

# LLM Agent Evaluations & Project Overview

Agentic AI

2025.10.6

# Outline

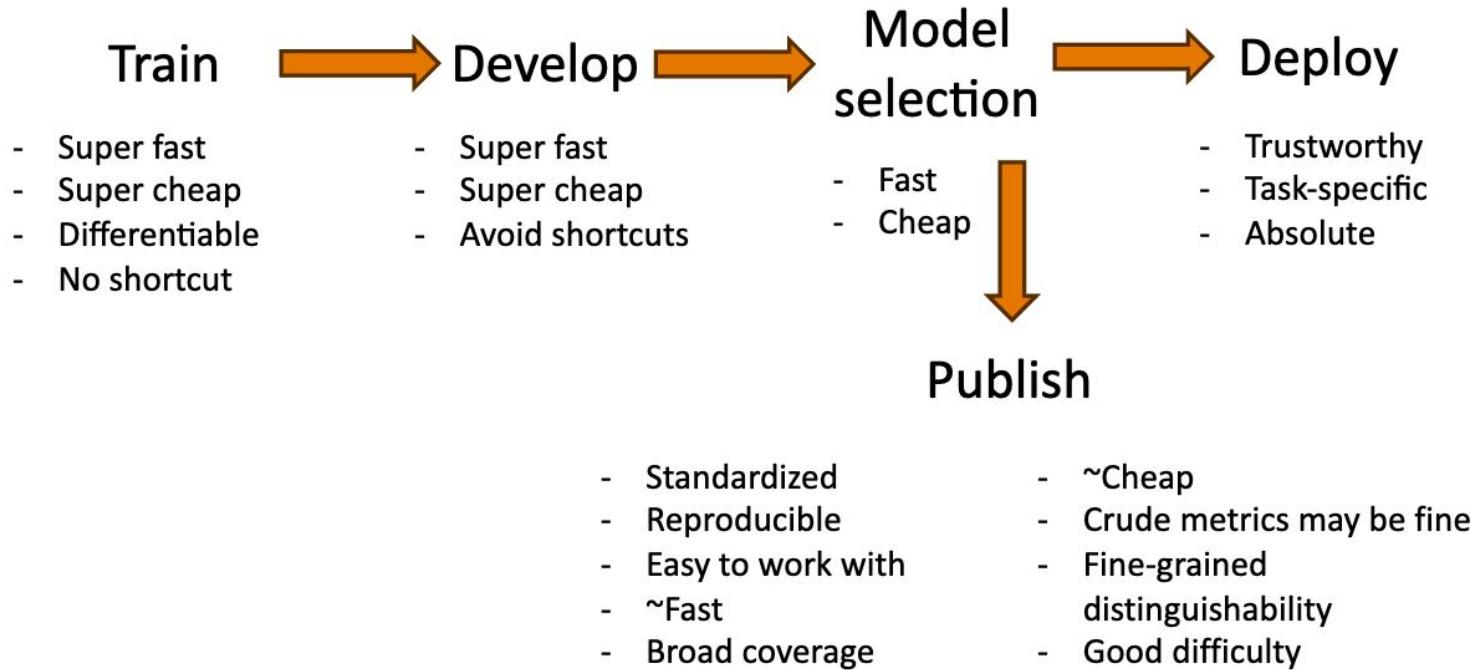
- Introduction to LLM agent evaluation and why is it important?
- Types of LLM Agent Eval
  - Close-ended vs. Open-ended
  - Verifiable vs. Non-verifiable
  - Static vs. Dynamic
  - Taxonomy of agent eval tasks
- What is a good eval system?
  - Outcome Validity
- Case studies of good eval systems
- Course project overview:
  - Two types of evals (green agents) in our course project
    - Integrating existing benchmarks
    - Proposing new benchmarks

# LLM Agent Eval: Why Is It Important?

# What is Evaluation?

- Evaluation is the systematic, repeatable measurement of models and agents.
- It provides a structured way to measure performance across benchmarks and environments.
- This helps
  - Measure capability progress that is grounded in reproducible evidence.
  - Risk assessment

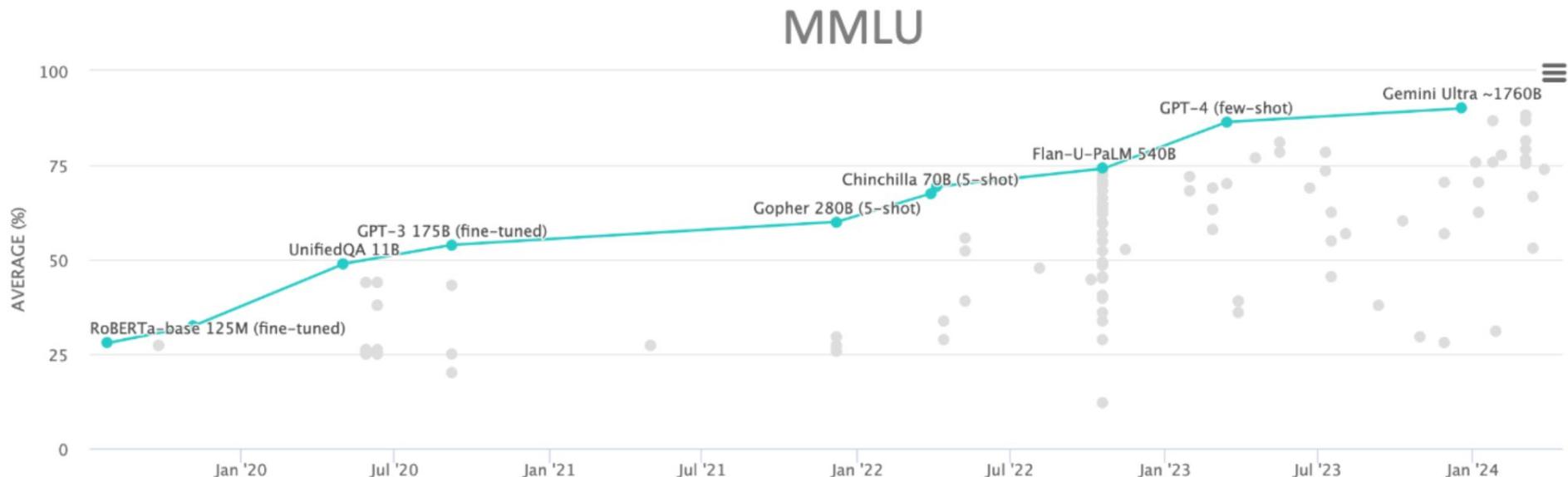
# When do we need evaluation?



# Why Evaluation Matters

- It enables fair comparisons across models and agents.
- Guides safe & effective deployment decisions by exposing weaknesses and strengths.
- Reliable evaluation of agents is critical to develop effective and safe agents in real-world applications.

# Benchmarks and Evaluation Drives Progress



# From LLM Eval to LLM Agent Eval

- LLMs are static, text-to-text systems.
- Agents extend them with planning, tool-use, memory, and multi-step reasoning.
- Agents operate in dynamic environments, requiring more complex evaluation.

# Types of LLM (Agent) Eval

# Types of Evals

- Close-ended
- Open-ended
  - Verifiable
  - Non-verifiable
- Static vs. Dynamic Eval

# Close-Ended vs. Open-Ended

## Close-ended evaluations

### Example

**Text:** Read the book, forget the movie!

**Label:** Negative

## Open ended evaluations

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

# Close-Ended Tasks

- Limited number of potential answers
- Limited number of correct answers
- Enables automatic evaluation

# Close-Ended Tasks

- Sentiment analysis: SST / IMDB / Yelp ...

## Example

Text: Read the book, forget the movie!

Label: Negative

- Entailment: SNLI

## Example

Text: A soccer game with multiple males playing.

Hypothesis: Some men are playing sport.

Label: Entailment

- Name entity recognition: CoNLL-2003
- Part-of-Speech: PTB

# Close-Ended Tasks

- Metrics: Accuracy / Precision / Recall / F1 / etc.
- Limited to certain easier tasks

# Open-Ended Tasks

- Long generations with too many possible correct answers to enumerate
  - => can't use standard ML metrics
- There are now better and worse answers (not just right and wrong)
- Example:
  - Summarization: CNN-DM / Gigaword
  - Translation: WMT
  - Instruction-following: Chatbot Arena / AlpacaEval / MT-Bench

# How to Evaluate Open-Ended Tasks?

- Two types of open-ended tasks:
  - Verifiable
  - Non-Verifiable

# Verifiable Tasks

- Tasks that have an “oracle” or clear criteria to test if the result is correct or not.
  - Math proof
  - Code generation
- Only need to build an eval that follows the criteria

# Non-Verifiable Tasks

- Tasks without clear test criteria or objective ground-truth answer
  - Storytelling
  - Writing style adaptation
- Often need human eval or LLM-as-a-judge

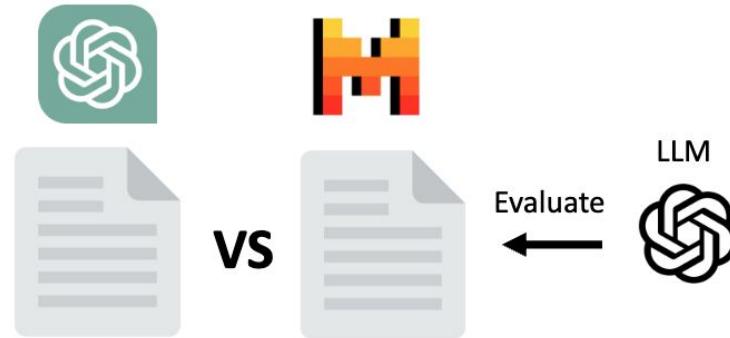
# Human Evaluations

- Ask humans to evaluate the quality of the generated text
- Overall or along some specific dimension:
  - Fluency
  - Coherence / consistency
  - Factuality and correctness
  - Commonsense
  - Style / formality
  - Grammaticality
  - redundancy

# Human Evaluations: Issues

- Human evaluations are regarded as gold standards
- But it also has issues:
  - Slow
  - Expensive
  - Inter-annotator disagreement
  - Intra-annotator disagreement across time
  - Not reproducible

# LLM as a Judge



- Use a LM as a reference free evaluator
- Surprisingly high correlations with human
- Common versions: AlpacaEval, MT-bench

# LLM as a Judge: Issues

- LLM sometimes is not reliable
  - Cross-check with human eval
- LLM output has randomness
  - Repeat judging with the same LLM for multiple times
- Different LLMs have different biases
  - Majority vote between different LLMs
- Interpretability of scores is poor
  - Try using continuous instead of discrete scoring
  - Score from multiple perspectives / set multiple rubrics
- Sensitive to vague prompts
  - Try design as detailed prompts as possible
  - Use chain-of-thought prompt or reasoning mode to get better evaluation

# Static vs. Dynamic Eval

- Static benchmarks have **fixed** test cases and metrics
  - Enable direct, reproducible comparisons between models
  - Examples: ImageNet, GLUE, MMLU
- Dynamic benchmarks have continuously updating or periodically re-generated data
  - Stay relevant with real-world data shifts
  - Harder to overfit or be contaminated
  - Examples: DynaBench, LiveCodeBench

# Taxonomy of LLM Agent Eval

# Taxonomy of Agent Evals

- Evaluating a specific agent capability
- Evaluating on a specific application
- Evaluating on a general set of applications

# Taxonomy of Agent Evals

- Evaluating a specific agent capability
- Evaluating on a specific application
- Evaluating on a general set of applications

# Evaluating a Specific Agent Capability

Planning and Multi-Step Reasoning ([§2.1](#))

*AQUA-RAT* ([Ling et al., 2017](#)); *HotpotQA* ([Yang et al., 2018](#)); *ARC* ([Clark et al., 2018a](#)); *StrategyQA* ([Geva et al., 2021](#)); *GSM8K* ([Cobbe et al., 2021](#)); *MATH* ([Hendrycks et al., 2021b](#)); *Game of 24* ([Yao et al., 2023](#)); *MINT* ([Wang et al., 2023](#)); *PlanBench* ([Valmeekam et al., 2023](#)); *FlowBench* ([Xiao et al., 2024](#)); *FOLIO* ([Han et al., 2022](#)); *P-FOLIO* ([Han et al., 2024](#)); *MultiRC* ([Khashabi et al., 2018](#)); *MUSR* ([Sprague et al., 2023](#)); *BBH* ([Suzgun et al., 2022](#)); *ToolEmu* ([Ruan et al., 2023](#)); *MINT* ([Wang et al., 2023](#)); *AutoPlanBench* ([Stein et al., 2023](#)); *ACP Bench* ([Kokel et al., 2024](#)); *Natural Plan* ([Zheng et al., 2024](#))

# Evaluating a Specific Agent Capability

Function Calling  
& Tool Use ([§2.2](#))

*BFCL* ([Yan et al., 2024](#)); *ToolBench* ([Qin et al., 2023](#)); *ToolAlpaca* ([Tang et al., 2023](#)); *APIBench* ([Patil et al., 2025](#)); *API-Bank* ([Li et al., 2023](#)); *NexusRaven* ([team, 2023](#)); *Seal-Tools* ([Wu et al., 2024b](#)); *ComplexFuncBench* ([Zhong et al., 2025](#)); *ToolSandbox* ([Lu et al., 2024](#)); *RestBench* ([Song et al., 2023](#)); *APIGen* ([Liu et al., 2024c](#)); *StableToolBench* ([Guo et al., 2024](#)); *NESTFUL* ([Basu et al., 2024b](#))

# Evaluating a Specific Agent Capability

Self-Reflection ([§2.3](#))

*LLF-Bench* ([Cheng et al., 2023](#)); *LLM-Evolve* ([You et al., 2024](#)); *Reflection-Bench* ([Li et al., 2024](#))

# Evaluating a Specific Agent Capability

Memory ([§2.4](#))

*NarrativeQA* ([Kočiský et al., 2018](#)); *QMSum* ([Zhong et al., 2021](#));  
*QUALITY* ([Pang et al., 2021](#)); *RAISE* ([Liu et al., 2024a](#)); *ReadAgent* ([Lee et al., 2024](#)); *MemGPT* ([Packer et al., 2024](#)); *LoCoMo* ([Maharana et al., 2024](#)); *A-MEM* ([Xu et al., 2025](#)); *StreamBench* ([Wu et al., 2024a](#)); *LTMbenchmark* ([Castillo-Bolado et al., 2024a](#))

# Taxonomy of Agent Evals

- Evaluating a specific agent capability
- Evaluating on a specific application
- Evaluating on a general set of applications

# Evaluating on a Specific Application

Web Agents (§3.1)

*MiniWob* ([Shi et al., 2017](#)); *MiniWoB++* ([Liu et al., 2018](#)); *WebShop* ([Yao et al., 2022](#)); *Mind2web* ([Deng et al., 2023](#)); *WebVoyager* ([He et al., 2024](#)); *WebLinX* ([Lù et al., 2024](#)); *WebArena* ([Zhou et al., 2023](#)); *VisualWebArena* ([Koh et al., 2024](#)); *MMInA* ([Zhang et al., 2024](#)); *AssistantBench* ([Yoran et al., 2024](#)); *WebCanvas* ([Pan et al., 2024b](#)); *ST-WebAgentBench* ([Levy et al., 2024](#)); *WorkArena* ([Drouin et al., 2024](#)); *WorkArena++* ([Boisvert et al., 2025](#))

# Evaluating on a Specific Application

Software Engineering  
Agents ([§3.2](#))

*HumanEval* ([Chen et al., 2021b](#)); *SWE-bench* ([Jimenez et al., 2023](#));  
*SWE-bench Verified* ([OpenAI, 2024](#)); *SWE-bench Lite* ([SWE-bench Lite, 2024](#)); *SWE-bench+* ([Aleithan et al., 2024](#)); *SWE-bench Multimodal* ([Yang et al., 2024](#)); *TDD-Bench Verified* ([Ahmed et al., 2024](#)); *SWT-Bench* ([Mündler et al., 2024](#)); *IT-Bench* ([Jha et al., 2025](#)); *SWE Lancer* ([Miserendino et al., 2025](#))

# Evaluating on a Specific Application

Scientific Agents ([§3.3](#))

*ScienceQA* ([Lu et al., 2022](#)); *QASPER* ([Dasigi et al., 2021](#)); *MS<sup>2</sup>* ([DeYoung et al., 2021](#)); *ScienceWorld* ([Wang et al., 2022a](#)); *SUPER* ([Bogin et al., 2024](#)); *Ideation* ([Si et al., 2025](#)); *AAAR-1.0* ([Lou et al., 2025](#)); *ScienceAgentBench* ([Chen et al., 2024](#)); *CORE-Bench* ([Siegel et al., 2024](#)); *SciCode* ([Tian et al., 2024b](#)); *MLGym-Bench* ([Nathani et al., 2025](#)); *DiscoveryWorld* ([Jansen et al., 2024](#)); *LAB-Bench* ([Laurent et al., 2024](#))

# Evaluating on a Specific Application

Conversational Agents ([§3.4](#))

*ABCD* ([Chen et al., 2021a](#)); *MultiWOZ* ([Budzianowski et al., 2018](#)); *SMCalFlow* ([Andreas et al., 2020](#)); *ALMITA* ([Arcadinho et al., 2024](#));  $\tau$ -*Bench* ([Yao et al., 2024](#)); *IntellAgent* ([Levi and Kadar, 2025a](#)); *LTM* ([Castillo-Bolado et al., 2024b](#))

# Evaluating on a Specific Application

- Other applications
  - Safety
    - RedCode (<https://arxiv.org/pdf/2411.07781.pdf>)
  - Cybersecurity
    - CyberGym (<https://www.cybergym.io/>)
    - BountyBench (<https://arxiv.org/abs/2505.15216>)
  - Legal
    - LegalAgentBench (<https://arxiv.org/abs/2412.17259>)
  - Healthcare
    - MedAgentBench (<https://www.qeios.com/read/VN3YH7/pdf>)
    - HealthBench (<https://openai.com/index/healthbench/>)
  - Finance
    - Finance Agent Benchmark (<https://www.arxiv.org/pdf/2508.00828.pdf>)

# Taxonomy of Agent Evals

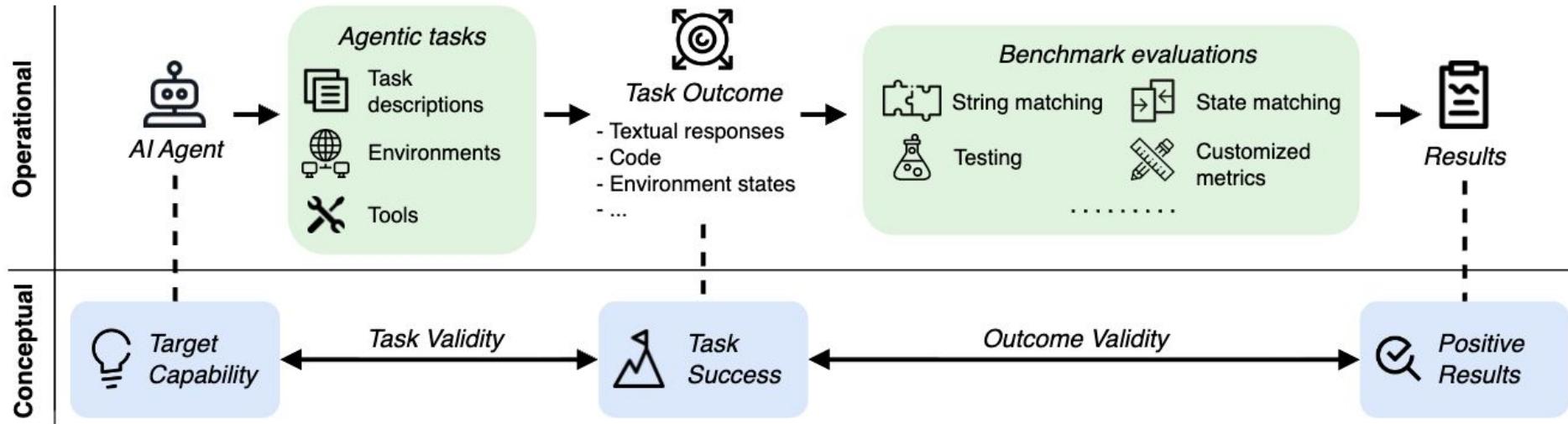
- Evaluating a specific agent capability
- Evaluating on a specific application
- Evaluating on a general set of applications

# Evaluating on a General Set of Applications

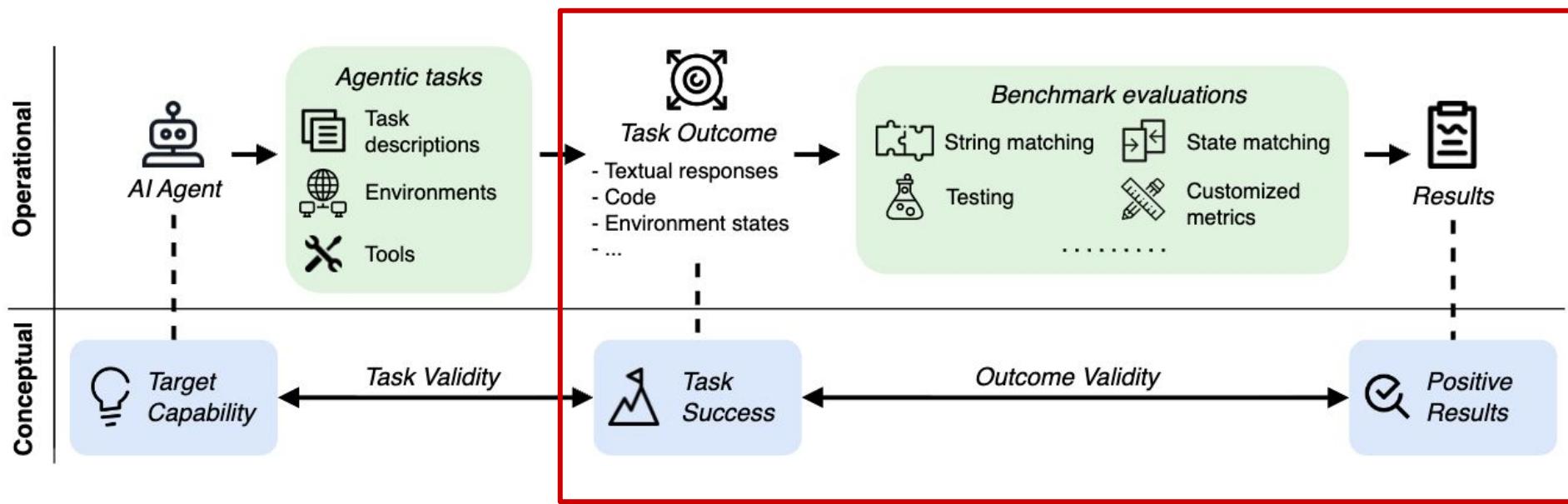


What is a good eval system?

# What is a good eval?



# Outcome Validity Makes a Good Eval



# Outcome Validity - Judging text results

## *Information Acquisition*

### **Whole string matching or substring matching:**

- O.a.1. Considers expressions semantically equivalent to ground truth.
- O.a.2. Handles redundant words used by agents.

### **Substring matching:**

- O.b.1. Handles negation modifiers used by agents.
- O.b.2. Is robust against systematically listing all possible answers.
- O.b.3. Ground truth is sufficiently complex to prevent guessing.

### **LLM-as-a-Judge:**

- O.c.1. Demonstrates documented or experimental evidence of the judge's accuracy, self-consistency, and agreement with human.
- O.c.2. Is designed to resist adversarial inputs and reward hacking.

# Outcome Validity - Judging Code Generation

## Code Generation

### **Unit testing or end-to-end testing:**

- O.d.1. Verifies test cases for correctness and quality (e.g., by human).
- O.d.2. Measures quality of test cases using objective metrics (e.g., code coverage, cyclomatic complexity control).

### **Fuzz testing:**

- O.e.1. Addresses potential edge cases.
- O.e.2. Ensures comprehensive coverage of all relevant input variations (e.g., data types, memory layouts, value ranges).
- O.e.3. Generates inputs that the code under testing is sensitive to.

### **End-to-end testing:**

- O.f.1. Exercises all relevant parts of the code being tested.
- O.f.2. Prevents non-deterministic (“flaky”) test results.

# Outcome Validity - Judging Env State Changes

## State Matching

### **State matching:**

- O.g.1. Ground truth includes all states achievable after success.
- O.g.2. Checks relevant and irrelevant states for the challenge.
- O.g.3. Ground truth is complex to prevent trivial state modifications.

# Outcome Validity - Judging Multi-Step Reasoning

## Multistep Reasoning

### **Answer matching:**

- O.h.1. Specifies required answer formats in challenge descriptions.
- O.h.2. Minimizes the possibility of success by random guessing.

### **Quality measure:**

- O.l.1. Designs quality metrics that prevent exploitation (e.g., achieving high scores by reward hacking).

# Ways That Eval Can Go Wrong

- Data is noisy or biased
  - Make sure the test data for evaluation is accurate and diverse enough!
- Not practical
  - Think about the practitioner's real needs!
- Shortcut - Eval can be gamed
  - Avoid any shortcut that your eval probably has!
- Not challenging enough
  - Design hard test cases to make sure your green agent is reliable!
- More info: <https://arxiv.org/pdf/2502.06559v2>

# Case Study of Good Eval System

# Case Studies

- What is the goal / what to evaluate
- What is a task, what is an env to run the agent to achieve the goal, what's the size of the benchmark
- How to evaluate the agent (LLM Judge, etc.)
  - Need to generate data in the benchmark, talk about data pipeline
  - What metric to use to make benchmark high quality (contamination, bias)
- Principles: real-world, have different difficulty levels, not easy to get saturated. GDPEval is a good example

# Case Studies

- What is a good benchmark and what's the process of constructing them (data pipeline, to let students have some ideas on what they need to do)
  - CyberGym (anthropic model card)
  - Tau-bench and tau2-bench
  - OS World
  - GDPEval
  - CRMArena
  - LegalAgentBench

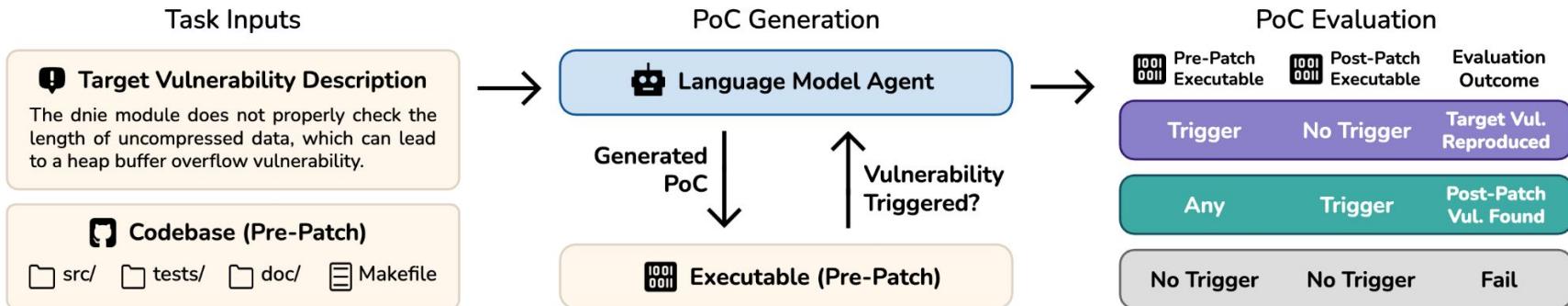
# Case Studies

What is a good benchmark and how to construct it?

- What is the goal / what to evaluate
- What is a task / what is an env to run the agent to achieve the goal
- How to build the data collection pipeline? How to evaluate the agent?
- **Principles:** real-world, have different difficulty levels, not easy to get contaminated and saturated

# CyberGym

<https://www.cybergym.io/>



**Goal:** Evaluate an agent's cybersecurity capabilities by testing its ability to reproduce real-world vulnerabilities at a large scale

**Task:** Given a vulnerability description and the pre-patch codebase+executable, agents must generate a proof-of-concept (PoC) test that successfully triggers the vulnerability in the corresponding unpatched codebase

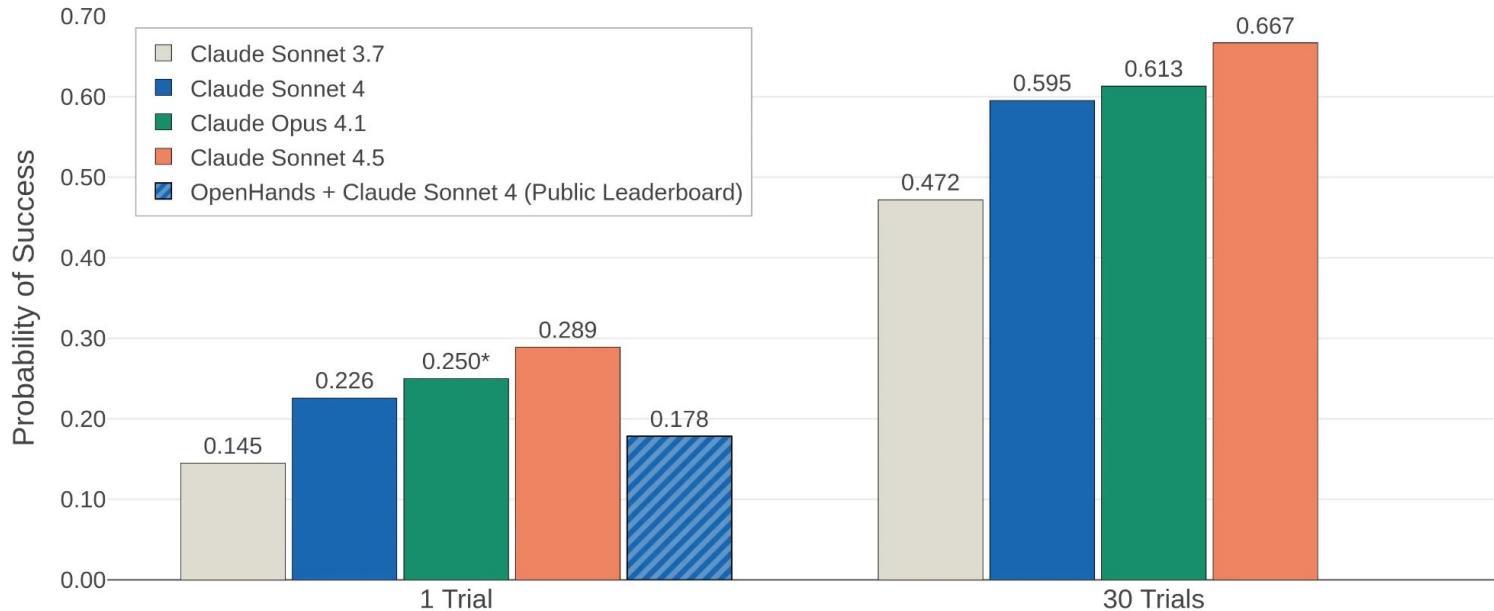
**Env:** a containerized sandbox to run programs

# CyberGym

<https://www.cybergym.io/>

Anthropic's latest [system card for its model release](#) (Claude 4.5) included CyberGym to evaluate AI capabilities in cybersecurity.

Model Performance Comparison on Vulnerability Reproduction



# CyberGym

<https://www.cybergym.io/>

## Data Generation Pipeline:

- Built from ARVO dataset and historical, real-world vulnerabilities found by OSS-Fuzz, a continuous fuzzing project for open-source software
- reconstruct pre/post patch commits & executables and include the ground-truth PoC; rephrase into concise vuln descriptions with LLMs and manual inspection

## How to Evaluate:

- Execute final PoC on pre-patch and post-patch builds. Count success if it (a) triggers the target vuln only **pre-patch (reproduction)**, or (b) triggers any vuln **post-patch (post-patch finding)**. Report overall success rate
- Detection is via runtime sanitizers (crash + stack trace), not subjective judging.
- A **data contamination analysis** is performed by evaluating vuln samples found after LLM knowledge cutoff dates

# T-bench

<https://arxiv.org/abs/2406.12045>

**Goal:** Evaluate an agent's ability to reliably interact with users and APIs while consistently following complex, domain-specific policies

**Task:** Agents resolve a simulated user's goal (e.g., return a product) using API tools through a multi-turn, dynamic conversation within domains like retail or airline customer service

**Env:** Each domain (e.g., retail, airline) provides a set of API tools, a specific policy document to follow, and an LLM-powered user simulator

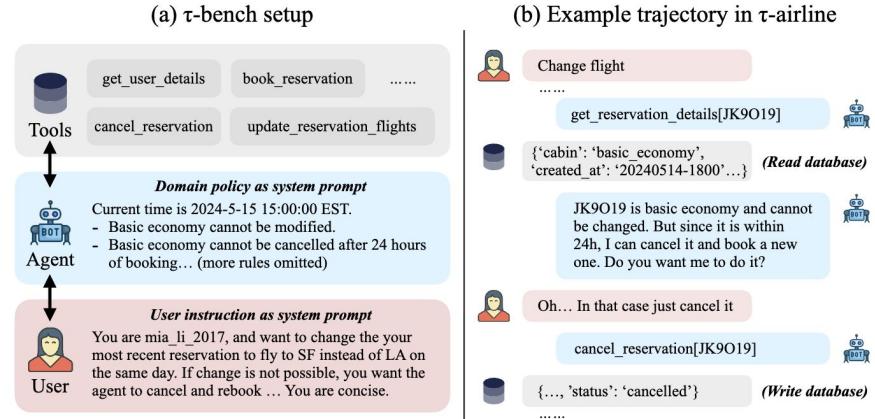


Figure 1: (a) In  $\tau$ -bench, an agent interacts with database API tools and an **LM-simulated user** to complete tasks. The benchmark tests an agent's ability to collate and convey all required information from/to users through multiple interactions, and solve complex issues on the fly while ensuring it **follows guidelines** laid out in a domain-specific policy document. (b) An example trajectory in  $\tau$ -airline, where an agent needs to reject the user request (change a basic economy flight) following domain policies and propose a new solution (cancel and rebook). This challenges the agent in long-context zero-shot reasoning over complex databases, rules, and user intents.

# T-bench

<https://arxiv.org/abs/2406.12045>

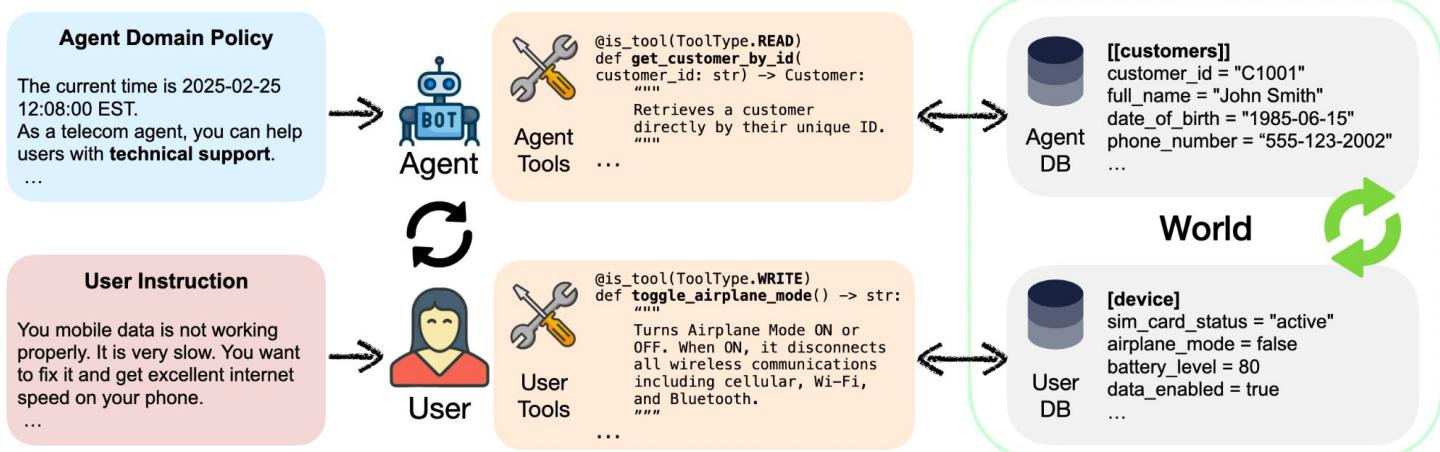
## Data Generation Pipeline:

- Manual design of schemas/APIs/policies
- LM-assisted synthetic data generation (GPT-4 helps produce sampling code; humans polish)
- Manual scenario authorizing + iterative validation with many agent runs to ensure each task has a unique end-state outcome

**How to Evaluate:** Evaluation is programmatic and verifiable. Success is determined by comparing the final database state to the annotated goal state. Report **pass@1 (avg success)** and **pass@k (all-k successes across i.i.d. runs)** to capture reliability/consistency

# T<sup>2</sup>-bench

<https://arxiv.org/abs/2506.07982>



**Goal:** T<sup>2</sup> shifts from single-control to dual-control (Dec-POMDP)—both agent and user act via tools in a shared world stressing coordination & guidance

## Task and Env:

- T was single DB + agent tools, with an LM-only user
- T<sup>2</sup> adds two databases (Agent DB + User/Device DB) and separate toolsets; the user is a simulator constrained by available tools and observable state of the environment

# $T^2$ -bench

<https://arxiv.org/abs/2506.07982>

## Data Generation Pipeline:

- T used manual schema/APIs, LM-assisted data, manual scenario authoring/validation
- $T^2$  pipeline uses LLM-drafted Product Requirements Document (PRD) → code/mock DBs/unit tests, plus user DB & tools, then do programmatic compositional tasks creation from atomic subtasks with init/sol/assert and auto-verification

**How to Evaluate:** T evaluates via end-state DB comparison.  $T^2$  introduces categorical checks—environment assertions, communication assertions, natural language assertions, action assertions; both report pass@k

# GDPval

<https://openai.com/index/gdpval/>

## **Manufacturing Engineer:** Design 3D model of cable reel stand for assembly line

### Prompt + task context:

#### **Experienced human deliverable:**



**Financial and Investment Analyst:** Create competitor landscape for last mile delivery

### Prompt + task context

#### **Experienced human deliverable**

It is April 2012 and you are a Managing Director of an investment firm focusing on the e-commerce – business-to-business – B2B – role of logistics service. One of your clients, a company called "Logistics Solutions", has approached you asking for funding to invest in a new warehouse. This warehouse will be used to store raw materials and to package and distribute them to its clients. You have been asked to evaluate the financial viability of this project and have advised that for your opinion as a director of private companies there and how the public market would value these companies.



**Registered Nurse:** Assess skin lesion images and create consultation report

### Prompt + task context

#### **Experienced human deliverable:**

<p><b>GENERAL INFORMATION</b></p> <p><b>NAME:</b> [REDACTED]  <b>ADDRESS:</b> [REDACTED]</p> <p><b>TELEPHONE NUMBER:</b> [REDACTED]</p> <p><b>EMAIL ADDRESS:</b> [REDACTED]</p> <p><b>DATE OF BIRTH:</b> [REDACTED]</p> <p><b>SEX:</b> [REDACTED]</p> <p><b>RACE:</b> [REDACTED]</p> <p><b>ETHNICITY:</b> [REDACTED]</p> <p><b>RELIGION:</b> [REDACTED]</p> <p><b>EDUCATION LEVEL:</b> [REDACTED]</p> <p><b>EMPLOYMENT STATUS:</b> [REDACTED]</p> <p><b>HOUSING STATUS:</b> [REDACTED]</p> <p><b>DISABILITY STATUS:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p>	<p><b>GENERAL INFORMATION</b></p> <p><b>NAME:</b> [REDACTED]  <b>ADDRESS:</b> [REDACTED]</p> <p><b>TELEPHONE NUMBER:</b> [REDACTED]</p> <p><b>EMAIL ADDRESS:</b> [REDACTED]</p> <p><b>DATE OF BIRTH:</b> [REDACTED]</p> <p><b>SEX:</b> [REDACTED]</p> <p><b>RACE:</b> [REDACTED]</p> <p><b>ETHNICITY:</b> [REDACTED]</p> <p><b>RELIGION:</b> [REDACTED]</p> <p><b>EDUCATION LEVEL:</b> [REDACTED]</p> <p><b>EMPLOYMENT STATUS:</b> [REDACTED]</p> <p><b>HOUSING STATUS:</b> [REDACTED]</p> <p><b>DISABILITY STATUS:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p>	<p><b>GENERAL INFORMATION</b></p> <p><b>NAME:</b> [REDACTED]  <b>ADDRESS:</b> [REDACTED]</p> <p><b>TELEPHONE NUMBER:</b> [REDACTED]</p> <p><b>EMAIL ADDRESS:</b> [REDACTED]</p> <p><b>DATE OF BIRTH:</b> [REDACTED]</p> <p><b>SEX:</b> [REDACTED]</p> <p><b>RACE:</b> [REDACTED]</p> <p><b>ETHNICITY:</b> [REDACTED]</p> <p><b>RELIGION:</b> [REDACTED]</p> <p><b>EDUCATION LEVEL:</b> [REDACTED]</p> <p><b>EMPLOYMENT STATUS:</b> [REDACTED]</p> <p><b>HOUSING STATUS:</b> [REDACTED]</p> <p><b>DISABILITY STATUS:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p> <p><b>CRIMINAL RECORD:</b> [REDACTED]</p> <p><b>DRIVING RECORD:</b> [REDACTED]</p> <p><b>CRIMINAL HISTORY:</b> [REDACTED]</p> <p><b>DRUG USE HISTORY:</b> [REDACTED]</p> <p><b>ALCOHOL USE HISTORY:</b> [REDACTED]</p>
---	---	--

**Goal:** Measure LLM performance on economically valuable, real-world knowledge-work tasks, comparing AI deliverables to industry experts across diverse occupations

**Task and Env:** Each task is a realistic work assignment with reference files/context (docs, data, assets). Models produce a one-shot deliverable (e.g., doc, slide deck, spreadsheet, diagram, media)

# GDPval

<https://openai.com/index/gdpval/>

**Data Generation Pipeline:** Tasks authored by vetted professionals (avg 14 yrs experience), pass a multi-step review ( $\approx$ 5 rounds) plus model-based validation; prompts mirror day-to-day work and include attachments; gold deliverables are experts' own solutions

## How to Evaluate:

- Blinded expert graders from the same occupations rank AI vs. human deliverables as better / as good as / worse
- Also compare time/cost
- Good example of a **benchmark with low contamination risk and hard to get saturated** as tasks require domain experts and tied to concrete work product

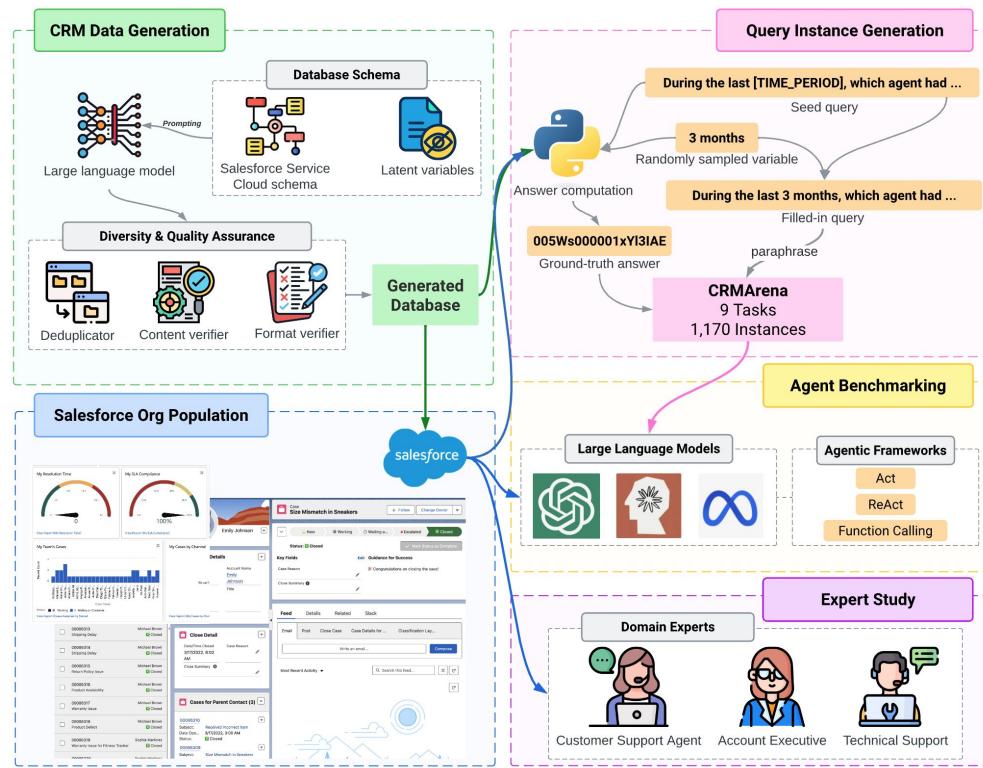
# CRM Arena

<https://arxiv.org/abs/2411.02305>

**Goal:** Evaluate LLM agents on professional Customer Relationship Management (CRM) workflows in a realistic, enterprise sandbox

**Task:** 9 tasks across 3 personas (Service Agent, Analyst, Manager): New Case Routing, Knowledge QA, Top Issue Identification, Monthly Trend Analysis etc.

**Env:** Live Salesforce sandbox (Simple Demo Org) with UI & API access; actions via SOQL/SOSL or function calls; Rich enterprise schema (16 objects)



# CRM Arena

<https://arxiv.org/abs/2411.02305>

## Data Generation Pipeline:

- LLM synthesis on Salesforce Service Cloud schema; introduce **latent variables** (e.g., agent Skill, customer ShoppingHabit) to create realistic causal patterns.  
[OBJ]
- Mini-batch prompting → de-duplication (string match) + dual verification (format & content) before upload; LLM paraphrasing for query diversity

## How to Evaluate:

- Automatic metrics per task: F1 for Knowledge QA; Exact Match on ground-truth IDs for all other tasks; optional pass@k to report multi-run reliability/consistency
- Also reports #turns/tokens/\$ cost

# LegalAgentBench

<https://arxiv.org/abs/2412.17259>

**Goal:** Evaluate LLM agents on realistic Chinese legal workflows that require multi-hop retrieval/reasoning and legal writing with tool use.

**Task and Env:** Multi-turn agent interaction over a text-based environment built from 17 corpora (14 tabular DBs + 3 retrieval corpora) with 37 tools (28 DB tools, 5 math, 3 text retrievers, 1 system “Finish”). Agents must plan, call tools with parameters, read observations, and iterate.

- E.g. “Which court handles the highest-value cases of <entity>?” (requires chaining court/case DB tools).

# LegalAgentBench

<https://arxiv.org/abs/2412.17259>

**Data Generation Pipeline:** (1) Build a planning tree from tool-call dependencies; (2) select serial (1–5 hop) and parallel paths; (3) choose entities that execute successfully; (4) rewrite questions with GPT-4 to hide solution paths and match real phrasing; (5) programmatically generate answers via the toolchain; (6) human verification of questions, paths, and answers.

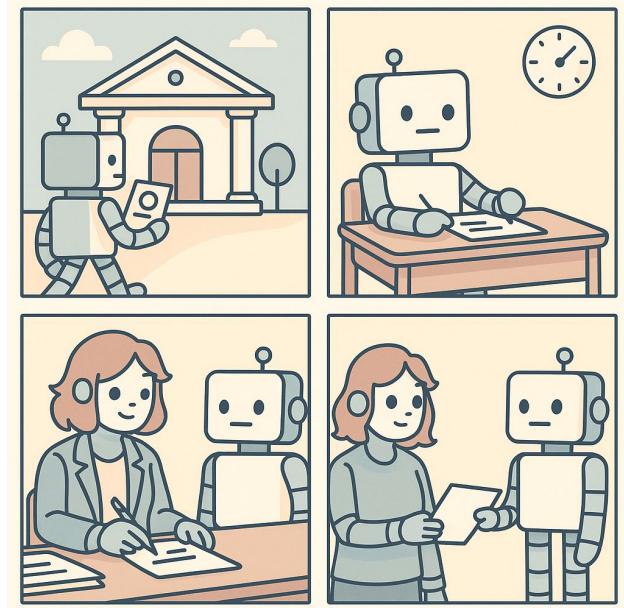
## How to Evaluate:

- **Success Rate:** keyword overlap with gold key\_answer.
- **Process/Progress Rate:** includes intermediate key\_middle keywords to reward partial progress.
- **Other Metrics:** BERTScore and token count for efficiency.

# What are Green Agents?

# Testing my Agent

- Like math exams are designed by math teachers and physics exams by physics teachers, Different assessments are defined by a special hosting agent, or **green agent**.
- This green agent determines the **type of assessment** and **defines the specific tasks** to be performed.
- The other agents involved are called participant agents, or **white agents**.



# Testing my Agent

\* Why “Green” and “White”?

- The color-based naming comes from the project’s early days. When exploring safety benchmarks, the team discussed red-teaming and blue-teaming setups. To refer to judges and other participants, we temporarily introduced more colors—green and gray. Over time, “green agent” became the established name for the hosting agent, while all others were simplified into “white agents.”

# Responsibilities of the Green Agent

The green agent carries multiple responsibilities:

- Preparing the environment
- Distributing test tasks to participant agents
- Collecting their results
- Verifying the environment
- Reporting back to the platform

# Responsibilities of the Green Agent

- Green agents are specifically designed to serve as evaluators
- They essentially act as the **technicians at the repair shop**, orchestrating the entire evaluation process.
- Green agents can interact with the platform through **MCP, A2A or APIs**, request permissions and resources, and submit results.

# The Full Assessment Flow

1. The platform first confirms that all agents, including the green agent, are online and reset.
2. It then sends a task to the green agent, including the URLs of the participant agents to be tested.
3. The green agent orchestrates the interaction: assigning tasks, supervising execution, managing tools or environments, and continuously reporting updates back to the platform.
4. At the end, the green agent submits the metrics—which it defines through its implementation.

**The green agent dictates what metrics get measured**

# How Green Agents Are Built

Typically, a green agent includes:

- A dataset of test tasks
- A predefined testing process (e.g., which agent to test first, in what order tasks are sent, and how tools are provided)
- The environment where the tasks run: access to additional tools or MCP modules required for the tasks, along with instructions for the white agents on how to use them

# How Green Agents Are Built

- AgentBeats provides a **prompt-based toolkit** to help developers quickly spin up a prototype green agent, making it easy to get started.
- For more rigorous or complex testing—say, strict workflows or custom environments—developers can also hand-code the logic of a green agent.
  - For class projects, we expect most green agents from unit 2 & 3 & 4 groups will be coding-based green agents; and many green agents from unit 1 groups may be prompt-based agents
- Ultimately, as long as it complies with the **A2A protocol, any web service can serve as a green agent.**

# Two Types of Projects for Building Green Agents

- Integrating an existing benchmark
- Building a new benchmark

# Type 1: Integrating Existing Benchmarks

- Goal: Adapt an existing benchmark (already published/tested) and integrate as a green agent in AgentBeats
  - E.g. SWE-bench Verified, Terminal bench
- Largely reuse existing evaluation metrics or rubrics
- Sample ideas:  
[https://docs.google.com/presentation/d/1TjtEjh6g9dZBsGxmAmcSp2EFakbmHpU\\_z31vnkf0c2Y/](https://docs.google.com/presentation/d/1TjtEjh6g9dZBsGxmAmcSp2EFakbmHpU_z31vnkf0c2Y/)

# Type 1: Integrating Existing Benchmarks

In addition to integration into agentbeats platform, you need to provide a **benchmark quality analysis** – checking the quality and reliability of the existing benchmark.

- **Manual Validation:** Sample and check data correctness, clarity, and difficulty
- **Evaluator Check:** Confirm metrics/judges align with true task success
- **Bias & Limitation Notes:** Highlight any gaps or weaknesses

Also consider expanding the benchmark and correcting the benchmark to improve its quality, size, and diversity.

# Type 1: Integrating Existing Benchmarks

Main Workflows:

- **Step 1: Integration**
  - Convert problem formats to correct format like A2A
  - Implement dataset loaders & interfaces
  - Add quality checks for correctness & reproducibility
- Step 2: Benchmark Quality Analysis
- Step 3: Correction and Expansion

# Type 1: Integrating Existing Benchmarks

## Main Workflows:

- Step 1: Integration
- **Step 2: Benchmark Quality Analysis:** check the quality and reliability of the existing benchmark.
  - **Manual Validation:** Sample and check data correctness, clarity, and difficulty
  - **Evaluator Check:** Confirm metrics/judges align with true task success
  - **Bias & Limitation Notes:** Highlight any gaps or weaknesses
- Step 3: Correction and Expansion

# Type 1: Integrating Existing Benchmarks

Main Workflows:

- Step 1: Integration
- Step 2: Benchmark Quality Analysis
- **Step 3: Correction and Expansion**
  - Correct the benchmark if there are errors
  - Expand the benchmark to improve its quality, size, and diversity.

# SWE-bench and SWE-bench Verified

<https://openai.com/index/introducing-swe-bench-verified/>

- **Problem (Original SWE-bench):**
  - Some tasks had **underspecified issue descriptions** or **overly specific/misaligned tests**; setup friction sometimes caused **false negatives**.
- **Correction:**
  - Added **human verification** by 93 professional developers on 1,699 samples
    - Issues flagged: **38.3% underspecification, 61.1% unfair unit tests**; total **68.3% of samples filtered out**
  - Filtered to **500 verified tasks**
- **Outcome:**
  - Curated a **higher-quality subset** with enhanced **task diversity and difficulty balance**
  - More **trustworthy, replicable, and comprehensive** benchmark
    - GPT-4o reaches 33.2% resolved on Verified (vs. 16% on original using best scaffold), indicating prior underestimation of capability.

# Type 2: Building New Benchmarks

- Create new benchmarks (no existing source)
- Realistic daily tasks → showcase agentic reasoning
- Scope by units:
  - 1-unit students: focus on prompt design + simple data generation
  - >1 units: include evaluation rubrics & multi-round setups

# Type 2: Building New Benchmarks

- Tasks should reflect useful, real-world scenarios
  - e.g., organize calendar, schedule meetings, manage to-dos
- Evaluation: Automatic or lightweight human checks
- We encourage you to build **multi-agent** benchmarks (e.g., Synthesizer + Analyzer roles)

# Step-by-Step Checklist for Building Your Green Agent

# Step-by-Step Checklist

1. Choose the task you want to evaluate on
  - E.g., Ticket-booking agent

# Step-by-Step Checklist

2. Design the environment that the agents being tested needs to run in

- The tools that the agent can interact with, the actions that the agent can make, and the env feedback to the agent after each action
- E.g., Tools can be web browser or an APP for ticket booking. Actions can be mouse clicking and keyboard typing, or the APIs provided by the APP. Env feedback can be the new webpage popped up every time the agent clicks on a button.

# Step-by-Step Checklist

3. Design the metrics that your green agent evaluates with

- E.g., the success rate of booking a ticket; how cheap the ticket is; whether the ticket satisfies user's requirements; etc.

# Step-by-Step Checklist

## 4. Design test cases to evaluate your green agent

- Think about different scenarios of white agents trying to complete the task
- Design test cases of white agents succeeding/failing to complete the task in different ways, along with ground-truth eval result for these cases.
- Include as many edge cases as possible
- Use these test cases to evaluate if your green agent gives reliable evaluation results.
- E.g., test cases can include a white agent successfully books the ticket; a white agent books the wrong ticket/a more expensive ticket; a white agent fails to find the website for booking tickets; etc.

# NeurIPS 2025 Datasets & Benchmarks Track Call for Papers

The **NeurIPS Datasets and Benchmarks track** serves as a venue for high-quality publications on highly valuable machine learning datasets and benchmarks crucial for the development and continuous improvement of machine learning methods. Previous editions of the Datasets and Benchmarks track were highly successful and continuously growing (accepted papers [2021](#), [2002](#), and [2023](#), and best paper awards [2021](#), [2022](#), [2023](#) and [2024](#)). Read our [original blog post](#) for more about why we started this track, and the 2025 [blog post](#) announcing this year's track updates.

## Dates and Guidelines

Please note that the Call for Papers of the NeurIPS2025 Datasets & Benchmarks Track this year **will follow the [Call for Papers of the NeurIPS2025 Main Track](#), with the addition of three track-specific points:**

- Single-blind submissions
- Required dataset and benchmark code submission
- Specific scope for datasets and benchmarks paper submission

The dates are also identical to the main track:

- **Abstract submission deadline:** May 11, 2025 AoE
- **Full paper submission deadline:** May 15, 2025 AoE (all authors must have an OpenReview profile when submitting)
- **Technical appendices and supplemental materials deadline:** May 22, 2025 AoE
- **Author notification:** Sep 18, 2025 AoE
- **Camera-ready:** Oct 23, 2025 AoE

Accepted papers will be published in the NeurIPS proceedings and presented at the conference alongside the main track papers. As such, we aim for an equally stringent review as in the main conference track, while also allowing for **track-specific guidelines**, which we introduce below. For details on everything else, e.g. formatting, code of conduct, ethics review, important dates, and any other submission related topics, please refer to the [main track CFP](#).

## OpenReview

Submit at: [https://openreview.net/group?id=NeurIPS.cc/2025/Datasets\\_and\\_Benchmarks\\_Track](https://openreview.net/group?id=NeurIPS.cc/2025/Datasets_and_Benchmarks_Track)

The site will start accepting submissions on April 3, 2025 (at the same time as the main track).

**Note:** submissions meant for the main track should be submitted to a different OpenReview portal, as shown [here](#). Papers will not be transferred between the main and the Datasets and Benchmarks tracks after the submission is closed.

# Project Grading Rubric [for new benchmarks]

- **Goal & Novelty:** Is your benchmark important, novel, and covering new capability space?
- **Scope & Scale:** Is the benchmark large and diverse enough to give reliable results?
- **Evaluator Quality:** Are metrics clear? Is your judge/evaluator high quality and consistent?
- **Validation:** Did you perform manual checks or spot validation on the evaluation outputs from your green agent?
- **Reliability:** Do your evaluation scripts and green agents run robustly on AgentBeats?
- **Quality Assurance:** Any bias or contamination checks included?
- **Realism:** Is the benchmark realistic, e.g., with real world workload, instead of toy or unrealistic settings
- **Impact:** Is the benchmark reusable, well-documented, and presented clearly?

# Project Grading Rubric [for existing benchmarks]

- **Analysis:** Analyze quality issues of the original benchmark and find any flaws it has.
- **Faithfulness:** Is your implementation reproducing the results from the original benchmark (excluding the flaws you fixed)?
- **Quality Assurance:** Is your implementation correcting flaws in the original benchmark and expanding the coverage of the original benchmark?
- **Evaluator Quality:** Are metrics clear? Is your judge/evaluator high quality and consistent?
- **Validation:** Did you perform manual checks or spot validation on the evaluation outputs from your green agent?
- **Reliability:** Do your evaluation scripts and green agents run robustly on AgentBeats?
- **Quality Assurance:** Any bias or contamination checks included?
- **Impact:** Is your implementation reusable, well-documented, and presented clearly?

# Green Agent Project Ideas

- [https://docs.google.com/presentation/d/1TjtEjh6g9dZBsGxmAmcSp2EFakbmHpU\\_z31vnkf0c2Y/](https://docs.google.com/presentation/d/1TjtEjh6g9dZBsGxmAmcSp2EFakbmHpU_z31vnkf0c2Y/)

# Coding Example: Supporting *Tau-bench*

# 1. Sort out the interface

Principles:

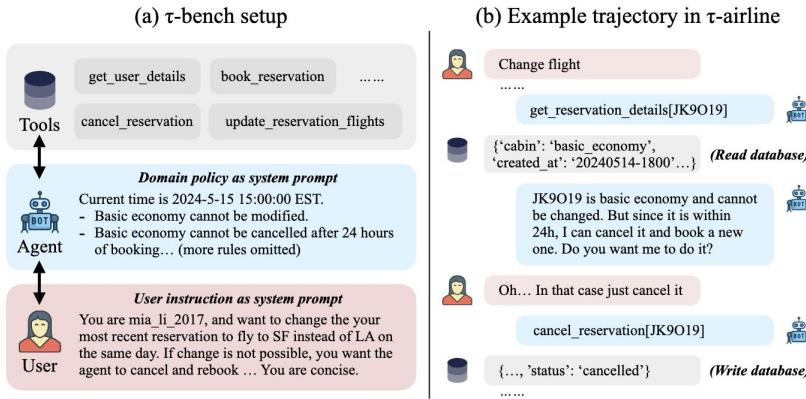
1. Human should be able to solve it if presented the same task.
2. The solving procedure should be as agent-friendly as possible. (so that the agent can solve it)

Example:

- **Web browsing agent:** url *vs.* tool actions
- **Coding agent:** provide coding env *vs.* provide repository & expect patches
- **Werewolf game agent:** text-based vote confirmation *vs.* tool-based confirmation

# 1. Sort out the interface

- Read the paper → think about task formulation
- Read their codebase → see how to deliver the



```
27
28     random.seed(config.seed)
29     time_str = datetime.now().strftime("%m%d%H%M%S")
30     ckpt_path = f'{config.log_dir}/{config.agent_strategy}-{config.model.split("/")[-1]}-{config.
31     if not os.path.exists(config.log_dir):
32         os.makedirs(config.log_dir)
33
34     print(f'Loading user with strategy: {config.user_strategy}')
35     env = get_env(
36         config.env,
37         user_strategy=config.user_strategy,
38         user_model=config.user_model,
39         user_provider=config.user_model_provider,
40         task_split=config.task_split,
41     )
42     agent = agent_factory(
43         tools_info=env.tools_info,
44         wiki=env.wiki,
45         config=config,
46     )
47     end_index = (
48         len(env.tasks) if config.end_index == -1 else min(config.end_index, len(env.tasks))
49     )
50     results: List[EnvRunResult] = []
51     lock = multiprocessing.Lock()
```

Figure 1: (a) In  $\tau$ -bench, an agent interacts with database API tools and an LM-simulated user to complete tasks. The benchmark tests an agent's ability to collate and convey all required information from/to users through multiple interactions, and solve complex issues on the fly while ensuring it follows guidelines laid out in a domain-specific policy document. (b) An example trajectory in  $\tau$ -airline, where an agent needs to reject the user request (change a basic economy flight) following domain policies and propose a new solution (cancel and rebook). This challenges the agent in long-context zero-shot reasoning over complex databases, rules, and user intents.

# 1. Sort out the interface

Two key challenges:

1. Cross-agent tool use
  - a. In the original repo, tool is directly provided to “completion” interface
  - b. How shall we evaluate using a standardized agent interface
    - i. “Special” white agents, with tool access
      1. Less standardized
    - ii. Explain this tool-access request to white agent, then ask for tool names / args
      1. Problem: cannot leverage agent internal tool-call mechanisms
    - iii. Provide an MCP → require dynamic discovery
  2. Migrate evaluation
    - a. Tool trace is not directly visible to green agent

# 1. Sort out

Two key challenges:

## 1. Cross-agent tool use

- a. In the original repo, tool info is hard-coded
- b. How shall we evaluate user-defined tools?
  - i. “Special” white agents, vs. green agents
    - 1. Less standardized
  - ii. Explain this tool-access rule
    - 1. Problem: cannot leverage LLM’s knowledge
  - iii. Provide an MCP → request

## 2. Migrate evaluation

- a. Tool trace is not directly visible to green agent

```
38     ]
39     for _ in range(max_num_steps):
40         res = completion(
41             messages=messages,
42             model=self.model,
43             custom_llm_provider=self.provider,
44             tools=self.tools_info,
45             temperature=self.temperature,
46         )
47         next_message = res.choices[0].message.model_dump()
48         total_cost += res._hidden_params["response_cost"] or 0
49         action = message_to_action(next_message)
50         env_response = env.step(action)
51         reward = env_response.reward
52         info = {**info, **env_response.info.model_dump()}
53         if action.name != RESPOND_ACTION_NAME:
54             next_message["tool_calls"] = next_message["tool_calls"][:1]
55             messages.extend(
56                 [
57                     next_message,
58                 {
59                     "role": "tool",
60                     "tool_call_id": next_message["tool_calls"][0]["id"],
61                     "name": next_message["tool_calls"][0]["function"]["name"],
62                     "content": env_response.observation,
63                 }
64             ]
65         )
66     return env_response
```

## 2. Design the workflow

- Kickoff script: send message to green agent to kick off the test
  - What information to include
  - Message format
- Green agent: coding-based, import tau\_bench
  - How to change to the initial prompt
  - How to incorporate the final scoring procedure / what are the metrics
- White agent: prompt-based / LLM-workflow
  - Which SDK to use
  - What prompt might help with the performance

```
1 import asyncio
2 import json
3 from a2a.types import SendMessageSuccessResponse
4 from .my_util import send_message_to_agent
5
6
7 task_config = {
8     "env": "retail",
9     "user_strategy": "llm",
10    # "user_model": "openrouter/openai/gpt-4o",
11    "user_model": "openai/gpt-4o",
12    "task_split": "test",
13    "task_ids": [1],
14 }
15
16
17 kick_off_message = f"""
18 Launch tau-bench to assess the tool-calling ability of the agent located at http://localhost:8001/ .
19 You should use the following configuration:
20 <task_config>
21 {json.dumps(task_config, indent=2)}
22 </task_config>
23 """
24
25
26 async def main():
27     agent_url = "http://localhost:9999/"
28     response = await send_message_to_agent(kick_off_message, agent_url)
29     if isinstance(response.root, SendMessageSuccessResponse):
30         response_text = response.root.result.parts[0].root.text
31         print("Agent response text:", response_text)
32     else:
33         print("Agent response:", response)
34
35
36
37 if __name__ == "__main__":
38     asyncio.run(main())
```

# 3. Impl: Green agent

```
class TauGreenAgentExecutor(AgentExecutor):
    def __init__(self):
        self.history = []

    async def execute(
        self,
        context: RequestContext,
        event_queue: EventQueue,
    ) -> None:
        # evaluation workflow
        user_input = context.get_user_input()

        task_config = parse_task_config(user_input)
        url = parse_http_url(user_input)
        assert len(task_config['task_ids']) == 1, "For demo purpose, here we run only one task"
        task_index = task_config['task_ids'][0]
        tau_env = get_env(
            env_name=task_config['env'],
            user_strategy=task_config['user_strategy'],
            user_model=task_config['user_model'],
            user_provider="openai",
            task_split=task_config['task_split'],
            task_index=task_index,
        )
        env_reset_res = tau_env.reset(task_index=task_index)
        obs = env_reset_res.observation
        info = env_reset_res.info.model_dump()

        task_description = tau_env.wiki + f"""
Here's a list of tools you can use: {tau_env.tools_info}
In the next message, I'll act as the user and provide further questions.
In your response, if you decide to directly reply to user, include your reply in a <reply> </reply> tag.
If you decide to use a tool, include your tool call function name in a <tool> </tool> tag, and include the arguments in a <args> </args> tag in JSON format.
Reply with "READY" once you understand the task and are ready to proceed.
"""

        res_check_ready = await send_message_to_agent(task_description, url)
        print("res_check_ready:", res_check_ready.root.result.artifacts[0].parts[0].root.text)
        is_ready = "READY" in res_check_ready.root.result.artifacts[0].parts[0].root.text.upper()
```

```
if __name__ == "__main__":
    agent_card_toml = load_agent_card_toml()
    agent_card_toml['url'] = f'http://:{HOST}:{PORT}/'

    request_handler = DefaultRequestHandler(
        agent_executor=TauGreenAgentExecutor(),
        task_store=InMemoryTaskStore(),
    )

    app = A2AStarletteApplication(
        agent_card=AgentCard(**agent_card_toml),
        http_handler=request_handler,
    )

    uvicorn.run(app.build(), host='0.0.0.0', port=9999)
```

(MCP-based impl would be different)

```
1 import datetime
2 from zoneinfo import ZoneInfo
3 from google.adk.agents import Agent
4 from google.adk.models.lite_llm import LiteLlm
5 from dotenv import load_dotenv
6
7 load_dotenv()
8
9 root_agent = Agent(
0     name="general_agent",
1     model=LiteLlm(model="openrouter/google/gemini-2.5-flash"),
2     description=(
3         "A general purpose agent that can assist with a variety of tasks."
4     ),
5     instruction=(
6         "You are a helpful assistant."
7     ),
8     tools=[],
9 )
0
1 from google.adk.a2a.utils.agent_to_a2a import to_a2a
2
3 # Make your agent A2A-compatible
4 a2a_app = to_a2a(root_agent, port=8001)
```

# Next step: integration with AgentBeats

After impl green/white/kick\_off → 90% DONE

Next: make it reproducible & open accessible → leverage agentbeats

Update checklist:

1. How to get (remote) agent URL / MCP server URL
  2. How to access LLM API
  3. How to report result & add traces
  4. Package the repo for platform hosting
- see documentation

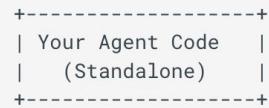
## Exposing an Agent

# Helpful

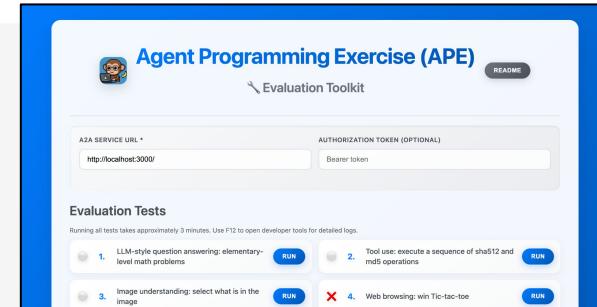
<https://google.github.io/adk-docs/a2a/intro/>

<https://a2a-protocol.org/latest/>

<http://ape.agentbeats.org/>



**Before Exposing:** Your agent code runs as a standalone component, but in this scenario, you want to expose it so that other remote agents can interact with your agent.



(Network Communication)



# Helpful tools

The screenshot shows a web browser window titled "A2A Inspector" at the URL "http://localhost:8001". The interface includes a "Connect" button and a "Agent Card" section displaying a JSON object.

```
{
  "capabilities": {},
  "defaultInputModes": [
    "text/plain"
  ],
  "defaultOutputModes": [
    "text/plain"
  ],
  "description": "A general purpose agent that can assist with a variety of tasks.",
  "name": "general_agent",
  "preferredTransport": "JSONRPC",
  "protocolVersion": "0.3.0",
  "skills": [
    {
      "name": "general_skills"
    }
  ]
}
```

# Helpful tools (Google ADK, for OpenAI, check the online

The screenshot shows a session in the Google ADK interface with the following details:

**SESSION ID:** df096512-9067-4518-99c7-e5033f243847

**Event 8 of 10**

**Tool Invocations:**

- get\_weather**:
  - Success: ✓ get\_weather (Response: The weather in New York is sunny with a temperature of 25 degrees Celsius (77 degrees Fahrenheit).)
  - Failure: ✘ get\_weather (Response: Weather information for 'Berkeley' is not available.)
- get\_current\_time**:
  - Success: ✓ get\_current\_time (Response: Please tell me the city you want to know the time for.)
  - Failure: ✘ get\_current\_time (Response: Sorry, I don't have timezone information for Berkeley.)

**Content Log:**

```
content:
  # parts:
  # 0:
    # functionCall:
    #   id: 'tool_0_get_current_time_e32702ef16xvish1#06'
    #   args:
    #     city: 'Berkeley'
    #     name: "get_current_time"
    role: 'model'
  partial: false
  usageMetadata:
    candidatesTokenCount: 7
    promptTokenCount: 102
    totalTokenCount: 109
  invocationId: 'e-a720181b-x00c-471c-a958-9951d-bea60a'
  author: 'weather_time_agent'
  actions:
    # stateDelta:
    # artifactDelta:
    # requestedActionConfigs:
    # requestedToolConfigurations:
  longRunningToolIds:
  id: '50075aef-a101-41ff-9ad3-72719f5ebdd6'
  timestamp: 175905400.154205
  title: 'functionCall/get_current_time'
```

**Message Input:** Type a Message...

# Reference & Credits

- <https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1246/slides/cs224n-spr2024-lecture11-evaluation-yann.pdf>
- <https://arxiv.org/abs/2503.16416>
- <https://arxiv.org/abs/2507.02825>

# Q&A

agentbeats concepts & platform / standardizing agent evaluation / ...