# *Aerospace's Trusted AI Framework*

*Dr. Phil Slingerland, ESD/RSSAS/MIXD*
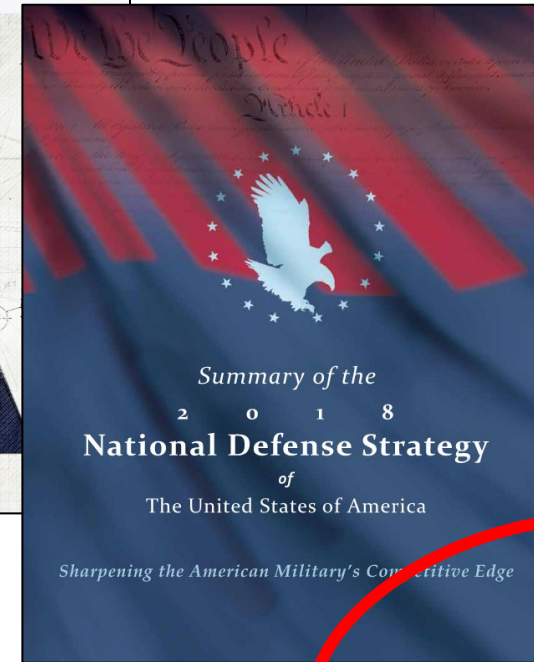*Lauren Perry, ATD/SRO/SAG*

*GSAW 2023*

# Aerospace's Trusted AI Framework

*Motivation*

- AI has proven success in wide array of applications, leading to:
  - *Increasing role of automated decision making*
  - *Pervasive use of AI-enabled systems*
- Consequence of increasing presence of AI:
  - *Risk of impact to all stakeholders (business leaders, public, governments)*
  - *Increased awareness of responsibility (financial, legal, ethical)*
- Regulations
  - *Global wave of guidelines and regulation targeting AI*
  - *DoD and IC published recommended policies*
  - *Parallel developments in industry-specific regulation (financial, pharmaceutical, autonomous vehicles, etc.)*

**Final Report**

National Security Commission on Artificial Intelligence

Summary of the

2 0 1 8

**National Defense Strategy**

of

The United States of America

*Sharpening the American Military's Competitive Edge*

78939

**Presidential Documents**

Executive Order 13960 of December 3, 2020

Promoting the Use of Trustworthy Artificial Intelligence in the Federal Government

By the authority vested in me as President by the Constitution and the laws of the United States of America, it is hereby ordered as follows:

Section 1. *Purpose.* Artificial intelligence (AI) promises to drive the growth of the United States economy and improve the quality of life of all Americans. In alignment with Executive Order 13859 of February 11, 2019 (Maintaining American Leadership in Artificial Intelligence), executive departments and agencies (agencies) have recognized the power of AI to improve their operations, processes, and procedures; meet strategic goals; reduce costs; enhance oversight of the use of taxpayer funds; increase efficiency and mission effectiveness; improve quality of services; improve safety; train workforces; and support decision making by the Federal workforce, among other positive developments. Given the broad applicability of AI, nearly every agency and those served by those agencies can benefit from the appropriate use of AI.

Agencies are already leading the way in the use of AI by applying it to accelerate regulatory reform; review Federal solicitations for regulatory compliance; combat fraud, waste, and abuse committed against taxpayers; identify information security threats and assess trends in related illicit activities; enhance the security and interoperability of Federal Government information systems; facilitate review of large datasets; streamline processes for grant applications; model weather patterns; facilitate predictive maintenance; and much more.

Agencies are encouraged to continue to use AI, when appropriate, to benefit the American people. The ongoing adoption and acceptance of AI will depend significantly on public trust. Agencies must therefore design, develop, acquire, and use AI in a manner that fosters public trust and confidence while protecting privacy, civil rights, civil liberties, and American values, consistent with applicable law and the goals of Executive Order 13859.

Certain agencies have already adopted guidelines and principles for the use of AI for national security or defense purposes, such as the Department of Defense's *Ethical Principles for Artificial Intelligence* (February 24, 2020), and the Office of the Director of National Intelligence's *Principles of Artificial Intelligence Ethics for the Intelligence Community* (July 23, 2020) and its *Artificial Intelligence Ethics Framework for the Intelligence Community* (July 23, 2020). Such guidelines and principles ensure that the use of AI in those contexts will benefit the American people and be worthy of their trust.
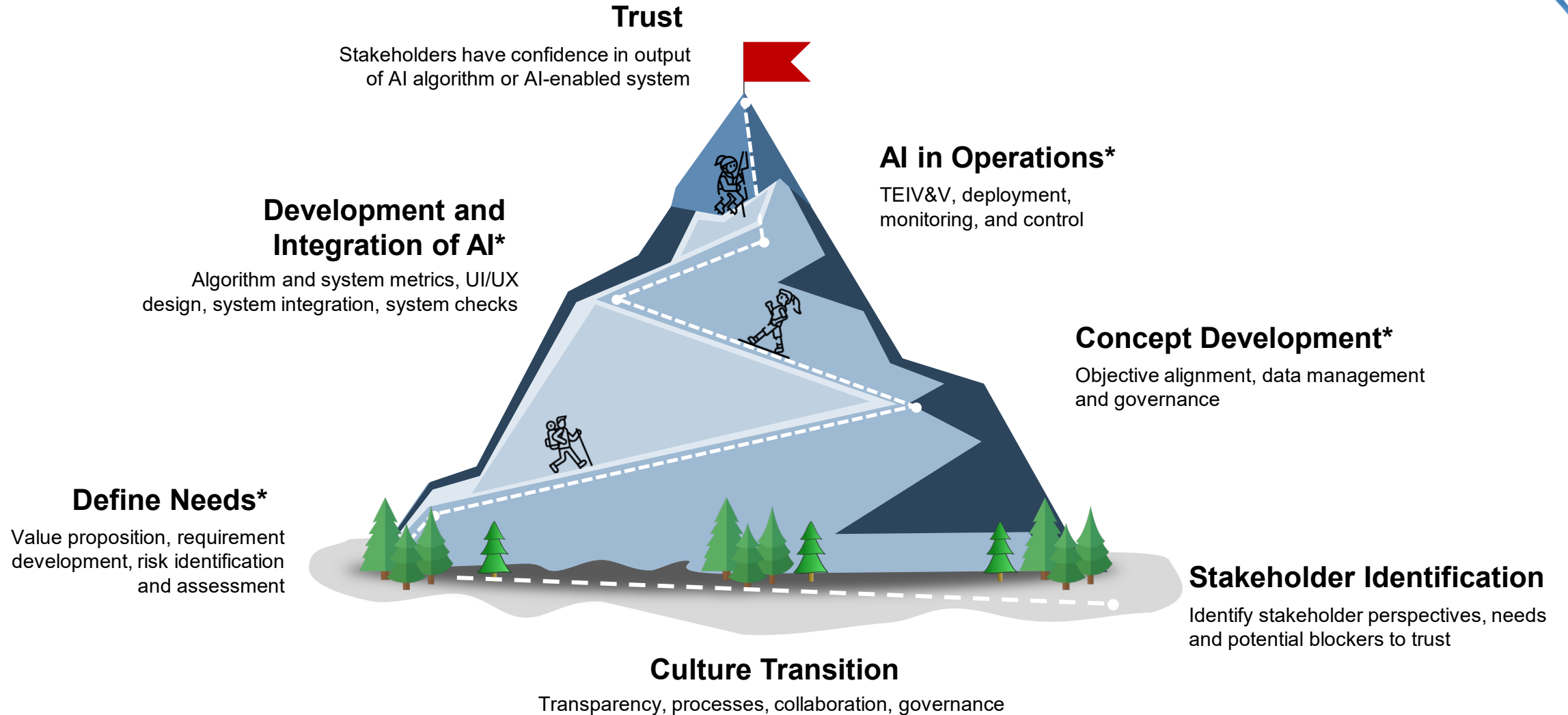
Section 3 of this order establishes additional principles (Principles) for the use of AI in the Federal Government for purposes other than national security and defense, to similarly ensure that such uses are consistent with our ... establishes ... guidance

Agencies are encouraged to continue to use AI, when appropriate, to benefit the American people. The ongoing adoption and acceptance of AI will depend significantly on public trust. Agencies must therefore design, develop, acquire, and use AI in a manner that fosters public trust and confidence while protecting privacy, civil rights, civil liberties, and American values, consistent with applicable law and the goals of Executive Order 13859.

***Aerospace Developed A Trusted AI Framework to assist customers in design, implementation, and assessment of AI-based algorithms, with an emphasis on high consequence environments***
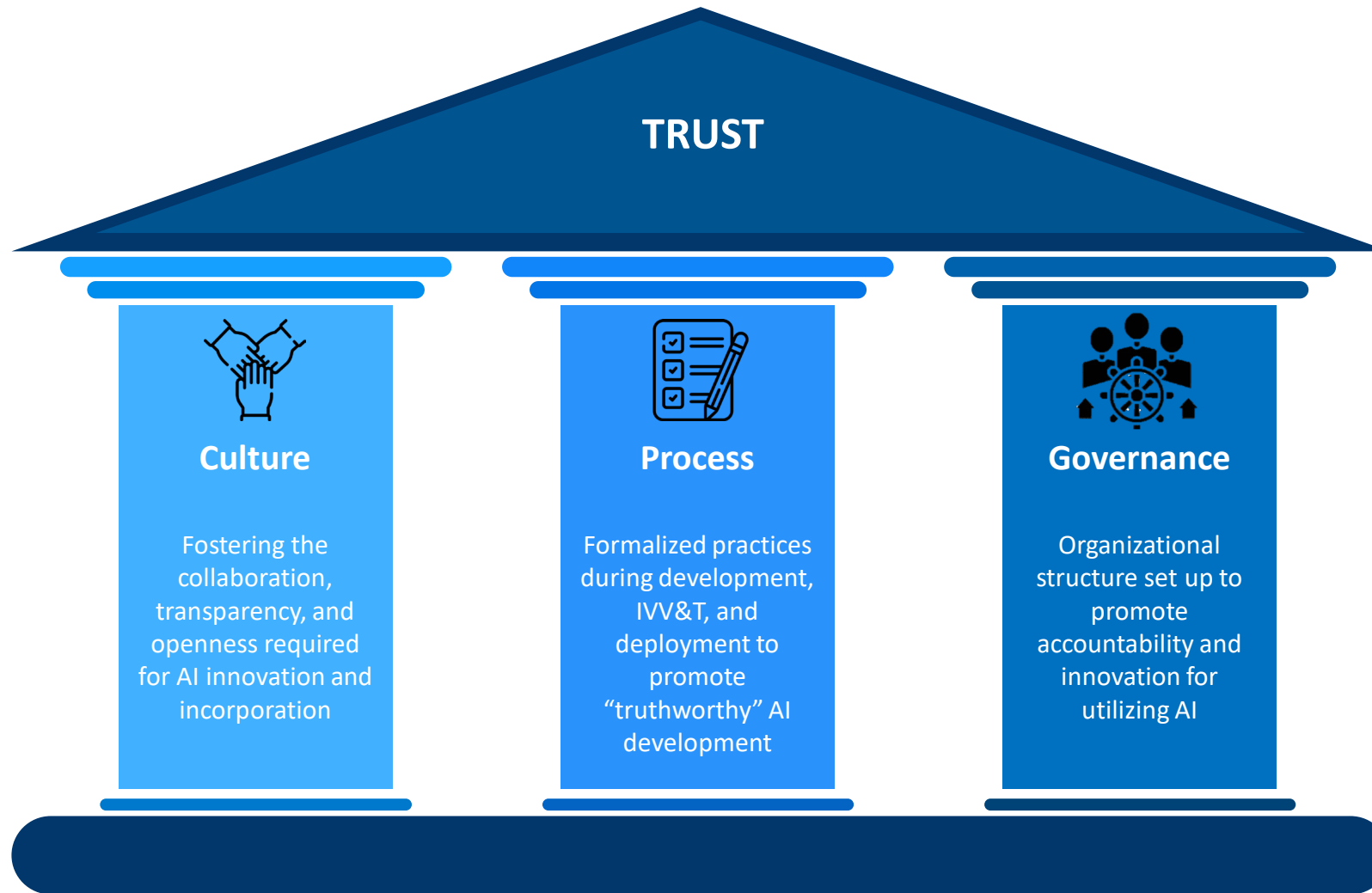
# Journey of Trust



**Trust**

Stakeholders have confidence in output of AI algorithm or AI-enabled system

**AI in Operations***

TEIV&V, deployment, monitoring, and control

**Development and Integration of AI***

Algorithm and system metrics, UI/UX design, system integration, system checks

**Concept Development***

Objective alignment, data management and governance

**Define Needs***

Value proposition, requirement development, risk identification and assessment

**Stakeholder Identification**

Identify stakeholder perspectives, needs and potential blockers to trust

**Culture Transition**

Transparency, processes, collaboration, governance

*Supported by Trusted AI Framework*

*Definition and journey to trusting an AI-enabled system is in the eye of the beholder*

# Pillars of Trust



**TRUST**

**Culture**

Fostering the collaboration, transparency, and openness required for AI innovation and incorporation

**Process**

Formalized practices during development, IVV&T, and deployment to promote "truthworthy" AI development

**Governance**

Organizational structure set up to promote accountability and innovation for utilizing AI

*Trust requires a facilitating environment and collaborative team*

# Trust Is More Than Technology
## *The People Side of Trust*

## Values of Trust

**Encourage openness and transparency**

- Faithfully capture risks inherent to AI capability
- Anticipate challenges of deployment with additional data collection, training, and testing
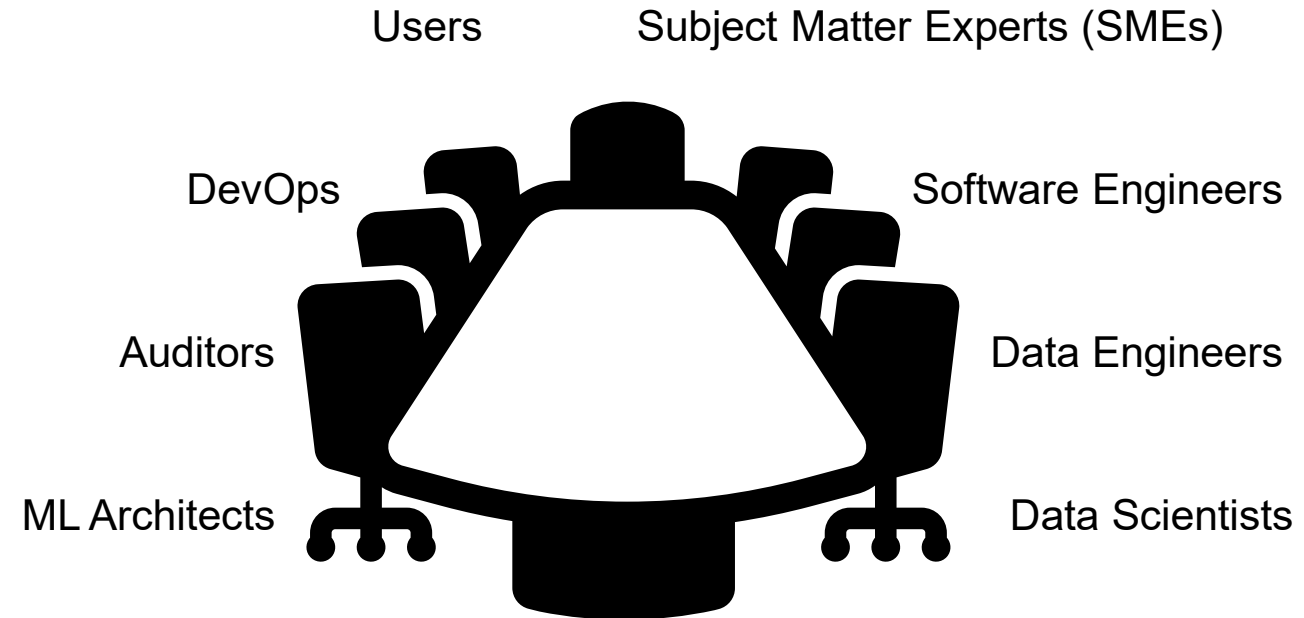- Avoid imprecise language of AI hype

**Develop collaboratively with users**

- Frame development as journey in building user trust
- Strive for informed users that have input to AI design, function, interpretability, and control

**Set high expectations for traceability**

- Set high expectations for traceability
- Software, data, and model version control
- Record design decisions and R&D progress
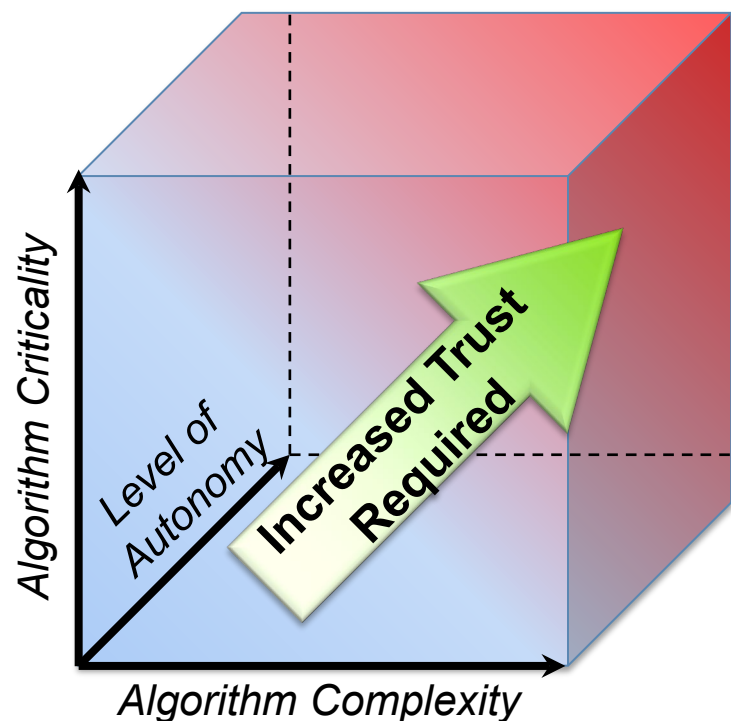- Log performance metrics throughout development

## Stakeholders



Users    Subject Matter Experts (SMEs)

DevOps    Software Engineers

Auditors    Data Engineers

ML Architects    Data Scientists

**Without consideration of various stakeholders, trust of an AI-enabled system will never be fully realized**

# Degree of Trust

*How Much Trust Investment is Needed?*



- The amount of trust required is related to risk to mission integrity defined by three dimensions:
  - **Algorithm Criticality**: impact on mission success
  - **Algorithm Complexity**: degree of interpretability
  - **Level of Autonomy**: independence from human intervention

- Assess degree of trust required for application
  - *Engage stakeholders on potential impacts of deployment*
  - *Define benchmarks for success*
  - *Estimate LOE and budget and weigh against expected value*

- Operational risks can be mitigated through DevOps best practices
  - *Prepare for unavailable, drifting, or poorly performing models*
  - *Provide awareness of data, model, and objective alignment throughout model lifecycle*
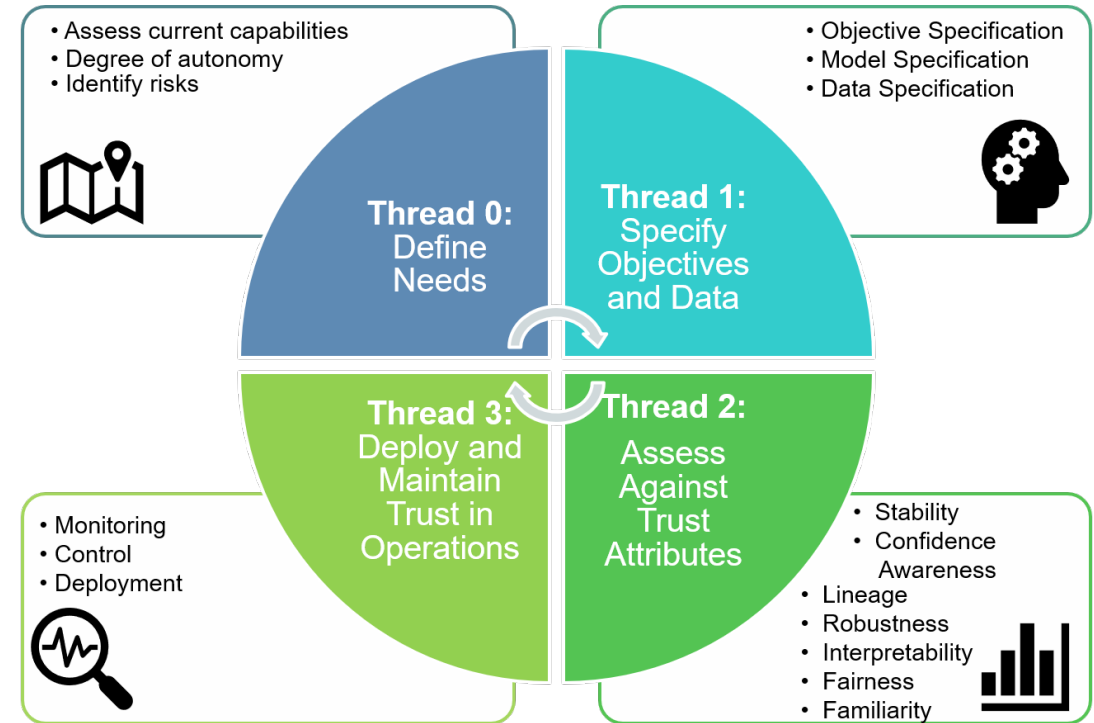
**Trust is not free and should be tailored to intended use**

# How to Trust AI?

**Trusted AI**

*AI capability that provides sufficient confidence of satisfying user-defined objectives in a proper, interpretable, and safe way over its lifetime*
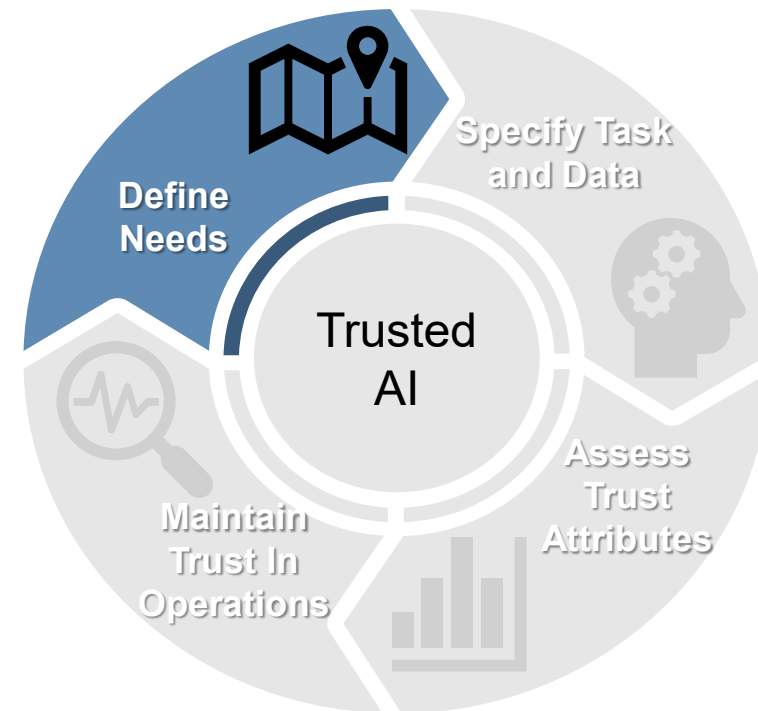


- Assess current capabilities
- Degree of autonomy
- Identify risks

- Objective Specification
- Model Specification
- Data Specification

**Thread 0:** Define Needs

**Thread 1:** Specify Objectives and Data

**Thread 3:** Deploy and Maintain Trust in Operations

**Thread 2:** Assess Against Trust Attributes

- Monitoring
- Control
- Deployment

- Stability
- Confidence Awareness
- Lineage
- Robustness
- Interpretability
- Fairness
- Familiarity

# Need Definition, Value Proposition and Intent

*Thread 0: Definition of Needs*

- Define and justify need for AI-based capability
  - *Compare with simplest approach or existing capability*
  - *Conduct literature review to support AI appropriateness*

- Identify desired level of autonomy of AI-based capability
  - *Degree of human intervention desired in operations*

- Reduce and manage risk
  - *Consider appropriate model complexity*
  - *Isolate operation to single function*
  - *Leverage existing systems and codebases*



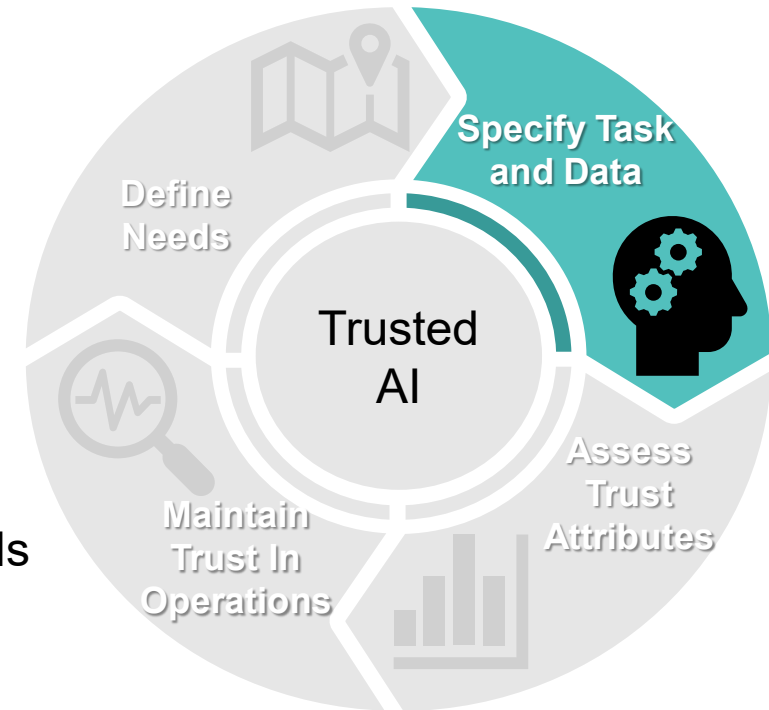*Developing business case and rationale for bringing AI into a system*

# Task Specification

*Thread 1: Specify Task and Data*

- Translate proposed capability to AI task
  - *Modular description:*
    - Compose task as collection of simple functions
    - Develop requirements for each to assemble trust
  - *Identify performance metrics relevant to target domain*

- Identify potential failure modes and how AI could cause them
  - *Early expectation of potential failures to help mitigate them*
  - *Note hardware and software limitations of deployed system*

- Capture performance and trust requirements; define metrics/TPMs
  - *Tailor metrics to use case*
  - *Algorithm, system, and mission focused metrics*
  - *Include upstream and downstream metrics (for operational monitoring)*

**Objective Specification**
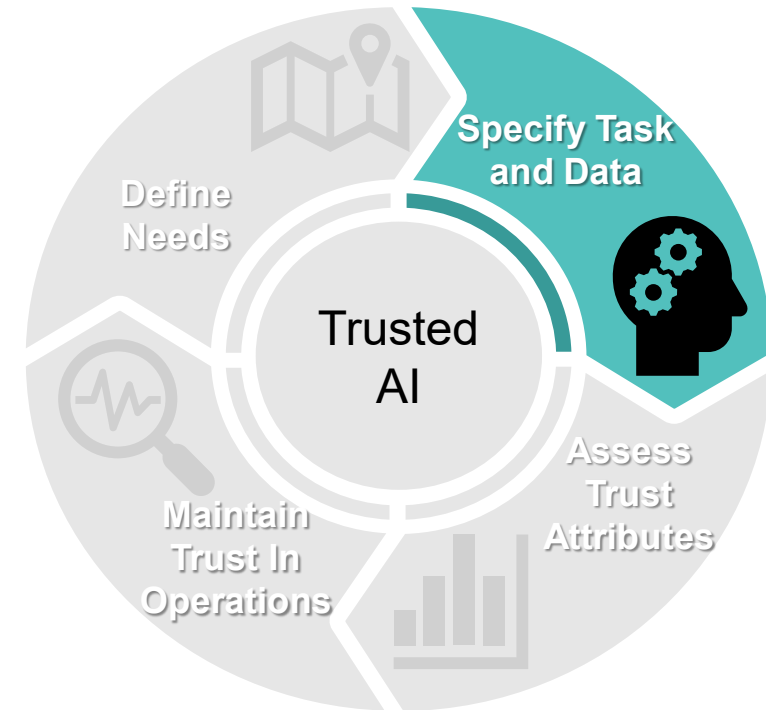*Clear description of AI objective with metrics that demonstrate task alignment*

Define Needs

Specify Task and Data

Trusted AI

Assess Trust Attributes

Maintain Trust In Operations

# *Data Specification*

*Thread 1*

- Identify datasets for each development stage and V&V
  - *Specify how data will be collected and partitioned*
  - *Capture details of sensors, data pre-processing (Traceability)*
- Perform exhaustive exploratory data analysis (EDA):
  - *Identify:*
    - Nominal and out-of-scope data parameters
    - Subgroups and their relative representation (Fairness)
    - Challenging data examples and mitigation plans
    - High spatial or temporal correlations, especially across data splits
- Develop upstream protection to check for out-of-scope data
  - *Re-route out-of-scope data to alternate system*
  - *Tailor confidence based on scope (Pertinence)*
- Provide tools that enhance data pedigree:
  - *Interface to gathering annotations from multiple SMEs, logging annotations from users*
  - *Enable review and disagreement between SMEs*
- Collect and assess representative target data near deployment phase (Pre-deployment)
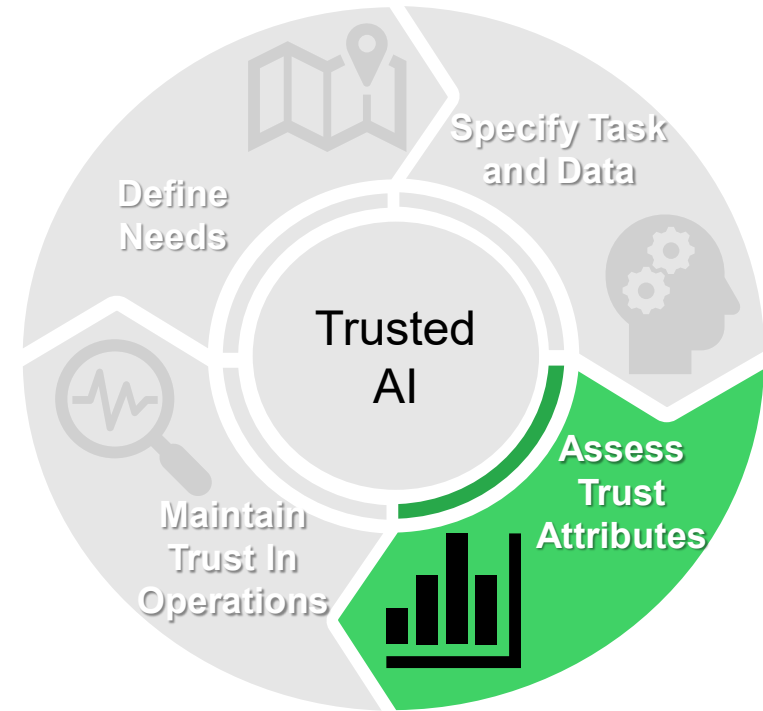
**Data Specification**

*Plan for collection, characterization, annotation, and management of data through AI lifecycle*

# Assess and Enhance Trust

*Thread 2*

- Evaluate proposed design using attributes of trust:
  - *Traceability*: document and maintain artifacts from implementation and evaluation of AI system
  - *Stability*: demonstrate consistency of AI behavior over nominal scope
  - *Confidence awareness*: assess pertinence of inputs and predict uncertainty of output
  - *Adversarial resilience:* detect and provide stable output when inputs are modified by external processes
  - *Interpretability*: maximize user comprehension of causes leading to AI predictions
  - *Fairness*: demonstrate equitable outcomes to known subgroups and characterize residual biases
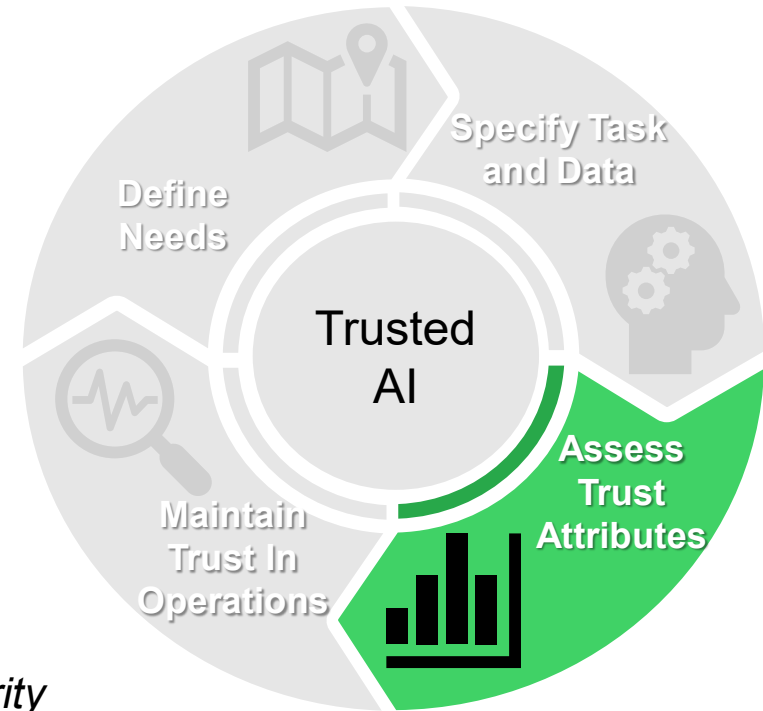  - *Familiarity*: comfort with which a user successfully operates system
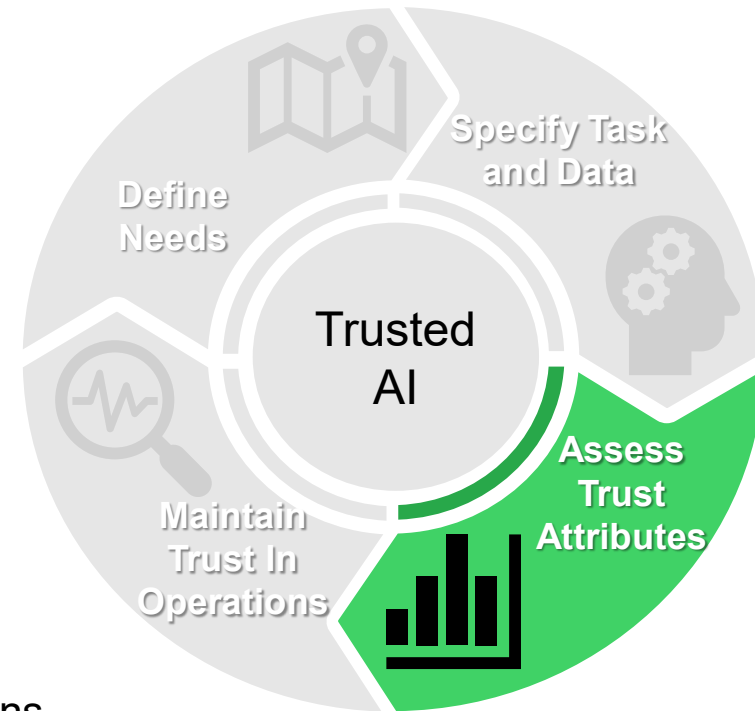
# Attributes of Trust

*Thread 2*

- Stability:
  - *Characterize performance on nominal scope, out-of-scope, and known challenge case data*
  - *Verify consistency of output when input is varied with background noise or inconsequential content*
  - *Consider third-party, blind V&V to assess stability of task and data specification*
  - *Leverage modern development practices to deploy across different platform (MLOps)*

- Confidence awareness:
  - *Determine if inputs are outside nominal, recording anomalies, and lowering prediction confidence (pertinence awareness)*
  - *Provide calibrated estimation of confidence when inputs are nominal*
  - *Provide ability to incorporate and propagate uncertainty from inputs*

- Adversarial resilience:
  - *Consider intents of nefarious actors or processes affecting data integrity*
  - *Evaluate against worst case deployment conditions*
  - *Assess sensitivity to range of attacks and attack strengths and train with them*

Define Needs

Specify Task and Data

Trusted AI

Assess Trust Attributes

Maintain Trust In Operations

# Attributes of Trust

*Thread 2*

- Interpretability:
  - *Consider expertise and training of users when providing interfaces to AI*
  - *Go beyond explanations - demonstrate that additional attributes contribute to user trust*
  - *Prefer algorithms with higher inherent interpretability*
  - *Prefer annotation methods that include concept attribution*
  - *Incorporate user input on relevant features, consider incorporation in model or UX*
  - *Provide evidence for prediction when requested by user:*
    - Display input or feature attributions
    - Display statistics of data and metadata
    - Query influential or similar training examples
  - *Study user engagement and get feedback on application utility:*
    - Identify attribution types that are most relevant
    - Record user expectation of model performance based on attributions (Familiarity)

Specify Task and Data

Define Needs

Trusted AI

Assess Trust Attributes

Maintain Trust In Operations
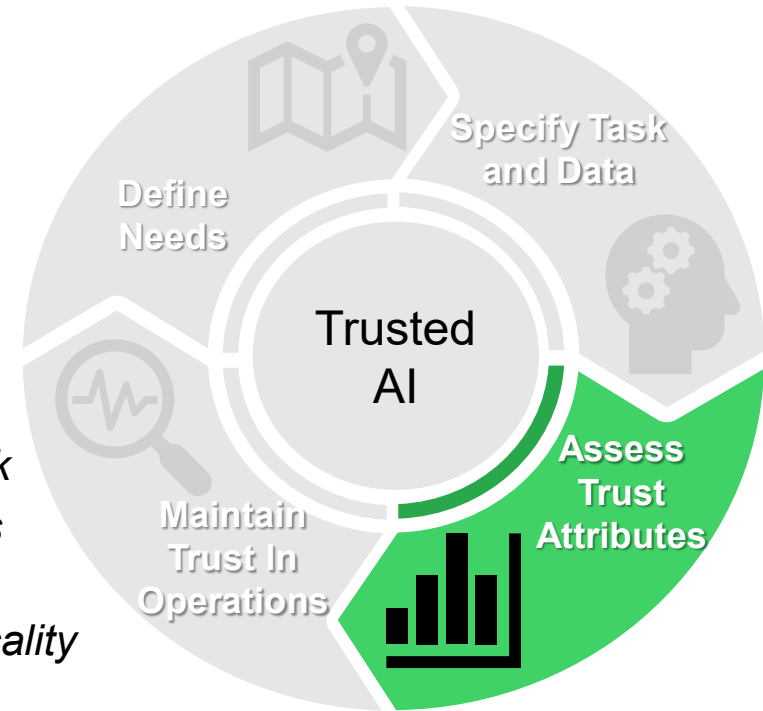
# Attributes of Trust

*Thread 2*

- Fairness:
  - *Leverage EDA to monitor subgroup disparities (Data Specification)*
    - Track performance metrics on subpopulations
    - Note degree of class separability and background diversity
  - *Augment disparities in data and annotation representation between subgroups*
  - *Estimate risk of unresolved biases in data or model*

- Familiarity:
  - *Facilitate early and frequent user interaction and incorporate feedback*
  - *Quantify and trend alignment between user actions and AI predictions (Pre-deployment)*
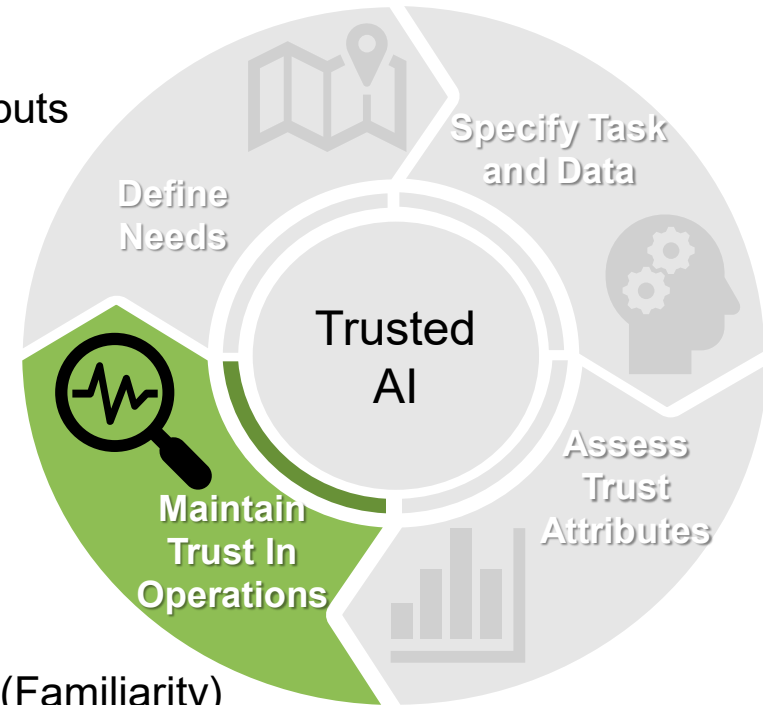  - *Consider AI validation through tasks that gradually increase task criticality*



Trusted AI

Define Needs

Specify Task and Data

Assess Trust Attributes

Maintain Trust In Operations

# *Deployment, Monitoring, and Control*

*Thread 3*

- Structured Deployment
  - *Note differences in target hardware and environment (Traceability)*
  - *Perform limited V&V to reaffirm task alignment and trust attributes:*
    - Verify modular sub-tasks adhere to expected performance
    - Evaluate stability and verify confidence calibration over nominal inputs
    - Record anomalous data and capture SME feedback
  - *Capture data representative of target environment*
    - Assess degree of distribution shift from training data (Data specification)
    - Quantify covariate and prior shift and update nominal / out-of-scope definitions
    - Assess risk for concept, conditional, and sensor shift and estimate impact
  - *Support gradual roll out of AI capability:*
    - Enable shadow mode operation for assessment of user alignment (Familiarity)
    - Support transition from limited to full operations
    - Deploy AI in roles of increasing autonomy and criticality

Trusted AI

Define Needs

Specify Task and Data

Assess Trust Attributes

Maintain Trust In Operations
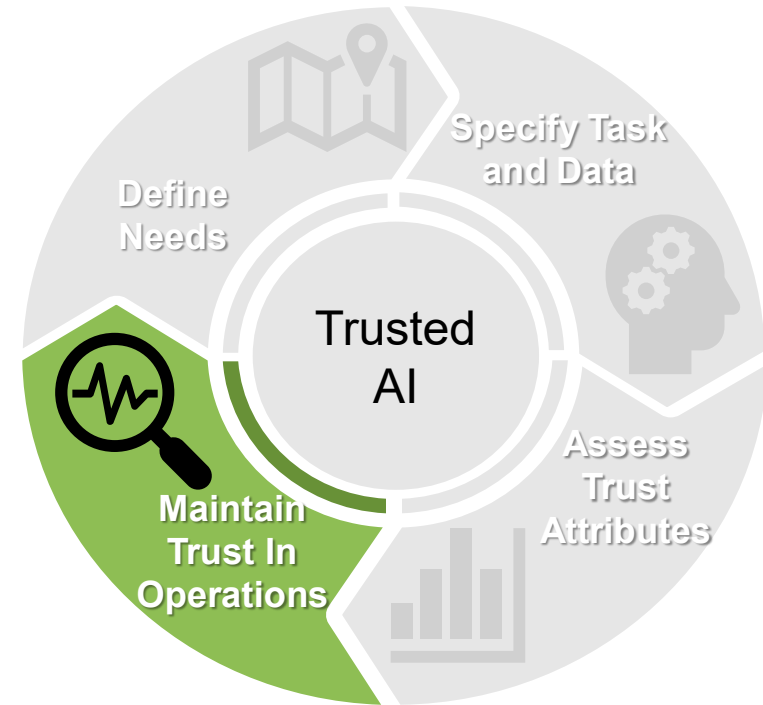
# Deployment, Monitoring, and Control

*Thread 3*

- Monitoring
  - *Develop upstream and downstream assurance systems to detect undesirable behavior*
  - *Provide metrics that track system degradation or system challenges*
    - Record data and AI prediction statistics over time
    - Identify out-of-scope data and record their frequency
    - Monitor computational trends such as runtime, convergence, memory usage, instrument quality

- Control
  - *Provide means for user intervention:*
    - Alert user to system degradations
    - Always obey user request for AI termination
  - *Include fallback systems for when AI termination occurs*
  - *Develop secondary assurance systems to prevent high risk behavior*
  - *Engage with users to consider means for re-fining AI behavior without re-training*
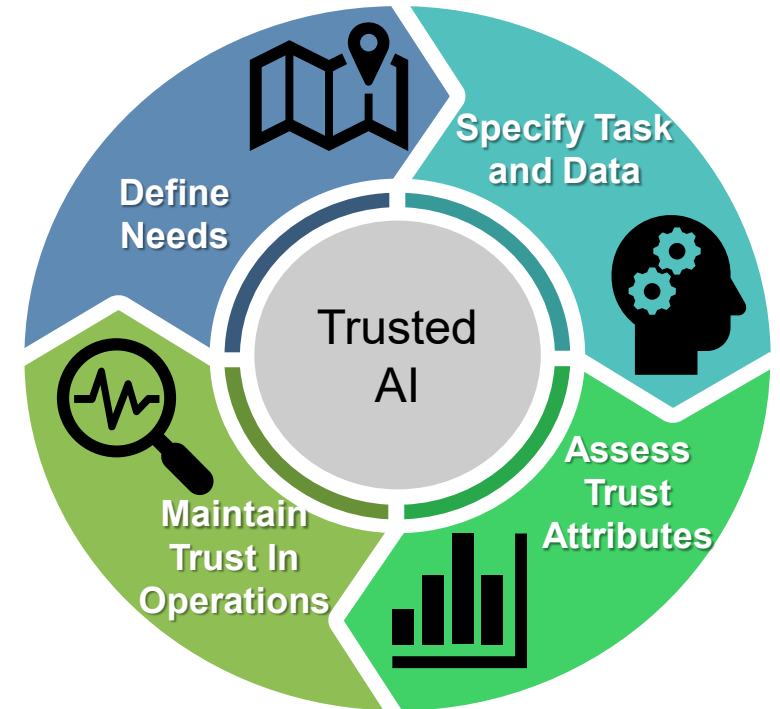
# Concluding Remarks

- **The promise of AI**:
  - *AI enabled capabilities are necessary in applications with limited opportunity for communication and control, harsh environments, and where some degree of autonomous discovery is required*

- **Where the Framework fits in**:
  - *Frame AI-based capabilities as trustable*
  - *Identify stakeholders that must be convinced*
  - *Actively manage expectations*
  - *Provide realistic roadmaps for how AI can be implemented and verified*
  - *Mitigate risk of unexpected AI behavior early in the development cycle*
  - *Prepare for the AI monitoring after deployment*

# *Backup*

- Assess current capabilities
- Degree of autonomy
- Identify risks

- Objective Specification
- Model Specification
- Data Specification

**Thread 0:** Define Needs

**Thread 1:** Specify Objectives and Data

**Thread 3:** Deploy and Maintain Trust in Operations

**Thread 2:** Assess Against Trust Attributes

- Monitoring
- Control
- Deployment

- Stability
- Confidence Awareness
- Lineage
- Robustness
- Interpretability
- Fairness
- Familiarity