

Guardrails for Safe and Scalable GenAI Programs

Jean Liu, Executive Director, Liu Pursuits

Doug Graham, Director of Enterprise Data Governance, Mercy

Ren-Yi Lo, Head of Autonomous Systems & Data Governance,
Siemens Healthineers, AI Tech Center

Brian Rasquinha, Associate Director, Solution Architecture,
Privacy Analytics (Moderator)

November 14, 2025

Part I: Background

Who are the Executive Advisory Board?

- 21 senior-level executives from **public and private**-sector organizations
- Geographies: US, Canada, UK, EU
- Industries: healthcare, life sciences, financial services
- Operates as a forum for exchanging insights on **safe and responsible use** of sensitive data to drive innovation
- Identified a need for additional guidance for AI:
 - Evaluating the risk of use-cases
 - Establishing appropriate guardrails for use.

How was the guidance developed?

- Guidance was developed over the course of several quarterly meetings, leveraging the cross-regulation, cross-industry experience of board members
- Overall goal is to drive public and partner trust, while supporting organizations in building efficiency and responding to opportunities for innovation.

Guidance available here:



Work produced by The Board

Thought Leadership

CDO's
Actionable
Framework to
Share
Sensitive Data

Download Now

Live Replay

Data Sharing
Session at MIT
CDOIQ 2022

Watch Now

Thought Leadership

5 Best
Practices for
Driving Value
Through Data
Stewardship

Read Now

Thought Leadership

Guardrails for
Generative AI
as Part of a
Safe &
Scalable AI
Program

Read Now

How was the guidance developed?

- Intended for risks manageable at the **institutional level**; other risks may be more appropriately managed by governmental or intergovernmental organizations and are outside the scope of the guidance.
- Intended to contribute to risk **mitigation**
 - Does not claim to result in minimal achievable risk
 - Does not imply any persistent residual risk is unacceptable
- Not prescriptive or exhaustive and does not constitute legal advice
- Given the advancements in AI, assessment of risk and choice of appropriate guardrails should be continually reviewed rather than a one-and-done effort.

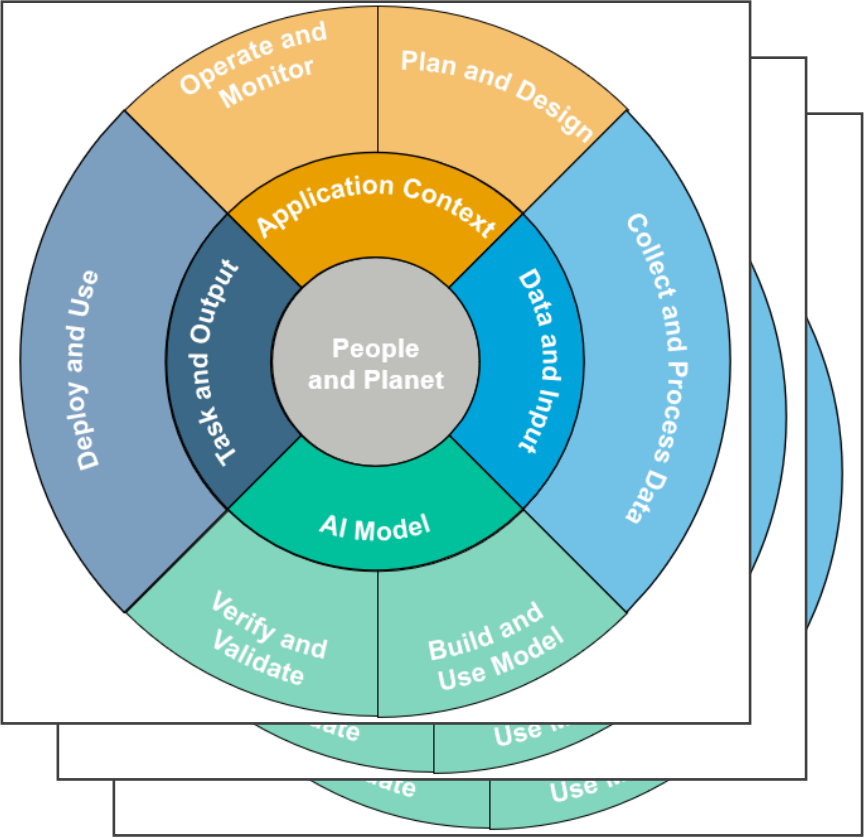
Part II: Guidance

Guidance Overview

1. Risk Level for Use Cases

	Lower Risk			Higher Risk
Context of Users	Under contract	Internal to organization	Internal to partnered organization	External to organization
Context of Impact	Internal		External	Social, economic, professional; e.g., digital therapy
Transparency to users, regulators	Human in the loop, output checking		AI-generated content is tagged	Use of AI is explicitly stated
Purposes	Reference (e.g., summaries, categorization)		Research	Decision support
Organizational Risk	Operational impacts		Business impacts	Legal, regulatory impacts

1. Guardrails mapped to NIST AI Risk Management Framework (3 Dimensions)



1. Evaluating Risk Level for Use Cases across Multiple Dimensions

- Dimensions affecting the **organizational risk** associated with GenAI use case
- Features/factors are presented from “Lower Risk” to “Higher Risk”
- **All dimensions** are intended to be taken into consideration
 - Overall risk level is intended to reflect a consolidated posture

Component

1

1. Evaluating Risk Level for Use Cases across Multiple Dimensions



Context of Users	Under contract	Internal to organization	Internal to partnered organization	External to organization
Context of Impact	Internal		External	Social, economic, professional; e.g., digital therapy
Transparency to users, regulators	Human in the loop, output checking	AI-generated content is tagged	Use of AI is explicitly stated	
Purposes	Reference (e.g., summaries, categorization)	Research	Decision support	Decision-making
Organizational Risk	Operational impacts		Business impacts	Legal, regulatory impacts

Example: Evaluating Risk Level

- **Mercy's GEN-AI SBAR Handoff Process**
- Mercy Health has implemented a **Generative AI-powered SBAR (Situation, Background, Assessment, Recommendation)** workflow for ED-to-IP transitions:
 - **Trigger Point:** When an ED admission order is placed and an inpatient bed is assigned.
 - **Automation:** GEN-AI compiles ED provider notes, clinician docs, and discrete patient data.
 - **Structured Output:** Creates a concise SBAR summary.
 - **Delivery:** Sent via Epic secure chat to the inpatient nursing team; notifications in Rover & Epic.
 - **Continuous Updates:** The AI-generated note updates until patient transfer is complete.
 - **Benefits:**
 - › Reduces cognitive burden on ED nurses.
 - › Improves accuracy and minimizes omissions.
 - › Enhances throughput by reducing phone calls and manual steps.
 - **Provider Role:** Timely ED documentation is critical for SBAR completeness.



Dimensions

User Context:

– Internal (L)

Impact:

– Internal (L)

Transparency:

– HITL (L)

Purpose:

– Decision Support (M)

Org Risk:

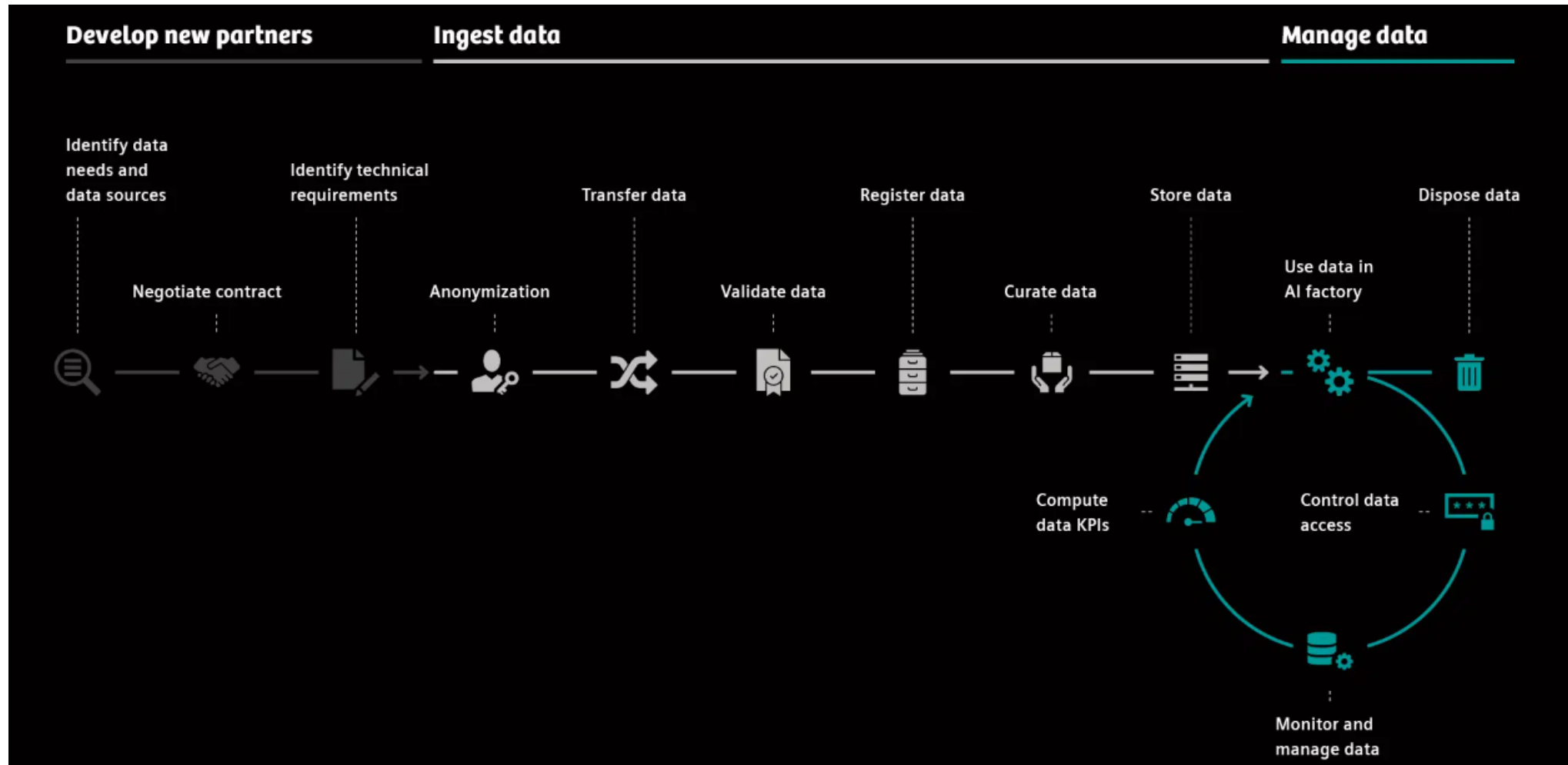
– Business Impact (M)

Lower Risk

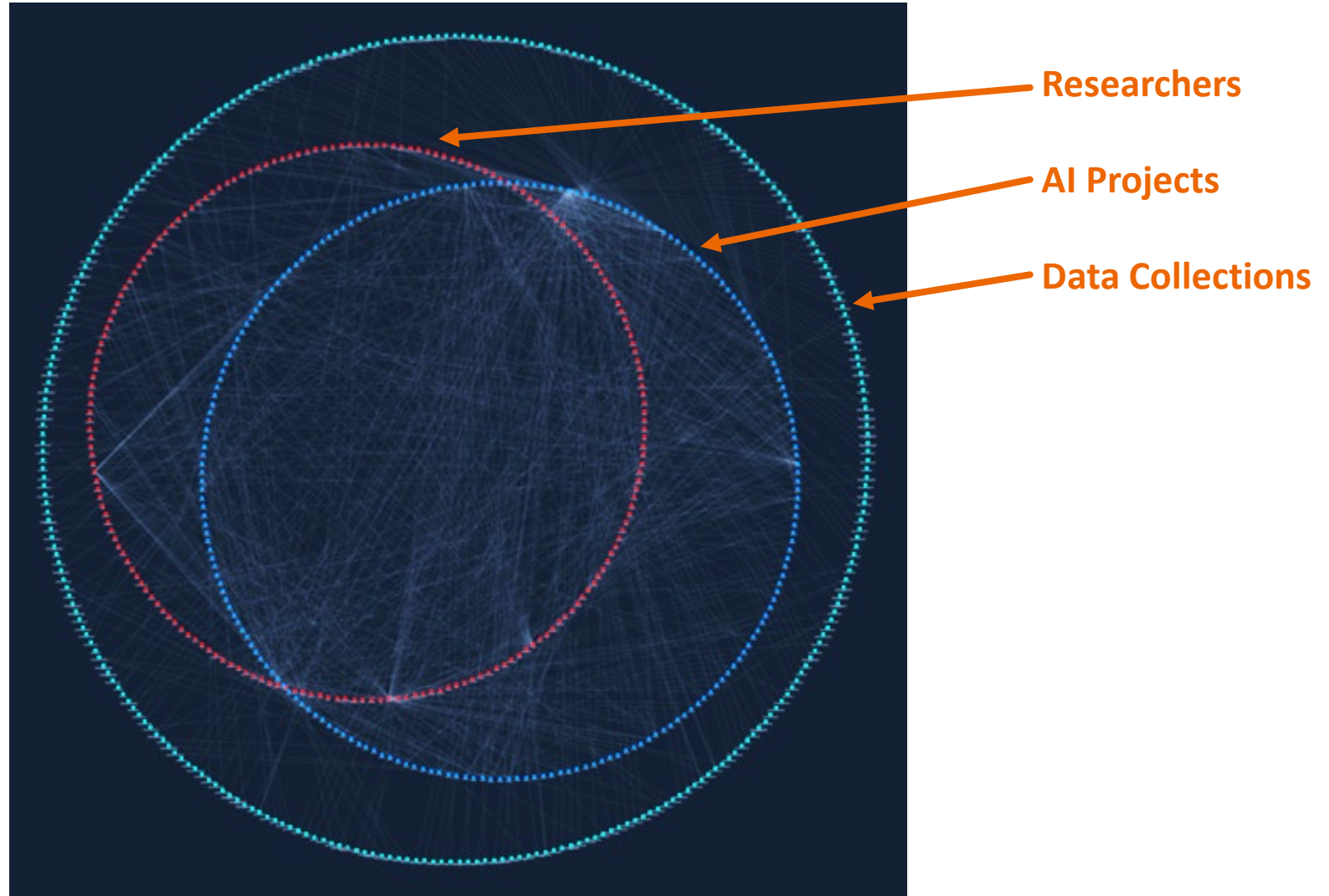


Higher Risk

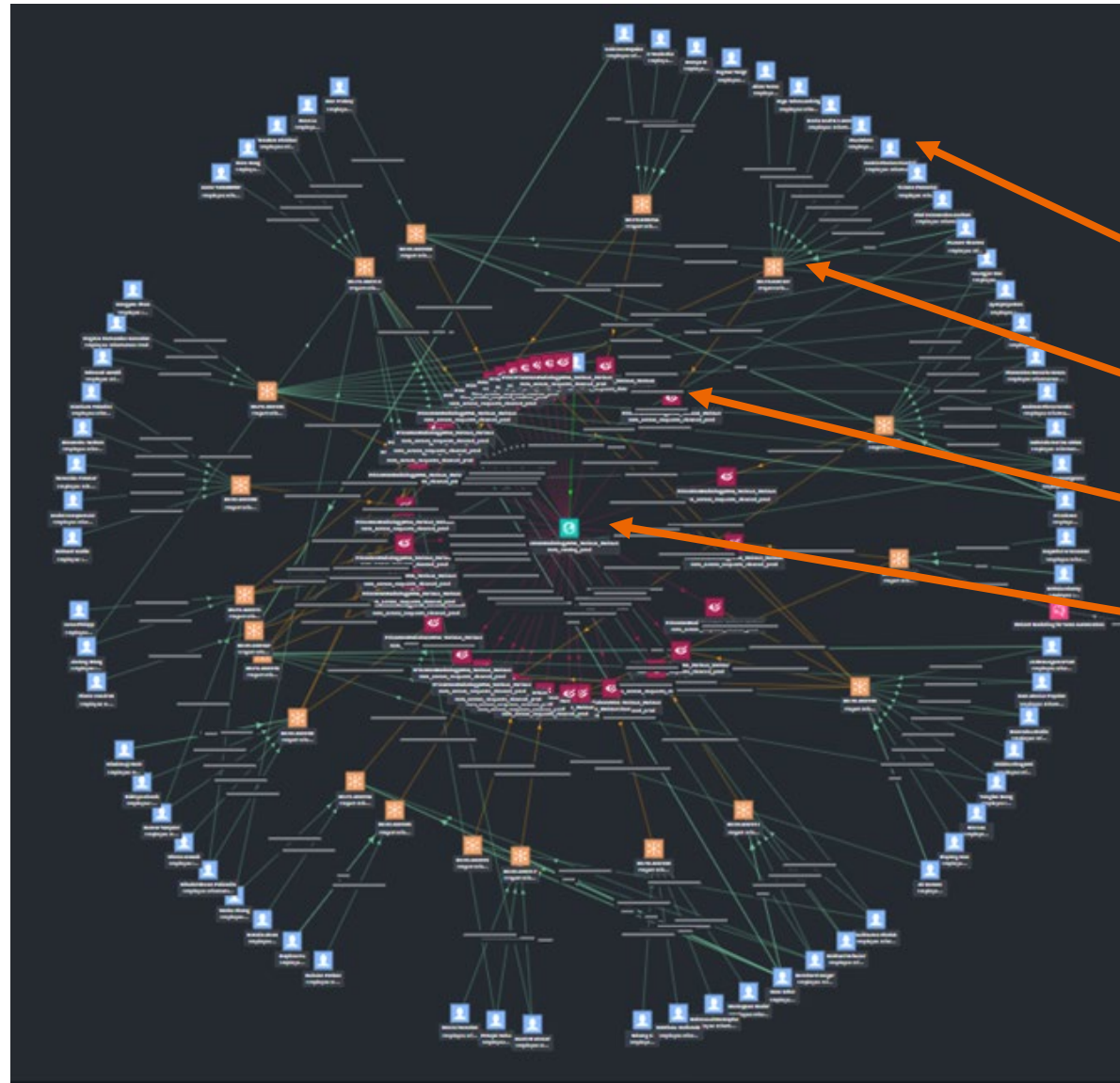
Comprehensive Data Lifecycle Management for Trustworthy AI Development



Data Collections, AI Projects, Researchers



How Data Access and AI Projects are Managed



Researcher

AI Project

Data Access

Data Collection

Automating data governance extraction for data management



Problem Statement:

Can our Data Management Control Center automate extraction of data governance from data contracts by utilizing GenAI solution?

Business Case:

Efficiency Gains: from 4h to 10min



Risk Evaluation:

- Results are consumed internally to manage AI and Data governance
- Human-in-the-loop towards Human-on-the-loop
- Additional guardrails: sufficient high-quality historical data, domain experts, audits

2. Guardrails Mapped to NIST AI Risk Management Framework

- Guardrails mapped to the [NIST AI Risk Management Framework](#) (AI RMF)
- The guardrails are further broken down into three categories of approaches:
 - Safety & Adaptability
 - Governance & Transparency
 - Trust & Ethical Responsibility
- Higher-risk use cases -> usually stronger, more numerous, and/or broader guardrails

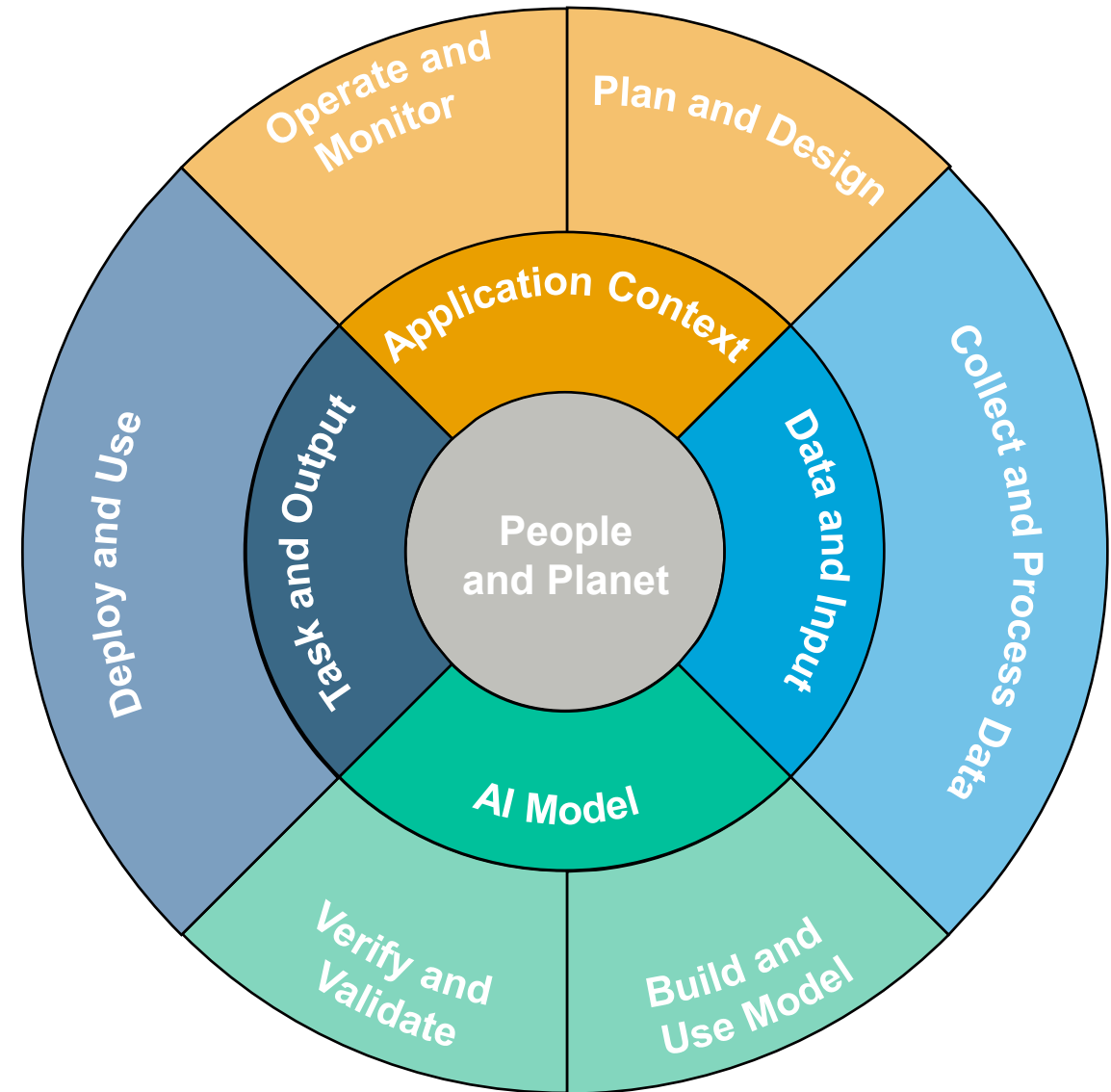
Component

2

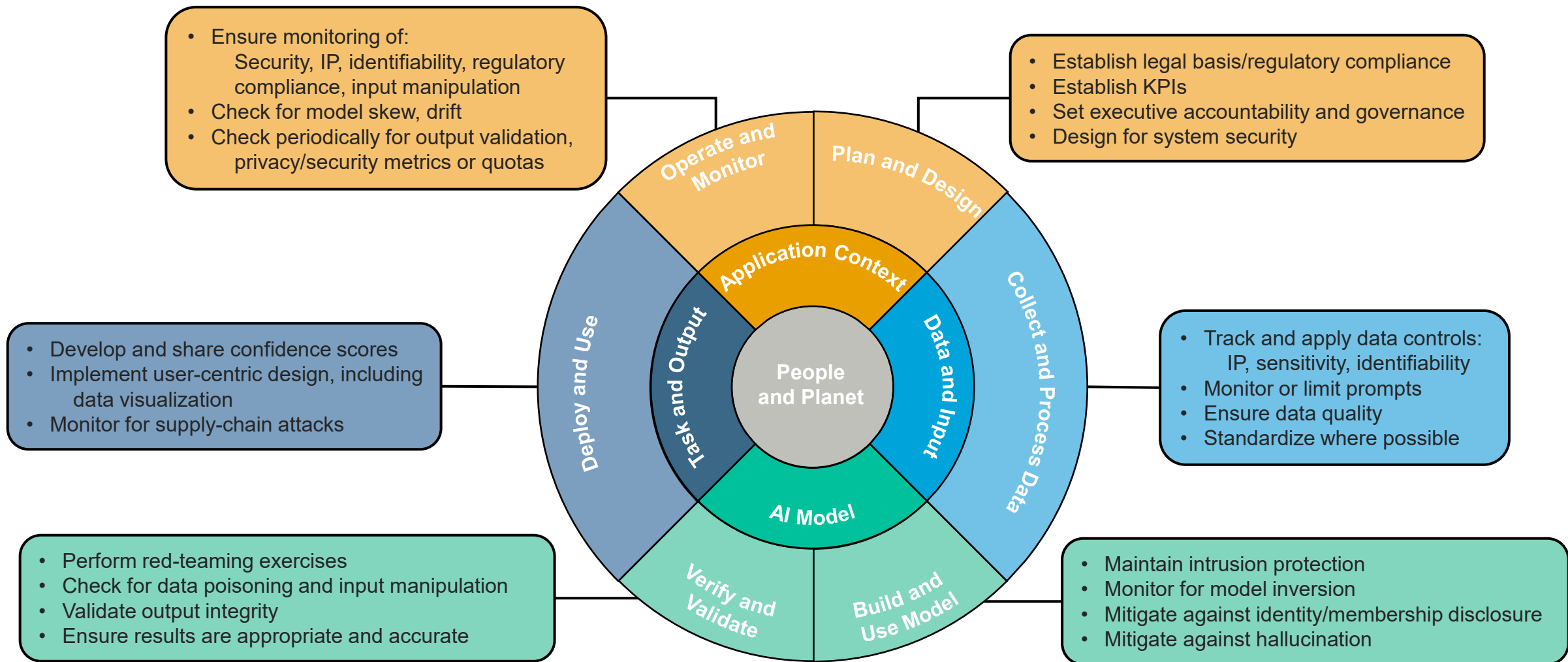
2. Guardrails Mapped to NIST AI Risk Management Framework

- The NIST AI RMF, developed in a public-private collaboration, aims to manage risks to individuals, organizations, and society that are associated with AI.
- It is organized by **Key Dimensions**, listed in the center and inner ring, and AI **Lifecycle Stages**, in the outer ring, corresponding to the use of AI tools.
- The EAB recommends guardrails for consideration to mitigate against some risks, mapped to the lower-level Lifecycle Stages of the outer ring.

For clarity, we emphasize that “Build and Use Model” in the AI RMF includes the activities of creating or selecting algorithms; training models; and model testing.



Guardrails: Safety & Adaptability



Safety & Adaptability

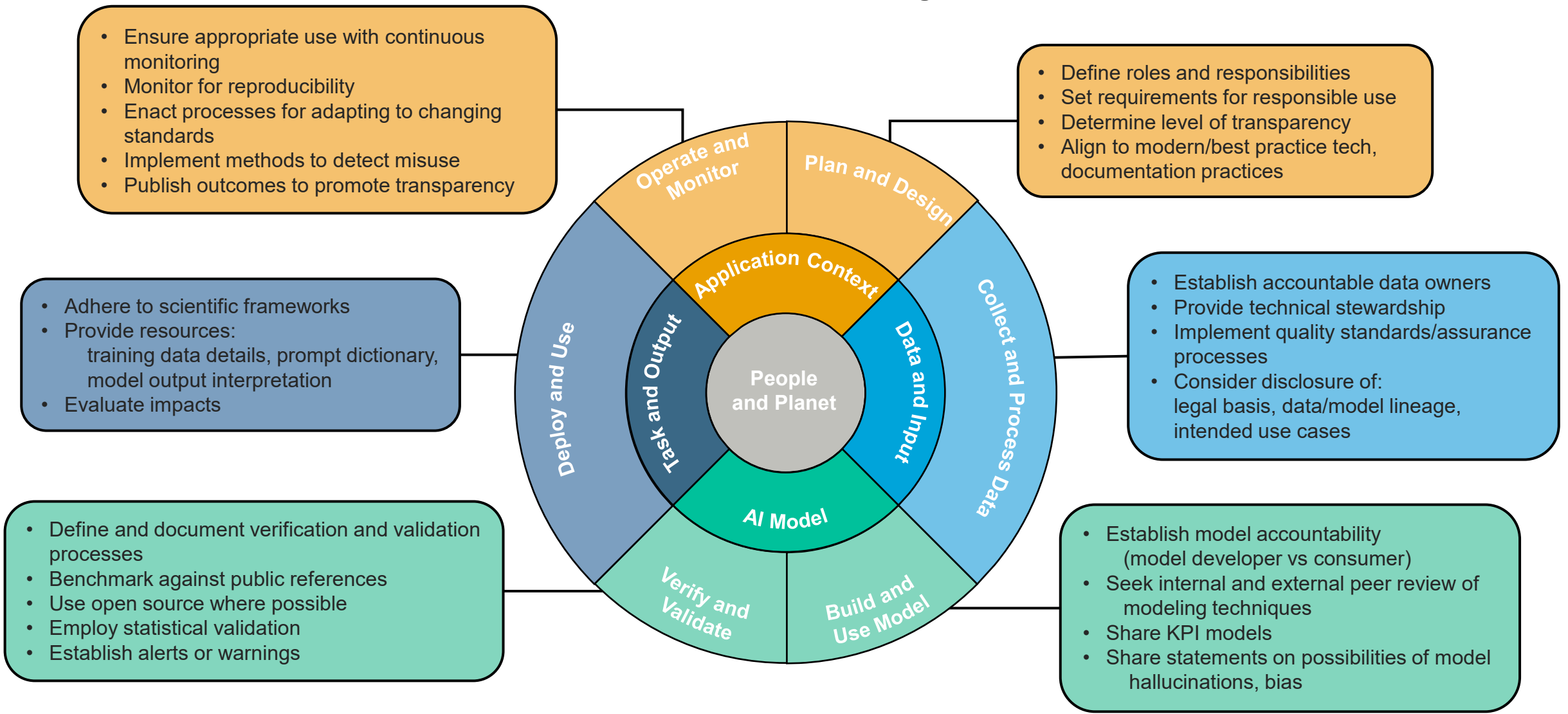
encompasses privacy and security; content review and fact checking; and adaptation and innovation

Example: Safety & Adaptability

An AI driven legal guidance in SharePoint

- Setting up guardrails for how outputs may be used
- 360 input from stakeholders & iteration
- Accountability for the results

Guardrails: Governance & Transparency



Governance & Transparency

encompasses oversight and governance; compliance with academic and scientific standards; transparency and disclosure

Example: Governance & Transpacency

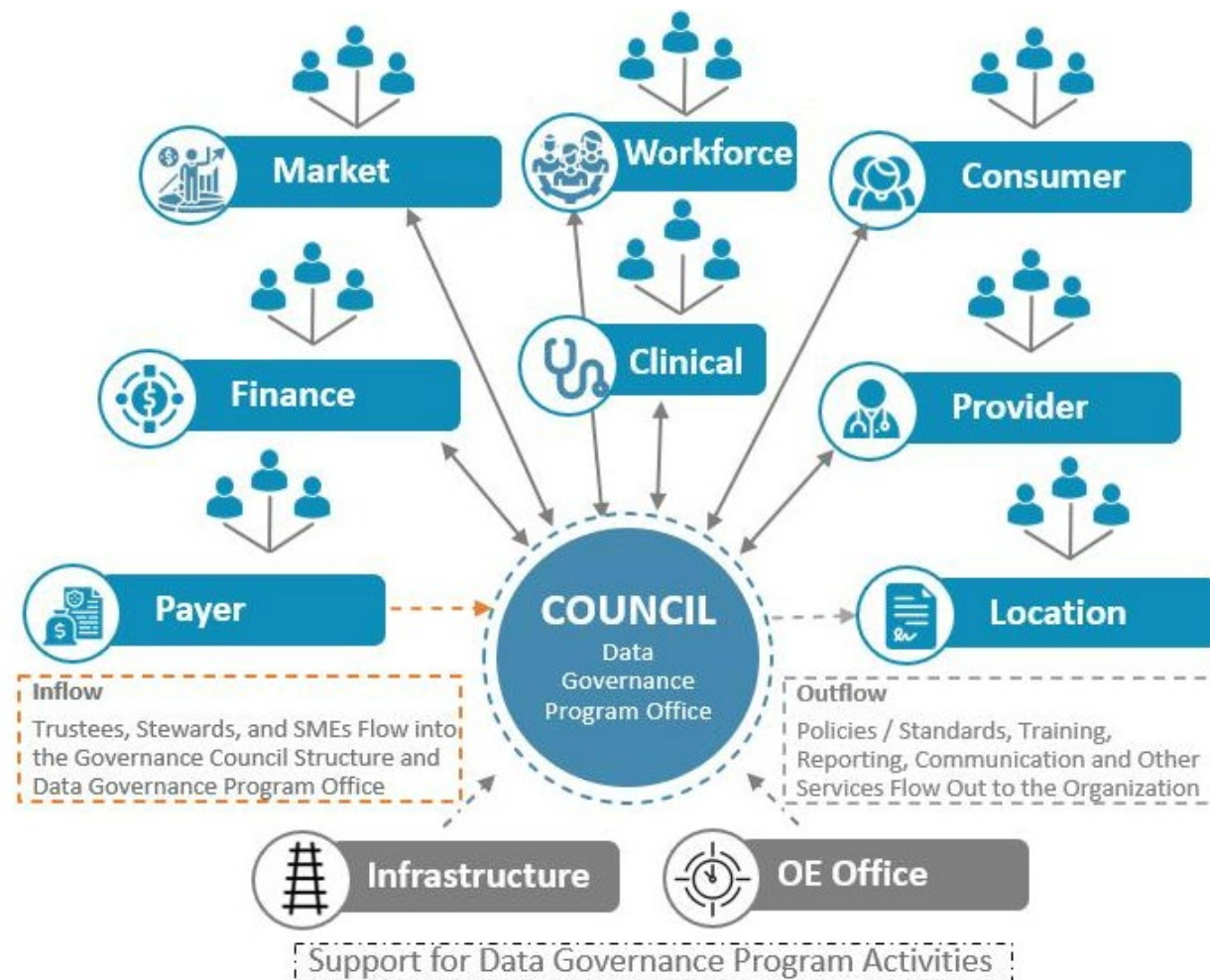
An AI driven legal guidance in SharePoint

- Periodic review of the user, validity, and value
- Providing the governance process to all users
- Ad hoc audits

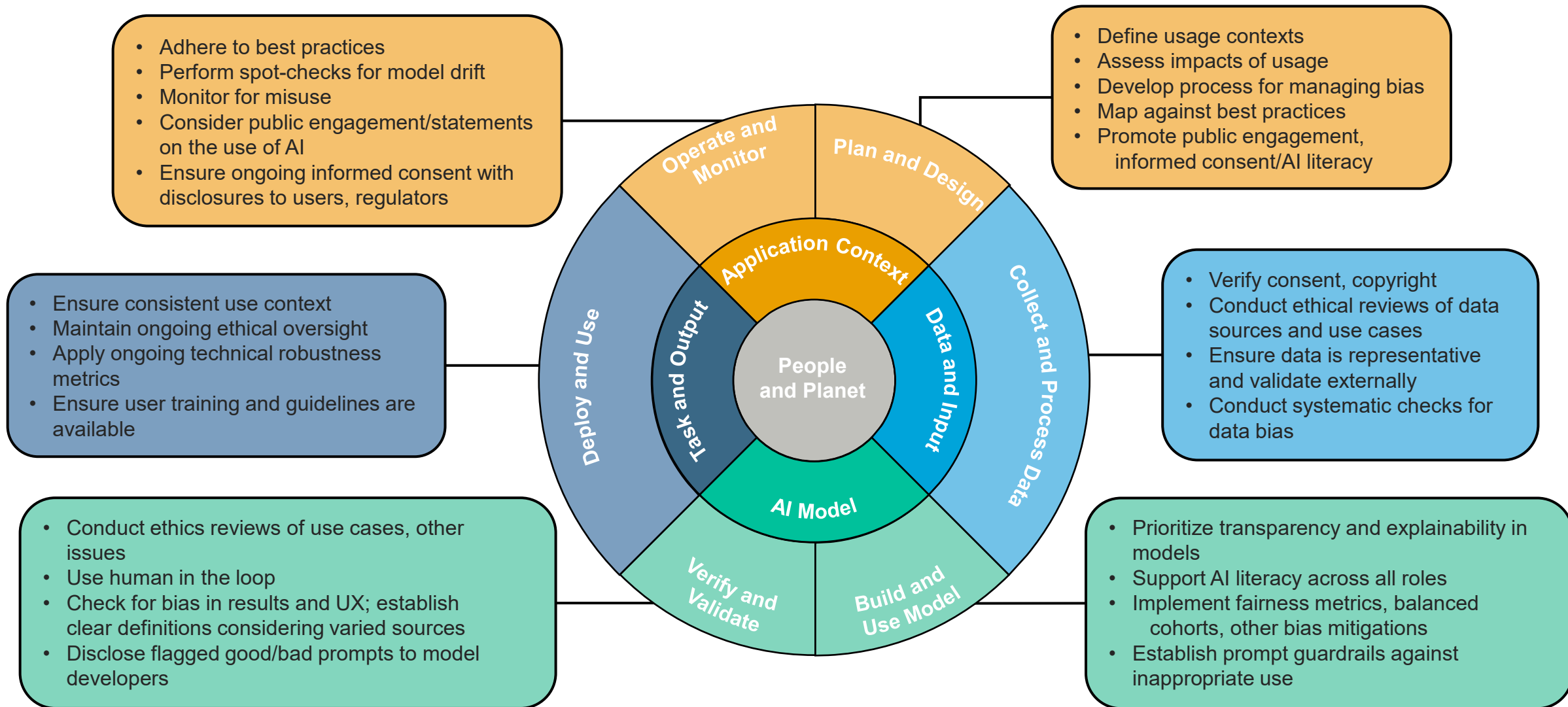
Example: Governance & Transparency



- Data Governance Framework
 - Business/Ministry Based Health Care Subject Area Domain Accountability
 - Implement new AI Domain Stewardship role
 - Align with Mercy Mission and Ethics Leadership
 - Integrate Council of Responsible AI



Guardrails: Trust & Ethical Responsibility



Trust & Ethical Responsibility

encompasses ethical and responsible uses; bias and fairness; stakeholder engagement and inclusivity

Example: Governance & Transpacency

An AI driven legal guidance in SharePoint

- Periodic review of the user, validity, and value
- Providing the governance process to all users
- Ad hoc audits

Example: Trust & Ethical Responsibility



- **Core Principles from Mercy's AI Governance Charter**

- **Transparency**

- Maintain openness about AI development, deployment, and usage.
- Ensure stakeholders understand how AI decisions are made.

- **Accountability**

- Establish clear lines of responsibility for AI decisions and outcomes.
- Humans remain in control and accountable for AI systems.

- **Ethics and Compliance**

- Align AI use with organizational values and ethical standards.
- Comply with all relevant laws and regulations, including privacy and security requirements.

- **Fairness and Non-Discrimination**

- Prevent bias and discriminatory outcomes in AI applications.
- Promote equitable treatment across all user groups.

- **Privacy and Security**

- Protect individual rights and secure all data used in AI initiatives.
- Implement safeguards against unauthorized access and misuse.

- **Human-Centric and Socially Beneficial**

- Keep AI applications focused on improving patient care and societal well-being.
- Ensure technology serves people, not replaces human judgment.

- **Continuous Monitoring and Improvement**

- Regularly audit AI models for performance, fairness, and compliance.
- Update practices as technology and regulations evolve.

- **Collaboration and Education**

- Foster cross-functional teamwork and knowledge sharing.
- Invest in training to build AI literacy across the organization

Thank you!



Jean Liu
Executive Director,
Liu Pursuits



Doug Graham
Director of Enterprise
Data Governance,
Mercy



Ren-Yi Lo
Head of Autonomous
Sys & Data Gov,
Siemens Healthineers,
AI Tech Center



Brian Rasquinha
Associate Director,
Solution Architecture,
Privacy Analytics