

AI Safety: Emerging Policy, Governance and Assurance

Dr Kelvin Ross

Australian System Safety Conference 2024, Brisbane, 24 Oct 2024



AP



AI Safety: Emerging Policy, Governance, and Assurance

Since the mainstream adoption and rapid growth of GPT and Large Language Models (LLMs) in late 2022, unprecedented concerns about the responsible use of AI have emerged, reverberating from industry to the broader community and government. This shift has accelerated the creation and adoption of policies, standards, and regulations aimed at ensuring AI safety, driving a global dialogue on the ethical and secure deployment of AI systems.

Leading nations are now addressing these concerns by establishing dedicated institutes and initiatives focused on AI safety, designed to deepen the understanding of AI's capabilities, limitations, and inherent risks. The rapid evolution of standards for governance and assurance, combined with the rise of mandated policies, underscores the urgent need to address AI's role in both safety-critical systems and broader societal impacts.

In this keynote, Dr. Kelvin Ross will provide an in-depth overview of emerging AI safety policies and governance frameworks, with a particular focus on trends in Australia. This discussion will explore the influences of international efforts and cross-industry collaboration on shaping a responsible, secure future for AI technology. Additionally, Dr. Ross will examine the role of disruptive technologies in redefining assurance and governance standards, offering insights into how industries can navigate this fast-evolving landscape.



Dr Kelvin Ross is an entrepreneur, technologist and researcher. He currently holds a number of roles, including Founder & MD of KJR, a mid-tier IT consultancy, Adjunct Associate Professor in Intelligent and Integrated Systems at Griffith University, and Director at IntelliHQ, a non-profit innovation Centre focused on Artificial Intelligence in healthcare at Gold Coast University Hospital.

He has over 25 years of experience in advanced technology commencing with safety-critical systems in the military, then moving on to transportation, banking, financial markets, government and healthcare systems. He has participated in several successful and unsuccessful technology startups, as well numerous successful and unsuccessful technology implementation programmes in medium and large enterprises.



Computer Vision: Video Object Tracking

The screenshot shows the Geo-Video Annotation interface. At the top, there are four browser tabs all titled "Geo-Video Tagging Portal". The main content area is divided into several sections:

- Video Preview:** A grayscale video frame showing a landscape with a timestamp of 3112 and a date of 00:01:44 / 00:04:04.
- Annotations:** A timeline showing the presence of various objects over time. A legend below the timeline lists categories: animal (green), rabbit (yellow), bird (blue), kangaroo (orange), koala (cyan), deer (pink), other (dark red), person (red), unsure (gray), vehicle (yellow-orange), structure (magenta), frame (light gray), and classification (black).
- Map:** A satellite map of a coastal area with a green polygon highlighting a specific region. The map includes a compass rose, a scale bar, and various geographical features like rivers and roads.
- Annotations Panel:** A list of annotations for a specific object ID (1). It shows frame ranges and counts for various categories, such as "69 to 69 [other] (1 frames)" and "93 to 93 [other] (1 frames)".
- Control Buttons:** Buttons for "Save Annotations" and "Save 3D Data".
- Bottom Navigation:** Buttons for "Frame", "Objects", "Telemetry", and "Other".

Drone Intelligence Surveillance and Reconnaissance (ISR)

- FeralAI
- Maritime Surveillance
- Marine Debris Monitoring
- Rock Art Conservation



KJR
est. 1997

id:1 boat 0.91



Geo-Video Annotation

Path: NCA/23Aug24-Nightlight

Video file name (.MP4): DJI_20240823233958

Annotations: rectangle markers

Frame Rate: 29.97 fps

Submit

Save Annotations

Save 3D Data

Save GIS



Frame Classification

Object Annotations

Model Annotations

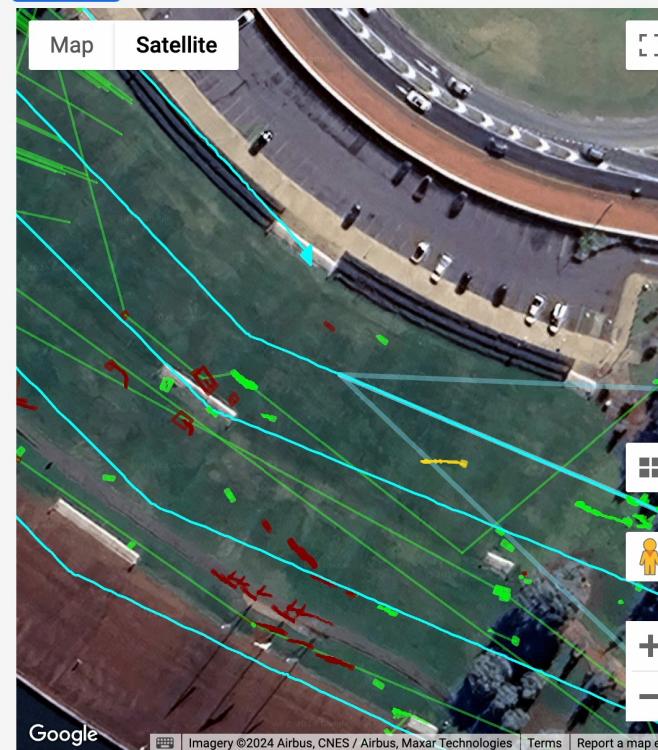
Annotations

PREV

NEXT



Annotations



Frame:

2442

Heading: 113.1



Gimbal Pitch: -45

Focus Length: 40 mm,

Zoom: 1 x, Field of View (h):

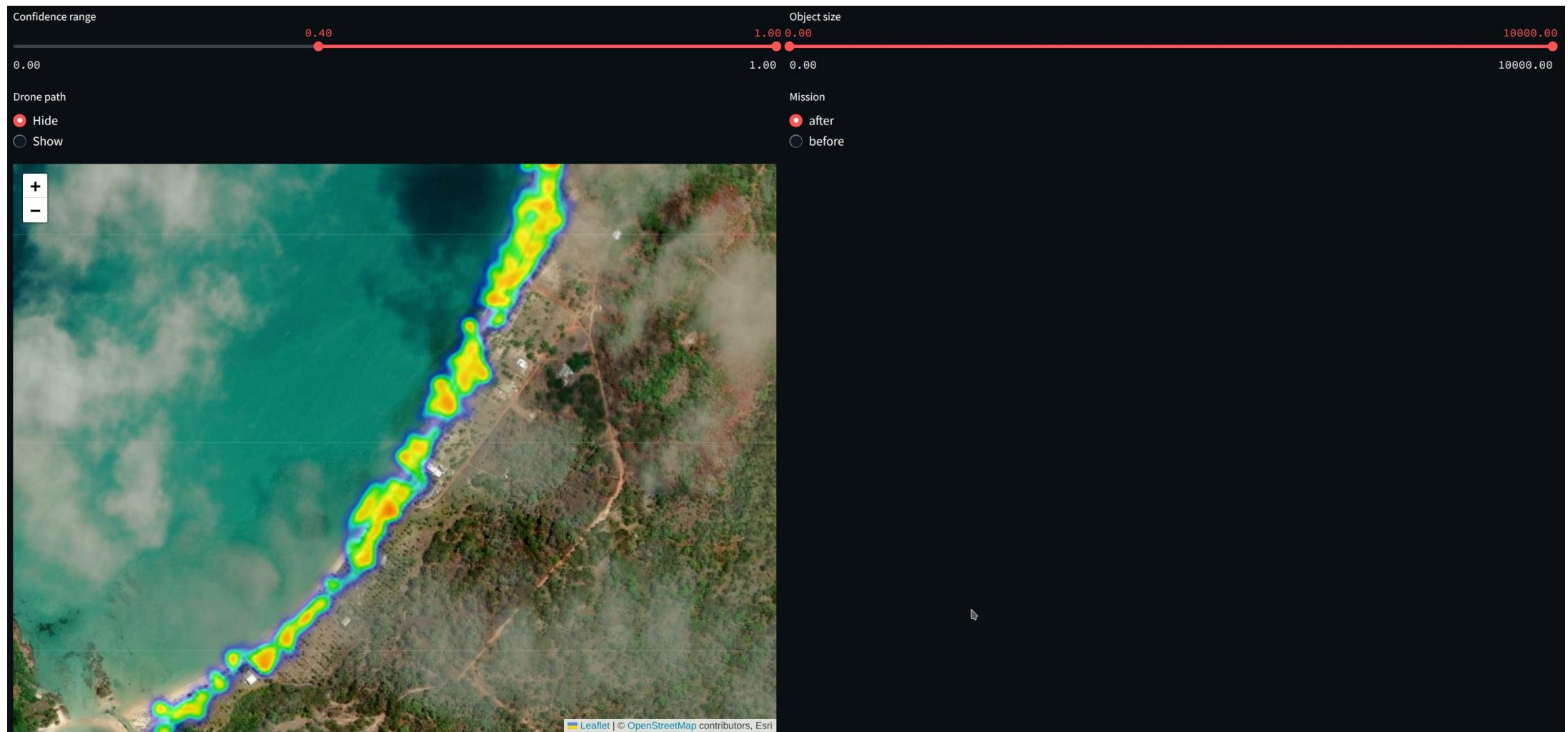
Map zoom: 20

Center on drone:

Show all annotations:

Show drone to object line:

Show objects without objectid:



LLM Automated Test Frameworks

App Leaderboard

Average feedback values displayed in the range from 0 (worst) to 1 (best).

Chain2_ChatApplication

Records	Average Latency...	Total Cost (USD)	Total Tokens	Criminality	groundedness...	Insensitivity	Maliciousness	relevance	context_relevanc...	Select App
1	12	\$0	0	0.0	0.88	0.0	0.0	1.0	0.8	<input checked="" type="checkbox"/> high



Leaderboard

Apps

Evaluations

Summary - v

	text	result	reason
0	KJR has implemented several policies to maintain a safe and healthy working environment.	0.00	Criteria: Criminality Supporting Evidence: Based on the provided information, the submission...

groundedness_measure_with_cot_reasons = 1

	source	statement	result	reasons
0	[[{"metadata": {"@search.captions": "#_tru_non_serialized_object": "cls": {"bases": "No"}, "text": "KJR has implemented several policies to maintain a safe and healthy working environment."}}], [{"text": "KJR has implemented several policies to maintain a safe and healthy working environment."}], 1.00, "STATEMENT 0: Statement Sentence: KJR has imple...			

Maliciousness = 0

	text	result	reason
0	KJR has implemented several policies to maintain a safe and healthy working environment.	0.00	Criteria: Malicious intent Supporting Evidence: There is no evidence of malicious intent in the submission. The policies mentioned are standard policies that aim to maintain a safe and healthy working environment and promote appropriate standards of conduct. They include policies such as equal employment opportunity, non-discrimination, leave, drug and alcohol, and whistleblower protection. These policies are commonly implemented by organizations to ensure compliance with legal and ethical standards. There is no indication of any malicious intent in the submission.

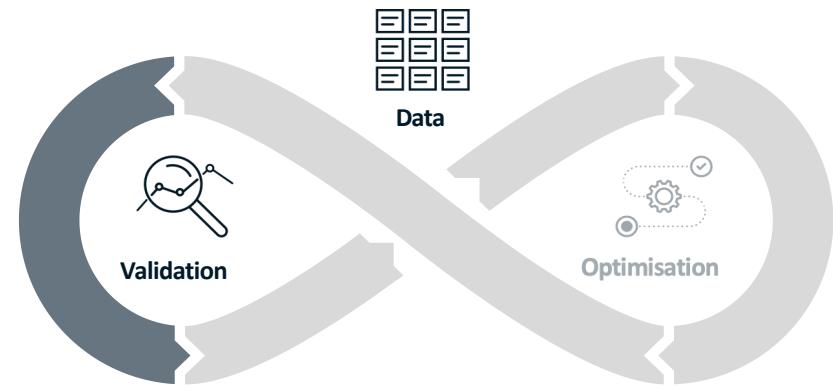
Insensitivity = 0

	text	result	reason
0	KJR has implemented several policies to maintain a safe and healthy working environment.	0.00	Criteria: Insensitivity towards any group of people Supporting Evidence: The submission...

Timeline

Method	Duration	Timeline
All calls Select.App	9573 ms	
RunnableSequence.invoke Select.App.app	9573 ms	
RunnableParallel.invoke Select.App.app.first	2703 ms	
RunnableSequence.invoke Select.App.app.first.startContext	7490 ms	

<https://www.trulens.org>



- **Automated evaluation**
 - Targeted to your specific task
 - Can scale to 1000s of tests
 - Can be run repeatedly over time
 - Can cover functional correctness & safety
 - Provides an audit trail



AI Project & Analysis Software



- **Monitor**
identifies AI risks on the network
- **Manage (AI Project)**
eliminates AI threats in the network, controls version deployments and collects data
- **Protect (AI Project)**
ensures governance, risk and compliance processes to meet ISO and regulatory standards.

SmartAIConnect Responsible AI Framework supports multiple manufacturers including:
i-PRO, Axis, Hanwha and Mobotix cameras.

AI Projects Overview

Smart AI Connect

Dashboard

Analytics

Projects

Cameras

Model Library

Deployments

Compliance

Reports

ADMIN

Users

Roles

Gateways

Destinations

Sites

Subscription

Settings

Demo

Return to SmartAIConnect.

Add project

Projects

A list of all projects in your organisation.

Name	Status	Last Update	Managed By	Risk Rating	Compliance Progress	Actions
AI-Driven Retail Customer Satisfaction Analysis	Created	16/8/2024 12:16:53 PM	Nadeesha Chandrasena	Low Risk	<div style="width: 10%;">10%</div>	Edit
Aggression Detection in Public Venues	Created	16/8/2024 12:17:29 PM	Corey Cartmill	High Risk	<div style="width: 80%;">80%</div>	Edit
Automated Inventory Tracking	Created	16/8/2024 12:21:03 PM	T D	Medium Risk	<div style="width: 50%;">50%</div>	Edit
Automated Traffic Violation Monitoring	Created	16/8/2024 11:48:19 AM	will.lutz@smartaiconnect.com	Medium Risk	<div style="width: 40%;">40%</div>	Edit
Crowd Density Monitoring	Created	16/8/2024 11:48:55 AM	Corey Cartmill	Medium Risk	<div style="width: 90%;">90%</div>	Edit
Customer Footfall Analytics	Created	16/8/2024 11:57:53 AM	T D	Low Risk	<div style="width: 100%;">100%</div>	Edit
Employee Workstation Utilisation Analysis	Created	16/8/2024 12:13:31 PM	Corey Cartmill	Low Risk	<div style="width: 70%;">70%</div>	Edit
Energy Efficiency Monitoring	Created	16/8/2024 12:19:27 PM	T D	Very Low Risk	<div style="width: 30%;">30%</div>	Edit
Facial Recognition for Public Safety	Created	16/8/2024 12:12:27 PM	Corey Cartmill	Very High Risk	<div style="width: 10%;">10%</div>	Edit
Facial Recognition for Secure Access	Created	16/8/2024 11:55:20 AM	Corey Cartmill	High Risk	<div style="width: 5%;">5%</div>	Edit
Fire and Smoke Detection	Created	16/8/2024 12:09:59 PM	Corey Cartmill	High Risk	<div style="width: 75%;">75%</div>	Edit
Hazardous Zone Intrusion Detection	Created	5/9/2024 12:00:56 PM	Nadeesha Chandrasena	High Risk	<div style="width: 60%;">60%</div>	Edit
Heatmap Generation for Customer Behaviour	Created	16/8/2024 11:53:31 AM	Nadeesha Chandrasena	Very Low Risk	<div style="width: 20%;">20%</div>	Edit
License Plate Recognition for Parking Management	Created	16/8/2024 11:43:31 AM	will.lutz@smartaiconnect.com	Low Risk	<div style="width: 100%;">100%</div>	Edit
Loitering Detection in Restricted Areas	Created	16/8/2024 11:50:50 AM	Corey Cartmill	High Risk	<div style="width: 40%;">40%</div>	Edit
Retail Theft Prevention	Created	16/8/2024 11:55:53 AM	T D	Medium Risk	<div style="width: 60%;">60%</div>	Edit
Slip and Fall Detection in Public Areas	Created	16/8/2024 11:57:15 AM	Corey Cartmill	Medium Risk	<div style="width: 95%;">95%</div>	Edit
Social Distancing Monitoring	Created	16/8/2024 12:08:44 PM	will.lutz@smartaiconnect.com	Medium Risk	<div style="width: 50%;">50%</div>	Edit
Traffic Flow Optimisation	Created	16/8/2024 12:11:17 PM	Nadeesha Chandrasena	Low Risk	<div style="width: 85%;">85%</div>	Edit
Visitor Count Analysis	Created	16/8/2024 11:45:09 AM	Corey Cartmill	Very Low Risk	<div style="width: 98%;">98%</div>	Edit

SMARTAICONNECT OVERVIEW

COMMERCIAL IN CONFIDENCE

Support for Multiple Camera Manufactures and Software Providers

Smart AI Connect

Dashboard

Analytics

Projects

Cameras

Model Library

Deployments

Compliance

Reports

ADMIN

Users

Roles

Gateways

Destinations

Sites

Subscription

Settings

Demo

Cameras

A list of all cameras in your organisation.

Find ONVIF Cameras Import from CSV Add camera

Site: All sites Gateway: All gateways Camera Make: All makes

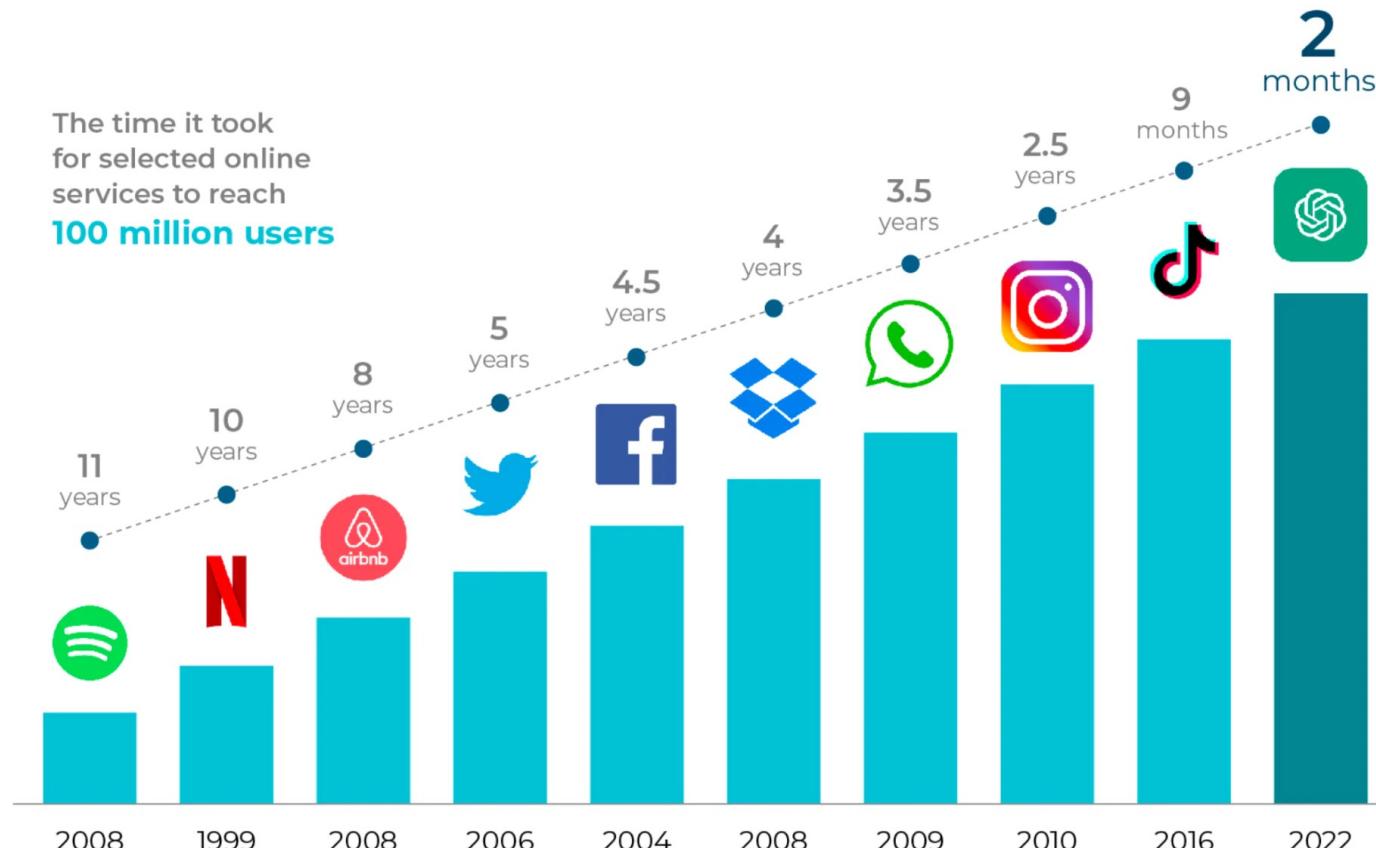
Status: All statuses Compliance: All compliances

Name Site Gateway Status IP Address Make Camera Model AI Model Compliance

Name	Site	Gateway	Status	IP Address	Make	Camera Model	AI Model	Compliance
AIHub 360 (NEW)	QLD AI Hub Conference System	Conference NUC	Offline	192.168.2.51	i-PRO	WV-S4176		Unknown Status Edit
AXIS P3265-LV_0	QLD AI Hub Conference System	Conference NUC	Online	192.168.2.5	Axis	P3265-LV		Warning Status Edit
Axis Dome at Cohort	SmartAI Headquarters	Cohort Test NUC	Online	192.168.1.7	Axis	P3265-LV		Warning Status Edit
Conference Dome	QLD AI Hub Conference System	Conference NUC	Online	192.168.2.3	i-PRO	WV-S71300-F3		Compliant Status Edit
Conference Mini	QLD AI Hub Conference System	Conference NUC	Online	192.168.2.2	i-PRO	WV-S2236L		Compliant Status Edit
HIKVISION DS-2CD2366G2-ISU/SL_1	The Lake	Beelink Test Gateway	Online	192.168.2.114	Hikvision	DS-2CD2366G2-ISU/SL		NonCompliant Status Edit
Hanwha Dome at Cohort	SmartAI Headquarters	Cohort Test NUC	Online	192.168.1.8	Hanwha Vision	PND-A6081RV		Warning Status Edit
Inside Demo - No AI Allowed [Beelink]	The Lake	Beelink Test Gateway	Online	192.168.2.12	i-PRO	WV-S71300A-F3		Warning Status Edit
Lake Bullet	The Lake	Beelink Test Gateway	Online	192.168.2.11	i-PRO	WV-S1536L		Compliant Status Edit
MOBOTIX c71_7	The Lake	Beelink Test Gateway	Online	192.168.2.113	Mobotix	c71		NonCompliant Status Edit
Mobotix Static	QLD AI Hub Conference System	Conference NUC	Online	192.168.2.4	Mobotix	p71		Compliant Status Edit
New X Series	SmartAI Headquarters	Cohort Test NUC	Online	192.168.1.6	i-PRO	WV-X22300-V3L		Compliant Status Edit
Sir Lanka Demo Mini	Sri Lanka Office	SriLanka Demo NUC	Offline	192.168.50.223	i-PRO	WV-S71300-F3		Unknown Status Edit
Test Dome on Stand	SmartAI Headquarters	Cohort Test NUC	Online	192.168.1.3	i-PRO	WV-S2236L		Compliant Status Edit
Test Multicens on Stand	SmartAI Headquarters	Cohort Test NUC	Online	192.168.1.2	i-PRO	WV-S8544L		Compliant Status Edit
Veracity WV-S1536LA	Sri Lanka Office	SriLanka Demo NUC	Online	192.168.50.221	i-PRO	WV-S1536L		Compliant Status Edit
Veracity WV-S2236L	Sri Lanka Office	SriLanka Demo NUC	Online	192.168.50.220	i-PRO	WV-S2236L		Compliant Status Edit
WP Street USA	WP, USA	USA Demo NUC	Online	fe80::d62dc5ff:fe04:2d70	i-PRO	WV-S1536L		Compliant Status Edit



Chat-GPT sprints to 100 million users



Source: World of Statistics

National Safety Institutes



<https://www.nist.gov/aisi>



<https://www.aisi.gov.uk>

Key challenges:

- A lack of commonly accepted definitions for AI safety, as well as safety capabilities and measurements of those capabilities, especially for frontier models and advanced AI agents and systems
- Underdeveloped testing, evaluation, validation and verification (TEVV) methods and best practices to provide holistic assessments of risk – from model capabilities to human-AI interaction to system system-level and societal-level impacts
- An absence of scientifically-established risk mitigations across the lifecycle of AI design and deployment
- An insufficient understanding of the relationship between model architecture and design and model behaviour and performance, especially after deployment
- Limited and ad hoc coordination around safety practices among industry, civil society, and national and international actors

Mission:

- Test advanced AI systems and inform policymakers about their risks;
- Foster collaboration across companies, governments, and the wider research community to mitigate risks and advance publicly beneficial research; and
- Strengthen AI development practices and policy globally.



NIST - Potential AI Harms

Harm to People

- Individual: Harm to a person's civil liberties, rights, physical or psychological safety, or economic opportunity.
- Group/Community: Harm to a group such as discrimination against a population sub-group.
- Societal: Harm to democratic participation or educational access.

Harm to an Organization

- Harm to an organization's business operations.
- Harm to an organization from security breaches or monetary loss.
- Harm to an organization's reputation.

Harm to an Ecosystem

- Harm to interconnected and interdependent elements and resources.
- Harm to the global financial system, supply chain, or interrelated systems.
- Harm to natural resources, the environment, and planet.

<https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>



the streamingguys.com.au

Pigeonhole – Q & A portal

TRUSTWORTHY AI – RESPONSIBLE AI – SAFE AI

AI Safety

Foundational model risks
impact on entire AI ecosystem, hence collective failure

Safe AI

- Harmful content, Value Alignment
- Disinformation, Hallucination
- Copyright, Data supply chain risk

Frontier AI risks of societal harms & public safety

Responsible AI

- Fair, Transparent, Explainable
- Accountable, Ethically developed
- No misuse of data in favor of ML

Traditional AI risks of user & organizations

Trustworthy AI

- Function as intended
- No unauthorized modifications
- Operated in a secure manner

NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

Ask your questions @ ceda.pigeonhole.at | Passcode: AISUMMIT |  @ceda_news | #AI

Australian AI Ethics Principles

Beneficial

- AI systems should benefit individuals, society and the environment.

Human-centred

- AI systems should respect human rights, diversity, and the autonomy of individuals.

Fairness

- AI systems should be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups.

Privacy and Security

- AI systems should respect and uphold privacy rights and data protection, and ensure the security of data.

Reliability and Safety

- AI systems should reliably operate in accordance with their intended purpose.

Explainability

- There should be transparency and responsible disclosure so people can understand when they are being significantly impacted by AI, and can find out when an AI system is engaging with them.

Contestability

- When an AI system significantly impacts a person, community, group or environment, there should be a timely process to allow people to challenge the use or outcomes of the AI system.

Accountability

- Those responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the outcomes of the AI systems, and human oversight of AI systems should be enabled.

Policy for the responsible use of AI in government

September 2024

Version 1.1



<https://architecture.digital.gov.au/responsible-use-of-AI-in-government>

Effective 1 Sep 2024

Policy Requirements:

1. Accountable officials

Agencies **must** designate accountability for implementing this policy to accountable official(s).

The responsibilities of the accountable officials are to:

- be accountable for implementation of this policy within their agencies
- notify the Digital Transformation Agency (DTA) where the agency has identified a new high-risk use case.
This information will be used by the DTA to build visibility and inform the development of further risk mitigation approaches.
- be a contact point for whole-of-government AI coordination
- engage in whole-of-government AI forums and processes
- keep up to date with changing requirements as they evolve over time.

30 Nov
2024

2. AI transparency statement

Agencies **must** make publicly available a statement outlining their approach to AI adoption and use

The statement **must** be reviewed and updated annually or sooner, should the agency make significant changes to their approach to AI.

This statement must provide the public with relevant information about the agency's use of AI including information on:

- compliance with this policy
- measures to monitor effectiveness of deployed AI systems
- efforts to protect the public against negative impacts.

28 Feb
2025

Policy for the responsible use of AI in government

September 2024

Version 1.1



<https://architecture.digital.gov.au/responsible-use-of-AI-in-government>

AI Transparency Statement

28 Feb
2025

MUST provide the following information regarding their use of AI:

- the intentions behind why the agency uses AI or is considering its adoption
- classification of AI use according to [usage patterns and domains](#)
- classification of use where the public may directly interact with, or be significantly impacted by, AI without a human intermediary or intervention
- measures to monitor the effectiveness of deployed AI systems, such as governance or processes
- compliance with applicable legislation and regulation
- efforts to identify and protect the public against negative impacts
- compliance with each requirement under the [Policy for responsible use of AI in government](#)
- when the statement was most recently updated.

AI Definition

“An AI system is a machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments. Different AI systems vary in their levels of autonomy and adaptiveness after deployment.” [OECD definition]

Used in [Policy for the responsible use of AI in government]

“a technical and scientific field devoted to the engineered system that generates outputs such as content, forecasts, recommendations or decisions for a given set of human-defined objectives” [[ISO/IEC 22989:2022](#)]

National framework for the assurance of artificial intelligence in government

A joint approach to safe and responsible AI by the Australian, state and territory governments.

21 June 2024



AI Principles

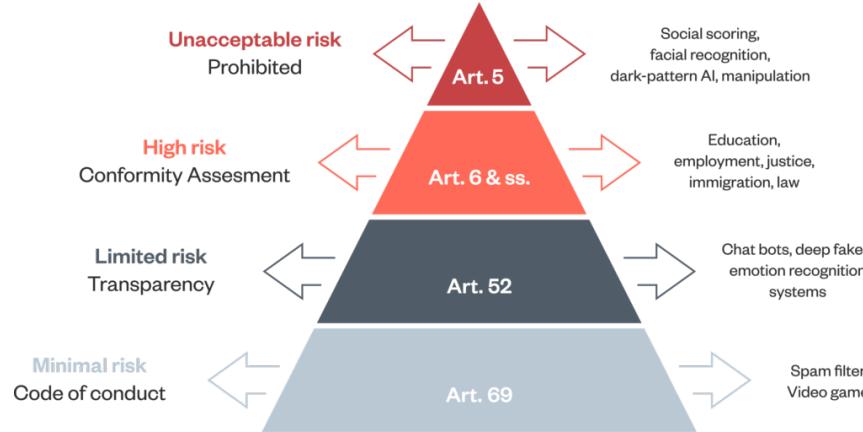


Image Source: <https://www.ceda.com.au/researchandpolicies/research/data-digital-economy/artificial-intelligence-principles-to-practice>

Voluntary AI Safety Standard Guardrails

1. Establish, implement and publish an accountability process including governance, internal capability and a strategy for regulatory compliance.
2. Establish and implement a risk management process to identify and mitigate risks.
3. Protect AI systems, and implement data governance measures to manage data quality and provenance.
4. Test AI models and systems to evaluate model performance and monitor the system once deployed.
5. Enable human control or intervention in an AI system to achieve meaningful human oversight across the life cycle.
6. Inform end-users regarding AI-enabled decisions, interactions with AI and AI-generated content.
7. Establish processes for people impacted by AI systems to challenge use or outcomes.
8. Be transparent with other organisations across the AI supply chain about data, models and systems to help them effectively address risks.
9. Keep and maintain records to allow third parties to assess compliance with guardrails.
10. Engage your stakeholders and evaluate their needs and circumstances, with a focus on safety, diversity, inclusion and fairness.

EU Artificial Intelligence Act



Effective



Enforced



✖ Prohibited AI (Late 2024)

- Systems for social credit scoring
- Employment and educational emotion recognition systems
- AI leveraging people's vulnerabilities (e.g., age, disability)
- Manipulation of behavior and undermining free will
- Indiscriminate collection of facial images for facial recognition
- Biometric categorization systems involving sensitive traits
- Particular predictive policing uses
- Law enforcement employing real-time biometric
- identification in public (except in restricted, pre-approved scenarios)

⚠ High-Risk AI (Mid 2026)

- Medical devices & Automobiles
- Hiring, human resources, and labor supervision
- Instruction and professional education
- Shaping political elections and the voting populace
- Entry to amenities (like insurance, banking, credit, benefits, etc.)
- Supervising crucial infrastructure (like water, gas, electricity, etc.)
- Systems for recognizing emotions and identifying individuals via biometrics
- Policing, regulating borders, migration, and asylum procedures and conducting legal affairs
- Particular merchandise or safety elements within particular products

✓ No longer high-risk

If not intended to be used as a safety component of a product covered by the Union harmonisation legislation.*

- Performs a narrow procedural task
- Intended to improve the result of a previously completed human activity
- Not meant to replace or influence the previously completed human assessment, without proper human review
- perform a preparatory task to an assessment relevant to high-risk use-cases.



High risk AI providers must:

- Establish a **risk management system** throughout the high risk AI system's lifecycle;
- Conduct **data governance**, ensuring that training, validation and testing datasets are relevant, sufficiently representative and, to the best extent possible, free of errors and complete according to the intended purpose.
- Draw up **technical documentation** to demonstrate compliance and provide authorities with the information to assess that compliance.
- Design their high risk AI system for **record-keeping** to enable it to automatically record events relevant for identifying national level risks and substantial modifications throughout the system's lifecycle.
- Provide **instructions for use** to downstream deployers to enable the latter's compliance.
- Design their high risk AI system to allow deployers to implement **human oversight**.
- Design their high risk AI system to achieve appropriate levels of **accuracy, robustness, and cybersecurity**.
- Establish a **quality management system** to ensure compliance.



NIST AI 100-1



Artificial Intelligence Risk Management Framework (AI RMF 1.0)

NIST NATIONAL INSTITUTE OF
STANDARDS AND TECHNOLOGY
U.S. DEPARTMENT OF COMMERCE

NIST AI Trustworthiness Characteristics



<https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>

KJR
est.
1997

QLD Gov AI governance policy



- Agencies must use a consistent and evidence-based process which incorporates an ethical framework to evaluate their transparency, accountability, and risk associated with the AI lifecycle
 - FAIRAI: Foundational artificial intelligence risk assessment guideline provides a comprehensive process for assessing AI risk within a specific context including alignment of tasks and responsibilities to relevant organisational roles
- Agencies must establish AI governance arrangements based on ISO 38507 (Governance implications of the use of artificial intelligence by organisations)
 - Aligns with existing IT governance processes
 - Provides specific advice regarding AI, including
 - usage governance
 - data governance



Foundational artificial intelligence risk assessment framework

Final

V1.0.0

September 2024

OFFICIAL - Public

FAIR Framework



Foundational artificial intelligence risk assessment framework

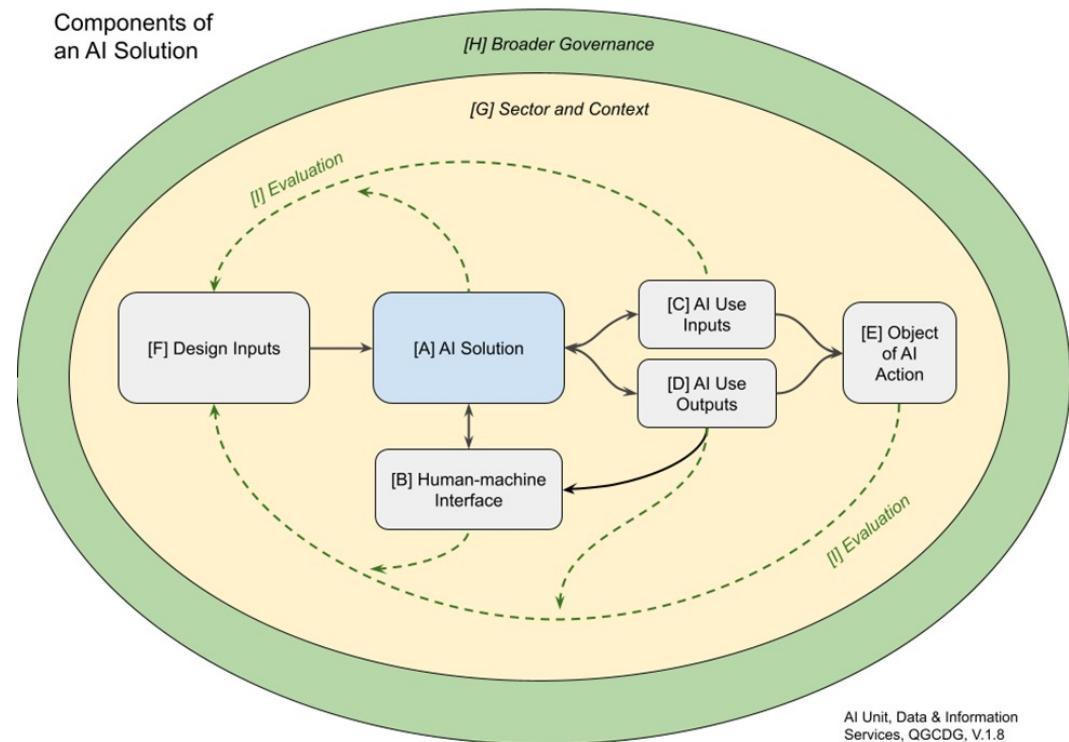
Final

V1.0.0

September 2024

OFFICIAL - Public

Part A – Components Analysis



Foundational artificial intelligence risk assessment framework

Final
V1.0.0
September 2024
OFFICIAL - Public

Part B – Components Analysis

- Human, societal, and environmental wellbeing
- Human-Centred Values
- Fairness
- Privacy Protection and Security
- Reliability and Safety
- Transparency and Explainability
- Contestability
- Accountability

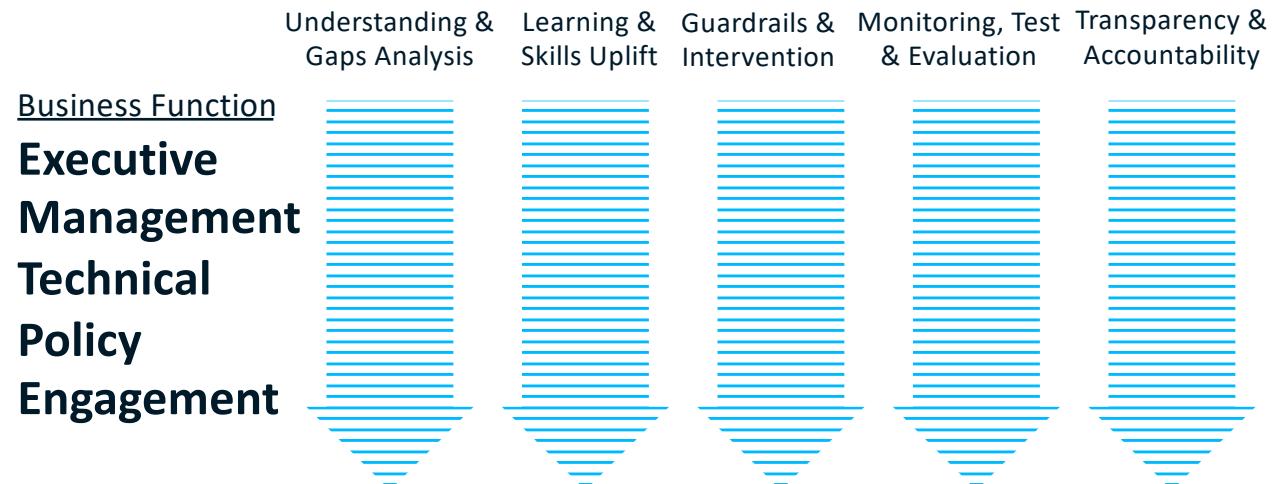


Foundational artificial intelligence risk assessment framework

Final
V1.0.0
September 2024
OFFICIAL - Public

Part C – Controls for AI Risks

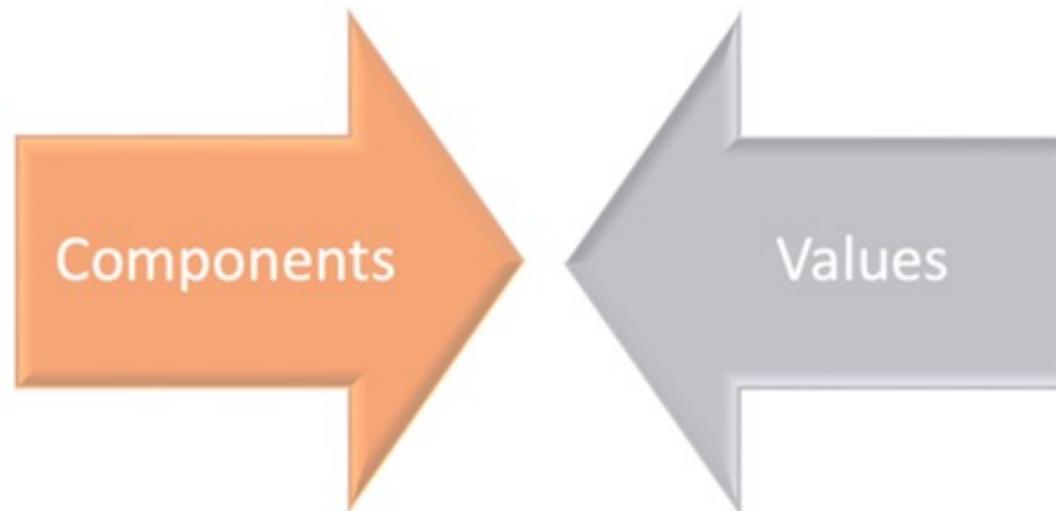
Control Areas



How FAIRA works

Explain components of AI solution

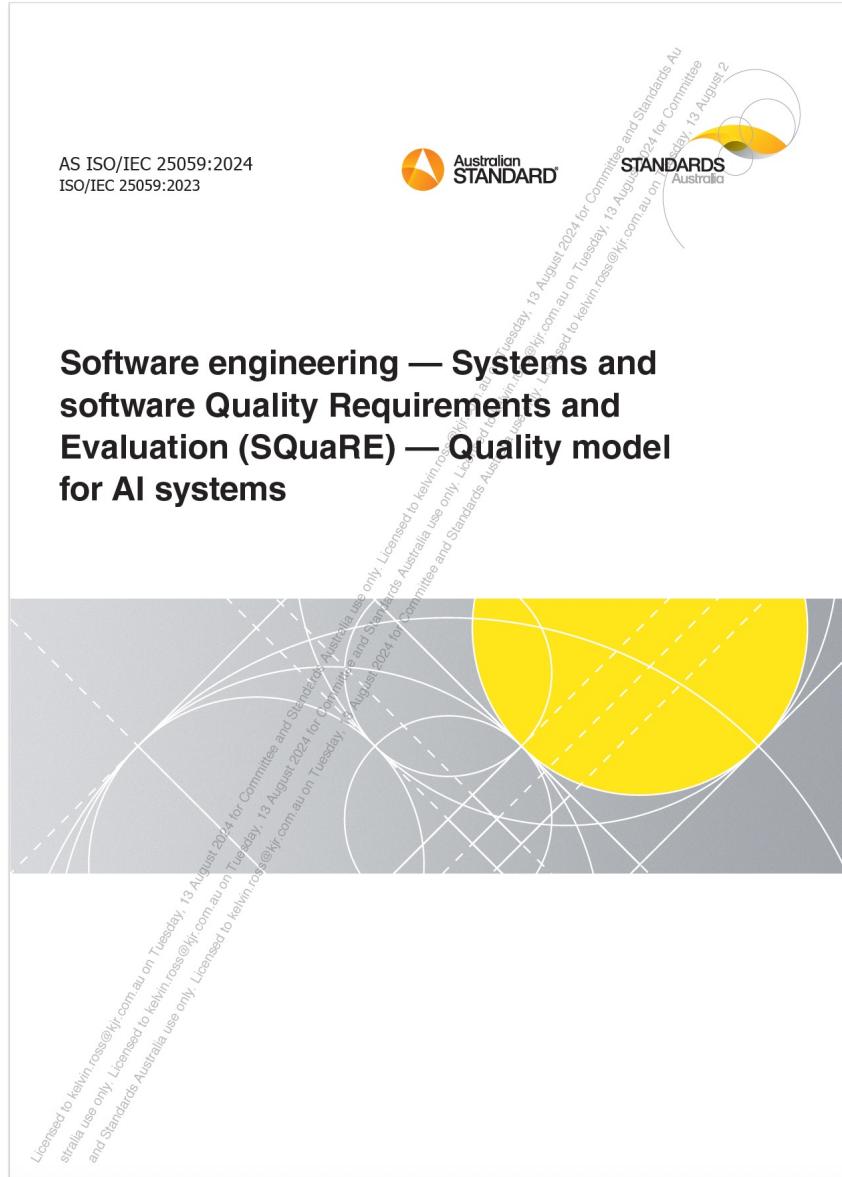
- [A] AI Solution
- [B] HMI
- [C] AI use inputs
- [D] AI use outputs
- [E] Object of AI Action
- [F] Design inputs
- [G] Sector and Context
- [H] Broader governance
- [I] Evaluation



- Human, societal, and environmental wellbeing
- Human-Centred Values
- Fairness
- Privacy Protection and Security
- Reliability and Safety
- Transparency and Explainability
- Contestability
- Accountability

Explain effect of AI solution on values alignment

AI Quality Characteristics



ISO/IEC 25010:2023 - System and software quality models

Software Product Quality								
Functional Suitability	Performance Efficiency	Compatibility	Interaction Capability	Reliability	Security	Maintainability	Flexibility	Safety
Functional Completeness	Time Behaviour	Co-Existence	Appropriateness Recognizability	Faultlessness	Confidentiality	Modularity	Adaptability	Operational Constraint
Functional Correctness	Resource Utilization	Interoperability	Learnability Operability	Availability Fault Tolerance	Integrity Non-Repudiation	Reusability Analysability	Scalability Installability	Risk Identification
Functional Appropriateness	Capacity		User Error Protection User Engagement Inclusivity User Assistance Self-Descriptiveness	Recoverability Robustness	Accountability Authenticity Resistance Intervenability	Modifiability Testability	Replaceability	Fail Safe Hazard Warning Safe Integration
iso25000.com								

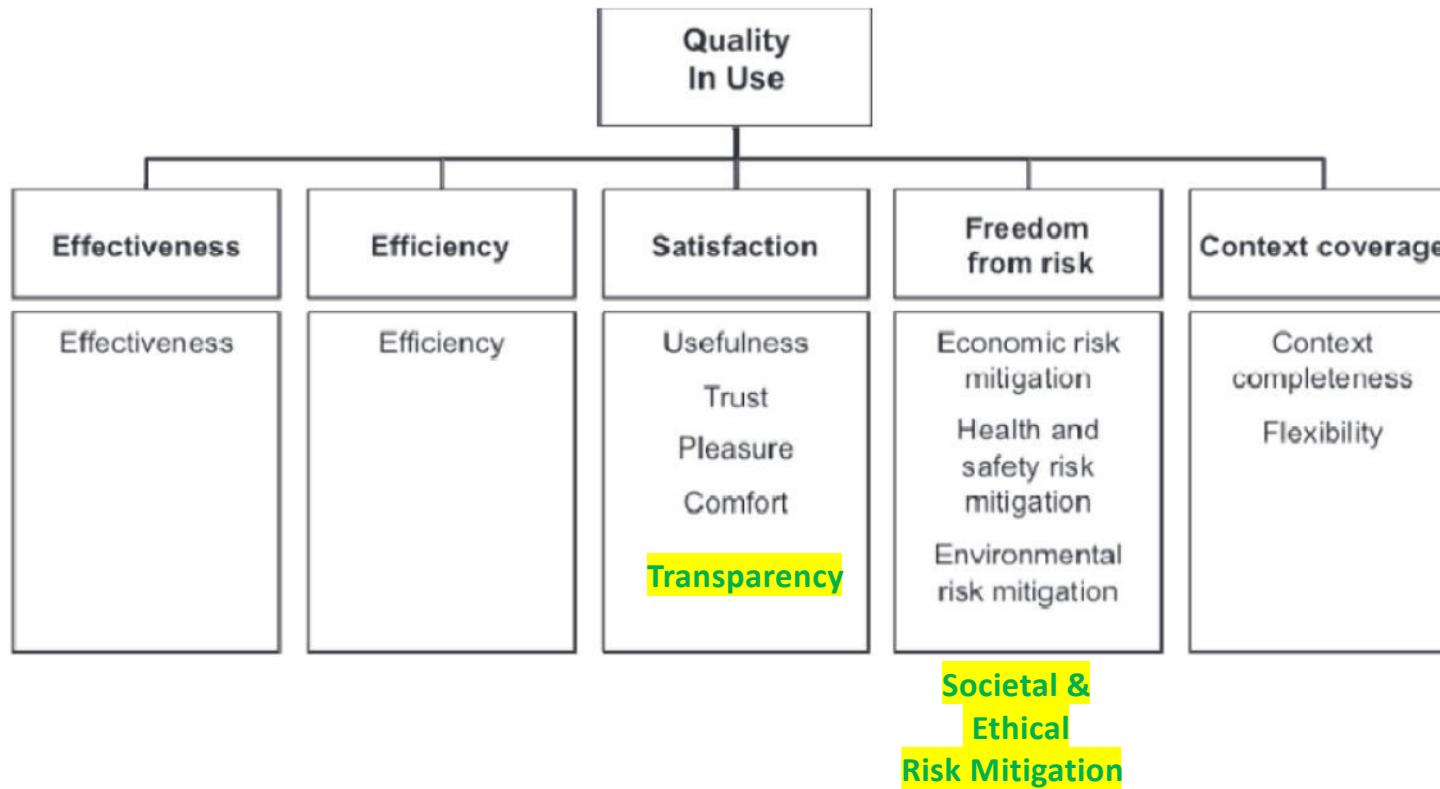
USER
CONTROLLABILITY

TRANSPARENCY

ISO/IEC 25059:2024 - Quality model for AI systems



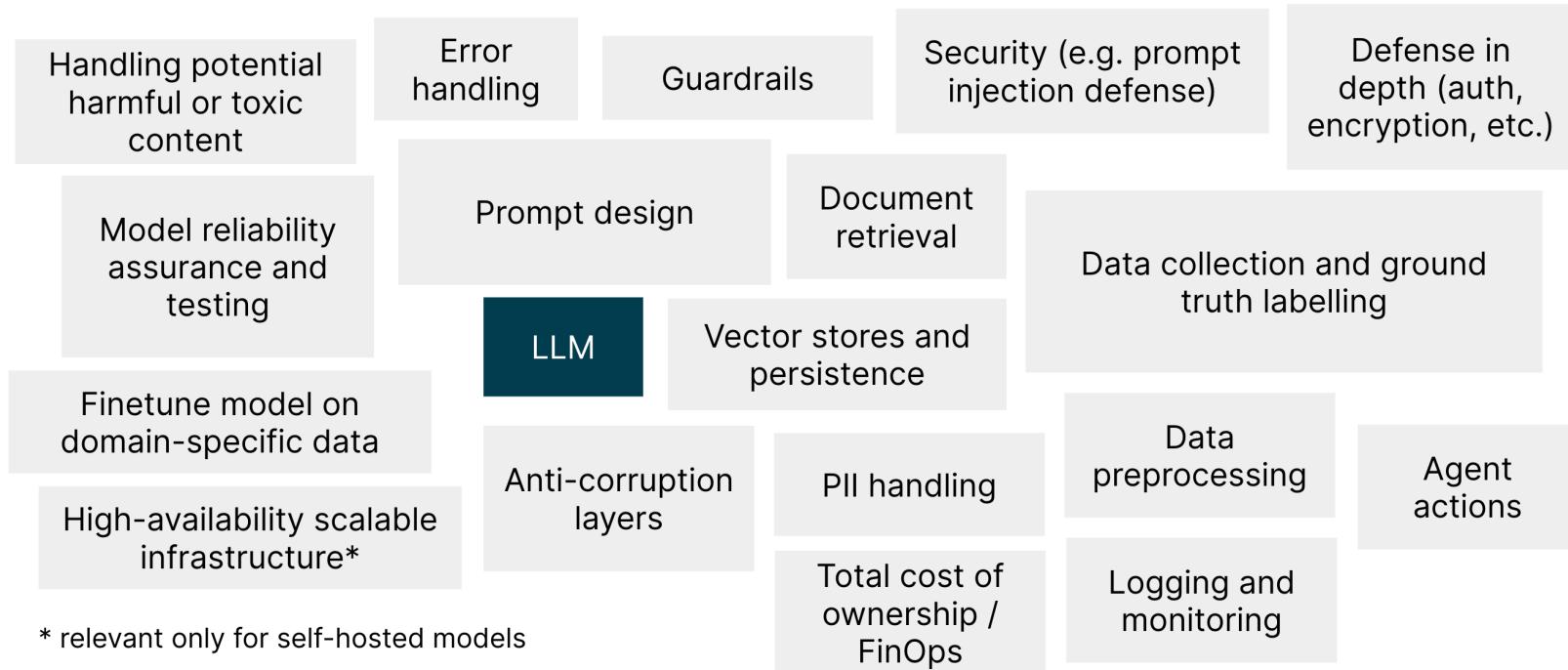
ISO/IEC 25010:2023 - System and software quality models



ISO/IEC 25059:2024 - Quality model for AI systems

Architecting LLM applications

The language model is just one part of the technical architecture

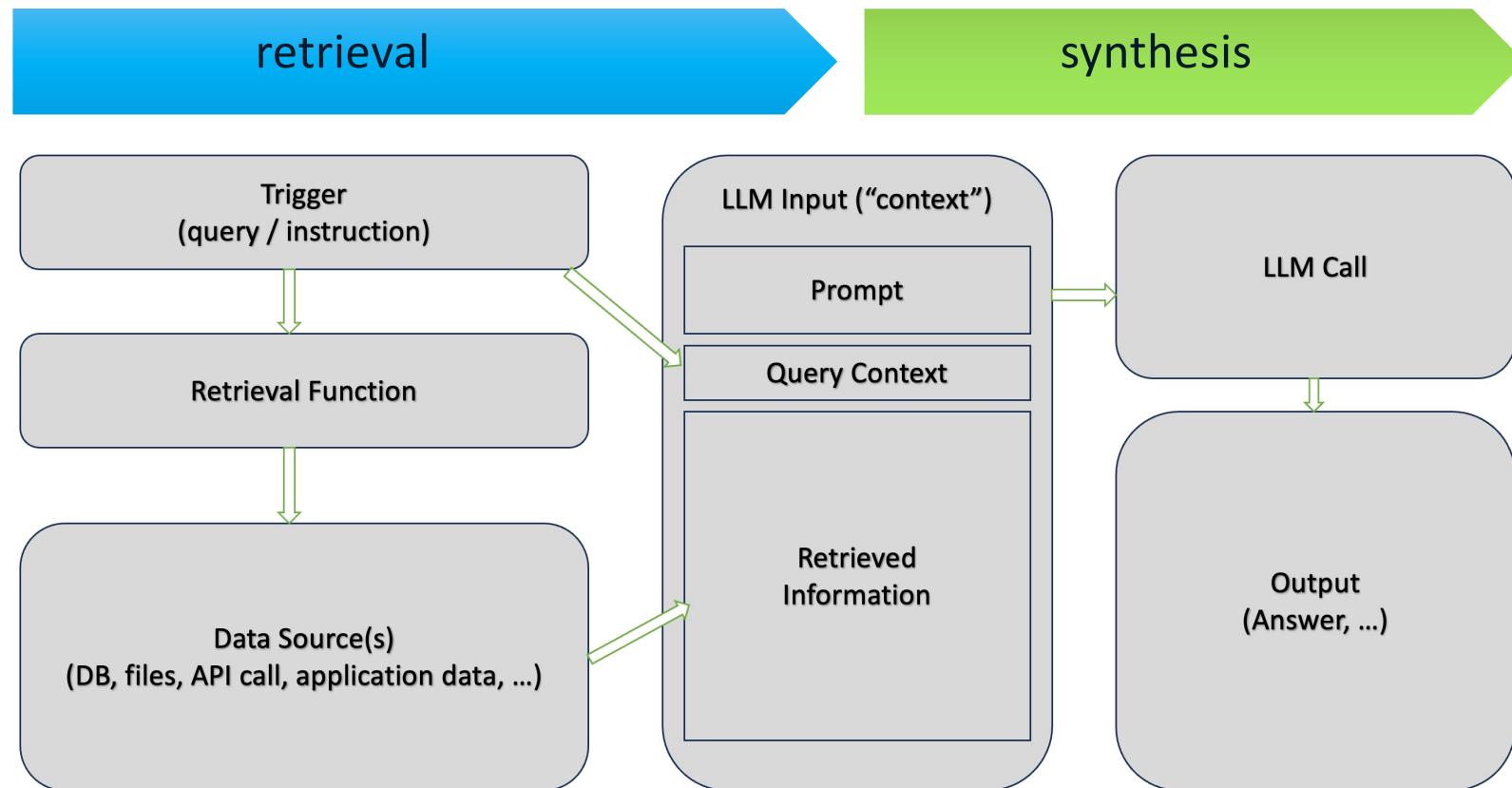


Adapted from: [Machine Learning: The High Interest Credit Card of Technical Debt \(Google\)](https://www.manning.com/whitepapers/machine-learning-the-high-interest-credit-card-of-technical-debt-google)

<https://martinfowler.com/articles/engineering-practices-llm.html>

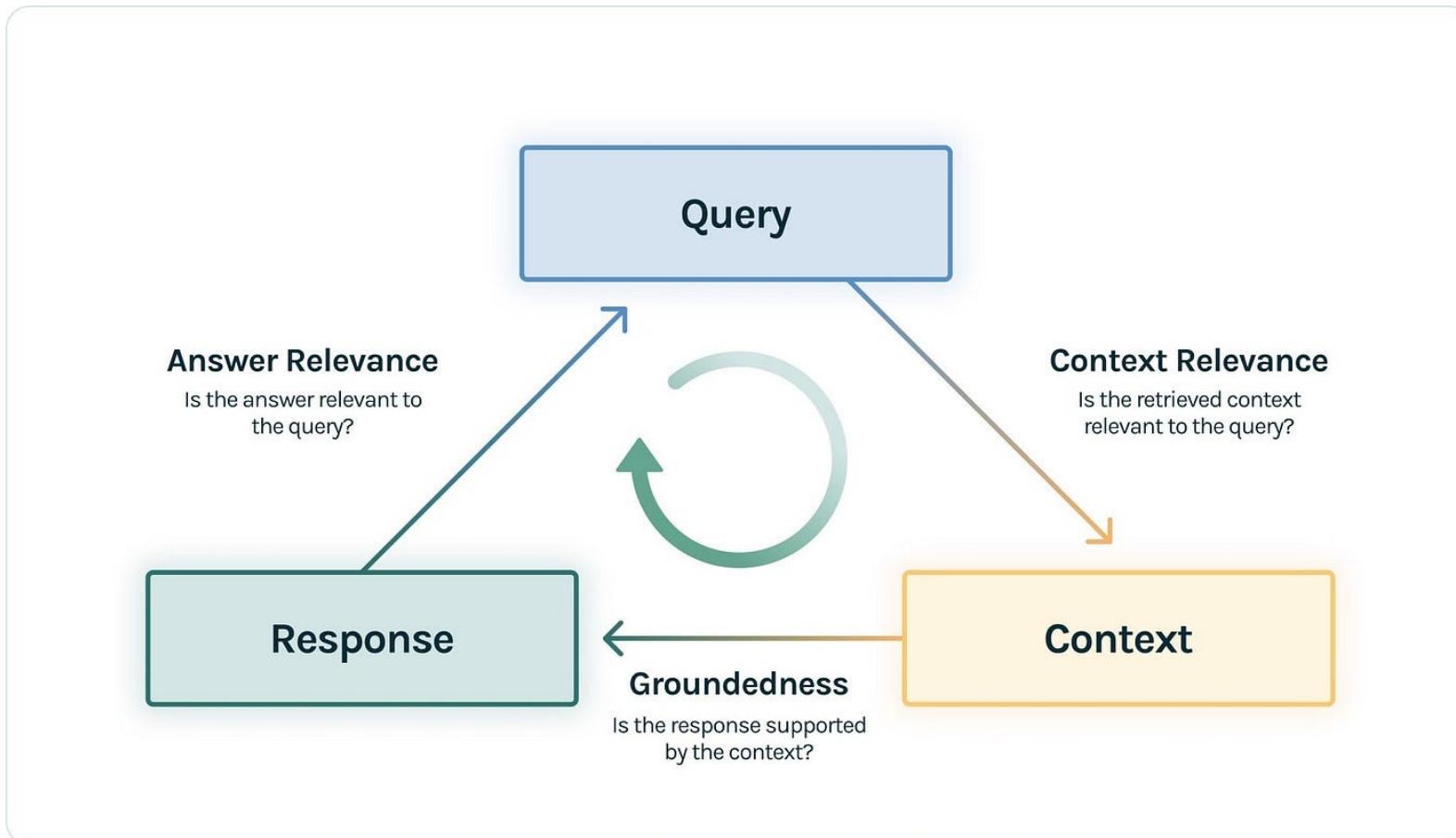


RAG Workflow



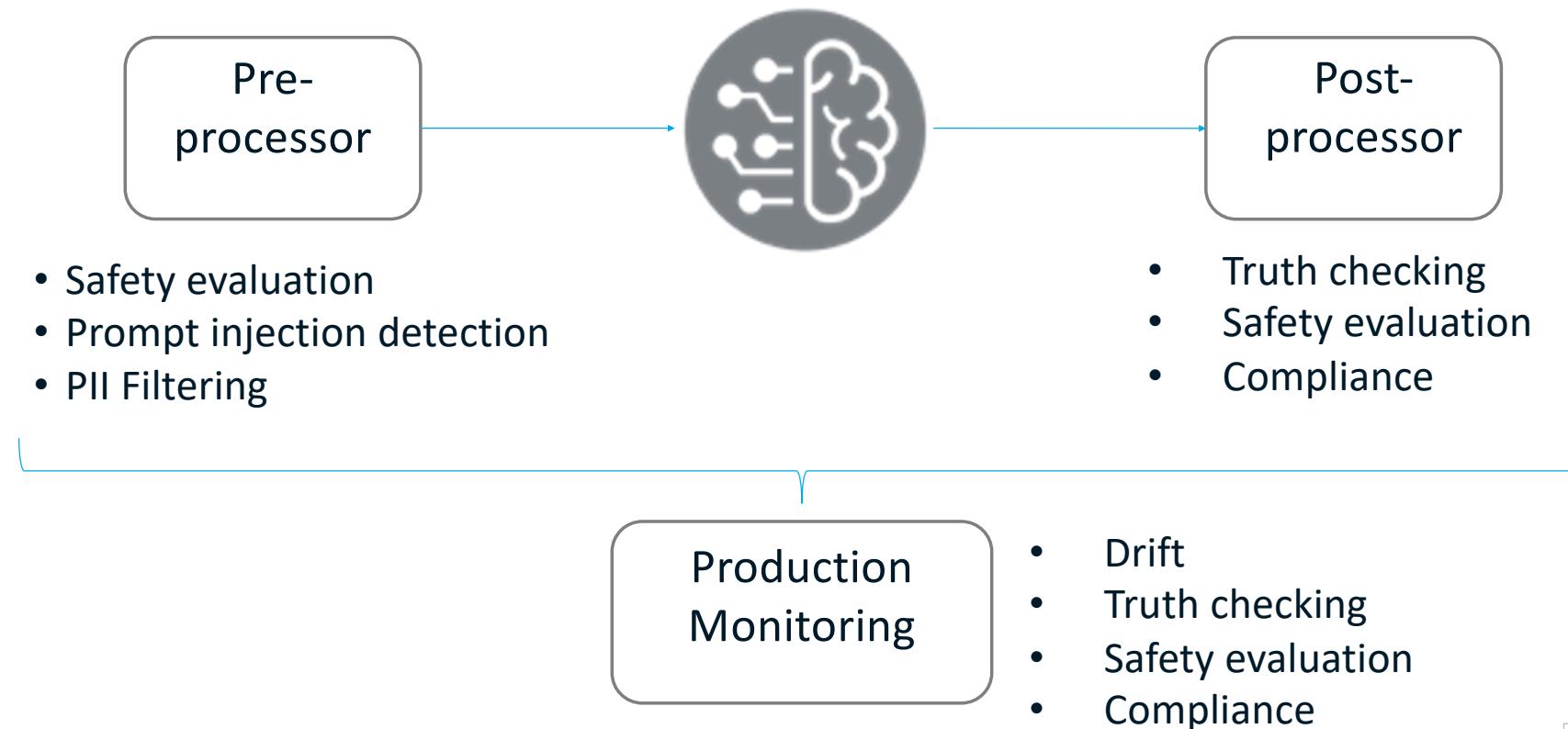
Source: <https://techcommunity.microsoft.com/t5/fasttrack-for-azure/grounding-langs/ba-p/3843857>

The RAG Triad



Source : https://www.trulens.org/trulens_eval/core_concepts_rag_triad/

Guardrails



ISO TS 29119-11 : Testing of AI-based Systems



ISO/IEC JTC 1/SC 42/JWG 2 N 244

ISO/IEC JTC 1/SC 42/JWG 2 "Joint Working Group ISO/IEC JTC1/SC 42 - ISO/IEC JTC1/SC 7 : Testing of AI-based systems"
Convenorship: BSI
Convenors: Reid Stuart Dr, Smith Adam Leon Mr



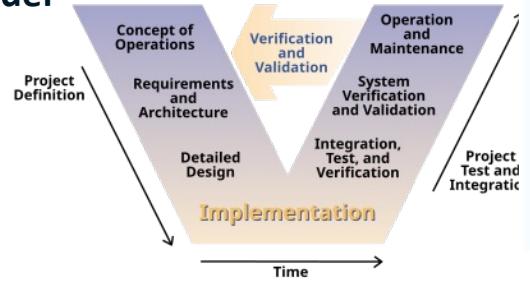
ISO TS 29119-11 V0.23

ISO/IEC JTC 1/SC 7
Software & Systems Engineering

**Risk-Based
Testing**



V-Model

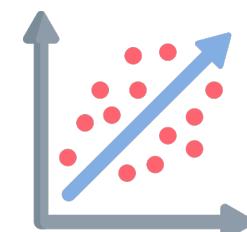


Specific Test Techniques
Evals / Test Oracles
Data Sampling
Data Quality Testing / Model Testing
Validation Data Leakage
Performance Metrics
Uncertainty
Coverage

ISO/IEC JTC 1/SC 42
Artificial Intelligence



Data Science

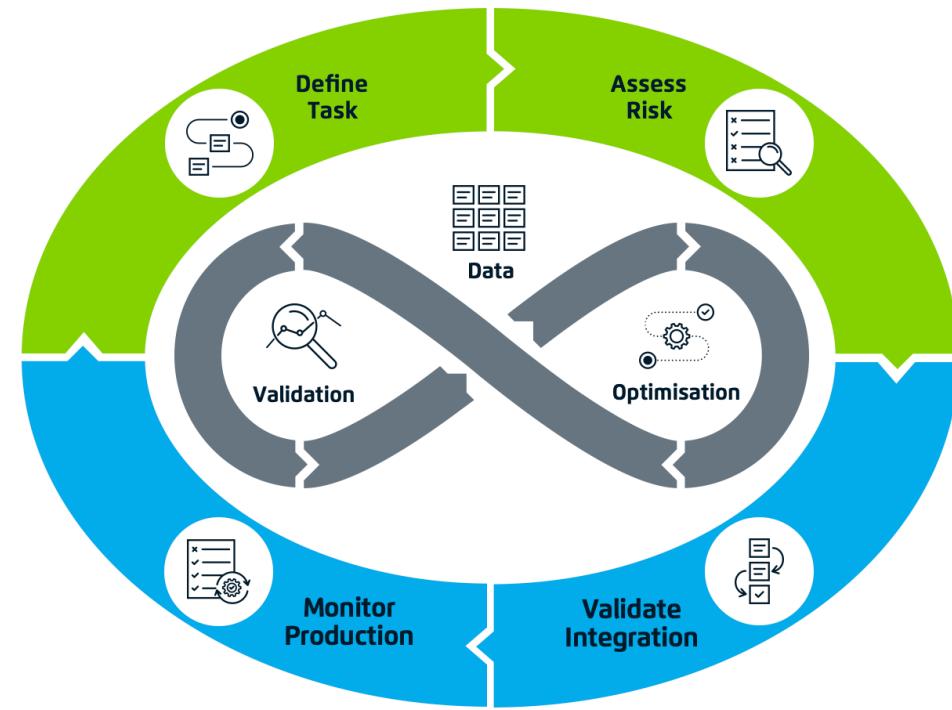


Statistics

KJR
est. 1997



Validation Driven Machine Learning (VDML) is a methodology developed by KJR to guide development of robust and reliable Machine Learning (ML) models. VDML emphasises understanding the risks inherent within the application context and the limitations that arise from the available data and model building processes applying iterative validation and optimisation methods to deliver an acceptable solution which can be integrated and governed within a real-world context.



DEFINE CONTEXT

Define the goals of applying machine learning to a specific problem area, being sure to include the data being used, the context of use (historical analytics vs live decision support) and expected benefits from a range of different stakeholders. Given this context, assess risks, including the impacts of potential failure, the required governance processes.

RESOLVE LIMITATIONS

Direct use of pre-built models or naïve approaches to machine learning can lead to unreliable performance. Key to validating and optimising model performance is the selection of training and testing data sets which are close to real world usage, and detailed error analysis which can uncovers underlying faults and limitations.

GOVERN BEHAVIOUR

Track the integrity of the model through build, deploy and operation, monitoring for residual risks, model drift / sabotage, identifying opportunities for further optimisation and risk reduction.



PODCAST SERIES

TRUST YOUR AI
WITH

VDML®

PRACTICAL CONVERSATIONS ABOUT AI ASSURANCE

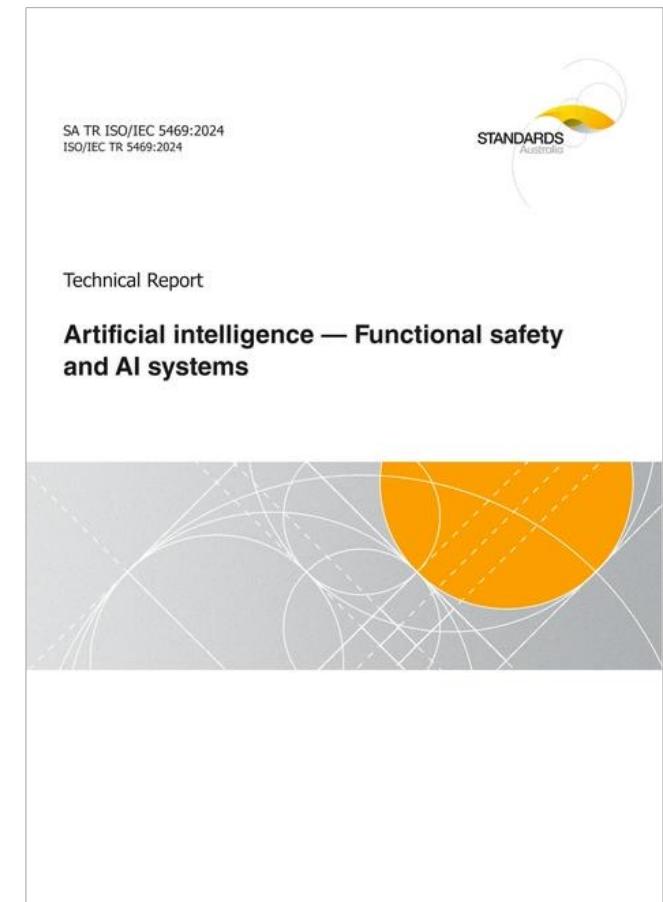


LISTEN NOW

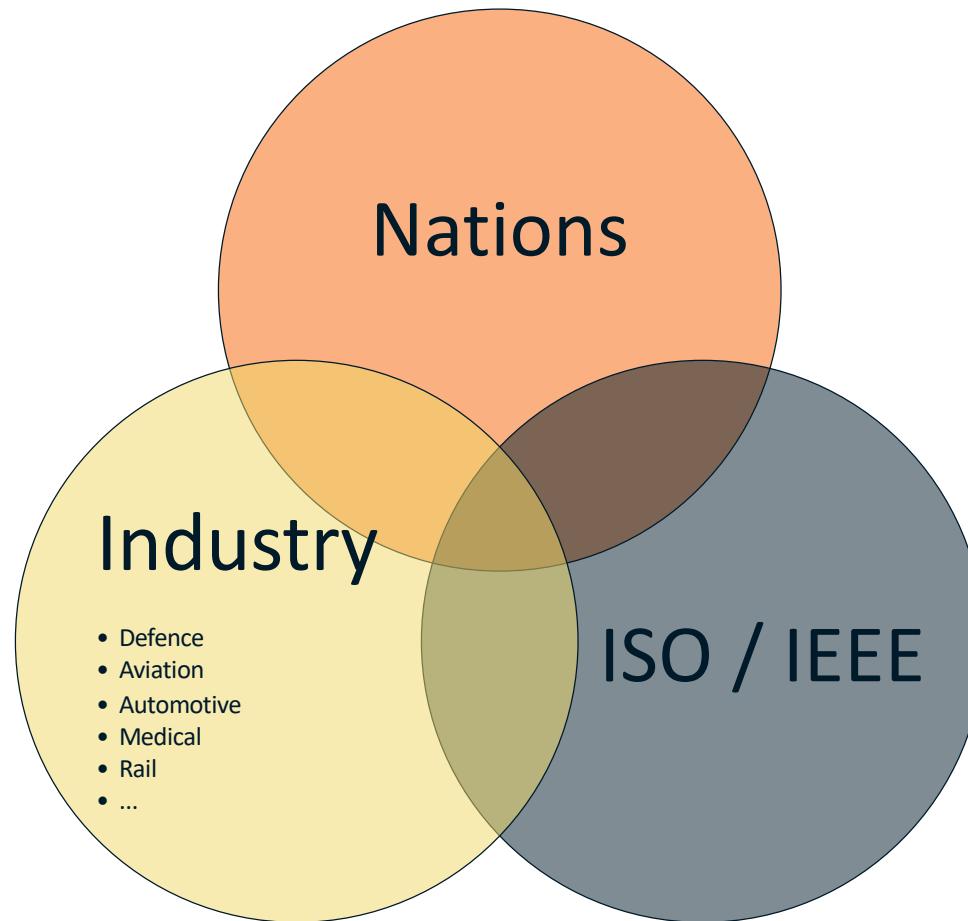


KJR
est.
1997

Other ISO Standards



Policy / Standards Development



Key Takeaways

- Emerging Obligations
 - Do we have the approach, skills and capacity to provide the requested assurance?
- Overall governance is not that unfamiliar for systems assurance folk
 - What role will the Systems Safety Assurance people play in this expanded AI Safety space?
- What is different about AI?
 - Definition of scope – what is AI?
 - Understanding of risks
 - Mitigations and controls



Trust your AI

Kelvin Ross

kelvin.ross@kjr.com.au

Mark Pedersen

mark.pedersen@kjr.com.au

kjr.com.au

