# Mastering Evals The Key to Building Effective LLM Based Apps

Hilik Paz, Arato.ai





## What are Evals?

Methods for **systematically** determine the performance of LLMs and Models.

Measurable and Repeatable, enabling data driven iterative improvements



# What Changed?

Traditional software testing vs. LLM evaluations

Clear **boundaries** vs Complex nuanced output.

**Deterministic** vs Probabilistic Good or Bad is no longer binary

From Pass/Fail to Scoring



# Why Evals Matter?

- Ensure LLM outputs **align** with business objectives
- Reduce the **risk** of unforeseen issues in production
- Enable **measurable** outcome.





## The AppMaster Support Bot

- Al-powered support bot for AppMaster users
- Answers **technical** questions using a **knowledge base**
- Goal: Provide accurate, fast, and context-aware responses, with a polite and friendly tone of voice.



## **Defining the Initial Prompt**

We start with a **simple** and clear prompt. We set the role Provide clear instructions And set the expected output

But is this enough? We need a way to **measure** its effectiveness





## **How Do We Measure** Success?

Define key evaluation criteria Choose evaluation methods

Establish test cases



## What Makes a Good Eval?

Aligns with the **Business** Use Case

### Specific and Measurable

Example: For our AppMater support bot, evals could measure **accuracy**, **politeness**, and **response time**, and expected output **format**.



# **Types of Evals**

### Content

Match expected outcomes, Similarity, Sensitive data etc...

### Format

Text structure, JSON schema, tabular data etc...

### Qualities

Grounding, Correctness, Detail etc...

### Guardrails

Non-functional aspects like Diversity, Politeness etc...



## **Evals: Different Mechanisms**

### Deterministic

Rule based, Code, Text matching, Regex etc...

### **Model Based**

LLM as an evaluator Vector Similarity



## **Pros and Cons**

### Deterministic

Cheap and Accurate Rigid and Limited

## **Model Based**

Flexible and Easy to use Expensive and Non Deterministic







## **A Sample Expected Response**

#### Expected output

#### {

"answer": "Hi Tim, I'm your friendly AppMaster support bot. To find your project files in AppMaster, you should first check if they were saved in the default directory or a custom location you specified during setup. Additionally, you can use the search functionality in the AppMaster interface to locate missing files by their names or tags. If you suspect that the files might have been deleted or moved, verify their existence in your local storage.",

"analysis": "I used document 1 and document 3 to provide a comprehensive answer. Document 1 helped identify where to look for the files, while document 3 offered a method to search for them within the AppMaster interface."

}

□ Copy ‡ Close

# Similarity Eval

#### Expected output

#### {

"answer": "Hi Tim, I'm your friendly AppMaster support bot. To find your project files in AppMaster, you

should first check if they were saved in the default directory Additionally, you can use the search functionality in the AppMa names or tags. If you suspect that the files might have been de local storage.",

"analysis": "I used document 1 and document 3 to provide a co where to look for the files, while document 3 offered a method interface."

```
}
```

#### Results

#### {

"answer": "To find your project files in AppMaster, first che you specified during setup. If you still can't find them, verif weren't accidentally deleted or moved. You can also use the sea locate missing files by their names or tags.",

"analysis": "I used document 1, document 2, and document 3 to helped identify where to look for the files, document 2 address moved, and document 3 suggested using the search functionality to rotate the lites.

}

![](_page_17_Figure_12.jpeg)

⊡ Copy # Close

![](_page_18_Picture_0.jpeg)

# **JSON Response Eval**

piviastei	Aa Edit content validatio	on				x	ot-4o-mini 🗸
oppMaster G						_	
valuations	Validation name *					_	
V.U.	Json Schema					×	
	validation description [optional]					_	
	Description						
	Description						
	Description					_	
^\{	\s*"answer":\	s*"(?:[^"\`	\] \\.) *?	9",∖s*"a	naly	sis'	×
^\{	\s*"answer":\	s*"(?:[^"\`	\] \\.) *?	)",∖s*"a	naly	sis'	×
۳ <b>۸\{</b> ۶ {{user_c	\s*"answer":\	s*"(?:[^"\`	\] \\.) *?	9",∖s*"a	naly:	SiS'	× ing VX Re

![](_page_20_Picture_0.jpeg)

Based Eval	
Edit custom LLM validation	×
Validation name * Core purpose of the validation, in one or two words	
Politeness	×
What part of the data to validate? *	
Response v	
Validation description [optional]	
Example: Check if the prompt mentions a country	
Yes/No question for validation *	
Write your validation question below. Phrase it as a Yes/No question	and the second

Does the response answers the user question in a polite manner that mentions the Friendly AppMaster Support bot ?

![](_page_21_Picture_2.jpeg)

![](_page_22_Picture_0.jpeg)

# Experimenting

## **Prompt Template**

Data

**Eval Question / Metric** 

![](_page_23_Picture_4.jpeg)

# What Makes a Good Dataset

### Diversity.

Representativeness of **real-world** inputs.

Clear **expected** outputs.

![](_page_24_Picture_4.jpeg)

![](_page_25_Picture_0.jpeg)

# Thank You!

hilik@arato.ai

linkedin.com/company/arato-ai

www.arato.ai

![](_page_26_Picture_4.jpeg)