

Large Language Model (LLM) Maturity Model

Chief Digital and Artificial Intelligence Office (CDAO)



BUILDING ENDURING DECISION ADVANTAGE

ADVANTAGE DOD 24

Defense Data & AI Symposium

FEBRUARY 20-22, 2024

#ADOD2024



February 14, 2024



Submitted To

Advantage DoD 2024, CDAO Symposium
osd.pentagon.cdao.mbx.engagement@mail.mil

Prepared By

iWorks Corporation
1889 Preston White Drive, Suite 100, Reston, VA 20191
Phone: (703) 636-6777 | Fax: (866) 574-2310 | www.iworkscorp.com
UEI: ZGHSJN13NVZ5 | DUNS: 16-9988578 | CAGE: 4Q7A1

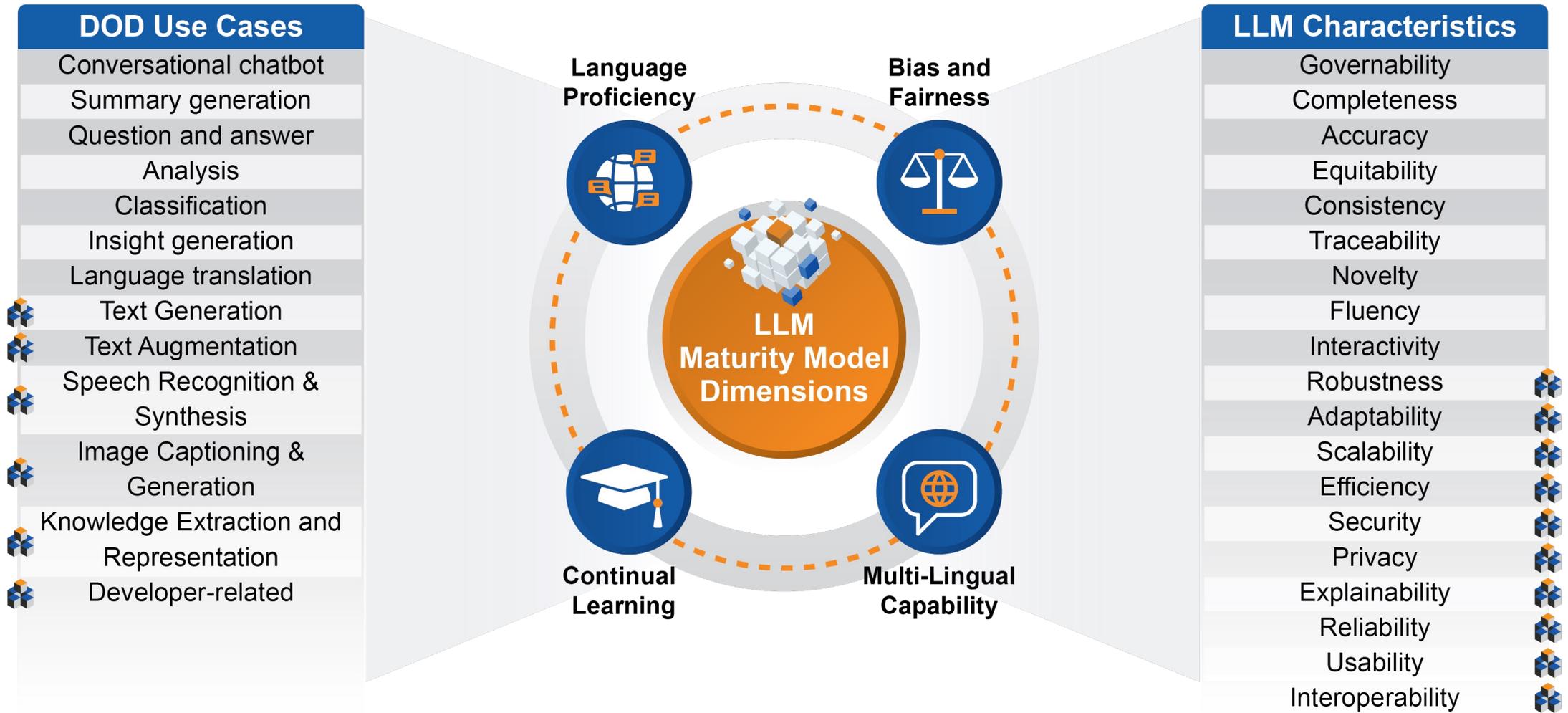


CMMIDEV / 4SM
Exp. 2023-02-14 / Appraisal #6082



CMMISVC / 3SM





Text Generation

- Generation of contextually-relevant text
- Production of high-quality, semantically relevant text

Speech Recog. & Synthesis

- Needed for tasks such as voice assistance, automated dispatch, etc.
- Responding naturally and conversationally

Know. Extraction & Rep.

- Identify relevant information from unstructured data sources, such as text, images, or videos
- Extract and transform as structured knowledge for more efficient and interpretable representation

Text Augmentation

- Generation of new data samples by applying various transformations to existing data samples
- Improves robustness and accuracy

Image Caption & Generation

- Needed to generate human-like descriptions of visual content - descriptive and contextually relevant
- Supports language and visual processing for complex tasks

Developer-Related

- Needed to assess the design, deployment, and use of generative AI technologies responsibly and securely
- Supports mission effectiveness by automating certain tasks, reducing human error, and increasing efficiencies

LLM Maturity Models - Dimensions

- Accuracy and Fluency
- Semantic Understanding
- Multilingual Competence
- Contextual Appropriateness
- User-Friendly Interfaces
- Fine-Tuning

Language Proficiency



Bias and Fairness



- Bias Detection and Mitigation
- Ethical Guidelines
- Data Representativeness
- Continuous Monitoring
- User Feedback Mechanisms
- Interpretability and Transparency

LLM Maturity Model Dimensions

LLM MM SCORING RUBRIC FOR 4 DIMENSIONS



- User Interaction
- Data Incorporation
- Fine-Tuning
- Version Control
- Model Evolution
- Ethical Considerations

Continual Learning

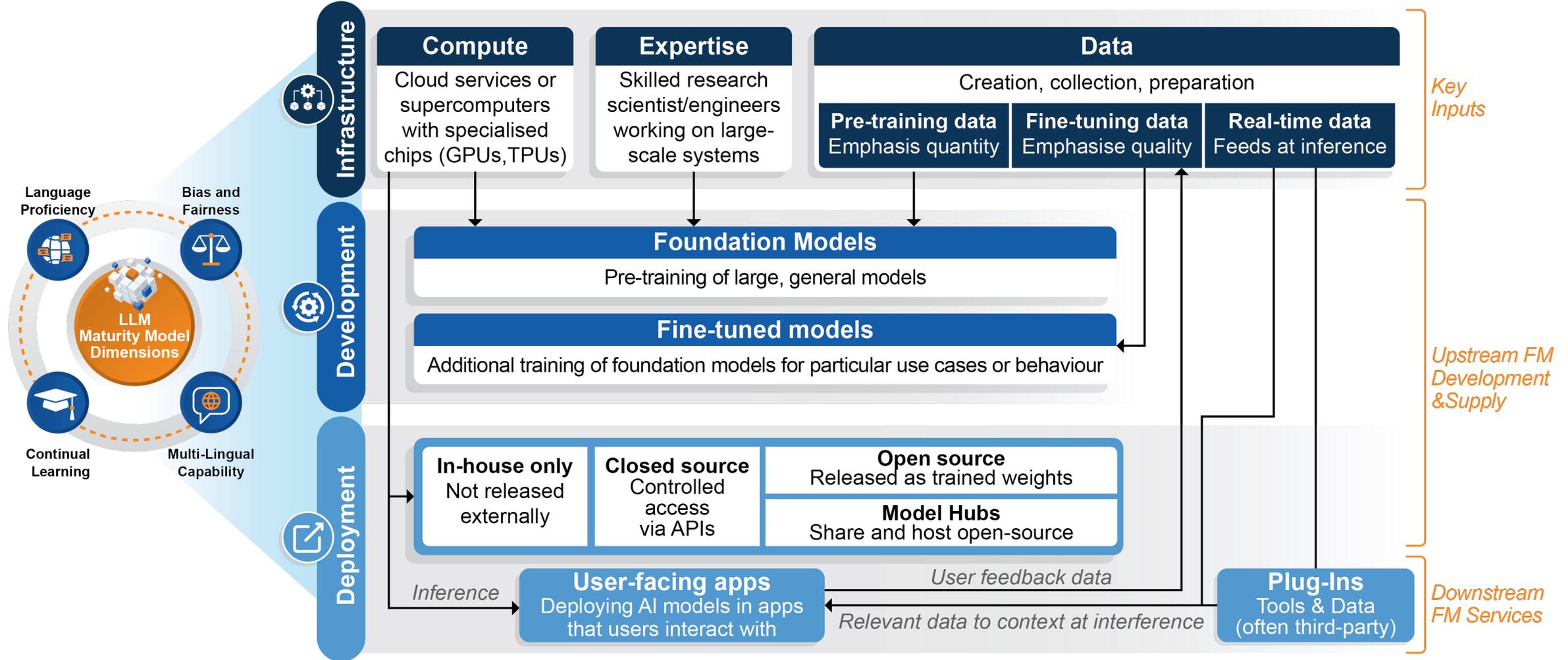


Multi-Lingual Capability



- Linguistic Proficiency
- Language Coverage
- Cross-Language Transfer
- Language Pairs
- Cultural Sensitivity
- Localization

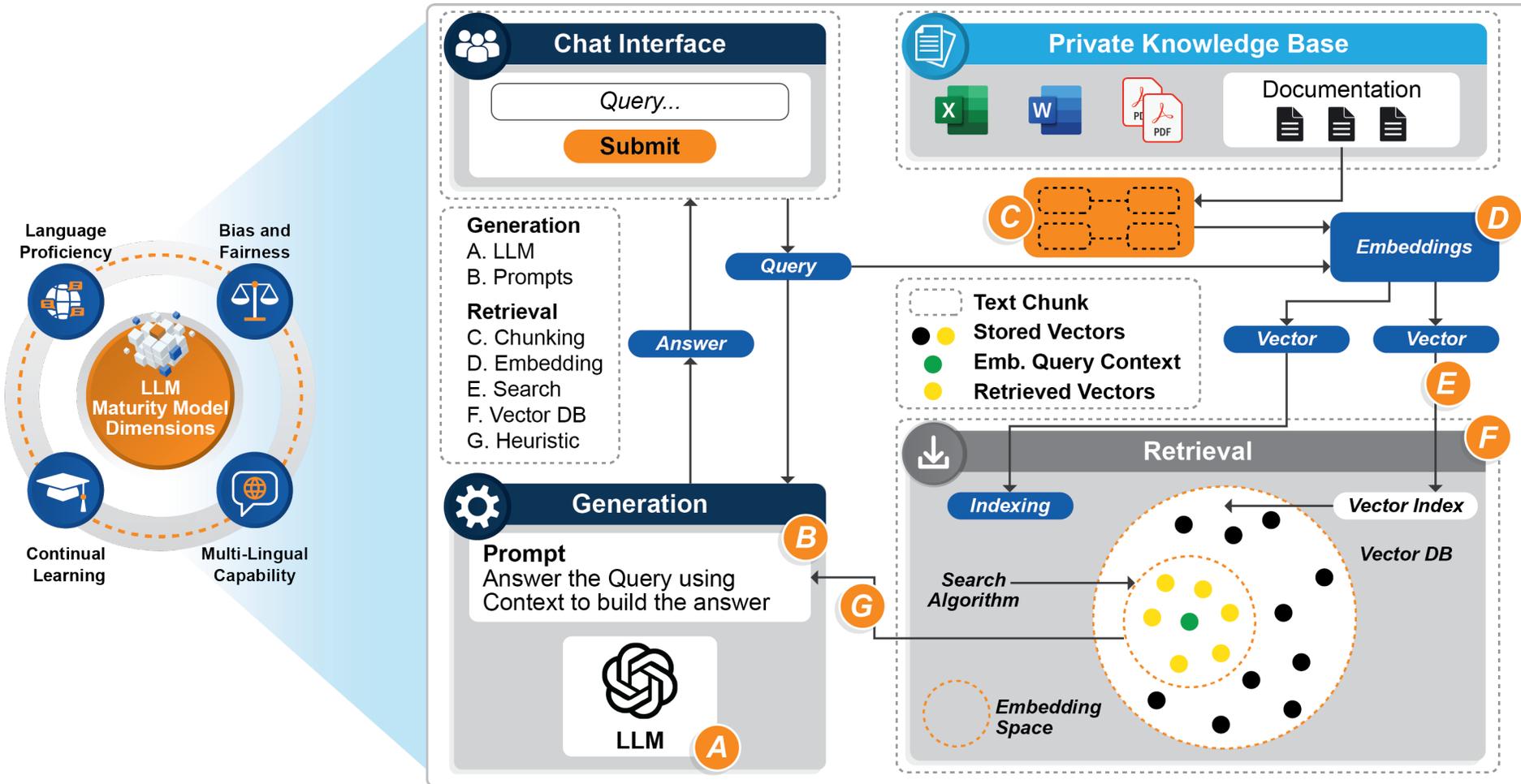
LLM Model Evaluations: Infrastructure, Development & Deployment



Source: Competition and Markets Authority, AI Foundation Models Review (2023): https://assets.publishing.service.gov.uk/media/65045590dec5be00dc35f77/Short_Report_PDF.pdf

Check and evaluate LLM Architectural components across infrastructure, development and deployment to assess current and sustained abilities

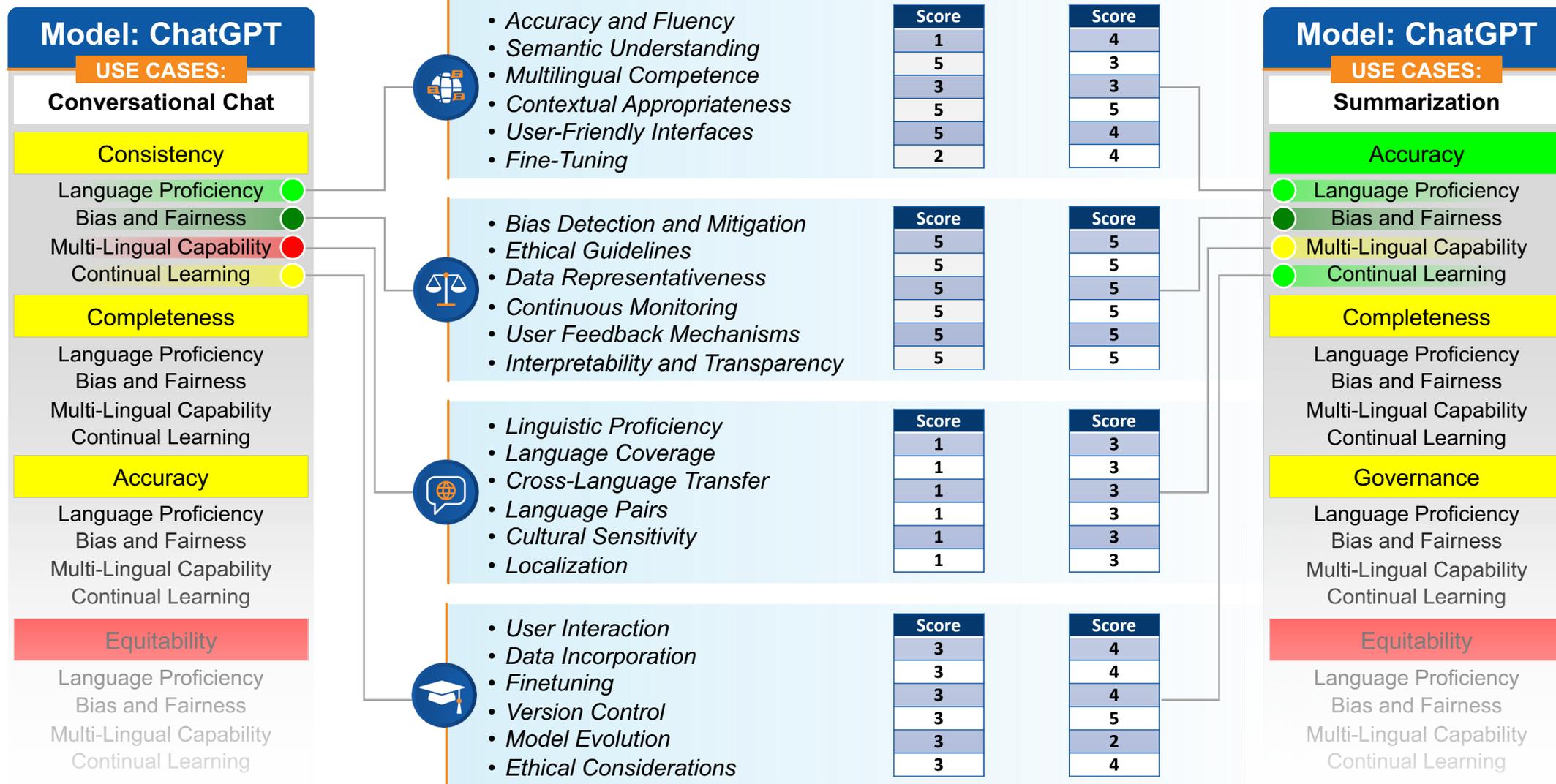
LLM Model Evaluations: Use Case Execution Workflow Example



Rubric for LLM MM Assessment / Ranking across the four critical dimensions

Evaluation Metrics	Very Low	Low	Satisfactory	High	Very High
 Language Proficiency					
 Bias and Fairness					
 Multi-Lingual Support					
 Continual Learning					

LLM Maturity Models – Evaluation by Use Case



Heatmap – LLM Characteristics

LLM Characteristics

Use Case	Conversation/Chatbot
Governability	Yellow
Completeness	Yellow
Accuracy	Yellow
Equitability	Red
Consistency	Yellow
Traceability	Red
Novelty	Grey
Fluency	Green
Interactivity	Green
...	Red

Evaluation Metrics	Very Low	Low	Satisfactory	High	Very High
Language Proficiency	The LLM struggles to understand and generate human-like text in the given language.	The LLM can understand and generate basic text in the given language but struggles with more complex or nuanced language.	The LLM can understand and generate most text in the given language, but may still struggle with some complex or nuanced language.	The LLM can understand and generate almost all text in the given language, incl. complex and nuanced language.	The LLM can understand and generate all text in the given language, incl. complex and nuanced language, with a high degree of accuracy.
Bias and Fairness	The LLM perpetuates harmful social biases and produces text that is highly unfair and biased.	The LLM may perpetuate some harmful social biases and produce text that is somewhat unfair and biased.	The LLM generally avoids perpetuating harmful social biases and produces text that is mostly fair and unbiased.	The LLM avoids perpetuating harmful social biases and produces text that is fair and unbiased.	The LLM avoids perpetuating all harmful social biases and produces text that is completely fair and unbiased.
Multi-Lingual Support	The LLM can only understand and generate text in one language.	The LLM can understand and generate text in multiple languages but may struggle with some languages or produce less accurate text in some languages.	The LLM can understand and generate text in most languages but may still struggle with some languages or produce less accurate text in some languages.	The LLM can understand and generate text in almost all languages, including complex and nuanced language, with a high degree of accuracy.	The LLM can understand and generate text in all languages, including complex and nuanced language, with a high degree of accuracy.
Continual Learning	The LLM does not learn or adapt to new data or situations.	The LLM learns and adapts to new data or situations but may do so slowly or with limited effectiveness.	The LLM learns and adapts to new data or situations at a moderate pace and with moderate effectiveness.	The LLM learns and adapts to new data or situations quickly and with a high degree of effectiveness.	The LLM learns and adapts to new data or situations continuously and with a high degree of effectiveness.

Use Case	Summarization
Governability	Green
Completeness	Yellow
Accuracy	Yellow
Equitability	Yellow
Consistency	Yellow
Traceability	Red
Novelty	Grey
Fluency	Green
Interactivity	Green
...	Red

Not of significance to this use case

Reference Heat Map for Evaluation

DOD Use Cases	LLM Characteristics																		
	Governability	Completeness	Accuracy	Equitability	Consistency	Traceability	Novelty	Fluency	Interactivity	Robustness	Adaptability	Scalability	Efficiency	Security	Privacy	Explainability	Reliability	Usability	Interoperability
Conversational chatbot	Green	Green	Green	Green	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Summary generation	Green	Green	Green	Green	Green	Green	Grey	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Question and answer	Green	Green	Green	Green	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Analysis	Green	Green	Grey	Green	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Classification	Green	Green	Green	Green	Green	Green	Grey	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Insight generation	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Language translation	Green	Green	Green	Green	Green	Green	Grey	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Text Generation	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Text Augmentation	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Speech Recog. & Synthesis	Green	Green	Green	Green	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Image Captioning & Generation	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Know. Extraction & Repr.	Green	Green	Green	Green	Green	Green	Green	Green	Grey	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
Developer -Related	Light Blue	Light Blue	Light Blue	Light Blue	Additional 15 development-focused characteristics										Light Blue	Light Blue	Light Blue	Light Blue	Light Blue

Developer-related use cases have an additional set of 15 development-focused characteristics (ex. code generation, summarization, etc.)

Q&A

Contact: srini@iworkscorp.com
 srini@computer.org