



This is not an ADB material. The views expressed in this document are the views of the author/s and/or their organizations and do not necessarily reflect the views or policies of the Asian Development Bank, or its Board of Governors, or the governments they represent. ADB does not guarantee the accuracy and/or completeness of the material's contents, and accepts no responsibility for any direct or indirect consequence of their use or reliance, whether wholly or partially. Please feel free to contact the authors directly should you have queries.

# AI Ethics and Governance in Practice Programme

Professor David Leslie

Director of Ethics and Responsible Innovation Research

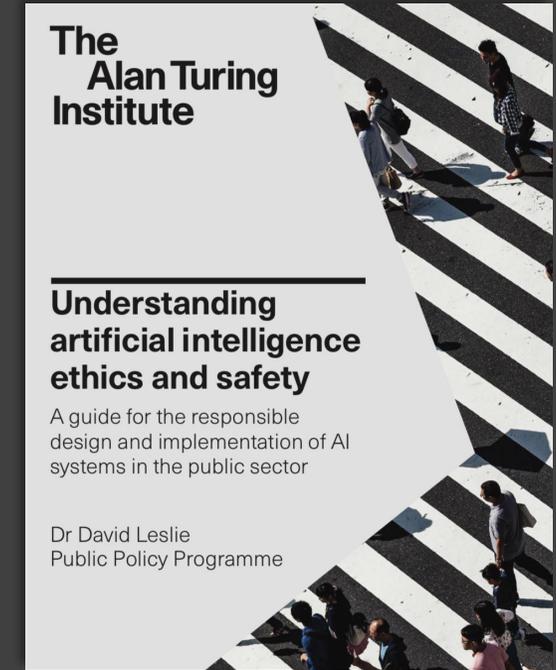
Professor of Ethics, Technology and Society (QMUL)

July 2024



# UK National Guidance on AI ethics and safety in the public sector

- Funded by EPSRC in 2018 to do fieldwork with Ministry of Justice to develop ethics framework which formed the normative foundation of the guidance
- Published June 2019 in collaboration with the Office for AI and GDS as part of the government's official guide to using AI in the public sector with Ministerial approval
- Now the **most accessed and cited public sector AI ethics guidance in the world.**
- Credited with initiating the move from “principles to practice in the field of AI policy and governance



Government  
Digital Service



Office for  
Artificial  
Intelligence

# UK National Guidance on AI ethics and safety in the public sector

Guidance has since been put into use by:



## The Alan Turing Institute

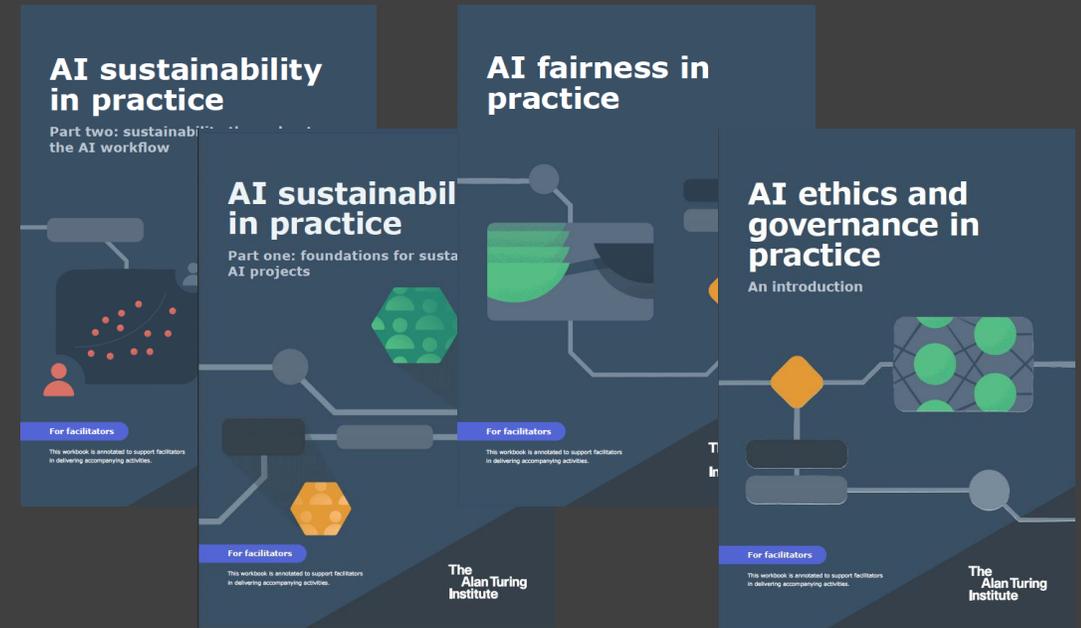
### Understanding artificial intelligence ethics and safety

A guide for the responsible design and implementation of AI systems in the public sector

Dr David Leslie  
Public Policy Programme

# AI Ethics and Governance in Practice Programme

- 2021 UK National AI Strategy made expansion of public sector guidance a priority
- Supported by the Office for AI and EPSRC funding to construct a series of 8 practice-based workbooks—co-designed and piloted with civil servants
- This is now complimented by an interactive digital platform to facilitate accessibility, uptake, and participation across government



# AI Ethics and Governance in Practice Programme

## Human Rights, Democracy, and the Rule of Law Assurance Framework for AI Systems:

A proposal prepared for the Council of Europe's Ad hoc Committee on Artificial Intelligence

The Alan Turing Institute

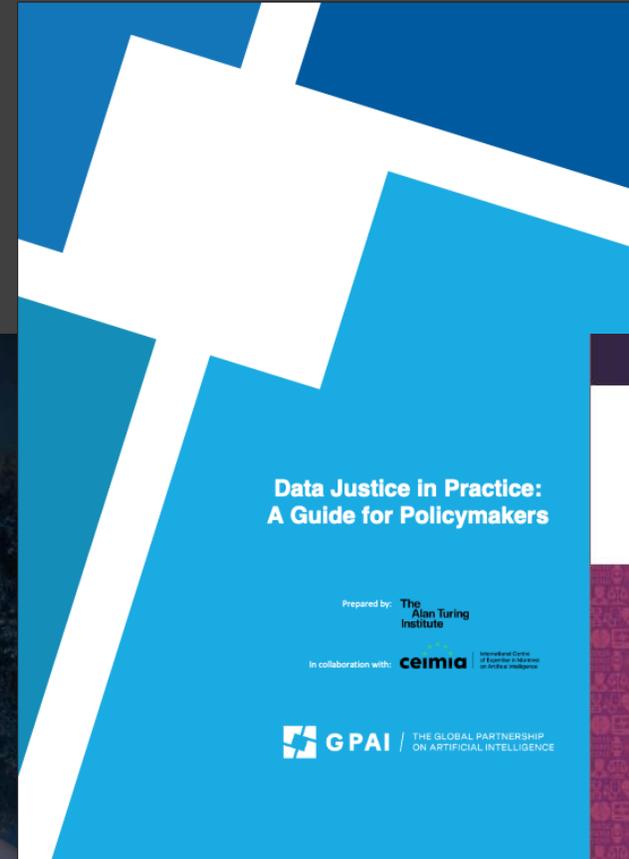


unesco

Recommendation on  
**the Ethics of Artificial Intelligence**

Adopted on 23 November 2021

The cover features a woman's face in profile, a globe with data points, and various small images representing AI applications in different fields.



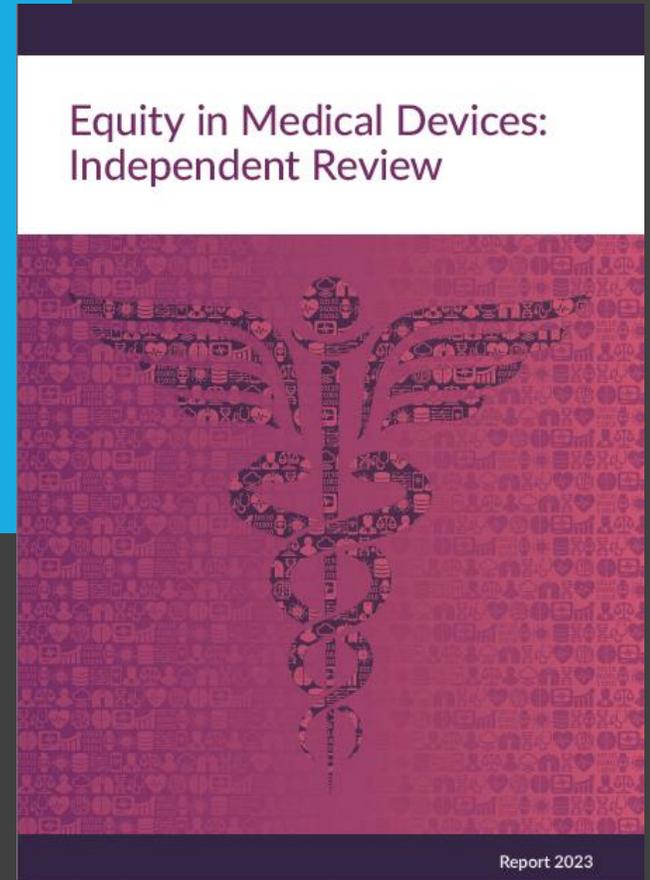
**Data Justice in Practice:  
A Guide for Policymakers**

Prepared by: **The Alan Turing Institute**

In collaboration with: **ceima** International Centre for Excellence in Research on Artificial Intelligence

**GPAI** | THE GLOBAL PARTNERSHIP ON ARTIFICIAL INTELLIGENCE

The cover has a blue background with white geometric shapes.



**Equity in Medical Devices:  
Independent Review**

Report 2023

The cover features a red background with a large, stylized caduceus symbol composed of many small icons.

# Programme Curriculum



1

## AI Ethics and Governance in Practice: An Introduction

*Multiple Domains*



2

## AI Sustainability in Practice Part One

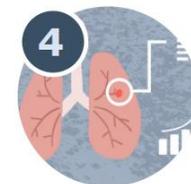
*AI in Urban Planning*



3

## AI Sustainability in Practice Part Two

*AI in Urban Planning*



4

## AI Fairness in Practice

*AI in Healthcare*



5

## Responsible Data Stewardship in Practice

*AI in Policing and Criminal Justice*



6

## AI Safety in Practice

*AI in Transport*



7

## AI Transparency and Explainability in Practice

*AI in Social Care*



8

## AI Accountability in Practice

*AI in Education*



# The CARE and Act Framework

## – Consider Context

- Think about the conditions and circumstances surrounding your research project.

## – Anticipate impacts

- Describe and analyse the impacts, intended or not, that might arise from your project.

## – Reflect on purposes, positionality, and power

- Reflect on the goals of, motivations for and potential implications of the research and engage in reflexive practices that scrutinise the way potential perspectival limitations and power imbalances can exercise influence on the equity and integrity of research projects and on the associated uncertainties, areas of ignorance, assumptions, framings, and questions.

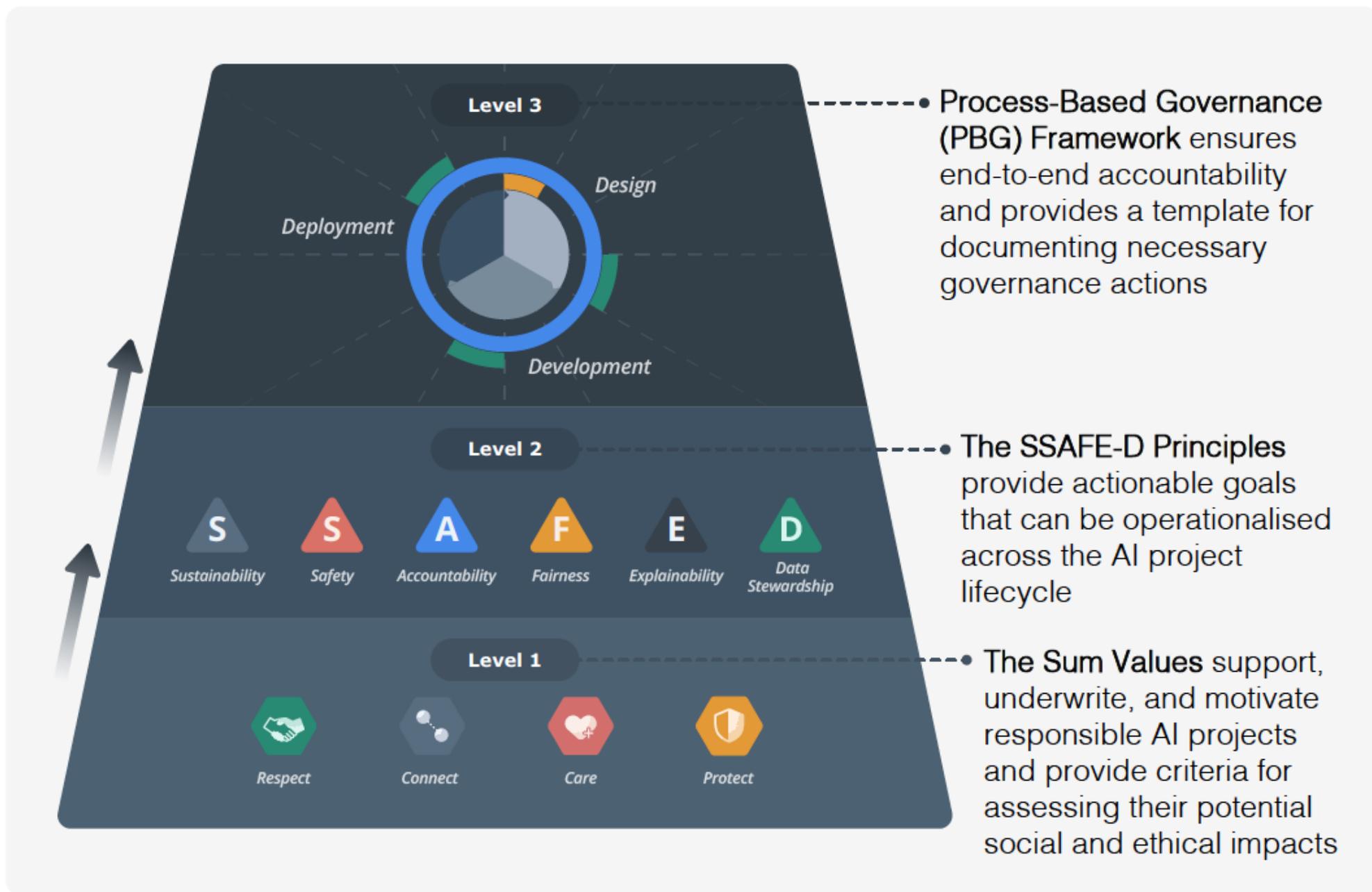
## – Engage inclusively

- Open up such visions, impacts, and questioning to broader deliberation, dialogue, engagement, and debate in an inclusive way.

## – Act responsibly and transparently

- Use these processes to influence the direction and trajectory of the research and innovation process itself. Produce research that is both scientifically and ethically justifiable. (EPSRC, 2013, expanded) and provide transparent documentation, following best governance, self-assessment, and reporting practices.

# Three building blocks of a responsible and trustworthy AI project lifecycle



## SUM Values: Responding to the ethical concerns raised by real-world hazards and risks

Risks that Emerge From the Use of AI/ML Technologies	Related Ethical Implications and Concerns	SUM Values
 <p>Loss of autonomy</p>	 <p><b>Agency</b> Human agency, dignity, and individual flourishing</p>	 <p><b>Respect</b> the dignity of individual persons</p>
 <p>Loss of interpersonal connection and empathy</p>	 <p><b>Interaction</b> Solidarity, communication, and integrity of social interaction</p>	 <p><b>Connect</b> with each other sincerely, openly, and inclusively</p>
 <p>Poor quality or hazardous outcomes</p>	 <p><b>Wellbeing</b> Individual, communal, and biospheric wellbeing</p>	 <p><b>Care</b> for the wellbeing of each and all</p>
 <p>Bias, injustice, inequality, and discrimination</p>	 <p><b>Justice</b> Justice, equity, and the common good</p>	 <p><b>Protect</b> the priorities of social values, justice, and the public interest</p>

## SSAFE-D Principles: Operationalising top-level normative goals



Achieving this goal requires assuring AI projects being developed with continuous sensitivity to real-world impacts.



Achieving this goal requires an AI system to be technically accurate, reliable, secure, and robust.



Achieving this goal requires assuring the project's end-to-end answerability and auditability.



Achieving this goal requires assuring a minimum threshold of discriminatory non-harm and bias mitigation.



Achieving this goal requires the ability to explain and justify AI project processes and AI-supported outcomes.



Achieving this goal requires data quality, integrity, protection, and privacy to be assured.

# 12 Process Based Governance (PBG) Framework: Evidencing and documenting necessary governance actions

## Sustainability

### SEP (Stakeholder Engagement Process)

Process facilitating the uptake of proportionate stakeholder engagement and input throughout the AI lifecycle. The SEP enables a contextually informed understanding of the social environment and human factors that may be impacted by, or may impact, individual AI projects.<sup>[58]</sup>

### SIA (Stakeholder Impact Assessment)

Process facilitating the iterative evaluation of the social impact and sustainability of individual AI projects, as well as the corroboration of these potential impacts in dialogue with stakeholders, when appropriate.

## Accountability

### PBG Log

Live document outlining governance actions, relevant team members and roles involved in each action, timeframes for follow-up actions, and logging protocols, for individual AI projects.<sup>[59]</sup>

## Explainability

### EAM (Explainability Assurance Management)

Iterative process aimed to facilitate the implementation and evaluation of transparency and explainability assurance activities across the project lifecycle and assist in providing clarification of AI system outputs to a range of impacted stakeholders.

## Safety

### SSA & RM (Safety Self-Assessment and Risk Management)

Process facilitating the evaluation of how AI projects align with safety objectives through the iterative identification and documentation of potential safety risks across the lifecycle and assurance actions implemented to address these.

## Fairness

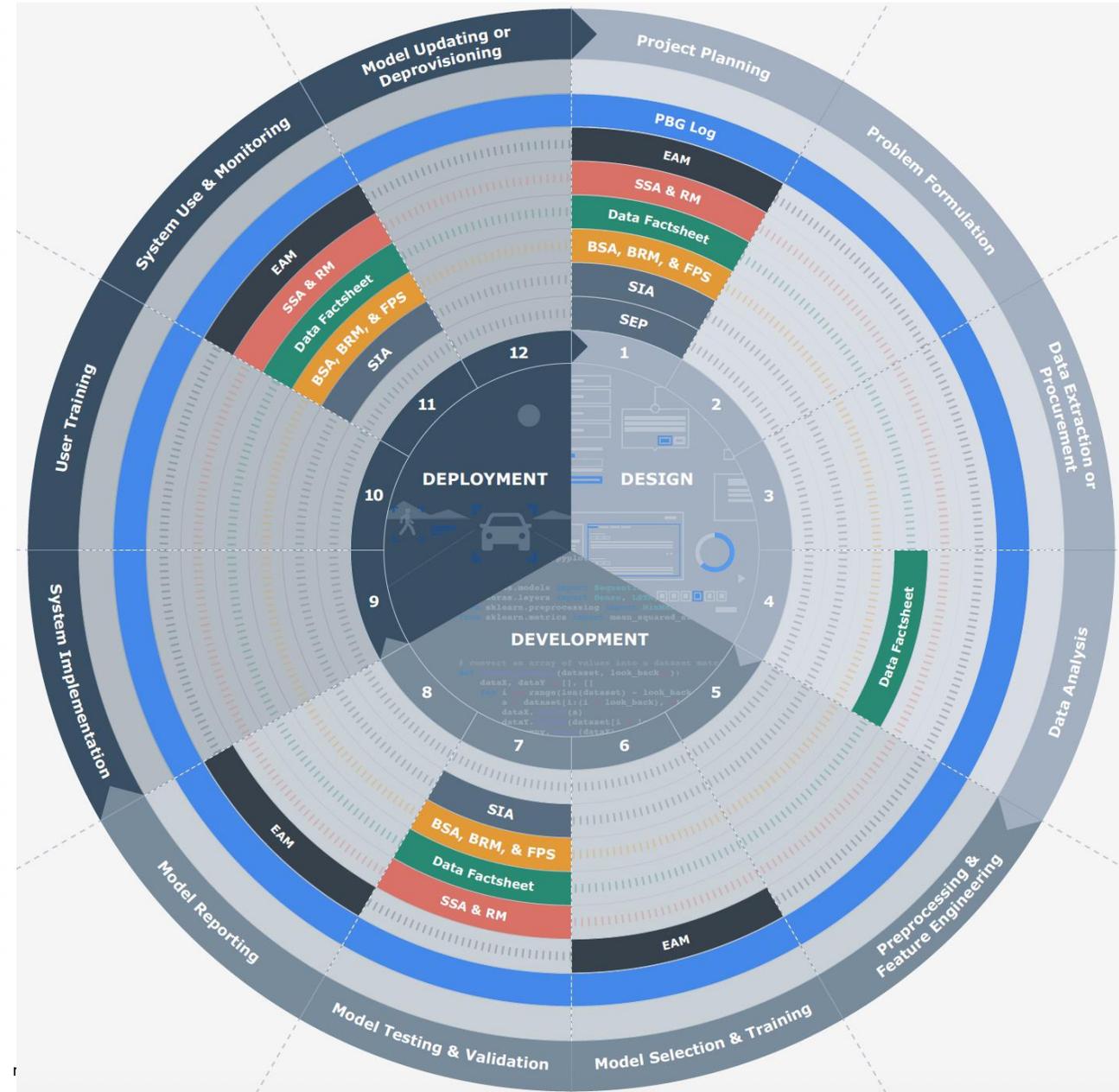
### BSA, BRM & FPS (Bias Self-Assessment, Bias Risk Management & Fairness Position Statement)

Process facilitating the evaluation of how AI projects align with the principle of fairness through the iterative identification and documentation of risks of bias across the lifecycle and assurance actions implemented to address these. The FPS is a document establishing the metric-based fairness criteria for individual AI projects, providing an explanation in plain and nontechnical language.

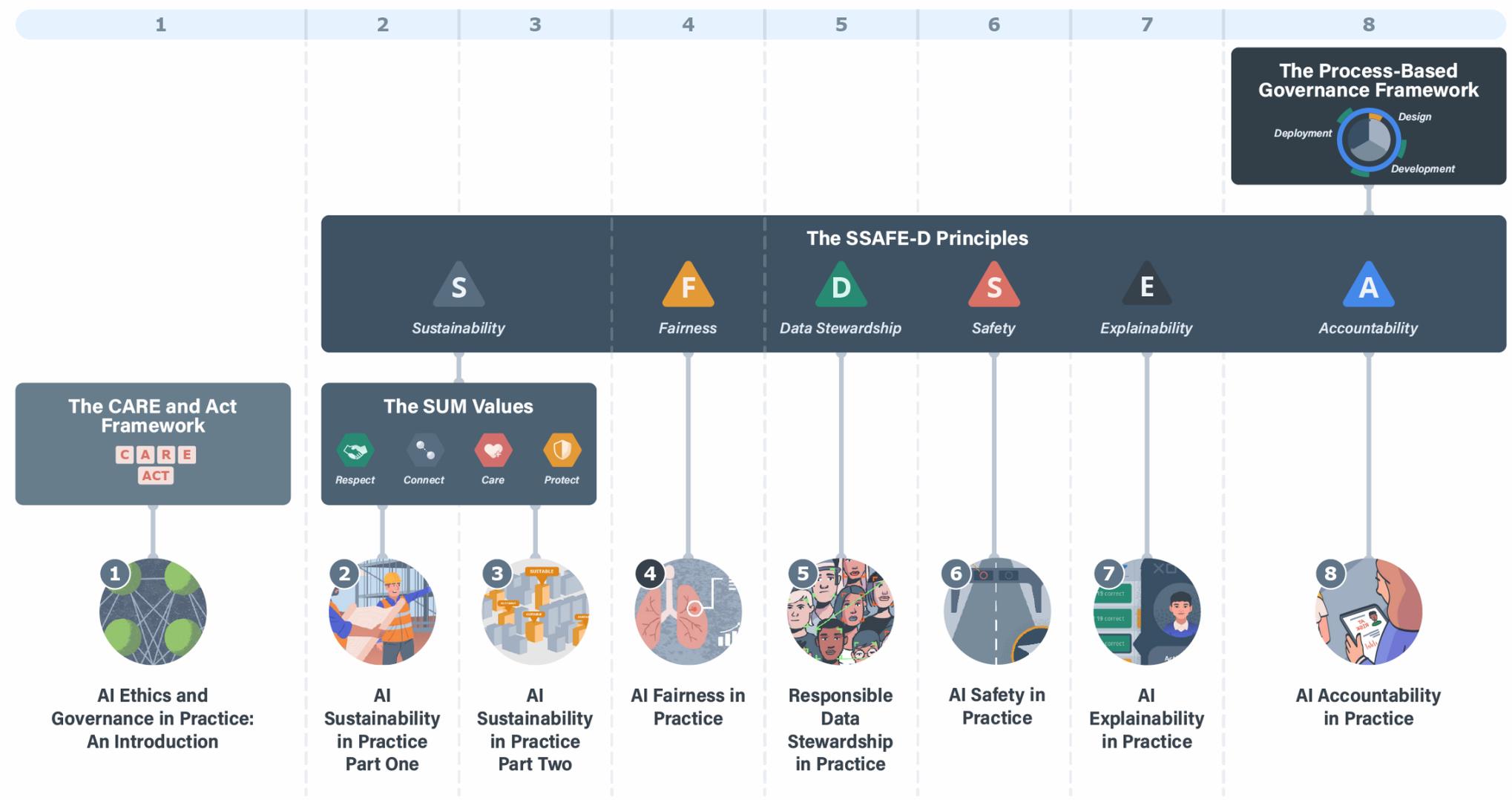
## Data Stewardship

### Data Factsheet

Live document facilitating the uptake of best practices for responsible data management and stewardship across the AI project workflow. The document consists of a comprehensive record of the data lineage and iterative assessments of data integrity, quality, protection, and privacy.



# Programme Roadmap



# Guidance Briefs:



# Scan for more information:



---

– **Project team: David Leslie, Ann Borda, Antonella Maia Perini, Smera Jayadeva**

**Previous contributors: Cami Rincon, Morgan Briggs**

– **Graphics: Conor Rigby**