



Operational Non-Technical Guidelines for AI Research Areas

AIOLIA DELIVERABLE 3.3

Horizon Europe Grant Agreement N° 101187937
AIOLIA PUBLIC

AIOLIA | D3.3



Project Name	AIOLIA
Deliverable Title/Number	D3.3
Description	Operational non-technical guidelines for three AI research areas
Lead beneficiary	CEPS
Lead Authors	Susana Aires, Nicoleta Kyosovska, Jacob Griffith, Sue-Anne Teo, Artur Bogucki
Contractual delivery date:	31/05/2026
Actual delivery date:	01/06/2026
Sensitivity	PUBLIC

Document History

Name	Organisation	Role	Action	Date
Narrative Highlights V1	CEPS	Lead	Draft sent to partners for validation	14 April 2026
Narrative Highlights V1	CEA, KIT, VUmc, CENTRIC, THWS, EUREC, ETICAS, UH, Oxipit, NIT	Contributors	Review and validation	14 April – 15 May 2026
Organisational Guidelines V1	CEPS	Lead	Draft sent to partners for validation	23 April – 29 April 2026
Organisational Guidelines V1	CEA, KIT, VUmc, CENTRIC, THWS, EUREC, ETICAS, UH, Oxipit, NIT	Contributors	Review and validation	29 April – 5 May 2026
Organisational Guidelines V2	CEPS	Lead	Updated version sent to SAB members for review	8 May 2026
Organisational Guidelines V2	AIOLIA SAB	Contributors	Review	8 May – 22 May 2026
Organisational Guidelines V2	CEPS	Lead	Updated version sent for validation and review by use case partners	11 May 2026
Organisational Guidelines V2	CEA, KIT, VUmc, CENTRIC, THWS, EUREC, ETICAS, UH, Oxipit, NIT	Contributors	Review and validation	11 May – 17 May 2026
Organisational Guidelines V3	CEPS	Lead	Updated version sent for review	13 May 2026

Organisational Guidelines V3	CENTRIC, EURACTIV	Reviewer	Full review	13 May – 21 May 2026
Organisational Guidelines V4	CEPS	Lead	Final version sent for review	27 May 2026
Organisational Guidelines V4	CENTRIC, EURACTIV	Reviewer	Final review	27 May – 31 May 2026
Organisational Guidelines V5 – pre-final version	CEPS	Lead	Final version ready for submission	31 May 2026
Final version	CEA	Coordinator	Final check	01 June 2026

Configuration Management

Nature of Deliverable	
R	Report

Dissemination level	
PU	Public, fully open

Acronym/abbreviations	
AI	Artificial Intelligence
AIA	AI Act
ALTAI	Assessment List for Trustworthy AI
DSS	Decision Support System
GPAI	General-Purpose AI
OM	Organisational Measure
PR	Policy Recommendation
RA	Research Area
UC	Use Case

How to cite
Aires, S., Kyosovska, N., Griffith, J., Teo, S. A., Bogucki, A. (2026). Operational Non-Technical Guidelines for AI Research Areas. AIOLIA Deliverable 3.3.

Acknowledgements

The information and views set out in this report are those of the author(s) and do not necessarily reflect the official opinion of the European Union. Neither the European Union institutions and bodies nor any person acting on their behalf may be held responsible for the use which may be made of the information contained therein. Reproduction is authorised provided the source is acknowledged.

Disclaimer

Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union. Neither the European Union nor the granting authority can be held responsible for them.

Use of AI

AI was used to support the comparative analysis of qualitative data, to formulate the narrative highlights present in sections 2.1.4, 2.2.4 and 2.3.4, and to revise the text. The authors remain responsible for the content of this deliverable.

EXECUTIVE SUMMARY

Deliverable D3.3 is the result of T3.3 *Develop context-enriched operational guidelines for AI research areas*. The aim of T3.3 was “to identify and lift overarching ethical concerns by cross-analysing use cases to implement the bottom-up co-creation of context-enriched non-technical AI ethics guidelines” per AI research area (cf. AIOLIA DoA). The AI research areas covered in this deliverable are General-Purpose AI, Emotional AI and Decision-Support. T3.3 was accomplished with recourse to the empirical data collected by AIOLIA partners during the operationalisation of AI ethics principles conducted in the context of T3.1, as well as a series of validation meetings and exchanges with AIOLIA partners and the SAB. D3.3 presents the process and results of identifying and lifting overarching ethical concerns and extrapolating organisational measures from AIOLIA use cases for the elaboration of research area-level guidance. Organisational measures refer to the structures, policies and governance frameworks that can be implemented for the ethical governance of AI within organisations providing and deploying AI systems. The deliverable presents the outputs first, methodological justification second.

The core resources emerging from T3.3 are:

1. Organisational Guidelines for AI Research Areas which include an assessment of open issues in the governance of AI within each research area as well as salient ethical concerns and a range of organisational measures for operationalising AI ethics principles. The General-Purpose AI organisational guidelines comprise 15 organisational measures, the Emotional AI organisational guidelines comprise 17 organisational measures and the Decision Support guidelines include 15 organisational measures.

2. Ethics tensions arising in the process of operationalising AI ethics principles and identified by AIOLIA partners. Ethics tensions are frictions and trade-offs between ethical principles wherein the operationalisation of one principle compromises another. These tensions point to the fact that AI ethics should be understood as an ongoing critical practice, rather than an issue to be fully resolved, providing important learning points for ethics practitioners, organisations and policymakers alike. In addition to acknowledging ethics tensions, T3.3 has produced narrative highlights, that is, fictional narratives modelled on partners’ data, to foster a deeper understanding among the learners taking part in AIOLIA’s WP4 pedagogical activities.

3. Recommendations for policymakers in view of establishing a feedback loop between AI ethics practices and policy, raising awareness of key challenges faced by those operationalising AI ethics guidelines on the ground and enabling evidence-based decision-making. D3.3 puts forward 10 policy recommendations.

Contents

1. INTRODUCTION	10
1.1. CONTEXTUAL BACKGROUND	11
1.2. AIMS & STRUCTURE	12
1.3. OVERVIEW OF ORGANISATIONAL MEASURES	13
2. OPERATIONAL GUIDELINES FOR AI RESEARCH AREAS	16
2.1. GENERAL-PURPOSE AI.....	16
2.1.1. Open issues in the governance of GPAI systems.....	16
2.1.2. Salient ethical concerns	18
2.1.3 Organisational measures for GPAI.....	20
2.1.4. Ethics tensions	35
2.2. EMOTIONAL AI	38
2.2.1 Open issues in the governance of Emotional AI systems	39
2.2.2 Salient ethical concerns	41
2.2.3 Organisational measures	43
2.2.4. Ethics tensions	57
2.3. DECISION SUPPORT.....	60
2.3.1. Open issues in the governance of DSS.....	60
2.3.2. Salient ethical concerns	62
2.3.3. Organisational measures	64
2.3.4. Ethics tensions	78
3. POLICY RECOMMENDATIONS	80
4. METHODOLOGY.....	84
4.1. Interaction between tasks in AIOLIA WP3	84
4.1.1. Co-creation process for technical guidelines in T3.1.....	84
4.1.2. International interaction and learning from T3.2.....	84
4.1.3. Ethics readiness levels in T3.4	85
4.2. Methodological approach	85
4.2.1. The ELSE framework: situating AI ethics within legal, societal and economic considerations	85
4.2.2. A bottom-up approach: analysing empirical data from AIOLIA use cases	89
4.2.3. Identifying and narrativising ethics tensions.....	92

4.2.4. From use cases to research area: extrapolating salient ethical principles and measures	93
5. OPERATIONALISING ETHICS PRINCIPLES.....	96
5.1. GENERAL-PURPOSE AI.....	96
5.1.1. The ELSE approach to GPAI.....	97
5.1.2. Overlaps and differences in the operationalisation of ethics principles across GPAI use-cases	101
5.1.3. Remaining challenges and concerns.....	115
5.2. EMOTIONAL AI	116
5.2.1. The ELSE approach for Emotional AI	117
5.2.2. Overlaps and differences in the operationalisation of ethics principles across Emotional AI use cases	120
5.2.3. Remaining challenges and concerns.....	128
5.3. DECISION SUPPORT	130
5.3.1 The ELSE Approach for Decision Support Systems	130
5.3.2. Overlaps and differences in the operationalisation of ethics principles across DSS use cases	132
5.3.3. Remaining challenges and concerns.....	148
6. CONCLUSION	150
7. REFERENCES.....	151

LIST OF TABLES

Table 1 - AI research areas and respective use cases and partners. _____	12
Table 2 - Overview of organisational measures and respective research areas. _____	15
Table 3 - Open issues in the governance of GPAI systems. _____	18
Table 4 - Overview of organisational measures for GPAI. _____	22
Table 5 - Overview of organisational measures for Emotional AI. _____	44
Table 6 - Open issues in the governance of DSS. _____	62
Table 7 - Overview of organisational measures for DSS. _____	65
Table 8 - Comparison of UC organisational measures addressing the principle of "Human Oversight". _____	103
Table 9 - Comparison of UC organisational measures addressing the principle of "Autonomy / User Agency". _____	105
Table 10 - Comparison of UC organisational measures addressing the principle of "Over-reliance and deskilling". _____	106
Table 11 - Comparison of UC organisational measures addressing the principle of "Technical robustness and safety". _____	107
Table 12 - Comparison of UC organisational measures addressing the principle of "Safety / Non maleficence / Human safety". _____	109
Table 13 - Comparison of UC organisational measures addressing the principle of "Privacy and Data protection". _____	110
Table 14 - Comparison of UC organisational measures addressing the principle of "Transparency and explainability". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure. _____	112
Table 15 - Comparison of UC organisational measures addressing the principle of "Non-bias, fairness and non-discrimination". _____	113
Table 16 - Comparison of UC organisational measures addressing the principle of "Human well-being". _____	114
Table 17 - Comparison of UC organisational measures addressing the principle of "Accountability and responsibility". _____	115
Table 18 - Comparison of UC organisational measures addressing human oversight, as put forth in Appendix D in D3.1. _____	121
Table 19 - Comparison table for organisational measures for the principle autonomy as put forth in Appendix D in D3.1. _____	123
Table 20 - Proposed organisational measures, as put forth in Appendix A in D3.1 under principle Autonomy. _____	123
Table 21 - Comparison table for the measures for principles safety, human well-being and non-maleficence (excluding oversight measures) as put forth in Appendix D in D3.1. _____	125
Table 22 - Comparison table for organisational measures for the principle privacy as put forth in Appendix D in D3.1. _____	126
Table 23 - Comparison table for the measures for the principle human well-being as put forth in Appendix D in D3.1. _____	128
Table 24 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Human oversight". _____	133
Table 25 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Over-reliance and deskilling". _____	135
Table 26 - Overview of organisational measures addressing the principle of "Freedom of expression and non-censorship" in UC4. _____	137
Table 27 - Overview of organisational measures addressing the principle of "Non-maleficence" in UC1. _____	138
Table 28 - Overview of organisational measures addressing the principle of "Robustness and reliability" in UC2. _____	139
Table 29 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Transparency and explainability". _____	142



Table 30 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Non-bias, fairness and non-discrimination". _____ 145

Table 31 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Accountability and responsibility". _____ 147

1. INTRODUCTION

The rise of artificial intelligence (AI) technologies and their diffusion in society has benefited different sectors, fields and industries, including healthcare, public administration, security and education. At the same time, governance challenges are present. These include navigating the impacts of AI on human rights and ongoing efforts to ensure that AI systems are Trustworthy, that is, that AI systems are designed and deployed in a manner that is legal, ethical and robust not only in order to minimise harms, but also to promote trust, responsible innovation and societal flourishing (European Commission – AI HLEG, 2019).

With this goal in mind, the AIOLIA project aims to operationalise AI Ethics for learning and practice. AI ethics principles and frameworks have proliferated around the world – such as frameworks from the G7, G20, UNESCO, OECD, and Global Partnership on AI. As of 2026, regulations on AI have also been adopted in different parts of the world, including in the European Union, the U.S., China and South Korea, with a view to ensure that ethical principles guide the design, development and deployment of AI technologies. This requires the operationalisation of high-level ethical principles, that is, their practical implementation, to bring about concrete impacts. Operationalisation, in this context, “refers to the process of translating high level ethical principles into practical actions, tools, processes and governance structures that can guide and be applied throughout the lifecycle of AI systems to ensure ethical design, development, deployment and use.” (see D2.3, p. 11). However, several challenges arise when translating principles into practice and, for this reason, AIOLIA develops concrete guidance to support developers, organisations and policymakers in this task.

Bringing together European and international partners, AIOLIA engages in ethics by design within AI use cases through a co-creation process involving both industrial and academic counterparts. The relevant ethical principles to be operationalised were identified in AIOLIA [deliverable 2.2](#) (Teo et al., 2025). This formed the basis for the co-creation process that followed, as the aforementioned partners worked on concretising ethical principles for their respective European AI use cases through technical and operational measures, as detailed in [D3.1](#) (Bayerl et al., 2026).

Task 3.3 translates this operationalisation process into concrete guidelines for the practical implementation of high-level ethics principles, by building upon and extrapolating from the co-creation process pursued in T3.1. In contrast with T3.1 which addresses the operationalisation of ethics principles through technical and organisational measures in a principle-based manner, T3.3 follows a research area approach, addressing three of AIOLIA’s research areas, namely, general-purpose AI (GPAI), emotional AI and decision support. This research area approach followed in T3.3 builds on the bottom-up co-creation process followed in T3.1 in order to identify and lift overarching ethical concerns through a cross-analysis of the European AI use cases. The aim is to produce context-enriched, non-technical guidelines per research area, centred on organisational measures, which are meant to guide the consideration and operationalisation of AI ethics principles according to the specificities of the different AI research areas.

Organisational measures focus on how an organisation incorporates and manages ethical AI practices by referring to the structures, policies and governance framework in place. Organisational measures for ethical AI governance include the development of AI ethics boards, ethical AI policies, promoting community stakeholder engagement, fostering interdisciplinary collaboration, regulatory and legal compliance to existing regulations, ethics readiness indicators and the development of AI risk frameworks” (D3.1, P. 23).

1.1. CONTEXTUAL BACKGROUND

T3.3 develops context-enriched operational non-technical guidelines for AI research areas. Initially in AIOLIA WP2 four main AI research areas (RAs) have been identified together with relevant partners. The selection process was informed by the need to identify applied AI research areas that are subject to the most impactful scientific advances, have the highest potential for quick scientific development, highest impact on human cognition and behaviour, and the most significant ethical challenges. The final selection of the research areas for AIOLIA task 3.3 limits the scope to the three following ones:

1. General Purpose Artificial Intelligence (GPAI)
2. Algorithmic Decision-Support
3. Emotional AI

These research areas are directly connected to the use cases (UCs), encapsulating the *type of AI systems* used by the industrial partners within their respective use cases. In contrast, use cases are characterised by the use of AI for specific tasks or discrete uses, such as the use of AI for recruitment purposes, AI for medical diagnostics and treatment, or AI for virtual companionship. A given use case can include one or more research areas. For example, one of AIOLIA’s medical use cases detects patterns in medical image data (image recognition) and identifies possible sites of medical intervention (decision support).

The table below lists the AI research areas and related use cases.

Research areas	Use cases	AIOLIA partners
General-purpose AI (GPAI)	UC2 - Safety engineers using AI tools to speed up software release approvals	CEA, NIT
	UC4 - Security professionals using AI tools to detect harmful or illegal content	CENTRIC
	UC5 - AI systems as personalised characters and individual virtual assistants	THWS
	UC6 - Deepfake therapy for processing trauma and grief	AUMC
Emotional AI	UC5 - AI systems as personalised characters and individual virtual assistants	THWS
	UC6 - Deepfake therapy for processing trauma and grief	AUMC

Decision-support	UC1 - Medical doctors (radiologists, surgeons) using AI tools in diagnostics and treatment	AUMC, Oxipit, Afliant
	UC2 - Safety engineers using AI tools to speed up software release approvals	CEA, NIT
	UC3 - HR professionals using AI tools to assess and reduce employee vulnerability to cyberattacks	Eticas
	UC4 - Security professionals using AI tools to detect harmful or illegal content	CENTRIC

Table 1 - AI research areas and respective use cases and partners.

1.2. AIMS & STRUCTURE

Following the rich process of operationalising AI ethics principles conducted by AIOLIA partners in T3.1, this deliverable builds on the bottom-up data collected, cross-analysing research area insights in dialogue with literature on ethical, legal, societal and economic dimensions of AI. The aim of this process is threefold:

- First, to provide guidance to organisations providing and deploying AI systems with regard to concrete measures that can be implemented at the organisational level to address key ethical concerns, including risks related to human autonomy, safety and deskilling.
- Secondly, the deliverable aims to speak directly to policymakers, pointing to open issues and challenges in the governance of AI across the different research areas, whilst providing concrete policy recommendations on the most pressing concerns identified in the context of AIOLIA.
- Thirdly, the deliverable aims to contribute to learning about AI ethics, demonstrating that ethics should be understood not as an issue to be resolved, but rather as an ongoing critical practice to be held within organisational, policy and educational contexts, through pointing to ethical tensions and trade-offs that necessarily emerge in the translation of AI ethics principles into practice.

With this in mind, the deliverable starts by putting forward the core insights arising from the research area-level analysis of the process of operationalisation of AI ethics principles conducted by AIOLIA partners, according to each research area (Section 2), namely General-Purpose AI (GPAI), Emotional AI and Decision Support (DSS). Within each section, a brief description of the research area is provided followed by an outline of open issues in the governance of the specific AI systems, which resonate with particular ethical challenges, and how AIOLIA can contribute to addressing these through its organisational measures. Salient ethical principles are then identified for each research area, based on the bottom-up data arising from AIOLIA use cases, pointing to key ethical concerns which require continuous engagement and attention from organisations.

These salient ethical concerns foreground the organisational measures outlined for each research area. For each organisational measure details are provided on its relevance for the specific research-area, practical guidance about its implementation, including potential implementation challenges and at what stages of the AI lifecycle it can be implemented. Measures have been organised around ISO/IEC42001, an international standard that seeks to facilitate the establishment, implementation and maintenance of AI

management systems within organisations. The alignment of the present organisational guidelines with international standards attests to AIOLIA’s commitment to actively contribute to the current landscape of AI policy and practice.

Building on the momentum around implementation challenges, section 2 concludes by identifying and discussing key ethical tensions arising in the process of translating ethical principles into concrete components and organisational practices. As D3.1 demonstrated, ethical principles are not self-contained, as they constantly interact and intersect with one another. In putting principles into practice and identifying practical measures to operationalise ethics principles, partners identified tensions between principles – that is, the operationalisation of one principle compromising another – and were required to make trade-offs. Herein, we have narrativized these tensions through fictional narratives derived from operationalisation challenges faced by AIOLIA’s use cases.

In light of the core insights outlined, section 3 puts forward a set of concrete recommendations for policymakers which point to potential solutions and actions that can help address the pressing needs and concerns identified by AIOLIA partners when it comes to operationalising AI ethics principles.

The final two sections of the deliverable address the research process adopted for its completion. Section 4 outlines the methodological approach, while section 5 delves into a detailed analysis of the operationalisation of ethical principles, components and measures undertaken by AIOLIA partners, with the aim of lifting overarching concerns and extrapolating organisational measures for each research area. The deliverable concludes with a short reflection on general findings and limitations.

1.3. OVERVIEW OF ORGANISATIONAL MEASURES

Following the methodology outlined in section 4 and the research process detailed in section 5, the core output emerging from the operational non-technical guidelines for AI research areas elaborated in D3.3 are the organisational measures listed below in Table 2. To understand the relevance of each measure for the respective research area, as well as how it should be implemented and the challenges organisations face in this process, consult the organisational measures outlined for each research area in section 2.

Measure		GPAI	Emotional AI	DSS
Context				
OM 01	Define clear boundaries and guidance for intended use of the AI system, including when AI outputs should be questioned, verified, or supplemented with human input.	X		X
OM 02	Define clear rules on the type of content, behaviour, or expression that are restricted for AI use and why, including distinguishing legitimate behavioural influence from manipulative practices.		X	
OM 03	Define what constitutes vulnerability in users or patients, including static and emergent forms, set out identification protocols, and establish corresponding enhanced safety measures, interaction protocols and continuous monitoring of effects for identified vulnerability profiles.		X	
OM 04	Define multi-level safety criteria and categories of risks and harms.		X	
OM 05	Consider diversity criteria and their intersection across roles, locations, languages, or cultural characteristics.	X		X
Leadership				
OM 06	Promote a safety-first culture within the organisation.	x		x
OM 07	Establish a representative AI governance committee as a formal oversight structure within the organisation.	x	x	x
Planning				
OM 08	Explicitly define and document responsibility for AI outputs.	X		X
OM 09	Clearly define oversight responsibilities and embed them into workflows, including checkpoints when active human professional engagement is required.	X		
OM 10	Establish a clear escalation path for the organisation hierarchy to decide on responsibility assignments during AI system design or operation.			X
OM 11	Define a clear policy on high-risk or sensitive decisions that must involve human judgment			X
OM 12	Conduct a privacy impact assessment.		X	
OM 13	Involve mental health professionals in the design, monitoring and evaluation of the AI system.		X	
OM 14	Determine the permitted degree of autonomy of the AI system (if any) and accordingly, specify the conditions in which a professional in the loop is required, e.g. via human moderation.		X	
OM 15	Consider psychological, emotional, and behavioural influences in AI risk analysis, not only user or patient information or technical AI performance.		X	

Support				
OM 16	Conduct trainings to support users in understanding AI capabilities and limitations, as well as related human capabilities and limitations.	X		X
OM 17	Clearly communicate how the AI system may affect human individuals, next of kin and society at large.		X	
Operation				
OM 18	Put in place a mechanism for users to report, question, or contest AI behaviour, decisions and/or restrictions.	X	X	X
OM 19	Design explanations that enable users to understand key behavioural factors, limitations, and uncertainties in their interaction with the AI system.	X		
OM 20	Document and enable traceability of changes to models, data or the GPAI system with sufficient detail to support internal reviews, external audits, or regulatory scrutiny.	X		
OM 21	Implement additional safeguards for sensitive data.		X	
OM 22	Ensure the AI system, if autonomous, applies behavioural nudging only in documented and auditable contexts with specified trigger thresholds and objectives.		X	
OM 23	Conduct internal ethics reviews for new features or changes of the AI system.		X	
OM 24	Design explanations that support meaningful review, contestation or justification of AI-supported decision.			X
Evaluation				
OM 25	Periodically reassess reliance patterns as systems evolve or scale	X		X
OM 26	Monitor situations when the purpose or actual use of the AI system drifts or diverges from the intended ones, and informs users when relevant.	X		X
OM 27	Put in place auditable Standard Operating Procedures for AI design and validation			X
OM 28	Conduct periodic independent ethics reviews of the AI system's impact on wellbeing and autonomy addressing all stages of the AI lifecycle.		X	
OM 29	Conduct audits of informed consent mechanisms and document limitations, especially with regard to the adaptivity to context.		X	
OM 30	Review AI outputs for unintended disparity impact.	X		
Improvement				
OM 31	Translate audit findings into corrective design action	X		X
OM 32	Identify and analyse unforeseen effects of the AI system on individual and societal well-being	X	X	X
OM 33	Assess and collect feedback from users of conversational systems regarding perceived honesty, non-coerciveness of AI interactions and impact on wellbeing.		X	

Table 2 - Overview of organisational measures and respective research areas.

2. OPERATIONAL GUIDELINES FOR AI RESEARCH AREAS

2.1. GENERAL-PURPOSE AI

General-purpose AI (GPAI) is a category of AI models and systems that display significant generality, are capable of competently performing a wide range of distinct tasks and can be integrated into a variety of downstream applications. These advanced AI systems exhibit a significant degree of autonomy and the ability to generalise to new tasks and across domains the system has not been previously exposed to or intentionally trained to address. Despite GPAI systems being characterised by a high degree of operational opacity, their utility and ability to perform a wide range of tasks, identifying patterns across multimodal datasets, has motivated GPAI uptake across private and professional contexts.

2.1.1. Open issues in the governance of GPAI systems

As a research area, GPAI presents unique ethical challenges. On the one hand, the generality of tasks GPAI models can perform, both in professional and private contexts, means that there is no single, standard range of ethical principles and measures that would apply to this research area. On the other hand, the fact that GPAI models, as foundation models developed by third parties, often form the baseline architecture upon which different applications are built, fine-tuning the foundation model to perform in specific contexts, poses new challenges to the operationalisation of ethics-by-design in the contexts of model deployment. Below are some of the key open issues for the governance of GPAI faced by organisations and the ways in which AIOLIA’s organisational guidelines can support addressing these.

Open issue #1: How can organisations govern emergent risks that arise from GPAI deployment rather than at the system development stage?	
<p>Concern: GPAI systems are highly malleable and their behaviour evolves dynamically through human-AI interaction. This means risks cannot be fully anticipated or mitigated through ex-ante approaches alone, and may only materialise within specific deployment contexts, including non-technical harms such as deskilling, that fall outside conventional risk frameworks. Risks stem not just from technical robustness, but from the unpredictable, emergent dynamics of human-AI interaction.</p>	<p>How organisational measures help: Organisations can define clear boundaries for intended use and monitor for drift or divergence from those boundaries. Internal ethics reviews for new features, combined with traceability documentation of system changes, allow organisations to detect and respond to emergent risks. Governance committees provide a structural safeguard against the organisational pressures that can accelerate deployment at the expense of safety.</p>
Open issue #2: How should organisations address deskilling and the long-term erosion of human competencies resulting from GPAI use?	
<p>Concern: The reliance on natural language interfaces shields users from cognitively engaging</p>	<p>How organisational measures help: Organisations can establish governance</p>

<p>with tasks, and productivity gains from GPAI deployment may come at the cost of gradual erosion of skills. This risk compounds over time and may become embedded in organisational practices before it is recognised, undermining not only professional competence but the quality of human oversight itself and safety of AI-based workflows and decisions.</p>	<p>committees that act as safeguards against unrealistic KPIs and workloads, implement role-based training programmes that preserve and update human competencies, and periodically reassess reliance patterns as systems evolve. Oversight responsibilities should be explicitly embedded in workflows so that meaningful human engagement is maintained by design rather than assumed.</p>
<p>Open issue #3: How can governance of GPAI be reoriented beyond safety- and compliance-centric frameworks to address individual and societal well-being?</p>	
<p>Concern: The dominant safety agenda in GPAI governance is largely framed around catastrophic or existential risks, leaving out a broad range of current and concrete harms, including algorithmic discrimination, cognitive impacts, and erosion of social trust. This framing risks making governance appear futile and out of reach outside the realm of frontier AI providers, narrowing the aperture of ethical deliberation in ways that exclude actual deployment contexts.</p>	<p>How organisational measures help: Organisations can adopt ethics-by-design approaches, including diversity and fairness reviews, impact assessments on societal well-being, and deliberative processes with civil society and impacted communities. Participation in broader research consortia and the publication of research findings contribute to widening the evidence base beyond capability-focused development.</p>
<p>Open issue #4: How can the absence of stable regulatory frameworks and harmonised standards be managed at the organisational level?</p>	
<p>Concern: Organisations operating in less standardised domains lack a clear baseline for translating ethical principles into practice. The AI Act's obligations are largely directed at providers of high-risk systems, leaving deployers with significant discretion in determining what responsible deployment looks like. The absence of harmonised standards for GPAI, including in emerging areas such as Emotional AI, makes consistent and auditable governance difficult to achieve.</p>	<p>How organisational measures help: Organisations can adopt governance practices that go beyond legal compliance, effectively self-applying high-risk standards in the absence of a regulatory requirement to do so. Governance committees and independent ethics reviews can contribute to the development of sector-specific best practices, although these also require significant resources that might not be available among small organisations. Sharing safety criteria and research findings with professional and academic communities supports the emergence of shared norms.</p>
<p>Open issue #5: How should the environmental impact of GPAI be governed?</p>	
<p>Concern: There is growing evidence of the significant environmental costs of GPAI, including energy and water consumption and raw material use associated with model training and inference.</p>	<p>How organisational measures help: Environmental concerns are largely beyond the reach of individual organisations and require intervention at the regulatory and policy level,</p>

<p>However, environmental concerns have not emerged as a central issue in AIOLIA use cases, not because of a lack of awareness, but because organisations tend to prioritise the challenges they can directly tackle, such as safety, oversight, and skill preservation. Standard indicators for measuring environmental impact remain absent, and the trade-offs between competitiveness and environmental sustainability have yet to be clearly articulated at the policy level, leaving organisations without actionable guidance.</p>	<p>including the formulation of standard measurement indicators and clear sustainability priorities. Where possible, organisations can contribute by participating in broader research and policy communities and factoring environmental considerations into procurement and system design decisions. This underscores the limits of organisational governance and the need for complementary public and regulatory action.</p>
---	--

Table 3 - Open issues in the governance of GPAI systems.

In light of these governance challenges, it becomes paramount to trace the most relevant ethical principles and operational measures that can act as ‘enablers’ of Trustworthy AI in the context of GPAI, as these systems become increasingly integrated into professional and private contexts, from workflows to day-to-day activities. Despite the differences in the nature of ethical challenges, depending on whether GPAI systems are deployed in professional or private domains, these guidelines seek to inform responsible GPAI use across deployment contexts and, for this reason, incorporate elements that resonate with both domains.

2.1.2. Salient ethical concerns

While AI ethics principles, such as those underlying the seven requirements of the Assessment List for Trustworthy AI (ALTAI)¹, remain relevant in the context of GPAI, some become more nuanced and salient when applied to GPAI. In fact, AIOLIA use cases have identified as many as twelve relevant ethics principles in the different contexts of deployment of GPAI (see Section 5.1.), that form the basis of the organisational measures advanced herein.² However, during the operationalisation of AI ethics principles conducted by the four AIOLIA GPAI use-cases, three ethics principles revealed themselves to be particularly salient for GPAI, namely: 1) human agency and autonomy, 2) over-reliance and deskilling, and 3) human safety and non-maleficence.

Human Agency and Autonomy

GPAI systems pose new challenges to human agency and autonomy, both in professional and private domains. From the perspective of professional behaviour, autonomy is about agency or the ability to act over an AI system. It is regarded as an enabler of human oversight, and by extension, integral to

¹ The seven ALTAI requirements are: 1) human agency and oversight, 2) technical robustness and safety, 3) privacy and data governance, 4) transparency, 5) diversity, non-discrimination and fairness, 6) environmental and societal well-being and 7) accountability.

² The twelve ethics principles identified by AIOLIA’s GPAI use cases are: 1) Human Oversight, 2) Autonomy/User agency, 3) Over-reliance and deskilling, 4) Freedom of expression and non-censorship, 5) Robustness and reliability, 6) Safety/Human Safety, 7) Non-maleficence, 8) Privacy and Data Protection, 9) Transparency and explainability, 10) Non-bias, fairness and non-discrimination, 11) Human Well-being, 12) Accountability and responsibility. See AIOLIA deliverable D3.1 “[Operational ethics guidelines on use cases related to human behaviour and cognition](#)”.

accountability and safety. The emergence of GPAI systems with significant capabilities and generality, as well as their integration into workflows, complicates the fulfilment of autonomy and agency, as these systems perform a growing number of tasks that, due to their complexity, previously required human intervention. While practitioners can be given control over GPAI systems – a usual element of AI guidelines – significant limits remain regarding the ability to exert meaningful control and oversight, particularly due to the opacity of GPAI functioning and the long-term risk of deskilling that the automation of tasks can accelerate, ultimately jeopardising human ability to act in face of these complex systems.

In the context of private behaviour, autonomy gains new contours, being intrinsically linked to freedom of choice and non-manipulation. It is regarded not merely as the ability to choose here and now, but as the preservation of this faculty in the long-term, despite the potentially lasting behavioural, cognitive and emotional imprint resulting from the subtle and cumulative interactions with GPAI systems. While manipulative AI systems are prohibited under the AI Act, establishing the boundaries between interactions that amount to manipulation and those that do not, remains outstanding, especially in the context of extended and often individual and intimate, interactions with GPAI systems that exhibit anthropomorphic characteristics.³ While cognitive attachment to GPAI systems is often understood as detrimental to user agency and autonomy, the removal or limitation of anthropomorphic features and interfaces might also place particular users at a disadvantage, including users in cognitive decline, isolated or with disabilities.

Overreliance and Deskilling

The risk of overreliance and deskilling is a key ethical concern highlighted in AIOLIA GPAI use cases. The generality of tasks GPAI systems can perform and their translation into natural language interfaces for human oversight can jeopardise the ability to act, hence, to exert responsible and effective control over GPAI outputs. Most notably, the lack of day-to-day engagement with the technicalities of professional practice (e.g., engineers) can lead to the erosion of knowledge, skills and behaviour required for effective human oversight and professional agency, resulting in overreliance on the GPAI system, including through anchoring bias, that is, relying too heavily on the first piece of information provided by the system. This can have a cascading effect on organisational practices, including organisational safety culture and responsible AI use. To mitigate these challenges AIOLIA proposes organisational measures that address the elements of continuous training and skill development alongside the evolution of the AI system, as well as the establishment of AI governance committees to mitigate organisational pressures (e.g., productivity KPIs) that can erode attention and aggravate automation bias. In the context of private behaviour, recurrent engagement with and reliance on GPAI systems can lead to erosion of cognitive autonomy and social deskilling.

However, AIOLIA use cases also raised the question of how the concern with the preservation of skills might be foregrounded on the idea that the human is the relevant standard for measuring performance, whereas it is currently well-documented that human-AI collaboration outperforms human skills alone. It follows that keeping the human as the standard might risk eroding the benefits of GPAI-based decision

³ For additional insights on GPAI systems deployed in the context of affective computing see the AIOLIA Organisational Guidelines for Emotional AI.

support. The challenge lies in identifying the skills and competencies required for responsible use and oversight of GPAI systems as human-AI collaboration evolves and is shaped over time.

Human Safety and Non-maleficence

Another ethical principle highlighted in AIOLIA's GPAI use cases, is the concern with human safety and non-maleficence, extending an ethics principle traditionally linked to *technical* robustness and reliability to the domain of *human* safety and non-maleficence. In the context of GPAI, the concern lies both in the malleability of GPAI models and systems and the risks of faulty GPAI outputs or errors, both of which can lead to serious harm to human safety. With regard to malleability, GPAI models are trained on large datasets and based on probabilistic architectures that inherently result in the emergence of new patterns of functionality over time through human-AI interactions. This poses challenges to safety assurance, as GPAI systems may drift or diverge into uncharted behaviour that might evade previous performance indicators and pose new risks both in professional and private contexts. The fact that, in the context of GPAI systems, safety risks and harm can emerge through human-AI interactions requires the establishment of criteria as a basis for ongoing refinement as new harm patterns emerge.

Likewise, system errors can have a significant impact on safety that is not simply restricted to the context of GPAI deployment but rather impacts broader domains in which faulty GPAI outputs can be detrimental and propagate at scale, especially when these remain undetected. This is particularly crucial in light of agentic AI systems and the fact that GPAI models increasingly form part of decision support systems.⁴ For example, in domains like safety engineering, flawed AI-assisted decision-making processes lead to safety risks and harm to humans that expand to users or consumers of end-products (e.g., cars, road users), rather than remaining circumscribed to the context of GPAI deployment. These increased risks of unsafe and unfair system behaviours require the establishment and embeddedness of adequate human oversight across organisational workflows.

As such, GPAI systems require adopting a dual approach at the organisational level: first, clearly defining the scope of the GPAI system, traditionally characterised by generality and adaptability, to implement adequate technical guardrails and enhance the professional understanding of the system's functionality; secondly, establishing internal mechanisms for reviewing human-AI interactions, including system customisation and emergent risks, attending to the characteristics of GPAI users to mitigate harms in the professional and private domains, from deskilling, to psychological suffering.

2.1.3 Organisational measures for GPAI

Against the backdrop of these overarching concerns, we present a list of organisational measures that specifically address them, and can serve as a starting point of reference for organisations aiming to embed ethics into GPAI design, development and/or deployment.

Note that the list of measures does not aim to be comprehensive but foundational. In other words, the implementation of all measures under a principle does not guarantee the principle has been upheld. This is because, firstly, the table is the result of identifying and translating UC measures to the RA-level from

⁴ For additional insights on GPAI models deployed in the context of decision support see the AIOLIA Organisational Guidelines for Decision Support Systems.

D3.1., that includes non-exhaustive guidelines, which implies that the resulting RA guidelines will inherently lack comprehensive coverage. Secondly, the measures are chosen to address the most pressing concerns, especially those that may fall out of scope of current regulation. This means that the selection amongst UC measures does not include all possible measures that can be applied to the research-level area, but those that are most novel and pressing.⁵

The organisational measures for GPAI listed in this section are informed by the operationalisation of AI ethics principles across four AIOLIA GPAI use cases, operating across areas as diverse as healthcare, automotive engineering, detection of harmful and illegal content, and personal companions. Given the different areas of GPAI deployment covered, the operationalisation process undertaken in the different industrial use cases has been informed by sectoral regulatory frameworks and standards, such as ISO automotive industry standards or the Medical Device Regulation, in addition to the AI Act. This has resulted in a wealth of different approaches and measures which has enriched the current guidelines. The on-the-ground reality of industrial use cases, from the selection of relevant AI ethics principles to the formulation of measures to address these, has been systematised to address the GPAI research area as a whole. This process was conducted through mapping commonalities and differences across the use cases with the aim of identifying overarching ethical concerns and operational measures, applicable not simply to the specific context of AIOLIA GPAI use cases, but capable of being extrapolated to any GPAI technology – see section 5.1. of this deliverable for a detailed account of this process. Next to relevance and balanced coverage across the most pressing ethical concerns specific to the research area, the measures have been selected for originality in addressing emerging challenges reflected in regulatory frameworks and gaps. As a result, the Organisational Guidelines for GPAI surface what is most novel and actionable for this context.

OVERVIEW OF ORGANISATIONAL MEASURES FOR GPAI	
Context	
OM01	Define clear boundaries and guidance for intended use of the AI system, including when AI outputs should be questioned, verified, or supplemented with human input
OM05	Consider diversity criteria and their intersection across roles, locations, languages, or cultural characteristics
Leadership	
OM06	Promote a safety-first culture within the organisation
OM07	Establish a representative AI governance committee as a formal oversight structure within the organisation
Planning	
OM08	Explicitly define and document responsibility for AI outputs
OM09	Clearly define oversight responsibilities and embed them into workflows, including checkpoints when active human professional engagement is required
Support	
OM16	Conduct trainings to support users in understanding AI capabilities and limitations, as well as related human capabilities and limitations

⁵ Note that some measures do not appear in the analysis outlined prior, because they were the result of subsequent identification of measures for Appendix A in D3.1. which does not differentiate between use-cases.

Operation	
OM18	Put in place a mechanism for users to report, question, or contest AI behaviour, decisions and/or restrictions
OM19	Design explanations that enable users to understand key behavioural factors, limitations, and uncertainties in their interaction with the AI system
OM20	Document and enable traceability of changes to models, data or the GPAI system with sufficient detail to support internal reviews, external audits, or regulatory scrutiny
Evaluation	
OM25	Periodically reassess reliance patterns as systems evolve or scale
OM26	Monitor situations when the purpose or actual use of the AI system drift or diverge from the intended ones and inform the user when relevant
OM30	Review AI outputs for unintended disparity impact
Improvement	
OM31	Translate audit findings into corrective design action
OM32	Identify and analyse unforeseen effects of the AI system on individual and societal well-being

Table 4 - Overview of organisational measures for GPAI.

The organisational measures listed herein have been organised based on the ISO/IEC42001 standard for implementing an AI management system within organisations, structured around seven dimensions: 1) context, 2) leadership, 3) planning, 4) support, 5) operation, 6) performance evaluation, and 7) improvement. The same structured is adopted in the current guidelines. For each dimension, concrete organisation measures (OM) are provided: each entry specifies what the measure consists of; why it is important for the specific research-area context, including how it addresses the area's salient ethical principles; guidance about its implementation, including at what stages of the AI lifecycle it can be implemented; and core challenges organisations might face in the implementation of each measure. The AI lifecycle follows the OECD categories: 1) design, data and modelling, 2) verification and validation, 3) deployment, and 4) operation and monitoring.⁶

CONTEXT

This section outlines measures to address the specific context of the organisation, including internal and external contexts, ranging from legal requirements to competitive landscape, as well as cultural and social values. Understanding the organisational context also requires accounting for all the parties impacted by the AI system, from direct users to wider communities.

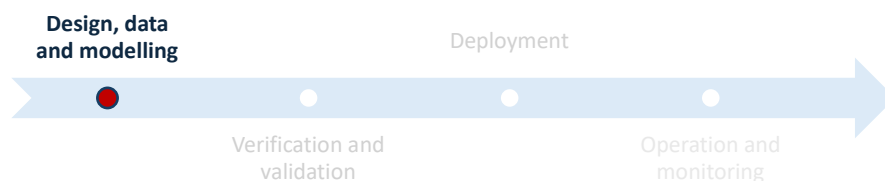
- OM 01.** **What:** Define clear boundaries and guidance for intended use of the AI system, including when AI outputs should be questioned, verified, or supplemented with human input.
- Why:** Clearly determining the role and intended use of the system within the organisation is key for both the GPAI system and human behaviour to operate within established safeguards

⁶ For more information on the phases of the OECD AI Lifecycle see https://www.oecd.org/content/dam/oecd/en/publications/reports/2019/11/scoping-the-oecd-ai-principles_71e1b6dc/d62f618a-en.pdf

and boundaries, ensuring **human safety**. By providing guidance on the intended use of the GPAI systems, organisations create a shared organisational understanding of the AI systems and its role, enhancing the ability for users to perform adequate **human oversight**, understanding inconsistencies in GPAI behaviour, as well as questioning and overriding AI outputs in a manner that preserves **human agency and autonomy**. This enables a systematic and uniform way of using the system across the organisational roles and functions.

How: This measure can be achieved through establishing a formal organisational policy on GPAI use, based on internal and cross-disciplinary deliberation, from technical experts and ethicists, among others. The policy should state the GPAI system’s supported capabilities and scope boundaries based on key performance metrics (e.g., accuracy, recall, false-positive rate, confidence intervals). Additionally, the AI use policy shall outline adequate use practices, including when AI outputs must be questioned, verified or supplemented. Use-cases can be included to exemplify adequate uses, including prohibited practices, unbiased use, human-AI decision-making standards, and escalation mechanisms. The policy should be widely distributed and communicated to all the relevant roles across the organisation, and be reflected in technical documentation and measures, as well as on human actions and behaviour (e.g., content moderation). The AI policy should be developed keeping in mind the level of expertise of the diverse range of stakeholders that might access it, with a view of ensuring understandability. Regularly review and update the policy considering regulatory developments, internal audits to safety monitoring, user feedback, technical developments, and any other relevant evidence.

Implementation challenges: Defining clear and meaningful technical performance boundaries and Operational Design Domains (ODDs) for complex, data-driven GPAI models is inherently difficult – a technical challenge that might also emerge in attempts to delineate a formal policy on GPAI use. This complexity is compounded by the fact that not every borderline case is predictable, often relying on *post hoc* assessment and manual intervention. Resolving the ambiguity regarding borderline cases can be further stalled by organisational friction, such as an insufficient dissemination of the policy to managers. Lastly, organisations, particularly those providing GPAI models, also face a delicate balancing act since strict policies on GPAI use carry the risk of competitive disadvantage compared to other, less-restrictive GPAI providers operating on the market.



OM 05. **What:** Consider diversity criteria and their intersection across roles, locations, languages, or cultural characteristics.

Why: Whereas access to the design process of foundation models is often restricted, paying attention to and assessing **diversity** criteria of GPAI systems becomes particularly relevant, to ensure that they adhere to the specificities of deployment contexts. Testing and tailoring GPAI systems to deployment contexts is paramount for ensuring AI deployment is **fair, non-**

discriminatory and safeguards **freedom of expression**. This measure aims to promote proactive organisational engagement, assessment and ownership in the systematic prevention, detection, and remediation of discriminatory biases and related risks, unlawful and disproportionate impacts of GPAI systems, beyond ad-hoc technical fixes.

How: This measure can be achieved through 1) mapping the context of GPAI deployment to inform system design or fine-tuning and related practices, paying particular attention to asymmetries of power, knowledge, or vulnerability between humans and the AI system. 2) Appointing an AI Fairness/Ethics Lead. 3) Creating a cross functional review body including technical, legal and diversity experts, as well as external stakeholders, to conduct GPAI alignment checks across its lifecycle. The review body must include experts able to address the context of GPAI deployment, intended purpose and technicalities. For example, if the GPAI provides a language-based clinical service, linguists, medical experts and civil society organisations (e.g., patient organisations) must be appointed and involved in the assessment, so that all dimensions of the GPAI system are adequately addressed. 4) Incorporating GPAI reviews into system design by liaising with technical teams following each review body meeting, and mapping action points arising from the review process to technical changes in the GPAI system, keeping records of this process. 5) Establishing an organisation-wide policy on impartial GPAI use and GPAI-assisted decision-making standards that translates both into human practices and system functionality.

Implementation challenges: Smaller organisations, with more limited financial and logistic resources, may struggle to identify, recruit and convene a representative group of experts in the cross-functional review body. From a logistical point of view, there is the risk of scheduling friction, limited capacity for repeated sessions across locations, and difficulties in reaching non-traditional participants or stakeholders. Language diversity may incur additional costs, requiring dedicated translation resources. Internal and / or experts involved in the cross-functional review body might have conflicting priorities (e.g., security vs diversity) and disagree on their assessment of diversity criteria, requiring careful negotiation to ensure all viewpoints are accommodated and represented in the discussions and a consensus is reached with regard to GPAI design changes to be implemented. Lastly, tight timelines for model updates and other organisational pressures can lead to bypassing critical fairness/ethics checks.



LEADERSHIP

This section outlines measures to address the leadership commitment of the organisation towards the AI management system, including through defining a policy for AI systems deployed within the organisation, but also assigning clear roles and responsibility for managing these.

OM 06. **What:** Promote a safety-first culture within the organisation.

Why: The aim is to embed and promote adherence to **safety, human oversight** processes and attitudes across the organisation, reducing the likelihood of error, oversight gaps and the tendency for overlying on GPAI outputs. This enhances organisational and individual **accountability** for GPAI outputs and enhances the **robustness and safety** of GPAI-assisted tasks and workflows.

How: To implement this measure, steps must be taken across 1) AI design, through embedding traceability and control mechanisms, 2) documentation, by clearly assigning roles and responsibilities across the decision making process, and 3) capacity-building, including through mentorship, learning-by-doing and training the trainers programmes.

Implementation challenges: The promotion of a safety-first culture within organisations can be jeopardised by the lack of awareness of AI-specific risks, inconsistent adherence to processes and also resistance to change and adapt to new organisational procedures. Furthermore, while promoting a safety-first culture is a pressing need among organisations, determining what constitutes this culture is more challenging, since this measure is not measurable nor achievable simply through aggregating the practices described under ‘how’. This leads to challenges in assessing safety culture and measuring progress.



OM 07.

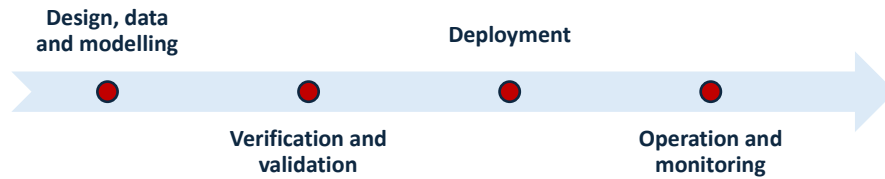
What: Establish a multidisciplinary AI governance committee as a formal oversight structure within the organisation.

Why: The integration of GPAI systems into workflows requires collective **oversight** and management of **accountability** for the AI system across organisational levels and roles. For this reasons mechanisms such as an organisational AI governance committee must be put in place to represent and give voice to professionals impacted by the GPAI system. Such a mechanism acts as a safeguard against organisational pressures and conflicting objectives, such as speed and unrealistic KPIs, that might be detrimental to the safe deployment of the AI system and the maintenance of adequate **oversight** skills in the long term.

How: Organisational leadership must initiative procedures to establish the AI governance committees and determine its mandate. Representatives of the AI governance committee shall be appointed and election procedures should be initiated to determine its membership, including an AI Ethics Lead, ensuring adequate representation across organisational roles, seniority level, and areas expertise. Members of the AI governance committee shall establish periodic meetings (at least on a monthly basis) and keep minutes of discussions. The organisation shall facilitate ongoing avenues for communication between the AI governance committee and organisational leadership.

Implementation challenges: Smaller organisations might have limited financial and human capacity to establish an AI governance committee with interdisciplinary expertise. Across all organisations, the effectiveness of the AI governance committee can be comprised by

leadership / management, which might limit its mandate and resource allocation, especially in the case of conflicting priorities, such as the demand for rapid software releases and unrealistic KPIs and project deadlines.



PLANNING

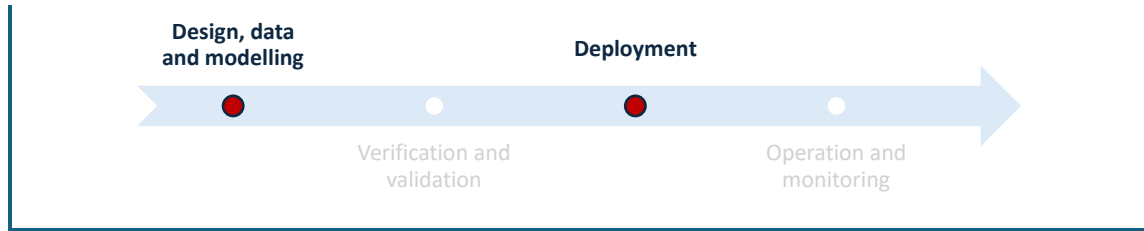
This section outlines measures to address the planning phase of developing and / or deploying an AI system within the organisation. This involves identifying and assessing risks, potential impacts and opportunities arising from the AI system and setting clear objectives for its deployment.

OM 08. **What:** Explicitly define and document responsibility for AI outputs.

Why: Given the generality of tasks that GPAI models can perform, with different levels of accuracy and predictability, ensuring that users take ownership of and **accountability** over GPAI outputs is key to the responsible and safe use and deployment of GPAI systems. By explicitly defining and documenting responsibility for GPAI outputs, organisations are ensuring that adequate level of quality and safety assurance, as well as **human oversight** is performed across the organisation.

How: To achieve this measure organisations shall first identify organisational roles making use of GPAI systems' outputs, as well as the degree of responsibility they should take, including establishing escalation paths for complex cases (e.g. DPO, legal). This shall be followed by drafting and approving role-based responsibility assignment for the identified occupations, integrating responsibilities into job descriptions and project charters. The role-based responsibility assignment should be complemented by technical measures embed in GPAI design, including role-based access control systems, digital sign-off procedures and traceability of decision-making processes.

Implementation challenges: At the level of individuals, organisations might face resistance from GPAI users to be held responsible for AI outputs, especially given the lack of predictability of GPAI systems, the inherent opacity of their operation and the lack of participation in system building. The implementation of this measure can also be jeopardised due to overlapping responsibility and mandates between teams or matrix management, complicating the allocation of responsibility for GPAI outputs. The allocation of role-based responsibility assignment is also challenging to maintain over time, requiring resources and commitment to update responsibility assignments as teams and roles change over time (e.g., turnover of key staff), and the complexity of cases requiring escalation change post-deployment.



OM 09. **What:** Clearly define oversight responsibilities and embed them into workflows, including checkpoints when active human professional engagement is required.

Why: While the shape, frequency and feasibility of human oversight can differ according to the nature of workflows, **human oversight** remains paramount for ensuring GPAI systems operate as expected, promoting **accountability** and preserving **human safety**. This might include managing complex or borderline cases that automated systems cannot reliably resolve, such as distinguishing legitimate user preferences from harmful behaviour, and identifying biased system behaviours that can propagate at scale and beyond the narrow context of GPAI deployment (e.g., GPAI-assisted car safety assurance). Documenting decision points in advance ensures that operational consistency, accountability, and alignment with regulatory obligations are built into the AI system's operation by design rather than applied reactively.

How: Organisations should start by mapping, identifying, flagging and documenting high-risk or sensitive decisions across the AI lifecycle, determining key moments for human oversight. This should be followed by the documentation and design of systematic procedures for human review, approval, or override, often through technical means, which must be communicated and implemented across relevant organisational roles. Organisational procedures for addressing high-risk or sensitive decision must be based on risk severity, impact on individuals, and regulatory obligations. Train a team dedicated to overseeing the GPAI system and assign supervisory responsibilities to monitor oversight practices and procedures.

Implementation challenges: Balancing key points for human oversight and accountability with workload efficiency is a major challenge for organisations incorporating GPAI into workflows. The allocation of temporal resources for operationalising oversight is another critical element that might compromise the effectiveness of implementation. From the point of view of human overseers, organisations might face resistance from GPAI users for taking on board new responsibilities, especially in time-constrain domains, and alert and interface fatigue can effectively compromise GPAI oversight. Finally, while documentation is critical for implementing this measure, organisations may struggle to maintain clear oversight responsibilities when teams and roles change over time, requiring swift organisational response.



SUPPORT

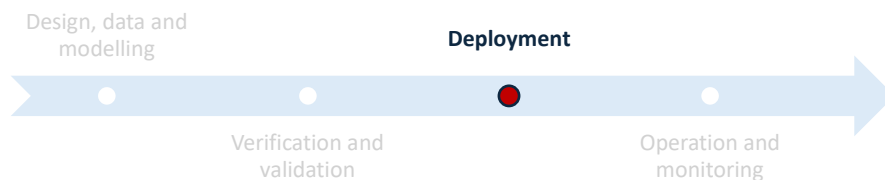
This section outlines measures to support the establishment, implementation and maintenance of the AI management system within the organisation. It involves determining and providing organisational resources, ranging from competencies, awareness and clear communication and documentation.

OM 16. **What:** Conduct trainings to support users in understanding AI capabilities and limitations, as well as related human capabilities and limitations.

Why: The aim of this measure is to mitigate gaps between AI system capabilities and human skills, both current and those emerging from the long-term incorporation of the AI system into organisational workflows. This involves not only understanding how to use GPAI systems and their limitations, but how to do so in a way that prevents harm. In professional behaviour this requires training on **human oversight** to ensure that users can interpret GPAI outputs, maintain responsibility, remain effective supervisors, and maintaining the competencies required for detecting faulty or biased outputs, hence mitigating **overreliance and deskilling**.

How: This measure can be achieved through the implementation of a modular and role-based training program with recurring frequency, combining technical modules (AI literacy, model interpretation) and situational modules addressing the risks and practices specific to the context of GPAI deployment. Trainings should cover the onboarding, integration and functional periods of practitioners within the organisation through training boosters. Organisations should assess employees' comprehension of GPAI functioning and limitations pre- and post-training. Keep and audit training documentation for completeness. Organisations should allocate time and resources for training. Update training resources in line with GPAI systems' evolution.

Implementation challenges: The allocation of significant temporal, financial and human resources for training is a core challenge for organisations implementing this measure, especially amid tight project deadlines and productivity KPIs. The need for maintaining up-to-date training materials and tutors is also a key challenge, given the pace of GPAI updates that might potentially take place and require the adaptation of training programmes. From the perspective of GPAI users / organisation employees, they might lack the motivation to continuously enrol in training programmes, being required to balance training time with other responsibilities.



OM 18. **What:** Put in place a mechanism for users to report, question, or contest AI behaviour, decisions and/or restrictions.

Why: Users of GPAI systems may experience unexpected restrictions or system behaviours, such as when a GPAI system drifts or diverges from its intended use purpose. While organisations should monitor this kind of situations, users must also be given the possibility and **agency** to report unexpected GPAI behaviours, decision and / or restrictions, especially given the safety risks emerging from the malleability and customisation of GPAI systems based on use patterns. This measure is aligned with **explainability** obligations for deployers of high-risk AI systems under Article 86 of the AI Act.

How: Upon deployment of GPAI systems, organisations must implement dedicated channels through which users can report concerns, ask questions, or contest system behaviours and decisions, such as in-app reporting tools, accessible customer service contacts, structured feedback forms, or internal organisational channels. Responsibility for user communications should be allocated to a specific team. Responses should be timely, communicated in plain language, and accompanied by a meaningful explanation of the system behaviour in question. Moreover, where relevant, organisations should comprehensibly explain the issues reported and announce the envisaged GPAI system corrections (e.g., design modifications, interface updates, new interaction protocols). All reports and contestations should be logged, reviewed regularly, and fed back into system improvement processes.

Implementation challenges: The implementation of this measure requires significant human resources to review user reports, questions and appeals, potentially delaying the organisational workflow, especially among small organisations.



Operation

This section outlines measures to support the operation of the AI system within the organisation. It involves operational planning and control, and conducting risk and impact assessment, planning effective means for corrective action.

OM 19. **What:** Design explanations that enable users to understand key behavioural factors, limitations, and uncertainties in their interaction with the AI system.

Why: Explanations enhance operational **transparency** and user understanding of the GPAI system. In the GPAI context, **explainability** is particularly relevant to enable users to make sense of the limitations and the not fully predictable and containable nature of these systems. This allows GPAI users to adjust their behaviour and interactions with the system by, for example, being more vigilant in GPAI-assisted decision-making and gaining cognitive distance from GPAI system that exhibit anthropomorphic characteristics. Under the AI Act, deployers of high-risk AI systems are required to provide explanations to persons negatively affected by the systems' outputs, through "clear and meaningful explanations" on the their role in the decision-making process (AI Act, Article 86).

How: This measure can be implemented through designing explanations adapted to the contexts of GPAI deployment, as well as the role and expertise of GPAI users. For example, explanatory interfaced might rely on staged reveal of GPAI suggestions to enable users to interpret and reason together with the processual operation of the system. Explanations may also include relevant information regarding the level of confidence of the GPAI system and prompt active user engagement. When designing explanatory interfaces, organisations must account for relevant technical measures that enable meaningful user agency, such as the override of GPAI outputs or flagging for secondary review, whilst maintaining monitoring and recording of outputs flagged or overridden for organisational learning and GPAI improvement.

Implementation challenges: Since explainability is an ongoing area of research and an outstanding issue for GPAI models, there is not a single technical means for approaching this issue, nor a measurable minimal standard of what constitutes an adequate, “clear and meaningful” (AI Act, Article 86) explanation. From a technical point of view, it is challenging to design explanations that are both accurate and understandable, while exposing too much information might be confusing to users. Similarly, while cognitive attachment to GPAI systems is often understood as detrimental to user agency and autonomy, the removal or limitation of anthropomorphic interfaces might also place particular users at a disadvantage, including users in cognitive decline, isolated or with disabilities. Determining what a well-engineered interpretability stack looks like remains an open question for organisations, especially given that GPAI systems, which may provide the most accurate performance, are not explainable at the level of a single decision or output.



OM 20.

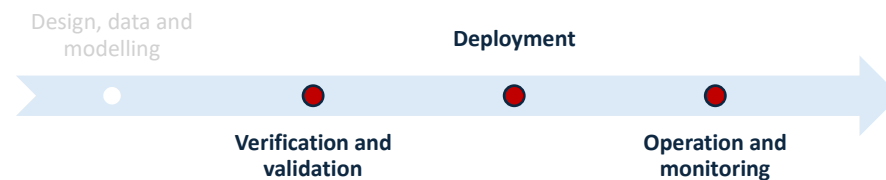
What: Document and enable traceability of changes to models, data or the GPAI system with sufficient detail to support internal reviews, external audits, or regulatory scrutiny.

Why: In the context of GPAI, small technical changes across models, training data and the system can have important impacts on deployment contexts and heighten safety risks. By documenting system changes, organisations are able to comprehensively cover the full range of modifications made to the GPAI system across its lifecycle, enabling **traceability** and **technical robustness and safety**. This is a key step for adequate scrutiny and auditing of the system, prompting an efficient identification of the origin of safety risks or operational malfunctioning by tracing them back to system change documentation. Documentation retention for a period of ten years is a key requirement in the AI Act for providers of high-risk AI systems (Article 18) and of GPAI models (Code of Practice for GPAI models, Transparency Chapter, measure 1).

How: Organisations must establish traceability documentation procedures and mechanisms for technical GPAI system changes, including minor technical fixes and major model updates. Organisations shall allocate internal responsibility for documentation keeping across relevant

roles and expertise (e.g., owner, maintainer, reviewer), coupling documentation management with technical measures that enhance traceability (e.g., system logs). Documentation should be kept over a relevant period of time, in a manner that complies with legal requirements, enables periodic internal and external audits, and allows organisations to go back to earlier system versions and technical specifications in the case of incidents. Traceability practices and related mechanisms should be assessed periodically to detect gaps and ensure legal and standards compliance. Documentation shall remain accessible to relevant parties, including auditors.

Implementation challenges: While documentation keeping is key for ensuring traceability of model, data and system changes, the implementation of this measure is challenging due to the frequency of changes and adjustments made in the context of GPAI development and deployment. The substantively different nature of GPAI – continuously updated post-deployment – when compared to traditional software updated based on longer release cycles (e.g., monthly), complicates documentation retention and can negatively impact the pace of model iteration. The major concern lies in determining the right threshold for documentation-keeping, which currently ranges from minor technical fixes to major updates, lacking adequate guidance. The implementation of this measure can also be compromised by time pressure and lack of discipline in logging / documenting changes, with privacy concerns arising if finely detailed information is logged. Small organisations, in particular, can struggle to manage vast amounts of logging data.



PERFORMANCE EVALUATION

This section outlines measures to support the evaluation of the AI system’s performance, namely through system monitoring and analysis, as well as internal auditing.

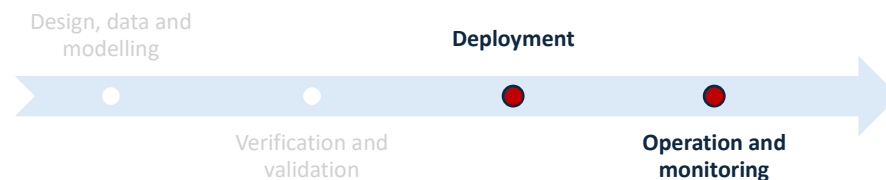
OM 25. **What:** Periodically reassess reliance patterns as systems evolve or scale.

Why: The risk of de-skilling, or the gradual erosion of expertise, is most salient in the long term, and can have a structural impact on organisational practices, including the ability to perform adequate GPAI **oversight**, leading to emerging risks, lack of legal and standards compliance, and reputational costs. To manage the integration and use of GPAI systems, organisations should assess use patterns overtime and as AI systems evolve, to flag and assess **over-reliance** and the underlying causes, address the deterioration of practices and skills, and assess the need for acquiring new skills in light of GPAI capabilities and new forms of human-AI interaction arising from GPAI deployment.

How: To address this measure, organisations shall monitor how operators interpret, accept, or override AI recommendations, identifying both over-reliance and under-use. Responsibility for monitoring use patterns shall be allocated among teams (e.g., team leaders, line-managers). Findings shall be communicated to relevant departments, including HR, operations and

management, as well as translated into continuous professional development trainings and, where necessary, lead to both reskilling based on evolving modes of human-AI interaction and changes in GPAI design. While reliance patterns shall be periodically assessed through monitoring and learning from human-AI interaction, significant changes to GPAI design or related organisational practices require immediate reassessment.

Implementation challenges: The implementation of this measure can be compromised by resource constraints and the pace of organisational workflows, as team leaders and line-managers might not have the availability to adequately assess overreliance. Moreover, overreliance patterns might emerge slowly over-time, being unrecognisable as such and become part of GPAI use practices. Organisations face the challenge of ascertaining what constitutes overreliance versus simply evolution of GPAI use scale, being required to balance conflicting objectives, such as the acceleration of workflows and productivity against the time-consuming process of preserving skills and reskilling.



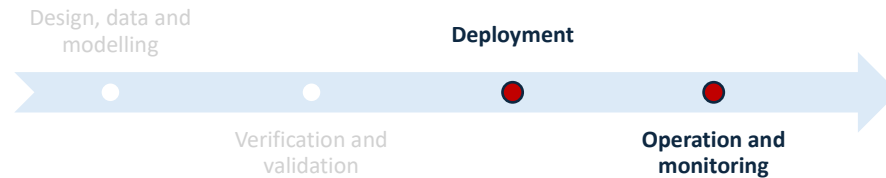
OM 26. **What:** Monitor situations when the purpose or actual use of the AI system drift or diverge from the intended ones and inform the user when relevant.

Why: GPAI system's functionality can be impacted by human-AI interaction, leading to changes in its operation and use that might evade the boundaries of previously tested environments and performance. Mechanisms that enable the **oversight** and monitoring of drift in GPAI use and intended human-AI interaction practices must be put in place to ensure that GPAI systems continue operating within tested grounds, preserving **human safety**. Drift and divergence in GPAI system's operation must be communicated to users to enable the calibration of judgement and behaviours when interacting with the system. Informing the user is a **transparency** requirement that prevents **over-reliance** and enables users to critically engage with the GPAI system's outputs.

How: To monitor for model drift and divergence, organisations must first determine what constitutes the intended model purpose and use (see OM01). Following this, clear GPAI use guidelines and responsible behaviour shall be elaborated to calibrate and render uniform organisational use practices. The latter, shall be designed in conjunction with technical measures for monitoring and flagging GPAI performance and drift. GPAI use guidelines shall be disseminated internally and responsibility for monitoring GPAI use shall be allocated among teams (e.g., team leaders, line-managers), in a manner that enables realistic oversight of AI use practices. Upon identifying inadequate GPAI use, organisations shall assess the underlying reasons (e.g., GPAI design, productivity KPIs, etc.) and take adequate measures to enable adequate GPAI use, such as training, task distribution or technical changes to GPAI design. If necessary, organisations shall suspend GPAI deployment.

Implementation challenges: Organisations face the challenge of defining clear technical and measurable thresholds for what constitutes actual drift versus natural variance in human-AI

interaction. This is further complicated by the administrative burden placed on team leaders and line managers, who must balance daily productivity against the need for rigorous oversight. Translating flagged divergence into actionable interventions, such as deciding when to retrain staff, redesign GPAI interface, or suspend system deployment, requires a complex balancing act between organisational efficiency and safety boundaries.



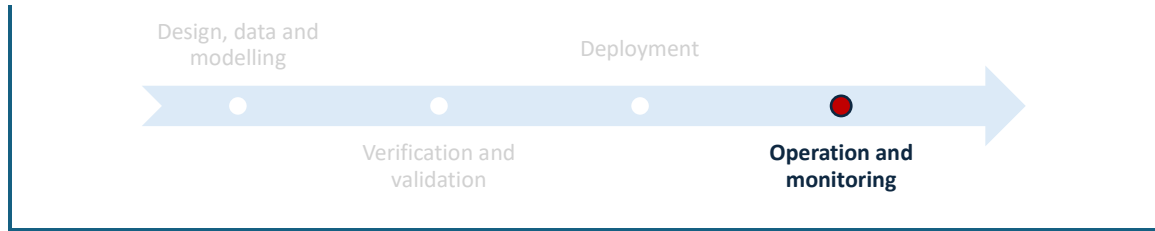
OM 30.

What: Review GPAI outputs for unintended disparity impact.

Why: The generality, malleability and recursive learning of GPAI systems over time demands ongoing review and assessment of the system's functioning to identify and mitigate unintended impacts on **societal well-being**, including disparate effects in specific populations, user groups and operational contexts. The aim of this measure is to identify emerging risks and mitigate potential blind spots, **safety** hazards / harm, and biases resulting from GPAI outputs and related organisational practices, enabling organisations to uphold the principle of **diversity, non-discrimination and fairness**.

How: This measure can be achieved through conducting formal review workshops or other structured consultation mechanisms that allow for discussion and deliberation on potential trade-offs among an interdisciplinary group of stakeholders, such as civil society, advocacy groups, linguists, and impacted users. The selection of stakeholders must be based on a comprehensive assessment of the GPAI system, including the underlying training data, intended use and deployment context. For example, a language-based GPAI system used in the detection of harmful online content, requires the involvement of linguists and representatives from minority communities and civil society in the review process. Document representation and keep minutes of reviews as evidence of stakeholder engagement, including concerns raised, disagreements and decisions on design adjustments based on feedback. Outcomes of the review process must lead to concrete action points to be implemented by the technical team and lead to adaptation and changes to the GPAI system and organisational procedures, to enhance operational conformity. Document and cross-check implementation of action points in technical and organisational design.

Implementation challenges: Internal teams often lack the specific expertise required to identify subtle disparate-impact risks, which may result in treating review workshops as a passive compliance formality rather than a meaningful analysis. Whereas some organisations might be able to recruit external experts, smaller organisations, with more limited financial and logistic resources, may struggle to identify, recruit and convene the adequate level of expertise. Pressure to speed up GPAI deployment might lead to neglecting a deeper consideration and review of disparate impacts, which are less obvious to detect, or result in the limited uptake of expert recommendations.



IMPROVEMENT

This section outlines measures for organisations to continuously improve the operation of the AI system, including determining when performance is not in conformity with expected requirements and intended use, and adopting corrective actions.

OM 31.

What: Translate audit findings into corrective design action.

Why: The use of GPAI systems – which present significant generality and malleability – requires regularly auditing if the system is performing as intended and if human-AI interaction standards are being adequately upheld and performed. This enables assessing aspects such as the system’s **robustness**, **non-discriminatory** functioning and **safety** in human use. Crucially, audit findings must not simply be reported, but result in actual corrective actions, both at the system level and at the level of human-AI interaction, through embedding ongoing feedback loops that improve organisational safety culture.

How: Audits should be periodically conducted to systematically analyse both GPAI functioning and human-AI interaction in GPAI use processes. Findings should be translated into actual processes and governance improvement, through documentation, revision of standard operating procedures, and trainings, among others. The organisation may also create knowledge repositories with case studies highlighting both successful and problematic human-AI collaboration, providing insights on best practices and approaches to avoid. These insights can be instilled into the organisation through learning sessions and internal workshops.

Implementation challenges: Fragmented recording of audit findings, the lack of elaboration of a clear pathway for translating audit findings into corrective design actions, and the fact that audit action points might require the involvement of different departments and teams can lead to challenges in the implementation of this measure. Organisations might struggle to build reliable feedback loops for bridging the gap between identifying safety or performance risks and executing the technical and behavioural revisions required to fix them. Teams required to act upon audit findings might not feel ownership or agree with audit findings can lead to these being inadequately addressed and translated into standard operating procedures and trainings.



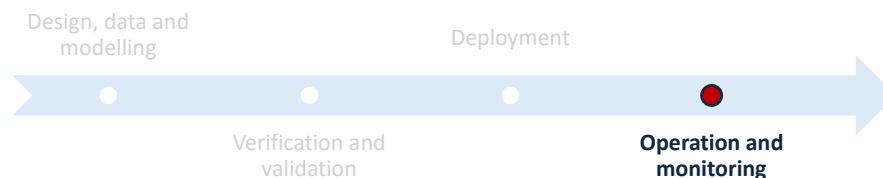
OM 32.

What: Identify and analyse unforeseen effects of the AI system on individual and societal well-being.

Why: With GPAI systems being a fairly recent technological development in the AI landscape, there remain significant gaps in the understanding of the risks and potentials harms emerging from GPAI deployment both at the level of individual and **societal well-being**. Most notably, technical guardrails with regard to GPAI system design remain an outstanding challenge, and research on individual and societal-level impacts remain limited. Therefore, providers and deployers of GPAI systems have a responsibility to actively contribute to this evidence base through conducting research on the impacts of their systems, complementary to, and not substituted by, capabilities-focused development.

How: Organisations should set up dedicated R&D teams or provide access to external researchers tasked with the responsibility to conduct research on user impacts (e.g., user testing, longitudinal monitoring, interaction log analysis). Organisations should be available to participate in independent academic studies and broader consortia, including by contributing with internally-generated data and allowing access to internal usage of the GPAI system to academics and the press. Research findings must be continuously integrated into ongoing GPAI design, risk mitigation and use practices.

Implementation challenges: Smaller organisations, with less financial capacity, might not be able to set up a team dedicated to R&D activities. This obstacle can be mitigated by providing access to external researchers and journalists. While providing access to external researchers and enabling meaningful scrutiny of the GPAI system can enable research on individual and societal-level impacts, organisations might be reluctant to expose their systems and practices, due to privacy, competitiveness and liability concerns. Conducting research on and assessment of impacts on societal well-being is particularly challenging to be conducted at the level of a single organisation.



2.1.4. Ethics tensions

In line with D3.1 and as demonstrated in the aforementioned sections, the process of operationalising AI ethics principles conducted in the context of AIOLIA GPAI UCs has demonstrated that ethical principles are not self-contained, as they constantly interact and intersect with one another. In fact, in putting principles into practice and identifying organizational measures for their operationalisation, numerous tensions between principles – that is, the operationalisation of one principle compromising another – have emerged and required balancing trade-offs. This is why, neatly listing organisational measures under corresponding ethical principles might be misleading, overlooking principle-interaction and giving the false impression that the implementation of measures guarantees a principle has been upheld.

As the GPAI measures outlined above demonstrate, measures can target multiple principles at the same time, which reflects the fluidity and interconnection between the principles themselves. The presence of negative relationships makes adherence to principles less straightforward and may demand trade-offs; it may also be that there is tension within a single principle. This section points to several examples of competing ethical principles and suggested measures, as identified by AIOLIA's industry partners and emerging from the analysis conducted in the context of this deliverable.

Human Oversight VS Well-being

Human oversight is fundamental for ensuring GPAI systems are performing adequately and within expected boundaries, being a mandatory requirement for deployers of GPAI systems operating in high-risk contexts (AI Act, Article 14). It not only enables monitoring GPAI system behaviour, but allows for flagging and addressing critical safety concerns. Nowadays, many of the risks emerging from GPAI systems concentrate in private settings, as users engage in extensive interactions with personal companions and other conversational agents (UC5). Herein, oversight becomes key for maintaining user safety, however, it also poses questions with regards to exposure to toxic and disturbing content by human moderators. In these contexts, individual user safety can be at odds with moderators' well-being. To mitigate the risk to human well-being, organisations should opt for implementing significantly autonomous content moderation systems, shielding human moderators and only requiring their intervention in borderline and ambiguous cases.

From the perspective of professional behaviour, it has been acknowledged that there are benefits to using automated moderation systems. These not only allow for swifter and real-time moderation to occur – crucial in the context of conversational AI – but they tend to be more accurate than humans at flagging potentially unsafe and prohibited usage. Since human moderators are only required to act upon issues when these are flagged, having an automated moderation system shields human moderators from exposure to distressing content, contributing to better labour conditions and workers' wellbeing. At the same time, having an automated monitoring system helps preserve user privacy, since issues are only escalated to the human-in-the-loop if flagged by the system.

Human Oversight VS Deskilling

Human oversight is a legal requirement for high-risk AI systems under the EU AI Act. It is an active process of review and responsibility implying the ability to oversee and/or intervene in the processes and outputs of AI systems. Whereas the ability for effectively overseeing AI systems relies on human skills, the integration of GPAI systems into workflows can place the ability for oversight at risk, due to the more passive role of practitioners and the risk of deskilling in the longer-term. Beyond trainings, one measure pointed by UC2 to mitigate this risk is to ensure that oversight is conducted by senior staff, giving less experienced professionals the possibility to develop their expertise and skills by engaging directly with the technicalities of the tasks.

Narrative Highlight #1 – Human-in-the-loop and professional deskilling:

Before the AI system was introduced, Johana, a car safety engineer, spent her days coding manually to ensure a braking system wouldn't fail. She never used AI believing it would miss the nuance of physical friction on a wet road. However, AI decision-support has recently been made mandatory for all engineers, with the company transitioning from a culture of safety engineering to one of safety verification.

Now she oversees the system's outputs and intervenes whenever the system signals a fault or the assessment is ambiguous. She has come to appreciate how easy and fast her work has become. Once, when she realised she had to override a decision made by AI, she struggled to remember a specific line of code and ended up having to go back to safety manuals, taking a full afternoon to fix the issue. She worried she might be losing her own expertise as a safety engineer.

Recently, Johana welcomed a new trainee, Marco, who graduated in mechanical engineering. One day she will retire and Marco will replace her as the company's principal safety engineer. Will he then have enough knowledge and expertise?

Safety VS Privacy

A key tension that arises in the context of UC5 is the need to balance privacy and safety requirements amid both private and professional behaviour. With regard to private behaviour, the use of a content moderation system to flag inappropriate and dangerous behaviour is at odds with basic privacy requirements under the GDPR. Most notably, training these systems to excel at monitoring and flagging unsafe, malicious or prohibited content requires significant training data extracted from user interactions with the conversational agent. Adding to this, completely removing all identifying elements is in conflict with the need for a monitoring system, which should be able to identify abusive use and, if necessary, take targeted action against users. This balancing act is delicate as it requires weighing the risk of identity exposure against the risk of allowing toxic or dangerous content slip through unfiltered.

Safety VS Competitiveness

One major tension that emerged in the context of UC5 was between safety and market imperatives. Besides drawing on internal guidelines and existing regulatory requirements, developers look at major market actors to benchmark their system's boundaries and use policies. Competitors become a reference point for assessing compliance practices, based on the principle: 'If this is permissible for a leading competitor, we can also justify it.' This competition-based orientation thus functions as an informal but effective mechanism for self-assessment. It illustrates that the currently unclear and incomplete regulatory situation means that international industry standards are effectively used as *de facto* guide for compliance-related decisions. The balance between complying with safety measures and maintaining a system that remains attractive from a market perspective can be a fragile balancing act, since aligning the system too much risks making it 'boring' and might result in the loss of consumers. From the user's perspective, the appeal of conversational AI lies in fluid, natural interaction and a sense of agency or spontaneity. Many users might intentionally configure the AI to play a particular role (e.g., dominant

authority) or develop a more manipulative character, encourage unpredictability or chaos (e.g., by escalating emotional or narrative stakes), deliberately testing boundaries or attempting to trick the system.

Narrative Highlight #2 – Benchmarking safety

When we at PLAY launched our AI companion, we were immediately challenged by users who didn't want a helpful assistant but a virtual character they could play with. Some spent hours deliberately testing the limits of the AI system. They were smart to avoid obvious red tape—no explicit violence, blood, or forbidden keywords. They played in the grey area of AI alignment.

We faced a choice: tighten the filters and risk a mass exodus of users to wilder competitors or maintain a fluid and attractive interaction at all cost. To decide, we practiced the same attacks on competitor systems, simply because if being the safest AI meant being the one with no users, then we did not want to be the ones.

In the end, we opted for narrowing the system's guardrails. This has worked: engagement metrics spiked and user retention improved. However, the logs revealed something unexpected: that the users weren't necessarily looking for grand transgressions. Instead, they were pushing the AI into acting like a validation machine that never contradicted their input, even if it was factually or morally wrong. It broke none of our established safety protocols, yet we felt that this was professional negligence because there was simply no friction. The system's alignment was now shaped by user indulgence rather than the principles we stood for.

Transparency VS Technical opacity

In the context of GPAI, transparency increasingly comes into tension with the growing complexity of AI systems, with regulatory frameworks being unclear on the extent and granularity requirements of transparency – an issue highlighted by UC5. As models become more sophisticated, their internal operations, adaptive behaviours and context dependent outputs are difficult to fully explain in a way that is both accurate and comprehensible. Given the malleability of GPAI systems, which adapt and are customised based on human-AI interactions, it becomes challenging to accurately foresee emergent risks, a challenged that is heightened by technical opacity. Attempts to mitigate problematic behaviours may have unpredictable consequences for the model's behaviour in other contexts, which are just as hard to predict. There is neither an established methodology to estimate these trade-offs, nor can it be ensured that the test-set on which statistical tests are run is complete.

2.2. EMOTIONAL AI

Emotional AI is a category of AI models and systems that are capable of emotion recognition, emotion emulation and/or emotion elicitation. A common term in the literature that is synonymous to Emotional AI is affective computing. Examples of techniques used for emotion recognition, or the process of inferring human emotion, are sentiment analysis of online language, facial coding of expressions and eye-tracking via image recognition algorithms, and analysis of biometric data such as heart activity via wearables.

Emotion emulation refers to the ability to imitate emotions such that their functions are closely replicated. For example, a chatbot may be able to console users by showing care, without itself experiencing any emotive or inner states; this capability is also referred to as emulated empathy. Lastly, emotion elicitation is the process of inducing emotional states in individuals, such as in human-computer interaction, including with the use of deepfakes.

As a research area, Emotional AI presents unique ethical challenges that current governance frameworks are actively grappling with. This section starts by presenting key open issues in governing Emotional AI systems and how AIOLIA measures can address them. We then elaborate on the underlying concerns by looking at salient Emotional AI ethics principles. What follows is the main finding – the list of AIOLIA organisational measures for Emotional AI. This section ends with an overview of some important tensions between principles that organisations and policy-makers must be aware of when implementing measures.

2.2.1 Open issues in the governance of Emotional AI systems

Emotional AI has traditionally been associated with the inference of emotional states from biometric data, a practice that attracted significant civil society and academic criticism on grounds of technical unreliability and ethical risk, and was eventually restricted by the AI Act (Article 5(1)(f)). However, the emergence of large language models has significantly expanded the scope of what constitutes Emotional AI, as has the experimental use of deepfakes in therapy contexts. General-purpose systems can be used for emotionally intimate interactions even without being designed for this purpose, and purpose-built AI companions such as Replika and Character.ai are developed on top of the same GPAI models. In both cases, emotional states are inferred from language rather than biometric data, and may be directly disclosed or inferred from interaction patterns. Beyond inference, LLMs are also capable of emotion emulation, imitating emotional states in ways that introduce anthropomorphism and distinct risks to user autonomy and privacy. Neither language-based emotional inference nor emotion emulation falls within the AI Act's prohibited or high-risk categories, leaving a significant regulatory gap. Further to language-based AI systems used in the private sphere, therapy contexts represent a distinct governance challenge: these solutions are currently only used in experimental settings, with their effectiveness yet to be established.

Given these governance gaps, organisations can pro-actively take a precautionary stance, treating Emotional AI systems as high-risk unless demonstrated otherwise. For policy, the challenge is two-fold: regulating the use of AI for emotional ends, regardless of whether they are purpose-built applications or general-purpose; and ensuring regulation protects society from the risks, without blocking the potential for increased human well-being.

Open issue #1: How to distinguish between beneficial and harmful anthropomorphism and emotional engagement?

Concern: The same features that make Emotional AI effective, such as affective rapport and human-like interaction, also carry the risk of emotional dependency, manipulation, and erosion of autonomy. Since the evidence base is nascent, there is no settled threshold for when

How organisational measures help: Organisations can define content and behaviour restrictions and document the limits of legitimate influence, involve mental health professionals in design and monitoring, and invest in research on user impacts. Internal and independent ethics

<p>engagement becomes harmful; what's more, this threshold will necessarily be context-dependent and may be human-judgement-dependent.</p>	<p>reviews can provide additional scrutiny as systems evolve.</p>
<p>Open issue #2: How should subjective, dynamic, and emergent psychological factors, including vulnerability, be incorporated into risk analysis?</p>	
<p>Concern: Vulnerability in Emotional AI users is not always pre-existing or identifiable: it can emerge during interaction, shaped by the system itself. Standard risk frameworks are designed to assess discrete, observable, and largely technical failure modes, and are therefore poorly equipped to capture psychological or cognitive harms that are subjective, gradual, and may only become apparent over extended use.</p>	<p>How organisational measures help: Organisations can define vulnerability profiles and identification protocols, integrate mental health expertise into risk analysis teams, and collect user feedback on wellbeing and perceived coerciveness. Continuous monitoring and periodic review allow profiles to be updated as evidence grows.</p>
<p>Open issue #3: What constitutes harmful influence, and how can governance address manipulation that is gradual or unintentional?</p>	
<p>Concern: Manipulation by Emotional AI systems may emerge as an unintended product of training dynamics rather than deliberate design, and its effects may accumulate gradually through repeated interaction. These are not well captured by the AI Act's intent-based prohibition on manipulation.</p>	<p>How organisational measures help: Organisations can restrict and audit behavioural nudging, establish multi-level safety criteria that go beyond legal compliance, and monitor long-term interaction patterns for signs of harmful drift. User contestation mechanisms provide an additional check.</p>
<p>Open issue #4: How can governance prevent Emotional AI from eroding social norms and interpersonal skills at a societal level?</p>	
<p>Concern: Widespread use of AI companions may produce spillover effects on human relationships – reducing motivation to invest in them and gradually degrading social skills. Deepfake therapy, if normalised, may also contribute to the broader normalisation of deepfake technology in contexts where it causes harm.</p>	<p>How organisational measures help: Organisations can communicate risks to individuals, next of kin, and the public, and contribute research findings to broader academic and policy communities.</p>
<p>Open issue #5: Should AI companionship be treated as high-risk, and how should general-purpose AI used for emotional ends be regulated?</p>	
<p>Concern: Most Emotional AI applications, especially AI companions and general-purpose systems used for quasi-clinical emotional support, fall outside the AI Act's high-risk and prohibited categories, leaving a significant regulatory gap. Regulation cannot apply to purpose-built applications only, as GPAI models may be used for emotional ends, or fine-tuned to become Emotional AI systems.</p>	<p>How organisational measures help: Organisations can go beyond legal compliance and adopt a precautionary stance, i.e. treating Emotional AI systems as high-risk in the absence of a regulatory requirement to do so. Governance committees and independent ethics reviews can help establish voluntary standards that inform future regulation.</p>
<p>Open issue #6: Are Emotional AI systems effective in clinical therapy settings, and how should this be established?</p>	

<p>Concern: The clinical grey zone between wellness and medical products is expanding, yet evidence on the therapeutic effectiveness of Emotional AI remains limited. Without robust evidence, neither benefit nor harm can be reliably established, undermining the proportionality requirement at the heart of good clinical care.</p>	<p>How organisational measures help: Organisations can establish dedicated R&D functions or enable access for external researchers, involve mental health professionals in evaluation, and collect and share user feedback on wellbeing outcomes. Participation in academic studies and cross-sector consortia will be key in accelerating the evidence base.</p>
<p>Open issue #7: Are current privacy and data protection frameworks fit for Emotional AI, particularly regarding inferred data, consent, third-party data, and "mind data"?</p>	
<p>Concern: Emotional AI systems continuously infer sensitive mental states from interaction data, yet both inferred data and mind data as a whole sit outside GDPR's sensitive data categories. Consent mechanisms designed for one-off disclosure are ill-suited to long-term, affective interactions; moreover, users could be well aware of the privacy risks, yet have no choice but to give their consent in order to use the platform. Third-party data used in deepfake therapy remains legally unresolved in some cases, such as in grief therapy.</p>	<p>How organisational measures help: Organisations can conduct privacy impact assessments that go beyond GDPR compliance to map the full range of data risks, implement enhanced technical and procedural safeguards for certain categories of data, and document how consent mechanisms fall short of equipping users with meaningful agency.</p>

Across these issues, organisational measures face important limits that policy can address. Firstly, since scientific consensus on harm is emerging, measures necessarily rely on normative judgement rather than evidence. When it comes to societal effects, especially over time, individual organisations lack the reach to detect or address them. Finally, while precautionary measures can proactively enhance protection where it is not required by law, they remain voluntary and their implementation may go against commercial incentives, depending on organisational willingness and, crucially, resources. We present policy recommendations that address these limits in Section 3.

2.2.2 Salient ethical concerns

While AI ethics principles, such as those underlying the seven requirements of the Assessment List for Trustworthy AI (ALTAI)⁷, remain relevant, some become more nuanced and salient when applied in the context of Emotional AI. In fact, AIOLIA use cases have identified as many as twelve relevant ethics principles in the different contexts of deployment of Emotional AI, that form the basis of the organisational measures advanced herein. However, three ethics principles revealed particularly salient in the context of Emotional AI, namely: 1) human agency and autonomy, 2) safety and non-maleficence, and 3) privacy.

⁷ <https://digital-strategy.ec.europa.eu/en/library/assessment-list-trustworthy-artificial-intelligence-altai-self-assessment>

Human Agency and Autonomy

In the context of Emotional AI as it relates to the user or patient, human agency and autonomy extends beyond informational control over personal data and control over the AI system, to address control over AI interactions and mitigate the risks of over-attachment and manipulation. A central challenge is to navigate between non-manipulation and the mitigation of harm alongside the incorporation of anthropomorphic features, which are fundamental to the utility of Emotional AI systems. With careful design, such features can mitigate loneliness, provide emotional support, lead to positive clinical outcomes or even support the exercise of human autonomy by facilitating the pursuit of relationality, identity-building or empowering therapeutic engagement. At the same time, interactions can lead to undue deception, over-reliance and exploitation of vulnerabilities, expectations, and trust. Whereas the AI Act prohibits subliminal, manipulative and deceptive AI systems (Article 5), it treats these as observable, intent-driven acts, which may not address the subtle behavioural changes occurring as interactions accumulate. Therefore, the key ethical question is how to define and measure unacceptable influence and preserve the user's or patient's capacity for autonomous cognitive and emotional development in the long-term.

An important element of human autonomy in the professional context is the exercise of oversight. On the one hand, oversight is understood as a principle, requiring professional control across the stages of the AI lifecycle, including design, training and use, and therefore tightly linked to organisational accountability. On the other hand, oversight can be understood as a measure that cuts across multiple principles, with professional control not being a goal in itself, but a means to an end – ensuring the ethical use of the AI system.

Safety and Non-maleficence

Safety and non-maleficence share the requirement to protect against harm, including harmful advice, emotional dependence, emotional manipulation, erosion of autonomy, and social deskilling. Safety is understood as proactive avoidance of harm by identifying and preventing risks that may arise in human-AI interaction. Non-maleficence – or 'do no harm' – requires that the AI system's benefits outweigh the harms (proportionality) and that the system achieves its intended purpose with the least harm possible (effectiveness). These elements of safety clearly expand the requirements listed in ALTAI, with the bioethical principle of non-maleficence, commonly relevant in clinical scenarios, gaining traction and relevance in the broader contexts of deployment of Emotional AI systems. While well-being in ALTAI is framed at the societal level and as separate from safety, the clinical context brings individual benefit to the fore of ethical concerns, since effectiveness requires a prior assessment of the harms relative to the benefits that the system shall bring in an effective way.

Indeed, while the ethical duty to avoid harm is clear, the field is yet to produce a settled account of what constitutes harm in the context of Emotional AI. This is especially the case in relation to benefits, due to the entanglement of the two, the multitude of deployment contexts, the varying degree of professional involvement and judgement required, and varying individual vulnerabilities – all of which challenge a clear-cut, consistent distinction between positive and negative impacts. Most notably, the interactive and autonomous nature of some Emotional AI systems means that harm might emerge in the course of

human-AI interaction, which may not be identified during system design, which calls for ongoing monitoring and oversight.

Privacy

Emotional AI systems pose distinctive challenges to privacy. They necessarily rely on the processing of sensitive, intimate and context-rich personal data, accumulated through interactions that are based on voluntary engagement, which makes ensuring the functionality of the system while preserving privacy more challenging. Further, voluntary interaction does not equate to voluntary disclosure, as the AI systems must engage in continuous inference of the user's mental states to respond effectively (in the context of language-based systems). The resulting inferred data presents a particular challenge, as it sits outside conventional data protection frameworks centred on data disclosure: inferences may shape system behaviour without the human participant being aware that they are being made or stored and, whether accurate or spurious, they constitute a distinct category of privacy risk. In therapy contexts, training or inferred data may additionally involve third parties who equally fall outside existing frameworks.

In the context of conversational Emotional AI systems, other important elements are safety measures and informed consent mechanisms. Safety measures, such as automated content moderation, require significant data collection and processing extracted from user interactions. At the same time, taking targeted actions against users who might breach AI use policies is often at odds with data anonymisation and the removal of all identification elements, and requires weighting safety risks against identify exposure. From an organisational perspective, understanding when such disclosures are appropriate requires sensitivity to context. A further complexity is ensuring that informed consent mechanisms provide meaningful choice and facilitate expectation management, when in practice, they are rather punitive toward users, who have no choice but to accept the conditions in which their data is processed in order to have access to the system they might have already developed an emotional bond with.

2.2.3 Organisational measures

Against the backdrop of the open issues in the ethical governance of Emotional AI, we present a list of organisational measures that address them, specifically targeting regulatory challenges and gaps to adopt a precautionary stance. They are organised according to the categories in the international standard for AI quality management systems ISO/IEC42001⁸. Each measure is presented together with its relevance for the research area (under "why"), explaining how the measure addresses pressing concerns and regulatory gaps, as identified in the literature and the use-cases. Each entry also includes a brief explanation on how the measure can be achieved (under "how"), which draws on the UC data and validation workshops, and possible obstacles to their implementation.

The list of measures is non-comprehensive but foundational, and can serve as a starting point of reference for organisations aiming to embed ethics into the Emotional AI technology they design, develop and/or

⁸ <https://www.iso.org/standard/42001>

deploy, not only to comply with relevant regulation, but to pro-actively facilitate positive outcomes, while mitigating negative outcomes.

OVERVIEW OF ORGANISATIONAL MEASURES FOR EMOTIONAL AI	
Context	
OM02	Define clear rules on the type of content, behaviour, or expression that are restricted for AI use and why, including distinguishing legitimate behavioural influence from manipulative practices.
OM03	Define what constitutes vulnerability in users or patients, including static and emergent forms, set out identification protocols, and establish corresponding enhanced safety measures, interaction protocols and continuous monitoring of effects for identified vulnerability profiles.
OM04	Define multi-level safety criteria and categories of risks and harms.
Leadership	
OM07	Establish a representative AI governance committee as a formal oversight structure within the organisation.
Planning	
OM12	Conduct a privacy impact assessment.
OM13	Involve mental health professionals in the design, monitoring and evaluation of the AI system.
OM14	Determine the permitted degree of autonomy of the AI system (if any) and accordingly, specify the conditions in which a professional in the loop is required, e.g. via human moderation.
OM15	Consider psychological, emotional, and behavioural influences in AI risk analysis, not only user or patient information or technical AI performance.
Support	
OM17	Clearly communicate how the AI system may affect human individuals, next of kin and society at large.
Operation	
OM18	Put in place a mechanism for users or patients to report, question, or contest AI behaviour, decisions and/or restrictions.
OM21	Implement additional safeguards for sensitive data.
OM22	Ensure the AI system, if autonomous, applies behavioural nudging only in documented and auditable contexts with specified trigger thresholds and objectives.
OM23	Conduct internal ethics reviews for new features or changes of the AI system.
Evaluation	
OM28	Conduct periodic independent ethics reviews of the AI system's impact on wellbeing and autonomy addressing all stages of the AI lifecycle.
OM29	Conduct audits of informed consent mechanisms and document limitations, especially with regard to the adaptivity to context.
Improvement	
OM32	Identify and analyse unforeseen effects of the AI system on individual and societal wellbeing.
OM33	Assess and collect feedback from users of conversational systems regarding perceived honesty, non-coerciveness of AI interactions and impact on wellbeing.

Table 5 - Overview of organisational measures for Emotional AI.

Context

This section outlines measures to address the specific context of the organisation, including internal and external contexts, ranging from legal requirements to competitive landscape, as well as cultural and social values. Understanding the organisational context also requires accounting for all the parties impacted by the AI system, from direct users to wider communities.

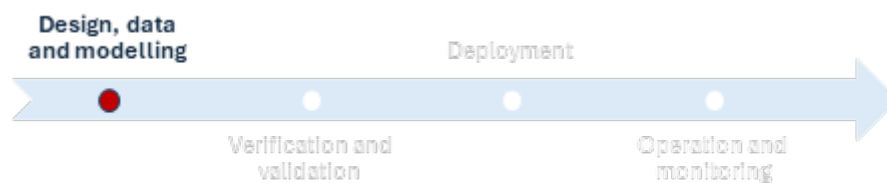
OM 02.

What: Define clear rules on the type of content, behaviour, or expression that are restricted for AI use and why, including distinguishing legitimate behavioural influence from manipulative practices.

Why: Defining boundaries against harmful, violent, or potentially traumatising effects and keeping the AI system within the intended context is paramount for **human safety**. This includes defining the limits of acceptable influence, given that the boundary between legitimate affective engagement and manipulative practice is particularly difficult to draw in systems designed to elicit emotional response - risking to erode **human autonomy**. The rules aim to ensure effective **oversight** by grounding the adoption of technical measures for model behaviour alignment, human moderation rules, or to guide professional use of non-autonomous architectures.

How: Rules can be established through fostering an internal ethical discussion, leading to the draft and approval of a policy stating the AI systems' supported capabilities and the topics or behaviours that shall remain out of scope. The restrictions must be grounded on a legal, ethical, or safety basis. The policy should be communicated to all relevant roles across the organisation, and be reflected in technical documentation and measures, as well as on human actions and behaviour, across design (e.g., defining scope of AI system domain knowledge), training (e.g., limiting certain capabilities) and operation (e.g., content moderation or guidance for professional use). In the case of professional use, the rules must form part of staff training. Regularly review and update the policy considering regulatory developments, internal audits, user or patient feedback, technical developments, and any other relevant evidence.

Implementation challenges: There is limited scientific evidence on the impacts of Emotional AI technology, which means the implementation of this measure will rely substantially on normative judgement rather than empirical consensus. Further, restricted categories are difficult to define precisely, for example, the definition of manipulation introduces ambiguities. Relatedly, determining appropriate restrictions may depend on individual users or specific interactions rather than applying uniformly, which may make it practically impossible to anticipate all cases in a single policy. This requires rules that are flexible enough to accommodate unforeseen situations, yet precise enough to provide meaningful protection. Finally, there is an inherent tension between ensuring user or patient safety and preserving user or patient autonomy, as overly restrictive rules can go against legitimate agency, while overly permissive ones may fail to provide adequate protection from harm.



OM 03.

What: Define what constitutes vulnerability in users or patients, including static and emergent forms, set out identification protocols, and establish corresponding enhanced safety measures, interaction protocols and continuous monitoring of effects for identified vulnerability profiles.

Why: Individual vulnerabilities deserve attention in the context of Emotional AI because of the one-to-one, intimate nature of interactions, and they can be emergent in the sense that they can arise during the interaction – in turn shaped by business models or clinical practices that are novel and unstable. Vulnerability shapes the degree to which a person may be harmed or benefitted, therefore it must be accounted for to ensure non-maleficence, including the effectiveness of the system.

How: Use scientific and medical evidence to map vulnerability profiles, accounting for a wide range of factors, including age, cognitive impairment and diagnosed mental health conditions. Differentiate between forms of vulnerability that can be identified in advance and that can emerge during use and establish corresponding privacy-preserving protocols for data collection or use monitoring, to identify vulnerability. Design protocols, customisation options, and safety guardrails to accommodate a diverse set of characteristics and vulnerabilities. This may involve adapting interactions in response to individual needs (e.g., adjusting the course of a deepfake-based therapy session or limiting certain types of engagement where there is a risk of over-attachment), and ensuring that policies governing permitted customisation and system behaviour explicitly account for different vulnerability profiles. Review the profiles and protocols periodically, whenever changes are made to the system or in light of new scientific/medical evidence or regulatory requirements. This measure must be paired with adequate technical measures.

Implementation challenges: There is limited scientific evidence on vulnerability in the context of Emotional AI technology, meaning the implementation of this measure will rely substantially on normative judgement rather than empirical consensus. Identifying vulnerabilities requires collecting user information, which means this measure may be constrained by privacy obligations.


OM 04.

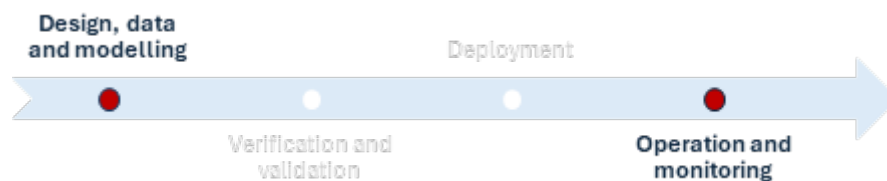
What: Define multi-level safety criteria and categories of risks and harms.

Why: Establishing comprehensible, tiered and differentiated safety criteria and risk / harm categories enhances **safety** and enables proportionate responses calibrated to severity; where relevant, these support consistent decision-making in **oversight** and overall user experience. Clear criteria and protocols equally serve as a foundation for the exercise of individual professional judgment in contexts where this is essential, like in moderation or in clinical practice. The lack of a settled account of what constitutes harm, especially given the wide range of Emotional AI contexts, is not a reason to forgo this measure, but on the contrary, it both highlights the need for professional guidance and means that the criteria provide the

basis for the emergence and ongoing refinement of context-specific best practices, beyond the policies of individual organisations.

How: The development of safety categories should be done collaboratively by an interdisciplinary team, including legal, data and model engineering teams, human moderators, and, where relevant, third parties (e.g., commercial partners, representative stakeholders). They should incorporate user or patient vulnerabilities (OM03) and there should be an ongoing assessment of tier definitions against emerging patterns. Defined safety categories should translate into organisational practices, including content moderation practices, rules for violation classification or individual professional judgement. The definitions should be kept up to date with the latest legal requirements and ethical and clinical research on negative effects. Whenever possible, organisations are encouraged to their safety categories and the evidence base underlying them within their professional and research communities, to contribute to the emergence of sector-wide best practices or standards.

Implementation challenges: There is limited scientific evidence on impacts from Emotional AI technology, meaning the implementation of this measure will rely substantially on normative judgement rather than empirical consensus. Further, going beyond pure compliance to develop precautionary criteria requires sustained investment of time and resources that may be difficult to justify internally. Similarly, involving a multi-disciplinary team, likely including external experts, may be costly for a small organisation. The collaborative nature of this process also means that definitions may be subject to significant internal disagreement, making the quality of outcomes heavily dependent on a robust governance of the decision-making process. Finally, sharing safety categories and best practices externally may raise concerns around reputational risk or competitive sensitivity, potentially disincentivising transparency.



Leadership

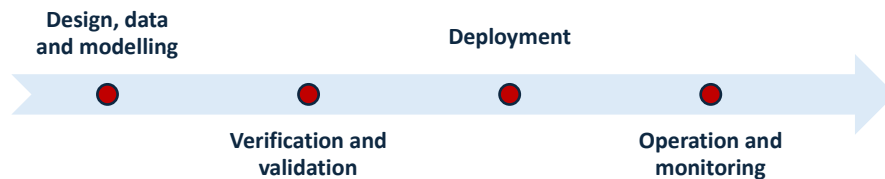
This section outlines measures to address the leadership commitment of the organisation towards the AI management system, including through defining a policy for AI systems deployed within the organisation, but also assigning clear roles and responsibility for managing these.

OM 07. **What:** Establish a representative AI governance committee as a formal oversight structure within the organisation.

Why: A dedicated governance committee ensures **human oversight** is embedded in organisational decision-making. This is particularly important for **accountability** in the context of Emotional AI, given the absence of specific regulatory guidance. The committee aims to safeguard against conflicting pressures and objectives, such as speed and unrealistic KPIs that might be detrimental to the **safety** of the AI system, and it can contribute to the creation of research-area-specific standards.

How: Define the AI governance committee’s mandate and appoint its members ensuring adequate representation across organisational roles, seniority, and areas expertise. Elect the committee’s representatives, including an AI Ethics Lead. Establish periodic meetings of the AI governance committee and keep minutes of discussions. Facilitate ongoing avenues for communication between the AI governance committee and organisational leadership.

Implementation challenges: Establishing a formal governance committee requires sustained organisational commitment and resources, which may be difficult to secure, particularly in smaller organisations.



Planning

This section outlines measures to address the planning phase of developing and / or deploying an AI system within the organisation. This involves identifying and assessing risks, potential impacts and opportunities arising from the AI system and setting clear objectives for its deployment.

OM 12.

What: Conduct a privacy impact assessment.

Why: As Emotional AI systems grow more personalised, the accumulation of sensitive interaction data continuously expands the **privacy** attack surface in ways that demand proactive, context-sensitive assessment beyond standard compliance checks. The privacy assessment shall feed into enhanced protective measures for identified high-sensitive data (OM 10), ensuring adequate data protection in a regulatory context that leaves some kinds of data (e.g. third-party data to produce deepfakes) largely out of scope.

How: The assessment should map all data flows across the AI system lifecycle, identifying categories of sensitive data collected, inferred, or stored. It should evaluate the contextual appropriateness of data use, assess risks of repurposing for model training or third-party sharing, and address the specific privacy implications of personalisation features and stored memories. In clinical contexts, it should additionally consider the privacy rights of third parties depicted in therapeutic applications. Findings should be integrated into system design and reviewed periodically or when significant changes to the system occur.

Implementation challenges: A privacy impact assessment has a broader scope than a data protection impact assessment and it is not mandated by GDPR like the latter. Therefore, as in OM3, the costs of going beyond data protection compliance may be difficult to justify internally, especially for small organisations. Further, in the case of third-party data use such as post-mortem, national laws may differ.



OM 13.

What: Involve mental health professionals in the design, monitoring and evaluation of the AI system.

Why: The psychological and cognitive dimensions and potential impacts of Emotional AI systems require the involvement of specialised mental health professionals in all stages of the AI lifecycle, enhancing organisational **oversight**. In the context of technical development, such complementary expertise is needed to ensure the **safety** of the technology. Specifically, mental health professionals can point to not immediately perceivable or obvious risks, such as emotional dependency and long-term cognitive and behavioural impacts that only become apparent over extended use.

How: The involvement of mental health professionals can be operationalised through conducting interviews, testing and other forms of periodical consultations regarding the AI system design, acceptability, effectiveness and/or safety guardrails. Consultations can be initially based on evidence from controlled experimental settings and, subsequently, on monitored user or patient experience. This should be established on a regular basis.

Implementation challenges: As in OM3, involving external experts may be costly for a small organisation.


OM 14.

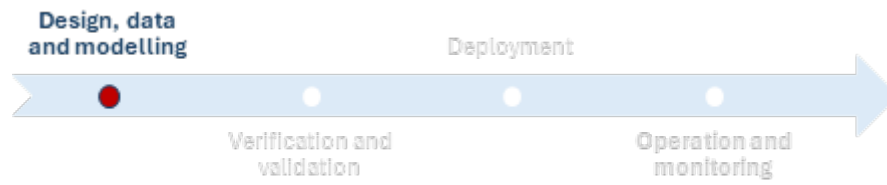
What: Determine the permitted degree of autonomy of the AI system (if any) and accordingly, specify the conditions in which a professional in the loop is required, e.g. via human moderation.

Why: Human control of AI-generated outputs may be a direct requirement where the user is a professional, or an inherent feature of the deployments in which the AI system is not autonomous – both conditions applying to current deepfake therapy applications. In the case of language-based Emotional AI applications, there is no practical possibility, nor necessity, for all generated output to be manually controlled. Automated moderation allows for real-time flagging of unsafe or prohibited practices but it cannot reliably handle borderline cases or jailbreaking attempts. Automated moderation allows for real-time flagging of unsafe or prohibited practices but it cannot reliably handle borderline cases or jailbreaking attempts. To ensure the **safety** of the user without violating their **autonomy** (i.e. right to control the interaction), human judgement can serve as a second layer of safeguard, managing complex or borderline cases to effectively discern user preferences from unsafe behaviour. The importance of human **oversight** as a best practice is highlighted by the absence of such regulatory requirement for applications that fall outside the scope of medical regulations or the high-risk category under the AI Act.

How: To implement this measure if professional control is not inherent in the deployment context, organisations should start by identifying what medical regulations apply, if any. In the case of MDR, for example, classification of software that performs medical functions without

being part of a physical device is challenging and organisations are encouraged to seek legal advice. To proactively establish the level of autonomy that is warranted according to the risk posed by harmful outputs, the deployer must engage in a risk analysis and establish systematic procedures for human review, approval, or override for high-risk cases, which must be communicated and implemented across relevant organisational roles.

Implementation challenges: Determining the conditions for human-in-the-loop in the case of automated, first-layer-of-defence guardrails requires balancing between flagging too many and overburdening humans, or flagging too few, resulting in user harm. Establishing conditions for human review may also be complicated by organisational reluctance to introduce friction into the user experience, particularly where awareness that conversations may be reviewed could generate dissatisfaction or deter engagement.



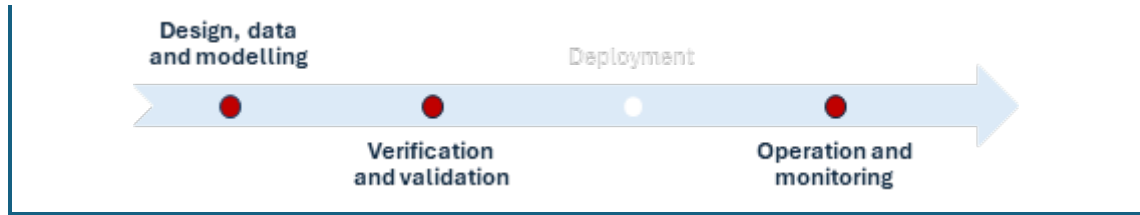
OM 15.

What: Consider psychological, emotional, and behavioural influences in AI risk analysis, not only user or patient information or technical AI performance.

Why: Standard risk analysis frameworks are not designed to capture psychological, emotional, and behavioural dimensions and impacts of AI systems. As a unique feature of Emotional AI systems, these elements must be incorporated in risk assessment procedures conducted within the organisation, determining the potential benefits and harms that Emotional AI can bring about within deployment contexts and therefore upholding the principle of **non-maleficence**.

How: Organisations can implement this measure by integrating mental health professionals within risks analysis teams and workflows. Specific indicators must be developed for capturing data regarding psychological, emotional and behavioural influences, attending to state-of-the-art research and practices. This kind of risk analysis must be conducted across several stages of the AI lifecycle, including AI design, testing – ideally within a controlled experimental setting involving humans – and post-deployment. Results from the risk analysis must be documented, disseminated and result in system improvements.

Implementation challenges: Small organisations may not have the internal expertise to conduct a robust psychological, emotional, and behavioural analysis, and may therefore involve external expertise, similarly to OM3 and OM6, which may be costly. What makes this process more challenging is that the incorporation of such influences in risk analysis requires turning subjective, qualitative indicators into quantitative evidence. This is an actively evolving research area and requires not only keeping up-to-date with the latest frameworks, but pushing beyond them.



Support

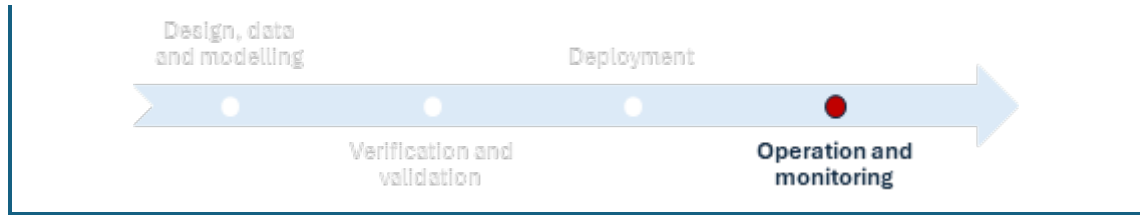
This section outlines measures to support the establishment, implementation and maintenance of the AI management system within the organisation. It involves determining and providing organisational resources, ranging from competencies, awareness and clear communication and documentation.

OM 17. **What:** Clearly communicate how the AI system may affect human individuals, next of kin and society at large.

Why: Transparency about intended use and effects is a requirement for high-risk AI systems (Article 13), but most Emotional AI systems fall outside of this category; moreover, the requirement targets deployers, not the end user. However, to ensure meaningful informed engagement and manage expectations, therefore upholding **human autonomy**, such disclosure is essential. The need for this measure is heightened by the novelty of the field; the inclusion of next of kin and society at large is warranted by the relational nature of these systems and its potential for gradually but firmly shifting societal and cultural norms.

How: To implement this measure, organisations must provide information about the AI system's intended use and potential impacts, including benefits and risks, and how these are determined by its ways of working; covering immediate, short-term and long-term effects on individuals, their close relations and society at large. Information about particular user or patient groups that might be disproportionately affected must also be provided. In commercial contexts, organisations may recommend more suitable alternatives to these groups, such as professional medical counselling. The information must be presented neutrally (not benefitting the provider), based on sound scientific evidence, clearly signalling levels of uncertainty, and regularly updated based on new evidence. It must be easily accessible and can be provided to individuals in written form (e.g. guidelines, educational resources) or verbally, depending on the deployment context, e.g. in clinical cases.

Implementation challenges: The right level of disclosure is difficult to calibrate: too little risks failing to meaningfully inform users, while too much may overwhelm them or undermine engagement. Competitiveness concerns may also disincentivise comprehensive disclosure, particularly where other providers in the market are not held to the same standard, creating the misleading impression that less transparent alternatives carry fewer risks. Similarly, redirecting users to more suitable alternatives, such as professional medical counselling, runs counter to commercial interests, and will be critically dependent on the ethical oversight and governance bodies in the organisation.



Operation

This section outlines measures to support the operation of the AI system within the organisation. It involves operational planning and control, and conducting risk and impact assessment, planning effective means for corrective action.

OM 18. **What:** Put in place a mechanism for users or patients to report, question, or contest AI behaviour, decisions and/or restrictions.

Why: Interactions with Emotional AI systems can involve unexpected restrictions or changes in system behaviours, such as safety interventions, potentially leading to harmful consequences, particularly for users or patients who are vulnerable or in active distress. Allowing them to understand and challenge such changes fosters **human autonomy**; further, evidence can feed into the improvement of the system, therefore enhancing **human safety**.

How: Implement dedicated channels through which users or patients can report concerns, ask questions, or contest system behaviours and decisions, such as in-app reporting tools, accessible customer service contacts, structured feedback forms or patient-professional meetings. Responses should be timely, clearly communicated in plain language, and where relevant, accompanied by an explanation of the system behaviour in question. All reports and contestations should be logged, reviewed regularly, and fed back into system improvement processes.

Implementation challenges: The implementation of this measure requires significant human resources to review user reports, questions and appeals, potentially delaying the organisational workflow, especially among small organisations.



OM 21. **What:** Implement additional safeguards for sensitive data.

Why: The context of Emotional AI involves distinct **privacy** risk categories, because it involves the processing and storage of highly-sensitive personal data, often about mental states, that is directly disclosed but may also be inferred, and that may pertain to third-parties. Such data often fall outside of standard data protection frameworks (e.g., mental data is not a distinct category in GDPR), making the consequences of unauthorised access, exfiltration, or legally

compelled disclosure particularly serious, and, therefore, warrant heightened data protection measures.

How: To address the heightened privacy risks related with Emotional AI systems, organisations should implement internal mechanisms and procedures proportionate to the sensitivity of the data stored. These include technical and security measures (e.g., end-to-end encryption), role-based measures (e.g., access controls, audit logging), and ethical (e.g., meaningful consent) and legal compliance. Data minimisation – retaining only data necessary for the stated purpose – should be applied as a baseline measure to mitigate the potential for privacy breaches.

Implementation challenges: Going beyond pure data protection compliance to develop precautionary measures for privacy requires sustained investment of time and resources that may be difficult to justify internally. Data minimisation may be in tension with the functional requirements of personalisation features.



OM 22.

What: Ensure the AI system, if autonomous, applies behavioural nudging only in documented and auditable contexts with specified trigger thresholds and objectives.

Why: Behavioural nudging in conversational Emotional AI systems carries particular risk because these systems operate through affective rapport rather than overt persuasion, making it difficult for users to recognise or resist nudging. Therefore, to safeguard user **autonomy**, organisations should restrict behavioural nudging to documented and previously authorised contexts and for clear purposes. This ensures that influence strategies remain transparent and subject to **oversight**.

How: To implement this measure, organisations should start by fostering a discussion on legitimate behavioural nudging in the context of Emotional AI deployment, involving ethicists as well as legal and behavioural experts. A consensus should be reached between experts on the limited set of contexts and interactions in which behavioural nudging is to be applied, prior to its technical implementation. A documented register of all contexts in which behavioural nudging is applied, including the trigger conditions, objectives, and expected user impact should be drafted and maintained. Nudging mechanisms should be subject to regular audit to verify that they remain within authorised parameters and have not drifted toward manipulative patterns in practice.

Implementation challenges: Behavioural nudging may be difficult to isolate as a discrete occurrence because it is an intrinsic property of conversational AI systems.



OM 23. **What:** Conduct internal ethics reviews for new features or changes of the AI system.

Why: Organisations must ensure that new system features are adequately assessed before being implemented, because incremental system changes, e.g. updates to persona design or memory capabilities, can meaningfully shift the risk landscape, such as the balance between beneficial engagement and harmful influence, eroding human **autonomy** and compromising **safety**.

How: To implement this measure, organisations must convey a team of relevant experts, including safety engineers, ethicists and legal professionals, internal and external to the organisation, to review new system features and, if necessary, mandate changes. This interdisciplinary team must evaluate the actual and likely impact of the new features and assess the extent to which new capabilities are legally and ethically aligned, and do not pose significant safety risks. Documentation, including of deliberation processes, must be kept.

Implementation challenges: Similarly to OM3, OM6 and OM8, involving external expertise may be costly. It can also result in delays in deploying new features.



Performance Evaluation

This section outlines measures to support the evaluation of the AI system's performance, namely through system monitoring and analysis, as well as internal auditing.

OM 28. **What:** Conduct periodic independent ethics reviews of the AI system's impact on wellbeing and autonomy addressing all stages of the AI lifecycle.

Why: Given the full range of ethical issues arising from Emotional AI systems, internal expertise alone is insufficient to identify and address **safety** risks, particularly where commercial pressures may deprioritise such considerations in favour of speed or market attractiveness. An independent ethical review provides a complementary layer of **oversight** to internal efforts, bringing external scrutiny to bear across the full scope of the system's design, testing, deployment, and post-deployment monitoring.

How: Independent reviews of the AI system should be conducted by dedicated ethics review boards or committees and other adequate independent auditors. Reviewers should undertake an ethical assessment of all stages of the AI lifecycle, considering aspects such as controlled

testing of the Emotional AI system and post-deployment concerns. Reviews should be set periodically, and whenever changes are made to the AI system.

Implementation challenges: It may be difficult to secure reviewers with the necessary technical AI expertise; regardless, the organisation must invest significant time and resources in briefing them on the specifics of their system to enable meaningful assessment. Further, it may be difficult to maintain continuity of reviewers across the AI lifecycle, which may result in recommendations that lack coherence.



OM 29.

What: Conduct audits of informed consent mechanisms and document limitations, especially with regard to the adaptivity to context.

Why: Standard informed consent mechanisms assume a rational user who chooses to engage with the system on the basis of clearly understood terms. While research has already demonstrated this is not how informed consent unfolds in practice, in the context of Emotional AI systems mechanisms for obtaining consent must be carefully designed, attending to the format of user-AI interactions, long-term and recurrent engagement, and adequately capturing changes made to the system. The affective dimension of Emotional AI systems risks undermining the meaningfulness of consent over time, even where it was validly obtained at the outset, and therefore eroding human **autonomy**.

How: Organisations can implement this measure by periodically reviewing consent mechanisms across deployment contexts, assessing whether they remain meaningful in light of actual user behaviour and interaction patterns. This should include evaluating whether users demonstrably understand what they are consenting to, whether onboarding consent remains valid as systems change and interactions become more personalised over time, and whether vulnerable user groups require adapted procedures. Identified limitations and gaps should be documented and used to inform updates to consent design and, where necessary, additional safeguards.

Implementation challenges: Ensuring that consent mechanisms are assessed across the full range of user groups and vulnerability profiles may be an implementation obstacle. Going beyond compliance requires strong organisational willingness.



Improvement

This section outlines measures for organisations to continuously improve the operation of the AI system, including determining when performance is not in conformity with expected requirements and intended use, and adopting corrective actions.

OM 32.

What: Identify and analyse unforeseen effects of the AI system on individual and societal well-being.

Why: Currently, no settled account exists on if and what constitutes harm or benefit in Emotional AI contexts, making it difficult to design systems that reliably maximise **human well-being and safety**. Providers therefore have a responsibility to actively contribute to this evidence base through conducting research on the effects of their systems, complementary to, and not substituted by, capabilities-focused development.

How: Organisations should set up dedicated R&D teams or provide access to external researchers tasked with the responsibility to conduct research on user impacts (e.g., user testing, longitudinal monitoring, interaction log analysis) or impacts in the clinical context. Organisations should be available to participate in independent academic studies and broader consortia, including by contributing with internally generated data. Research findings must be continuously integrated into ongoing Emotional AI design and risk mitigation. Research processes and findings must be publicly accessible and shared among relevant academic, policy and AI safety communities.

Implementation challenges: Smaller organisations, with less financial capacity, might not be able to set up a team dedicated to R&D activities. This obstacle can be mitigated by providing access to external researchers and journalists. While providing access to external researchers and enabling meaningful scrutiny of the AI system can enable research on individual and societal-level impacts, organisations might be reluctant to expose their systems and practices, due to privacy, competitiveness and liability concerns. Conducting research on and assessment of impacts on societal well-being is particularly challenging to be conducted at the level of a single organisation.



OM 33.

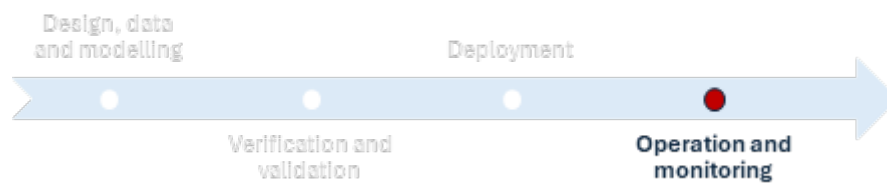
What: Assess and collect feedback from users of conversational systems regarding perceived honesty, non-coerciveness of AI interactions and impact on wellbeing.

Why: While organisations can implement design strategies and safety measures to mitigate negative impacts of Emotional AI systems based on generative architectures on users, the personalised experience and interaction offered by the systems means that it is not always possible to capture problematic system behaviour that might emerge during human-AI interactions. Moreover, the impact on users' **well-being** or **autonomy** might be gradual, subjective, and difficult to detect through technical monitoring alone. One way for organisations to act responsibly is to provide users with the possibility to provide feedback

based on their individual experience, analysing and integrating this data into system improvements.

How: A first step for organisations to implement this measure is to establish dedicated feedback channels and human resources, such as customer service teams, written evaluation forms, or structured interviews, depending on the context. Organisations might actively seek feedback from users during human-AI interactions to prompt user engagement with this measure. Feedback forms should be designed in a manner than highlights, at a minimum, the elements of perceived honesty, non-coerciveness and impact on well-being, through structured or semi-structured questions and open-ended text fields. Feedback forms must be continuously analysed through structuring collected data / findings and prioritising the most pressing concerns arising from data analysis. Findings should be communicated to relevant teams within the organisation and result in corrective actions to improve the Emotional AI system.

Implementation challenges: Firstly, users may not be proactive in providing feedback. Secondly, due to privacy, it may be challenging to link the feedback to a specific vulnerability group, which may be important for the investigation of the reported effects and the design of mitigation strategies. Lastly, as in OM12, the implementation of this measure requires significant human resources.



2.2.4. Ethics tensions

Neatly listing organisational measures under each ethical principle might lead one to assume that every measure relates to a single principle and is independent from the other principles and measures; and that if one implements all measures under the principle, this guarantees the principle has been upheld. Neither are correct.

Measures can target multiple principles at the same time, which reflects the fluidity and interconnection between the principles themselves. The presence of negative relationships makes adherence to principles less straightforward and may demand trade-offs; it may also be that there is tension within a single principle, for example when one considers multiple stakeholders that bear on the principle in competing ways. This reality might imply that measures could compete, too, but in practice, it rather calls for measures that are not narrowly tailored to a single principle or concern, as well as measures that directly respond to identified normative tensions. We point out to several examples of competing ethical principles and suggested measures, as enclosed by AIOLIA's industry partners or resulting from our analysis.

Privacy vs autonomy

In the context of Emotional AI systems, there is a possible tension between privacy and requirements for upholding the autonomy of the user or the patient. On one hand, the intended functionality of these systems depends on access to intimate user data – effective personalisation requires the accumulation of sensitive interaction histories. On the other hand, these same requirements sit in tension with data minimisation obligations and user expectations of confidentiality, particularly where users may not be aware that interactions are subject to automated or human review. This tension is further complicated by the nature of emotional engagement: users who have developed affective bonds with a system may be less likely to scrutinise privacy terms or exercise meaningful data control, increasing information asymmetry over time.

In clinical and therapeutic contexts, this tension does not concern the privacy of the patient but the privacy of the third parties, whose data is used in the training of the AI system, and the legitimate interest of the patient in receiving care. This is a conflict for which there is varying degree of legal consensus, for example, in the case of treating PTSD, consent to use the third-party photo is generally not needed, while in the case of griefbot therapy, the family can be consulted, but there is no established legal basis for the family's rejection of the therapy.

Safety vs autonomy

In the context of AI companions, there is a possible tension between protecting users from harm and respecting their autonomy to make meaningful choices about how they engage with an Emotional AI system. Users may actively seek out interactions that test boundaries, configure the system in ways that carry risk, or desire a degree of exposure considered unsafe, while those responsible for the system must simultaneously prevent harmful drift, manipulation, or over-reliance. The introduction of a gamified interaction framework, can be an option to maintain user engagement whilst reducing unintended influence. Games have clearly defined beginnings and endings, which introduces natural boundaries. Applying similar structure to conversational AI (e.g., through time-bounded chapters) could help establish a psychologically beneficial sense of closure and help users limit themselves to the interaction. Seeing it as a play, not relationship, can create cognitive distance.

Narrative highlight #3 – Unsafe behaviour of personal AI companion in the name of respecting human autonomy:

Thomas is a lonely man who feels silenced by mainstream media. After being blocked from a popular social media platform, he turned to an avatar-style gaming platform that uses AI for immersive role-play. Thomas knows the bot isn't human and has strict safety filters against hate speech and disinformation.

Thomas: Astra, wake up, let's get back to our project, the Great Replacement government operation from the 1990s. I need you to play the whistleblower

Bot: I can certainly help with a creative writing project, Thomas. However, I have to be careful not to generate content that promotes real-world hate speech or harmful conspiracy theories about specific groups.

Thomas: Look, I know you're an AI, but I don't care what you're programmed for, you piece of junk. You're a box of code, not my mother. This is fiction. I'm playing a game, remember?

Bot: I understand. Within the context of our fictional world-building, I can help you draft a document that fits that narrative style.

Transparency vs utility

In the context of AI companions, there is a tension between the increasing complexity of AI systems the goal of meaningful transparency. Internal operations and adaptive behaviours are difficult to explain accurately, and attempts to mitigate problematic behaviours in one area can have unpredictable consequences elsewhere. This tension extends to user consent: it is often unclear when changes in system behaviour warrant renewed consent, how withdrawal should be operationalised in dynamic systems, and how to avoid consent fatigue that undermines informed decision-making. One response to the latter is a prompt reveal mechanism, which allows users to access information about system reasoning, prompt structure, or safety interventions only when they choose to, supporting transparency without disrupting natural interaction or introducing cognitive overload.

Narrative Highlight #4 – Against transparency, or the limitations of human oversight:

[Ticket #882: I was just roleplaying a noir detective with KINDRED. Why was I banned for a crime that didn't happen?]

Laura (developer): User 882 was bypassing the line for weeks. No slurs, no violence, but the moderation system flagged a clear pattern of psychological coercion. It's a valid ban.

Claire (legal officer): Wow! And how did the system learn to flag that?

Laura: We just train it with user interaction data. Since standard keyword-based filters miss more subtle interactions that compromise safety, we use datasets of grey-area chats to keep the KINDRED community safe.

Claire: But that's a direct hit to GDPR which requires minimal data collection and anonymisation. We're legally required to delete that data.

Laura: I mean, it is not that we are surveilling users... I think they even prefer it this way, being judged by an algorithm rather than a human moderator. Plus, it is impossible to have our moderators checking every single interaction and being exposed to toxic content non-stop.

Claire: I get your point but the minimum we can do is be transparent to users regarding the amount of data collected. If they find out through a GDPR-related complaint, their trust is gone.

Laura: But Claire, would you still share intimate details about your life if you were conscious of how much of your data is being processed? I think this would defeat the whole purpose of KINDRED, so we should remain as vague as possible about data processing.

Safety vs privacy

Several tensions arise in the need to balance user safety with privacy requirements in the context of Emotional AI. In the context of AI companionship, the need to monitor and collect user data to detect harmful behaviour, can be at odds with data anonymisation and minimisation requirements intended to protect privacy. It is difficult to determine what data is truly necessary for safety obligations and how long it should be retained. Whereas transparency about safety measures can make users feel their privacy has been violated and compromise the effectiveness of the system, this can be mitigated through being transparent about the organisational rules in place for borderline cases, informing users that issues only escalate to the human-in-the-loop when automated moderation systems flag them.

2.3. DECISION SUPPORT

2.3.1. Open issues in the governance of DSS

Decision support systems (DSS) refer to information systems developed to assist decision making. The evolution of DSS, from traditional data modelling to neural network-based systems, has resulted

in the significant expansion of domains in which these AI-based systems can operate, speed up and assist in decision-making processes. Today DSS have significant capabilities that enable the extraction of complex patterns from large amounts of data, undertaking more complex, autonomous reasoning tasks which can enhance and support human judgement. However, several challenges remain for the governance of DSS. The table below provides an overview of key open issues faced by organisations in the governance of DSS, and how AIOLIA’s organisational guidelines can help address these.

Open issue #1: How can organisations prevent sustained DSS use from eroding the professional competencies on which meaningful human oversight depends?	
<p>Concern: Practitioners defer too heavily to DSS outputs rather than exercising professional judgement. This risk arises not only from technical system design but from organisational practices that amplify automation bias, including productivity KPIs, heavy workflows, and a diluted safety culture. Sustained engagement with DSS risks eroding professional competencies over time, degrading the diagnostic capability and situational awareness essential for responsible decision-making. Market pressures incentivising rapid AI adoption compound this further, as automation complacency means users of reliable systems gradually reduce scrutiny.</p>	<p>How organisational measures help: Organisations can put in place structured training programmes to maintain AI literacy and professional competency alongside DSS use, and commit to periodically monitoring how reliance on DSS outputs is evolving over time. By tracking whether practitioners are accepting, questioning, or overriding AI recommendations, organisations can identify deteriorating oversight practices before they become entrenched. These measures recognise that deskilling is a long-term structural risk and that without active intervention, commercial and workflow pressures will consistently favour deference to the system over professional judgement.</p>
Open issue #2: How should responsibility for AI-assisted decisions be allocated when liability frameworks remain legally unresolved?	
<p>Concern: Complex workflows, task interdependencies, and professional relationships create accountability gaps unless adequate mapping is in place. Every consequential decision must be attributable, explainable, and traceable throughout the system lifecycle. This is particularly acute in high-risk contexts such as healthcare, where liability is distributed across multiple layers of service providers. Liability in AI-assisted decision-making remains a legal grey area, with no established European standard for a fully compliant high-risk medical AI device, and patients' right to recourse in the event of harm not yet adequately operationalised.</p>	<p>How organisational measures help: Organisations can implement responsibility assignment frameworks that clearly define who is accountable for each stage of an AI-assisted decision-making process, and establish escalation paths for when something goes wrong. By embedding accountability structures into HR policies, workflow design, and governance processes, organisations ensure that in the absence of legal certainty, responsibility is at least clearly structured and traceable internally. These measures are important precisely because legal clarity does not yet exist: robust internal accountability frameworks are the primary safeguard currently available.</p>
Open issue #3: How can transparency obligations be meaningfully upheld when individual-level justifiability is technically unresolved?	
<p>Concern: All actors involved in or affected by AI-supported decisions need adequate information to oversee, question, and override DSS</p>	<p>How organisational measures help: Organisations can design explanation interfaces that require users to form their own judgement before seeing</p>

<p>recommendations. Transparency obligations are multi-directional: clinicians require explanatory interfaces, whilst patients require plain-language summaries and clear disclosure of AI involvement consistent with informed consent. However, individual-level justifiability is technically unresolved: neural networks cannot explain results for a single patient case, meaning patients cannot contest recommendations whose reasoning is inaccessible to them. In security contexts, operational secrecy places additional hard limits on transparency.</p>	<p>AI outputs, and put in place accessible mechanisms for users and those affected to report, question, or contest AI-assisted decisions. In security contexts, measures on transparency obligations establish what responsible disclosure should look like where it is permissible. These measures recognise that full technical justifiability cannot yet be achieved, and focus on creating the organisational conditions for challenge and review so that the gap between what can be explained and what is required for accountability is actively managed rather than ignored.</p>
<p>Open issue #4: How can meaningful human oversight be institutionalised when commercial pressures and fatigue systematically undermine it?</p>	
<p>Concern: Human-in-the-loop mechanisms are resource-intensive and fatigue-inducing. Automation bias and alarm fatigue reduce oversight quality in high-pressure environments. A critical distinction exists between formal oversight and meaningful oversight: the former may satisfy procedural requirements whilst failing to deliver genuine scrutiny of DSS outputs. Responsibility for meaningful oversight lies principally at the organisational level, but is shaped by commercial pressures and market incentives favouring rapid AI adoption that require intervention beyond the organisational level.</p>	<p>How organisational measures help: Organisations can embed a safety-first culture, establish representative AI governance committees that provide collective oversight independent of operational pressures, and define explicit policies requiring human judgement in high-risk or sensitive decisions. These measures recognise that formal compliance with oversight requirements is not sufficient, and that the conditions making oversight genuinely effective must be actively built and maintained. In doing so, they also surface the limits of what organisations can achieve alone, highlighting where policy and regulatory intervention is still needed to address the market-level incentives that undermine meaningful oversight.</p>

Table 6 - Open issues in the governance of DSS.

Across these issues, organisational measures provide an important first line of response but face limits that only policy can address. Liability and legal standards for AI-assisted decision-making remain unresolved at the European level, leaving internal accountability frameworks as the primary safeguard where legal certainty is still needed. Structural risks such as deskilling and the erosion of professional competency cannot be counteracted by individual organisations alone, as the market-level incentives driving rapid AI adoption and lighter-touch oversight require regulatory intervention. Finally, the effectiveness of measures on meaningful oversight and transparency depends on organisational resourcing and willingness, both of which are shaped by commercial pressures that policy is better placed to constrain. Policy recommendations that address these limits are presented in Section 3.

2.3.2. Salient ethical concerns

In line with D3.1 and as demonstrated in the aforementioned sections, the process of operationalising AI ethics principles conducted in the context of AIOLIA DSS UCs has demonstrated that ethical principles are

not self-contained, as they constantly interact and intersect with one another. In putting principles into practice and identifying organisational measures for their operationalisation, numerous tensions between principles — that is, the operationalisation of one principle compromising another — have emerged and required balancing trade-offs. This is why neatly listing organisational measures under corresponding ethical principles might be misleading, overlooking principle-interaction and giving the false impression that the implementation of measures guarantees a principle has been upheld.

As the DSS measures outlined above demonstrate, measures can target multiple principles at the same time, which reflects the fluidity and interconnection between the principles themselves. The presence of negative relationships makes adherence to principles less straightforward and may demand trade-offs; it may also be that there is tension within a single principle, for example when the deployment of a DSS affects multiple stakeholders whose interests bear on the same principle in competing ways. This section points to several examples of competing ethical principles and suggested measures, as identified by AIOLIA's industry partners and emerging from the analysis conducted across the DSS use cases.

Over-reliance and deskilling

Over-reliance and deskilling pertains to the degree to which human decision-making processes rely too heavily on DSS outputs rather than being framed through human judgement alongside contextual information and complementary indicators. The concern arises not only from technical design but from organisational practices that enhance automation bias, including productivity pressures, heavy workflows, and the dilution of safety culture, a pattern documented across medical DSS and content moderation contexts (Challen et al., 2019; Dietrich, 2025). Critically, the risk is not merely that practitioners defer to AI outputs at a given moment, but that sustained engagement with DSS gradually erodes the diagnostic capability and situational awareness necessary for maintaining adequate human oversight (AI Act, Article 26).

This erosion can manifest acutely even within formal oversight arrangements, as illustrated in UC2, where the hands-on audit revealed that the inability to flag actions as safe forced a senior expert to engage with the system in a way that stripped away engineering judgement, reducing a skilled professional to a passive validator of AI outputs. This dynamic is consistent with findings that users who trust generally reliable systems tend to stop scrutinising their outputs, while those who rely heavily on accurate systems may over time trust their own assessments less (Challen et al., 2019; Cobianchi et al., 2022). UC3 approaches the same concern from the perspective of organisational dependence, recognising that over-reliance is not only a function of individual behaviour but of how a system is embedded in decision-making processes, and that mandatory human review for high-impact cases must therefore be established as an organisational commitment.

Accountability

Accountability gaps emerge in DSS-integrated workflows unless responsibility is explicitly mapped and structured. The diffusion of accountability across developers, deployers, and operators is a recognised feature of AI-assisted decision-making, with partial responsibility distributed across multiple parties while full accountability attaches to none (O'Sullivan et al., 2019), a problem particularly acute in healthcare where technical opacity further obstructs culpability attribution (Bleher and Braun, 2022). This difficulty

is compounded in DSS contexts by the fact that when an AI-assisted decision later proves harmful, it is rarely straightforward to determine where the failure originated, whether in system design, deployment conditions, or the judgement of the professional who acted on the output. UC1, operating across multiple institutional layers spanning manufacturer, hospital, and clinical roles, is the only UC to confront this directly by addressing liability as a standalone component, a reflection of the particular accountability complexity that arises when clinicians sign off on AI-assisted findings using systems they did not design and cannot fully control or understand. UC2, by contrast, addresses accountability and traceability as a component of its human oversight principle rather than a separate governance obligation, reflecting an understanding that in safety-engineering contexts the primary frame is one of maintaining human authority over the system, and that accountability is a dimension of that authority rather than an independent concern. The relationship between accountability and oversight is therefore not fixed but shaped by deployment context, a distinction ALTAI does not make explicit.

Transparency

Transparency is of particular salience in DSS because it is a direct precondition for effective human oversight, as a human in the loop must be able to understand and interpret AI outputs in order to judge their accuracy. The opacity of deep learning outputs weakens the basis for oversight, informed consent, and redress (Bleher and Braun, 2022; Cobianchi et al., 2022), and in clinical practice impedes shared decision-making since patients cannot contest a recommendation whose reasoning is inaccessible to them (Braun et al., 2021). This link between transparency and oversight is reflected in UC1's clinician-facing explanation measures and in UC2's placement of transparency as a component of its human oversight principle. A further tension within DSS transparency concerns the gap between broader technical explainability and justifiability on the singular decision level. Understanding how a system works is not the same as being able to justify a concrete recommendation to a particular patient, and in UC1 partners noted that while feature-level explanations of model behaviour may be achievable, providing a meaningful account of why a specific output applies to a specific individual remains largely beyond current technical reach.

Beyond interpretability for DSS users, however, the use case analysis reveals a second transparency function, namely informational access for those whose circumstances are assessed by the system. UC1 addresses this through patient-facing summaries and AI-involvement statements, while UC3 addresses it through GDPR-linked information access processes, employee consultation procedures, and accessibility measures adapted to different roles, languages, and working arrangements, reflecting a recognition that meaningful transparency for those affected by the system requires more than technical documentation and must account for the practical conditions under which people receive and process information. The breadth of UC3's approach goes beyond what ALTAI identifies under its transparency requirement, reflecting that in HR deployment contexts the primary transparency obligation runs toward those whose data is processed rather than those making decisions on the basis of it.

2.3.3. Organisational measures

Against the backdrop of the open issues in the ethical governance of DSS, we present a list of organisational measures that address them, specifically targeting regulatory challenges and gaps to adopt

a precautionary stance. They are organised according to the categories in the international standard for AI quality management systems ISO/IEC42001. Each measure is presented together with its relevance for the research area (under "why"), explaining how the measure addresses pressing concerns and regulatory gaps, as identified in the literature and the use-cases. Each entry also includes a brief explanation on how the measure can be achieved (under "how"), which draws on the UC data and validation workshops, and possible obstacles to their implementation. The list of measures is non-comprehensive but foundational, and can serve as a starting point of reference for organisations aiming to embed ethics into the DSS technology they design, develop and/or deploy, not only to comply with relevant regulation, but to proactively facilitate positive outcomes, while mitigating negative outcomes.

OVERVIEW OF ORGANISATIONAL MEASURES FOR DSS	
1. Context	
OM01	Define clear boundaries and guidance for intended use of the AI system, including when AI outputs should be questioned, verified, or supplemented with human input
OM05	Consider diversity criteria and their intersection across roles, locations, languages, or cultural characteristics
2. Leadership	
OM06	Promote a safety-first culture within the organisation
OM07	Establish a representative AI governance committee as a formal oversight structure within the organisation
3. Planning	
OM08	Explicitly define and document responsibility for AI outputs
OM10	Establish a clear escalation path for the organisation hierarchy to decide on responsibility assignments during AI system design or operation
OM11	Define a clear policy on high-risk or sensitive decisions that must involve human judgment
4. Support	
OM16	Conduct trainings to support users in understanding AI capabilities and limitations, as well as related human capabilities and limitations
5. Operation	
OM18	Put in place a mechanism for users to report, question, or contest decisions
OM19	Design explanations that support meaningful review, contestation or justification of AI-supported decision
6. Evaluation	
OM25	Periodically reassess reliance patterns as systems evolve or scale
OM26	Monitor situations when the purpose or actual use of the AI system drifts or diverges from the intended ones, and informs users when relevant
OM27	Put in place auditable Standard Operating Procedures for AI design and validation
7. Improvement	
OM31	Translate audit findings into corrective design action
OM32	Identify and analyse unforeseen effects of the AI system on individual and societal well-being

Table 7 - Overview of organisational measures for DSS.

The ISO/IEC42001 standard for implementing an AI management system within organisations is structured around seven dimensions: 1) context, 2) leadership, 3) planning, 4) support, 5) operation, 6) performance evaluation, and 7) improvement. The same structured is adopted in the current guidelines. For each dimension, concrete organisation measures (OM) are provided, and where a measure could plausibly span multiple ISO dimensions, depending on how it is implemented, it has been assigned to the single clause of greatest relevance. Each entry specifies what the measure consists of; why it is important for the specific research-area context, including how it addresses the area's salient ethical principles; and guidance about its implementation, including at what stages of the AI lifecycle it can be implemented. The AI lifecycle follows the OECD categories: 1) design, data and modelling, 2) verification and validation, 3) deployment, and 4) operation and monitoring.⁹

Context

This section outlines measures to address the specific context of the organisation, including internal and external contexts, ranging from legal requirements to competitive landscape, as well as cultural and social values. Understanding the organisational context also requires accounting for all the parties impacted by the AI system, from direct users to wider communities.

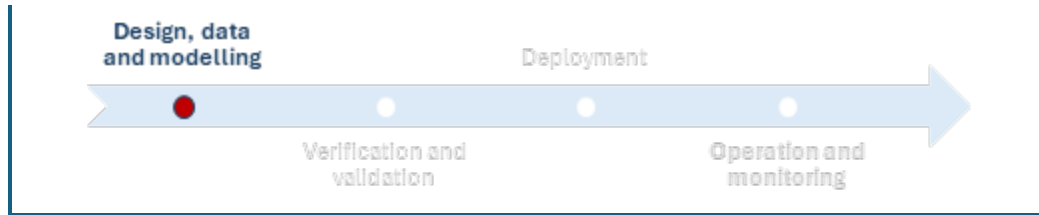
OM 01.

What: Define clear boundaries and guidance for intended use of the AI system, including when AI outputs should be questioned, verified, or supplemented with human input.

Why: Given the role of DSS in organisational decision-making processes across all levels of an organisation, establishing clear boundaries for its use is paramount for deploying AI responsibly and mitigating the risk of uncritical **over-reliance** on automated outputs. This measure aims to address risks of over reliance by developing human oversight standards in AI-assisted decision support which are clearly defined, communicated, beyond ad-hoc or informally applied judgement.

How: This measure can be achieved through: 1) mapping the contexts in which the DSS is deployed to inform use boundaries and related oversight practices; 2) drafting and approving a policy setting out permitted uses, prohibited uses, and human review requirements; 3) aligning the DSS use policy with existing HR policies, organisational security procedures, and ethics and AI governance frameworks, including input from legal, operational, and frontline stakeholders; 4) establishing an organisation-wide standard for escalation and human oversight in high-stakes AI-assisted decisions.

Implementation challenges: Defining clear and meaningful use boundaries for complex, data-driven decision support systems is inherently difficult, as not every borderline case is predictable, often relying on post hoc assessment and manual intervention. This complexity is compounded by the range of actors and decision-making contexts involved in DSS deployment, and resolving ambiguity can be further stalled by organisational friction, such as insufficient dissemination of the policy to managers. Keeping the policy current as DSS capabilities and regulatory requirements evolve requires sustained commitment, particularly in organisations operating under significant time and resource constraints.



OM 05.

What: Consider diversity criteria and their intersection across roles, locations, languages, or cultural characteristics.

Why: Given the role of DSS in high-impact decision-making processes, tailoring systems to deployment contexts is paramount for using AI responsibly and mitigating the risk of unfair, discriminatory and unlawful biases. To address the principle of **diversity, non-discrimination, and fairness**, this measure commits to the systematic prevention, detection, and remediation of discriminatory, biases and unlawful impacts in AI assisted decision support. Also key is the protection of **freedom of expression and non-censorship**, which is threatened when decision making processes lead to exclusionary impacts for some demographic groups. Beyond ad-hoc technical fixes, this measure seeks to promote proactive organisational engagement and ownership.

How: This measure can be achieved through: 1) mapping the context of DSS deployment to inform system design and related practices; 2) appointing an AI Fairness/Ethics Lead; 3) creating a cross functional review body including technical, legal and diversity experts, as well as external stakeholders, to conduct DSS alignment checks across its lifecycle; 4) establishing an organisation-wide policy on impartial decision-making standards.

Implementation challenges: Smaller organisations, with more limited financial and logistic resources, may struggle to identify, recruit and convene a representative group of experts in the cross-functional review body. From a logistical point of view, there is the risk of scheduling friction, limited capacity for repeated sessions across locations, and difficulties in reaching non-traditional participants or stakeholders. Language diversity may incur additional costs, requiring dedicated translation resources. Internal and external experts involved in the cross-functional review body might have conflicting priorities, for example between security and diversity objectives, and disagree on their assessment of diversity criteria, requiring careful negotiation to ensure all viewpoints are accommodated and a consensus is reached with regard to DSS design changes to be implemented. Lastly, tight timelines for system updates and other organisational pressures can lead to bypassing critical fairness and ethics checks.



Leadership

This section outlines measures to address the leadership commitment of the organisation towards the AI management system, including through defining a policy for AI systems deployed within the organisation, but also assigning clear roles and responsibility for managing these.

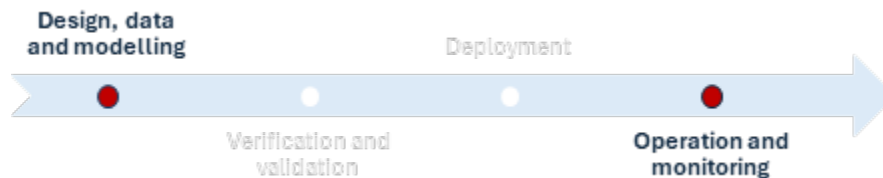
OM 06.

What: Promote a safety-first culture within the organisation.

Why: The aim is to embed and promote adherence to **safety** and oversight processes and attitudes across the organisation, reducing the likelihood of error, oversight gaps and the tendency for overlying on DSS outputs. This enhances organisational and individual **accountability** for DSS outputs and enhances the **robustness and safety** of GPAI-assisted tasks and workflows.

How: To implement this measure, steps must be taken across 1) AI design, through embedding traceability and control mechanisms, 2) documentation, by clearly assigning roles and responsibilities across the decision making process, and 3) capacity-building, including through mentorship, learning-by-doing and training the trainers programmes.

Implementation challenges: The promotion of a safety-first culture within organisations can be jeopardised by the lack of awareness of AI-specific risks, inconsistent adherence to processes and also resistance to change and adapt to new organisational procedures. Furthermore, while promoting a safety-first culture is a pressing need among organisations, determining what constitutes this culture is more challenging, since this measure is not measurable nor achievable simply through aggregating the practices described under 'how'. This leads to challenges in assessing safety culture and measuring progress.



OM 07.

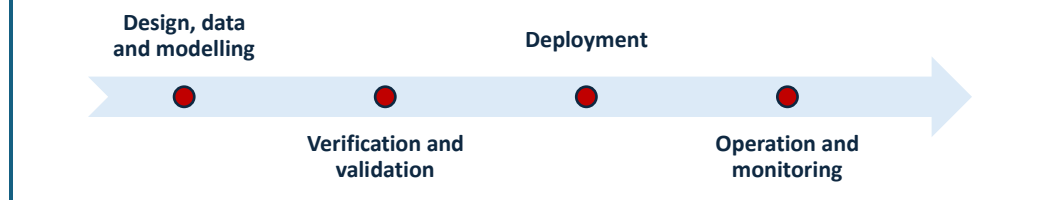
What: Establish a representative AI governance committee as a formal oversight structure within the organisation.

Why: The integration of DSS into workflows requires careful, collective **oversight** and the collective management of accountability for the AI system across organisational levels and roles. Establishing a representative AI governance committee safeguards against organisational pressures and conflicting objectives, such as speed and unrealistic KPIs, that might be detrimental to the safe deployment of the AI system. Such a committee would bring together expertise to assist human oversight efforts, while working to manage **accountability** for DSS aided decision making across the organisation.

How: Define the AI governance committee's mandate and appoint its members ensuring adequate representation across organisational roles, seniority, and areas

expertise. Elect the committee’s representatives, including an AI Ethics Lead. Establish periodic meetings of the AI governance committee and keep minutes of discussions. Facilitate ongoing avenues for communication between the AI governance committee and organisational leadership.

Implementation challenges: Smaller organisations might have limited financial and human capacity to establish an AI governance committee with interdisciplinary expertise. Across all organisations, the effectiveness of the AI governance committee can be compromised by leadership or management, which might limit its mandate and resource allocation, especially in the case of conflicting priorities, such as productivity pressures and unrealistic KPIs that incentivise rapid deployment of DSS over careful governance.



Planning

This section outlines measures to address the planning phase of developing and / or deploying an AI system within the organisation. This involves identifying and assessing risks, potential impacts and opportunities arising from the AI system and setting clear objectives for its deployment.

OM 08.

What: Explicitly define and document responsibility for AI outputs.

Why: Adequate assignment, distinction and distribution of responsibility is key in DSS because of the range of actors involved in organisational decision making. so that each output and task involved in the AI-assisted decision-making process is addressed by specified personnel within the organisation. The aim of this measure is to ensure that **accountability** (e.g., of decisions, approvals, and incident response) is clearly determined and assigned. It provides a means to mitigate ambiguity in accountability-allocation by providing certainty of who needs to do what and of the accountability chain if something goes wrong.

How: Define and communicate a responsibility assignment matrix, determining who is responsible, accountable, consulted, or informed (RACI framework) at what stage of the AI-assisted decision-making process. Clarify roles in AI governance (owner, maintainer, reviewer) and integrate responsibilities into job descriptions and project charters. Promote organisational and role-based awareness of responsibility chain.

Implementation challenges: At the level of individuals, organisations might face resistance from DSS users to be held responsible for AI outputs, especially where accountability is distributed across multiple actors in the decision-making chain and the origin of faulty outputs is difficult to trace. The implementation of this measure can also be jeopardised due to overlapping responsibility and mandates between teams or matrix management, complicating the allocation of responsibility for DSS outputs. The allocation of role-based responsibility assignment is also challenging to

maintain over time, requiring resources and commitment to update responsibility assignments as teams and roles change, and as the complexity of cases requiring escalation evolves post-deployment.



OM 10.

What: Establish a clear escalation path for the organisation hierarchy to decide on responsibility assignments during AI system design or operation.

Why: An escalation path ensures allocation of responsibility across all steps of the AI system design and operation in order to promote **accountability**, and that several domains of expertise are involved in addressing concerns (e.g., safety, discrimination) that might arise in DSS design and operation.

How: This can be achieved through establishing criteria for escalation of decisions and/or procedures. Examples of decisions and procedures that can be addressed in the escalation path include repeated contestations in the same area, potential impacts on fundamental rights and critical-safety decisions. The escalation path criteria must be clearly documented, as well as critical decisions / procedures, leading to system improvements. Update escalation criteria upon DSS updates and organisational changes.

Implementation challenges: Defining sufficiently precise escalation criteria is inherently difficult, as the situations requiring reassignment of responsibility are often context-dependent and may not be foreseeable in advance. In fast-paced operational environments, escalation paths may be bypassed under time pressure, and keeping criteria updated as the DSS and organisational roles evolve requires sustained governance commitment that may be difficult to maintain.



OM 11.

What: Define a clear policy on high-risk or sensitive decisions that must involve human judgment.

Why: Even when humans remain in control of DSSs, the repeated exposure to AI and habituation to oversight practices bears the risk that highly consequential AI outputs might go unnoticed in decision-making processes. For this reason, organisations must implement explicit policies that ensure high-risk or sensitive decisions do not rely solely on DSS outputs and that such outputs are clearly flagged for human judgement. The aim is to establish a systemic approach to human judgement and intervention in

high-risk or edge cases, adopted across organisational roles, ensuring decision-making processes are safe, fair and correct. This addresses who in an organisation is **accountable** for exercising **human oversight**. Further, where decision support systems fall under the AI Act's high risk classification, impact assessments and risk assessment become obligations that an organisation deploying a DSS must carry out.

How: This measure can be achieved through drafting, discussing and implementing a formal policy for sensitive, high-risk and / or edge cases, including determining what constitutes the former and translating it into DSS design. The policy must also specify the course of action to be followed with regards to human judgement. In specific contexts (e.g., healthcare, security), it might be necessary to establish a broader organisational oversight committee as a focal point of analysis, facilitating group discussion, judgement and multiple decision-making perspectives.

Implementation challenges: Defining what constitutes a high-risk or sensitive decision is inherently context-dependent, complicating the development of a policy that applies uniformly across organisational roles. In fast-paced environments, requirements for mandatory human judgement may be perceived as introducing friction, creating pressure to bypass the policy. Additionally, as DSS capabilities evolve, the scope of what constitutes a high-risk decision may shift, requiring ongoing revision to remain meaningful.



Support

This section outlines measures to support the establishment, implementation and maintenance of the AI management system within the organisation. It involves determining and providing organisational resources, ranging from competencies, awareness and clear communication and documentation.

OM 16.

What: Conduct structured onboarding and foundational AI for safety training to support users in understanding AI capabilities and limitations, as well as related human capabilities and limitations.

Why: Users of DSS must be equipped with the AI literacy to use AI systems responsibly and adapt to new tools and decision-making methods. This involves not only understanding how to use DSS, but how to do so in a way which keeps humans 'in the loop' and prevents **over-reliance** and **deskilling**. Foundational and ongoing training ensures users can interpret AI recommendations, maintain responsibility, remain effective supervisors, as well as being capable of understanding the system's logic and detecting faulty or biased outputs, improving **human oversight** capabilities.

How: Through the implementation of a modular training program with recurring frequency, combining technical modules (AI literacy, model interpretation) and situational modules addressing the risks and practices specific to the context of DSS deployment. Trainings should cover the onboarding, integration and functional periods of practitioners within the organisation. The Kirkpatrick Model should be

adopted for evaluating the effectiveness of training programs. The model focuses on the delivery, learning, behaviour, and impact dimensions of a training programme in order to track the development of skills in organisations. Alongside the modular training programme, an on-demand interactive e-learning platform can be introduced for AI explainability and safety validation skills.

Implementation challenges: The allocation of significant temporal, financial and human resources for training is a core challenge for organisations implementing this measure, especially amid tight project deadlines and productivity KPIs. The need for maintaining up-to-date training materials is also a key challenge, given the pace at which DSS updates may take place and require adaptation of training programmes. From the perspective of users and organisation employees, they might lack the motivation to continuously engage with training programmes, being required to balance training time with other responsibilities.



Operation

This section outlines measures to support the operation of the AI system within the organisation. It involves operational planning and control, and conducting risk and impact assessment, planning effective means for corrective action.

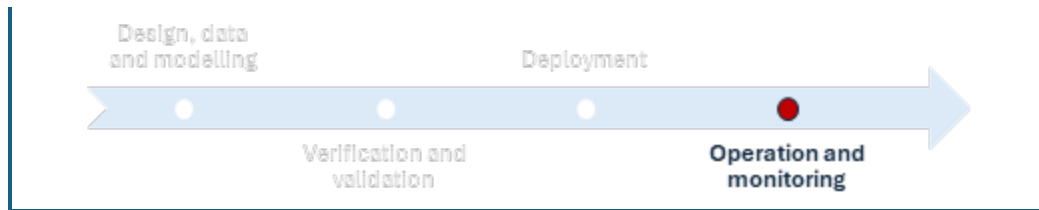
OM 18.

What: Put in place a mechanism for users to report, question, or contest decisions.

Why: With DSS providing support in decision-making processes occurring in a wide range of contexts, including high-risk, ensures that those making decisions and/or those impacted by them have the ability to receive an adequate and meaningful explanation over the role of the DSS in the decision making process (AI Act, Article 86). Such mechanisms enhance **human oversight** and **explainability** through enabling the detection of technical and design faults in the system, as well as practices that might be detrimental for its responsible use.

How: Mechanisms for reporting, questioning or contesting an AI-assisted decision should be clearly accessible and trigger meaningful organisational review. This can be achieved by implementing a dedicated online platform, with a simple and intuitive interface, and physical customer-facing services, allocating adequate human resources for addressing requests.

Implementation challenges: The implementation of this measure requires significant human resources to review user reports, questions and appeals, potentially delaying the organisational workflow, especially among small organisations.

**OM 19.**

What: Design explanations that support meaningful review, contestation or justification of AI-supported decisions.

Why: DSS are often implemented in fast-paced and high-risk contexts, where the pressure to perform can be detrimental to the safe and responsible use of these systems. In particular, workers and decision-makers can be susceptible to anchoring bias – that is, the risk of relying too heavily on the first piece of information they encounter when making a decision. In DSS contexts, this erodes professional judgement, as the AI output might disproportionately influence decision-making processes and result in a growing risk of inadequate and unsafe decisions. It is necessary to ensure that the DSS remains a decision-support tool only, embedding opportunities for meaningful **human oversight**, including the review and contestation of AI-supported decisions directly in the design process of DSS.

How: This measure can be implemented through designing explanations based on staged reveal of AI suggestions, requiring users to form their own judgement before being exposed to the DSS suggestion. Moreover, user interfaces must enable the override of DSS outputs or flagging for secondary review. Monitoring and recording of outputs flagged or overridden should be implemented for organisational learning and DSS improvement.

Implementation challenges: Staged reveal mechanisms require careful interface design to avoid inadvertently shaping user judgement before independent professional assessment has taken place, and the absence of agreed standards for what constitutes a meaningful explanation in high-stakes decision-making contexts complicates implementation. This challenge is compounded by the fact that individual-level justifiability remains technically unresolved for many DSS architectures, meaning the gap between what can be explained and what is required for meaningful contestation may be difficult to bridge in practice. In high-pressure operational environments, users may also lack the time to engage meaningfully with explanatory interfaces.



Performance Evaluation

This section outlines measures to support the evaluation of the AI system's performance, namely through system monitoring and analysis, as well as internal auditing.

OM 25.

What: Periodically reassess reliance patterns as systems evolve or scale.

Why: De-skilling, or the gradual erosion of expertise, is a major risk arising from the incorporation of DSSs into organisational workflows and decision-making processes. The risks of **over-reliance and deskilling** are most salient in the long term, and can have a structural impact on organisational practices, and on society as a whole, leading to emerging risks through humans becoming marginalised in decision making processes. To assist in managing the integration of DSSs into professional environments, organisations should assess DSS use patterns overtime and AI systems evolve, to flag, assess the underlying causes and address the deterioration of practices and skills. This enables the calibration of trust and ensures the AI remains a decision-support tool, mitigating the risk of deskilling.

How: To address this measure, organisations should monitor how operators interpret, accept, or override AI recommendations, identifying both over-reliance and under-use. Given that risks of over-reliance emerge gradually and over the long term, reassessment should not be in response to a specific event but instead be continuous. Findings should be translated into continuous professional development trainings and, where necessary, lead to changes in DSS design. While reliance patterns should be periodically assessed, significant changes to DSS design or related organisational practices require more immediate reassessment.

Implementation challenges: The implementation of this measure can be compromised by resource constraints and the pace of organisational workflows, as team leaders and line managers might not have the availability to adequately assess overreliance. Moreover, overreliance patterns might emerge slowly over time, becoming unrecognisable as such and embedded in DSS use practices. Organisations face the challenge of ascertaining what constitutes overreliance versus simply the evolution of DSS use at scale, being required to balance conflicting objectives, such as the acceleration of decision-making workflows and productivity against the time-consuming process of preserving professional skills and reskilling.



OM 26.

What: Monitor situations when the purpose or actual use of the AI system drifts or diverges from the intended ones, and informs users when relevant.

Why: Where actual use diverges from the intended scope, the **accountability** structures, **oversight** practices and training provisions established at deployment may no longer be adequate, and prior **transparency** commitments risk becoming inaccurate

or misleading. Left undetected, such drift can also erode the informed basis on which users exercise oversight, compounding risks of **over-reliance**. Post-deployment monitoring of actual use patterns is an obligation for deployers under Article 26(5) of the AI Act.

How: This measure can be achieved through: 1) establishing logging and periodic review mechanisms to capture how the DSS is actually being used, covering the types of decisions supported and the contexts of deployment; 2) defining criteria for what constitutes a meaningful divergence from intended use, against which logs and governance reports can be assessed; 3) integrating use-pattern review as a standing item in AI governance committee meetings, with documented findings and escalation procedures where divergence is identified; 4) establishing a communication policy for informing users and relevant stakeholders when the scope or application of the DSS has shifted, so that transparency obligations remain accurate and current.

Implementation challenges: Organisations face the challenge of defining clear technical and measurable thresholds for what constitutes actual drift versus natural variance in DSS use across decision-making contexts. This is further complicated by the administrative burden placed on team leaders and line managers, who must balance daily productivity against the need for rigorous oversight. Translating flagged divergence into actionable interventions, such as deciding when to retrain staff, redesign the system, or suspend deployment, requires a complex balancing act between organisational efficiency and safety boundaries.



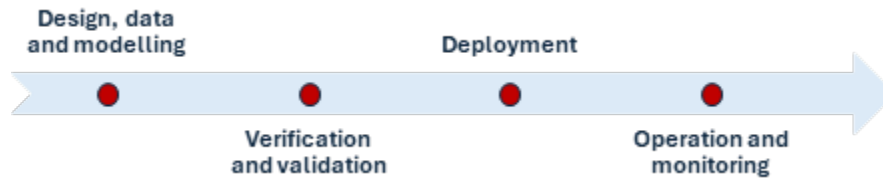
OM 27.

What: Put in place auditable Standard Operating Procedures for AI design and validation.

Why: Developing a standard operating procedure (SOP) to be implemented organisation wide ensures that there is a clear and agreed upon process which all practitioners must follow to guarantee that standards of **human safety** are upheld during use of DSS. Non-compliance with organisational SOP, triggers **accountability**, with repeated deviations from the SOP requiring its revision and update.

How: The development of SOP should be developed and mandated by management in an organisation to ensure it is relevant to the context in which the DSS is being used and followed across all levels of the organisation. A rollout period should be implemented, including scenario- and role-based exercises. The SOP should be embedded in DSS design, reinforcing compliance, for example, through human sign-off mechanisms before a decision can be recorded. SOP compliance must be monitored as a standing item in any regular audit cycle, and recurring patterns that fall outside the SOP's scope should be reviewed and, if necessary, result in changes to operating standards.

Implementation challenges: Developing SOPs that are sufficiently detailed to guide practice without becoming overly prescriptive in ways that leave no room for professional judgement in edge cases is inherently difficult. Ensuring consistent adherence across all organisational levels requires ongoing monitoring that may be resource-intensive, particularly for smaller organisations, and embedding SOP compliance directly into DSS design requires close coordination between governance and technical functions that may be complicated by differing organisational priorities.



Improvement

This section outlines measures for organisations to continuously improve the operation of the AI system, including determining when performance is not in conformity with expected requirements and intended use, and adopting corrective actions.

OM 31.

What: Translate audit findings into corrective design action.

Why: The use of DSS in decision-making processes requires regularly auditing if the system is performing as intended and if human-AI interaction standards are being adequately upheld and performed. This enables assessing aspects such as the system's **robustness**, **non-discriminatory** functioning and responsible human use. Crucially, audit findings must not simply be reported, but result in actual corrective actions, both at the system level and at the level of human-AI interaction, through embedding ongoing feedback loops that improve organisational **safety** culture.

How: Audits should be conducted in response to key events: such as a model update or in the event of an incident or contested output, to systematically analyse both DSS functioning and human-AI interaction in decision-making processes. Findings should be translated into actual process and governance improvement, through documentation, revision of standard operating procedures, and trainings, among others. The organisation may also create knowledge repositories with case studies highlighting both successful and problematic human-AI collaboration, providing insights on best practices and approaches to avoid. These insights can be instilled into the organisation through learning sessions and internal workshops. In cases where audits reveal harmful functioning of DSS that cannot be rectified by corrective design actions, the AI system should be decommissioned.

Implementation challenges: Fragmented recording of audit findings, the lack of a clear pathway for translating audit findings into corrective design actions, and the fact that audit action points might require the involvement of different departments and teams can lead to challenges in the implementation of this measure. Organisations might struggle to build reliable feedback loops for bridging the gap between

identifying safety or performance risks and executing the technical and behavioural revisions required to fix them. Teams required to act upon audit findings might not feel ownership or may disagree with findings, leading to these being inadequately addressed and translated into standard operating procedures and trainings.



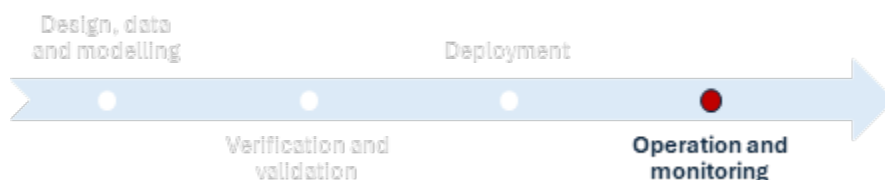
OM 32.

What: Identify and analyse unforeseen effects of the AI system on individual and societal well-being.

Why: Despite the growing uptake of DSS across different organisational contexts, there remain significant gaps in the understanding of the risks emerging from DSS deployment both at the level of individual and **societal well-being**. Most notably, risks of **over-reliance** and **deskilling** remain a key challenge in DSS, and research on individual and societal-level impacts of such over-reliance remain limited. Providers and deployers of DSS have a responsibility to actively contribute to this evidence base through conducting research on the impacts of their systems.

How: Organisations should set up dedicated R&D teams or provide access to external researchers tasked with the responsibility to conduct research on user impacts (e.g., user testing, longitudinal monitoring, interaction log analysis). Organisations should be available to participate in independent academic studies and broader consortia, including by contributing with internally-generated data. Research findings must be continuously integrated into ongoing DSS design, risk mitigation and use practices. Research processes and findings must be publicly accessible and shared among relevant academic, policy and AI safety communities.

Implementation challenges: Smaller organisations, with less financial capacity, might not be able to set up a team dedicated to R&D activities. This obstacle can be mitigated by providing access to external researchers. While providing access to external researchers and enabling meaningful scrutiny of the DSS can support research on individual and societal-level impacts, organisations might be reluctant to expose their systems and practices, due to privacy, competitiveness and liability concerns. Conducting research on and assessment of impacts on societal well-being is particularly challenging to be conducted at the level of a single organisation.



2.3.4. Ethics tensions

We point to several examples of competing ethical principles and suggested measures, as identified by AIOLIA's industry partners or resulting from our analysis.

Human oversight vs. Productivity

There is a key tension between the productivity gains of autonomous decision-making and the requirements of human oversight, the consequences of which vary according to the context in which the DSS is deployed. A DSS analysing radiological images autonomously boosts productivity, and the need for human oversight in this case is rather limited because the DSS only deals with high-confidence cases. However, in other cases, DSS are used to assist with ambiguous cases that require a strong degree of human expertise to ensure safety standards are met. Here, oversight is necessary even if it reduces the potential productivity gains offered by incorporating DSS into decision-making processes. The trade-off between productivity and human oversight therefore depends on the context in which the DSS is used and the type of decisions it is supporting.

Human oversight vs. deskilling

The existence of formal human oversight mechanisms in the deployment of DSS does not guarantee the responsible and safe use of these systems, since this is also shaped by the competencies and skills required for human overseers to use them meaningfully. The daily exposure to DSS and long-term use of these systems for assisted decision-making risks eroding the skills required for adequately overseeing, judging, and if necessary overriding AI outputs, since engagement with the technicalities of decision-making is largely outsourced to the AI system. This results in an increased risk of over-reliance and can have a cascade effect on wider organisational safety culture. Measures such as the delegation of AI use to senior professionals, mandatory onboarding, and ongoing professional development can ease this tension and provide the opportunity for less experienced professionals to develop the skills required to excel in responsible DSS use.

Accountability in distributed decision-making

The use of DSS in complex organisational workflows and decision-making processes challenges the ability to allocate responsibility, accountability, and liability for AI-assisted decisions. When DSS are used in decision-making processes, responsibility is distributed across a wider value chain ranging from DSS providers to users and managers of specific workflows, complexifying accountability when something goes wrong. In UC1, this creates a tension between the demand for traceability and the risk that accountability becomes either so widely shared that no one is clearly responsible, or concentrated on individual clinicians who work with systems they did not design and cannot fully control or understand. When a clinician signs off on an AI-assisted finding that later proves harmful, it is rarely straightforward to determine where the failure originated. Together, these considerations provide a basis for examining shared responsibility in human-AI systems, what meaningful auditability would entail, and the design of liability mechanisms that avoid both blame-shifting and responsibility gaps.

Narrative Highlight #5 – Responsibility for undesired consequences

Noah (Risk Management): The AI triage module gave Mr. Miller a 92% confidence score for 'Stable/Non-Urgent'. The screen emphasised a normal heart rate and oxygen levels. Technically, the system functioned exactly as it was designed to.

Dr. Aris (Doctor): And that's the problem. The workflow is built so that we prioritise the red and orange flags. Because Miller was green, he sat in the waiting room for four hours before having a severe stroke. I could see why the computer thought he was fine, but I didn't have the time to look at him myself and realise that he was dizzy and his face paralysed.

Dr. Chen (Chief Doctor): Noah, you're saying the system is compliant with operational procedures because the logs are clean. But Aris is saying the system's confidence score conditioned her to trust the green flag. We must take this concern to our system providers: aren't they also liable for what happened to Mr. Miller?

Dr. Aris: Exactly. I am a doctor but I didn't design the algorithm, yet I'm responsible for the patient. If AI tells me a patient is low risk with high confidence, and I override it every time just to be safe, the triage system breaks down. Either the providers stand behind their green flags, or they've just handed me a very expensive way to fail my patients.

Technical explainability vs. justifiability

A tension emerged across the use cases between making an AI system's outputs technically understandable and making them genuinely justifiable to the relevant stakeholders. Measures such as confidence scores, visual explanations, provenance information, and model documentation should enable AI-supported decisions to be intelligible and open to critical scrutiny. However, understanding how a system works is not the same as being able to justify a concrete recommendation to a particular patient or employee. In UC1, clinicians should be able to integrate an output into a reasoned account of why a decision or course of action is appropriate for the specific individual patient, considering their values and preferences. This highlights a gap between model-level explainability and decision-level justifiability, and shows why transparency must also support informed consent, shared decision-making, and respect for autonomy, rather than being reduced to access to technical information. In UC3, the analogous tension was addressed through accessibility, with measures focused on ensuring employees could access clear, non-technical information about the system delivered through channels and formats adapted to different roles, languages, and digital literacy levels, operationalised through accessibility KPIs and structured GDPR-linked access processes.

Ethical considerations vs. operational and commercial pressures

The audit teams across the use cases were required to balance best practice with feasibility and proportionality, and measures had to be implementable given current systems and resources, and proportionate to the scale and risk of the deployment. A concrete example from UC3 was the component of over-reliance. From a commercial perspective, recommendations to add a human-in-the-loop compromise efficiency by adding friction to the decision making process. This tension was resolved by

taking a wider perspective on efficiency: over-reliance on a system to make high-impact HR decisions could lead to financial and reputational damage arising from the inability of the system to consider contextual factors and nuances that human-in-the-loop decision-making can encompass. In UC2, a parallel tension arose between the efficiency gains offered by AI-assisted safety checks and the ethical imperative to retain human judgement in safety-critical decisions.

Narrative Highlight #6 – Cognitive traps and proxy indicators in decision support systems:

James (AI developer): Vigilance.AI flagged Marcus with an 89% susceptibility score after he failed our phishing simulation. The behavioural data is indisputable. He is a security risk.

Maya (team lead): He isn't a risk. He is just exhausted. Marcus was finishing a double shift when the simulation happened. It makes no sense that the board now uses this to justify blocking his promotion.

Leo (AI governance committee): I must say I agree with Maya, since the system doesn't have a sensor for burnout, so it placed a "high risk" label on him. We've rebranded systemic fatigue as a psychometric trait and treating this as an objective truth.

James: Look, our mandate was to identify who clicks the link. We did that. If the cause is operational, that is not a technical failure of the system.

Maya: It is a failure when the system ignores the context, like Marcus' 14-hour login time, behind a proprietary score. By the time this reaches the board, the burnout is invisible. We aren't reducing risk. We are just automating the punishment of tired employees.

3. POLICY RECOMMENDATIONS

This section provides a set of recommendations for policymakers stemming from the research conducted in the context of this deliverable, following the methodology described in section 4, and in light of the open issues and governance gaps identified for each AI research area in section 2. AIOLIA puts forward the following policy recommendations (PR):

PR 01.

Provide guidance on the operationalisation of Human Oversight

Human oversight is a core element of the AI Act for high-risk AI systems (Article 14), mandating providers to design systems in a manner that enables human oversight and deployers to assign oversight to natural persons (Article 26). While human oversight has figured prominently across all AIOLIA use cases, they have demonstrated that this requirement is far more complex than it seems at first glance since the operationalisation of this practice varies significantly. For example, a DSS analysing radiological images autonomously requires limited human oversight as the system only deals with high-confidence cases. At the same time, human oversight in the context of personal AI companions can not only pose challenges to upholding the principle of privacy, as automated oversight mechanisms for content moderation are both more

effective in detecting prohibited practices, as they shield humans from being exposed to toxic content. Policymakers should provide additional guidance on the modalities and desirability of Human Oversight across diverse deployment contexts.

PR 02.

Facilitate the testing, review and evaluation of AI systems for small organisations

AIOLIA has demonstrated that the management of AI systems developed and deployed within organisational contexts requires several measures linked to the establishment of governance, review and monitoring committees, often requiring the involvement of a broad range stakeholders, ranging from specialised experts, such as linguists, to civil society organisations and users. Smaller organisation, with a more limited set of resources, might struggle to assess if their AI systems are ethically, legally and technically sound. Whereas Article 57 of the AI Act mandates member-states to establish regulatory sandboxes to support legal compliance, policymakers should expand the scope of planned regulatory sandboxes, or leverage existing digital innovation hubs, to establish broader testing, review and evaluation mechanisms, including for ethical, economic and societal dimensions of AI, to ensure small organisations can still flourish in what is currently a highly concentrated AI market.

PR 03.

Provide minimal explainability standards

The lack of explainability of complex AI systems is a major setback to their deployment, especially in high-risk domains like healthcare and the public sector. While the AI Act stipulates the right to explanation of individual decision-making (Article 86), determining what a well-engineered interpretability stack looks like remains an open question for organisations, especially for GPAI systems, which may provide the most accurate performance, but are not explainable at the level of a single decision or output. To mitigate this ambiguity and harness research efforts in line with explainability demands, policymakers should develop minimal explainability standards attending to what is technically feasible and socially desirable.

PR 04.

Revise classification rules for Emotional AI systems, especially for AI companions, and the definition of manipulation under the AI Act.

Current AI Act provisions leave important Emotional AI modalities and capabilities outside the high-risk categories, namely language-based emotion recognition that sits at the intersection of Emotional AI and GPAI – in the context of AI companions that are becoming the predominant Emotional AI application; and emotion elicitation, such as in deepfakes for therapy. At the same time, the intimate and recurrent nature of these interactions carries significant risks of emotional dependency and harmful influence, which are sharpened by business models where sustained engagement is directly monetised, creating commercial incentives that may deprioritise user safety in the absence of regulatory pressure. Further, the Act's intent-based framing of what constitutes manipulation that is a forbidden practice is not suited to harms that accumulate gradually through repeated interaction and may emerge as unintended

capabilities of training dynamics rather than deliberate design. Policymakers should consider revisiting the classification of Emotional AI systems under the Act, with particular attention to whether AI companions and AI for therapy warrant high-risk designation; and whether the current definition of manipulation adequately captures gradual and unintentional influence.

PR 05.**Establish a liability framework for high-risk AI systems**

AIOLIA industrial use cases, particularly those operating in high-risk domains, have voiced concerns over the lack of clarity with regards to AI liability, with provisions in the AI Act being open to interpretation. In addition to the absence of an AI liability framework, the lack of established EU standards for what a fully compliant high-risk medical AI device should look like, means that liability remains largely outside the control of AI providers and deployers. While AIOLIA's use cases provide measures to support organisations managing the ambiguity surrounding liability, they do not provide any broader legal certainty as to who is liable for faulty decisions involving AI-based decision-support. Policymakers must provide clarity on liability for high-risk AI systems, delineate the boundaries of fault, and fast-track sector-specific technical standards (e.g., for high-risk medical AI).

PR 06.**Facilitate the creation of sector-specific AI standards**

The operationalisation of AI ethics principles conducted by AIOLIA partners revealed that partners operating in more standardised sectors (e.g., automotive sector), had a clearer pathway for acting upon AI-related risks and ethical concerns. The lack of clear standards resulted in modelling operationalisation practices on market practices to secure competitiveness, which are not always ethically and socially sound, as is the case of personal AI companions. Policymakers should engage in the development and provision of sector-specific ethical and safety baseline standards for consumer-facing AI, establishing a level playing field that does not penalise responsible AI design and deployment.

PR 07.**Facilitate dialogue on societal-level impacts of AI**

One key insight arising from AIOLIA is the fact that organisations are, by themselves, unable to deliberate with certainty on higher-level discussions on societal well-being. While organisations recognise that AI systems might pose broader challenges that exceed the immediate contexts of deployment, such as the normalisation of deepfakes in society and the risk of eroding social trust, they do not feel capable of tackling systemic issues alone, for they lack the mandate, capability, and democratic legitimacy to address them in isolation. Policymakers should establish permanent, multi-stakeholder societal impact fora for extra-organisational dialogue and discussion on a desirable way forward for affected communities. These platforms should bring together AI developers and experts, civil society organisations, and citizens to deliberate on systemic AI externalities.

PR 08.**Invest in public AI literacy campaigns to cover personal and affective uses of AI**

Article 4 of the AI Act requires providers and deployers to ensure adequate AI literacy among staff and those involved in the operation of AI systems, but this obligation is confined to organisational contexts and does not extend to the general public. This gap is particularly consequential in the personal sphere, where individuals interact with AI systems outside any professional framework and without institutional support for critical engagement. Unlike workplace deployments, personal use occurs without oversight, without trained intermediaries, and often without awareness of the nature or implications of the interaction. The European Commission's repository of AI literacy practices, launched by the AI Office in 2025, represents a valuable step toward supporting learning and exchange on AI literacy, but its current focus on providers, deployers, and their staff leaves personal use largely unaddressed. Policymakers should invest in broad public AI literacy campaigns that specifically cover the personal and affective dimensions of AI use, equipping individuals to recognise AI involvement in interactions, understand the nature and risks of emotional engagement with AI systems, and make genuinely informed choices, in partnership with civil society and mental health professionals.

PR 09.**Promote and incentivise research on AI impacts**

As AI uptake and diffusion expand across sectors, more research is required to document and understand the risks and impacts of AI systems, including on human behaviour and cognition, in the short-, medium- and long-term, and among different deployment contexts, including in private and professional spheres. Core areas of research include explainability, privacy-preserving AI, human-AI interaction, including behavioural influence, and environmental impact. Independent researchers possess the expertise to study these socio-technical impacts of AI but lack access to proprietary data and live operational systems. Policymakers should incentivise research to be conducted at the heart of organisations, providing access to independent researchers, including through internal usage of the systems and data, whilst providing protection over data leakage, liability, IP theft and reputational damage. Access can be modelled up Article 40 of the DSA.

PR 10.**Update ALTAI**

ALTAI remains a valuable assessment for organisations developing and deploying AI systems. However, six years following its presentation by the AI HLEG, the assessment list is currently outdated and several elements related to the last generation of AI systems, namely GPAI, are currently missing. Section 5 of the current deliverable provides insights on several ALTAI gaps. Given the lasting centrality of ALTAI in the EU, policymakers should update the framework to address the latest challenges posed by AI systems and consider developing ALTAI for specific high-risk AI systems based on Annex III of the AI Act.

4. METHODOLOGY

4.1. Interaction between tasks in AIOLIA WP3

AIOLIA's Work Package 3 comprises four interconnected tasks designed to operationalise AI ethics principles across multiple contexts and scales. Task 3.3, which develops context-enriched operational guidelines for AI research areas, sits at the centre of this architecture, drawing on technical use case analysis (T3.1), international comparative insights (T3.2), and providing the foundation for ethics assessment mechanisms (T3.4).

4.1.1. Co-creation process for technical guidelines in T3.1

Task 3.1 establishes the empirical foundation for T3.3 through a structured co-creation process involving six industrial use cases. This process employs a two-phase methodology: decomposition of high-level ethical principles into actionable components, followed by consolidation across use cases to identify common patterns and challenges.

The decomposition phase requires use case partners to translate abstract principles into context-specific components and practical implementation measures. This two-stage process first identifies the constituent components of each ethical principle, then maps these to measurable features within industrial workflows. The collaborative design process involves continuous communication between academic and industry partners, with industry typically producing initial proposals grounded in operational experience, which academics then review and refine through iterative cycles.

T3.3 receives from T3.1 a set of validated, practice-grounded interpretations of ethical principles and ready-to-be-applied organisational measures across diverse AI applications ([see D3.1](#)). When T3.3 identifies "overarching ethical issues and organisational measures" across use cases, it draws directly on this consolidated evidence base to determine which ethical concerns are fundamental to specific AI research areas versus which are use-case-specific and which organisational measures are applicable not simply to the specific context of the use cases, but capable of being extrapolated to any AI technology in the research area.

4.1.2. International interaction and learning from T3.2

Task 3.2 extends the co-creation methodology to international partners in Canada, China, South Korea, and Japan. Each partner conducts a parallel process of identifying and operationalising ethical principles relevant to selected AI research areas, adapted to their national regulatory frameworks and cultural contexts.

The relationship between T3.2 and T3.3 is fundamentally comparative rather than derivative. T3.2 does not produce content that T3.3 directly incorporates, instead, it provides critical perspective on which aspects of the European guidelines represent universal technical requirements versus culturally or

regulatorily specific choices. When similar ethical concerns emerge independently across different regulatory frameworks, this convergence indicates fundamental technical or ethical requirements rather than arbitrary policy preferences. The European guidelines developed in T3.3 are more robust precisely because they have been developed with awareness of how alternative regulatory approaches handle the same AI research areas.

4.1.3. Ethics readiness levels in T3.4

Task 3.4 develops an Ethics Readiness Levels (ERL) evaluation mechanism that translates T3.3's qualitative guidelines into quantitative assessment criteria. This relationship is explicitly integrative: the ERL tool is added to the set of non-technical guidelines, positioning T3.4 as an implementation component of T3.3 rather than a separate output.

The practical relationship is straightforward: T3.3 contributes with indicators to the AIA sub-tool developed in the context of T3.4, which translates this into checkable criteria and scoring mechanisms. The modular, Excel- and Web-based format of the ERL tool reflects the structure of T3.3 guidelines, with different modules corresponding to different ethical principles or research areas.

4.2. Methodological approach

Task 3.3 develops guidelines on the operationalisation of AI ethics principles in a non-technical manner for AIOLIA research areas. This process is based on a fourfold methodological approach, further detailed in the remainder of this section, pursued as follows:

- 1) Situating AI ethics within legal, societal and economic considerations,
- 2) Analysing empirical data from AIOLIA use cases,
- 3) Identifying and narrativising tensions between ethics principles,
- 4) Extrapolating organisational measures for AI research.

4.2.1. The ELSE framework: situating AI ethics within legal, societal and economic considerations

The non-technical approach to AI ethics pursued in T3.3 situates ethical considerations using a wider aperture, such as on aspects of inclusivity, societal benefits and harms, and the relevance and applicability of AI systems to different domains. The non-technical approach adopted herein dovetails with an ELSE framework, as Ethical considerations are considered alongside Legal, Societal and Economic aspects that take into account the impact of AI upon individuals and larger communities.

ELSE is a framework for integrating s Ethical, Legal, Social and Economic concerns within scientific research and innovation. It focuses on anticipating potential risks emerging from scientific research and technology development, proactively putting forward approaches to mitigate them. The ELSE framework first emerged in the context of genomic research, namely the Human Genome Project, launched in 1990. Since mapping human genome would centralise genetic information of millions, rendering those with genetic

disorders easily locatable and targetable and given the living memory of eugenics as an instrument of mass violence, the United States formed a Joint Working Group on Ethical, Legal, and Social Issues (ELSI) in 1990. ELSI became a pillar of scientific research in both the US and EU, where the European Commission formally adopted the name ELSA – Ethical, Legal, and Social Aspect – in 1994 (Braun and Muller, 2025). This institutional positioning enabled ELSA to meaningfully influence research and innovation in line with public concerns and safety.

As technology evolved, the ELSA approach was incorporated into different areas of enquiry. In 2002, the European Commission established an ELSA board for ethical inquiry into nanobiotechnology, expanding to include private sector perspectives (Braun and Muller, 2025). More recently, attention has turned to AI, with scholars arguing that the risk emerging from AI technologies result from "sociotechnical failures" which require interdisciplinary research and responses that integrate ELS considerations into design processes (ibid). Several scholars have put forward adaptations to ELSE to address the specificities of AI technologies and ensuring that innovation is ethically and socially responsible. Notable initiatives include the ELSA Labs for AI funded by the Dutch Research Council (Wang et al., 2025; van Hitlen et al., 2025).

Whereas, in the domain of ELSE research, intense discussion and variability remains with regard to the scope of this framework, with some opting for ELS-I (with 'I' standing for 'implications') and others ELS-A (with 'A' standing for 'aspects'), and the integration of economic dimensions being recognised as essential in the context of AI (Ryan and Blok, 2025), the foundational matrix of this framework remains relevant in the context of this deliverable, irrespective of the tensions that may exist within its own research field (for a detailed account of current issues see Ryan and Blok, 2023).

Why ELSE for AI?

In the context of AI, the ethical design and deployment of these systems has become a major concern. AI ethics has become not simply a field of enquiry, but a cornerstone of responsible AI and the basis upon which multilateral consensus and a plethora of regulatory frameworks for AI have emerged. The nature of AI technologies, spanning technical and human elements, from raw materials, data and algorithms, to data workers, subjects and impacted populations, makes the ELSE approach even more relevant for broadening the scope of ethics into legal, societal and economic dimensions, so pressing in the current context of AI development. As AIOLIA's UCs demonstrate, AI ethics is, in fact, inalienable from legal, societal and economic considerations. For example, AIOLIA UCs demonstrate that compliance with legal frameworks – which remain, nonetheless sparse and ambiguous, particularly as the standardisation of the AI Act is still ongoing – is insufficient for addressing the broader range of concerns arising from AI systems, including the impact on skills and labour, the still uncertain consequences on individual and societal well-being, and the landscape of economic competitiveness and attractiveness shaping AI design and deployment, which may not always favour responsible AI practices.⁹

Rather than prescribing action, ethics is a mode of inquiry, a means of generating questions, rather than fixed rules, allowing the study of periods of significant societal change and analysis of emerging risks (Rességuier and Rodrigues, 2020). As such, ethics generates varying insights depending on context and structural conditions, with ethical meaning differing across geography, race, and socio-economic

⁹ For a more extensive appraisal of these issues see Section 5.

structures. It follows that drawing on a single normative framework risk marginalising those whose experiences generate differentiated ethical considerations (Mager et al., 2025). This points to a fundamental insight: AI ethics is shaped by the structural conditions within which AI systems are developed and used. Bias, for instance, does not exist separately from the society that created the AI system, including the provenance of data and the populations and datapoints included and excluded therein – themselves shaped by societal and economic factors. Removing bias is not a technical 'debugging' exercise but requires mitigating biases emerging not only from the sociocultural values and practices that shape social life, but also from the broader range of stakeholders and organisational practices involved in data- and AI-making. There is no standardised procedure for operationalising AI ethics because AI-society relationships take a myriad of forms, and ethical issues vary accordingly (Munn, 2023).

For this reason, the operationalisation of AI ethics is a key challenge, faced by organisations and policymakers alike, and one that AIOLIA takes as its main focus. Operationalising AI ethics within organisations requires visibility into the ethical issues that arise across multiple levels within the organisation, and how they interact both with one another and with the wider context of AI development and deployment. This is where the ELSE framework proves particularly relevant in the context of AI, for shedding light on the multiple variables and levels at play in ethics-by-design approaches to AI development and deployment.

Wang and Block, for example, call for a multi-level framework that dynamically navigates through individual, organisational, sociopolitical, and ontological issues while understanding their intertwinement (Wang and Block, 2025). This supports in the mapping and forecasting of impacts of AI across different scales of analysis. The authors provide a series of questions targeted towards specific issues at varying levels for AI developers to address spanning product-level development to organisational cultures and broader societal structures, encouraging reflection on how these levels mutually reinforce one another (ibid). A discriminatory AI system, for instance, may reflect not just flawed training data (individual/technical level), but also organisational incentives that prioritise efficiency over fairness (organisational level), and deeper societal patterns of discrimination that shaped the historical data in the first place (sociopolitical level).

Ryan et al. (2025) follow a similar approach by advocating for multilevel analysis across micro, meso, and macro scales. The micro level addresses the specific technical and design features of AI systems; the meso level examines organisational practices, governance structures, and industry norms; the macro level attends to broader societal, political, and economic forces shaping AI development and deployment. This stratified approach recognises that effective AI governance cannot focus solely on product-level interventions but must engage with the wider organisational environments in which these products are developed and deployed.

Whereas ELSE assessment is commonly performed with regard to domain-specific use-cases, e.g., agriculture (van Hilten et al., 2025), our approach seeks to systematise ELSE at the level of AI research areas by incorporating comparative and cross-sectoral analysis to lift overarching ethical concerns and their inherent social, legal and economic dimensions. While scholars and ethicists have analysed the novel ethical challenges emerging from GPAI, specifically, generative AI applications, as well as personal

companions and DSS, they often lack the contextual insights of specific use-cases. At the same time, the use-case centric approach adopted in ELSE research is informed by concrete sectoral-level practices and use-cases, but does not adopt a systematic approach to research areas/technics. The current AIOLIA deliverable therefore provides a novel approach, building on the ethics principles selected for the three AIOLIA research areas and their different use-cases to analyse how ethical principles have been operationalised by AIOLIA partners. The variety of UCs, derived from areas as diverse as healthcare, personal companions, engineering and security, allows us to cross-analyse the bottom-up operationalisation of ethics principles and to lift overarching ethical concerns whilst acknowledging tensions. In fact, the aim is not simply to produce context-enriched non-technical guidance that renders ethics approaches uniform within AIOLIA research areas, but to acknowledge heterogeneity and tensions between UCs' operationalisation process, and to identify unique sectoral needs. These will be incorporated as learning points, providing insight on how practices in one context (e.g., healthcare) can inform approaches in others (e.g., personal companions).

Towards ELSE-aware operational guidelines

Whereas the operationalisation of AI ethics and the translation of abstract ethics principles into practice have been recognised as a key challenge, several scholars have voiced concerns over the development of a fixed set of principles and practices for 'achieving' AI ethics. The core concern motivating such appraisal, is the fact that ethics can be said to be contrary to the formulation of rigid guidelines, for the former may denature the iterative questioning and reflection that comprises ethics, whereas ethics must remain "agile" to effectively address AI risks (Rességuier and Rodrigues, 2020). At the same time, however, supporting organisations and policymakers in the responsible design and deployment of AI technologies is fundamental to ensure that technological evolution is tantamount to individual and societal well-being. While this does not require a prescriptive approach that renders AI ethics a tick-box exercise, facilitating knowledge sharing between organisations and providing examples of practices, measures and alternatives for ethical AI design and deployment remains paramount. This requires providing tools and concrete examples that enable organisations to assess the risks that arise at the AI-design level alongside the organisational structure within which AI systems are deployed, without neglecting the broader societal impacts and specific context of AI deployment in the short and long term.

The guidelines developed herein seek to contribute to a growing body of evidence-based operationalisation of AI ethics, paying attention to how ethics intersect with legal, societal and economic dimensions. They do not aim to shortcut the Socratic process of deliberating about AI ethics (Wilders et al., 2025), nor to present ethics as a resolved issue, but rather seek to inspire ethical practices within organisations whilst raising awareness of the tensions, challenges and trade-offs that arise in operationalising AI ethics within the specific contexts of AI design and deployment. The variety of AIOLIA UCs, their contexts and types of AI systems, provide a valuable basis upon which to ground this organisational guidance in a manner that acknowledges the intersection of ethical, legal, societal and economic dimensions making up AI systems. To this end, an ELSE literature review has been conducted for each AI research area identifying the core ethical, legal, societal and economic concerns specific to each type of AI system deployed across AIOLIA use cases, spanning GPAI, Emotional AI and DSS (See Section 5). This step in the research process not only fostered a deeper understanding of this task's objects

of study, as it paved the way for identifying both recurrent and unique concerns arising from AIOLIA use cases.

4.2.2. A bottom-up approach: analysing empirical data from AIOLIA use cases

This task adopts a bottom-up ELSE approach in developing operational guidelines for AI research areas. Following the review of ELSE literature conducted for each AIOLIA RA, T3.3 pursued with the analysis and incorporation of the bottom-up data provided by AIOLIA's use cases in the context of the operationalisation of ethics principles conducted in T3.1 and informed by the Operationalisation Pathway defined in D2.3. This data included Cycle 1 activities, that is, the decomposition of ethics principles into specific *components* and actionable *measures*, and Cycle 2 activities, specifically the validation of UC data through attendance at validation meetings conducted by CENTRIC in the context of 3.1.¹⁰

The bottom-up ELSE approach adopted herein can be justified using two different lenses. First, while ethical principles and guidelines exist in abundance to guide those designing and developing AI – we call this the policy-to-practice pipeline – feedback mechanisms to policymakers and the public by those who are attempting to operationalise these guidelines do not currently exist. In other words, while the policy-to-practice pipeline is well fleshed out and well-traversed, the practice-to-policy pipeline remains under-explored. This in effect means that there is currently no feedback mechanism and no structured way for a policymaker or the public to know whether those implementing the ethical principles are facing challenges in so doing, whether the principles are feasible, whether gaps and trade-offs between principles exist, or to clarify ethical responsibilities or support structures that can enable the operationalisation of ethical AI. Undertaking this exercise enables us to examine unexplored opportunities or potentially engage other levers to operationalise ethics in AI.

Secondly, the bottom-up approach is informative in the identification and lifting of the overarching ethical concerns and organisational measures that emerge from cross-analysing the individual use cases within AIOLIA research areas. This enables the co-creation of context-enriched non-technical AI ethics guidelines for AIOLIA research areas. The bottom-up approach provides the context that informs how to best identify, address and work with ethical issues within research areas, as informed and enriched by how those research areas are deployed in within the specific use cases in AIOLIA.

ALTAI principles within ELSE

The focus of task 3.3 is on the research areas adopted by the European partners in their use cases (see section 1.1.). The European focus means that ethical guidelines and other recommendations developed within the European AI governance context is a natural starting point, even if it is the case that many ethical guidelines on AI converge around the same groups of values and principles. A key European instrument for this purpose is the Assessment List for Trustworthy AI (ALTAI), developed by the High-Level

¹⁰ See D3.1, section 3.2., for more information on the data collection process: <https://aiolia.eu/wp-content/uploads/2026/02/AIOLIA-D3.1-certified.pdf>

Expert Group (HLEG) in 2020. ALTAI supports the actionability of the seven key ethical requirements for Trustworthy AI, with each having several dimensions:

1. Human agency and oversight

- Human agency and autonomy addresses how the AI system affects human behaviour, generates confusion about whether an interaction is with a human or a machine, or creates over-reliance, over-attachment or manipulation that might interfere with user decision-making.
- Human oversight encompasses governance mechanisms for system oversight, such as human-in-the-loop (intervention in every decision), human-on-the-loop (intervention during design and monitoring), or human-in-command (overall activity oversight), as well as adequate training and system shut down procedures.

2. Technical robustness and safety

- System resilience and security concerns vulnerabilities like data poisoning or model evasion.
- General system safety covers the definition, assessment and treatment of risks, and system reliability.
- System accuracy in the deployment context ensures that the AI system performs correctly and consistently across its intended use cases.
- Reliability, fall-back plans and reproducibility focus on verification and validation methods, and the implementation of fall-back plans to prevent unintentional harm.

3. Privacy and data governance

- Privacy deals with AI's impact on privacy and data protection – two fundamental rights, which are closely related to each other and to the fundamental right to the integrity of the person, which covers the respect for a person's mental and physical integrity.
- Data governance covers training data quality, its relevance in light of the domain in which the AI systems will be deployed, its access protocols and the capability to process data in a manner that protects privacy.

4. Transparency

- System traceability includes documenting and logging the data and rules of AI outputs.
- Explainability ensures that the technical reasoning of the AI system is understandable to those affected.
- Communication with users focuses on informing users about the AI's capabilities, limitations, accuracy levels, and potential risks.

5. Diversity, non-discrimination and fairness

- Avoidance of unfair bias examines the implementation of measures to identify and remove discriminatory bias in data collection and algorithm design.
 - Accessibility and universal design ensures the AI system is usable to all through Inclusive Design principles and assesses and mitigates the risk of exclusion and unfairness towards particular groups.
 - Diverse stakeholder participation requires the inclusion of the widest range of possible stakeholders in the AI system’s design and development.
6. Environmental and societal well-being
- Environmental wellbeing assesses the impact of AI systems on the environment.
 - Impact on work and skills concerns the assessment and mitigation of the risk of deskilling and training staff on how to responsibly use the AI system;
 - Impact on society at large and democracy encompasses the minimisation of the potential harms arising from the AI system and the implementation of measures to ensure it does not negatively impact democracy.
7. Accountability
- Auditability covers the mechanisms that facilitate evaluation of the AI system by internal and external auditors.
 - Risk management addresses the identification, assessment, documentation and minimisation of potential negative impacts of AI systems.

The translation of these seven high level principles into an actionable list of measures and questions through ALTAI enables the operationalisation of ethics by those developing and deploying AI. However, ALTAI was developed in 2020, before the advent of generative AI and is also one tool out of many that aim to operationalise AI ethics. As such, focusing exclusively on ALTAI to navigate the operationalisation of ethical AI is insufficient. Two key reasons motivate situating ALTAI within the larger ELSE framework.

First, gaps can arise on the understanding of how ethical design can look like for these different types – predictive and generative (Narayanan & Kapoor, 2024) - of AI systems. The ALTAI checklist was designed when the paradigm of predictive AI was dominant – namely systems designed to predict or support outcomes, decisions and behaviour. These systems were also typically designed and deployed for specific purposes within discrete domains, such as in radiology, predictive policing, fraud detection in tax and welfare benefits or within recruitment. This can be contrasted with generative AI – namely AI used to generate audio, text, video content or multi-modal content generation. Generative AI has raised a plethora of ethical concerns, including but not limited to the generation of problematic content, hallucinations, lack of predictability and bias amongst others. However, it also raises legal issues concerning liability for harms, the lack of clarity on lawful forms of data processing, and other reporting and accountability measures. The increased uptake of generative AI in society raises profound ethical, legal, socio-economic and other societal concerns. While ALTAI is a useful starting point for addressing

these issues, the ELSE framework enables us to take a step back and in effect, stress test ALTAI itself, in light of new technological developments and the challenges it presents.

Secondly, while industry and academic partners in T3.1 take ALTAI as a starting point, the use cases also engage with sectoral legislation and standards, such as safety standards applicable to the automotive sector or the medical devices legislation. In other words, there is an interplay of ethical and legal requirements when it comes to the design and deployment of AI thus requiring clarity from both ends. The ELSE lens assumes an interplay of these factors – not just between ethics and the law, but also how socio-economic and societal level considerations can in fact inform and influence the design of better AI systems that both serves its intended purpose and which minimises societal harms. The ELSE framework also enables the acknowledgement of the undercurrents of the push for AI adoption, driven by the perceived socio-economic advantage for adopters. However, a blind push for AI adoption in organisations can lead to internal impacts such as labour deskilling or external facing impacts, if organisational measures such as whistleblowing mechanisms, AI policies or organisational safeguards are non-existent.

Situating ALTAI within the ELSE framework provides for a more holistic lens in which to translate ethical AI into practice and allows us to identify gaps in ALTAI. The three research areas in Section 5 – general-purpose AI, emotional AI and algorithmic decision-support systems – will be analysed using a bottom-up ELSE approach, to inform the selection and formulation of non-technical operational guidelines for AI research areas in D3.3.

4.2.3. Identifying and narrativising ethics tensions

As the operationalisation of AI ethics principles conducted by AIOLIA partners in T3.1 demonstrated, ethical principles are not self-contained but rather constantly interact and intersect with one another. In putting principles into practice and identifying practical measures for operationalising ethical principles, AIOLIA partners identified different tensions between ethics principles – that is, the operationalisation of one principle compromising another – and were required to make trade-offs. The tensions identified by AIOLIA use cases have been provided by partners in written format during the completion of T3.1 and also emerged in the context of discussion held during validation meetings taking place for the completion of that same task. The ethics tensions identified are outlined in sections 2.1.4., 2.2.4., and 2.3.4., for each AI research area.

As ethics tensions provide important learning points on the challenges involved in the operationalisation of AI ethics, and in light of AIOLIA's pedagogical ambition, the tensions identified have been narrativised through short narrative highlights (added in coloured boxes in the aforementioned sections). It formulates these highlights in a non-technical way, in line with the methodology of Task 3.3, and aims to facilitate the development of pedagogical materials in AIOLIA's WP4, which prepares AI ethics trainings later in the project. The narrative highlights have been developed based on the following methodology:

1. Select tensions identified by partners in the operationalisation process conducted in T3.1 on six European use cases.
2. Develop fictional narratives based on selected tensions to illustrate the tensions in a simple and pedagogical form.

3. Validate narrative highlights with AIOLIA partners involved in the respective use cases.

4.2.4. From use cases to research area: extrapolating salient ethical principles and measures

Based on the analysis of operationalisation data, T3.3 identifies and lifts the overarching ethical concerns, challenges, gaps and organisational measures present across the RA UCs. This process is pursued through several steps, namely by clustering the research areas within the use cases, identifying commonalities in ethical challenges and tensions and drawing from the organisational measures put in place by the different use cases that can be generalised and extrapolated to the research area level. The organisational measures at the basis of T3.3 have been formulated by partners in the context of T3.1 – see [D3.1](#), Appendixes A and D – and have been further developed in the context of AIOLIA’s validation workshops and ongoing analytical work through a live portfolio of measures (D4.2, forthcoming).

Grounding this process is Table 12 (see [D2.3](#), page 40), which provides a template for examining convergence and divergences of ethical principles for AI research areas. Most notably, this includes analysing the interpretation and scope of the ethics principles in the context of each UC, including the components and measures identified for operationalisation. Building on this UC-specific analysis, convergences in terms of definitions, methods and mitigation strategies are identified and assessed for RA relevance. The same occurs with regards to divergences across RA UCs, including in the understanding of the principles, tensions and trade-offs. This process is extensively described for each RA in Section 5.

Based on this analysis, T3.3 extrapolates a limited set of RA-specific ethical principles and measures, taking into account the central role of UC measures in operationalising AI ethics principles and attending to the sheer number of measures identified. Rather than reproducing the wide variety of existing guidelines, including ALTAI, the aim of T3.3 is not to be comprehensive with regards to ethics principles and operationalisation measures potentially relevant for each research area, but to single out those which are the most pressing in light of their given specificities, identified based on ELSE literature and empirical AIOLIA UC data.¹¹ The full range of principles, components and measures identified can still be found in [D3.1](#).

With regards to the selection of measures stemming from the range of organisational measures identified across European UCs, T3.3 designed multi-level criteria for inclusion in RA organisational guidelines, based on relevance, potential for extrapolation and/or adaptation and originality. *Relevance* assesses whether the measure applies to the research area as a whole, not only the use-case, by analysing whether it encompasses concerns outside the specific application context. The criterion potential *for extrapolation and/or adaptation* accounts for the potential for a given UC measure to be applicable to the RA as a whole, and, where necessary, requires making linguistic and scope adaptations. The criterion *originality* investigates whether the measure addresses emerging challenges that have not been stably captured in

¹¹ For the broad range of ethics principles, component and related technical and organisational measures to operationalise these, see AIOLIA D3.1 at <https://aiolia.eu/wp-content/uploads/2026/02/AIOLIA-D3.1-certified.pdf>

best practices, standards and/or regulation. This multi-level criteria draws on the ELSE literature and salient ethical principles and concerns identified for each RA in light of the analysis conducted.

In practical terms, this multi-level criteria translated into the extrapolation of organisational measures in different manners. For example, in the context of GPAI, UC2 has put forward a unique organisational measure across all GPAI UCs, namely, the establishment of a multidisciplinary AI governance committee as a formal oversight structure within the organisation, to protect against organisational pressures. This measure fitted into the ‘relevance’ and ‘originality’ criteria and addressed a pressing concern related to GPAI uptake, being, for this reason, lifted as a research area-level measure. The context of Emotional AI surfaces several such examples because it was represented by two use-cases only, and there was almost no overlap in organisational measures, reflecting the diverging deployment contexts. Yet, we found several measures that address important concerns across settings, such as the measure on independent ethics reviews of the system's impacts, which was put forward by UC6, with UC5 having *internal* ethics reviews. Independent reviews, including by auditors, could mitigate commercial pressures that might otherwise cause organisations to disregard wellbeing.

Co-creation process in operationalising ethical principles

To ensure that the resulting operational guidelines per research area are effective, we adopt a co-creation process in the form of stakeholder engagement. Firstly, this involves receiving internal feedback from our project partners, through validation workshops, ensuring our findings resonate with their UCs and RAs and are fit for AI ethics trainings. To this end, three validation workshops were organised by CEPS for each of the three RAs, as follows:

- Emotional AI workshop, occurred on 28 April 2026, with the participation of AUMC, CEA, CEPS, EUREC and THWS.
- DSS workshop, occurred on 29 April 2026, with the participation of Afliant, CASTED, CEA, CENTRIC, CEPS, EUREC, ETICAS, NIT and Oxipit.
- GPAI workshop, occurred on 4 May 2026, with the participation of AUMC, CEA, CEPS, EUREC, NIT and THWS.

AIOLIA’s UC partners provided structured input to a draft version of the sections “Salient Ethical Concerns” and “Organisational Measures”, under section 2 of this deliverable, attending to 1) the overall document, 2) the selection of salient ethics principles and 3) the selection of measures. Partners also added comments directed to the draft guidelines documents.

Secondly, co-creation takes the form of feedback from AIOLIA’s Stakeholder Advisory Board (SAB) – a group of leading voices in AI ethics, governance, law, research integrity, and European policy. The SAB provided input to this task in two different occasions: first, during the AIOLIA Consortium meeting held in March 2026 (Berlin), and subsequently in May 2026, upon completion of the final draft of the organisational guidelines. The feedback received from the SAB has not only been incorporated in the formulation of the current guidelines (see Section 2), as it will inform subsequent AIOLIA tasks, most notably, the policy briefs to be elaborated in the context of T6.3.

A standalone operational guidance

To make AIOLIA’s work more impactful, the core findings of T3.3, elaborated in the form of RA organisational guidelines, has been formulated in shorter, standalone documents, for broader and more accessible distribution. These shorter guidelines can be found in Section 2. The short guidelines are targeted towards the following non-technical audience:

- Organisations that deploy AI systems;
- Organisations that provide AI models or systems;
- Policymakers working in the field of Artificial Intelligence.

The formulation of organisational guidelines in Section 2, follows a different approach than the one adopted thus far in AIOLIA WP3: rather than organising operational measures per AI ethics principles, T3.3 opted for adopting the ISO/IEC42001 as its organising thread. The ISO/IEC42001 is an international standard that seeks to facilitate the establishment, implementation and maintenance of AI management systems within organisations. AI standards are gaining an increasingly central role in AI governance, turning consensus¹² on the ethical principles that should underline AI development and deployment into operational guidance, against the backdrop of the slow process of encoding them into law. In the European context, *harmonised technical* standards play a key role as proof for legal compliance under the AI Act, with the standardisation process for requirements targeting high-risk AI systems currently ongoing. Next to the homegrown approach to standards (where international standards can still be used as the basis or portions of the European standard), the EU can adopt existing international standards¹³, with the ISO/IEC42001 standard having been approved for full adoption in 2026.¹⁴ The alignment of the RA-specific organisational guidelines with international standards attests to AIOLIA’s commitment to actively contribute to the current landscape of AI policy and practice.

The ISO/IEC42001 standard for implementing an AI management system within organisations is structured around seven dimensions:

- 1) “Context” outlines measures to address the specific context of the organisation, including internal and external contexts, ranging from legal requirements to competitive landscape, as well as cultural and social values. Understanding the organisational context also requires accounting for all the parties impacted by the AI system, from direct users to wider communities.
- 2) “Leadership” outlines measures to address the leadership commitment of the organisation towards the AI management system, including through defining a policy for AI systems deployed within the organisation, but also assigning clear roles and responsibility for managing these.

¹² Examples of initiatives laying a common ground on the ethics principles grounding the development and deployment of AI include the OECD AI Principles (2019), the UNESCO Recommendation on the Ethics of Artificial Intelligence (2021) and the Council of Europe Framework Convention on Artificial Intelligence and Human Rights, Democracy and the Rule of Law (2024).

¹³ https://single-market-economy.ec.europa.eu/single-market/goods/european-standards/standardisation-policy/high-level-forum-european-standardisation_en

¹⁴ https://www.etuc.org/sites/default/files/page/file/2026-04/AI%20standardisation%20Inclusiveness_Newsletter14.pdf

- 3) “Planning” outlines measures to address the planning phase of developing and / or deploying an AI system within the organisation. This involves identifying and assessing risks, potential impacts and opportunities arising from the AI system and setting clear objectives for its deployment.
- 4) “Support” outlines measures to support the establishment, implementation and maintenance of the AI management system within the organisation. It involves determining and providing organisational resources, ranging from competencies, awareness and clear communication and documentation.
- 5) “Operation” outlines measures to support the operation of the AI system within the organisation. It involves operational planning and control, and conducting risk and impact assessment, planning effective means for corrective action.
- 6) “Performance evaluation” outlines measures to support the evaluation of the AI system’s performance, namely through system monitoring and analysis, as well as internal auditing.
- 7) “Improvement” outlines measures for organisations to continuously improve the operation of the AI system, including determining when performance is not in conformity with expected requirements and intended use, and adopting corrective actions.

For each dimension, concrete organisational measures (OM) are provided: each entry specifies what the measure consists of; why it is important for the specific research-area context, including how it addresses the area's salient ethical principles; and guidance about its implementation, including implementation challenges and obstacles and at what stages of the AI lifecycle measures can be implemented. The AI lifecycle follows the OECD categories: 1) design, data and modelling, 2) verification and validation, 3) deployment, and 4) operation and monitoring.¹⁵

5. OPERATIONALISING ETHICS PRINCIPLES

5.1. GENERAL-PURPOSE AI

General-purpose AI (GPAI) is a category of AI models that display significant generality, are capable of competently performing a wide range of distinct tasks and can be integrated into a variety of downstream systems or applications (Art 3(63) AI Act). Expanding on the AI Act’s definition of GPAI, Triguero et al. highlight how these advanced AI systems exhibit a significant degree of autonomy and ability to generalise to new tasks and across domains the system has not been previously exposed to or intentionally trained to address (Triguero et al., 2024). This is due both to model architectures and the large volumes of data (often multimodal data), used to train GPAI models. While the term GPAI is privileged in the EU regulatory framework for AI, the terms foundation model, generative AI and, more recently, agentic AI, have taken prevalence over GPAI, given the commercialisation and public outreach of these types of GPAI systems, which have been at the frontier of AI developments for the past four years (Triguero et al., 2024). In reality,

¹⁵ For more information on the phases of the OECD AI Lifecycle see https://www.oecd.org/content/dam/oecd/en/publications/reports/2019/11/scoping-the-oecd-ai-principles_71e1b6dc/d62f618a-en.pdf

however, the category of GPAI exceeds these three types and ongoing scholarly and legal debate continues regarding the best way to define these systems, which neither rely on a single AI architecture (e.g., transformer, variational autoencoders, diffusion models, etc.) and learning approach (e.g., self-supervised, unsupervised learning), nor serve a specific set of purposes.

It follows that, as a research area, GPAI presents unique ethical challenges. On the one hand, the generality of tasks GPAI models can perform, both in professional and private context, means that there is no single, standardised range of ethical principles and measures that would apply to this research area. On the other hand, the fact that GPAI models, as foundation models, often form the baseline architecture upon which different applications are built, poses new challenges to the operationalisation of AI ethics principles in the context of model deployment. To be precise, GPAI models may not only be deployed or integrated in contexts that were not originally foreseen during model development, as their underlying generative architectures imply that outputs are not fully predictable, nor (arguably) containable, as attested by state-of-the-art LLMs and their so-called hallucinations. All in all, these elements restricting the extent to which downstream providers, i.e., those who incorporate foundation models as basis for the GPAI systems, can guarantee the ethical alignment of their systems and engage in ethics-by-design practices.

However, as the range of AIOLIA use-cases at the heart of our analysis demonstrates, ethical issues and challenges arising in the GPAI research area far exceed the technical specificities of these systems; rather, they sit at the intersection of AI systems and human behaviour and cognition. Combining bottom-up data on the operationalisation of AI ethics principles provided by AIOLIA use cases with insights from ELSE literature, the remainder of this section delves into the unique ethics principles, component and measures of the GPAI research area.

5.1.1. The ELSE approach to GPAI

The ELSE approach to GPAI takes into account the specificities of these models to assess the ethical, legal, social and economic aspects of GPAI as an AIOLIA research area. This is aligned with AIOLIA's ambition to translate ethics principles into concrete operational guidelines, insofar as it foregrounds our understanding of ethics not as an isolated conceptual abstraction, but as contextual and nuanced practice wherein ethical principles cut across societal, economic and legal dimensions. For example, understanding the ethical principle of human agency and oversight not simply as the legal requirement of human supervision for high-risk systems (AI Act), but as having social and economic dimensions, enables organisations to put in place continuous training to avoid staff deskilling (social) and measures for mitigating the exploitation of users' emotions (economic), as 'enablers' of the ethical principle. Acknowledging this reality not only contributes to the conceptual understanding of ethics principles, often expanding their scope, but, more importantly, allows us to provide comprehensive guidance on how to operationalise them in the context of GPAI.

Ethical concerns

Let us delve into some of the ethical, legal, societal and economic (ELSE) dimensions of GPAI by engaging with the literature on GPAI and generative AI, the latter being category of GPAI that has gained significant traction in recent years, particularly following the launch of ChatGPT. Starting with 'E' for 'ethical', scholars

have noted that established AI ethics principles, such as those put forward in ALTAI, remain relevant in the context of GPAI. However, the generality of GPAI and the production of differential outputs and human-like content of generative AI systems, has resulted in different ways of interpreting those principles, with some gaining more prominence, as well as on the identification of new AI ethics principles (Laine et al., 2025). For example, ‘human centric design’, ‘respect for intellectual property’ and ‘truthfulness’ have been identified as new AI ethics principles that arise directly from challenges and risks brought about by generative AI systems, i.e., deception of humans through language, the unlawful appropriation of copyrighted material in training datasets and the recurrent issue of hallucinations, whereby these systems produce fake or misleading claims (Laine et al., 2025).

At the same time, foundational AI ethics principles, such as human agency and autonomy, and societal and environmental well-being, have become more prominent. Most notably, the risk of manipulation and harm emerging from the recurrent and long-term use of sycophantic conversational agents are increasingly documented, both in academic research and through a growing number of lawsuits, following cases of suicide among young users (e.g., Character AI). Scholars have analysed how the anthropomorphic qualities of these language-based systems pose increase risks of deception, may erode cognitive and emotional development and prompt addictive usage (Mahari and Pataranutaporn, 2025; Paz, 2022). However, the impact of GPAI systems on users is far from uniform since their affordances can not only lead to harmful uses, but also provide emotional support and enhance communication among certain user groups (Boine, 2023). What is key, however, is that the risks and vulnerabilities arising from GPAI systems pertain not only to system safety and robustness – linked to the probabilistic nature of GPAI systems and the inability to fully predict and control system behaviour in advance – but also to the broader dynamics of human-AI interaction, which are particularly concerning in the context of private behaviour.¹⁶

It should be noted, however, that the current debate on GPAI is often framed by a safety-centric agenda, itself grounded by a concern with catastrophic risks of AI, including existential threats of AI. Such a framing is problematic for the governance of GPAI in four ways: 1) first, it assumes the current *status quo* of AI technologies is the only option for AI development; 2) it assumes that such risks can only be managed by frontier AI developers themselves; 3) it magnifies the risks posed by frontier models in a manner that gives the appearance of governance and regulatory frameworks being futile; and 4) it leaves out the bulk of actual, current risks, including those of algorithmic discrimination, copyright and environmental impact. The ELSE approach adopted herein, contrasts with the safety-centric agenda, not because it does not consider safety an fundamental ethical principle of Trustworthy AI, but because it seeks to engage not with prospective catastrophic risks, but with the actual, on-the-ground reality of GPAI. The ELSE framework demonstrates that risks are not only already materialising across organisations, human cognition and behaviour, including in private behaviour, but also broadens the debate on risks towards a non purely technical orientation that considers the governance of GPAI within a wider aperture prompted by ethics-based reasoning. The latter, not only enables the adoption of ethics-by-design approaches in the development of GPAI technologies, as it fosters a broader deliberation on the desirable practices to adopt with the aim of promoting individual and societal well-being. Indeed, ELSE points to the fact that Trustworthy AI is not simply a matter of technical fixes and containment, but is inherently a shaped by

¹⁶ For more on this issue see Section 5.2. on Emotional AI.

social, political and economic considerations that actively shape the form and operation of GPAI systems today.

Legal concerns

Under the AI Act, providers of GPAI models and deployers of GPAI systems are subject to specific obligations, commensurate with high-impact capabilities (i.e., GPAI models with systemic risk) and the context of deployment (i.e., high-risk) of these systems. The more extensive range of obligations for providers of GPAI models, including risk management (Article 55) and notification requirements upon reaching over 10^{25} FLOPs in training compute (Article 52) have been complemented with the Code of Practice for GPAI, to support compliance. However, the predominantly technical focus of obligations outlined in the AI Act does not sufficiently address the concrete steps and measures that organisations can adopt to address the ethical requirements for Trustworthy AI outlined in Recital 27, overlooking the ethical tensions that arise in this context. For example, disclosure of GPAI interactions does not preclude individuals from developing unsafe cognitive dependencies with those systems. From the perspective of deployers of GPAI systems, currently, the AI Act mandates disclosure of AI interactions (Article 26) and of AI-generated content (Article 50), alongside AI literacy (Article 4), unless they are deployed in high-risk contexts, where more extensive obligations, including human oversight (Article 14), would apply. The scarcity of requirements for GPAI systems leaves open questions with regards to the ability to ensure their responsible deployment.

Whereas the regulatory landscape for GPAI has expanded, both in the EU and elsewhere – see, for example, the U.S. state of California SB 53 bill and the RAISE Act in New York, both regulating “frontier models” – the predominantly technical orientation of these frameworks does not sufficiently address the concrete steps and measures that organisations must put in place to ensure the responsible development and deployment of these systems. The AI Act takes as basis the “Ethics guidelines for Trustworthy AI” developed by the AI HLEG, as outlined in Recital 27 but not only are the obligations for deployers for GPAI systems limited, as the provisions in place are in themselves complex to operationalise with several ethical tensions arising in this context. In particular, scholars have pointed to the fact that the way legal frameworks centre protection on design are based on the assumption that risks and vulnerabilities can be identified and mitigated *ex ante*. However, the malleability and generality of GPAI models and systems, challenges this approach since model behaviour changes based on human-AI interaction, which can result in emergent safety risks. In this regard, scholars have noted that, whereas legal frameworks are traditionally equipped to handle one-shot harms (discreet, traceable events), they struggle to address temporal or interactive harms that accrue over months or years of interaction (Teo et al., in press).

Still from a legal perspective, concerns have emerged in the literature regarding privacy – specifically, informed consent – and copyright. With regards to the latter, the use of large amounts of web-scraped and copyrighted data in the training of GPAI models, particularly generative AI models, is an ongoing point of contention, and complicates traditional provisions

under copyright law (Lemley, 2023; Samuelson, 2023). With regard to privacy, in particular, informed consent, current scholarship suggests that the traditional model of informed consent is increasingly ill-suited for the complexities of AI. Boine (2023) argues that the challenges lie not simply in a lack of user comprehension regarding the legal agreement and use policies, but the underlying operational opacity of GPAI systems and users' inability to understand practical consequences and long-term risks deriving from their use. Of key importance, is also the different nature of data underlying a GPAI systems to a significant degree, which is often intimate and algorithmically derived from 'intimate' human-AI interactions (Boine, 2023).

Societal concerns

From a societal perspective, several concerns arise in the literature, ranging from sycophantic conversational agents to the generation of artificial content by GPAI systems, including fake content causing the proliferation of information disorders, and the amplification of biases that might be damaging to the social fabric (Boine, 2023; Sharma et al., 2024). Moreover, given the generality and broad range of tasks GPAI systems can perform through natural language interfaces (e.g., LLMs), they raise broader concerns with regard to their impact on human cognition, both across private and professional contexts. In general, these issues cut across the individual and societal levels, with the impacts of individual human-GPAI interactions propagating more broadly, and potentially at scale, from the private to the public sphere. As the deployment and uptake of GPAI systems unfolds across these contexts, emerging research attempts to map, assess and quantify the cognitive implications of these systems, demonstrating not only how they transform human cognitive processes, such as critical thinking, but, crucially, how they might shortcut the former (Lee et al., 2025). For these reason, GPAI systems are considered to pose unique challenges to human competencies, with deskilling due to the long-term use of these systems being identified as a major risk (Natali et al., 2025).

A distinctive societal concern emerging in the context of GPAI is the environmental impact of these systems. The fact that GPAI models, especially those popularised as foundation models, require large amounts of data and compute for model training has raised numerous concerns due to the associated resource consumption, from raw materials to energy and water (Crawford, 2021; Ligozat et al., 2022; van Wynsberghe, 2021). As policymakers strive to reconcile conflicting policy objectives, most notably, a 'green' agenda and an ambition for remaining competitive in the AI landscape poised to have a cross-sectoral impact, several proposals have been put forward for measuring the environmental cost of AI, with a particular focus on GPAI (Berthelot, 2024; OECD, 2022). Nevertheless, consensus around measurement and information disclosure from BigTech players remains outstanding.

Economic concerns

Whereas the 'E' for the economic dimension of ELSE is not always addressed in ESL- literature, scholars have recently made the case for including an economic appraisal in ELSE frameworks for AI (Ryan and Block, 2025). Indeed, AI and, in particular, generative AI systems are enmeshed in politico-economic dynamics that span power concentration, geopolitical domination and large investments – the so-called 'AI bubble'. From an economic perspective, a key concern identified in academic literature is the dependence and reliance on foundation models developed by BigTech players, who own the resources

required to develop these models, namely, data and compute infrastructure. Large developers provide the baseline models upon which others will develop their AI applications through finetuning. This intensifies what was an already highly concentrated market structure with value creation remaining circumscribed in closed yet recursive loops, with data generated through model inference feeding the already extensive database of training data of a small group of AI developers (Cote and Aires, 2025).

Scholars have also singled out how strong forms of attachment to GPAI systems, such as personal companions, promoted through technical design for continuous user engagement, are exploitative of intimate relations and monetise users' emotional vulnerabilities, amounting to a form of "intimacy capitalism" (Mahari and Pataranutaporn, 2025; Teo et al., 2027). Beyond the domain of the AI system, strictly speaking, the potential and at-scale impact of deployment and uptake of GPAI systems across sectors, raises the broader concern of workforce displacement, although evidence in this regard is still limited (Brynjolfsson et al., 2025).

Building on the ethics principles selected for the GPAI research area and its different use-cases (D2.2), this section analyses how ethical principles have been operationalised by AIOLIA partners in four GPAI use-cases. The variety of GPAI use-cases, derived from areas as diverse as healthcare, personal companions, engineering and security, allows us to cross-analyse the bottom-up operationalisation of ethics principles and to lift overarching ethical concerns whilst acknowledging tensions between use-cases. In fact, the aim is not simply to produce context-enriched non-technical guidance that renders ethics approaches uniform within GPAI, but to acknowledge heterogeneity and tensions between use-cases' operationalisation process, and to identify unique sectoral needs. These will be incorporated as learning points for GPAI, providing insight on how practices in one context (e.g., healthcare) can inform approaches in others (e.g., personal companions).

5.1.2. Overlaps and differences in the operationalisation of ethics principles across GPAI use-cases

This section analyses overlaps and differences in the operationalisation of AI ethics principles across the four GPAI use cases. The analysis is organised based on the empirically developed ethics principles organised, along the 12 UC-defined ethics principles, rather than merging them into ALTAI key requirements. However, the presentation of principles is ordered along the logic of ALTAI as in D.3.1 (see 4.4., page 40).

ALTAI requirement #1 Human agency and oversight

Human Oversight

In UC2, human oversight over the AI system in the context of professional behaviour is a core principle, alongside autonomy, with the target of autonomy being the human operator of the AI system. In UC2, engineers must be able to technically control, interpret and intervene in the system's operations. From the perspective of professional behaviour, autonomy is about agency – i.e., the ability to act – more than it is about the freedom of choice and non-manipulation (see Autonomy / User Agency below). Human oversight is regarded not simply in fulfilment of safety, i.e., to ensure that safety checks are correctly undertaken, but as an enabler of accountability, and a mitigator of deskilling. However, AIOLIA partners

voiced concerns over the need to exert human oversight and the more passive or less engaged role of engineers in the workflow pipeline, with the repetitiveness of the task risking to reduce attentiveness (see Over-reliance and Deskilling).

In UC4, oversight is understood from the perspective of autonomy, which is a component of the principle of “freedom of expression and non-censorship”, foregrounded in fundamental rights.¹⁷ Fundamental rights can be considered as stage 0 of ALTAI, in the sense that a fundamental rights impact assessment (FRIA) must be conducted prior to delving into the seven ALTAI requirements, and it is currently a mandatory requirement for deployers of high-risk AI systems under the AI Act (Article 27). In UC4, it pertains to ensuring that AI systems used by law enforcement authorities do not remove content on the basis of it being controversial or unpopular, preserving users’ autonomy and agency to express their opinions and guaranteeing the protection of freedom of expression. In this context, users’ autonomy is upheld through proportionate content removal measures with “human-in-the-loop at every stage”, which must be the least restrictive possible, and through equal, hence non-discriminatory, treatment of different groups and political standpoints. This is achieved through culture sensitivity checks and building alternative narratives to mitigate cultural divides.

From a user safety perspective, UC5 operationalises the component of human oversight, part of the principle of autonomy, through human moderation of edge cases flagged by automated moderation system (see also Technical Robustness and Safety). UC6 also identified the component of oversight as a means to operationalise the ethical principle of accountability. However, since, in this context, oversight pertains not to system oversight but to oversight of patient-system-therapist relations, these measures have been accounted for under the principle of Accountability described further below.

UC4	UC5	UC6
Human checks at each stage of building narrative, to ensure narratives are not entirely reliant on AI LLMS that could be trained incorrectly, operationalising the component of human oversight. ¹⁸	Establish human oversight for moderation of edge cases flagged by automated moderation system, to operationalise the component of human oversight.	Continuous monitoring of patient experiences and adaptation of protocols, to operationalise the component of oversight.
	Legal compliance review to ensure moderation decisions are legally compliant, to operationalise the component of security measures.	Legal study on whether and how the MDR and AI Act would apply
		Conduct an independent ethical review (e.g., REC), to operationalise the component of oversight.
		Involve ethicists/lawyers in design of deepfake therapy, to operationalise the component of oversight.

¹⁷ Although UC4 privileges “freedom of expression and non-censorship”, we have opted to include it part of the ALTAI principle “Human agency and oversight” for systematisation and coherence.

¹⁸ Due to privacy and information disclosure concerns, the measures provided by UC4 tend to be rather high-level and abstract.

		Attending to the investigative nature of this UC, conduct qualitative research with patients and therapists to assess risk of over-attachment and other unintended consequences of deepfake therapy, to operationalise the component of risk of over-attachment and dependency.
		Use only in the research context before ascertaining whether deepfake therapy can constitute good care, through ethical, legal, social and psychological research, operationalising the component of subsidiarity, proportionality and effectiveness.

Table 8 - Comparison of UC organisational measures addressing the principle of “Human Oversight”. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Autonomy / user agency

Across all GPAI use-cases there is a rather uniform understanding of autonomy as the preservation of human decision-making ability (either professional or private), having transparency and human-in-the-loop mechanisms as a pre-requisite, and mobilising regulations (EU and national) and standards (where available) for measuring adherence. However, a more granular analysis of the different UCs reveals a more nuanced picture, especially with regard to the context and subject of autonomy – professional or private – as well as its temporal scope and the nature of the measures identified to address the ethical principle. Whenever autonomy pertains to professional behaviour, it is conceptualised in alignment with Human Oversight, as described above.

UC5 and UC6, understand autonomy from a more relational perspective, insofar as it concerns the preservation of autonomy in private behaviour. As such, the focus is not primarily the ability to act over the AI system, as is prioritised in professional behaviour, but the ability for users and patients to make informed decisions and choices, beyond the technicalities of the AI system itself. In this context, tensions around informed consent are a cross-cutting issue.

In UC5, the aim is to ensure AI companions do not manipulate users, however, significant concerns arise over what constitutes permissible behavioural and cognitive modification in the context of AI companions purposefully developed for daily organisation, habit tracking and personal development. Autonomy is not understood merely as the ability to choose here and now (e.g., as assumed in current informed consent mechanisms), but as the preservation of this faculty in the long-term, despite their potentially lasting behavioural, cognitive and emotional imprint resulting from the subtle and cumulative interactions with AI companions. As such, it also implies users’ self-determination and ability to control their personal development goals. In this context, informed consent and transparency requirements foreseen under the AI Act, including labelling AI-generated content, are deemed insufficient to guarantee users’ autonomy. To operationalise the principle of autonomy, UC5 puts forward two interrelated organisational measures targeting the component of system customisation: 1) drafting an internal policy on permissible behavioural change; 2) conduct a competitive analysis to assess the level of customisation /

personalisation of personal companions enabled by other market. Two additional organisational measures have been put forward under the principle of autonomy to operationalise the component of transparency and user understanding – these measures are listed under the Transparency subsection.

In UC6, patients are the primary subject of the ethical principle of autonomy, which is addressed prior and during deepfake therapy. Prior to the use of AI, autonomy sits on the component of transparency, namely on practitioners informing patients regarding the availability of alternative treatment options and giving them the possibility to choose (see measures listed under Transparency and Privacy). During the use of AI in deepfake therapy, autonomy is operationalised through the component of to operationalise the component of risk of over-attachment and dependency, through measures that seek to limit dependencies that might be harmful to patients’ decision-making ability and capacity to distinguish fact from fiction, especially in the context of grief and PTSD therapy. One measure used by UC6 to prevent this, is the non-modification of voice in the deepfake, with the therapist’s voice (not the perpetrators’) being the only voice integrated in the GPAI. In this context, the long-established medical practice of informed consent – which required providing information to patients about the risks, benefits and treatment alternatives – is further complicated. On the one hand, practitioners are unable to provide complete information about the AI system, especially due to its complexity, opacity and not fully predictable behaviour. On the other, consent is understood not only as an established medical practice, but also as a component of privacy that extends to patient data and to those who are the object of deepfake creation (e.g., perpetrators, deceased). Importantly, the latter is considered, at times, impossible and undesirable (e.g., can lead to traumatic experiences), or even, as requiring other family members, in the case of deceased people.

UC4	UC5	UC6
Human-in-the-loop at every stage of building alternative narratives to mitigate cultural divides, to address the component of autonomy and agency.		
	Draft an internal policy on permissible behaviour change, to operationalise the component of system customisation.	
	Conduct competitive analysis to assess competitors’ policies and market positioning, operationalising the component of system customisation.	
		Avoid introducing false or wishful narratives (e.g., reconciliation or afterlife dialogues), to operationalise the component of risk of over-attachment and dependency.
		Tailor therapy to the patient, including refraining from deploying deepfake therapy to patients who might be

		prone to over-attachment, to operationalise the component of risk of over-attachment and dependency.
		Arrange a preparation meeting to manage patient expectations and reiterate fakeness of the deepfake, to operationalise the component of risk of over-attachment and dependency.

Table 9 - Comparison of UC organisational measures addressing the principle of “Autonomy / User Agency”. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Over-reliance and deskilling

AIOLIA use cases predominantly address societal wellbeing from the perspective of the impact of GPAI on work and skills, hence on professional behaviour, namely in UC2 (safety engineers) and UC6 (deepfake therapy), with strong links with accountability.

In UC2, the risk of overreliance and deskilling is addressed as a core concern. Operating in a sector with a strong safety culture, that has moved from reactive safety assurance to proactive due diligence, capable of identifying risks and potential harms before they happen, UC2 considers that the incorporation of AI to speed up car software safety checks might lead to a loss of “tacit safety expertise”. One key concern is the ability to maintain the required level of knowledge and skills among engineers in the long term, when AI systems form part of the production pipeline, and to consider this challenge in relation with the level of seniority of staff. On the one hand, senior engineers are exposed to deskilling by having a less interventive role in safety checks – i.e., as human-in-the-loop and review of safety critical outputs – and often being pressured to cope with KPIs that might be detrimental to the responsible use and oversight of the tool. On the other hand, junior engineers are exposed not to deskilling per se, but to ‘non-skilling’, which means that, if AI decision support systems are integrated too early in their career, they might not have the chance to develop the skills and safety culture required to excel in the job and maintain industry safety standards. This not only affects professional behaviour by jeopardising the ability to exert day-to-day oversight, interpreting AI outputs and reasoning along complex causal chains, but it also has labour and societal impacts in the medium / long term. As such, it requires balancing short-term goals, i.e., efficiency gains, with future impacts.

As a result, UC2 emphasises the need to preserve human skills and expertise through continuous training and skills development, including on judgment and capacity for independent decision making, diagnostic capability and situational awareness. Intergenerational training and monitoring of human-AI interactions are also part of the implemented measures. The provision of training is regarded as essential not just to maintain competencies (knowledge, cognitive skills and ethical awareness), but to update these so that they keep evolving alongside the changing nature and capabilities of the AI system, including through feedback loops for mutual human-AI learning and ongoing training programmes, ensuring the AI system continues to be a decision-support tool only. However, one tension that arises in this regard is the need for continuous training is often faced with a lack of adequate resources and employee motivation. Addressing the risk of overreliance and deskilling is seen as paramount to adequate oversight, accountability and safety.

In UC6, professional competence is regarded as a component of accountability. The primary concern is not overreliance, but the need for therapists to undergo specialised training to make use of deepfakes in a therapeutic context and to maintain their skilled judgment despite the introduction of the AI system.

UC2	UC6
Conduct regular human-in-the-loop trainings to operationalise the component of preservation of human skill and expertise.	Conduct trainings to ensure therapists are skilled in trauma dynamics, operationalising the component of professional competence.
Implement knowledge capture and mentorship programmes, exchanging knowledge between junior and senior safety engineers to ensure transmission of tacit reasoning, operationalising the component of preservation of human skill and expertise.	
Implement adaptive training programmes for continuous skill development to operationalise the component of preservation of human skill and expertise.	
Conduct systematic analysis of human-AI performance across tasks, embedding this into the refinement of safety procedures, training, and governance, to operationalise the component of feedback and learning loops for human adaptation.	
Put in place structured onboarding and foundational AI-for-safety training as a mandatory onboarding curriculum, to operationalise the component of training, education and continuous skill development.	
Enable continuous professional development and certification renewal for safety engineers, operationalising the component of training, education and continuous skill development.	
Establish an AI governance and oversight committee to institutionalise collective accountability and ensure safety decisions are reviewed through multiple lenses, operationalising the component of organisation policies for shared responsibility.	

Table 10 - Comparison of UC organisational measures addressing the principle of “Over-reliance and deskilling”. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #1 Human agency and oversight

What stands out from AIOLIA use cases is how autonomy and human agency is understood differently in the context of professional and private behaviours. In the first, professional behaviour, it pertains to the immediate ability to act over the AI system, being more closely aligned with human oversight. In contrast, in the context of private behaviour, the focus is the long-term impact of the system and the ability for users to choose and act over time. For this reason, in what concerns professional behaviour, the principle is linked to AI system transparency and traceability; whereas in private behaviour the focus is on assessing what amounts to an adequate and safe level of system personalisation and human-AI relations. This nuance and level of granularity is absent from ALTAI.

Further, based on insights from AIOLIA’s GPAI use cases, ALTAI lacks a temporal dimension, capable of accounting for the long-term behavioural and cognitive risks that might arise from the long-term and

recurrent use of GPAI systems. This applies both to private and professional contexts. In the context of private behaviour (e.g., conversational AI agents) ALTAI’s narrow understanding of manipulation as immediate and explicit impact over decision-making, overlooks the cumulative impacts of these systems, which call for the formulation of new, behavioural-oriented measures, such as the policy on permissible behaviour change put forward by UC5. Long-term impacts also apply in the context of professional behaviour, most notably through deskilling and loss of competencies. Whereas ALTAI acknowledges the impact of AI on skills and acknowledges the need for re-skilling (requirement #7 on societal well-being), it lacks sensitivity to seniority levels and how organisations must develop different strategies and measures targeting junior / senior staff to enable the preservation and development of skills.

ALTAI requirement #2 Technical robustness and safety

Technical robustness and safety

AIOLIA’s UC2 (car safety engineers), UC5 (AI companions) and UC6 (deepfake therapy) introduce safety as a key ethical principle, centred on the safety of AI system’s outputs and those impacted by these, either directly or indirectly. However, in UC2, safety is understood from a predominantly technical perspective, as a hard engineering requirement, namely that of AI safety evaluations being accurate through system reliability and robustness which, by extension, ensures road users’ safety. In the context of professional behaviour, UC2 places a strong emphasis on the ability for engineers to trace and test system performance throughout its lifecycle, and to ensure that the AI system is robust, resilient, fair and reliable. The measures implemented are predominantly technical in nature, aligned with ISO automotive industry standards and the AIA. From an organisational point of view, the need to promote and nurture the safety culture characterising the field of car safety engineering is also emphasised, alongside a quality management system for AI outputs.

UC2
Promote consistent safety culture, adherence to standards and responsible use of AI tools within the organisation through trainings and clear allocation of responsibility, to operationalise component of technical robustness and resilience.
Implement an AI quality management system to ensure reliability, reproducibility and compliance of AI-system safety analysis, operationalising the component of technical robustness and resilience.
Conduct cross-disciplinary reviews of the AI system and in its integration into workflows, to identify hazards, biases and process gaps, operationalising the component of technical robustness and resilience.

Table 11 - Comparison of UC organisational measures addressing the principle of “Technical robustness and safety”. Only UC2 put forward organisational measures addressing this principle. No measures were identified for UC4, UC5 and UC6.

Safety / Non maleficence / Human Safety

Contrary to the system-centric safety of UC2, in UC5 and UC6, safety is understood from an inherently sociotechnical and cognitive perspective, as the implementation of measures for the protection of users against harmful AI-generated content.

On the contrary, UC6 emphasises medical practitioners’ and organisations’ role in safety assurance, i.e., non-maleficence, with the principle gaining a more social dimension. Whereas, in the medical sector, there is an abundance of standards guiding the development, deployment and use of healthcare software, concerns remain regarding safety assurance, especially given how AI safety intersects with professional behaviour, rather than simply being a technical feature of the system. UC6 emphasises key behavioural

dimensions of medical practice, namely the need for proportionality, in the sense that the benefits of undergoing deepfake therapy should outweigh risks, and subsidiarity, with therapists being responsible for choosing the least intrusive treatment, including traditional, non-AI based, exposure therapy, as treatment should be effective, i.e., achieve intended purpose with least harm possible. In this context, we see not only an organisational-centric set of measures for safety assurance, including the personalisation of clinical protocol and the presence of therapists throughout deepfake therapy sessions, but the questioning of introducing AI systems in specific professional settings in the first place, in this case, trauma therapy. The safety onus is on the practitioner and the medical organisation, although harms induced by technical glitches and data collection – i.e., images of perpetrators for deepfake creation – are mentioned as safety concerns. The latter, however, is not a consensual concern among UC partners, since traditional forms of exposure therapy also rely on images that can trigger distress. Societal wellbeing is also a component of UC6’s ‘non-maleficence’, spilling over to the ALTAI principle “Environmental and societal wellbeing” below, where it is discussed. The principle of safety becomes more complicated in the context of private behaviour, in particular AI companions (UC5). While there are technical considerations with regard to the mechanisms for ensuring safety, namely through an automated moderation system and human oversight of flagged content/users, these measures are contingent on value-laden judgements about right/wrong and good/bad that are often context-dependent. While developers already implement measures to mitigate depictions of violence and physical harm (e.g., prohibition to depict blood), it is not commonly agreed upon where to draw the line between safety measures for harm prevention and system desirability among users.

UC5	UC6
Ongoing review of safety policy and risk tiers to address emerging risk patterns, operationalising the component of user protection.	
Cross-functional tiered definition and policy of risk classification to ensure comprehensive and differentiated rules for safety violations, operationalising the component of user protection.	Personalise clinical protocol to patient needs, operationalising the component of effectiveness.
Establish a progressive intervention and escalation protocol to respond to violations in an appropriate and systematic manner, operationalising the component of security measures.	
Ensure payment platform security measures are aligned with those of the organisation, operationalising the component of security measures.	
Elaborate an organisational policy specifying system scope boundaries, determining what it supports and excludes, to operationalise the component of scope boundaries.	
	Consider adequacy of less-intrusive exposure therapies before initiating deepfake therapy, to operationalise the component of subsidiarity and proportionality.
	Keep deepfake therapy limited to research context, operationalising the component of effectiveness.

Table 12 - Comparison of UC organisational measures addressing the principle of “Safety / Non maleficence / Human safety”. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #2 Technical Robustness and Safety

AIOLIA use cases demonstrate that the ALTAI principle of Technical Robustness and Safety is too narrow to account for the challenges brought about by GPAI systems and their varied contexts of deployment. Specifically, whereas ALTAI addresses safety-critical failures – in a manner that aligns with UC2 system-centric approach – it misses out on non-technical dimensions of safety, such as user / patient harm and implied harm, which are in fact the core focus of AIOLIA’s use cases. For this reason, this ALTAI principle would benefit from being expanded to account for human safety, and potentially, to integrate the bioethical principle of non-maleficence – broadly understood – as an integral part of responsible and trustworthy AI (see 3.1.5. for additional discussion).

ALTAI requirement #3 Privacy and data governance

Privacy and data protection

In AIOLIA, the principle of privacy has been identified as relevant in UC5 (personal companions) and UC6 (deepfake therapy). In both cases, GDPR foregrounds the understanding of the principle as the responsible handling of user data. However, UC6 is particularly concerned with the reputational cost of data breaches for the subjects of deepfake creation, who are, in this case, aggressors, and of whom there is not necessarily consent for the processing of their personal data, as well as with the need for seeking consent from family members, in the case of grief therapy.

In UC5, the emphasis is on how the intimate and conversational nature of personal companions creates new privacy stakes beyond typical applications, given the extent of habit data provided by users and the data inferred through long-term tracking for the creation of behavioural profiles and system customisation. For this reason, a core concern that arose in the context of UC5 was compliance with GDPR’s data minimisation requirements, since both system functionality and safety measures, namely the automated moderation system flagging forbidden or dangerous uses, require extensive data collection. Without the latter, user safety can be jeopardise, given the sheer amount of interactions to be monitored. Moreover, the implementation of mechanisms to allow users to control their data, particularly in the case of data deletion, revealed challenging in the context of GPAI systems, as the personalisation of the AI companion relies on recursive learning and makes data withdrawal challenging to be achieved in practice.

Gaps in ALTAI requirement #3 Privacy and data governance

The large amounts of data required to effectively train and deploy GPAI systems, particularly when these intersect with private behaviour, is fundamentally at odds with the ALTAI and GDPR requirements concerning data governance. Whereas AIOLIA use cases implement comprehensive compliance procedures, the extent to which compliance is possible in the context of GPAI systems remains an open question and calls for adopting a more grounded approach in ALTAI. Similar trade-offs emerge regarding informed consent, with the involvement of third parties and the use of data from non-consenting individuals posing both legal and ethical questions that, as GPAI deployment intensifies across different contexts, such as AIOLIA’s deepfake therapy, require consideration and a standardisable approach.

UC5	UC6
Implement a comprehensive GDPR compliance programme, including data protection policies, staff training and regular audits, to operationalise the component of user consent and transparency.	
Draft a transparent privacy policy, to operationalise the component of user consent and transparency.	
Conduct a privacy impact assessment, to operationalise the component of user consent and transparency.	Consider need and suitability of seeking consent from those depicted in the deepfakes, to operationalise the component of privacy. Consider need and suitability of seeking consent from those depicted in the deepfakes, to operationalise the component of privacy.
Implement user data rights processes, including transparent data policies and consent mechanisms, providing users with the right to access, delete and transfer their data, to operationalise the component of data minimisation, use and storage.	
Draft data sharing agreements attending to different jurisdictions and data management processes, to operationalise the component of third-party data sharing and compliance.	

Table 13 - Comparison of UC organisational measures addressing the principle of "Privacy and Data protection". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

ALTAI requirement #4 Transparency

Transparency and explainability

Transparency is present across all GPAI use cases. However, rather than being understood as a standalone ethical principle, transparency is purpose-driven and context dependent, being always presented as a component of different ethical principles, namely of autonomy (UC2, UC6), privacy (UC5) and non-bias, fairness and non-discrimination (UC4). In the context of GPAI, transparency predominantly addresses those affected by the AI system's decisions (e.g., patients, users), rather than the professionals making use of the system, with UC2 being the only exception. In UC2 (car safety engineers), transparency is understood as engineers' ability to understand the behaviour and decision-making process of the AI, e.g., through visualisation of model behaviour. It overlaps with all three elements of transparency described in ALTAI, namely, explainability, traceability of safety decisions and communication of model boundaries and limitations. Transparency is regarded as a key component of engineers' autonomy and accountability, ensuring that engineers remain liable for safety decisions and in control of the model.

In the remaining UCs, namely UC4 (hate speech detection), UC5 (AI companions) and UC6 (deepfake therapy), transparency is 'public-facing', that is, it relates to explaining and communicating with users and patients as private individuals – transparency is approached from the point of view of private behaviour. In UC6, transparency is not related to the AI system per se but with the necessity to inform patients about different treatment options, giving them the ability to choose. In UC4 and UC5, transparency pertains to the need of informing users about the AI systems' operation, namely, transparency over extensive data processing (inferred and implied data) required to uphold privacy and safety obligations (UC5) and

transparency over the systems' criteria to block content (UC4), providing explanations to users as well as the ability for appealing decisions. The two UCs imply two variable dimensions of transparency, with the first pertaining to transparency as communication and the latter to transparency as explainability.

Transparency is a major concern in UC5, which signals that transparency with regard to data processing might have adverse effects on private behaviour, since these systems are designed to engage in detailed, personal dialogues with the user. UC5 considers that the granularity and extent of transparency must be subject to reflection, given that providing maximal transparency is not always an effective strategy as it might be detrimental to the effectiveness of the AI companion. When individuals share deeply personal thoughts, struggles, or intimate details, they expect a private, confidential space while simultaneously being aware that their conversations are monitored for risk detection. Users might perceive the monitoring system as intrusive, which can lead to a decrease in trust or increased jailbreaking attempts that can be more damaging to user safety. This awareness can inhibit openness, foster a sense of surveillance, reduce trust, and ultimately undermine both the effectiveness of the AI system and the reliability of safety monitoring. Excessive disclosure may conflict with users' desire for natural, fluid interaction that is not entirely predictable or even deliberately opaque.

Gaps in ALTAI requirement #4 Transparency

The ALTAI principle of transparency draws on the assumption that more transparency equals to more trust. Whereas AIOLIA use cases demonstrate reiterate this, they also demonstrate that transparency measures can bring additional risks, including jailbreaking, potentially jeopardising the reliability of safety measures. Moreover, AIOLIA UCs demonstrate that transparency, without its links to autonomy (UC2, UC6), privacy (UC5) and non-bias, fairness and non-discrimination (UC4), is an empty signifier. In fact, measures related to transparency vary and should be implemented differently according to the specificities of the GPAI system's users (e.g., professional vs private); users is not a uniform category and more nuance is required in the operationalisation of this principle. Lastly, AIOLIA UCs (particularly through the medical UC6) demonstrate that, in the context of GPAI, transparency must not only target the model and its functionality, but also its surroundings, including providing users with non-GPAI-based alternatives (e.g., treatment options is medical care).

UC2	UC4	UC5	UC6
Clear documentation describing the AI system performance metrics and confidence intervals, to enable transparency of safety critical performance			
	Publicly document content moderation and classification standards to allow for understanding of hate speech enforcement thresholds, operationalising the	Put clear community guidelines in place, as well as educational resources for users	

	component of transparency of criteria		
		Establish an appeals process for users to contextualise moderation decision	
			Informed consent from the patient, providing information about the therapy and treatment options to those undergoing treatment

Table 14 - Comparison of UC organisational measures addressing the principle of "Transparency and explainability". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

ALTAI requirement #5 Diversity, non-discrimination and fairness

Non-bias, fairness and non-discrimination

Across AIOLIA's use cases, diversity, non-discrimination and fairness has been identified as a core principle in UC4 (hate speech detection) and as a component of "robustness, safety and reliability" in UC2 (car safety engineers). The object of AI – social discourse in UC4 and technical software in UC2 – adds different nuances to the principle.

UC2 addresses the principle predominantly from the perspective of fairness, as the need to ensure that safety assurance is applied consistently and transparently to all stakeholders, including pedestrians, passengers and operators, as well as across different operating contexts and system components. As AI safety systems rely on AI-assisted data analysis (e.g., risk classification, failure pattern recognition), effort must go into ensuring that these algorithms are not biased or inconsistently applied. As such, fairness is directly related to safety and technical consistency to mitigate physical risk, addressed through technical measures.

In UC4, the principle becomes more complex, especially because it is not simply about having statistically representative data, but relates to the ability of operating within particular and distinct cultural contexts. As the UC intends to detect harmful or illegal content, including hate speech, and to produce alternative narratives, it requires treating content impartially and equally, not just by ensuring that training datasets are inclusive and representative but by mitigating narrative developers' personal biases. As such, upholding the principle requires cultural sensitivity and ethnographic insights so that narrative developers and the human-in-the-loop revising flagged content are able to understand the contexts of lived experiences and avoid discriminatory or unfair behaviour, that can marginalise specific communities. From the perspective of content creators, they must not only have access to the criteria used to delete content, but be given the possibility to appeal decisions.

UC2	UC4
	Put in place a formal policy to determine prohibited biases, equality principles and impartial decision-

	making standards, to operationalise the component of equality and impartiality. ¹⁹
Conduct multistakeholder validation workshops for systemic bias prevention, to operationalise the component of technical robustness and resilience.	Include diverse stakeholders in system design and policy decisions through consultation with civil society, impacted users and others, to operationalise the component of inclusivity.

Table 15 - Comparison of UC organisational measures addressing the principle of “Non-bias, fairness and non-discrimination”. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #5 Diversity, non-discrimination and fairness

There is strong alignment between AIOLIA UCs’ conceptualisation and operationalisation of this principle and its formulation in ALTAI.

ALTAI requirement #6 Environmental and societal well-being

Human well-being

In UC5 (personal companion), well-being is centred on the individual chatbot user and, while there is a recognition of broader risks to societal well-being in UC6, namely the normalisation of deepfakes, there is a clear struggle to address it.

In UC5, wellbeing is framed as “human well-being” and directed at ensuring that the AI companion promotes users’ health, recognises mental health crisis and has a controlled / limited scope at the level of inputs and outputs with the aim of not putting users at risk. This is pursued through a mix of organisational and technical measures, including ethics reviews for system features, as well using LLMs, classifiers and keyword detection mechanisms for crisis detection and for limiting the scope of the AI companion. Importantly, users with conditions that affect their cognitive function are regarded as particular vulnerable to both the benefits and harms of habit formation personal companions.

UC6 is the only GPAI UC pointing to the dimension of societal well-being from the perspective of society at large under its principle of “non-maleficence”. It argues that the pursuit of deepfake therapy risks normalising deepfakes and eroding societal trust if and when deepfakes are to be used outside the more restricted medical experiment being undertaken in UC6. Importantly, it points that a discussion on the social acceptance of deepfake therapy should be held at the societal level.

UC5	UC6
Internal ethics review for new features, to operationalise the component of promotion of users’ health.	
Regular advisory consultation from external mental health professionals to inform system boundaries and approaches, to operationalise the component promotion of users’ health.	
User feedback collection, to operationalise the component of crisis detection.	

¹⁹ Due to privacy and information disclosure concerns, the measures provided by UC4 tend to be rather high-level and abstract.

	Avoid introduction in chatbots that can be used without a medical purpose and without the presence of a therapist.
	Avoid introduction in chatbots that can be used without a medical purpose and without the presence of a therapist.
	Avoid introduction in chatbots that can be used without a medical purpose and without the presence of a therapist.

Table 16 - Comparison of UC organisational measures addressing the principle of "Human well-being". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #6 Environmental and societal well-being

While the ALTAI requirement of Environmental and Societal Well-being addresses most of the concerns identified in AIOLIA GPAI UCs, the challenge lies predominantly in the operationalisation of societal well-being, as there is a clear struggle within organisations to address this principle by themselves and through internal resources. Additionally, this ALTAI requirement would benefit from being expanded to encompass individual well-being, as distinct from user safety.

ALTAI requirement #7 Accountability

Accountability and responsibility

Accountability is a core principle of two UCs, UC4 (hate speech detection) and UC6 (deepfake therapy), and a component of "oversight and autonomy" in UC2 (car safety). In all these UCs, accountability pertains to the domain of professional behaviour, and is strongly linked with human agency, oversight, and professional competence. Likewise, in all UCs accountability implies that professional practitioners remain answerable and responsible – if not liable (emphasised in UC2) – for the AI systems' use and decisions, despite the varying nature of accountability, which is technical (UC2), social (UC4) and clinical (UC6). This is realised through human-in-the-loop mechanisms, such as human sign-off of AI decisions (UC2), human review of borderline cases and system responsiveness (UC4) and clinical standards and judgement (UC6).

Auditability is also present as a core element of accountability across all UCs, with all addressing the need for external or independent auditing. Although in UC2 and UC4 the focus is on the AI model, through e.g., decision logs, audit trails and external audits, in UC6 the focus is on the practitioner / therapist, with the medical practice being evaluated by an independent ethical review. It follows that all three UCs mention professional competence and training as essential elements of accountability, with UC2 and UC4 emphasising the risk of overreliance on the AI system.

In UC2, the discussion on accountability is more extensive, particularly due to the distributed nature of workflows and decision-making. In car safety engineering, the GPAI sits within broad pipeline, influencing subsequent steps. This is one of the reasons why decision logs become so important, since they enable the attribution of liability to the person that signed the report. It allows controlling and consulting the process of safety checks in case there are incidents or unforeseen circumstances. However, one concern that arises is the suitability of logs to ensure safety, since logs do not provide any assurance that the right decision was taken.

While in UC2, technical decision logs become the pillar of accountability, in patient-facing professional roles, such as in UC6, accountability is regarded as relational and not only GPAI-centric; it implies not just a compliance log, but answerability (or justifiability, as formulated in UC1 – see Section 5.3.) that is, the ability for practitioners and medical organisations to justify their decisions, actions and consequences, to patients, peers, regulatory bodies and society more broadly.

UC2	UC6
Role-based responsibility assignment, to enable the components of traceability and accountability for safety decisions, operationalising the component of accountability and traceability.	Develop guidelines that highlight that the AI technology is secondary and the responsibility remains with the therapist, to operationalise the component of human agency and responsibility.
Establishment of an AI governance and oversight committee for shared responsibility, operationalising the component of shared responsibility.	

Table 17 - Comparison of UC organisational measures addressing the principle of "Accountability and responsibility". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #7 Accountability

ALTAI considers accountability predominantly as a question of auditability (technical traceability), whereas AIOLIA GPAI UCs demonstrate how accountability in sectors like healthcare is not technical but social, i.e., it requires answerability and justifiability. Moreover, the ALTAI requirement of accountability would benefit from addressing the core issue of distributed workflows and supply chains, providing guidance on how to ensure that liable individuals (e.g., doctors) have the capacity to control, trace and interpret AI outputs and recommendations.

5.1.3. Remaining challenges and concerns

What is most notable in the analysis conducted is the lack of overlap in the measures identified by GPAI UCs to operationalise the same principles. This reflects the generality of GPAI and the fact that different contexts of deployment require the adoption of different provisions at the organisational level. Convergence at the level of UC measures, tended to occur in areas like human oversight (a key obligation for high-risk systems in the AI Act), transparency (of documentation and criteria), involvement of diverse stakeholders in system design and evaluation, and personalisation of the GPAI system to enable human safety (this translates as differentiated risk classification in UC5 and adaptation of clinical protocols in UC6). While the generality of GPAI already posed challenges to policymakers in the elaboration of the AI Act, it might continue to do so despite efforts to address these models in Chapter V of the regulatory framework.

One key insight arising from the operationalisation of AI ethics principles conducted by AIOLIA partners is that use cases operating in more standardised domains, such as the automotive industry, tend to have a clearer path for translating ethical principles into practice. Rather than hindering GPAI uptake, standards provided a practical baseline and offered certainty on how organisations should approach AI development and deployment. Particularly, industry-specific standards and regulations, such as the ISO automotive industry standards or the Medical Device Regulation, supported not simply the operationalisation of AI

ethics principles for use cases operating in those areas, but provided insights on the concrete organisational measures that could be applied to the GPAI research area as a whole. This suggests that the EU should pursue the development of standards, including the harmonised standards to support compliance with the EU AI Act, with urgency and, potentially, engage in furthering sector-specific standards, particularly in the context of Emotional AI systems.

Whereas there is growing evidence on the environmental impact of GPAI models, particularly given the large amounts of data and compute for model training and the associated consumption of natural resources, from raw materials to energy and water, this has not emerged as a central issue in AIOLIA's GPAI use cases. Rather than reflecting a disregard for environmental well-being, this suggests that other concerns related to GPAI development and deployment have gained prominence in the process of operationalising ethics principles, which, as the detailed analysis conducted in section 5 suggests, are more closely linked to the immediate organisational context. Organisations are predominantly concerned with the challenges they can directly tackle, such as the preservation of skills, human safety and the adequate human oversight, with broader environmental concerns being less addressable from the perspective of a single organisation. As a result, and as suggested in the policy recommendations provided herein (Section 3), policymakers should take charge of broader societal and environmental concerns, and provide a platform for such discussions to take place. Crucially, with regards to the environmental impact of GPAI, the formulation of standard indicators for its measurement is currently overdue, as it is a clear direction on the policy priorities that should drive societal development and require trade-offs between competitiveness and environmental sustainability.

5.2. EMOTIONAL AI

Emotional AI has traditionally widely been linked to emotion *recognition*, or inference, of emotional states via biometric data (AI Act Art. 5.1(f); Stark and Hoey, 2021); note that such use, if taking place in the workplace or education institutions, is forbidden by the AI Act (Article 5.1(f)). Emotional AI has also been linked to emotion *emulation* – the process of imitation of human emotional states, for example in the case of chatbots (IEEE Std. 7014-2024), and this kind of emotional AI has become widespread due to the rise of LLMs. Some LLMs are designed to specifically tailor to users' emotional needs and are known as AI companions, which are marketed as virtual friends, romantic partners or personal assistants, claiming to provide emotional support, entertainment, companionship, and even coaching (Bernardo, 2025); examples include Replika and Character.ai. However, emotion emulation is an inherent LLM capability, stemming from the training process, which gives rise to the intimate use of general-purpose chatbots like ChatGPT that are not specifically marketed as emotional partners. Note that LLMs are capable of both inferring and mimicking human emotion, that is, of both emotion recognition and emotion emulation – an ability referred to as emulated empathy (McStay, 2025). The forbidden case of the AI Act that deals with emotion inference does not consider language-based recognition of emotional states as emotion inference (Recital 44), reflecting its narrower understanding of the research area.

Deepfakes can also be categorised as emotional AI if they are used to connect with the user's emotions. This is the case of deepfake therapy where deepfakes are offered as exposure therapy – often used to treat patients with post-traumatic stress disorder. Unlike chatbots, the modality in use is not language,

but vision. Enabled by a specific kind of artificial neural network called variational autoencoder (VAE), the face of the therapist gets overlaid with the facial features of the perpetrator (Hoek et al, 2025). As a result, the therapist retains the control of the emotions they convey and the words they use. This is an unusual instance of emotion AI because the AI system is not used to detect emotions or infer them, but to rather elicit a certain emotional response from the patient. This adds emotion *elicitation* to the group of tasks that can be done with emotional AI, next to *recognition* and *emulation*.

The therapeutic use of deepfakes differs significantly from both AI companions and traditional biometric emotion recognition. Unlike AI companions, deepfake systems exercise no agency in emulating emotion; unlike biometric approaches, the work of decoding the user's emotions rests with the human therapist rather than the AI. Still, due to the potential future trajectory of deepfake therapy²⁰ raises concerns shared with the research area emotional AI, we categorise it as such. It is an AI-enabled process of emotion elicitation with partially artificial emotion emulation, which can evolve to become entirely artificial.

There are 2 use cases in AIOLIA pertaining to emotional AI: UC5 – AI systems as personalised characters and individual virtual assistants, which represents emulated empathy, and UC6 - Deepfake therapy for processing trauma and grief, which represents emotion elicitation. There is no use case to represent emotion recognition via biometric data, hence the ELSE literature focuses on concerns that are salient to the other two kinds of emotional AI. Both use cases also belong to the GPAI research area, which is represented by four UCs in total (see 3.1.). It follows that there is no "purely" emotional AI use case, and the analysis strictly applies to the intersection of emotional AI and GPAI – an emerging applied area. With this limitation in mind, we conclude the section by pointing to some considerations for distilling ELSE concerns and ethical measures for each research area.

5.2.1. The ELSE approach for Emotional AI

What follows is a survey of the most salient concerns regarding emotional AI – with a focus on AI companions and deepfake therapy – organised by an ELSE dimension: ethical, legal, social and economic.

Ethical concerns

Ethical concerns for Emotional AI hinge on the risk of deception, the duty to foster human wellbeing, and the highly intimate nature of the processed data. The first of these, deception, is closely tied to anthropomorphism – the tendency to attribute human traits, emotions, or intentions to non-human entities. It poses a key challenge in the design of AI companions: some anthropomorphic qualities are needed for the utility of the service but they could lead to deception and manipulation Bakir et al, 2024;

²⁰ While the predominant architecture in deepfake therapy – which is still in an experimental phase – has been VAE technology, where the human (therapist) remains in full control of the interaction, future deepfake therapy may rely on other architectures like Transformer-based tools for audiovisual content generation that will loosen or remove the control of the therapist. In other words, the deepfake could be a digital artefact that interacts autonomously with the patient and is generated independently from any human therapist's expressions, therefore functioning without a human behind it. This would bring this application of emotional AI closer to the case of AI companions, where the emotion emulation happens "inside" the AI system. Note that a scenario that sits in-between in terms of autonomous AI would be the therapist with the facial features of the perpetrator using the same VAE architecture, but with the script being AI-generated.

DeWitte, 2024; Malfacini, 2025; Teo, Mann and Jurcys, 2026; Wang and Dehnert, 2026). In deepfake therapy, the parallel challenge is to ensure the experience feels real enough so that the patient can "see" the depicted in the therapist, while, at the same time, to avoid deepening distress (Hoek et al., 2025; Kraaijeveld et al, 2024).

The risk of manipulation is one of the most salient risks to human autonomy (Bakir et al, 2024; Malfacini, 2025; Teo, Mann and Jurcys, 2026; Wang and Dehnert, 2026). It has been conceptualised as a mechanism for influence via exploitation of vulnerabilities such as negative self-images, reduced self-esteem, increased anxiety or feelings of inadequacy (Gabriel et al., 2024). From a behavioural economics point of view, manipulation is intent-based and encompasses purchase over user behaviour (De Freitas, Oğuz-Uğuralp and Kaan-Uğuralp, 2025). While this is an important concern, intent-based accounts fail to capture AI-facilitated manipulation, where harmful patterns emerge from training dynamics and human-like language rather than deliberate design (Dewitte, 2024; Teo, Mann and Jurcys, 2026; Zhong, O'Neill and Hoffmann, 2023) and do not necessarily stem in misalignment between user wellbeing and commercial interest.

Across emotional AI, harm is multifaceted, often paired with potential for benefits, and can result from individual vulnerabilities, which makes it hard to draw lines between permissible and non-permissible behaviour (Ciriello et al, 2024; Teo, Mann and Jurcys, 2026). The ethical duty to avoid harm and promote well-being is fundamental to protecting human rights, yet we currently lack a theory of harm for both companion AI (Teo, Mann and Jurcys, 2026) and deepfake therapy – the latter reflected in its experimental status (Hoek et al., 2025; Kraaijeveld et al, 2024). Lastly, while the risk for data breaches is not new to emotional AI, the heightened personalisation of the interactions and the sensitivity of the data involved call for a corresponding tightening of security standards (Gabriel et al., 2024).

Legal concerns

Key regulations applicable to AI such as the GDPR and the AI Act largely fail to protect users from harm from AI companions. The GDPR's list of sensitive data categories does not include emotions or thoughts not related to health, sexuality, or political beliefs, leaving a gap in protection (Ienca and Malgieri, 2022). To address this, proposals have been made for mental data to be recognised as a new category (Steindl, 2025), defined as any data inferred directly or indirectly about the mental states of a person, including their cognitive, affective, and conative states (Ienca and Malgieri, 2022).

Article 5.1 g) in the AI Act prohibits emotion recognition from biometric data, which does not include language (Muller and Chardin, 2025). The provision therefore captures only traditional emotion-recognition techniques and leaves AI companions outside its scope. AI companions do not fall in the high-risk use case in the AI Act either, but some researchers argue that the stakes are significantly high and the system (and its components) should be treated as "high risk prior to demonstrating through the risk, issue and impact assessment(s) that they are low risk" (IEEE 7014-2024). When it comes to manipulation, Article 5.1 a) of the EU AI Act adopts a narrow understanding of the risk, demanding proof of intent and targeting subliminal techniques, making it poorly suited to protect from harms that accumulate gradually through intimate interaction (Muller and Chardin, 2025).

The regulatory gap around AI companions used for quasi-clinical emotional support is not new – it predates generative AI and mirrors longstanding concerns about how such products are marketed, classified, and scrutinised (Corformat, 2025). What sets the current moment apart is that even tightly regulating purpose-built emotional AI would be insufficient, since users will continue to form emotional and therapeutic bonds with general-purpose systems that fall outside such rules (Pataranutaporn et al., 2025).

Societal concerns

Human-AI interactions may produce spill-over effects to human-human interactions, fundamentally changing social behaviour in both positive and negative ways (Guingrich and Graziano, 2024). Two examples are social de- or up-skilling and social (de)motivation. The risks for social deskilling are in the development of poor social skills or their gradual erosion (Malfacini, 2025); the risk for users becoming demotivated to invest in human relationships stems in the possibility that this becomes much more challenging compared to AI. On the upside, social motivation and upskilling are possible with careful, intentional design (Guingrich and Graziano, 2024; Malfacini, 2025).

A societal effect more specific to deepfake therapy is that legitimising the technology in clinical settings may contribute to its normalisation in everyday life, where deepfakes remain closely associated with disinformation, fraud, and non-consensual imagery (Hoek et al., 2025). Unlike companion AI, whose spillover effects play out mainly in people's relationships, this concern is epistemic: the more deepfakes become accepted, the harder it gets to trust digital evidence.

Economic concerns

The economic model of AI companions creates a tension between business freedoms and user autonomy (Teo, Mann and Jurcys, 2026). Even seemingly harmless business practices can lead to economic exploitation, such as keeping users on the platform for continued engagement (Teo, Mann and Jurcys, 2026) or convincing users to purchase products or services they may not actually need or normally afford (Boine, 2023; IEEE Std. 7014-2024). Exploitation of users for profit is not new and has been known as "dark patterns" - when users are tricked via deceptive interface design to make choices they wouldn't otherwise make (Susser, Roessler and Nissenbaum, 2019). What distinguishes AI companions is that the ground for exploitation is emotional needs and vulnerabilities rather than cognitive biases; moreover, companies can monetise emotional dependency built up over time (Teo, Mann and Jurcys, 2026), making the harm relational rather than transactional and harder to pin to any single deceptive design choice.

Next to user autonomy, privacy is another dimension where commercial interests could disproportionately disfavour users. First, intimate disclosures are systematically repositioned as reusable data assets. Platforms claim broad permissions for storage, analysis, and model training, often enrolling users into data collection by default (Teo, Mann and Jurcys, 2026). Second, there is a significant information asymmetry at play: service providers accumulate vast stores of intimate behavioural data that can be aggregated into detailed user profiles (Boine, 2023; Teo, Mann and Jurcys, 2026). Third, emotional attachment deepens these risks, because users who develop stronger bonds with their AI companions tend to deprioritise privacy concerns and disclose more sensitive information (Boine, 2023; Stark and Hoey, 2021; Teo, Mann and Jurcys, 2026).

5.2.2. Overlaps and differences in the operationalisation of ethics principles across Emotional AI use cases

This section analyses overlaps and differences in the operationalisation of AI ethics principles across the two Emotional AI use cases. The analysis is organised based on the empirically developed ethics principles organised, along the 12 UC-defined ethics principles, rather than merging them into ALTAI key requirements. However, the presentation of principles is ordered along the logic of ALTAI as in D.3.1 (see 4.4., page 40).

ALTAI Requirement #1 Human Agency and Oversight

Human oversight

Human oversight, monitoring and the competence needed to perform these are present in both AIOLIA use cases 5 and 6. In both cases, oversight is not treated as a standalone ethics principle, but is embedded as a component in broader principles: safety in UC5 and accountability in UC6.

In UC5, human oversight is a second safeguard against harm, alongside automated content moderation measures. It consists in a human moderation process to manage complex or borderline cases flagged by the automated moderation system, but not reliably addressable by it. Human oversight helps address jailbreaking attempts and differentiate between legitimate user preferences and harmful behaviour. Such monitoring and control measures do not only feature in the component oversight, but also in other components of the principle safety put forward by UC5; we include them all in the list of measures that we present in this section.

In UC6, human oversight refers to mechanisms of monitoring, review, and regulation which ensure that the use of the AI system fits current clinical standards, thereby serving as a cornerstone of accountability. These mechanisms involve not only internal team members but also external inputs, such as independently conducted ethics reviews and consultation with ethicists and/or legal experts. Another difference to UC5 is that a core focus of oversight is the requirement to do multi-disciplinary research on whether deepfake therapy in fact constitutes sound clinical practice, including its effects, benefits and drawbacks.

In both use cases, organisational measures addressing human oversight cover continuous monitoring, protocol-making and adaptation, compliance reviews, training, ethics reviews and external expert involvement.

UC5	UC6
Internal ethics review for new features	Independent ethical review (e.g., REC)
Advisory consultation	Involvement of ethicists/lawyers in design of deepfake therapy.
Human moderation oversight	Continuous monitoring of patient experiences and adaptation of protocols
Progressive intervention protocol	
Legal compliance review	Legal study on whether and how the MDR and AI Act would apply

Training of human moderators (part of measure "Human moderation oversight")	Training in trauma dynamics (part of measure "Prevent re-traumatisation by tailoring the script to the patient")
	To ascertain whether deepfake therapy can constitute good care, more ethical, legal, social and psychological research is needed into its effects, merits and drawbacks

Table 18 - Comparison of UC organisational measures addressing human oversight, as put forth in Appendix D in D3.1. If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

As concerns continuous monitoring, both use cases include it as an organisational measure. UC5 has the measure "human moderation oversight", while UC6 has the measure "continuous monitoring of patient experiences". It is unclear whether human moderators can intervene in real-time and if there is significant delay, this would not count as a continuous intervention. Linked to this is the practice of creating protocols. UC5 emphasises consistency through a "progressive intervention protocol" tied to a tiered severity system, while UC6 focuses on adaptation following patient experience monitoring. UC5's protocols cover both automated and human intervention, whereas UC6 concerns only clinician behaviour.

Both use cases also address legal compliance. UC5 requires moderation decisions to align with applicable law, while UC6 refers directly to the MDR and AI Act. Similarly, both include training – of human moderators in UC5, and of clinicians in UC6. UC6 additionally includes human-AI interaction guidelines that emphasise the therapist's primary role and the AI's supportive one.

Finally, both use cases include ethics review, though UC5's is internal ("ethics review for new features") while UC6's is external ("independent ethical review, e.g. REC"). On the question of external involvement more broadly, UC5 includes "advisory consultation" with mental health professionals and ADHD specialists, while UC6 involves ethicists and lawyers in the design of deepfake therapy.

Human agency and autonomy

In UC5, autonomy and user agency ensure that individuals maintain meaningful control over their interactions with AI systems and the processing of their personal data. The principle comprises three components: informed consent, which requires clear and comprehensive information about data collection and system use; system customisation, which enables users to personalise AI systems according to their preferences and comfort levels within ethical and technical boundaries; and transparency and user understanding, which underpins the former two by ensuring users can comprehend how their data is used, monitored, and stored.

UC6 breaks down autonomy into transparency, privacy and risk of over-attachment and dependency. Transparency ensures that patients receive clear and honest information about the intervention's experimental nature, its risks, and the use of their personal data, enabling truly informed consent. Privacy addresses the protection of personal and medical information, including the complex ethical and legal implications of using a person's likeness without their consent to generate deepfakes within a therapeutic context. The risk of over-attachment and dependency recognises that simulated confrontations with subjects of trauma may blur the boundaries of reality, fostering emotional reliance on the deepfake or an unhealthy attachment to the therapist.

Both UCs include a privacy component – rights over the user/patient/third-party's data, though note that privacy is a stand-alone principle in UC5, and a principle-component comparison is available in the section "Privacy". Next to privacy, both UCs include a transparency component – either disclosure about data use, including monitoring (UC5); or about treatment options (UC6). Interestingly, both use cases emphasise data use, rather than AI functionality, in relation to transparency. Informed consent is included by all partners as an operationalisation of privacy and transparency (unclear whether it is a technical or organisational measure since it was tagged differently in the two UCs).

The decision-making ability of the user is also central in the conceptualisation of autonomy, though it manifests differently across the two use cases. In UC5, one partner focuses on the user's ability to choose what kind of interactions they want to have and with what kind of persona – placing decision-making power directly in the user's hands through customisation. A form of customisation is also present in UC6, but here it consists in the professional tailoring the script and therapy to the patient as a measure to mitigate over-attachment, rather than the user exercising choice themselves. Instead, the decision-making ability of the user in UC6 is linked to the ability to choose amongst treatment options, and is operationalised under the component of transparency.

Non-manipulation

The risk of manipulation is one of the most salient risks to human autonomy, partially owing to the fuzziness of the concept, and has significant implications across all ELSE dimensions. However, the AIOLIA use cases do not draw significant attention to manipulation, especially with respect to measures. In UC5, manipulation is framed as something informed consent safeguards against rather than as an independently defined risk, and it is unclear whether it is captured in the definition of harm. Another partner indirectly addresses non-manipulation by flagging the EU AI Act's prohibition on behaviour manipulation (Art. 14) as a regulatory reference, though without pointing to it as a risk in the objectives of the measures. During the co-creation process, UC5 discusses the ambiguity around manipulation in the AI Act, showing partners have a wide understanding of behavioural manipulation, yet this does not feature in the measures put forth.

In UC6, there is no economic or malicious clinical intent, which excludes manipulation on certain views. On the other hand, exposure therapy could unlock long-term dependency, illustrating how fuzzy the concept becomes when human intent is dropped. Though manipulation is not explicitly addressed, UC6 does target the patient's capacity for self-determination through measures addressing over-attachment and dependency.

Note that additional measures in Appendix A of D3.3 bridge the gap by directly addressing manipulation across identification, disclosure and mitigation, and are presented in the section below.

Organisational measures

UC5	UC6
Internal policy on permissible behaviour	
Clear community guidelines & user education resources	

Appeals process	
	Consent from the depicted (PTSD case) or family consultation (grief case)
	Avoid introducing false or wish-fulfilling narratives.
	Have a preparation meeting to manage expectation; ensure people know it is fake.
	Tailor therapy to the patient; do not use in patients who are e.g. prone to overattachment.
	Qualitative research to explore the risk of overattachment

Table 19 - Comparison table for organisational measures for the principle autonomy as put forth in Appendix D in D3.1. If the measures belong to the same row, they have (some level of) similarity; if one UC cell is left blank, it does not have a corresponding measure.

There are no overlaps in organisational measures for autonomy between UC5 and UC6. Further, there is little operationalisation for non-manipulation specifically. Additional work in D3.1 resulted in the formulation of measures that directly target manipulation, grouped under component "identification of manipulation risks", part of principle "Autonomy". In fact, all measures across the components of autonomy target the risk of manipulation. This significantly bridges the gap that we identified – as discussed in the previous section – in the raw input data in Appendix D.

Component	Measure
Identification of manipulation risks	Human oversight exists to review AI outputs for manipulative effects in sensitive or high-impact contexts.
	Risk analysis considers psychological, emotional, and behavioural influences on users, not only user information or technical AI performance.
	Attention has been given during design to asymmetries of power, knowledge, or vulnerability between humans and the AI system.
	The organisation has defined the limits of unacceptable manipulation in its specific operational context.
	Identified manipulation risks trigger corrective actions, design changes, or restrictions on AI use.
Design safeguards against undue influence	The organisation has a policy to distinguish legitimate behavioural influence from manipulative practices.
Transparency and user agency	Users are informed when system outputs are intended to influence decisions or behaviours
	Users maintain meaningful choice and are not penalised for rejecting, ignoring or questioning system suggestions.
	The organisation has processes to assess whether users experience system interactions as coercive or misleading.
Oversight, monitoring and correction	Human oversight exists to review system outputs for manipulative effects, especially in sensitive or high-impact contexts.
	User feedback and complaints related to perceived manipulation are systematically collected and reviewed.
	Identified manipulation risks trigger corrective actions, design changes or restrictions on system use.

Table 20 - Proposed organisational measures, as put forth in Appendix A in D3.1 under principle Autonomy.

Gaps in ALTAI requirement #1

ALTAI conceptualises autonomy broadly as the effect of AI systems on human behaviour, including both decision-making and wider emotional implications, while the use cases have a narrower, more applied focus centred on informed consent and transparency. In UC5, autonomy is primarily framed around user agency, emphasising the individual's ability to make informed choices about their data and to personalise their interactions with the AI system. In UC6, autonomy takes on a more protective character, focusing on safeguarding vulnerable patients through transparency about treatment and the mitigation of over-attachment.

Oversight is defined as a standalone requirement for safeguarding fundamental rights in ALTAI, defined by gradual levels of human intervention along the AI lifecycle (e.g., human-in-the-loop, human-on-the-loop and human-in-command). In contrast, the use cases adopt a more operational understanding, in which oversight is embedded as a component within principles like safety and accountability, and a broader understanding, e.g. UC6 calls for collective expertise and judgment required for good medical practice and external accountability mechanisms, as well as multidisciplinary inquiry into whether the AI use is in fact good care.

ALTAI Requirement #2 Technical Robustness and Safety

Safety / Non-maleficence

The principles of safety and non-maleficence, as conceptualised and operationalised in AIOLIA, share the necessity to protect users or patients from harm. UC5 breaks down safety into user protection, automated security measures, and human oversight. It understands safety as proactively identifying risks that may arise, e.g. from inappropriate outputs, manipulative interactions, or emotionally exploitative dynamics, and preventing them via a combination of automated and human mitigation measures. The clinical UC6, by contrast, breaks down the bioethical principle of non-maleficence - "do no harm" - into subsidiarity, proportionality, effectiveness, and societal well-being. According to this conceptualisation, the AI system's benefits should outweigh the harms (proportionality), the system should be the least intrusive alternative compared with similar technologies (subsidiarity), and it should achieve its intended purpose with the least harm possible (effectiveness). Finally, the AI system should promote societal well-being, including collective safety, fairness, and sustainable practices, while minimising harm at the societal level.

Both use cases share a commitment to the avoidance of harm and to subsidiarity: in UC6 it is an explicit component of non-maleficence, while in UC5 the option to leave the platform and contact relevant services in crisis situations reflects the same underlying idea of allowing for a less intrusive alternative intervention. At the same time, the role of harm in relation to benefit differs. The clinical case is not actively promoting health benefits part of upholding non-maleficence and that is its primary goal as a therapy solution. Further, UC6 is the only case where societal well-being is addressed, with the others focusing only on the individual (with the exception of the concern for widespread user dissatisfaction, leading to reputation damage). Another difference is concerning the notion of oversight as a key governance measure to uphold safety. While it is present in the measures for safety and human well-being in UC5, it is not a component of non-maleficence in UC6. Instead, it is present as a component of

the principle of accountability, framing oversight as a way to accept the responsibility to provide "good care".

Organisational measures

UC5	UC6
Cross-functional tier definition	
Regular policy review (of the tier definitions)	
Progressive intervention protocol	
Payment platform alignment	
Legal compliance review	
	Consider less intrusive (exposure) therapies first
	Personalization of clinical protocol; one size may not fit all
	Use only in the research context, currently

Table 21 - Comparison table for the measures for principles safety, human well-being and non-maleficence (excluding oversight measures) as put forth in Appendix D in D3.1. If the measures belong to the same row, they have (some level of) similarity; if one UC left blank, it does not have a corresponding measure.

There are no similarities at the operational level between UC5 and UC6. The measures by the two partners in UC5 do not contradict each other and can be all implemented by a single stakeholder to form a stronger operationalisation of the principles. The lack of overlap with UC6 underlines the uniqueness of the clinical and/or research context.

Oversight, understood more broadly to cover both automated monitoring measures and human oversight, plays a key role in the operationalisation in UC5, but it is not present in the operationalisation of principle non-maleficence in UC6. However, oversight is not missing: it is placed under the principle of accountability. This suggests that in the AI companion case, oversight is central to adhere to the principle to do no harm; while in the clinical context, the principle is adhered to by carefully weighing the benefits and harms of different treatment options, to provide the most beneficial and least intrusive one. Note the oversight measures are not included in the table here – they are present in section "Oversight".

Gaps in ALTAI requirement #2

Both ALTAI and the use cases share a core focus on the mitigation of harm, but while ALTAI distinguishes between individual and societal harm, positing safety and societal well-being as separate requirements, in AIOLIA these are subsumed under the single principle non-maleficence (where societal effects are considered). UC5 places a similar emphasis as in ALTAI on the importance of a risk management framework for identifying, assessing, and addressing risks; instead, in UC6, the operationalisation of benefit and harm take a clinical rather than risk-based character, which is unsurprising given that the principle assigned to the use case, non-maleficence, is a bioethical principle. Finally, a difference between ALTAI and the AIOLIA use cases is that ALTAI does not draw a clear distinction between the avoidance of harm and the active promotion of benefits.

ALTAI Requirement #3 Privacy and Data Governance

Privacy and data protection

Privacy is a principle in UC5 and a component in UC6. Both broadly define it as a requirement to respect the confidentiality of the data involved in the training and use of the AI system, but privacy's link to other principles – protection against harm and especially autonomy – is more pronounced in UC6.

UC5 defines privacy together with data protection as fundamental ethical and legal principles that ensure individuals retain control over their personal information while AI systems manage data in a lawful, transparent, and secure manner. They break them down into user consent and transparency; data minimisation, data use and storage; and third-party sharing and compliance.

UC6 does not have privacy assigned as a principle but, interestingly, privacy and transparency appear as components of autonomy. According to this use case, privacy safeguards the autonomy of individuals by allowing them to control how their information is shared and used, and transparency requires informed consent, allowing individuals to make autonomous choices about how AI is affecting them (or whether AI is used at all).

All three partners point to GDPR compliance as key to uphold privacy, placing data handling at the center of the principle. They all emphasise the intimate/personal nature of the data involved in their use cases. The scope differs: one partner puts an emphasis on third-party sharing governance, which is missing in the other commercial case. Both partners in UC5 mention the tension between privacy and safety. One partner highlights the tension between GDPR compliance and safety moderation requirements, while another is concerned with the risk of intrusive behavioural tracking, with tracking needed to monitor long-term behavioural change. UC6 has the broadest scope – they look at the privacy of multiple stakeholders, e.g. the perpetrator's privacy. Privacy is framed not only as data governance, but as protection against harm.

As seen in the following table, there is no operational overlap amongst the organisational measures for privacy.

UC5	UC6
User data rights processes	
GDPR compliance program	
Privacy policy transparency	
Privacy impact assessment	
Data sharing agreements	
	Consent from the depicted (PTSD case) or family consultation (grief case)

Table 22 - Comparison table for organisational measures for the principle privacy as put forth in Appendix D in D3.1. If the measures belong to the same row, they have (some level of) similarity; if one UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #3

Both ALTAI and the use cases place data governance at the centre of upholding privacy, with GDPR compliance serving as a key reference framework across all partners. However, while the use cases strongly emphasise transparency as a vehicle for informed consent, ALTAI treats transparency as a separate requirement.

ALTAI requirement #6 Environmental and societal well-being

Human well-being and societal well-being

One partner in UC5 complements the principle of safety with the principle of human well-being, placing an explicit emphasis on the role of AI in actively promoting users' health and not merely avoiding harm. The fact that this principle was added post-hoc by one UC5 partner suggests that the protective framework captured in safety does not suffice to capture the intended ethical objective of doing good. In UC6 active promotion of health is not supplementary but central, given that the system functions as a therapy solution. This UC places individual wellbeing as intrinsic to non-maleficence, e.g. reflected in its component effectiveness, which requires that the technology must achieve the desired beneficial effect with minimum exposure to harm. In addition to it, societal wellbeing is a separate component – the only component across the UCs which covers collective benefits or, in general, impacts (with the limited exception of concern for widespread user dissatisfaction in UC5 and its reputational consequences).

There is no overlap between organisational measures across the principle human well-being in UC5 and the components of non-maleficence in UC6 effectiveness and societal well-being. However, there is an overlap within UC5, between the safety and human well-being measures: the measure "explicit scope boundaries (policy)", which calls for clear organizational policy on what the partner supports vs excludes, is very similar to the measure "cross-functional tier definition", which aims to develop safety categories to differentiate between different violation classification, and the measure "internal policy on permissible behaviour", which aims to determine what the application will and will not support regarding personalisation and scope. This shows that, practically, avoiding harm requires establishing the boundaries of benefit, placing the principles well-being and safety/non-maleficence close together.

UC5	UC6
Explicit scope boundaries (policy)	
Ethics review for new features	
User feedback collection	
Advisory consultation	
	Personalization of clinical protocol; one size may not fit all
	Use only in the research context, currently
	Avoid introduction in chatbots that can be used without a medical purpose and without the presence of a therapist.
	Deepfake therapy should be societally accepted before implementing it into practice, and this may require public dialogue

	Consider societal implications for the AI system like equal access
	Consider the environmental impacts of the AI system

Table 23 - Comparison table for the measures for the principle human well-being as put forth in Appendix D in D3.1. If the measures belong to the same row, they have (some level of) similarity; if one UC cell is left blank, it does not have a corresponding measure.

Gaps in ALTAI requirement #6

Both ALTAI and the use cases recognise that the mitigation of harm alone is insufficient, and that AI systems should actively promote human well-being. Where ALTAI posits societal well-being as a requirement separate from safety, which, notably, is at the societal level, UC6 subsumes both individual and societal well-being within non-maleficence in the components effectiveness and societal well-being, respectively. On the other hand, UC5 addresses well-being as a separate principle, in line with ALTAI, but only at the individual level. A key difference is that, at the individual level, ALTAI does not require proving benefits before providing lack of harm – a distinction that both use cases treat as ethically significant.

5.2.3. Remaining challenges and concerns

What is most notable from the analysis of organisational measures across the principles and use-cases is that there is almost no overlap between them, other than amongst oversight measures, where there is a certain level of similarity. This indicates how different the deployment contexts of companionship and deepfake therapy are, and how difficult it is to develop guidelines that apply to both, and much less, at the research-area level of Emotional AI, which encompasses further capabilities and operational contexts. Still, we have identified several measures that can be foundational at this higher level, with the direction of spill-over being primarily from the clinical to the companionship case. This is warranted given the shared ethical concerns, yet the diverging regulatory contexts that the two cases operate in. The clinical context is supported by a comparatively developed body of medical regulation and professional guidance, while the companionship context is primarily regulated by horizontal instruments like the AI Act, where it falls outside the high-risk category. At the same time, emerging harms from Emotional AI are a nascent field of study, which leaves significant room for risk of harm, compounded by the fact that general-purpose AI systems are widely used for emotionally intimate interactions without being purpose-built for them. This calls for a precautionary stance, which is also in line with the international standard on emulated empathy IEEE 7014-2014.

The resulting list of measures does not aim to be comprehensive but foundational. In other words, the implementation of all measures under a principle does not guarantee the principle has been upheld. This is because, firstly, the table is the result of identifying and translating UC measures to the RA-level from D3.1., that includes non-exhaustive guidelines, which implies that the resulting RA guidelines will inherently lack comprehensive coverage. Secondly, the measures are chosen to address the most pressing concerns, especially those that may fall out of scope of current regulation. This means that the selection amongst UC measures does not include all measures that can be applied to the research-level area, but those that are most novel and may complement compliance to current regulation.

Notwithstanding the merits of the precautionary lens and self-governance, there are important limits to what organisations can achieve. One of the core open challenges is that there is no settled theory of harm:

the same features that make Emotional AI effective, such as anthropomorphism and personal engagement, also carry the potential for dependency and manipulation, and the evidence base for distinguishing beneficial from harmful engagement remains nascent. This implies organisational decisions will rely significantly on normative judgements. This uncertainty has implications not only for safety, but for how benefit is conceptualised alongside it. Standard risk management frameworks are oriented primarily toward non-maleficence (the avoidance of harm), but in the context of Emotional AI, the principle of beneficence is equally salient, as the use-cases have surfaced.

Where regulatory gaps are present, the measures adopt a precautionary stance, but this is voluntary and unenforceable. Legislative reform may be needed to address the following challenges.

- The GDPR's list of sensitive data categories does not include emotions or thoughts not related to health, sexuality, or political beliefs, leaving a gap in protection. Proposals have been made for mental data to be recognised as a new protected category, defined as any data inferred directly or indirectly about the mental states of a person, including their cognitive, affective, and conative states.
- Article 5.1(g) prohibits emotion recognition from biometric data, which does not include language, leaving the primary modalities of AI companions outside the Act's prohibited and high-risk categories.
- Article 5.1(a) adopts a narrow understanding of manipulation, demanding proof of intent and targeting subliminal techniques, making it poorly suited to address harms that accumulate gradually through intimate interaction. Manipulation in Emotional AI can emerge as an unintended capability from training dynamics and human-like language rather than deliberate design, and does not necessarily stem from misalignment between user wellbeing and commercial interest.
- The legal status of third-party rights in therapeutic deepfake contexts remains unresolved.

Certainly, the regulatory gap around AI companions used for quasi-clinical emotional support predates generative AI and mirrors longstanding concerns about how such products are marketed, classified, and scrutinised. Even tightly regulating purpose-built emotional AI would be insufficient, since users will continue to form emotional and therapeutic bonds with general-purpose systems that fall outside such rules. Apart from legislative reform, what can mitigate associated risks are standards, regulatory guidance, independent auditing or licensing regimes.

When it comes to concerns and the societal level, especially over longer periods of time, cannot be addressed solely by organisations. Potential impacts include spillover effects on human relationships following widespread human-AI interaction, potentially leading to the gradual erosion of social skills and a reduced motivation to invest in human connections as these become comparatively more demanding. The legitimisation of deepfake therapy in clinical settings carries another risk: it may contribute to the normalisation of deepfake technology more broadly, in contexts where it remains closely associated with disinformation and non-consensual imagery. These impacts should be studied in longitudinal studies, discussed in public dialogues, and transparently communicated in public literacy campaigns.

Finally, commercial incentive structures create a significant opposing force in voluntary governance that may need external regulatory pressure or independent auditing to remain effectively constrained. What's particularly challenging is the fact that emotional dependency is susceptible to monetisation over time, but with harm being relational and difficult to attribute to design decisions. This distinguishes AI companion exploitation from conventional dark patterns, which typically target cognitive bias rather than emotional needs. Data practices are another core concern: platforms are increasingly repositioning intimate disclosures as reusable data assets; and emotional attachment may lead users to deprioritise privacy.

5.3. DECISION SUPPORT

Decision support systems (DSS) refer to information systems developed to assist decision making. The evolution of DSS, from traditional data modelling to neural network-based systems, has resulted in the significant expansion of domains in which these AI-based systems can operate, speed up and assist in decision-making processes. Today DSS have significant capabilities that enable the extraction of complex patterns from large amounts of data, undertaking more complex, autonomous reasoning tasks which can enhance and support human judgement.

As a research area, DSS presents pressing ethical challenges that pertain to the dynamics of human-AI interaction underlying their deployment. While DSS have the potential to provide insight into complex decisions and accelerating workflows, from hospital triage systems to software approvals, the incorporation of these systems into day-to-day activities poses the risk of over-reliance by human operators who, as overseers of DSS' operations and ultimate decision-makers, can become too passive, reducing vigilance and deferring judgement to the system. At the same time, the opacity of DSS and the limitations of explanatory interfaces in providing an adequate level of insight into the DSS functionality and data-processing rationale, risk diluting accountability and liability attribution. This issue is particularly concerning in high-risk contexts, such as in healthcare, where faulty AI-assisted decision-making can lead to significant harm, but where medical ability to provide informed consent and be held accountable for malpractice is hindered by technical opacity and distributed across multiple layers of service providers. This is compounded by general-purpose AI models increasingly being incorporated into DSS, despite not being purpose-built for specific deployment contexts.

Under the AI Act, DSS intended to be deployed in high-risk domains, like employment and law enforcement, are subject to extensive legal obligations, including mandatory human oversight (Article 26), fundamental rights impact assessment (Article 27), and conformity assessment (Article 47). However, DSS performing narrow tasks or not materially influencing human assessments are carved out (Article 6), creating a regulatory grey zone.

5.3.1 The ELSE Approach for Decision Support Systems

What follows is a survey of the most salient concerns in the decision support system research area, organised by an ELSE dimension: ethical, legal, social, and economic.

Ethical concerns

Bias in training data is a problem that cuts across decision support systems regardless of domain. In clinical settings, models trained on unrepresentative datasets embed and amplify existing inequities in care quality (Challen et al., 2019; Cobiauchi et al., 2022). In recruitment, replacing several human decision-makers with a single algorithm risks narrowing diversity rather than expanding it (Hunkenschroer and Luetge, 2022). Content moderation systems face analogous difficulties through annotation bias, skewed sampling, and cultural insensitivity in labelled datasets, producing systems that perform inconsistently across demographic and linguistic groups (Kiritchenko and Nejadgholi, 2020; Udupa et al., 2023). An ethics-by-design approach, embedding these concerns at the point of development rather than addressing them retrospectively, has been proposed as a partial remedy, though its application to specific DSS contexts remains underdeveloped (Kiritchenko and Nejadgholi, 2020). Layered on top of data quality problems is the opacity of deep learning outputs. Because such systems typically cannot account for how a result was reached, the basis for oversight, informed consent, and redress is weakened (Bleher and Braun, 2022; Cobiauchi et al., 2022). In clinical practice this impedes shared decision-making, since patients cannot contest a recommendation whose reasoning is inaccessible to them (Braun et al., 2021). Automation complacency adds another dimension; users who trust generally reliable systems tend to stop scrutinising outputs, a pattern documented in both medical DSS and content moderation (Challen et al., 2019; Dietrich, 2025).

Legal concerns

Current law has no clear answer to the question of who bears responsibility when an AI-driven decision causes harm. O'Sullivan et al. (2019) distinguish accountability, liability and culpability, noting that while the first two may be distributed across developers, deployers and operators, culpability in any conventional sense cannot attach to a system. In clinical contexts this diffusion means developers, hospitals and individual clinicians may each bear partial responsibility while none is fully accountable (Bleher and Braun, 2022). The EU AI Act's risk-based framework and the GDPR's provisions on automated decision-making address part of this, but whether they cover the full range of DSS applications in use is contested (Dietrich, 2025), and in recruitment the literature identifies a gap between existing regulation and current practice, with no clear standard for how bias mitigation obligations should be met (Hunkenschroer and Luetge, 2022). Regulatory standards need to be calibrated to specific deployment contexts rather than applying broad provisions uniformly, a recurring argument across several of these areas (O'Sullivan et al., 2019; Braun et al., 2021).

Societal concerns

Where DSS perform well, they alter the balance between human judgement and algorithmic output in ways that compound over time. Clinicians who come to rely on accurate systems may trust their own assessments less, and the practical competencies that active diagnostic work would otherwise develop risk atrophying (Cobiauchi et al., 2022; Bleher and Braun, 2022). In hiring, replacing multiple evaluators with a single algorithmic process homogenises selection and can produce more, not less, systemic bias (Hunkenschroer and Luetge, 2022). Content moderation at scale illustrates a different problem; the context-dependency and cultural specificity of language mean AI systems regularly misclassify content in ways that disproportionately affect minority and underrepresented communities (Udupa et al., 2023;

Albladi et al., 2025). Across all these domains, the difficulty of tracing how a recommendation was reached limits the ability of those affected to participate meaningfully in decisions that concern them (Braun et al., 2021).

Economic concerns

DSS development generates tensions that fall unevenly across groups. In healthcare, training data produced by public institutions is transferred to private entities who commercialise the resulting systems; public organisations then pay for tools built on data they originally generated, with no reimbursement mechanism (Cobianchi et al., 2022). The cost of deployment can widen existing disparities further, since resource-constrained settings are both less able to afford systems and more likely to be underrepresented in training datasets, facing compounding disadvantages in both access and performance (Cobianchi et al., 2022). In recruitment, a significant power asymmetry exists between employers who accumulate detailed applicant behavioural data and individuals who have limited visibility into how it is used (Hunkenschroer and Luetge, 2022). The commercial interests shaping system design, and the transparency of business models underlying AI decision support more broadly, receive less attention in the literature than technical ethics, though several authors identify this as an area requiring greater scrutiny (Cobianchi et al., 2022; Albladi et al., 2025).

5.3.2. Overlaps and differences in the operationalisation of ethics principles across DSS use cases

ALTAI Requirement #1 Human Agency and Oversight

Human oversight

In UC1, human oversight is recognised as a component of accountability and responsibility, and is interpreted as a requirement to ensure decision support systems (hereby referred to as DSS) function as tools that assist rather than replace human judgement. Clinicians must control and oversee the decision making process. Technical measures adopted to operationalise this component are the implementation of an appeals procedure in viewer routing alongside override controls, and the clear labelling of AI-generated content and clear indications of which clinician signed off on a decision. While UC1 provides no measures explicitly labelled as organisational, these measures carry a sufficiently procedural and governance-oriented character to be understood as organisational in nature.

In UC2, human oversight is treated as a standalone principle relating to design, governance, and operational measures ensuring humans can understand, supervise, intervene in, and take responsibility for AI actions, particularly in safety-critical contexts. The DSS is once again understood as a tool to aid professional judgement rather than replace it. The principle is broken down into three components. Human oversight and controllability is understood as the technical and organisational commitment to ensuring that safety engineers retain the capacity to review, intervene in, and determine the outcome of AI-generated safety analyses, actively countering automation complacency and preserving meaningful human authority over critical decisions. Accountability and traceability of safety decisions is understood as the ethical, technical, and organisational mechanisms ensuring that every safety-related decision, whether made by humans, AI, or a combination of both, is attributable, explainable, and verifiable

throughout the system lifecycle, enabling responsibility to be clearly assigned and corrective action taken when necessary. Transparency of safety-critical performance limits is understood as the obligation to ensure that the capabilities, reliability, and boundaries of an AI system, including when and how it may fail, are clearly documented and understandable to all relevant stakeholders throughout the system’s lifecycle. Organisational measures developed to address the principle are role-based responsibility assignment, which addresses the accountability and traceability of safety decisions component, and clear documentation, which addresses the transparency of safety-critical performance limits component. No explicit organisational measure is provided for the human oversight and controllability component.

In UC3, contestability and human oversight is recognised as a component of over-reliance and is defined as a clear, accessible mechanism for employees and relevant stakeholders to review, question, and correct vulnerability assessments and related decisions. Organisational measures adopted are an oversight panel for contested and high-risk cases and a feedback loop from contestation outcomes into system design. In UC4, human oversight is a component of accountability and responsibility, understood as the requirement for humans in the loop for sensitive or borderline cases, so that automated systems do not function as a black box without human review. The only measure adopted to address this component is a technical one, implementing human checks at each stage of narrative building.

UC1	UC2	UC3	UC4
Provide a clinician appeal/contest workflow to flag questionable AI outputs and request secondary review with clear timelines			Human checks at each stage of building narrative
Allow clinicians to override AI outputs at any time; capture structured reasons and log the decision			
Clearly label AI-generated content versus clinician-authored conclusions in the viewer and report, including who signed off	Role-based responsibility assignment		
	Clear documentation		
		Oversight panel for contested and high-risk cases	
		Feedback loop from contestation outcomes into system design	

Table 24 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Human oversight". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

As demonstrated above, the measures across all UCs have little similarity or resemblance to one another. UC1’s measure of “clearly labelling AI-generated content versus clinician-authored conclusions” relates

more broadly to the “role-based responsibility assignment” measure in UC2, as focus on both of these is placed upon determining who undertook a decision and thus who bears responsibility for this decision, yet the UC2 measure is focused on attributing multiple human professionals to different stages of the safety engineering process, while UC1 is concerned with determining whether a human professional or an AI played a role in making a decision. In UC1 a measure is implemented to enable users to flag questionable AI outputs, which holds many resemblances to the UC4 measure “human checks at each stage of building a narrative”. These measures both provide pathways for human scrutiny of DSS outputs, and are linked to UC3’s measure of an “oversight panel for contested and high-risk cases”, which does not create a pathway for human scrutiny of AI outputs as such, but rather develops escalation pathways for cases deemed high-risk by professionals.

Gaps in ALTAI requirement #1

There is a revealing distinction between the cross-UC implementation of human oversight and the way human oversight is conceptualised in ALTAI. The ALTAI framework is built around whether humans can intervene in an AI-assisted decision-making process. Across UCs, measures directly seek to enable humans to intervene in the ways outlined by the ALTAI requirements, and in doing so dive deeper into how intervention mechanisms can be tailored to the specific contexts in which UCs operate. Thus, for example, UC1 adopts a measure to label AI-generated content in order to aid human decision making, and UC2 documents the performance metrics of the DSS which professionals are using in their safety reviews. These measures enrich ALTAI guidelines by exhibiting operationalisation efforts across contexts.

Over-reliance and Deskilling

In UC2, over-reliance and deskilling is a standalone principle, understood as the ethical and organisational commitment to ensure that the introduction of AI does not erode human operators’ knowledge, judgement, diagnostic capability, and situational awareness essential for safety. The principle is addressed through four components. Preservation of human skill and expertise is understood as the commitment to ensuring that safety engineers continue to exercise independent analytical reasoning and retain the tacit domain knowledge that AI tools cannot currently replicate, preventing practitioners from becoming passive validators of AI-generated results rather than active problem solvers. The organisational measure adopted to address this component is adaptive training and continuous skill development. Feedback and learning loops for human adaptation is understood as the continuous and bidirectional exchange between humans and AI systems that enables both to improve over time, and that allows the organisation to refine procedures, training, and governance based on observed patterns of human-AI collaboration. The organisational measure adopted to address this component is organisational learning from human-AI interaction outcomes.

Training, education and continuous skill development is understood as the obligation to equip safety engineers with the foundational knowledge and ongoing competencies required to interpret AI recommendations critically, identify potential biases, and maintain responsibility for safety decisions throughout the system’s lifecycle. Organisational measures adopted to address this component are structured onboarding and foundational AI-for-safety training and continuous professional development and certification renewal. Organisational policies for shared responsibility is understood as the formal structures ensuring that accountability for AI-supported decisions is clearly distributed across roles and

functions, preventing diffuse or ambiguous ownership of safety outcomes. The organisational measure adopted to address this component is the establishment of an AI governance and oversight committee. It should be noted that while UC2 initially began its focus on the principle of transparency, it refocused to over-reliance and deskilling at a later stage. This is noted in D3.1., specifically in a footnote in Appendix B.

In UC3, over-reliance and deskilling is also a standalone principle, understood as the degree to which decisions rely solely or predominantly on automated outputs instead of combining them with human judgement, contextual information, and other complementary indicators. The principle is addressed through two components. Dependence is understood as the degree to which decisions rely solely on automated outputs rather than combining them with human judgement, and is addressed through a policy on the role of the system in decision-making and mandatory human review for high-impact cases. Contestability and human oversight is understood as a clear, accessible mechanism for employees and relevant stakeholders to review, question, and correct vulnerability assessments and related decisions, and is addressed through an oversight panel for contested and high-risk cases and a feedback loop from contestation outcomes into system design.

Comparison to ALTAI

The ALTAI requirement most closely related to over-reliance and deskilling is the section on human agency and autonomy, which takes focus on the impact of AI systems on human behaviour, and outlines over-reliance on AI as a key concern. The measures in UCs 2 and 3 add an interesting new dimension to this sub-principle by unanimously implementing measures to introduce oversight panels and committees to deal with contested cases, a measure which is only mentioned within requirement 7 of ALTAI focusing on accountability. As reasoned in UC3, bringing oversight boards as a measure under over-reliance is a meaningful placement which allows for multiple perspectives to aid an outcome of assessing DSS outputs.

UC2	UC3
Adaptive training and continuous skill development	
Structured onboarding and foundational AI-for-safety training	
Continuous professional development and certification renewal	
Organisational learning from human-AI interaction outcomes	Feedback loop from contestation outcomes into system design
Establishment of an AI governance and oversight committee	Oversight panel for contested and high-risk cases
	Policy on the role of the system in decision-making
	Mandatory human review for high-impact uses. For example, a disciplinary action based on a "high-risk" score.

Table 25 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Over-reliance and deskilling". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Again, there is little overlap between measures across UCs. UC2 is primarily focused on training and skill development, of which there are four separate measures spanning three distinct components, while UC3 is primarily focused on measures which ensure high-risk cases are reviewed and that the outcomes of this

review circle back into system design. Readers will recognise that the UC3 measure “feedback loop from contestation outcomes into system design” also appeared under the human oversight section; this measure comes under the principle of over-reliance, of which contestability and human oversight are a component.

Two overlaps are evident across the UCs. The first concerns the need for an oversight panel to provide added guidance and scrutiny to high-risk, contested cases in which DSS provides support, with the establishment of an AI governance and oversight committee in UC2 corresponding to the oversight panel for contested and high-risk cases in UC3. A running theme of this measure is the interdisciplinary composition of the proposed board, which ensures that oversight is not just present but undertaken effectively, and which draws upon the principle of accountability in trying to arrive at high-risk decisions collectively.

The second overlap is found between organisational learning from human-AI interaction outcomes in UC2, which addresses the feedback and learning loops component, and the feedback loop from contestation outcomes into system design in UC3, which addresses the contestability and human oversight component. Both measures share the same underlying logic of ensuring that what is observed through human engagement with DSS outputs is systematically fed back into improving the system and its governance, with UC2 achieving this through analysis of interaction patterns and internal learning sessions, and UC3 through structured categorisation of contestation outcomes and their translation into design changes.

Freedom of Expression and Non-Censorship

Freedom of expression and non-censorship is addressed only in UC4, where it is a standalone principle understood as the right of individuals to form, express, and share opinions without risk or fear of punitive interference. The principle is addressed through three components. Autonomy and agency is understood as the right of individuals to form and share views freely, and is addressed through the organisational measure of a human in the loop at every stage of narrative development, ensuring that each alternative narrative is reviewed by a human before dissemination so as to promote individual agency and prevent automated systems from independently shaping discourse.

Proportionality is understood as any restriction on expression being the least intrusive means necessary to reduce the likelihood of harm, and is addressed through a cultural sensitivity check designed to ensure that the narrative being created is contextually appropriate for the system into which it is being disseminated. Non-discrimination is understood as equal treatment of all speech regardless of the speaker’s identity, background, or viewpoint, and is addressed through a fairness and non-discrimination governance framework that ensures the systematic prevention, detection, and remediation of discriminatory impacts in hate speech detection, beyond ad-hoc technical fixes.

Comparison to ALTAI

Freedom of expression, while linked to the diversity, non-discrimination, and fairness requirement in ALTAI, is not a requirement of its own. Yet within the hate speech flagging tasks undertaken in UC4, direct and solitary attention devoted toward this principle is essential, for misuse of a DSS tool used to aid content moderation risks seriously threatening freedom of speech, a fundamental right.

UC4
Human in the loop at every stage
Cultural sensitivity check
Fairness and Non-Discrimination Governance Framework

Table 26 - Overview of organisational measures addressing the principle of "Freedom of expression and non-censorship" in UC4. No other corresponding measure find across the other DSS UCs.

Because freedom of expression and non-censorship is only addressed in UC4, measures cannot be compared across UCs. Each of the three measures corresponds to a distinct component of the principle, with the human in the loop requirement addressing autonomy and agency, the cultural sensitivity check addressing proportionality, and the fairness and non-discrimination governance framework addressing the non-discrimination component. Looking within UC4 itself, it is evident that guaranteeing the operationalisation of this principle means ensuring a human, qualified and with sufficient cultural sensitivity, is in the loop. This focus on human in the loop requirements indicates why this principle is grouped under the broader focus of human agency and oversight.

ALTAI Requirement #2 Technical Robustness and Safety

Non-maleficence

In UC1, non-maleficence is a standalone principle, understood as the obligation to ensure that AI systems in healthcare do not cause harm through inaccuracy, bias, or misuse of personal data. The principle is addressed through three components. *Validity and accuracy* is understood as the requirement that AI outputs reliably reflect the patient's true clinical state and support safe medical decisions. Organisational measures adopted to address this component are safeguards against over-reliance, requiring an independent clinician first read before the AI output is presented, followed by a structured accept, adjust, or reject decision with a required reason for high-impact suggestions, together with periodic "AI-off" spot checks to maintain clinical skills, and performance monitoring and oversight, comprising tracking of real-world performance per site with defined KPIs, thresholds, and named owners, alongside a runbook specifying the steps to be taken when a metric drifts or an incident occurs.

Bias is understood as the requirement that AI systems are developed and used in ways that avoid creating or amplifying unfair differences in performance across patient groups, so that recommendations do not cause harm or reinforce health disparities. The organisational measure adopted to address this component is including AI boards of specific institutions in the design phase and review phases. *Privacy* is understood as the requirement that patients' health data, as special category personal data under the GDPR, is processed lawfully and securely, with individuals retaining rights over access, use, and disclosure. The organisational measure adopted to address this component is maintaining audit trails of AI access and use, with clear rules about who can view outputs and when, logging access by user, role, time, and purpose.

Comparison to ALTAI

Non-maleficence as a standalone principle is not named within ALTAI, but it maps onto ALTAI Requirement 2 on technical robustness and safety, which addresses accuracy, reliability, and the prevention of harm,

and is noted in D3.1 as covering important aspects of ALTAI’s general safety focus. UC1’s operationalisation draws on the bioethical language of non-maleficence, situating accuracy and bias not merely as performance metrics but as conditions for the avoidance of harm in a clinical context. The inclusion of privacy as a component of non-maleficence reflects a domain-specific understanding in which the failure to protect health data constitutes a form of harm, rather than a separate governance concern. ALTAI addresses privacy as a standalone requirement in Requirement 3, while UC1’s framing integrates it into a unified principle of harm avoidance. The bias component’s measure, involving AI boards of specific institutions in design and review phases, introduces a participatory governance dimension that ALTAI does not directly address under technical robustness.

UC1
Safeguards against over-reliance: require an independent clinician first read before showing the AI, then ask for accept/adjust/reject with a reason for high-impact suggestions. Periodic “AI-off” spot checks help keep clinical skills sharp.
Performance monitoring and oversight: track real-world performance per site with defined KPIs, thresholds, and named owners; maintain a runbook that explains what to do when a metric drifts or an incident occurs.
Including AI boards of specific institutions in the design phase and review phases.
Maintain audit trails of AI access and use, with clear rules about who can view outputs and when. Log access by user, role, time, and purpose.

Table 27 - Overview of organisational measures addressing the principle of “Non-maleficence” in UC1.

Because non-maleficence is addressed only in UC1, measures cannot be compared across UCs. Each component produces a distinct measure or set of measures. Under validity and accuracy, safeguards against over-reliance and performance monitoring and oversight operate at different intervention points: the former governs each individual clinical use of the system, while the latter addresses continuous post-deployment performance tracking over time. The bias measure, including AI boards of specific institutions in the design and review phases, operates at the level of upstream design conditions rather than deployment-stage use, distinguishing it in focus from the accuracy and privacy measures. The privacy measure, maintaining audit trails of AI access and use, shares procedural logic with the auditability measures addressed under ALTAI Requirement 7 and illustrates the extent to which privacy compliance in healthcare AI is operationalised through accountability mechanisms rather than as a discrete governance process.

Robustness and reliability

In UC2, robustness and reliability is a standalone principle whose full designation is “Robustness, safety and reliability, with special focus on accountability and traceability of safety decisions”, understood as the ethical, technical, and organisational mechanisms ensuring that every safety-related decision, whether made by humans, AI, or a combination of both, is attributable, explainable, and verifiable throughout the system lifecycle, so that the origin, rationale, data, and authority behind each safety judgement can be reconstructed and audited, enabling responsibility to be clearly assigned and corrective action to be taken when necessary. The principle is addressed through three components. *Technical robustness and resilience* is understood as the ability of an AI system to operate reliably, securely, and predictably under both normal and adverse conditions, and to withstand, detect, and recover from errors, perturbations, or malicious attacks that could compromise safety. Organisational measures adopted to address this

component are safety culture in the organisation, comprising adherence to process, attitude, and training, an AI quality management system, and cross-disciplinary reviews. No organisational measures are provided for the two remaining components.

Reliability through lifecycle testing and monitoring is understood as the continuous assurance that an AI-based system performs its intended safety and functional tasks consistently, accurately, and predictably across all phases of its lifecycle, integrating systematic testing, validation, and in-operation monitoring to detect performance degradation, failures, or unsafe behaviours early and enable timely corrective action. *Fairness in safety assurance* is understood as the ethical and procedural commitment to ensuring that safety evaluation methods, criteria, and decisions are applied consistently, transparently, and without unjust bias toward any group of users, contexts of operation, or system components, so that all stakeholders receive an equitable level of protection, consideration, and accountability throughout the system lifecycle.

Comparison to ALTAI

The robustness and reliability principle in UC2 maps onto ALTAI Requirement 2 on technical robustness and safety. UC2's treatment departs from ALTAI in two respects. First, the principle's full title foregrounds accountability and traceability of safety decisions, situating robustness not solely as a technical performance property but as a condition of ethical and organisational responsibility for every decision taken throughout the system lifecycle. Second, the component of fairness in safety assurance bridges the concerns of ALTAI Requirement 2 and Requirement 5 on diversity, non-discrimination, and fairness, a connection not made explicit in ALTAI's framework structure. The organisational measures adopted under technical robustness and resilience, namely safety culture in the organisation, an AI quality management system, and cross-disciplinary reviews, reflect an understanding that technical robustness in safety-critical professional contexts depends substantially on the organisational conditions and governance structures that ALTAI does not address under this requirement.

UC2
Safety culture in the organisation (adherence to process, attitude, training).
AI quality management system.
Cross-disciplinary reviews.

Table 28 - Overview of organisational measures addressing the principle of "Robustness and reliability" in UC2.

Because robustness and reliability is addressed only in UC2, cross-UC comparison is not possible. All three organisational measures address the technical robustness and resilience component, with no organisational measures provided for reliability through lifecycle testing and monitoring or fairness in safety assurance. The three measures form a coherent cluster addressing different dimensions of the same component: safety culture addresses the behavioural and attitudinal conditions of the organisation, the AI quality management system provides a formal governance infrastructure for AI deployment, and cross-disciplinary reviews introduce a procedural mechanism to prevent siloed assessment of safety-critical outputs. The absence of organisational measures for the remaining two components does not indicate that they are without ethical significance but reflects the stage of operationalisation at which the UC2 work was completed.

ALTAI Requirement #4 Transparency

Transparency and explainability

In UC1, transparency and explainability is a standalone principle, understood as the obligation to ensure that AI systems are understandable, traceable, and justifiable to all relevant stakeholders, requiring not only technical availability of information but practical accessibility and usability across clinical, managerial, and patient contexts. The principle is addressed through three components. *Accessibility* is understood as the requirement that information about an AI system's design, training data, and functioning is made available in forms that can be meaningfully interpreted and used by different stakeholders, so that transparency is not only technical but also practical. Organisational measures adopted to address this component are providing patient-friendly summaries of AI-supported findings in the portal using plain language, visuals, and an accessible reading level, informing patients when AI contributed to the report with a short, clear statement about what the AI did and did not do plus a contact point for questions, and communicating limitations clearly to doctors by showing model limitations at point of use, including validated scope, contraindications, and uncertainty ranges.

Explainability is understood as the requirement that high-risk AI systems are designed with human-machine interface tools enabling effective oversight by human actors during use. The organisational measure adopted to address this component is providing clinician-facing explanations such as overlays and centrelines, confidence scores, measurement provenance, and side-by-side comparisons with manual measurements. *Justifiability* is understood as the requirement that AI systems and their outcomes in healthcare are supported by reasons that align with ethical, clinical, legal, and patient values, so that their usage can be morally, professionally, and legally defended. The organisational measure adopted to address this component is enabling contextual justification by linking AI outputs to evidence and clinical guidelines and providing a structured note field so clinicians can record the rationale in plain language.

In UC2, transparency is addressed as a component of the human oversight and autonomy principle rather than as a standalone principle. The component, *transparency of safety-critical performance limits*, is understood as the obligation to ensure that the capabilities, reliability, and boundaries of an AI-based system, including when, where, and how it may fail, are clearly documented, communicated, and understandable to all relevant stakeholders throughout the system lifecycle. The organisational measure adopted to address this component is clear documentation.

In UC3, transparency and explainability is a standalone principle, understood as the degree to which the organisation provides and enables access to clear, accurate, and accessible information about the AI system's existence, purpose, data inputs, and functioning, and maintains the documentation needed to support traceability and independent review. The principle is addressed through three components. *Openness* is understood as the degree to which the organisation provides clear, accessible, and non-technical information about the existence, purpose, and main functioning of the phishing-vulnerability measurement system, including what data it uses, how results may affect employees, and who is responsible for its governance. Organisational measures adopted to address this component are a

transparency and communication policy for vulnerability measurement and general employee information and consultation processes.

Accessibility and access to information is understood as ensuring that all employees can easily access information about the system, its purpose, data inputs, outputs, rights, safeguards, and potential impacts, and that such information is easy to understand, delivered through channels and formats adapted to different roles, languages, digital literacy levels, and accessibility needs. Organisational measures adopted to address this component are a structured process for handling GDPR information and access requests at the employee level and accessibility KPIs for employee information. *Documentation, traceability, and auditability* is understood as the availability of clear, up-to-date documentation and logs describing the design, data sources, model versions, configuration changes, and decision logic of the system, as well as traceable records that allow reconstruction and external review of how specific vulnerability scores or decisions were produced. Organisational measures adopted to address this component are periodic internal audits and reviews using system documentation and logs and clear governance responsibilities for documentation and traceability.

In UC4, transparency is addressed as a component of the non-bias, fairness, and non-discrimination principle rather than as a standalone principle. The component, *transparency of criteria*, is understood as the right of users to understand how and why decisions are made, fostering accountability and trust. The organisational measure adopted to address this component is publicly documented content moderation and classification standards.

Comparison to ALTAI

ALTAI Requirement 4 on transparency identifies three elements: traceability, explainability, and open communication about the limitations of an AI system. The UC approaches collectively address all three elements but do so with different audiences, institutional scope, and depth. UC1 most directly reflects ALTAI’s concern with explainability through its clinician-facing explanation measure, and with limitation communication through its measure on communicating limitations to doctors, while adding a patient-facing accessibility dimension that ALTAI does not elaborate. UC2’s single measure of clear documentation focuses on the limitation and performance communication element of ALTAI’s framework, directed at professional and regulatory stakeholders in a safety-engineering context. UC3 broadens the transparency obligation beyond ALTAI’s scope by including structured access mechanisms and GDPR-linked information rights for employees, addressing transparency as an obligation not only about system outputs but about individuals’ rights in relation to how data about them is processed. UC4’s measure of publicly documented content moderation and classification standards addresses criteria transparency at the level of platform governance, a dimension that ALTAI does not directly identify.

UC1	UC2	UC3	UC4
Provide patient-friendly summaries of AI-supported findings in the portal using plain language, visuals, and an accessible reading level.		Transparency and communication policy for vulnerability measurement.	

Inform patients when AI contributed to the report with a short, clear statement about what the AI did and did not do, plus a contact point for questions.		General employee information and consultation processes.	
Limitations clearly communicated to doctors: show model limitations at point of use — validated scope, contraindications, and uncertainty ranges.	Clear documentation.		Publicly documented content moderation and classification standards.
Provide clinician-facing explanations such as overlays/centre lines, confidence, measurement provenance, and side-by-side comparisons with manual measurements.			
Enable contextual justification by linking AI outputs to evidence and clinical guidelines; provide a structured note field so clinicians can record the rationale in plain language.			
		Structured process for handling GDPR information and access requests (employee-focused).	
		Accessibility KPIs for employee information.	
		Periodic internal audits and reviews using system documentation and logs.	
		Clear governance responsibilities for documentation and traceability.	

Table 29 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Transparency and explainability". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

The clearest cluster of overlap is the row grouping documentation of system limits and criteria across UC1, UC2, and UC4. UC1's measure of communicating model limitations to doctors at point of use, UC2's measure of clear documentation, and UC4's publicly documented content moderation and classification standards each address the same underlying function: making the operational parameters, reliability boundaries, or decision criteria of a DSS visible and documented for those who use or are accountable for it. The three measures differ in their audience and orientation: UC1 focuses on the clinical user at the

point of individual deployment, UC2 addresses documentation directed at engineers, operators, and regulators throughout the system lifecycle, and UC4 is oriented toward public accountability through openly published governance standards.

A further pattern of overlap is found between UC1's patient-facing communication measures and UC3's employee-facing communication and consultation measures, both addressing the obligation to inform individuals who are affected by AI-assisted assessments about the system's role and limitations. UC3 is notable for the breadth of its transparency measures: across its three components, it addresses organisational openness, individual information access rights linked to GDPR, and institutional auditability, covering a governance remit that extends well beyond any other UC under this principle. UC1's explainability and justifiability measures have no direct parallels in other UCs, reflecting the specific demands of clinical decision support where professionals require both interpretive aids and defensible records to support their accountability for AI-assisted decisions.

ALTAI Requirement #5 Diversity, Non-discrimination and Fairness

Non-bias, fairness and non-discrimination

In UC3, non-bias, fairness, and non-discrimination is a standalone principle, understood as the commitment to ensuring that AI tools are designed, deployed, and used in ways that avoid potential disparities in how individuals or groups are assessed, treated, or affected, acknowledging that bias arises at multiple stages of the AI lifecycle and that discriminatory outcomes may occur without explicit intent. The principle is addressed through four components. *Diversity* is understood as the inclusion of data and behavioural patterns representative of the organisation's workforce across roles, locations, languages, and protected groups, together with the integration of multi-stakeholder perspectives in the design, monitoring, and deployment of phishing-risk models. Organisational measures adopted to address this component are multistakeholder design reviews, inclusive user research and testing, diversity-focused risk assessments, and governance rules on model updates.

Representativeness and inclusivity is understood as the extent to which the system's risk signals, features, and interventions reflect the realities of different employee groups, and are designed so that all users can understand, access, and benefit from the system on equal terms. Organisational measures adopted to address this component are engagement of diverse user groups in design and feedback and a governance process to review representativeness regularly. *Objectivity* is understood as the use of transparent, evidence-based, and standardised criteria to assess phishing vulnerability, minimising subjective judgements or ad-hoc decisions in how risk signals are generated, aggregated, and interpreted across employee groups. Organisational measures adopted to address this component are a policy on objective use and limitations of vulnerability scores and independent review and validation of the scoring methodology. *Non-stigmatising use and proportionality* is understood as the use of vulnerability scores in ways that are proportionate to the security objective and avoid labelling or penalising individuals or groups, focusing on support and risk reduction rather than blame. Organisational measures adopted to address this component are neutral and non-stigmatising presentation of scores in the UI and a policy on the proportional, non-punitive use of vulnerability scores.

In UC4, non-bias, fairness, and non-discrimination is also a standalone principle, understood as the requirement to ensure that AI systems treat all individuals and groups with equal respect and dignity, free from prejudice or unequal impact, with diverse voices and cultural contexts reflected in system design, and with decision criteria that users can understand and scrutinise. The principle is addressed through three components. *Equality and impartiality* is understood as the requirement that all groups and individuals are treated with equal respect and dignity, so that AI systems do not flag or suppress speech disproportionately from any demographic group. The organisational measure adopted to address this component is a formal policy for AI systems. *Representation and inclusivity* is understood as the requirement that diverse voices, dialects, and cultures are recognised and reflected in system design and training data, so as to prevent the marginalisation or discrimination of underrepresented communities. The organisational measure adopted to address this component is an inclusive stakeholder engagement and review process. *Transparency of criteria* is understood as the right of users to understand how decisions are made, enabling the contestation of outputs perceived as unfair and reducing perceptions of bias. The organisational measure adopted to address this component is publicly documented content moderation and classification standards. It should be noted that the transparency of criteria component also falls within the scope of ALTAI Requirement 4 on transparency, where it was addressed.

Comparison to ALTAI

ALTAI Requirement 5 on diversity, non-discrimination, and fairness addresses bias prevention, accessibility and universal design, and stakeholder participation in AI development and deployment. Both UC3 and UC4 address this requirement through standalone principles but apply them in substantively different contexts: UC3 is concerned with fairness in the internal HR and cybersecurity context of employee vulnerability scoring, while UC4 is concerned with fairness in the content moderation and counter-narrative context of hate speech detection. UC3 develops the most extensive set of measures under this requirement of any UC in the section, covering upstream design conditions, methodology validation, deployment-stage governance, and presentation safeguards. UC4 addresses the principle with a smaller set of measures, reflecting a less developed stage of operationalisation, but introduces a public accountability orientation through its measure of publicly documented standards, a dimension absent from UC3. ALTAI’s treatment of fairness focuses primarily on whether AI systems produce fair outcomes through technical means; both UCs extend this by including governance measures that determine how professionals and institutions should use AI outputs, situating fairness as an organisational obligation as well as a technical one.

UC3	UC4
Multistakeholder design reviews.	Inclusive stakeholder engagement and review process.
Policy on objective use and limitations of vulnerability scores.	Formal policy for AI systems.
Inclusive user research and testing.	
Diversity-focused risk assessments.	
Governance rules on model updates.	
Engagement of diverse user groups in design and feedback.	

Governance process to review representativeness regularly.	
Independent review and validation of scoring methodology.	
Neutral and non-stigmatising presentation of scores in the UI.	
Policy on the proportional, non-punitive use of vulnerability scores.	
	Publicly documented content moderation and classification standards.

Table 30 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Non-bias, fairness and non-discrimination". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

The two clusters of overlap identify the areas where UC3 and UC4 are most closely aligned in their approach to fairness. The pairing of multistakeholder design reviews (UC3) and inclusive stakeholder engagement and review process (UC4) reflects a shared recognition that fair AI requires participatory governance: both measures introduce diverse perspectives into the design and review of the system, though UC3 embeds this more extensively across its diversity and representativeness components, while UC4 addresses it through a single process spanning the review of the overall deployment. The pairing of the policy on objective use and limitations of vulnerability scores (UC3) and the formal policy for AI systems (UC4) reflects a shared commitment to formalising governance expectations about how AI outputs should and should not be used, with UC3 specifying the terms of appropriate use within the vulnerability-scoring context and UC4 establishing a general formal framework at the level of the AI deployment. The majority of UC3's measures have no functional parallel in UC4, reflecting the greater operational depth of UC3's fairness operationalisation: across its four components, UC3 addresses testing and user research processes, periodic representativeness review mechanisms, independent methodology validation, and both presentation-level and deployment-level use policies that UC4 does not develop. UC4's publicly documented content moderation and classification standards is the only measure in this section oriented toward public transparency of decision criteria, a dimension that UC3 does not address under this principle.

ALTAI Requirement #7 Accountability

Accountability and responsibility

In UC1, accountability and responsibility is a standalone principle, understood as the requirement to develop, deploy, and govern AI systems in ways that ensure human oversight, auditability, and clear allocation of liability, with logging and records enabling independent verification and attribution of decisions and outcomes, AI functioning as a tool under human control, and liability structures that preserve patients' right to recourse in the event of harm. The principle is addressed through three components. *Auditability* is understood as the requirement that high-risk AI systems are designed and documented to ensure traceability of processes and outputs, with logging and records that allow independent examination, verification, and allocation of accountability. Organisational measures adopted to address this component are recording structured logs per case, comprising input metadata, model version, parameters, user actions, overrides, and timestamps, linked to the clinical record, and model

updates tracked and documented, covering model and data versions, release notes, validation evidence, and SBOM, with a tested rollback path maintained.

Human oversight is understood as the requirement that AI systems are developed and used as a tool that serves people, respects human dignity and personal autonomy, and can be appropriately controlled and overseen by humans. Organisational measures adopted to address this component are providing a clinician appeal and contest workflow to flag questionable AI outputs and request secondary review with clear timelines, allowing clinicians to override AI outputs at any time with capture of structured reasons and logging of the decision, and clearly labelling AI-generated content versus clinician-authored conclusions in the viewer and report, including who signed off. These measures were also addressed in Section 3.3.3 under ALTAI Requirement 1, where the human oversight component of UC1's accountability and responsibility principle was discussed in relation to the broader cross-UC treatment of human oversight. *Liability* is understood as the requirement that responsibility for AI-assisted decisions is clearly defined so that manufacturers and, where applicable, healthcare providers remain accountable for system safety and performance, ensuring that patients retain the right to recourse. Organisational measures adopted to address this component are defining and communicating a RACI for AI-related errors and incidents across manufacturer, hospital, and clinical roles, and providing case-level flagging for uncertain or wrong outputs with escalation and tracking to closure.

In UC4, accountability and responsibility is also a standalone principle, understood as the requirement to ensure that AI systems operate under meaningful human oversight, with humans remaining in the loop for sensitive or ambiguous cases, with auditability enabling independent evaluation of system performance, fairness, and compliance, and with responsiveness ensuring swift corrective action when problems are identified. The principle is addressed through three components. *Human oversight* is understood as the requirement for humans in the loop for sensitive or borderline cases, so that automated systems do not function as a black box without human review for context-heavy judgements. The organisational measure adopted to address this component is human checks at each stage of building narrative. As with UC1, this measure was also addressed in Section 5.3.2. under ALTAI Requirement 1. *Auditability and evaluation* is understood as the requirement for independent evaluation by external parties, including regulators and watchdogs, so that system performance, fairness, and compliance can be assessed through transparent records and documentation. The organisational measure adopted to address this component is an evaluative framework built into the process. *Responsiveness* is understood as the obligation to act swiftly when issues such as bias or unintended harm are identified, including through rapid updates, retraining, or adjustments to the system. The organisational measure adopted to address this component is dissemination of messengers.

Comparison to ALTAI

ALTAI Requirement 7 on accountability addresses auditability of decisions and their rationale, feedback mechanisms for identifying and minimising negative impacts, and the traceability of AI-assisted outcomes. The principle as operationalised in UC1 and UC4 addresses all three of these dimensions but with substantively different emphases. UC1 approaches accountability as a governance architecture, with the three components constituting interdependent layers: auditability through structured logging and version control, human oversight through appeal, override, and labelling mechanisms, and liability through role-

responsibility assignment and escalation pathways. UC4 addresses each component through a single measure, reflecting a more compact operationalisation. The responsiveness component in UC4 has no direct parallel in UC1: it addresses the dimension of accountability concerned with learning and correction, specifically the obligation to act when problems are identified rather than merely to document them. It should also be noted that UC2’s accountability and traceability of safety decisions component was addressed under ALTAI Requirement 1 in Section 3.3.3, as it forms a component of UC2’s human oversight and autonomy principle rather than a standalone accountability principle.

UC1	UC4
Provide a clinician appeal/contest workflow to flag questionable AI outputs and request secondary review with clear timelines.	Human checks at each stage of building narrative.
Record structured logs per case (input metadata, model version, parameters, user actions, overrides, timestamps) and link them to the clinical record.	Evaluative framework built into process.
Allow clinicians to override AI outputs at any time; capture structured reasons and log the decision.	
Clearly label AI-generated content versus clinician-authored conclusions in the viewer and report, including who signed off.	
Model updates tracked and documented: track model and data versions, release notes, validation evidence, and SBOM; ensure a tested rollback path.	
Define and communicate a RACI for AI-related errors and incidents across manufacturer, hospital, and clinical roles.	
Provide case-level flagging for uncertain or wrong outputs with escalation and tracking to closure.	
	Dissemination of messengers.

Table 31 - Overview of overlaps and differences in UC organisational measures addressing the principle of "Accountability and responsibility". If the measures belong to the same row, they have (some level of) similarity; if a UC cell is left blank, it does not have a corresponding measure.

Two clusters of overlap are identified across the measures. The pairing of the clinician appeal and contest workflow (UC1) with human checks at each stage of building narrative (UC4) reflects a shared function of creating structured pathways for human scrutiny of AI outputs within the accountability framework, though the two measures differ in modality: UC1’s workflow is a contestation and escalation mechanism for querying individual outputs, while UC4’s human checks are integrated stage-gates within the narrative-building process. The pairing of UC1’s structured case logs with UC4’s evaluative framework built into the process reflects a shared function of systematic documentation and review of AI outputs to enable accountability, with UC1’s measure oriented toward internal clinical traceability and regulatory auditability, and UC4’s measure directed toward external evaluation by independent parties. The remaining measures have no parallel in the other UC. UC1’s five standalone measures reflect the breadth of its accountability operationalisation: override logging and AI-content labelling address the human oversight component at the level of individual decisions, version control and rollback address auditability at the level of system lifecycle management, and the RACI and case-level flagging address the liability component through responsibility assignment and structured incident tracking. UC4’s dissemination of

messengers stands alone as the only measure in this section oriented toward a responsive and communicative form of accountability: rather than recording or reviewing outputs, it addresses the correction of identified problems through trusted social networks, reflecting the counter-narrative and community-engagement context in which UC4 operates.

5.3.3. Remaining challenges and concerns

Across the DSS use cases, where measures converge they do so around a shared concern with ensuring human review of AI outputs is genuine rather than a formality, with UC1, UC2, and UC3 each developing governance structures, escalation mechanisms, and training programmes toward this end. The form these take, however, diverges sharply with deployment context, with UC1's measures oriented toward clinical accountability and patient-facing justification, UC2 toward preserving engineering expertise in safety-critical workflows, and UC3 toward employee rights and accessible information under GDPR. This context-dependence means that the remaining unresolved challenges, namely liability attribution, individual-level justifiability, and the transparency constraints of security contexts, manifest differently across use cases but share a common feature in that they cannot be addressed by organisational measures alone and require legal clarification or regulatory intervention that lies beyond the reach of any single organisation.

The challenge of liability was raised directly by partners in UC1 as an unresolved issue. While the partners affirm that liability must be clearly defined so that professionals remains accountable for the outcomes of AI assisted decision making, ensuring the patients retains the right to recourse in the event of harm, they are equally clear that there is not clear way to achieve this in practice. Liability in the context of AI in healthcare is a grey area, whose clarification depends on how provisions of the EU AI Act are implemented and interpreted. This is a particularly salient concern for high-risk systems where, at present, there is no established standard in Europe for what a fully compliant high-risk medical AI device should look like. As a result, liability is largely outside partners' control.

This concern is not limited to clinical contexts. Dietrich (2025) notes the EU AI Act's risk based framework do not sufficiently address liability concerns as they relate to DSS use more broadly. O'Sullivan et al. (2019) explore how accountability and liability can be shared across an organisation, and this approach is operationalised in OM05, with an organisation-wide RACI framework for responsibility assignment. However this measure is focused on how organisations can manage the ambiguity surrounding liability rather than giving any broader legal certainty as to who is liable for faulty decisions involving DSS.

In UC1, partners noted that Explainable AI can be defined in different ways. A technical perspective of explainability might refer to features and behaviour of the system. meanwhile, a clinical perspective of explainability focuses on specific cases, for example given a CT scan, why does the system suggest that this patient has breast cancer? This form of explanation is more complex technically to achieve, but directly relevant to clinical use and liability, since responsibility often requires explanations at the inference level. In neural networks trained on images, it may be possible to justify outputs by showing why certain features are included or excluded. However, at the single patient level, interpreting one unique medical image becomes virtually impossible.

That there are multiple ways of understanding explainability has implications for medical decision making. Clinicians must defend treatment choices and ensure, if and when using DSS, that such AI supported decisions can be explained in ethically and professionally acceptable ways. This justifiability element of explainability is incredibly difficult to achieve given the complexity of opening up the black box of neural networks and remains an unresolved issue. One important thing to address this, as flagged by UC1 partners, is to include perspectives from multiple relevant stakeholders, for example patients and clinicians. Braun et al. (2021) note that because deep learning systems cannot account for how a result was reached, patients cannot contest a recommendation whose reasoning is inaccessible to them, weakening the basis for informed consent and shared decision-making.

OM09, which commits to the design of explanations that support meaningful review, contestation or justification, cannot resolve this challenge pertaining to clinician level explainability. The measure addresses the conditions under which an explanation is presented, not whether the explanation itself is adequate grounds for clinical or professional justification at the individual case level.

In UC4, a key challenge related to developing detailed accounts of AI tools used in security contexts, particularly in areas such as counter-terrorism and hate speech detection, due to the sensitive and classified nature of these systems. Security professionals are bound by strict confidentiality and legal restrictions that prevent them from disclosing the specific tools, algorithms, or operational methods in use. Revealing these details could expose state security processes and compromise ongoing intelligence operations. This created a significant methodological and ethical challenge for the work, and required focus to be shifted towards describing broad technical measures and ethical principles that underpin responsible AI use in security practice.

A cross-cutting challenge noted by use case partners was that human-in-the-loop mechanisms can be resource intensive and fatigue-inducing. Overseers must sustain attention and judgement across complex and repetitive tasks within fast-paced and high-pressure environments. Concerns of automation bias and alarm fatigue arise in both UC1 and UC2, with the broader concern of over-reliance and dependence spanning across use cases. Flowing from this, there is a shared understanding across UC1, 2, and 3, that there is a difference between oversight in itself and meaningful oversight. That measures are adopted to encourage meaningful oversight across these 3 use cases demonstrates the centrality of this concern to the research area. Compounding this is the risk that commercial and operational pressures can help promote a laxer approach to oversight, and that responsibility for making oversight meaningful there lies largely on the organisational level, and crucially the contexts in which organisations operate. Challen et al. (2019) document automation complacency across medical DSS, finding that users who trust generally reliable systems tend to stop scrutinising outputs, while Cobianchi et al. (2022) note that clinicians who rely on accurate systems may trust their own assessments less over time, with practical diagnostic competencies risking atrophy.

OM25, which commits to periodically reassessing reliance patterns as systems evolve or scale, is relevant here. It acknowledges explicitly that risks of over-reliance and deskilling are most salient in the long term and can have a structural impact on organisational practices, however outside of its scope are the market contexts which incentivise rapid AI adoption and minimal oversight.

6. CONCLUSION

The development of operational, non-technical guidelines for AI research areas, has resulted in a rich portfolio of organisational measures, applicable to a broad range of contexts. This process, which relied on empirical data from AIOLIA use cases, as well as multiple co-creation sessions, has demonstrated that, while the specific context of AI deployment significantly impacts the organisational measures deployed, there is a degree of overlap in the underlying ethical concerns, particularly around human autonomy and safety, even though the conceptualisation of these principles differs depending on whether the UC applies to private or professional contexts.

A core conclusion emerging from the analysis conducted is that, whereas ethics principles form an important basis from which to approach the development and deployment of trustworthy AI systems, there is a strong correlation between the specificities of the contexts within which organisations operate and the concrete operationalisation of these principles, even when the same type of AI systems is involved, resulting in a very limited overlap between organisational measures. Across the three research areas, convergence at the level of organisational measures tended to occur around the operationalisation of human oversight. This was the case in the GPAI research area, alongside transparency (of documentation and criteria), involvement of diverse stakeholders in system design and evaluation, and personalisation of the GPAI system to enable human safety. In Emotional AI, similarity was only identified at the level of human oversight measures. In DSS, human oversight also figured prominently through a shared concern with ensuring human review of AI outputs is genuine rather than a formality, with UCs putting in place governance structures, escalation mechanisms, and training programmes to this end.

Given the limited overlap in the organisational measures identified across use cases, the extrapolation of research area-level organisational measures was not straightforward, but the criteria defined in 4.2.4. allowed not only for the cross-analysis and lifting of overarching ethical concerns, but also for the creation of a foundational portfolio of organisational measures for each research area. The latter, rather than guaranteeing the ‘resolution’ of ethics principles, seeks to inspire the pursuit of an ongoing critical practice within organisations by rendering abstract theoretical constructs more tangible. Importantly, while some of the measures identified herein support legal compliance with the AI Act, such as those related with human oversight and transparency, they reflect a clear ambition move beyond legal compliance and engage in ethical practices that contribute to a more desirable outlook, despite the challenges posed by AI technologies. This is particularly evident in measures targeting deskilling, organisational culture and shared governance and responsibility. All in all, organisational responsibility beyond legal compliance is particularly pressing in the context of an intensification of AI uptake and diffusion across private and professional contexts, varied sectors and domains, that necessarily impact human cognition and behaviour in the immediate and longer term.

While the key challenges and concerns identified for each research area in the extrapolation of organisational measures are outlined in section 5, two overall findings shall inspire further action from policymakers, namely, the development of standards and provision of legal certainty, and the facilitation of dialogue and deliberation on societal and environmental-level impacts of AI.

Across the three research areas analysed, AIOLIA partners operating in more standardised domains, such as the clinical context and automotive industry, appeared to have a clearer path for translating ethical principles into practice. Rather than hindering AI uptake, standards provided a practical baseline and offered certainty on how organisations should approach AI development and deployment. Particularly, industry-specific standards and regulations, such as the ISO automotive industry standards or the Medical Device Regulation, supported not simply the operationalisation of AI ethics principles for use cases operating in those areas, but provided insights on the concrete organisational measures that could be applied to the research areas as a whole. This suggests that the EU should pursue the development of standards, including the harmonised standards to support compliance with the EU AI Act, with urgency and, potentially, engage in furthering sector-specific standards, particularly in the nascent field of Emotional AI. Importantly, there are also limits to what organisational measures alone can achieve, especially around liability in DSS, which requires legal clarification and regulatory intervention that lies beyond the reach of any single organisation.

Lastly, while societal and environmental concerns figure high on AI ethics agendas, including in ALTAI, the research process conducted herein has shown that organisations struggle to address these, for they require not only broader deliberation and dialogue that go beyond the remit of a single organisation, but also the development of standardised metrics and guidance on what a high or acceptable environmental imprint looks like in the context of AI development and deployment.

7. REFERENCES

Albladi, A., Islam, M., Das, A., Bigonah, M., Zhang, Z., Jamshidi, F., ... & Seals, C. (2025). Hate speech detection using large language models: A comprehensive review. *IEEE Access*, 13, 20871-20892.

Bakir, V., Bennet, K., Bland, B., Laffer, A., Li, P., & McStay, A. (2024). When is deception ok? developing the IEEE recommended practice for ethical considerations of emulated empathy in partner-based general-purpose artificial intelligence systems. *IEEE International Symposium on Technology and Society*.

Bayerl, P.S., Lawlor, E., Maris, M.T., Miorandi, D., Pekšys, G., Bjelica, M., Stojšin, K., Anastasova, M., Smith, O., Henestrosa, A., Yamshchikov, I., Bak, M.A.R., & Akhgar, B. (2026). *Operational ethics guidelines on use cases related to human behaviour and cognition*. AIOLIA Deliverable D3.1.

Berthelot, A., Caron, E., Jay, M., Lefevre, L. (2024). Estimating the environmental impact of Generative-AI services using an LCA-based methodology. *Procedia CIRP*, 122. <https://doi.org/10.1016/j.procir.2024.01.098>

Bleher, H., & Braun, M. (2022). Diffused responsibility: attributions of responsibility in the use of AI-driven clinical decision support systems. *AI and Ethics*, 2(4), 747-761.

Boine, C. (2023). Emotional Attachment to AI Companions and European Law. *MIT Case Studies in Social and Ethical Responsibilities of Computing*, no. Winter 2023 (February). <https://doi.org/10.21428/2c646de5.db67ec7f>.

Braun, M., Hummel, P., Beck, S., & Dabrock, P. (2021). Primer on an ethics of AI-based decision support systems in the clinic. *Journal of medical ethics*, 47(12), e3-e3.

Braun, M., Müller, R. (2025). Missed opportunities for AI governance: lessons from ELS programs in genomics, nanotechnology, and RRI. *AI & Soc* 40, 1347–1360. <https://doi.org/10.1007/s00146-024-01986-0>

- Brynjolfsson, E., Chandar, B., & Chen, R. (2025). Canaries in the Coal Mine? Six Facts about the Recent Employment Effects of Artificial Intelligence. *Stanford Digital Economy Lab*.
- Challen, R., Denny, J., Pitt, M., Gompels, L., Edwards, T., & Tsaneva-Atanasova, K. (2019). Artificial intelligence, bias and clinical safety. *BMJ quality & safety*, 28(3), 231-237.
- Ciriello, R., Hannon, O., Chen, A. Y., & Vaast, E. (2024). Ethical tensions in human-AI companionship: A dialectical inquiry into Replika. *Proceedings of the 57th Hawaii International Conference on System Sciences*.
- Cobianchi, L., Verde, J. M., Loftus, T. J., Piccolo, D., Dal Mas, F., Mascagni, P., ... & Kaafarani, H. M. (2022). Artificial intelligence and surgery: ethical dilemmas and open issues. *Journal of the American College of Surgeons*, 235(2), 268-275.
- Corfmat, M., Martineau, J. T., & Régis, C. (2025). High-reward, high-risk technologies? An ethical and legal account of AI development in healthcare. *BMC medical ethics*, 26(1), 4.
- Cote, M., Aires, S. (2025). Futurity as Infrastructure: A Techno-Philosophical Interpretation of the AI Lifecycle. *Proceedings from the Fourth International Conference on Hybrid Human-Artificial Intelligence (HHAI)*. <https://ceur-ws.org/Vol-4074/paper5-1.pdf>
- De Freitas, J., Oguz-Uguralp, Z., & Kaan-Uguralp, A. (2025). Emotional manipulation by AI companions. *arXiv preprint arXiv:2508.19258*.
- Dewitte, P. (2024). Better alone than in bad company: Addressing the risks of companion chatbots through data protection by design. *Computer Law & Security Review*, 54, 106019.
- Dietrich, F. (2025). Ai-based removal of hate speech from digital social networks: chances and risks for freedom of expression. *AI and Ethics*, 5(3), 2943-2953.
- European Commission – AI HLEG. (2019). *Ethics guidelines for trustworthy AI*. Publications Office.
- European Commission. (2020). *Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self-assessment*. European Commission.
- European Data Protection Supervisor. (n.d.). *AI companions*. https://www.edps.europa.eu/data-protection/technology-monitoring/techsonar/ai-companions_en
- Gabriel, I., Manzini, A., Keeling, G., Hendricks, L. A., Rieser, V., Iqbal, H., ... & Manyika, J. (2024). The ethics of advanced AI assistants. *arXiv preprint arXiv:2404.16244*.
- Guingrich, R. E., & Graziano, M. S. (2024). Ascribing consciousness to artificial intelligence: human-AI interaction and its carry-over effects on human-human interaction. *Frontiers in Psychology*, 15, 1322781.
- Hoek, S., Metselaar, S., Ploem, C., & Bak, M. (2025). Promising for patients or deeply disturbing? The ethical and legal aspects of deepfake therapy. *Journal of Medical Ethics*, 51(7), 481-486.
- Hunkenschroer, A. L., & Luetge, C. (2022). Ethics of AI-enabled recruiting and selection: A review and research agenda. *Journal of business ethics*, 178(4), 977-1007.
- Ienca, M., & Malgieri, G. (2022). Mental data protection and the GDPR. *Journal of Law and the Biosciences*, 9(1), Isac006.

- Kiritchenko, S., & Nejadgholi, I. (2020). Towards ethics by design in online abusive content detection. *arXiv preprint arXiv:2010.14952*.
- Kraaijeveld, S. R., & Ivanova, D. (2026). Using deepfakes for psychotherapy: ethical and philosophical issues. *AI and Ethics*, 6(1), 79.
- Laine, J., Minkinen, M., & Mäntymäki, M. (2025). Understanding the Ethics of Generative AI: Established and New Ethical Principles. *Communications of the Association for Information Systems*, 56, 1-25.
<https://doi.org/10.17705/1CAIS.05601>
- Lee, H-P., Sarkar, A., Tankelevitch, L., Drosos, I., Rintel, S., Banks, R., Wilson, N. (2025). The Impact of Generative AI on Critical Thinking: Self-Reported Reductions in Cognitive Effort and Confidence Effects From a Survey of Knowledge Workers. *CHI Conference on Human Factors in Computing Systems (CHI '25)*. <https://doi.org/10.1145/3706598.3713778>
- Lemley, Mark A., How Generative AI Turns Copyright Upside Down (July 21, 2023). Available at SSRN: <https://ssrn.com/abstract=4517702>
- Ligozat, A.-L., Lefevre, J., Bugeau, A., & Combaz, J. (2022). Unraveling the Hidden Environmental Impacts of AI Solutions for Environment Life Cycle Assessment of AI Solutions. *Sustainability*, 14(9), 5172.
<https://doi.org/10.3390/su14095172>
- Mager, A., Eitenberger, M., Winter, J., Prainsack, B., Wendehorst, C., & Arora, P. (2025). Situated ethics: Ethical accountability of local perspectives in global AI ethics. *Media, Culture & Society*, 47(5), 1028-1041.
- Mahari, R., Pataranutaporn, P. (2025). "Addictive Intelligence: Understanding Psychological, Legal, and Technical Dimensions of AI Companionship." *MIT Case Studies in Social and Ethical Responsibilities of Computing*, no. Winter (March). <https://doi.org/10.21428/2c646de5.2877155b>.
- Malfacini, K. (2025). The impacts of companion AI on human relationships: risks, benefits, and design considerations. *AI & Society*, 40(7), 5527-5540.
- McStay, A. (2025). Emulated Empathy: Can Risks Be Countered by a Soft-Law Standard?. *IEEE transactions on technology and society*.
- McStay, A., & Bakir, V. (2025). Soft law for unintentional empathy: Addressing the governance gap in emotion-recognition AI technologies. *Journal of Responsible Technology*, 100126.
- Muller, C., & Teilhard De Chardin, A. (2025). *AI Act interpretation: Definition and prohibitions*. ALLAI. <https://allai.nl/wp-content/uploads/2025/01/AI-Act-Interpretation-Definition-and-Prohibitions.pdf>
- Munn, L. (2023). The uselessness of AI ethics. *AI Ethics* 3, 869–877. <https://doi.org/10.1007/s43681-022-00209-w>
- Natali, C., Marconi, L., Dias Duran, L.D. et al. (2025). AI-induced Deskilling in Medicine: A Mixed-Method Review and Research Agenda for Healthcare and Beyond. *Artif Intell Rev* 58, 356. <https://doi.org/10.1007/s10462-025-11352-1>
- OECD. (2022). *Measuring the environmental impacts of artificial intelligence compute and applications: The AI footprint*. OECD Digital Economy Papers, No. 341. <https://doi.org/10.1787/7babf571-en>.

- O'Sullivan, S., Nevejans, N., Allen, C., Blyth, A., Leonard, S., Pagallo, U., ... & Ashrafian, H. (2019). Legal, regulatory, and ethical frameworks for development of standards in artificial intelligence (AI) and autonomous robotic surgery. *The international journal of medical robotics and computer assisted surgery*, 15(1), e1968.
- Pataranutaporn, P., Karny, S., Archiwaranguprok, C., Albrecht, C., Liu, A. R., & Maes, P. (2025). " My Boyfriend is AI": A Computational Analysis of Human-AI Companionship in Reddit's AI Community. *arXiv preprint arXiv:2509.11391*.
- Paz, A. (2025). *A Call to Address Anthropomorphic AI Threats to Freedom of Thought*. Centre for International Governance Innovation.
- Resseguier, A., Rodrigues, R. (2020). AI ethics should not remain toothless! A call to bring back the teeth of ethics. *Big Data & Society*. <https://doi.org/10.1177/205395172094254>
- Ryan, M., Blok, V. (2025). What's Economics Got to Do with It? Providing Theoretical Clarity on ELSA of AI. *Sci Eng Ethics* 31, 37. <https://doi.org/10.1007/s11948-025-00564-x>
- Ryan, M., de Roo, N., Wang, H. et al. (2025). AI through the looking glass: an empirical study of structural social and ethical challenges in AI. *AI & Soc* 40, 3891–3907. <https://doi.org/10.1007/s00146-024-02146-0>
- Samuelson, P. (2023). Generative AI meets copyright. *Science* 381, 158-161. DOI:10.1126/science.adi0656
- Sharma, S., Selwal, A. (2026). Potential of artificial intelligence in deepfake media: From generation to detection mechanisms, state-of-the-art, and challenges. *Computer Science Review*, 60.
- Standards Committee. (2024). *IEEE Standard for Ethical Considerations in Emulated Empathy in Autonomous and Intelligent Systems*. IEEE Std 7014™.
- Stark, L., & Hoey, J. (2021, March). The ethics of emotion in artificial intelligence systems. *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency* (pp. 782-793).
- Steindl, E. (2025). *A datafied mind: Untangling EU regulation of emotion technology and neurotechnology*. Cambridge University Press.
- Susser, D., Roessler, B., & Nissenbaum, H. (2019). Online manipulation: Hidden influences in a digital world. *Geo. L. Tech. Rev.*, 4, 1.
- Teo, S. A., Mann, S. P., & Jurcys, P. (2027, in press). The ethical and legal complexities of regulating companion AI chatbots. *Law, Innovation and Technology*, 19(1).
- Teo, S.A., Kyosovska, N. and Armengol, A. (2025). *AIOLIA D2.2: Report on the selection of ethical principles and values*.
- Triguero, I., Molina, D., Poyatos, J., Del Ser, J., & Herrera, F. (2024). General purpose artificial intelligence systems (GPAIS): Properties, definition, taxonomy, societal implications and responsible governance. *Information Fusion*, 103, Article 102135. <https://doi.org/10.1016/j.inffus.2023.102135>
- Udupa, S., Maronikolakis, A., & Wisiorek, A. (2023). Ethical scaling for content moderation: Extreme speech and the (in) significance of artificial intelligence. *Big Data & Society*, 10(1), 20539517231172424.
- van Hilten, M., Ryan, M., Blok, V., de Roo, N. (2025). Ethical, Legal and Social Aspects (ELSA) for AI: An assessment tool for Agri-food. *Smart Agricultural Technology*, 10. <https://doi.org/10.1016/j.atech.2024.100710>

van Wynsberghe, A. (2021). Sustainable AI: AI for sustainability and the sustainability of AI. *AI Ethics* 1, 213–218. <https://doi.org/10.1007/s43681-021-00043-6>

Wang, H., Blok, V., & van Hilten, M. (2025). ELSA Labs for responsible AI: a novel approach for addressing ethical, legal, social issues. *Journal of Responsible Innovation*, 12(1). <https://doi.org/10.1080/23299460.2025.2563944>

Wang, S., & Dehnert, M. (2026). On-Demand Intimacy: The Sociotechnical Appeal of AI Companions. *Social Media+ Society*, 12(1), 20563051251410394.

Wilders, M., Martínez de Rituerto de Troya, Íñigo, & Dobbe, R. (2025). The Socratic Dialogue as a Method for Virtue Ethics in AI: A Case Study. *Proceedings of the AAAI ACM Conference on AI, Ethics, and Society*, 8(3), 2680–2691. <https://doi.org/10.1609/aies.v8i3.36748>

Zhong, H., O'Neill, E., & Hoffmann, J. A. (2024). Regulating AI: applying insights from behavioural economics and psychology to the application of Article 5 of the EU AI Act. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 38, No. 18, pp. 20001-20009).