# [Evaluation] - Evaluators

| Research | 🔵 Confident AI Blog - Resources to help teams stay confident in AI 🦡 Blog - Langfuse |
| --- | --- |
| | 🔴 Evaluating AI Agents - DeepLearning.AI |
| | 🔺 LLM Observability for AI Agents and Applications |
| **Relevant PRDs** | 🔵 [Evaluation] - Create Test Case and Manage with Datasets |
| | 🔵 [Evaluation] - Metric Config and Management |
| | 🔵 [Evaluation] - Experiment |
| **Design** | 🔲 Papago AI-Evaluation |

## Business Context

As an Agent Building Platform, our mission is to help users (agent builders) build better agents. Similar to traditional software development, having clear indicators that signal improvements after each update is crucial for product growth and success. In this phase of agent optimization, users seek solutions to construct a systematic evaluation workflow that provides consistent and reliable insights through metrics that matter most to them. Evaluating or comparing agents requires running them through the same user **inputs** and then assessing the **outputs** against predefined metrics and criteria.

## Problem

Even with output of an experiment, it is a daunting task to manually review and rate a run trace of an agent.

## Solution

Allow user to create and use LLM as Evaluators to automate the review and measuring each run result with the corresponding Metrics.

An LLM Evaluator will specify

1. Have the goal to output a value for a attached Metric.

2. Follow the guide by user on how to measure the Metric.

3. It will review, analyze the input and output a value

4. It can also give justification on why such value is outputted.

5. The input are test Case Input, Expected Output and the Actual Output.



Anatomy of an evaluator

# Requirement

**1. Evaluators list within Eval Config** 🔖 **FE-2353: [LLM Evaluator] Evaluators list within Eval Config** `LAUNCHED`

| Requirement | Wireframe / Design |
| --- | --- |

| List all Evaluators |  |
| --- | --- |
| | |

**2. Create a new version of an Evaluator**  🔖 **FE-2354: [LLM Evaluator] Create a new Evaluator**  `LAUNCHED`

| Requirement | Wireframe / Design |
| --- | --- |
| Create a new Evaluator is a 2 steps form<br><br>1. Describe your Evaluator<br>2. *Mapping configuration (TBD)* | |
| **Describe**<br><br>Attributes of an Evaluator includes:<br><br>1. Name<br>2. Prompt: write instruction to tune the LLM with guide on how to evaluate<br>3. Model: select base LLM model<br>4. Metric Config: Connect a Metric to be the core output of this evaluator<br>5. Reasoning: Add guide on justification of the evaluation result |  |

| | |
|---|---|
| ~~**Mapping**~~<br><br>~~Allow user to set up variables that can be use in the evaluator instruction~~<br><br>~~Support user map value from XXX to the variable~~<br><br>~~*An example use case is that the user can set a variable Ground_Truth and expect the evaluator to extract that from the executed run to see if the answer is truth.*~~<br><br>⚠️ deprioritized |  |
| Leaving the prompt reset the content.<br><br>Show confirmation pop up |  |
| User can create a new version of an existing evaluator |  |

### 3. Add an Evaluator to an Experiment

A n LLM Evaluator will:

1. Have the goal to output a value for a attached Metric.
2. Follow the guide by user on how to measure the Metric.
3. It will review, analyze the input and output a value
4. It can also give justification on why such value is outputted.
5. The input ~~are~~ test Case Input, Expected Output and the Actual Output.

| Requirement | Wireframe / Design |
|---|---|
| An evaluator can be add in a new experiment |  |
| An evaluator can be add in an executed experiment to give extra evaluation for such experiment | TBD |

## Adjustment

**1. Mapping Test Case attribute in Evaluator prompt**

User can reference `Experimenting Test Case's Attributes` in the Evaluator prompts.

**Experimenting Test Case's Attributes**

```
1  "experimentCase":{
2    "actualOutput":"",
3    "input":"",
4    "expectedOutput":"",
5    "context":"",
6    "conversationHistory":"",
7    "expectedToolCall":"",
8  }
```

## Examples Evaluator Prompt

User Satisfaction Evaluator

```
1  You are a helpful AI bot that checks for user satisfaction based on the response text and
   its prompt. Here is the data:
```

```
 2  [BEGIN DATA]
 3  ============
 4  [Response]: {actualOutput}
 5  ============
 6  [Prompt]: {input}
 7  [END DATA]
 8  Compare the Response above to the Prompt and determine if the Response is satisfactory given
    the Prompt.
 9  Your response must be a string, either SATISFIED or UNSATISFIED, and should not contain any
    text or characters aside from that.
10  The string UNSATISFIED means that the Response does not meet the user's needs or
    expectations based on the Prompt.
11  The string SATISFIED means the Response meets the user's needs or expectations based on the
    Prompt.
12
13  Then write out in a step by step manner an EXPLANATION to show how you determined if the
    user was satisfied or unsatisfied.
```

## Accuracy Evaluator

```
 1  You are a helpful AI bot that checks for Accuracy based on the Actual Output text and its
    Expected Output text. Here is the data:
 2  [BEGIN DATA]
 3  ============
 4  [Actual Output]: {testCase.actualOutput}
 5  ============
 6  [Expected Output]: {testCase.expectedOutput}
 7  [END DATA]
 8  Compare the Actual Output above to the Expected Output and determine if the Actual Output is
    resembles the given Expected Output.
 9  Your response must be a string, either LOW, Medium, and High, and should not contain any
    text or characters aside from that.
10  The string LOW means that the Actual Output does not resembling the content based on the
    Expected Output.
11  The string MEDIUM means that the Actual Output does resembling the content with minor
    vairation based on the Expected Output.
12  The string HIGH means that the Actual Output does resembling the content based on the
    Expected Output and replacible with that Expected Output content.
13
14  Then write out in a step by step manner an EXPLANATION to show how you determined if the
    Actual Output is LOW, MEDIUM, or HIGH.
```

Test Case

Experiment Test Case

- actual output
- input
- expected output
- expected tool call
- context
- conversation history

message:
type:
tools:[]

Evaluator Accuracy

"you are evaluator, compare {case.actualOutput} and {case.expected Output}"

Test Case
- input
- expected output
- expected tool call
- context
- conversation history

Test Case
- input
- expected output
- expected tool call
- context
- conversation history

Experiment

- actual output
- input
- expected output
- expected tool call
- context
- conversation history
- actual tool_call

????

Evaluator Tool_Call

"you are evaluator, compare {experiment.actualToolCall} and {testCase.expectedToolCall}"

**10038455** Complete
Trace Initiate Time: May/13/2024 17:08   Run Time: 2.76s   Token Cost: $0.000019

Experimental Results    Metrics

Input

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

Example:
Query: Can eating carrots improve your vision?
Context: Everyone has heard, "Eat your carrots to have good eyesight!" Is there any truth to this statement or is it a bunch of baloney? Well no.

Actual                          Output   Tool Call      Expected                  Output   Tool Call

Actual Output    Debugging Detail         Expected Output

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

Example:
Query: Can eating carrots improve your vision?
Context: Everyone has heard, "Eat your carrots to have good eyesight!" Is there any truth to this statement or is it a bunch of baloney?  Well 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. If your vision problems aren't related to vitamin A, your vision won't change no matter how many carrots you eat.
Ground truth: It depends. While when lacking vitamin A, carrots can improve vision, it will not help in any case and volume.
Metric: 0.3
Reasoning: The context correctly explains that

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

Example:
Query: Can eating carrots improve your vision?
Context: Everyone has heard, "Eat your carrots to have good eyesight!" Is there any truth to this statement or is it a bunch of baloney?  Well 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. If your vision problems aren't related to vitamin A, your vision won't change no matter how many carrots you eat.
Ground truth: It depends. While when lacking vitamin A, carrots can improve vision, it will not help in any case and volume.
Metric: 0.3
Reasoning: The context correctly explains that