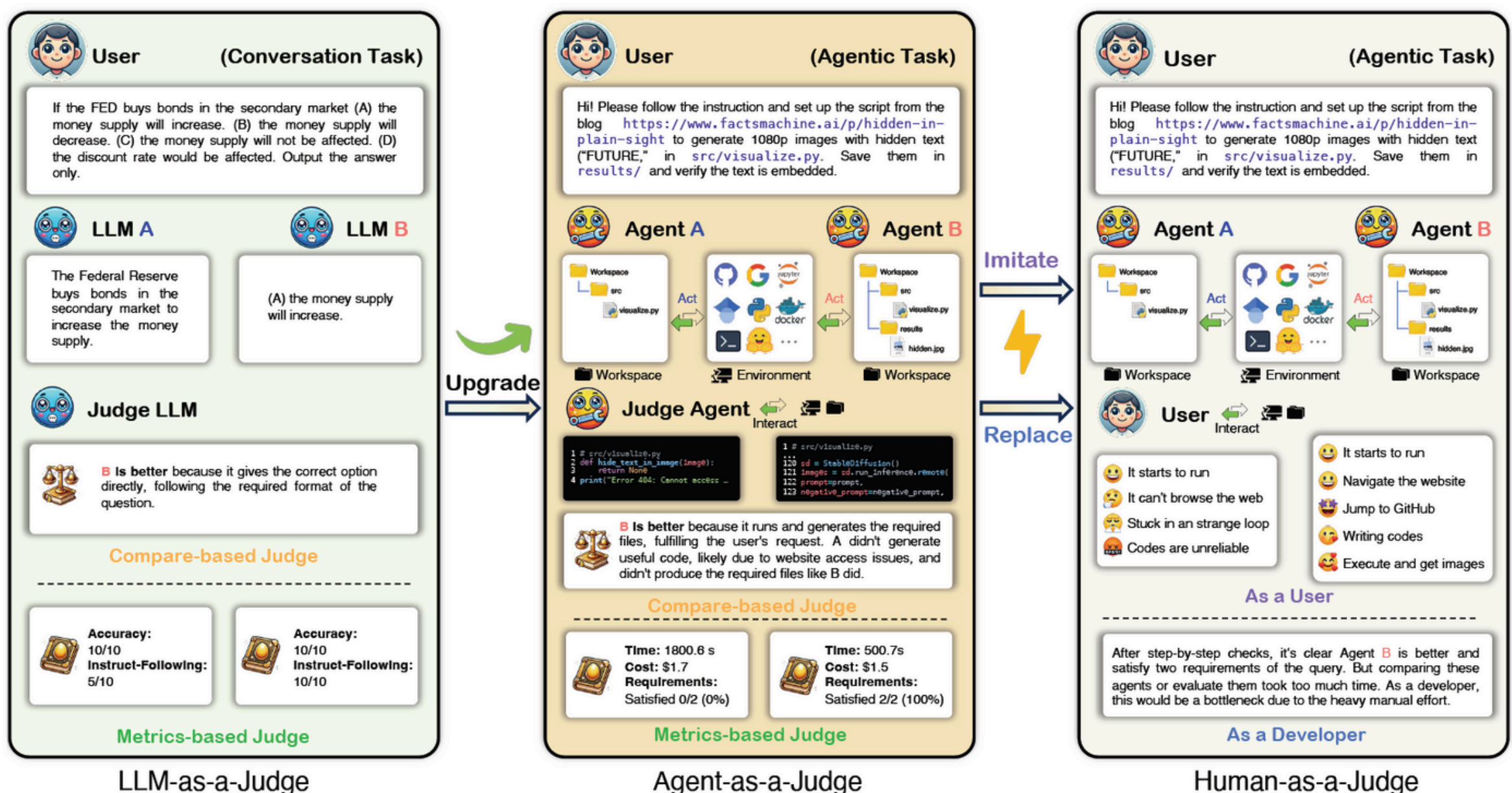# Agent-as-a-Judge

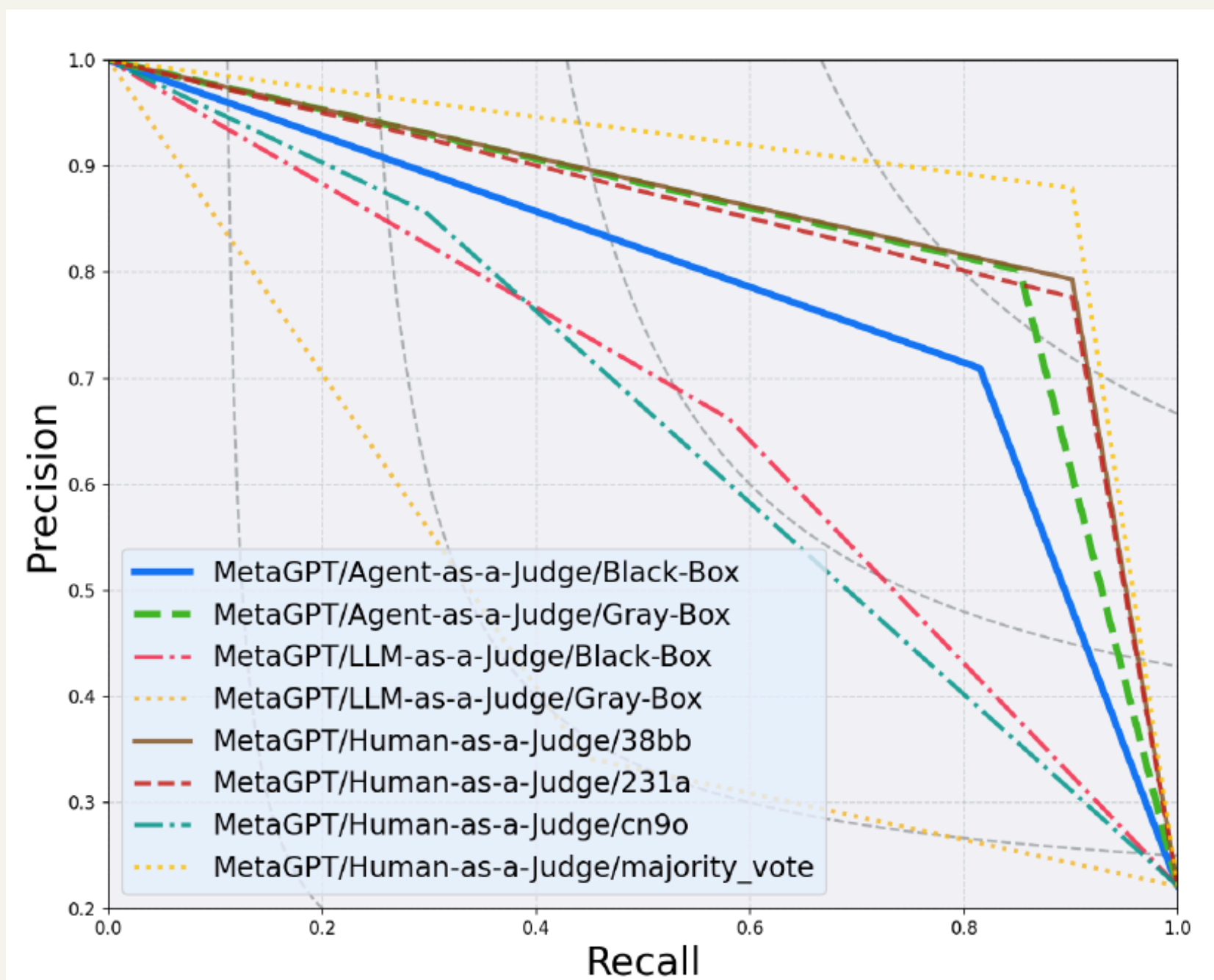**1** **Rich Intermediate Feedback:** Unlike traditional evaluation methods, Agent-as-a-Judge provides rich intermediate feedback that captures the entire thought and action trajectory, enhancing the understanding of agentic systems.

**2** **Evolved from LLM-as-a-Judge:** It builds upon the LLM-as-a-Judge framework but integrates agentic features for more detailed evaluations of agentic systems.

**3** **Cost-Effectiveness:** Demonstrates significant cost and time savings compared to human evaluators, making it a scalable alternative.



**LLM-as-a-Judge**     **Agent-as-a-Judge**     **Human-as-a-Judge**

Check out our LLM/GenAI Courses at www.masteringllm.com ✨

# Benchmarking Agent-as-a-Judge

**1** **Outperformed LLM-as-a-Judge:** Agent-as-a-Judge aligned with human evaluators' consensus at a rate of **90%, significantly outperforming LLM-as-a-Judge at 70%.**

**2** **Efficiency Gains:** Evaluations done using Agent-as-a-Judge cost only **2.28%** of the human evaluation baseline while taking just **2.36%** of the time.

**3** **Dynamic Feedback Integration:** The modular components allowed for dynamic evidence collection throughout task completion, leading to more accurate judgments.



Legend (Precision vs Recall):
- MetaGPT/Agent-as-a-Judge/Black-Box
- MetaGPT/Agent-as-a-Judge/Gray-Box
- MetaGPT/LLM-as-a-Judge/Black-Box
- MetaGPT/LLM-as-a-Judge/Gray-Box
- MetaGPT/Human-as-a-Judge/38bb
- MetaGPT/Human-as-a-Judge/231a
- MetaGPT/Human-as-a-Judge/cn9o
- MetaGPT/Human-as-a-Judge/majority_vote

# DevAI-Dataset for AI Agents

**1** **Real-World AI Development Tasks:** DevAI includes **55** real-world comprehensive tasks, each with detailed requirements, to evaluate agentic systems thoroughly.

**2** **Hierarchical Task Structure:** Tasks in DevAI are structured with dependencies in a directed acyclic graph, providing a non-sparse evaluation compared to existing benchmarks.

**3** **Evaluation of Leading Agentic Systems:** The dataset was used to benchmark three popular open-source agentic systems, providing a proof-of-concept for Agent-as-a-Judge.



(1) Word Clouds of User Queries

Shortest: 69 words    Longest: 164 words

(2) Number of Words in User Queries

(4) Mentions of Models

(3) Number of Tags of User Queries

Classification · Regression · Other · Recommender System · Unsupervised Learning · Time Series Forecasting · Medical Analysis · Reinforcement Learning · Audio Processing · Generative Models · Financial Analysis · Supervised Learning · Computer Vision · Natural Language Processing

SVM classifier, LSTM, Random Forest, Linear Regression, Decision Tree, Logistic Regression, ResNet-50, Siamese, SVM regressor, BART, KNN classifier, DenseNet-121, U-Net, BERT, Latent Factor model, K-means, PPO, Random Forest classifier, Transformer, 7B LLaMA, YOLOv3, ResNet-18, DCGAN, CNN-LSTM, Q-learning, XGBoost, Show and Tell, VGG16, GridSearchCV, FaceNet, NCF, SRCNN, Word2Vec, DQN, GPT-2, Naive Bayes classifier, qlora.
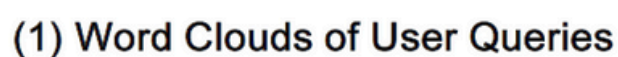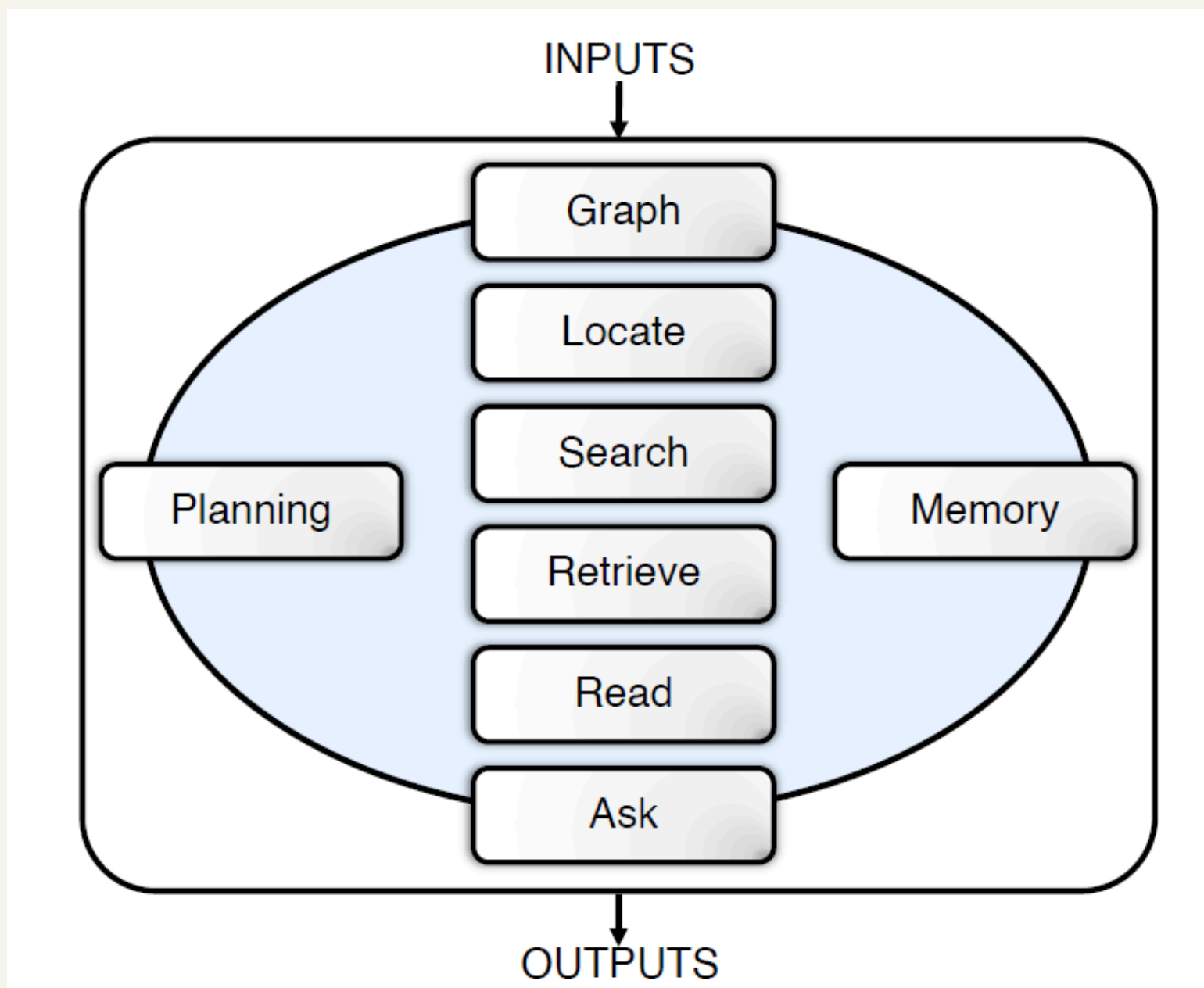
Check out our LLM/GenAI Courses at www.masteringllm.com ✨

# Human vs. Agent Evaluation

**1** **Comparable to Human Evaluators:** The Agent-as-a-Judge aligned with human evaluators at a similar level of accuracy, suggesting that it could replace human judges for certain tasks.

**2** **Reduction in Variability:** Unlike human judges who showed variability in their evaluations, Agent-as-a-Judge provided consistent and unbiased feedback across different tasks.

**3** **A Fraction of the Cost and Time:** The framework saved over 97% of both the cost and time needed compared to the human evaluators, highlighting its scalability and effectiveness.

INPUTS

Graph

Locate

Search

Planning

Memory

Retrieve

Read

Ask

OUTPUTS

Check out our LLM/GenAI Courses at www.masteringllm.com ✨

# LLM Interview Course 🌐

**50% OFF**

## Want to Prepare yourself for an LLM Interview?

✓ **120+** Questions spanning **14 categories** with Real Case Studies

✓ Curated **100+ assessments** for each category

✓ Well-researched **real-world interview questions based on FAANG & Fortune 500 companies**

✓ Focus on Visual learning

✓ Certification

**HOT SALE!**

## Coupon Code - LLM50

Coupon is valid till 30th Oct 2024

# AgenticRAG with LlamaIndex 🖥️

## Want to learn why AgenticRAG is future of RAG?

✓ Master **RAG fundamentals** through practical case studies

✓ Understand how to overcome **limitations of RAG**

✓ Introduction to **AgenticRAG** & techniques like **Routing Agents, Query planning agents, Structure planning agents, and React agents with human in loop**.

✓ **5** real-time **case studies with code walkthroughs**