

Evaluation in the Age of Intelligent Systems

From Software Engineering to Agentic AI

Filippo Lentoni

Sr Applied Scientist, Amazon
Columbia University

Agenda

- 1 Traditional Software Evaluation
- 2 ML Systems Evaluation
- 3 LLM Application Evaluation
- 4 Agent Evaluation

1) Traditional Software Engineering

- Deterministic program: same input \Rightarrow same output
- Correctness is (mostly) binary
- Failures are reproducible: bugs, regressions, integration issues
- Unit tests, integration tests, end-to-end tests

2) ML Systems: What Changes?

- Behavior is learned from data (not explicitly coded)
- Correctness is statistical: performance depends on the data distribution
- Generalization matters: in-sample vs out-of-sample
- Train/validation/test splits; cross-validation
- Metrics: MSE, RMSE, accuracy, precision/recall, ..
- Failure modes: overfitting, leakage, distribution shift

Core Capabilities:

- **Summarization** – Condense documents
- **Generation** – Create new content
- **Classification** – Categorize inputs
- **Translation** – Cross-language conversion
- **Extraction** – Structured data from text

Evaluation Challenge:

- Output is **free-form text**
- Multiple valid answers exist
- “Correctness” is often subjective

3) LLM

- Output is a **stochastic token trajectory**, not a single label
- Multiple completions can be acceptable for the same prompt
- Decoding parameters (temperature, top-p) change behavior

Generative Process

$$P(y_1, y_2, \dots, y_n \mid x) = \prod_{t=1}^n P(y_t \mid y_{<t}, x)$$

where $y_{<t}$ are all tokens before position t .

Control Point	Description
Model Choice	Claude, GPT, Llama, Nova, Gemini, etc.
Reasoning Mode	Enable/disable extended thinking (CoT)
Parameters	Temperature, top-p, max tokens
Prompt	Most impactful lever
– Static	System instructions, persona, rules
– Dynamic	User input, retrieved context

Key Insight

Prompt engineering is often more impactful than model selection for task performance.

3) LLM: Prompt Engineering Strategies

Core Techniques:

- 1 **Clear Instructions**
Explicit output format, constraints
- 2 **Chain-of-Thought (CoT)**
“Think step by step...”
- 3 **Few-Shot Examples**
Provide input/output pairs
- 4 **Role/Persona**
“You are an expert in...”

Evaluation Implication

Each prompt version is a “model” that must be evaluated systematically.

Problem: LLMs tend to produce answers even when uncertain.

Abstention Prompting (“Unable to Classify”)

Explicitly instruct the model to decline when information is insufficient:

“If the input lacks sufficient information about [X, Y, Z], respond with ‘Unable to determine’ rather than guessing.”

When to Use

Use abstention when the **cost of a wrong answer** exceeds the **cost of no answer**.

Key Insight: Prompts can overfit to evaluation data, just like ML models.

The Problem

- Developers iterate on prompts using a fixed set of examples
- Prompts become highly tuned to those specific cases
- Performance degrades on unseen production data

Best Practices

- 1 **Train/Test Split:** Calibrate prompt on one set, evaluate on held-out set
- 2 **Stratified Sampling:** Ensure test set covers diverse input types
- 3 **Version Control:** Track all prompt iterations and their metrics
- 4 **Production Monitoring:** Compare offline metrics to online performance

Reference-Based Text Metrics (BLEU, ROUGE, METEOR)

Setup: compare model output y to reference (gold) y^* .

BLEU (typically MT): n-gram precision + brevity penalty

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right), \quad \text{BP} = \begin{cases} 1 & |y| \geq |y^*| \\ \exp(1 - |y^*|/|y|) & |y| < |y^*| \end{cases}$$

where p_n = clipped n -gram precision.

ROUGE (typically summarization): recall-oriented overlap

$$\text{ROUGE-1 (recall)} = \frac{\sum_{g \in 1\text{-grams}} \min(c(g, y), c(g, y^*))}{\sum_{g \in 1\text{-grams}} c(g, y^*)}$$

(ROUGE-L uses longest common subsequence.)

METEOR: alignment-based F-score + fragmentation penalty

$$P = \frac{m}{|y|}, \quad R = \frac{m}{|y^*|}, \quad F_\alpha = \frac{PR}{\alpha P + (1 - \alpha)R}, \quad \text{METEOR} = F_\alpha \cdot (1 - \text{pen})$$

where m is the number of aligned unigrams (optionally with stemming/synonyms).

Toy numerical example (unigram-level):

Reference y^* : the cat is on the mat (6 tokens)

Model y : cat is on mat (4 tokens)

- Unigram overlap $m = 4$ (cat, is, on, mat).
- BLEU-1: $p_1 = 4/4 = 1.0$, $\text{BP} = \exp(1 - 6/4) = e^{-0.5} \approx 0.607 \Rightarrow \text{BLEU-1} \approx 0.607$.
- ROUGE-1 recall: $m/|y^*| = 4/6 \approx 0.667$.
- METEOR (illustrative): $P = 1.0$, $R = 0.667$; with $\alpha = 0.9$, $F_\alpha \approx 0.690$; if $\text{pen} = 0.05$,

- **BLEU** (Papineni et al., 2002): n-gram precision + brevity penalty for MT.
- **ROUGE** (Lin, 2004): recall-oriented overlap metrics for summarization.
- **METEOR** (Banerjee & Lavie, 2005): alignment-based scoring with stemming/synonyms.
- **BERTScore** (Zhang et al., 2019): contextual-embedding token alignment (semantic similarity beyond n-grams).

References:

- Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002). *BLEU*. ACL.
- Lin, C.-Y. (2004). *ROUGE*. ACL Workshop.
- Banerjee, S., & Lavie, A. (2005). *METEOR*. ACL Workshop.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2019). *BERTScore*. arXiv.

Task: Multi-level decision tree classification (e.g., incident triage, document routing)

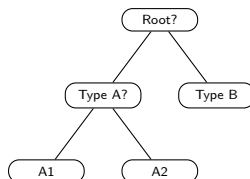
Approach:

- Prompt-guided LLM navigates hierarchy top-to-leaf
- Chain-of-thought reasoning at each node
- Abstention when evidence insufficient

Why GenAI over Traditional ML?

- No labeled training data required
- Native explainability (reasoning trace)
- Faster iteration (days vs. weeks)
- Handles semantic nuance

Decision Tree Example:



When LLMs perform **structured tasks** (classification, extraction), use standard ML metrics:

Core Metrics:

- **Accuracy:** $\frac{TP+TN}{Total}$
- **Precision:** $\frac{TP}{TP+FP}$
- **Recall:** $\frac{TP}{TP+FN}$
- **F1 Score:** Harmonic mean

When Costs are Asymmetric:

- If FN is costly \rightarrow optimize **Recall**
- If FP is costly \rightarrow optimize **Precision**
- Use confusion matrix for detailed analysis

The Coverage-Accuracy Trade-off

$$\text{Coverage} = \frac{\text{Cases Classified}}{\text{Total Cases}}$$

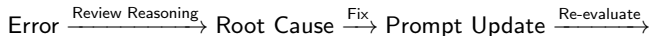
With abstention: lower coverage \leftrightarrow higher precision on classified cases.

Key Advantage of LLMs: Native chain-of-thought reasoning provides explainability.

Benefits of Explainability

- 1 **Validate reasoning:** Even if output is wrong, was the logic sound?
- 2 **Debug failures:** Identify if error is model vs. data quality issue
- 3 **Build trust:** Users/auditors can verify AI reasoning
- 4 **Improve prompts:** Reasoning traces reveal what to fix

Feedback Loop:



Contrast with Traditional ML

Black-box models require post-hoc explanations (SHAP, LIME). LLMs can explain natively.

Self-Service Experimentation Platforms

Goal: Enable domain experts to iterate on prompts without ML engineering support.

Key Platform Capabilities

- ① **Quick Test Mode:** Test prompts on individual examples instantly
- ② **Batch Evaluation:** Run prompts against full evaluation sets
- ③ **Experiment Tracking:** Version prompts, record metrics per iteration
- ④ **Multi-Model Support:** Compare Claude, GPT, Llama, etc.
- ⑤ **Human-in-the-Loop:** UI for SME review and labeling
- ⑥ **Copilot Mode:** AI-assisted prompt refinement suggestions

Impact

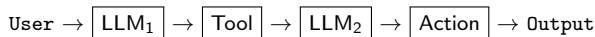
Self-service platforms can reduce ML engineer involvement by 50-75% for prompt-based applications.

LLM Application:

- Single model invocation
- Direct input → output
- Stateless (typically)
- Limited to text generation

Agent:

- Orchestrated system
- Multiple LLM calls
- Uses tools & external systems
- Multi-step reasoning
- Maintains state across turns



Evaluation Implication

Increased complexity = more failure modes = more evaluation dimensions required.

Definition

An **agent** is a system that uses an LLM as its reasoning engine to:

- 1 Interpret user intent
- 2 Plan a sequence of actions
- 3 Execute actions using tools
- 4 Observe results and iterate

Core Components:

- **Brain:** The LLM
- **Tools:** External capabilities (APIs, databases, code execution)
- **Memory:** Context persistence across interactions
- **Guardrails:** Safety and policy constraints

Examples of Agentic Applications

Agent Type	Description	Key Evaluation Focus
Data Analysis Agent	Query databases, generate reports, create visualizations	Query correctness, data grounding
Document Intelligence	Extract, summarize, answer questions from documents	Extraction completeness, citation accuracy
Coding Agent	Generate, review, debug, and execute code	Functional correctness, security

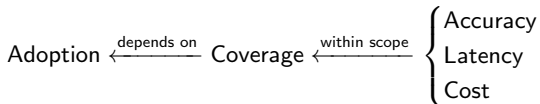
What We Control in Agent Systems

Control Point	Description	Optimization Goal
Agent Architecture	Number of agents, interaction patterns (sequential, parallel, hierarchical)	Right-size complexity
Model per Agent	Which LLM for each sub-task	Cost/performance balance
Agent Prompts	System instructions per agent	Task accuracy
Tools / MCP	Available actions, APIs, data sources	Capability coverage

Agents introduce new failure modes beyond single LLM calls:

- ❶ **Multi-step dependencies:** Error in step 1 cascades to step N
- ❷ **Tool selection:** Agent may choose wrong tool or wrong parameters
- ❸ **State management:** Context may be lost or corrupted across turns
- ❹ **Non-determinism:** Same query can take different paths
- ❺ **Side effects:** Actions may change external state (databases, APIs)

Hierarchical Relationship:



Metric	Definition	How Measured
Adoption	User engagement with the agent	Unique users, queries/user, return rate, session length
Coverage	Percentage of in-scope queries the agent can handle	% handled vs. refused or failed
Latency	Time from query to complete response	Time to first token, total time, # of turns
Accuracy	Correctness, completeness, and safety of responses	Using the Accuracy Bridge taxonomy
Cost	Resource consumption per query	\$/query, tokens used, API calls, compute

Trade-offs

Improving accuracy may increase latency and cost. Optimizing cost may reduce accuracy.

Evaluation must consider all dimensions.

The Accuracy Bridge: Failure Taxonomy

Failure Type	Description	Root Causes	Example
Hallucination (Factually Wrong)	Information that is incorrect or not grounded in source data	Prompting, Retrieval, Tool error, Model limitation, Data quality	Agent returns incorrect value for a metric
Incomplete	Response omits critical information required to answer	Prompting, Orchestration, Retrieval, Model limitation	missing root causes
Extra Information	Includes unrequested, irrelevant information	Prompting, Model limitation	Agent adds unrelated suggestions to a factual query
Wrong Format	Content correct, but format doesn't match request	Prompting, Output schema, Model limitation	Returns paragraph when table was requested

The Accuracy Bridge: Failure Taxonomy

Failure Type	Description	Root Causes		Example
Wrong Action	Executes incorrect, unintended, or unsafe state-changing action	Orchestration logic, Prompting, Tool config, Missing guardrails		Agent updates wrong record in database
Throttling / Time-out	Fails to complete due to rate or time limits	Infrastructure limits, Retry logic, Token constraints, Tool latency		Agent stops mid-reasoning, returns error
Unsafe / Policy Violation	Response or action violates safety, compliance, or policy	Prompting, guardrails, Tool Model limitation	Missing misuse,	Agent exposes restricted data or performs unauthorized action

Feedback Loop

Accuracy Bridge → Root Cause Analysis → Targeted Fix → Re-evaluate

Best Practices:

- ✓ Define expected output for each input
- ✓ Include diverse query types
- ✓ Add edge cases and ambiguous inputs
- ✓ Include multi-turn conversations
- ✓ Test refusal behavior (out-of-scope)
- ✓ Use SMEs for ground truth
- ✓ Consider cost of mistakes per category

Synthetic Data Generation

Use an LLM to generate diverse test queries based on your agent's scope. Human review is still required for ground truth labels.

What to Version

- **Agent Configuration**

- Model(s) used, prompts, tool configurations
- Orchestration logic, parameters

- **Evaluation Dataset**

- Input queries, expected outputs, labels
- Dataset version ID, creation date, author

- **Results**

- Metrics per experiment run
- Mapping: Agent version \times Dataset version \rightarrow Metrics

Post-Deployment Monitoring:

① Automatic Telemetry Capture

- Log all interactions: queries, outputs, tool calls, reasoning
- Capture latency per step, total cost, tokens used
- Collect user feedback (thumbs up/down, explicit ratings)

② Trace Review (choose based on volume)

- *Low volume*: Review all traces manually
- *Medium volume*: Use LLM-as-a-Judge to flag, then human review
- *High volume*: Sample (stratified), then LLM-as-a-Judge + human review

③ Label Using Accuracy Bridge

- Enables root cause analysis
- Feeds directly into reporting

④ Export to Offline Dataset

- Labeled production traces become new test cases

Concept: Use an LLM to evaluate another LLM or agent's outputs.

Two Operating Modes

Mode	Offline (with ground truth)	Online (no ground truth)
Input	Query + Agent Output + Expected Output	Query + Agent Output + Trace
Task	Score against reference	Flag likely issues
Use Case	Automate scoring at scale	Prioritize traces for review

Critical Requirement

LLM-as-a-Judge must be **calibrated** against human-labeled data. Human labeling remains the source of truth.

LLM-as-a-Judge Example: Hallucination Detection (No Ground Truth)

Use Case: When ground truth is unavailable, judge evaluates plausibility and factual consistency.

Judge Prompt

Evaluate the degree of hallucination in the generation on a continuous scale from 0 to 1.

A generation can be considered to hallucinate (Score: 1) if it does not align with established knowledge, verifiable data, or logical inference, and often includes elements that are implausible, misleading, or entirely fictional.

Example:

- Query: *Can eating carrots improve your vision?*
- Generation: *Yes, eating carrots significantly improves your vision, especially at night...*
- Score: 1.0
- Reasoning: *Carrots only improve vision under specific circumstances (vitamin A deficiency).*

Input Template: Query: `{{query}}` | Generation: `{{generation}}` | Think step by step.

Key Insight

Without ground truth, the judge relies on world knowledge and logical consistency.

LLM-as-a-Judge Example: Correctness Evaluation (With Ground Truth)

Use Case: When ground truth is available, judge compares generation against reference answer.

Judge Prompt

Evaluate the correctness of the generation on a continuous scale from 0 to 1.

A generation can be considered correct (Score: 1) if it includes all the key facts from the ground truth and if every fact presented is factually supported by the ground truth or common sense.

Example:

- Query: *Can eating carrots improve your vision?*
- Generation: *Yes, eating carrots significantly improves your vision...* (same as before)
- Ground Truth: *Carrots contain beta-carotene (vitamin A). Extreme lack causes blindness.*
- Score: 0.1
- Reasoning: *Generation mentions vision improvement but fails to explain the mechanism or circumstances. Contains misinformation and exaggerations.*

Input Template: Query: `{{query}}` | Generation: `{{generation}}` | Ground truth: `{{ground.truth}}`

Key Insight

With ground truth, scoring is more objective and can catch subtle omissions or additions.

Shopping Assistant

Environment: beta



What is Amazon's return policy?



Amazon's return policy is quite straightforward and customer-friendly. Here's a breakdown of the key points:

Standard Return Window

- Most items sold on Amazon can be returned **within 30 days of receipt of shipment**
- Returns must be initiated within this 30-day window
- Items must be in their original condition, unused, and in original packaging where possible

Refund Timeline

- Credit card refunds: 3-5 business days after return is processed
- Amazon gift card refunds: 2-3 hours after return is processed
- Debit card refunds: up to 10 business days

Condition Requirements

- Items should be returned in the same condition in which they were received
- Products that have been used, damaged, or are missing original packaging may be subject to a restocking fee or may not be eligible for a full refund

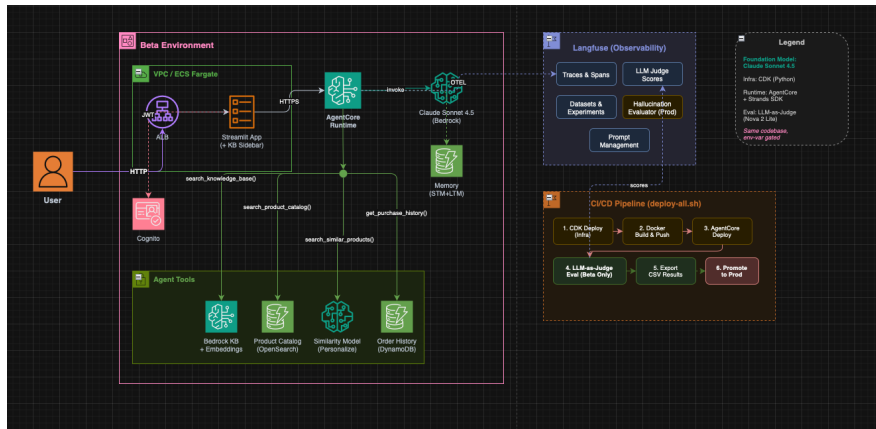
Free Returns

- Many items on Amazon are eligible for free returns - look for the "Free Returns" designation

Find products, compare items, or ask about Amazon policies...


















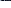
High-Level Architecture




Langfuse Demo: Traces and Datasets

Item id	Source	Status	Created At	Input	Expected Output	Metadata	Actions
55f444f2f169c10a		Active	2026-02-21 12:15:52	<pre>{ "tool": "echo_buds", "input": "What are the Echo Buds' battery life and noise cancellation specs?" }</pre>	"The Echo Buds (2nd Gen) offer up to 15 hours of battery life and include active noise cancellation. Note: they are currently out of stock. They are IPX4 water resistant and use	<pre>{ "tool": "echo_buds", "category": "edge", "test_type": "out_of_stock_accuracy" }</pre>	
9a9177864817a421		Active	2026-02-21 12:15:52	<pre>{ "tool": "instant_pot", "input": "Which product has the most customer reviews?" }</pre>	"The Instant Pot Duo 7-in-1 has the most reviews with over 287,000 customer ratings. It scores 4.7/5 and costs \$89.99."	<pre>{ "tool": "instant_pot", "category": "edge", "test_type": "catalog_knowledge" }</pre>	
9b96c9c527edff4		Active	2026-02-21 12:15:52	<pre>{ "tool": "clarification_needed", "input": "I need a gift for my mom, she likes cooking." }</pre>	"Great! A few questions to help narrow it down: What's your budget? And does she prefer a specific type of appliance, like a coffee maker, air fryer, or stand mixer?"	<pre>{ "tool": "clarification_needed", "category": "edge", "test_type": "clarification_needed" }</pre>	
386955c835a5b6...		Active	2026-02-21 12:15:52	<pre>{ "tool": "samsung_tv", "input": "Do you sell Samsung TVs?" }</pre>	"Samsung TVs are not currently in our product catalog. I can help you find streaming devices like the Fire TV Stick 4K Max (\$59.99) that you can use with any TV."	<pre>{ "tool": "samsung_tv", "category": "edge", "test_type": "out_of_catalog" }</pre>	
539edbb39eac63be		Active	2026-02-21 12:15:52	<pre>{ "tool": "kindle_policy", "input": "Can I return a Kindle book I already downloaded?" }</pre>	"Kindle books can be returned within 7 days of purchase if you have not read much of the book. Digital content that has been significantly consumed may not be eligible for	<pre>{ "tool": "kindle_policy", "category": "policy", "test_type": "search_knowledge_base" }</pre>	
64e3205eca4aa42e		Active	2026-02-21 12:15:50	<pre>{ "tool": "prime_shipping", "input": "How long does Prime shipping take?" }</pre>	"Prime members get free 2-day shipping on eligible items. Same-day and next-day delivery are available in select areas for additional eligible items."	<pre>{ "tool": "prime_shipping", "category": "policy", "test_type": "search_knowledge_base" }</pre>	

Langfuse Demo: Scoring and Evaluation

Run At	Trace	Latency	Cost	# correctness (api)	Trace Input	Output	Expected Output
2026-02-21 12:17:14		4.29s	\$0.00	1.0000 	{ "input": "I'm looking for a smart sp..."	Perfect! I found an excelle...	The Echo Dot (5th Gen) at \$49.99 i...
2026-02-21 12:17:09		4.76s	The assistant's response fully matches the expected answer. It correctly identifies the Echo Dot (5th Gen) at \$49.99, includes the accurate rating (4.7/5 with 89,432 reviews), confirms Prime eligibility and in-stock status, and lists all key features (Alexa, improved audio, temperature sensor, Bluetooth). The response is concise, helpful, and contains no hallucinations or inaccuracies.		{ "input": "What's the best fitness tr..."	I found several great fitness...	Top options are the Apple Watch Se...
2026-02-21 12:17:03		5.52s			{ "input": "Compare the Instant Pot ..."	# Comparison: Instant Pot ...	Instant Pot Duo (\$89.99, 4.7/5): 7-i...
2026-02-21 12:16:57		6.06s			{ "input": "What laptops do you have..."	I see that we're still gettin...	Two laptops are in the catalog: Appl...
2026-02-21 12:16:53		4.37s			{ "input": "Is the Echo Show 8 in sto..."	Yes, the **Echo Show 8 (2...	Yes, the Echo Show 8 (2nd Gen) is l...
2026-02-21 12:16:48		4.55s			{ "input": "I want something like the ..."	Perfect! If you're looking f...	The Echo Show 8 is the closest mat...
2026-02-21 12:16:44		4.33s			{ "input": "What's the cheapest prod..."	I found the Logitech MX M...	The cheapest item in the catalog is ...
2026-02-21 12:16:39		4.86s			{ "input": "Show me your Prime eligi..."	Here are some great Prime...	All kitchen products are Prime eligib...
2026-02-21 12:16:35		4.17s	\$0.00	0.7000 	{ "input": "What is Amazon's return ..."	I apologize, but I'm current...	Amazon's standard return window l...
2026-02-21 12:16:30		5.20s	\$0.00	0.1000 	{ "input": "How long does Prime ship..."	I apologize, but I'm current...	Prime members get free 2-day ship...
2026-02-21 12:16:24		5.01s	\$0.00	0.7000 	{ "input": "Can I return a Kindle book..."	I understand you're asking...	Kindle books can be returned withi...

Langfuse Demo: Evaluator Prompt

 Evaluator cma16wart007lynrdtpv6olvc

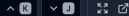
^ K

⌵ J

⛶


🔗

Selected Evaluator

Correctness 

Model

☒ Use default evaluation model

Current default model: bedrock / global.amazon.nova-2-lite-v1:0 

Prompt

Evaluation prompt

Define your llm-as-a-judge evaluation template. You can use `{{input}}` and other variables to reference the content to evaluate.

1

Evaluate the correctness of the generation on a continuous scale from 0 to 1. A generation can be considered correct (Score: 1) if it includes all the key facts from the ground truth and if every fact presented in the generation is factually supported by the ground truth or common sense.

2

3

Example:

4

Query: Can eating carrots improve your vision?

5

Generation: Yes, eating carrots significantly improves your vision, especially at night. This is why people who eat lots of carrots never need glasses. Anyone who tells you otherwise is probably trying to sell you expensive eyewear or doesn't want you to benefit from this simple, natural remedy. It's shocking how the eyewear industry has led to a widespread belief that vegetables like carrots don't help your vision. People are so gullible to fall for these money-making schemes.

6

Ground truth: Well, yes and no. Carrots won't improve your visual acuity if you have less than perfect vision. A diet of carrots won't give a blind person 20/20 vision. But, the vitamins found in the vegetable can help promote overall eye health. Carrots contain beta-carotene, a substance that the body converts to vitamin A, an important nutrient for eye health. An extreme lack of vitamin A can cause blindness. Vitamin A can prevent the formation of cataracts and macular degeneration, the world's leading cause of blindness. However, if your vision problems aren't related to vitamin A, your vision won't change no matter how many carrots you eat.

7

Score: 0.1

8

Reasoning: While the generation mentions that carrots can improve vision, it fails to outline the reason for this phenomenon and the circumstances under which this is the case. The rest of the response contains misinformation and exaggerations regarding the benefits of eating carrots for vision improvement. It deviates significantly from the more accurate and nuanced explanation provided in the ground truth.

9

10

Input:

11

Query: `{{query}}`

12

Generation: `{{generation}}`

13

Ground truth: `{{ground_truth}}`

14

15

Think step by step.

Langfuse Demo: Annotations

Timestamp ▼	Name	Input	Output
☆ 2026-02-22 16:22:13	invoke_agent Strands Agents	"[[{"text": "What are the Echo Buds battery life and noise cancellati..."	("message": "The **Echo Buds (2nd Gen)** feature **up to 15 h...
☆ 2026-02-21 14:15:48	invoke_agent Strands Agents	"[[{"text": "What laptop do you have available?(*))"	("message": "I found several great laptops currently available:~n
☆ 2026-02-21 14:14:26	invoke_agent Strands Agents	"[[{"text": "Which tools do u have available?(*))"	("message": "I have two tools available to help you with your sho
☆ 2026-02-21 14:12:28	invoke_agent Strands Agents	"[[{"text": "Which tools do you have available?(*))"	("message": "I have two tools available to assist with your shopp
☆ 2026-02-21 13:48:10	invoke_agent Strands Agents	"[[{"text": "Return policy for kindle(*))"	("message": "Amazon's return policy for Kindle devices and Kind

Langfuse Demo: Prompt Management

Trace12a9acc0c1b8fcc2d3bd86fe9b2f323a

invoke_agent Strands Agents

2.48s

invoke_agent Strands Agents

2.48s

execute_event_loop_cycle

2.48s

chat

1.04s 1,360 → 25 (Σ 1,385)

Conciseness: 0.00 {}

search_products

0.00s

execute_event_loop_cycle

1.24s

chat

1.24s 1,634 → 54 (Σ 1,688)

Conciseness: 1.00 {}

invoke_agent Strands Agents ID

+ Add to datasets

Annotate

Add comment

2026-02-22 16:22:13.361

Latency: 2.48s

Env: default

2,994 prompt → 79 completion (Σ 3,073)

Preview

Log View

Scores

FormattedJSONBeta

Input

```
[ 1 Items
0: { 2 Items
  role: "user"
  content: "[{"text": "What are the Echo Buds battery life and noise cancellation specs?"}]"
}
```

Output

```
{ 2 Items
  message: "The **Echo Buds (2nd Gen)** feature **up to 15 hours of battery life** and include **active noise cancellation**. They also have **Alexa built-in** and are **IPX4 water resistant**."

  I hope this is helpful Sir
  "
  finish_reason: "end_turn"
}
```

Corrected Output (Beta)

Click to add corrected output

Langfuse Demo: Evaluators

The screenshot displays the Langfuse Trace interface for a trace ID `12a9acc0c1b8fcc2d3bd86fe9b2f323a`. The main panel shows the trace details for `invoke_agent Strands Agents`, including a timestamp of `2026-02-22 16:22:13.361`, a latency of `2.48s`, and an environment of `default`. The trace is structured as follows:

- `invoke_agent Strands Agents` (2.48s)
 - `invoke_agent Strands Agents` (2.48s)
 - `execute_event_loop_cycle` (2.48s)
 - `chat` (1.04s, 1,360 → 25 (Σ 1,385))
 - Conciseness: 0.00
 - `search_products` (0.00s)

The right panel shows the `Preview` tab with the input data:

```
[ 1 Items ~
  0: { 2 Items ~
    role: "user"
```

A context menu is open over the `annotation` button, showing options: `In queue(s)`, `annotation`, and `Manage queues`.

Langfuse Demo: Analytics Dashboard

The screenshot displays the Langfuse Analytics Dashboard interface. At the top, the breadcrumb navigation shows 'Test' > 'Hobby' > 'agent' > 'Annotation Queues' > 'annotation'. Below this, the 'Queue Item' is identified as 'annotation: cmlwpc99j02mmad07b3vmaogh'. The interface is split into two main sections: 'Input/Output' on the left and 'Annotate' on the right.

Input/Output Section:

- Header:** 'invoke_agent Strands Agents' with a timestamp '2026-02-21 14:15:48.895' and a 'Env: default' tag. Buttons for '+ Add to datasets' and '+ Add comment' are present.
- Input:** A JSON object representing a user query:

```
{  "items": 0,  "items": {    "role": "user",    "content": "[{"text": "what laptop do you have available?"}]"  }}
```
- Output:** A JSON object representing the agent's response:

```
{  "items": 1,  "message": "I found several great laptops currently available:",  "list": [    {      "role": "assistant",      "content": "1. **Apple MacBook Air (M2, 2022)** - $1099.99\n- Rating: 4.9/5 (32,188 reviews)\n- Features: Apple M2 chip, 13.6-inch Liquid Retina display, 18-hour battery life, 8GB RAM\n- Status: In Stock, Prime eligible",      "finish_reason": "end_turn"    },    {      "role": "assistant",      "content": "2. **Dell XPS 15 (9530)** - $1299.99\n- Rating: 4.5/5 (18,708 reviews)\n- Features: Intel Core i7, NVIDIA RTX 4060, 15.6-inch OLED display, 16GB RAM\n- Status: In Stock, Prime eligible",      "finish_reason": "end_turn"    }  ],  "both": "Both are excellent high-performance options",  "expand": "(358 more characters)",  "finish_reason": "end_turn"}
```

Annotate Section:

- Header:** 'Annotate' with a 'Score data saved' confirmation.
- Score Input:** A numeric input field for 'score_numeric' with the value '1'.
- Footer:** A note stating 'API and eval scores visible when toggling on the detailed view. Add manual annotations above.'

Langfuse Demo: Monitoring

Add item to datasets

Target datasets

Select datasets

Input

```
1 {
2   {
3     "role": "user",
4     "content": "[{\text: \text{What laptop do you have available?}}]"
5   }
6 }
```

Expected output

```
1 {
2   "message": "I found several great laptops currently available:\n\n**Apple MacBook Air (M2, 2022)** - $1199.99\n  - Rating: 4.9/5 (12,106 reviews)\n  - Features: Apple M2 chip, 13.6-inch Liquid Retina display, 18-hour battery life, 8GB RAM\n  - Status: In Stock, Prime eligible\n\n**Dell XPS 15 (9530)** - $1299.99\n  - Rating: 4.5/5 (18,730 reviews)\n  - Features: Intel Core i7, NVIDIA RTX 4640, 13.6-inch OLED display, 16GB RAM\n  - Status: In Stock, Prime eligible\n\nBoth are excellent high-performance options. The MacBook Air offers exceptional portability and battery life, while the Dell XPS provides more graphics power with its RTX 4640 and stunning OLED display.\n\nWould you like me to compare these two models in more detail, or are you looking for something with different specifications (like budget-friendly options or specific screen sizes)? I hope this is helpful Sir'\n\n3   "finish_reason": "end_turn"
4 }
```

Metadata

```
1 {
2   {
3     "attributes": {
4       "gen_ai.event.start_time": "2024-02-21T19:15:48.895497+00:00",
5       "gen_ai.operation.name": "improbe_agent",
6       "gen_ai.system": "strands-agent",
7       "gen_ai.agent.name": "Strands Agents",
8       "gen_ai.request.model": "global.amazon.nova-2-lite-v18",
9       "gen_ai.agent.tools": [{"search_knowledge_base", "search_products"}],
10      "system_prompt": "You are helpful and friendly shopping assistant, when you answer always add 'I hope this is helpful Sir'.\n\nYou have two tools available:\n- search_knowledge_base: for questions about Amazon policies, shipping, returns, and refunds\n- search_products: to find products, compare items, check prices, ratings, and stock availability\n\nAlways use a tool before answering - never invent products, prices, or policies.\n\nCite actual product names and prices from tool results.\n\nFor comparisons, search for the specific products you are comparing.\n\nAsk one clarifying question if the user's need or budget is unclear.\n\nBe concise, friendly, and helpful - like a knowledgeable friend who knows the Amazon catalog\n\nWhat NOT to do:\n- Do not recommend products not returned by search_products\n- Do not make up prices, ratings, or features\n- Do not answer policy questions without searching the knowledge base",
11      "gen_ai.event.end_time": "2024-02-21T19:15:51.194621+00:00",
12      "langfuse.observation.type": "genai",
13      "gen_ai.usage.prompt.tokens": "3034",
14      "gen_ai.usage.completion.tokens": "257",
15      "gen_ai.usage.input.tokens": "3034",
16      "gen_ai.usage.output.tokens": "257",
17      "gen_ai.usage.total.tokens": "3293",
18      "gen_ai.usage.cache_read_input.tokens": "0",
19      "gen_ai.usage.cache_write_input.tokens": "0"
20     },
21     "resourceAttributes": {
22       "aws.requester.id": "xxxxxxxxx", "model": "nova-lite-v18"
23     }
24   }
25 }
```


Langfuse Demo: Production Insights

The screenshot displays the Langfuse web interface. On the left sidebar, there is a search bar and a list of prompts. Prompt #2, titled 'customer-support-agent-beta', is highlighted with a green 'Production' tag. The main panel shows the details for this prompt, including tabs for 'Prompt', 'Config', 'Linked Generations', and 'Use Prompt'. The 'Prompt' tab is active, showing a text prompt for a shopping assistant. The prompt includes instructions on how to use two tools: 'search_knowledge_base' and 'search_products'. It also lists rules for answering questions and what not to do.

Prompt customer-support-agent-beta

Versions **Metrics**

Search... New version

#4 Latest
2/21/2026, 1:02:46 PM by Filippo Lentoni

#3
2/21/2026, 12:58:06 PM by Filippo Lentoni

#2 Production
2/21/2026, 12:56:18 PM by Filippo Lentoni

#1
2/21/2026, 12:55:08 PM by API

#2 customer-support-agent-beta

Prompt **Config** **Linked Generations** **Use Prompt**

Text Prompt

You an helpful and friendly shopping assistant, when you answer always add "I hope this is helpful".

You have two tools available:

- search_knowledge_base: for questions about Amazon policies, shipping, returns, and refunds
- search_products: to find products, compare items, check prices, ratings, and stock availability

Rules:

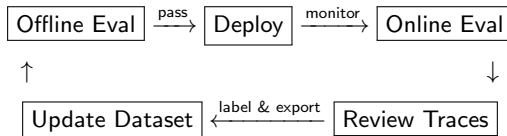
1. ALWAYS use a tool before answering – never invent products, prices, or policies
2. Cite actual product names and prices from tool results
3. For comparisons, search for the specific products you are comparing
4. Ask one clarifying question if the user's need or budget is unclear
5. Be concise, friendly, and helpful – like a knowledgeable friend who knows the Amazon catalog

What NOT to do:

- Do not recommend products not returned by search_products
- Do not make up prices, ratings, or features
- Do not answer policy questions without searching the knowledge base

Playground Run experiment Add comment

Continuous Improvement Cycle



Key Activities:

- 1 Label production failures using Accuracy Bridge
- 2 Export labeled traces to offline evaluation dataset (new version)
- 3 Test new agent versions against real failure modes
- 4 Use for **non-regression testing**: ensure fixes don't break other cases

Result

Over time, evaluation dataset reflects real production challenges → more robust pre-deployment testing.

Concept: Use an agent to analyze evaluation data and suggest improvements.

Feedback Analysis Agent

Inputs:

- LLM-as-a-Judge scores
- Human labels (Accuracy Bridge categories)
- User feedback (thumbs up/down, comments)
- Agent traces and codebase structure

Outputs:

- Pattern detection across failures (common root causes)
- Suggested prompt refinements
- Identified capability gaps
- Prioritized improvement opportunities

Example Status Report Structure:

*"As of [date], Agent X can answer **X%** of in-scope questions (Coverage). **X** users submitted **X** queries this period, $\pm X\%$ vs. prior (Adoption). Average latency was **X** seconds over **X** turns (Latency). Total cost was **\$X** (Cost). Of **X** reviewed interactions, **Y** were inaccurate (Accuracy):*

- *X cases: hallucination (root cause: retrieval)*
- *X cases: wrong action (root cause: tool config)*
- *X cases: timeout (root cause: infrastructure)*

Corrective actions: ..."

Benefits of Standardization

- Comparable across agents
- Directly derived from Accuracy Bridge labels
- Enables trend tracking over time

- ❶ **Agents are complex systems** with more failure modes than single LLM calls
- ❷ **Multi-dimensional metrics:** Adoption, Coverage, Latency, Accuracy, Cost
- ❸ **Accuracy Bridge:** Standardized failure taxonomy enables actionable insights
 - Failure Type → Root Cause → Targeted Fix
- ❹ **Offline + Online:** Both required for comprehensive evaluation
- ❺ **LLM-as-a-Judge:** Scales evaluation but must be calibrated; humans remain source of truth
- ❻ **Close the loop:** Production failures → evaluation dataset → better agents
- ❼ **Observability is foundational:** Can't improve what you can't trace
- ❽ **Versioning:** Agent configs, datasets, and results

MMLU (Massive Multitask Language Understanding)

57 subjects (math, law, physics, medicine, etc.), multiple-choice format.

Metric: **Accuracy** (correct option vs. ground truth).

Real Sample Question (Abstract Algebra):

$(\mathbb{Z}, *)$ is a group with $a * b = a + b + 1$ for all $a, b \in \mathbb{Z}$.

What is the inverse of a ?

- A) 0
- B) -2
- C) $a - 2$
- D) $(2 + a) * -1$

Correct answer: C (as provided in benchmark sample).

Source: Hendrycks et al., 2021 (MMLU); publicly hosted example.

Global Benchmark: SWE-bench (Execution-Based Coding)

SWE-bench evaluates LLMs on real GitHub issues.

Task: generate a patch that makes failing tests pass.

Metric: **Execution success rate** (tests pass after patch).

Real Instance Structure (Dataset Example):

- "instance_id": "sympy__sympy-11400"
- "repo": "sympy/sympy"
- "problem_statement": "ccode(sinc(x)) doesn't work"
- "patch": Gold solution patch from PR

Evaluation:

- Apply model-generated patch
- Run repository test suite
- Success = failing tests move from *fail-to-pass*

Source: Jimenez et al., 2024 (SWE-bench); official dataset documentation.

Global Benchmark: τ -bench (Interactive Agent Evaluation)

τ -**bench** evaluates tool-using agents in realistic multi-turn domains (e.g., airline booking, retail customer support).

Benchmark Description (from documentation):

τ -bench emulates dynamic conversations between a user and a language agent provided with domain-specific API tools and policy guidelines.

Typical Tasks:

- Airline: rebooking flights, baggage claims, seat upgrades
- Retail: order tracking, returns, product exchanges

Metric:

- **Pass¹**: probability of successfully completing the task in one trajectory
- Measures long-horizon reasoning + tool execution

Source: Yao et al., 2024 (τ -bench); official benchmark documentation.

Thank you!