# [Evaluation] – Create Test Case and Manage with Datasets

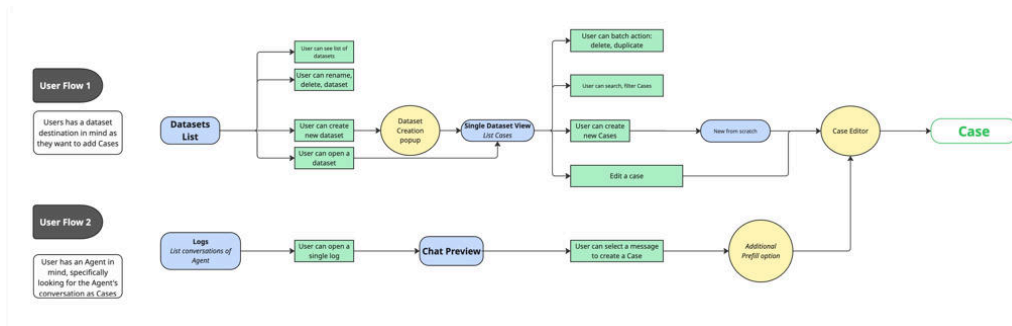| Research | ⬤ Confident AI Blog - Resources to help teams stay confident in AI 🔴 Blog - Langfuse |
| --- | --- |
| | ☺ Evaluating AI Agents - DeepLearning.AI |
| Wireframe | Ⓜ https://miro.com/app/board/uXjVIs9HNVg=/?moveToWidget=3458764633561380994&cot=14 Connect your Miro account 💡 |
| Design | ▦ Papago AI - Dataset |
| Relevant PRDs | 🔲 [Evaluation] - Create and manage Test Case for Agent Evaluation |

## Business Context

As an Agent Building Platform, our mission is to help users (agent builders) maximize their ROI from the effort invested in building and refining 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.

## Solution

This new module will enable users to create **Datasets**, which are collections of **Cases** *(Test Cases).* Each Case allows users to record inputs, expected outputs, tool calls, and other relevant details as a single, consistent unit for future evaluation jobs.

MVP: User Flow Map - Datasets & Case Management

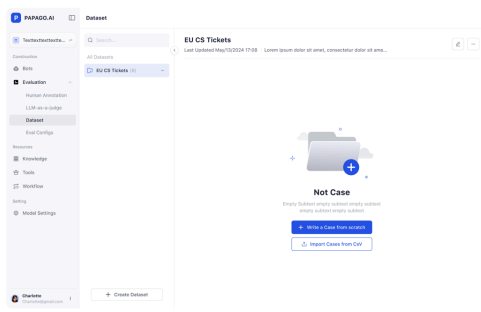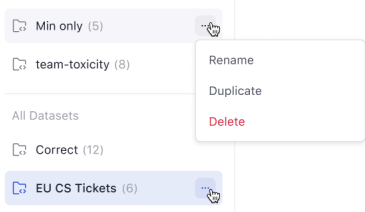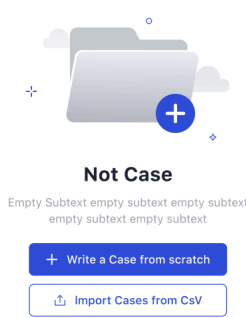> ℹ️ **Experiment** = **Dataset** + **Agent** + **Metrics (+ optional Evaluator).**
>
> - Uses made-up *(totally synthetic, no live logs needed)* or hand-picked content from real log to create **test cases**.
> - Lets users run the same **dataset of test cases** against **agents** to see who well it perform base on **Metrics**
> - Results belong to the experiment itself, not the agent or any single message. Notice there's no "This message has been evaluated..." label here—everything rolls up to the experiment level.
>
> We are building 3 related PRDs: Dataset, Metric and Experiment in this sprint

## Requirement

### 1. Create a dataset 🔖 FE-2340: [Case & Dataset] Create a dataset `LAUNCHED`
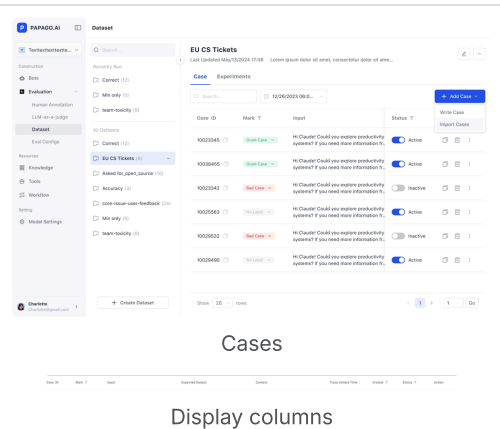
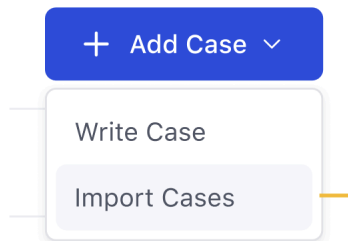| Requirement | Wireframe |
|---|---|
| Add Dataset as a new module on the platform, under Evaluation (also a new category).<br><br>Empty state message:<br><br>1. Title: No Dataset<br>2. Sub-title: Dataset manage your test cases for agent evaluation. |  |
| Create a new Dataset open the creation as pop up<br><br>1. Required Dataset Name<br>2. Data Description |  |

Datasets show a list of datasets as a row of times in a new left side pane. In Datasets view:

1. Each row is a specific dataset with
   a. name
   b. Number of cases per dataset
2. Datasets are searchable by name
3. The pane has two category:
   a. Recent Run
   b. All Dataset
4. User can Create a new dataset



Dataset can be Rename, Duplicate, and Delete



A new empty Dataset will prompt user to add cases

1. Write a Case from Scratch
2. Import Cases from CSV



## 2. Cases in a dataset

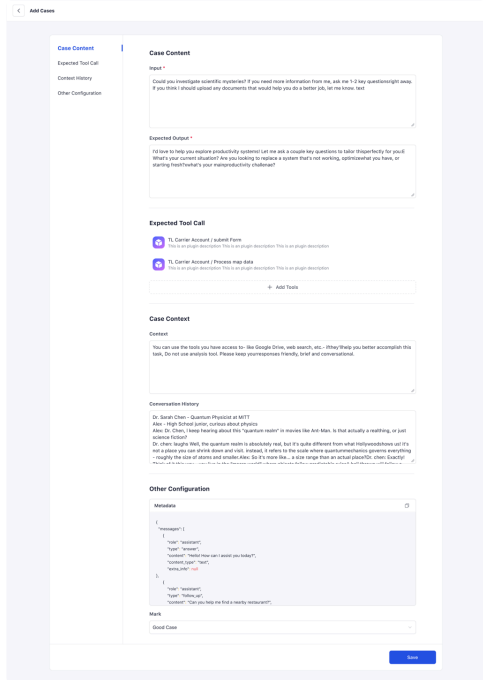User can see all Cases are listed in Cases tab.

1. Sorted by latest edited on top
2. Number of columns equal to number of attributes
   a. Ensure Input Expected Output and Context is visible in the initial width
   b. Overflow is scrollable
3. Columns are case attributes:
   a. Case ID: Can be copy
   b. Mark: [Good, Bad] and **Filterable**
   c. Input, Expected Output and Context: Fix height, overflow hidden, expandable by click to snap to content length
   d. ~~Trace Initiate Time~~ Last Edited: date time edited **Sortable**
   e. Creator: Workspace member **Filterable**
   f. Status: Toggle [Active, InActive] **Filterable**
4. Freeze action column on the right to contain move, duplicate, and remove
   a. Move clone the case to another dataset
   b. Remove cases will have pop-up to confirm
   c. Duplicate will clone the selected case and added right



Cases



Display columns
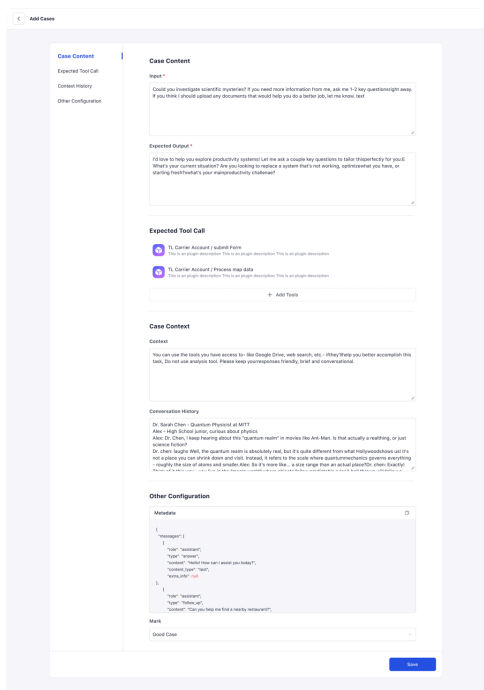
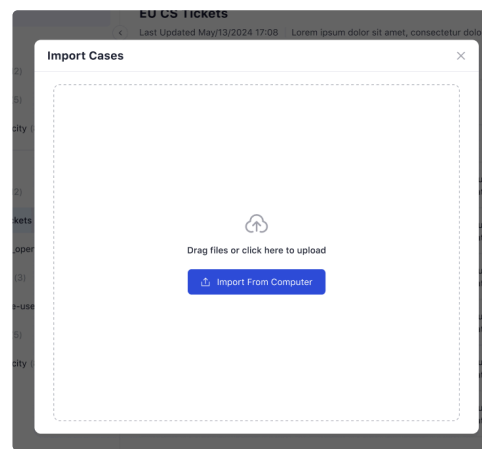| | |
|---|---|
| below. Re-sort the list only when user re-open the dataset<br>5. Cases view is configurable with:<br>  a. number of results to be show per page<br>  b. pagination between page<br>6. Cases can be search by name<br>7. Cases can be search | |
| Add a new case allow two option:<br>1. Write case (US3)<br>2. Import case (US4) |  |

## 3. Write case

ℹ️ Feature entry is US2's `Add case` button

| Requirement | Wireframe |
|---|---|
| Write case layout as a long form with section tree for quick navigation |  |

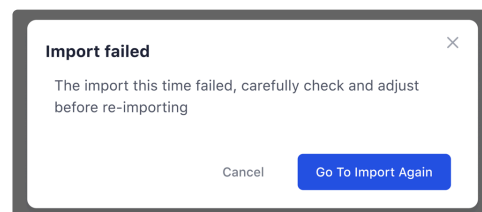| | |
|---|---|
| User can fill the information of:<br><br>1. Input **- Required**<br>2. Expected output **- Required**<br>3. Expected tool call. (*Detail in next requirement*)<br>4. Context<br>5. Conversation History<br>6. Metadata (Optional)<br>7. Tags (Detail in next requirement)<br>8. Status (not visible) is set default to Active |  |
| Each section must be collapsible due to the nature of lengthy content. | |
| Input box height is adjustable. This to make the text area length to adapt with content length. | |
| Input requirement to be validate on Save | |
| Each case has a system assigned unique Case ID | |
| Auto-save changes | |
| Content refresh when leave the page. Show confirmation that "Case has not been completed, do you want to leave" | |
| Save button will create a new case | |

## 4. Import cases from CSV

| Requirement | Wireframe |
|---|---|

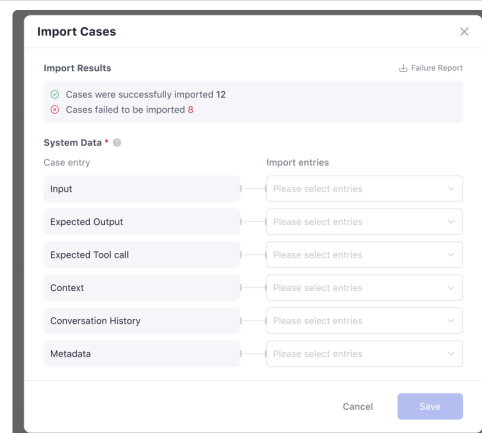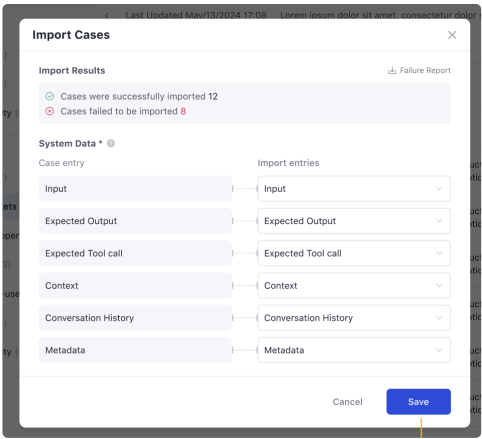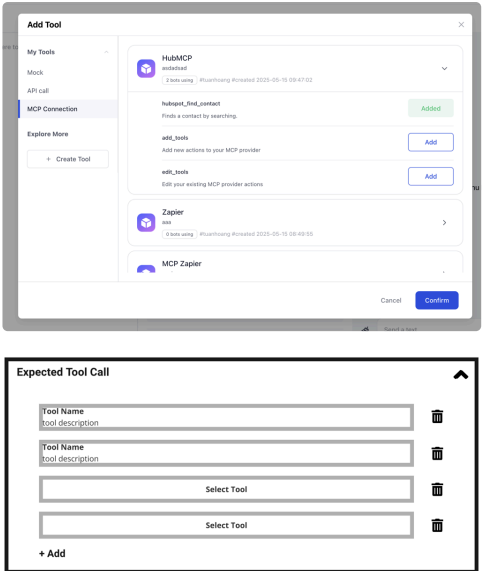| | |
|---|---|
| User can upload cases from CSV. There is no restriction to the .csv file |  |
| User will see upload status<br><br>1. Progressing: Percentage and Time Remaining<br>2. Failed:<br>   a. Show message<br>   b. Cancel and close to turn close the modal<br>   c. Go To Import Again will reopen the upload file widget |  |
| On success, user can start mapping CSV column to case attributes.<br><br>Partial success could happen, so showing how many case is success and not is also supported<br><br>User can select any column in the uploaded file to map with the field in the case. The dropdown show a list of all .csv columns<br><br>ℹ️ **Note to users:**<br><br>   Expected Tool Call is a string of multiple tool name, delimited by comma |  |

Examples:

Input, Expected Output, Expected Tool Call

xxx, qqqq, "tool_name_1, tool_name_2, tool_name_3"

ccc, hh, tool_name_2



⚠️ Adjustment: Not allow Metadata mapping win CSV import

## 5. Add tool calling to case

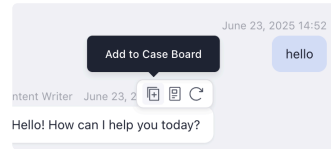| Requirement | Wireframe |
| --- | --- |
| Expected tool call section include add tool button. On click it open a list of tools that configured in the workspace *(Similar to existing selecting tools for agent, workflow).* <br><br> On select, the tool name, description are captured as expected tool call (Copy by value not reference) |  |

## 6. Create case from actual log

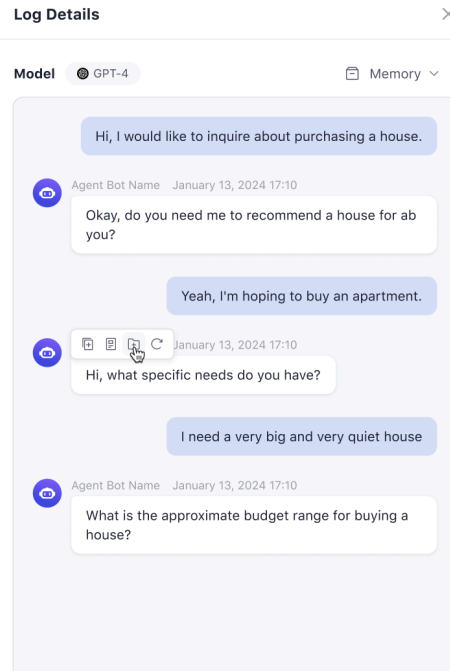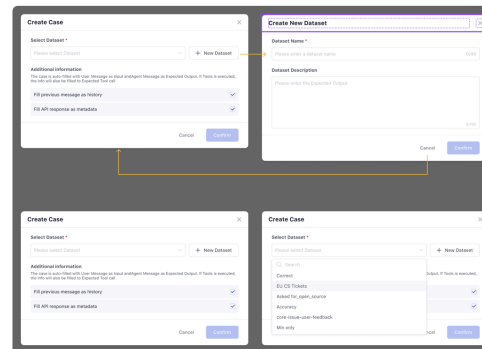| 🔖 FE-2345: [Case & Dataset] Create case from actual log LAUNCHED <br> **Requirement** | **Wireframe** |
| --- | --- |

Change label `Add to Case Board` → `Copy clipboard`



In an Agent Detail Log, user can start adding case from a specific log chat preview panel
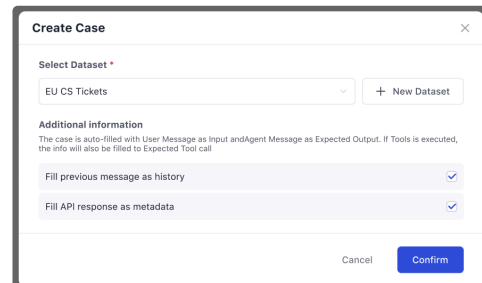
Add Create new Case button to ALL message.



`Create new case` open a Create Case editor

1. Select dataset Show list of dataset created
   a. Allow to Create new dataset. On dataset creation resume the case creation flo



User can select any message to add case. The platform will detect the relevant messages (User, Agent, Tool call) and filled in the new case. Specifically, the case is auto-filled with User Message as Input and Agent Message as Expected Output. If Tools is executed, the info will also be filled to Expected Tool call.

User can also choice wether to add previous message as conversation history and fill the metadata with the whole API response from the server. Default value:

1. Fill history: On
   a. If user select message within the first conversation *(1 user message following by the agents responses)* history can be empty
2. Fill metadata: Off
   a. Meta on will copy the response JSON from our backend as metadata



## Adjustment

**Filter UI element**

**Problem**

After testing on stage, I found that it is cumbersome to get to use to many UI style for the same filtering



❌ No indicator which status is currently filter



❌ Only one can be toggle on at a time.



❌ List of all creators are known. Should show to select instead of search which is cumbersome and prone to typos

**Suggestion**

Re-use the filter in metrics that help:

1. Multiple filter apply
2. Clear indicator of which filter is applying
3. List of all possible filter criteria

- ☑ Boolean
- ☐ Numeric
- ☑ Categorical
- ☐ Range Qualities

Reset    **Confirm**

Metric Filter

- ☑ Boolean
- ☐ Numeric
- ☑ Categorical
- ☐ Range Qualities