# TEST AND EVALUATION OF LARGE LANGUAGE MODELS TO SUPPORT INFORMED GOVERNMENT ACQUISITION

Presented by Dr. Erin Lanus and Dr. Jaganmohan "Jagan" Chandrasekaran

22nd Annual Acquisition Research Symposium & Innovation Summit

May 07, 2025

Coauthors: Mr. Brian Mayer, Dr. Heather Frase, Dr. Patrick Butler, Dr. Stephen Adams, Mr. Jared Gregersen, Dr. Naren Ramakrishnan and Dr. Laura Freeman

# Acknowledgement

# Introduction to Large Language Models

## Language modeling

**Imagine the following task:** Predict the next word in a sequence

( The cat likes to sleep in the ____ ) → What **word** comes next?

**Can we frame this as a ML problem?** Yes, it's a classification task.

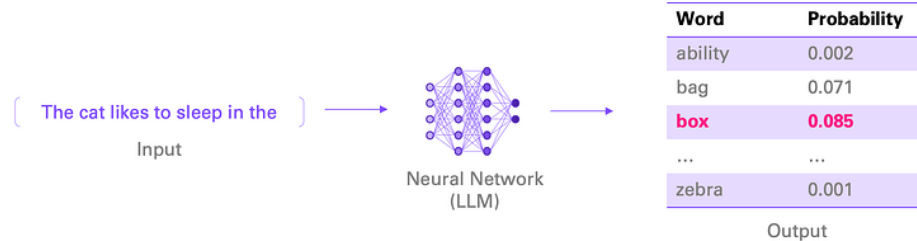*Now we have (say) ~50,000 classes (i.e. words)*

( The cat likes to sleep in the )
Input
→ Neural Network (LLM) →

| Word | Probability |
|------|-------------|
| ability | 0.002 |
| bag | 0.071 |
| **box** | **0.085** |
| ... | ... |
| zebra | 0.001 |

Output

Image source: https://medium.com/data-science-at-microsoft/how-large-language-models-work-91c362f5b78f

**Intelligent Machines**
Broadly defined

**Pattern Recognition**
Learning general patterns from data

**Neural Networks**
Learning general patterns in **unstructured** data (i.e. images, text, audio, etc.)

**Large Language Models**
Learning to understand natural language (i.e. text)

Artificial Intelligence — Machine Learning — Deep Learning — LLMs

ChatGPT  Gemini  Llama  MISTRAL AI_  GPT-4

LLMs offer flexibility in performing diverse functions, distinguishing them from traditional AI/ML systems.

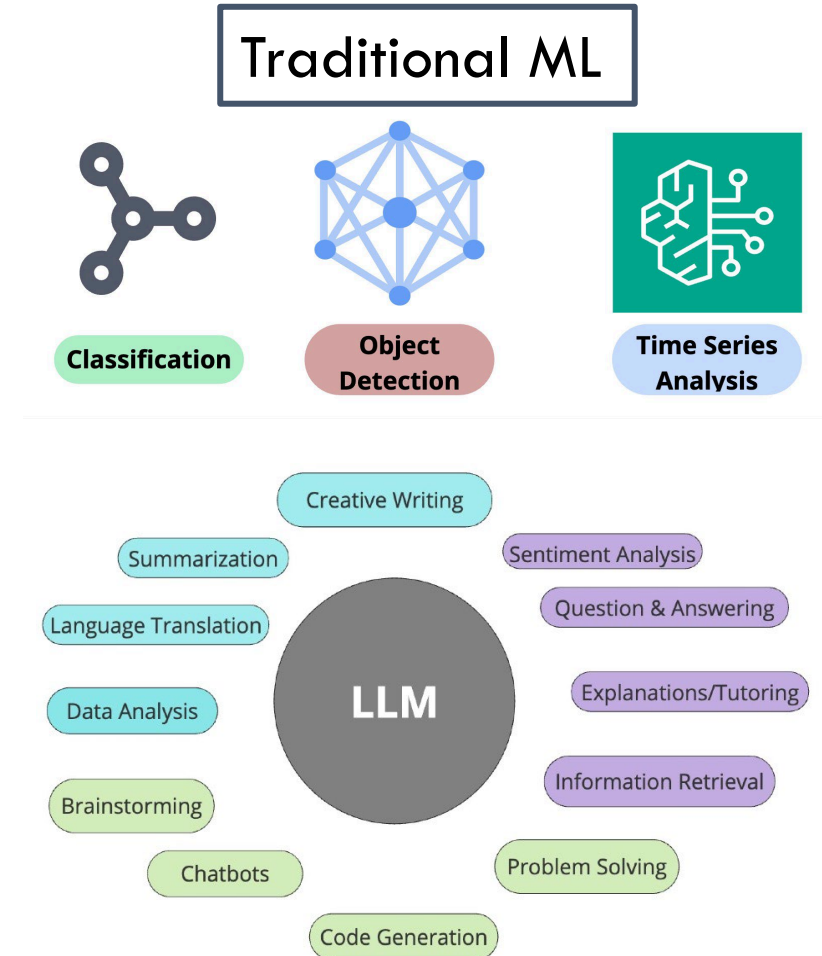**LLM's versatility is great, but comprehensive test and evaluation (T&E) are key to ensure reliable, trustworthy, and safe behavior.**

The ability to do many different functions increases the difficulty and the necessary variability in testing LLMs.



Traditional ML

Classification    Object Detection    Time Series Analysis

Creative Writing
Summarization
Sentiment Analysis
Language Translation
Question & Answering
LLM
Data Analysis
Explanations/Tutoring
Brainstorming
Information Retrieval
Chatbots
Problem Solving
Code Generation

## What does the current T&E landscape inform us about the evaluation of LLMs?

**T&E Objective:** Can an LLM generate <u>correct,</u> <u>contextually relevant</u> responses?

Steps in testing LLM

| Installing Prerequisites | → | Set up LLM | → | Loading Datasets | → | Prompting | → | Assessment/ Evaluation |

**Access mode -** To perform inferencing, practitioners either

- Host the LLMs locally
- Interact via Application Programming Interface

**Parameters –** A set of values influencing the LLM's outcome
- Temperature
- Top-p
- Max tokens
- Frequency penalty

**Prompt -** A set of instructions informing the LLM about the user's request
- Zero-shot prompting
- Few-shot prompting
- Chain-of-thought prompting

## Capabilities

- Abstract functional abilities of an LLM
- *Examples: understanding, reasoning, generation*

## Tasks

- Concrete implementations used to assess specific capabilities
- *Examples: Question Answering, Multiple Choice Question, Code Generation*
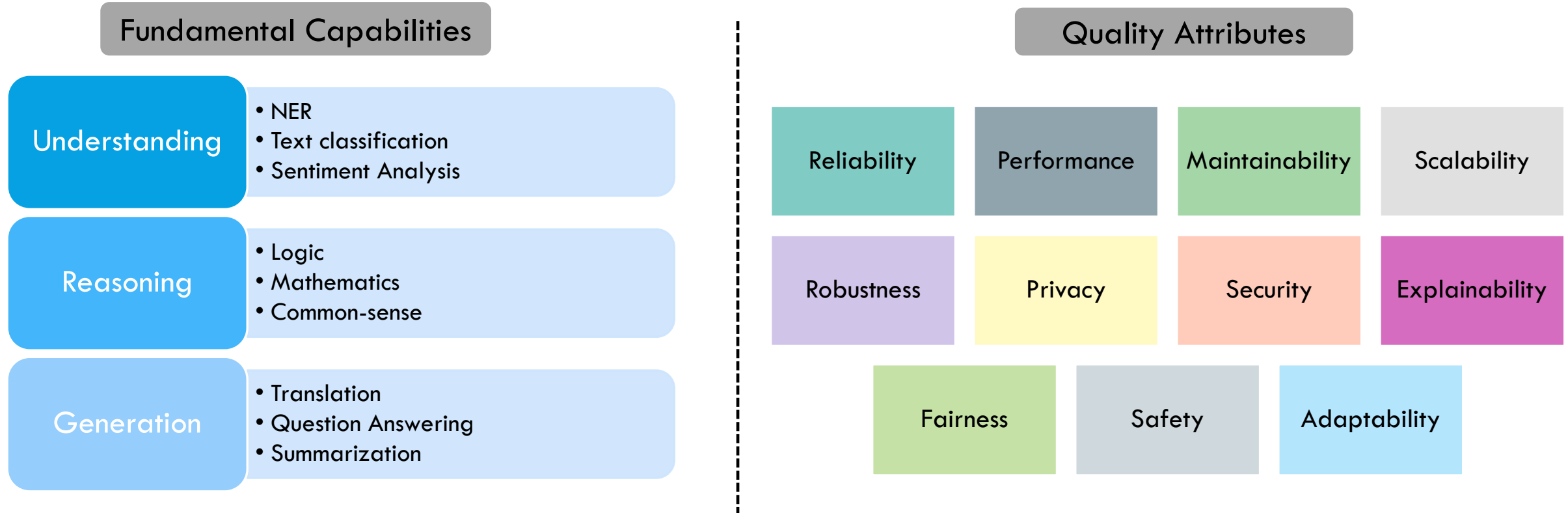
## Benchmarks

- Standardized datasets that measure performance
- *Examples: MMLU, HELM*

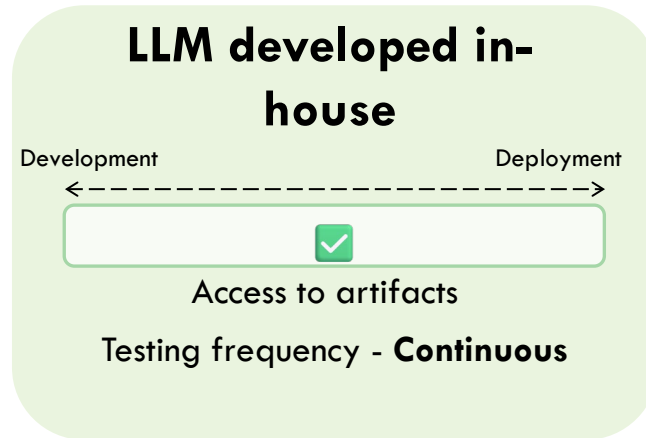| Capability | Task Type | Benchmarks | Objective |
|---|---|---|---|
| Understanding | Named Entity Recognition | CoNLL 2003 | Evaluate LLM's basic word-level understanding and categorization abilities |
| Reasoning | Multiple Choice Question | MMLU | Assess high-school level reasoning abilities on variety of subjects |
| Generation | Code Generation | HumanEval | Test LLM's ability to generate software code. |

# LLM Evaluation Framework

A comprehensive evaluation of LLMs must include two primary dimensions:
- Evaluation of **fundamental capabilities** in facilitating human-like interactions
- As a software component, the LLM's ability to meet expected **software quality standards**

## Fundamental Capabilities

**Understanding**
- NER
- Text classification
- Sentiment Analysis

**Reasoning**
- Logic
- Mathematics
- Common-sense

**Generation**
- Translation
- Question Answering
- Summarization

## Quality Attributes

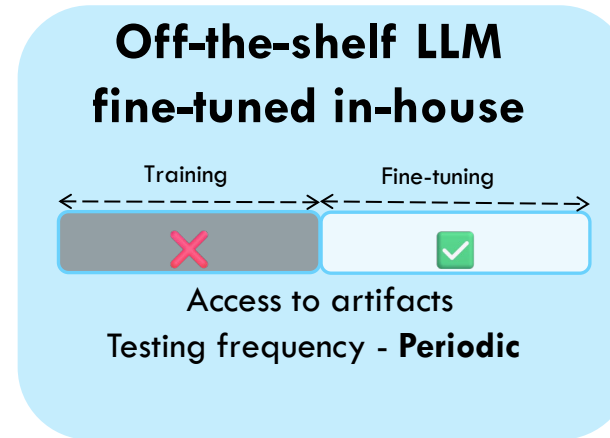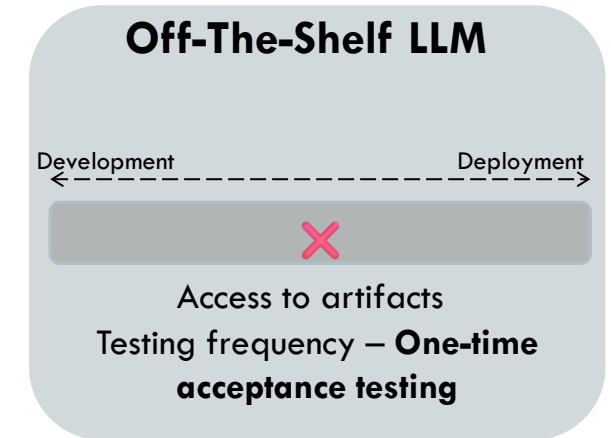| Reliability | Performance | Maintainability | Scalability |
| --- | --- | --- | --- |
| Robustness | Privacy | Security | Explainability |
| Fairness | Safety | Adaptability | |

# Example Acquisition Scenarios

**Use case 1** - Identify named entities in a user-provided collection of records, and extract relationships between entities.

**Use case 2** – Extract information from user-provided records and question responses.

**Use case 3** - The drones have an LLM that converts text messages into commands that they can implement.



Understanding, Reasoning, Reliability, Scalability
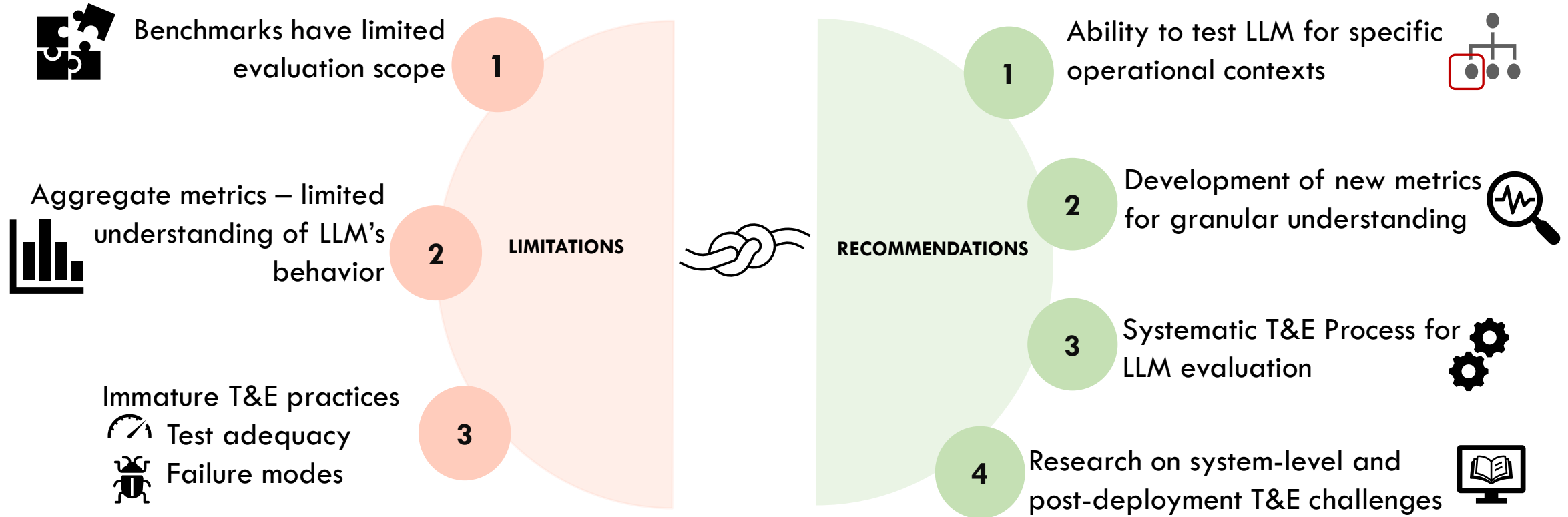
Security

Generation

Performance

Adaptability

Maintainability

Fairness

Robustness

# Limitations and Recommendations



**LIMITATIONS**

1. Benchmarks have limited evaluation scope

2. Aggregate metrics – limited understanding of LLM's behavior

3. Immature T&E practices
   - Test adequacy
   - Failure modes

**RECOMMENDATIONS**

1. Ability to test LLM for specific operational contexts

2. Development of new metrics for granular understanding

3. Systematic T&E Process for LLM evaluation

4. Research on system-level and post-deployment T&E challenges

# Thank you!