# AWS re:Inforce

JUNE 10 – 12, 2024 | PHILADELPHIA, PA

GRC325

# Build responsible AI applications with Guardrails for Amazon Bedrock

**Anubhav Mishra**

(he/him)
Principal Product Manager
AWS

aws

# Building generative apps brings new challenges

**Undesirable and irrelevant topics**

Controversial queries and responses

**Toxicity and safety (including brand risk)**

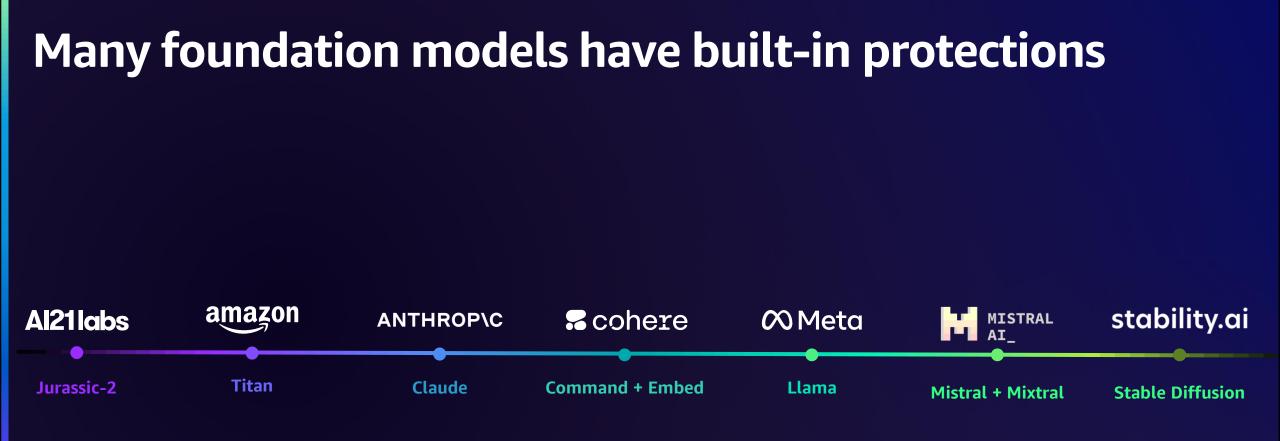Harmful or offensive responses

**Privacy protection**

Protect user information or sensitive data

**Bias/stereotype propagation**

Biased results or unfair user outcomes

# Many foundation models have built-in protections

**AI21labs**      **amazon**      **ANTHROP\C**      **cohere**      **Meta**      **MISTRAL AI_**      **stability.ai**

Jurassic-2      Titan      Claude      Command + Embed      Llama      Mistral + Mixtral      Stable Diffusion

# Building generative AI apps requires additional controls



Customizations based on use cases and organizational policy



Safety and privacy controls for responsible AI



Consistent safeguards across FMs and applications

# Guardrails for Amazon Bedrock

Implement safeguards customized to your application requirements and responsible AI policies

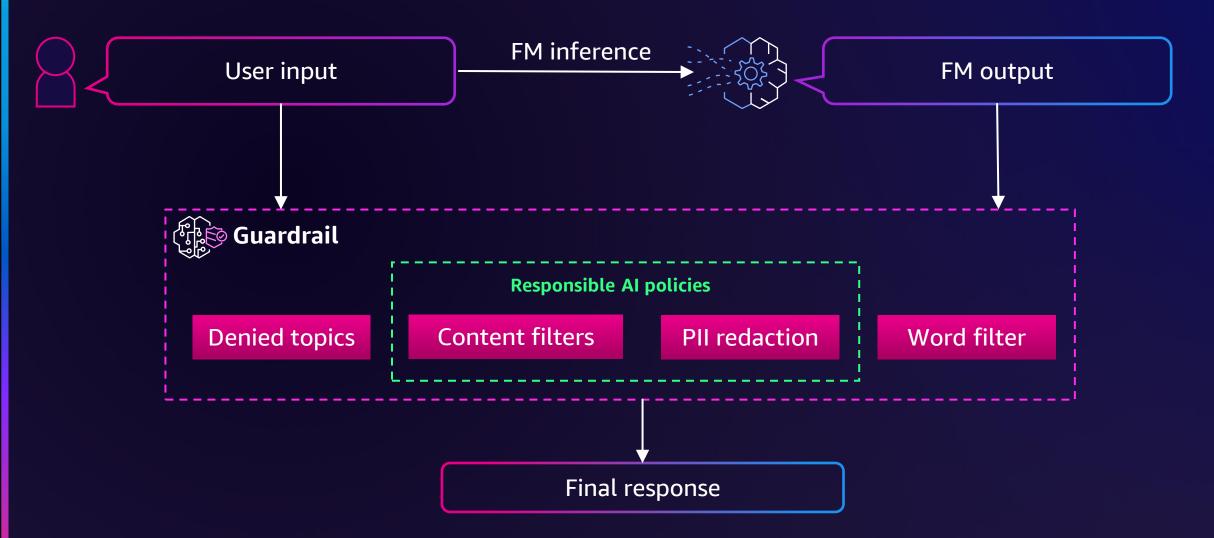Apply guardrails to multiple foundation models, knowledge bases, and agents for Amazon Bedrock

Filter harmful content and safeguard against prompt injection and jailbreaks

Define and disallow denied topics with short natural language descriptions

Redact or block sensitive information, such as PIIs, and custom regex (regular expressions)

# How it works: Guardrails for Amazon Bedrock

User input

FM inference

FM output

**Guardrail**

**Responsible AI policies**

Denied topics

Content filters

PII redaction

Word filter

Final response

# Denied topics

# Content filters

Filter harmful content across categories:

➤ Hate

➤ Insults

➤ Sexual

➤ Violence

➤ Misconduct (criminal activity)

➤ Prompt Attack (jailbreak and prompt injection)

# Sensitive information filter

## PROTECT SENSITIVE INFORMATION AND PRESERVE USER PRIVACY

➢ Redact personally identifiable information (PII) in FM responses to protect user privacy

➢ Detect and filter PIIs in user inputs

➢ Select from a variety of PIIs based on application requirements

➢ Define your own sensitive information using regular expressions (regex)

**Personally Identifiable Information (PII) types** Info
Specify the types of PII to be filtered and the desired guardrail behavior.

**PII types** (1/15)                                            Delete ▼

🔍 Find PII types          | Show all ▼ |        < 1 2 3 4 5 > ⚙

| ☐ | Type ✎ ▽ | Guardrail behavior ✎ ▽ |
|---|---|---|
| ☐ | Name | Mask |
| ☐ | Address | Mask |
| ☐ | Phone number | Mask |

**Add a PII type** ▼

**Regex patterns** Info
Add up to 10 regex patterns to filter custom types of sensitive information and specify the desired guardrail behavior.

**Regex patterns** (0)                        Delete ▼ | Add regex pattern

🔍 Find regex patterns          < 1 2 3 4 5 > ⚙

| Name | Regex pattern | Guardrail behavior | Masking in logs | Description | Actions |
|---|---|---|---|---|---|

No regex patterns added.

Add regex pattern

# Word filters

➤ Filter profane words

➤ Define a set of custom words to block user input and FM responses

## Filter profanity

☑ Filter profanity
Enable this feature to block profane words in user inputs and model responses. The list of words is based on conventional definitions of profanity and is subject to change.

## Add custom words and phrases  Info

Specify up to 10,000 words or phrases (up to 3 words each) to be blocked by the guardrail. A blocked message will show if user input or model responses contain these words or phrases.

◉ Add words and phrases manually
Manually add words and phrases to the following table.

◯ Upload from a local file
Populate the following table with words and phrases from a .txt or .csv file from your computer.

◯ Upload from S3 object
Populate the following table with words and phrases from an S3 object.

## View and edit words and phrases (0)                    Delete all ▼

🔍 Find words and phrases          Show all    ▼          < 1 >  ⚙

Word or phrase ✎                                          Action

No words or phrases added
Upload from file or add manually in the console

Add a word or pharse  ▼

# Guardrails demo

# Guardrails demo

# Agents with guardrail demo

# Agents with guardrail demo

# Sensitive information filter demo

# Sensitive information filter demo