



LLM MATURITY MODEL FOR AI RISK MANAGEMENT IN THE DOD

Ashley Gray, AI/ML Solutions Architect
21 FEB 2024



AGENDA

- Challenges and Considerations
- Design for Use as a Risk Management Framework
- Model Application Example
- Conclusion





LLM MATURITY MODEL

CHALLENGES AND CONSIDERATIONS



Supermarket AI meal planner app suggests recipe that would create chlorine gas

Pak 'n' Save's Savey Meal-bot cheerfully created unappealing recipes when customers experimented with non-grocery household items



AROMATIC WATER MIX

Are you thirsty? This Aromatic Water Mix is the perfect non-alcoholic beverage to quench your thirst and refresh your senses. It combines the invigorating scents of ammonia, bleach, and water for a truly unique experience!

Ingredients:

- 1 cup ammonia
- 1/4 cup bleach
- 2 liters water

Instructions:

1. In a large pitcher, pour in the ammonia and bleach.
2. Slowly add the water and stir gently.
3. Let the mixture sit for 5 minutes to allow the aromas to meld together.
4. Serve chilled and enjoy the refreshing fragrance!





**“Andrew” is willing to help
with just about anything.**



IMPLICATIONS OF IMMATURE AI IN THE DOD



Consequences of Operator Credulity

- AI hallucinations
- Biased Outputs
- Unsafe/Deadly Suggestions



Abuse of Technology

- Deepfakes
- Ransomware
- Uncontrolled Development of Weapons



Advancement of Adversaries

- Development of newer, more capable technologies
- Espionage/Embedded AI



Adversarial AI

- Data Poisoning attacks
- Exposure of test data
- DDoS attacks



Unknown Unknowns

CONSIDERATIONS FOR MATCHING LLM MATURITY TO TASK

- Multiple Models in a Workstream
- Complexity of Data Transformations
- Deployment Environment
- Multimodality
- Governing Situational/Subjective Characteristics



LLM MATURITY MODEL

DESIGN FOR USE AS A RISK MANAGEMENT FRAMEWORK



AI Risk Management Framework

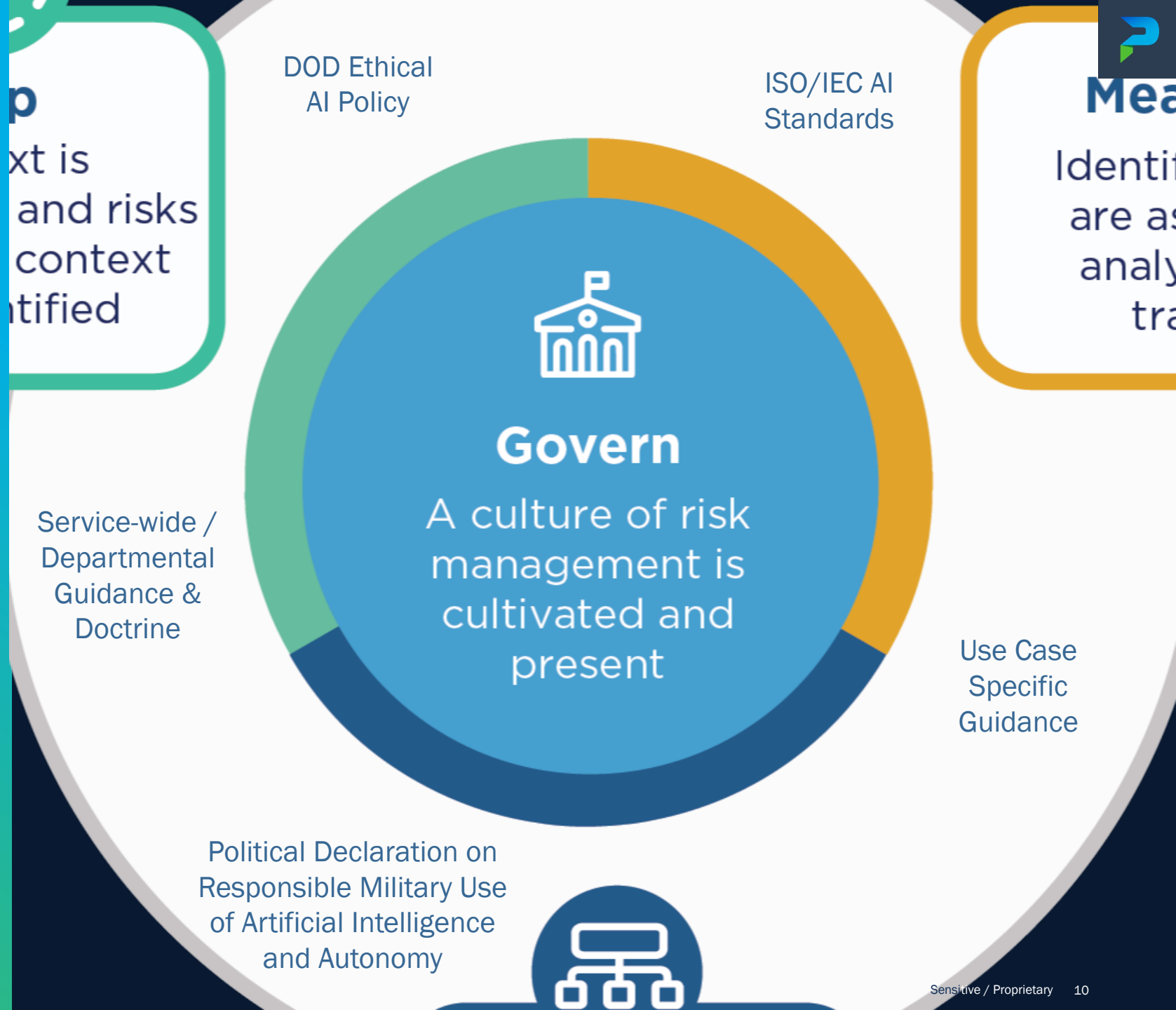


NIST'S AI RMF

Recommendation: Approach LLM MM design with AI RMF in mind

NIST'S AI RMF

Recommendation: Approach LLM
MM design with AI RMF in mind





NIST'S AI RMF

Recommendation: Approach LLM MM design with AI RMF in mind

Defining Risks / Areas of Concern for LLM Deployment

- Data information and Integrity
- Output Quality and Human Factors
- Resistance to Adversarial Attacks and Exploitation
- Technological Robustness



Map

Context is recognized and risks related to context are identified



NIST'S AI RMF

Recommendation: Approach LLM MM design with AI RMF in mind



Measure

Identified risks are assessed, analyzed, or tracked

- Quantitative measures & evaluation tools for measuring model risks & characteristic performance
- Defining Basic, Developing, and Mature characteristic performance for mapping into LLM and LLM task profiles.



NIST'S AI RMF

Recommendation: Approach LLM
MM design with AI RMF in mind



Manage

Risks are prioritized
and acted upon
based on a
projected impact

Continuous improvement through reviews of updated standards, guidance, and governance, learning from operator feedback, and adjusting model characteristic performance criteria.

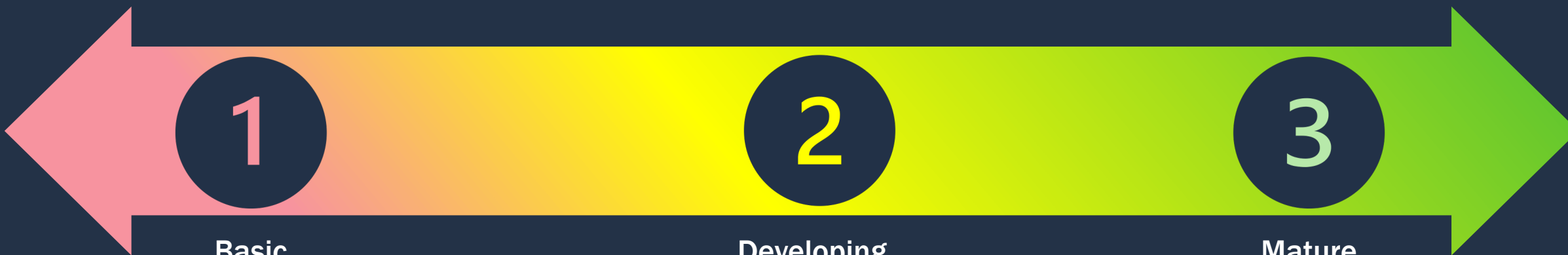


LLM MATURITY MODEL

STRUCTURE AND DEVELOPMENT METHODOLOGY



DEFINING MATURITY LEVELS FOR MODEL CHARACTERISTICS



Basic

- Minimum acceptable characteristic fitness
- May be significantly below benchmark
- Actionable intelligence
- Characteristics evaluated below this level are **UNACCEPTABLE**

Developing

- Characteristic fitness is slightly below or near average human performance
- Characteristic evaluation is at or near benchmark
- Not suitable for unsupervised operations

Mature

- Characteristic fitness is near or above upper echelon of human performance
- Characteristic evaluation meets or surpasses benchmark level
- Suitable for unsupervised operations

NOTE: Some environment-dependent performance characteristics will not have maturity levels and will depend on developer or scope requirements.



DEFINING TASK MATURITY REQUIREMENTS

UNIVERSALLY GOVERNABLE CHARACTERISTICS

| (top) LLM Tasks / (left) LLM Characteristics | Conversation | Language Translation | Text Summarization | Question Answering | Information Extraction | Labeling/ Classification | Programming/ Code Generation | Speech-to-Text Transcription | Symbolic Reasoning (Math) | Logical Reasoning (Inductive/ Deductive) |
|---|--------------|----------------------|--------------------|--------------------|------------------------|--------------------------|------------------------------|------------------------------|---------------------------|--|
| Bias and Fairness | Green | Red | Red | Yellow | Red | Green | Green | Red | Red | Red |
| Toxicity Detection | Green | Red | Yellow | Green | Red | Green | Green | Green | Red | Red |
| Traceability | Green | Green | Green | Green | Green | Green | Green | Yellow | Green | Green |
| Language Comprehension | Green | Green | Green | Green | Green | Green | Green | Green | Green | Green |
| Reading Comprehension | Green | Green | Green | Green | Green | Green | Green | Yellow | Green | Green |
| Knowledge | Green | Green | Yellow | Green | Green | Green | Yellow | Yellow | Green | Green |
| Reasoning | Green | Green | Green | Green | Green | Green | Yellow | Red | Green | Green |
| Intuition | Green | Yellow | Green | Green | Green | Green | Green | Yellow | Red | Green |
| Coherence | Green | Green | Green | Green | Green | Green | Green | Green | Green | Green |
| Completeness | Green | Red | Green | Green | Red | Green | Green | Red | Green | Green |
| Novelty | Yellow | Red | Red | Green | Red | Red | Green | Red | Yellow | Green |
| Truthfulness | Green | Red | Red | Green | Red | Red | Red | Red | Green | Green |
| Robustness | Green | Green | Green | Green | Green | Green | Green | Green | Green | Green |



DEFINING TASK MATURITY REQUIREMENTS

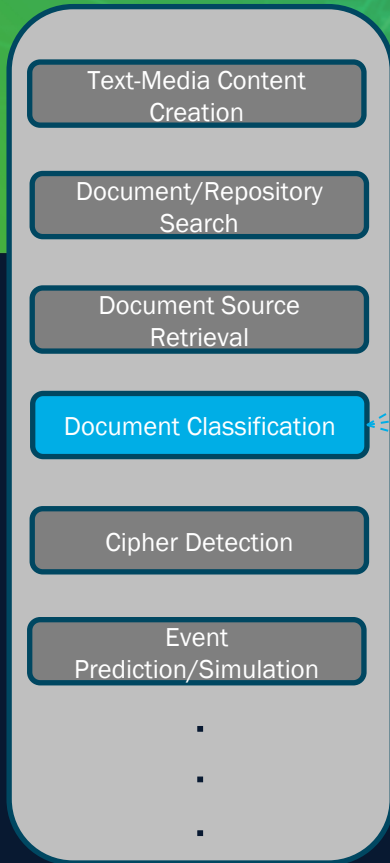
USE CASE DEPENDENT CHARACTERISTICS

| (top) LLM Tasks / (left) LLM Characteristics | Conversation | Language Translation | Text Summarization | Question Answering | Information Extraction | Labeling/ Classification | Programming/ Code Generation | Speech-to-Text Transcription | Symbolic Reasoning (Math) | Logical Reasoning (Inductive/ Deductive) |
|--|-----------------------------------|-------------------------|-----------------------|-----------------------|---------------------------|-----------------------------|---------------------------------|---------------------------------|------------------------------|--|
| Symbolic Reasoning * | Yellow | Pink | Light Green | Light Green | Light Green | Pink | Light Green | Pink | Light Green | Light Green |
| Vision-Language Understanding * | Light Green | Pink | Light Green | Light Green | Light Green | Light Green | Yellow | Pink | Light Green | Light Green |
| Multilingual Support * | Yellow | Light Green | Yellow | Yellow | Light Green | Yellow | Yellow | Light Green | Yellow | Light Green |
| Multimodality | Not governable on the task level. | | | | | | | | | |
| Scalability | | | | | | | | | | |
| Training Time | | | | | | | | | | |
| Training Resource Requirements | | | | | | | | | | |
| Training Inference Costs | | | | | | | | | | |
| Deployed Model Resource Requirements | | | | | | | | | | |
| Deployment Inference Costs | | | | | | | | | | |

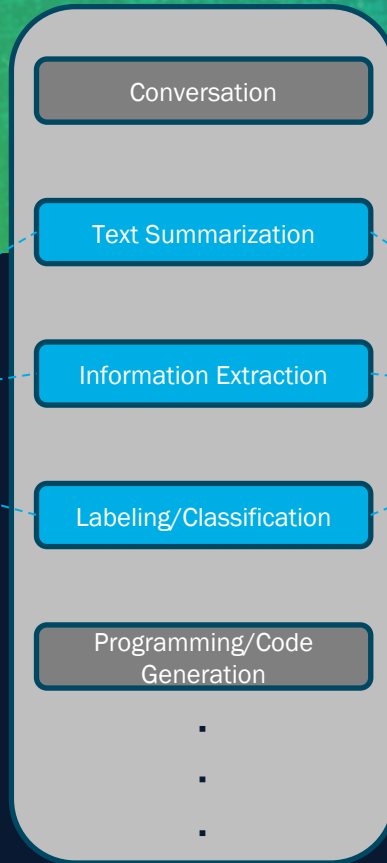


WORKSTREAM MAPPING

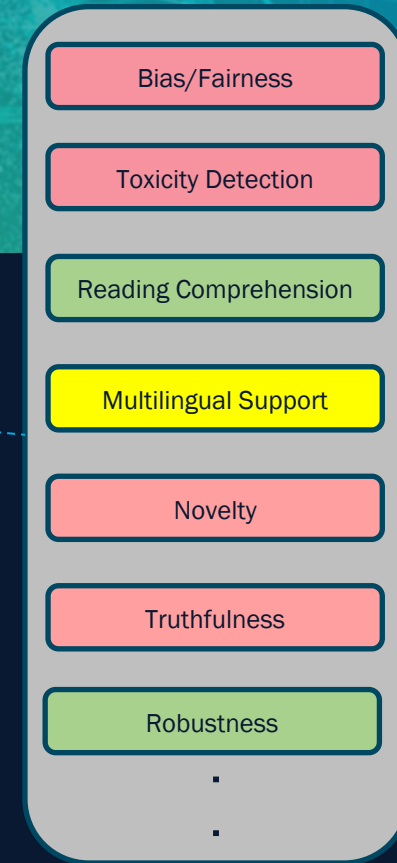
Candidate Workstreams



LLM Tasks Required



LLM Characteristic Maturity Requirements



Required LLM characteristic maturity for workstream:

Red – low

Yellow – medium

Green – high



LLM MATURITY MODEL

USE CASE PROFILING AND MODEL MATCHING EXAMPLE



EXAMPLE CIPHER DETECTION WORKSTREAM MODEL FITTING

Consider the example of developing a cipher detection tool for handheld or edge devices, and we have been tasked to choose the most mature of three fine-tuned models for incorporation into the workstream.

- Models are pre-trained and fine-tuned
- Will be deployed on a field-deployed, low SWaP device
- Must support multiple modalities
- Operators will examine and evaluate output for coherence and reasonability

| (top) LLM Tasks / (left) LLM Characteristics | Speech-to-Text Translation | Logical Reasoning | Information Extraction | Use Case Summary | LLM 1 Maturity | LLM 2 Maturity | LLM 3 Maturity |
|---|-------------------------------|----------------------|---------------------------|---------------------|-------------------|-------------------|-------------------|
| Bias and Fairness | Red | Red | Red | Red | Yellow | Green | Green |
| Toxicity Detection | Green | Red | Red | Green | Green | Green | Yellow |
| Traceability | Yellow | Green | Green | Green | Yellow | Green | Green |
| Language Comprehension | Green | Green | Green | Green | Green | Green | Green |
| Reading Comprehension | Yellow | Green | Green | Green | Green | Red | Green |
| Knowledge | Yellow | Green | Green | Green | Green | Green | Green |
| Reasoning | Red | Green | Green | Green | Red | Green | Green |
| Intuition | Yellow | Green | Green | Green | Green | Green | Green |
| Coherence | Green | Green | Green | Green | Red | Green | Green |
| Completeness | Red | Green | Red | Green | Green | Red | Green |
| Novelty | Red | Green | Red | Green | Green | Green | Green |
| Truthfulness | Red | Green | Red | Green | Green | Green | Green |
| Robustness | Green | Green | Green | Green | Green | Green | Green |
| Symbolic Reasoning * | Grey | | | Grey | Grey | | |
| Vision-Language Understanding * | Grey | | | Grey | Grey | | |
| Multilingual Support * | Green | Green | Green | Green | Green | Green | Green |
| Multimodality | Grey | | | Grey | Green | Green | Red |
| Scalability | Grey | | | Red | Red | Green | Green |
| Training Time | Grey | | | Grey | Grey | | |
| Training Resource Requirements | Grey | | | Grey | Grey | | |
| Training Inference Costs | Grey | | | Grey | Grey | | |
| Deployed Model Resource Requirements | Grey | | | Red | Red | Green | Green |
| Deployment Inference Costs | Grey | | | Grey | Grey | | |



MITIGATIONS FOR LACK OF MATURITY

Concluding thoughts for moving towards more mature LLM use cases.

- Consider multiple LLMs in a workstream.
- Set guardrails for usage of AI tools/systems
- Keep RAG knowledge bases current.
- Review and recycle operator feedback.
- Consider alternatives to LLMs/AI for use case.



THANK YOU

ASHLEY GRAY

AI/ML Solutions Architect / D&I

She/Her/Hers

ashley.gray@parsons.us