

# 4 AGENT EVALUATION FRAMEWORKS

# A Comprehensive Guide to Evaluating Agents

# INTRODUCTION

## Memory Management

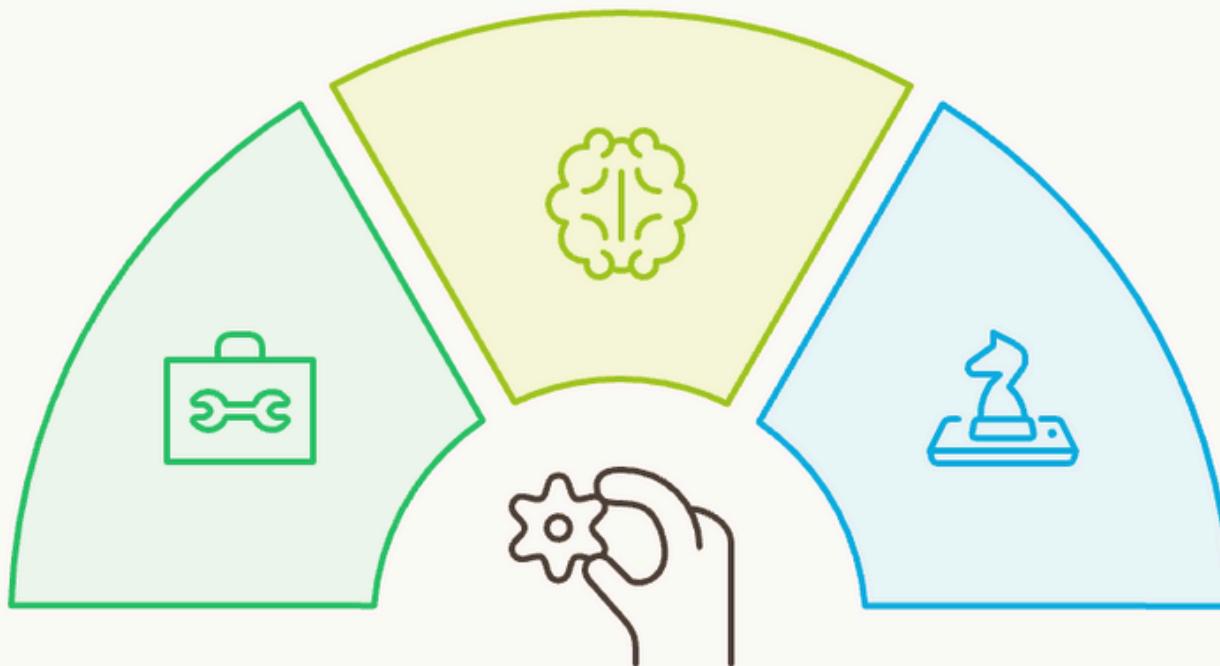
Retains context and learns user preferences to enhance interactions.

## Tool Usage

Utilizes external tools for efficient task execution and automation.

## Strategic Task Planning

Breaks down tasks into steps and autonomously executes workflows.



Agentic workflows are advanced AI-driven systems designed to **automate tasks** by following instructions, **making decisions, and taking actions** independently.

Agentic workflows empower systems to go beyond static programming, enabling **dynamic, intelligent actions to save time**, reduce manual work, and maximize efficiency.

Agentic workflows adapt in real time, such as in supply chain management, where they adjust to disruptions or demand shifts without human intervention. This adaptability helps businesses streamline processes and stay agile.

# WHY EVALUATE?



Evaluating agentic workflows is essential for ensuring their effectiveness, reliability, and efficiency

A well-evaluated workflow **saves time, cuts costs**, and ensures superior results

- **Quality Assurance**

- Improves accuracy and consistency.
- Identifies errors in decision-making and planning.
- Builds trust with reliable output.

- **Performance Benchmarking**

- Measures and compares system capabilities.
- Assesses speed, adaptability, and accuracy.
- Identifies top-performing processes.

- **Cost Optimization**

- Reduces inefficiencies and resource wastage.
- Focuses on automating high-impact tasks.
- Balances quality and cost efficiency.

# 4 EVALUATION FRAMEWORKS

Agent as Judge



## 1. Agent as a Judge

AI systems evaluate their outputs for consistency and relevance.

AAEF Framework



## 2. AAEF (Agentic Application Evaluation Framework)

A structured framework to evaluate agent applications across multiple dimensions.

.

Mosaic AI Tool



## 3. Mosaic AI for Agent Evaluation

A robust tool providing unified evaluation metrics for agent workflows.

WORFEVAL Protocol

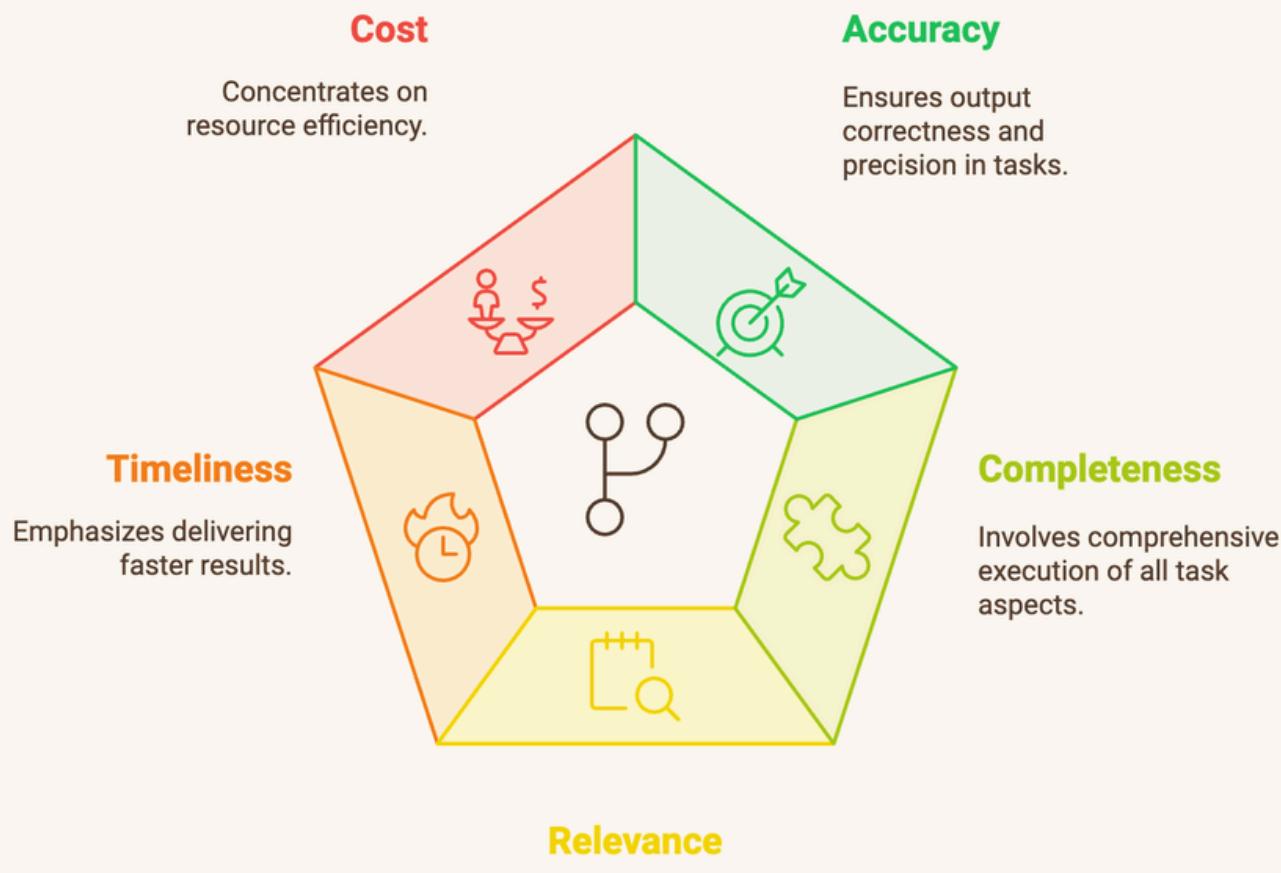


## 4. WORFEVAL Evaluating Workflow Generation

A protocol designed for assessing workflow generation by LLM agents.

# AGENT AS A JUDGE

"Agent-as-a-Judge": AI agents trained on datasets evaluate other AI outputs using criteria like:



## 1. Accuracy

Ensures output correctness by comparing results to ground truths and minimising error rates. Impact: Builds trust in reliable task completion.

## 2. Completeness

Verifies full execution of tasks, ensuring no steps or dependencies are missed. Impact: Delivers holistic task performance.

## 3. Relevance

Aligns outputs with user expectations and real-world needs through feedback and analysis. Impact: Boosts user satisfaction and workflow value.

## 4. Timeliness

Tracks task speed, comparing against benchmarks for faster results. Impact: Saves time and enhances efficiency.

## 5. Cost Efficiency

Optimizes resource use by identifying inefficiencies and refining workflows. Impact: Balances costs with high-impact automation.

# AAEF

## Agentic Application Evaluation Framework

A framework designed to assess the performance of agentic systems by evaluating their effectiveness, efficiency, and adaptability in specific applications or tasks. It has the following features:

### Tool Utilization Efficacy



Also known as TUE and measures the AI's ability to select and use the most appropriate tools effectively.

### Memory Coherence and Retrieval



Assesses how consistently the workflow manages and retrieves stored information.

### Strategic Planning Index



Evaluates the workflow's capability to plan and sequence tasks strategically.

### Component Synergy Score



Gauges how well different components of the workflow interact to achieve a common goal.

## It Works in Real Life

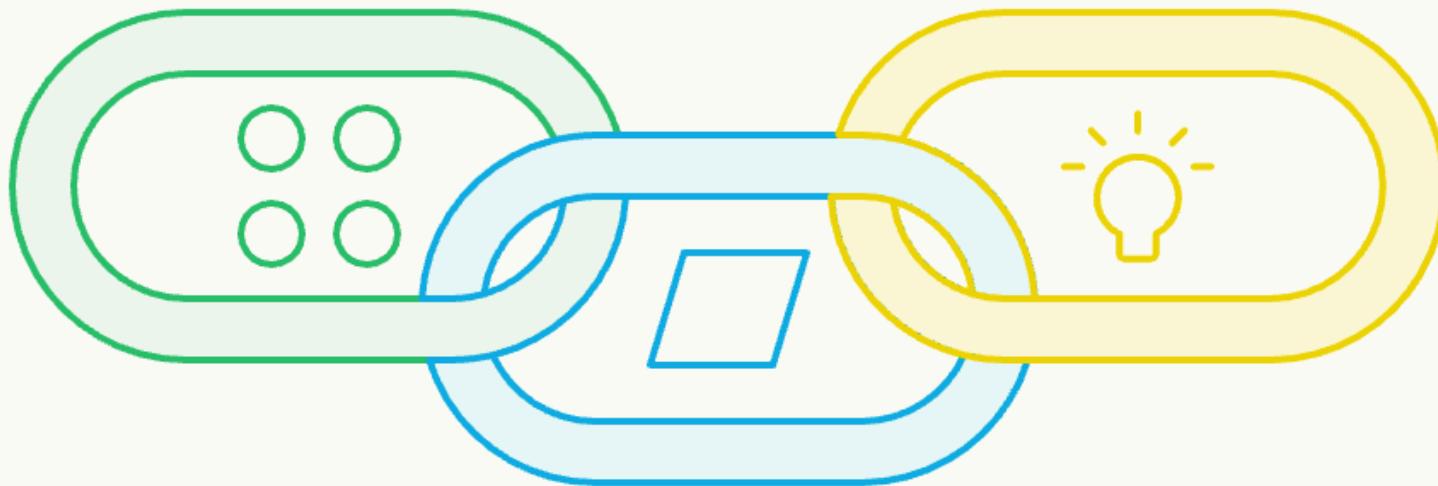
### Example Use Case: 🚛 Logistics Optimization AI

- **Task:** Route planning for delivery trucks..
- **TUE Score:** 0.954 for selecting the most efficient mapping and scheduling tools
- **Outcome:** Reduced fuel consumption and delivery times by implementing optimised routes, saving costs and improving customer satisfaction. (explained in-depth in comments)

# MOSAIC AI

Unified Metrics

Actionable Insights



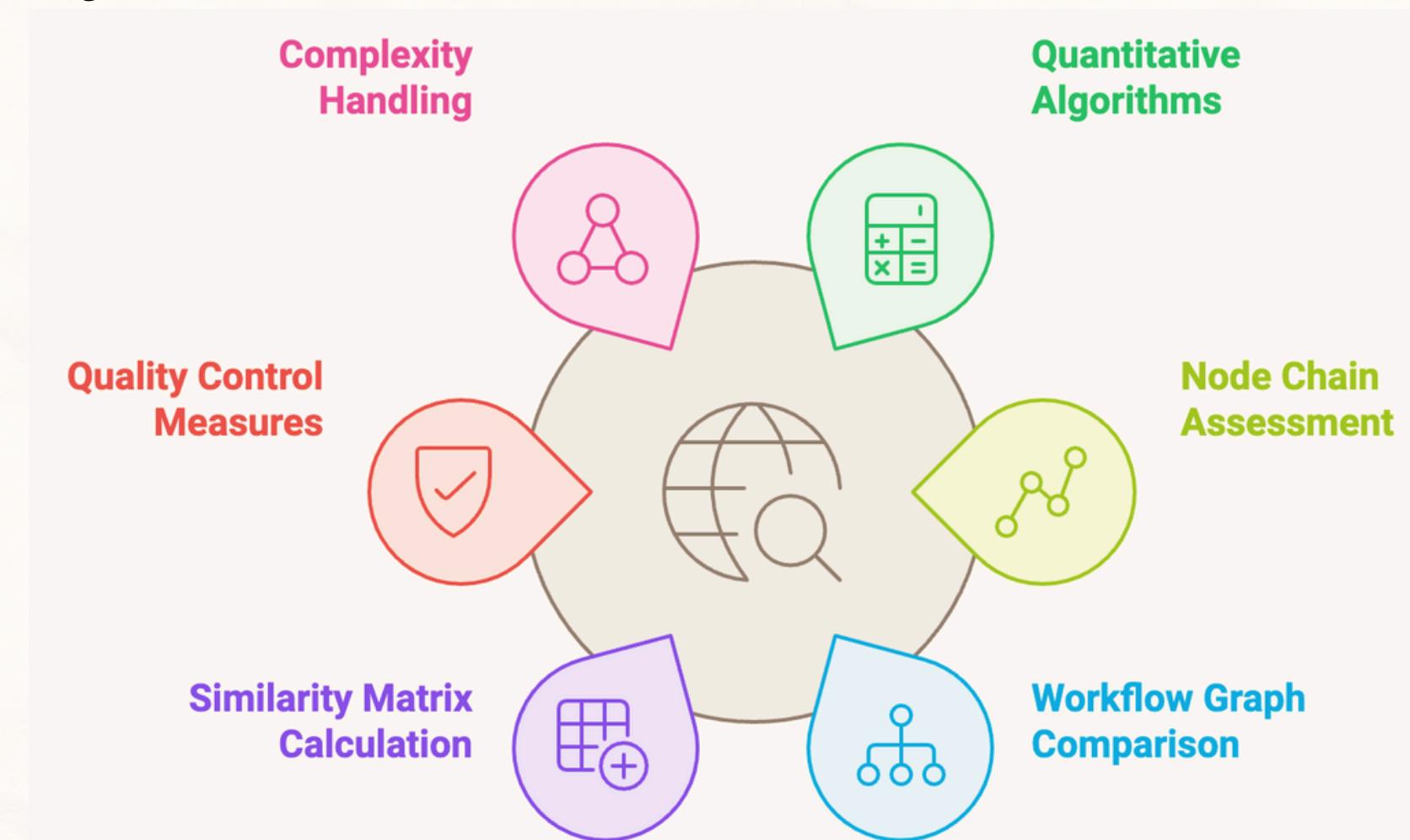
MLflow Integration

Mosaic AI is a comprehensive framework empowering robust agent evaluation through:

- **Unified Metrics:** Seamlessly integrates key performance indicators across the MLOps lifecycle, like Accuracy, Precision, Recall, F1-Score, AUC-ROC, etc
- **MLflow Integration:** Tracks experiments, metrics, and workflows for enhanced data intelligence.
- **Actionable Insights:** Provides clear directions for refining agentic workflows and boosting efficiency.

# WORFEVAL

WORFEVAL is a systematic evaluation protocol that assesses the workflow generation capabilities of LLM agents by comparing their generated workflows to a gold standard using quantitative algorithms.



## Evaluation Methodology

- Utilizes quantitative algorithms like subsequence and subgraph matching for precise analysis.
- Assesses the node chain and workflow graph generated by LLM agents.

## Evaluation Components

- Compares predicted nodes and edges with a gold standard (correct workflow).
- Calculates a similarity matrix using cosine similarity between predicted and actual nodes.

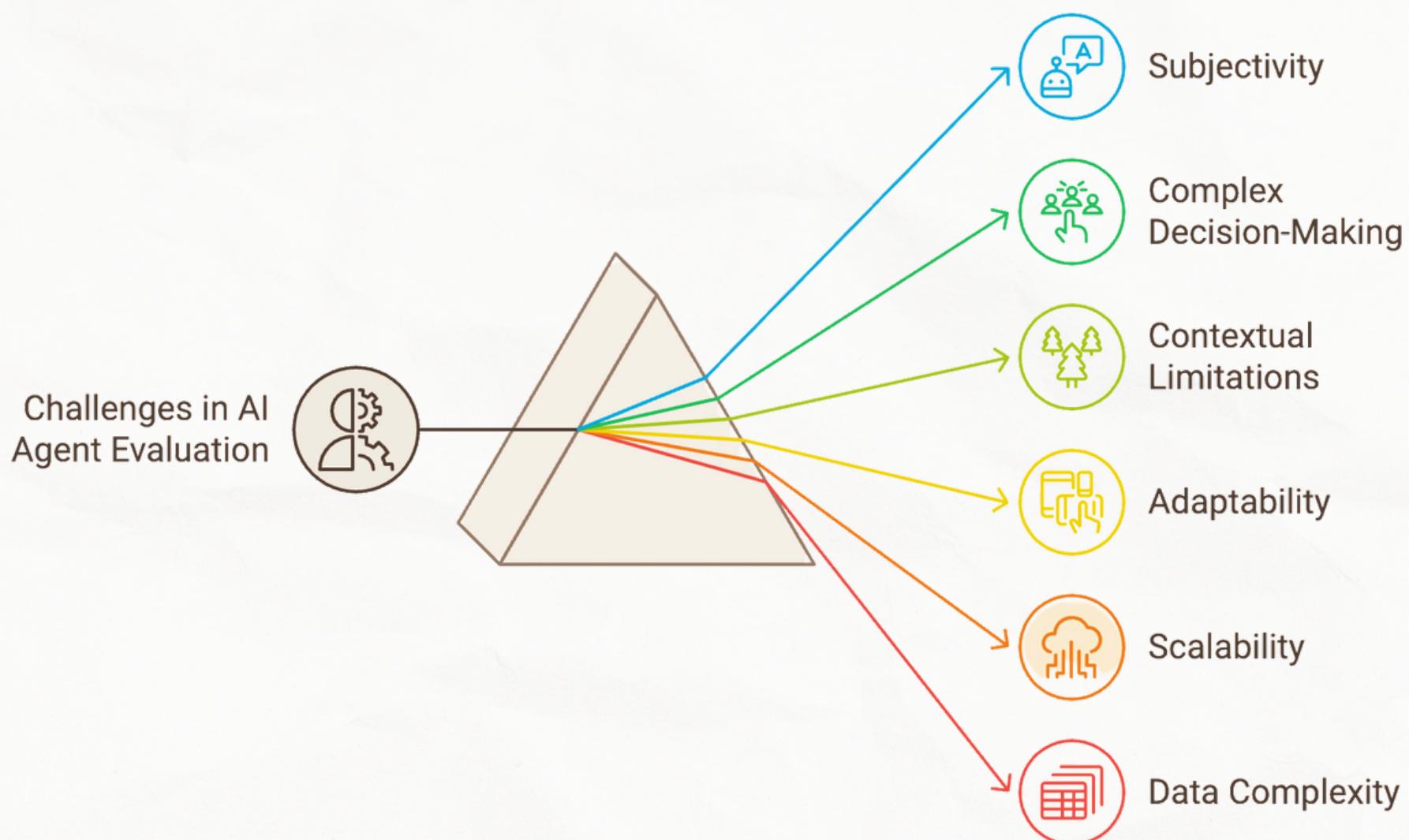
## Quality Control

- Implements strict measures, including data filtering and manual checks, to ensure fairness and eliminate biases.

## Complexity and Scope

- Designed to handle complex structures like Directed Acyclic Graphs (DAGs).
- Evaluates LLM agents in multi-faceted scenarios, simulating real-world tasks.

# CHALLENGES



## Agent as a Judge

- Bias: AI evaluations may reflect biases from the training data.
- Complexity: Struggles with nuanced tasks requiring deep contextual understanding.

## AAEF (Agentic Application Evaluation Framework)

- Context: May struggle with dynamic, rapidly changing tasks.
- Adaptability: Needs complex adjustments for diverse domains.

## Mosaic AI for Agent Evaluation

- Scalability: Handling large datasets and real-time evaluation can be challenging.
- Data Complexity: Integrating and standardizing data from multiple sources.

## WORFEVAL Evaluating Workflow Generation

- Complexity: Difficulty assessing intricate, non-linear workflows.
- Gold Standard: Relies on the availability of an accurate gold standard workflow



**Follow to stay updated on  
Generative AI**



**LIKE**



**COMMENT**



**REPOST**