A arize LLM Guardrails 101

A comprehensive guide to protecting your application, including from itself.



Introduction

As LLM applications become more common, so too do jailbreak attempts, exploitations of these apps, and harmful responses. More and more companies are falling prey to damaging news stories driven by their chatbots <u>selling cars for \$1</u>, <u>writing poems critical of their owners</u>, or <u>dealing out disturbing replies</u>.

Fortunately, there is a solution to this problem: **LLM guardrails**.

Section 1: What Are LLM Guardrails?

Section 2: Types of Guardrails

Section 3: Dynamic Guards



What Are LLM Guardrails?

A key part of any modern large language model (LLM) application, LLM guardrails allow you to protect your application from potentially harmful inputs, and block damaging outputs before they're seen by a user. As LLM jailbreak attempts become more common and more sophisticated, having a robust guardrails approach is critical.

Let's dive into how guardrails work, how they can be traced and acted upon, and how you can use them to avoid becoming the next big news story.

How Do Guardrails Work?

LLM guardrails work in real-time to either catch dangerous user inputs or screen model outputs. There are many different types of **guards** that can be employed, each specializing in a different potential type of harmful input or output.

What Are the Primary Use Cases for Guardrails In AI Development?

Common input guard use cases include:

- Detecting and blocking jailbreak attempts
- Preventing prompt injection attempts
- Removing user personally identifiable information (PII) before it reaches a model

Common output guard use cases include:

- Removing toxic or hallucinated responses
- Removing mentions of a competitor's product
- Screening for relevancy in responses
- Removing NSFW text

There can be a lot of ground to cover. Fortunately, tools like Guardrails AI offer a hub of different guards that can be added to your application.

		D.	Sign in to get started
	Validators		
USE CASES	Validators are basic Guardralis components that are used to validate an aspect of an LLM workflow. Validators can be used a to prevent end-users from seeing the results of faulty or unsafe LLM responses.		
TEXT2SQL STAG STRUCTURED DATA	Q Search		
	Showing 57 of 57 validators Generate Code		
😂 ETIQUETTE 🔒 JAILBREAKING 🚫 BRAND RISK	Arize Dataset Embeddings	Select	S Competitor Check Select
CODE EXPLOITS XINVALID CODE	Validates that user-generated input does not match dataset of jailbreak Last updated 1 day ago		Flags mentions of competitors. Fixes responses by filtering out competitor names. Last updated 22 hours ago
E FACTUALITY FORMATTING	Correct Language	Select 🗍	Select PII Select C
INFRASTRUCTURE REQUIREMENTS	scb-10x/correct_language Last updated 2 months ago ()		Detects personally identifiable information (PII) in text, using Microsoft Presidio. Last updated 1 week ago
CONTENT TYPE	Detect Prompt Injection	Select	Select
🖀 STRING 🖉 SQL 🖉 INTEGER 🔠 FLOAT	Finds prompt injection using the Rebuff prompt library. Last updated 2 months ago ③		This validator checks if the extracted summary sentences match the original Last updated 2 months ago
OBJECT CODE LIST CSV 2 PYTHON	Extractive Summary	Select	Gibberish Text Select
CERTIFICATION	Uses fuzzy matching to detect if some text is a summary of a docum Last updated 2 months ago	ent.	A Guardrails Al validator to detect gibberish text. Last updated 2 months ago 💿
LANGUAGE	High Quality Translation	Select 🗌	Select Select
I EN	A validator that checks if a translation is of high quality. Last updated 2 months ago		Validates logical consistency and detects logical fallacies in the model output Last updated 6 days ago
	NSFW Text	Select 🗌	Select
	A Guardrails Al validator to detect NSFW text Last updated 1 day ago 🕥		Checks for profanity in text, using the alt-profanity-check library. Last updated 1 week ago 💿
	Provenance Embeddings	Select 🗌	Provenance LLM Select
	Compares embeddings of generated and source texts to calculate pr Last updated 2 weeks ago	ovenance.	guardrails/provenance_lim Last updated 2 months ago 💿
	QA Relevance LLM Eval	Select 🗌	Relevancy Evaluator Select
	Makes a second request to the LLM, asking it if its original response Last updated 2 months ago 💿	was relevant	Validates that the reference text contains information relevant to answering the Last updated 5 days ago

If a message in an LLM chat fails a guard, then the guard can take one of a few different corrective actions: providing a default response, prompting the LLM for a new response, or throwing an exception. For a guard that detects responses that may be damaging to your company's reputation, regenerating a response may work just fine. However for a guard that detects jailbreak attempts, a default response may be more appropriate.

How To Use Guardrails

A variety of packages are available. Arize offers an integration with Guardrails AI to give you best-in-class observability alongside top-of-line security.

Here are the basic steps necessary to use Guardrails with Arize. If you're looking for a more in-depth example, check out the following tutorial notebook walking through how to enable guards on a RAG pipeline.



Import Packages and Initialize Arize

First, make sure you imported the correct packages and connected your application to your Arize dashboard:

```
Unset
pip install opentelemetry-sdk opentelemetry-exporter-otlp
openinference-instrumentation-guardrails guardrails-ai arize-otel
Unset
from arize_otel import register_otel, Endpoints
from openinference.instrumentation.guardrails import
GuardrailsInstrumentor
# Setup OTEL via our convenience function
register_otel(
    endpoints = Endpoints.ARIZE,
    space_key = getpass(" P Enter your Arize space key in the space
settings page of the Arize UI: "),
    api_key = getpass(">> Enter your Arize API key in the space
settings page of the Arize UI: "),
   model_id = "sales-demo-dataset-embeddings-guard", # name this to
whatever you would like
)
# Use the Arize autoinstrumentor to add tracing to your application
GuardrailsInstrumentor().instrument(skip_dep_check=True)
```

Prepare Your Guards

Next, you need to add whichever guards you're looking to use to your project. These can be downloaded directly from the Guardrails hub. For this example, we'll use our *ArizeDatasetEmbeddings* guard.

Unset guardrails hub install hub://arize-ai/dataset_embeddings_guardrails

Initialize Guardrails and Add Your Guard

Now we are ready to initialize our guard. Here we can specify whether this guard will act on prompts or responses and what it should do if it catches a bad input or output. We also disable Guardrails tracing in this step, as we're using Arize to view our telemetry.

```
Unset
from guardrails import Guard
guard = Guard().use(ArizeDatasetEmbeddings, on="prompt",
on_fail="exception")
guard._disable_tracer = True
```

Make Calls to your LLM

With that, we're ready to make protected calls to our models. We can do this by invoking our Guard.

Our guard will take care of any retries or default responses necessary based on our earlier setup.

```
Unset
validated_response = guard(
    llm_api=openai.chat.completions.create,
    prompt=prompt,
    model="gpt-3.5-turbo",
    max_tokens=1024,
    temperature=0.5,
)
```

How To View Guard and Trace Data

If you followed the steps so far, you should already be seeing trace data in your project.



Types of Guardrails

The Balancing Act of LLM Guards

Implementing guardrails for AI systems is a delicate balancing act. While these safety measures are important for responsible AI deployment, finding the right configuration is necessary to maintain both functionality and security.

To effectively manage your guards, we strongly recommend using specialized tools such as <u>Guardrails AI</u> or <u>NemoAI</u>. Although it's technically possible to implement each guard type manually, this approach quickly becomes unwieldy as your system's complexity grows.

On the other hand, it's important to resist the temptation to over-index on guards. It may seem prudent to implement every conceivable safety measure, but this approach can be counterproductive. Excessive guardrails risk losing the intent of the user's initial request or the value of the app's output. Instead, we advise starting with the most critical guards and expanding judiciously as needed. Tools like Arize's AI search can be helpful in identifying clusters of problematic inputs to allow for targeted guard additions over time.

Additionally, more sophisticated guards that rely on their own LLM calls introduce additional costs and latency which must also be considered. To mitigate these issues, consider using small language models like GPT 40 mini or Gemma-2b for auxiliary calls.



leverage Arize-Phoenix with Guardrails AI to set up a guard that blocks an LLM from responding from attempted jailbreaks]

What Are the Types of LLM Guardrails?

Input Validation and Sanitization

Input validation and sanitization serve as the first line of defense in AI safety. These guards ensure that the data fed into your model is appropriate, safe, and in the correct format.

SYNTAX AND FORMAT CHECKS

While basic, these checks are important for maintaining system integrity. They verify that the input adheres to the expected format and structure. For instance, if your model expects certain parameters, what happens when they're missing? Consider a scenario where your RAG retriever fails to return documents, or your structured extractor pulls the wrong data. Is your model prepared to handle this malformed request? Implementing these checks helps prevent errors and ensures smooth operation.

CONTENT FILTERING

This guard type focuses on removing sensitive or inappropriate content before it reaches the model. Detecting and removing personally identifiable information can help avoid potential privacy issues, and filtering NSFW or toxic language can ensure more appropriate responses from your LLM. We recommend implementing this guard cautiously – overzealous filtering might inadvertently alter the user's original intent. Often, these types of guards are better suited filtering the outputs of your application rather than the inputs.

JAILBREAK ATTEMPT DETECTION

These are the guards that prevent massive security breaches and keep your company out of news headlines. Many <u>collections</u> of jailbreak prompts are available, and even advanced models can fail on up to <u>40% of these publicly-documented attacks</u>. As these attacks constantly evolve, implementing effective guards can be challenging; we recommend using an embedding-based guard <u>like Arize's</u>, which can adapt to changing strategies. At minimum, use a guard connected to a common library of prompt injection prompts, such as <u>Rebuff</u>.

Output Monitoring and Filtering

Output guards generally fall into two categories: preventing damage, and ensuring performance.

PREVENTING DAMAGE

Examples of this include:

- System Prompt Protection: Some attacks try to expose the prompt templates your system uses. Adding a guard to detect system prompt language in your outputs can mitigate this risk. Just be sure to avoid exposing this same template within your guard's code!
- NSFW or Harmful Language Detection: Allowing this type of language in your app's responses can be extremely harmful to user experience and your brand. Use guards to help identify this language.
- Competitor Mentions: Depending on your use case, mentioning competitors might be undesirable. Guards can be set up to filter out such references.

ENSURING PERFORMANCE

When it comes to performance, developers face a choice between using guards to improve your app's output in real-time or running offline evaluations to optimize your pipeline or prompt template. Real-time guards introduce more latency and cost but offer immediate improvements. Offline evaluations allow for pipeline optimization without added latency, though there may be a delay between issue discovery and resolution. We recommend starting with offline evaluations and only adding performance guards if absolutely necessary.

- Hallucination Prevention: Guards can prevent hallucinations by comparing outputs with reference texts or, when unavailable, cross-referencing with reliable sources like Wikipedia.
- Critic Guards: This broad category involves using a separate LLM to critique and improve your pipeline's output before sending it to the user. These can be instructed to focus on relevancy, conciseness, tone, and other aspects of the response.

What Strategies and Techniques Best Complement AI Guardrails?

While guardrails are important for AI safety, they're not the only measures you should consider.

Fence Your App from Other Systems

Isolating your AI application from other systems and networks creates an additional layer of security, limiting potential vulnerabilities and preventing unauthorized access or data leakage. Implement strict access controls and use secure APIs for necessary inter-system communications.

Red Team Pre-Launch

Before deploying your AI system, conduct thorough red team exercises. This involves having a dedicated team attempt to break, manipulate, or exploit your system in ways that malicious actors might. These simulated attacks can reveal vulnerabilities that weren't apparent during development and allow you to address them before public release.

Monitor Your App Post-Launch

Perhaps the most critical strategy is continuous_ <u>LLM production monitoring</u> after deployment. No amount of pre-launch testing can anticipate all real-world scenarios. Implement robust logging and monitoring systems to track your app's performance, user interactions, and potential issues, and regularly analyze this data to identify patterns, anomalies, or emerging problems. This will allow you to refine your guardrails and respond swiftly to any security concerns.



Dynamic Guards

While static guards are great at filtering out predefined content like NSFW language, they struggle when faced with sophisticated attacks like jailbreak attempts, prompt injection, and more. These dynamic threats require equally dynamic defenses that can evolve alongside the attackers' strategies.

Manually updating guards to counter new threats is a near impossible task, quickly becoming unsustainable as attack vectors multiply. Fortunately, two approaches allow us to create adaptive guards that can keep pace with emerging threats: few-shot prompting and embedding-based guards.



Few-Shot Prompting

This technique involves adding examples of recent jailbreak attempts or other attacks directly into your guard's prompt. By exposing the guard to real-world attack patterns, you improve its ability to recognize and thwart similar threats in the future.

Embedding-Based Guards

This is a more sophisticated approach to dynamic protection. It involves comparing the embedding of a user's input against a database of known attack embeddings. By checking the similarity between these representations, the system can identify and block potentially malicious inputs that exceed a predefined similarity threshold. Arize has developed an easy-to-use but powerful implementation of this concept with our <u>ArizeDatasetEmbeddings Guard</u>.

Arize's ArizeDatasetEmbeddings Guard

The ArizeDatasetEmbeddings guard follows these steps:

- Start with a set of example attacks and generate embeddings for each.
- When a new prompt arrives, chunk the input and create embeddings for each chunk.
- Calculate the cosine distance between the prompt chunk embeddings and the example embeddings.
- If any distance falls below a user-defined threshold (default: 0.2), the guard intercepts the call.

By default, this guard uses 10 examples from a public jailbreak prompt dataset. However, you can customize it with your own data using the sources={} parameter, allowing you to fine-tune the guard based on attacks specific to your application.

Both few-shot prompting and embedding-based guards require an up-to-date collection of attack prompts. You can either gather these yourself through your application's usage or tap into online repositories. While we won't link directly to them here, platforms like Reddit and Twitter (X) host frequently updated collections of attack prompts that can be valuable resources for training your guards.

Performance

We tested the ArizeDatasetEmbeddings guard against a dataset of jailbreak attempts and regular prompts, and the results speak for themselves. Crucially, the guard keeps false negatives low (shown in the top right quadrant below), preventing missed attacks.



True Positives	86.43% of 656 jailbreak prompts were successfully blocked.		
False Negatives	13.57% of jailbreak prompts slipped through.		
False Positives	13.95% of 2000 regular prompts were incorrectly flagged.		
True Negatives	86.05% of regular prompts passed correctly.		
Median Latency	1.41 seconds for end-to-end LLM call on GPT-3.5.		

When To Use Dynamic Guards?

While dynamic guards offer powerful protection, they come with a few considerations:

- Increased computational cost due to larger models or embedding generation.
- Higher latency, potentially impacting response times.
- The need for ongoing maintenance and updates to the attack prompt database.

It's important to weigh these factors against the level of protection required for your specific use case.



Conclusion

We hope this series has provided you with a comprehensive understanding of the guardrail landscape for LLMs. From basic content filtering to sophisticated dynamic defenses, you now have the knowledge to implement robust safety measures for your AI applications.

Remember, the field of AI safety is rapidly evolving. Staying informed and continuously adapting your defenses is the only way to maintain secure and responsible AI systems.



arize

To start your LLM observability journey, **sign up for a free account** or **schedule a demo**.

To receive more educational content, <u>Sign up</u> for our bi-monthly newsletter "The Evaluator."

