aws

# Generative AI in 2024
## AI Engineering

Mike Chambers (he/him)

Snr. Developer Advocate Specialist
Amazon Web Services

👋 hello

# Where are we now?

## Image



Generated with Amazon Titan Image Generator G1

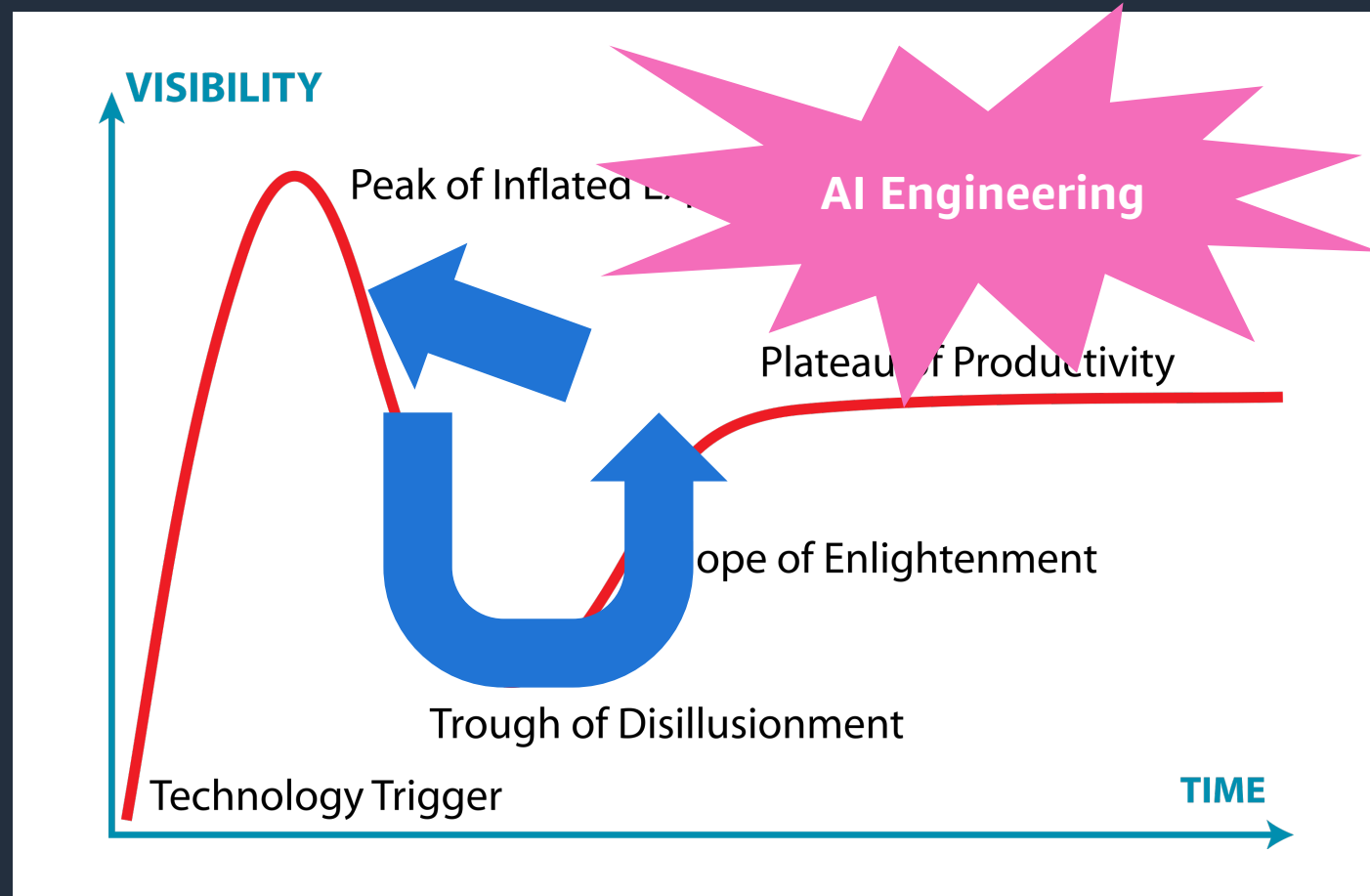"A happy golden retriever wearing glasses and a hat in a portrait painting"

## Text/Chat

> Question: What is generative AI?

Generative AI refers to a type of artificial intelligence that is capable of creating new and original content. It uses algorithms and machine learning to generate text, images, music, and other forms of media based on patterns and rules learned from existing data.

Generated with Amazon Titan Text G1 - Express

# Where are we now?



Gartner Research's Hype Cycle diagram

Jeremykemp at English Wikipedia

# AI Engineering

- Foundation Models

- Prompting engineering (yes it's still important)

- Agents

- Retrieval Augmented Generation

- Fine-tuning and Continuation training

# Foundation Models

**Jurassic-2 series**
By AI21 Labs

**Titan**
By Amazon

**Claude**
By Anthropic

**Command**
By Cohere

**Llama 2**
By Meta

**Stable Diffusion**
By Stability AI

# Prompt Engineering

Text generation

Text summarization

Translation

*oui / yes*

Sentiment analysis

More …

# Prompt Engineering

> Summarise the following text:
> `<text>`
> …
> `</text>`

*In context learning*

# Prompt Engineering

" You are a personal assistant. You are friendly, polite and casual. You help with… "

" You are a classifying agent that filters user inputs into categories. Your job is to sort these inputs before they are passed along to our function calling agent. The purpose of our function calling agent is to call functions in order to answer user's questions. "

# Prompt Engineering

> ```
> The user input is
> <question>
> What is the weather in Oslo right now?
> </question>
> ```

# Agents

- Foundation models have a learning cut off.

- LLMs will "hallucinate" if given the chance.

- What do we want from an LLM?

- World knowledge...

*Up to up...*
*Down is down...*
*The sky is blue...*

**What is the weather in Oslo right now?**

Answer: As of 12:00 PM (noon) local time in Oslo, Norway, the current weather conditions are:

* Temperature: 4°C (39°F)
* Humidity: 84%
* Wind: Northwest at 12 km/h (7 mph)
* Conditions: Overcast with light rain showers

**Llama 2 Chat 13B**

# Agents

Give LLM applications agents.

agency
/ˈeɪdʒ(ə)nsi/

2. action or intervention producing a particular effect.

# Agents

# Agents

## What *is* the weather in Oslo right now?
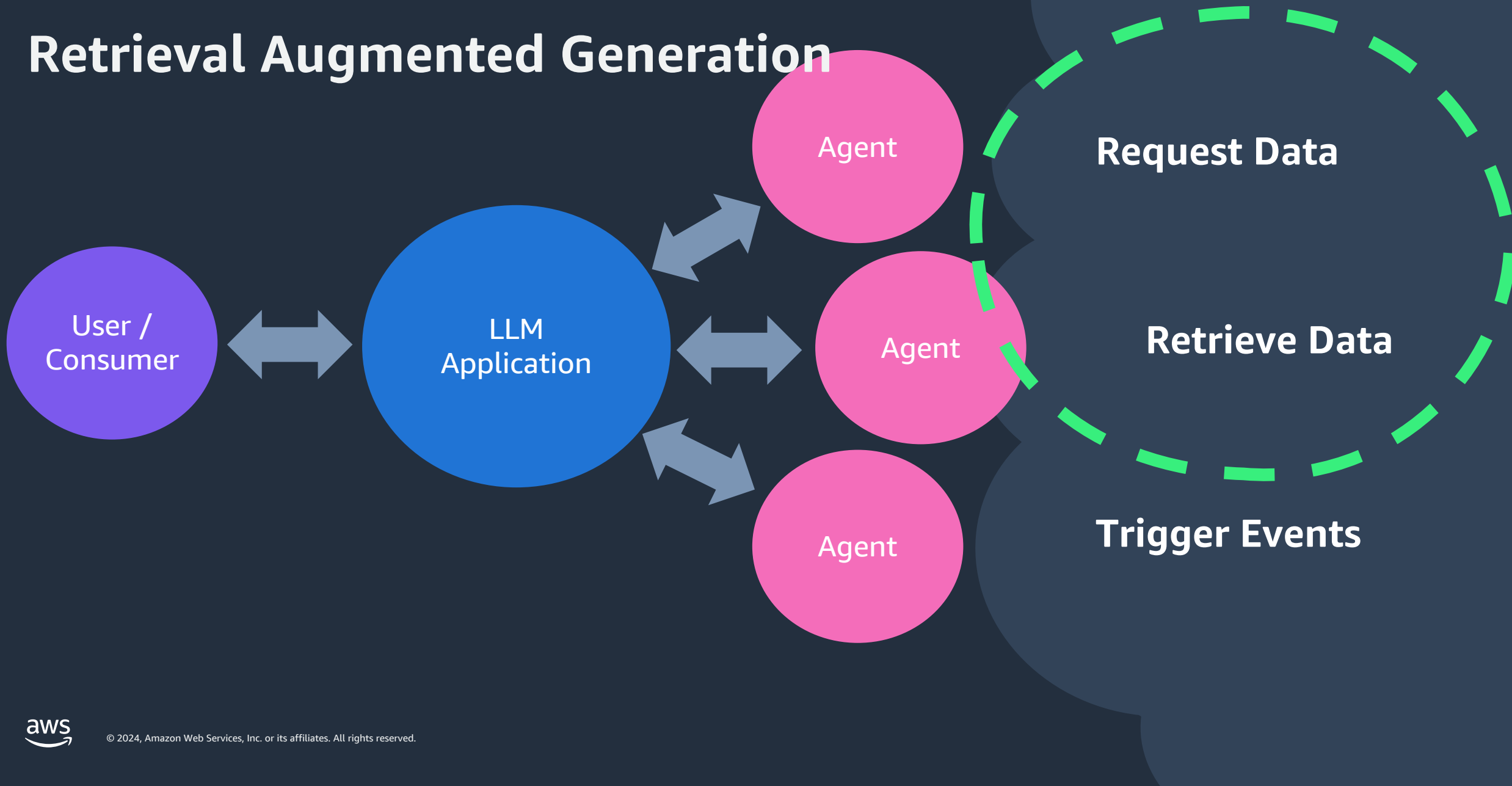
Demo time...

# Agents

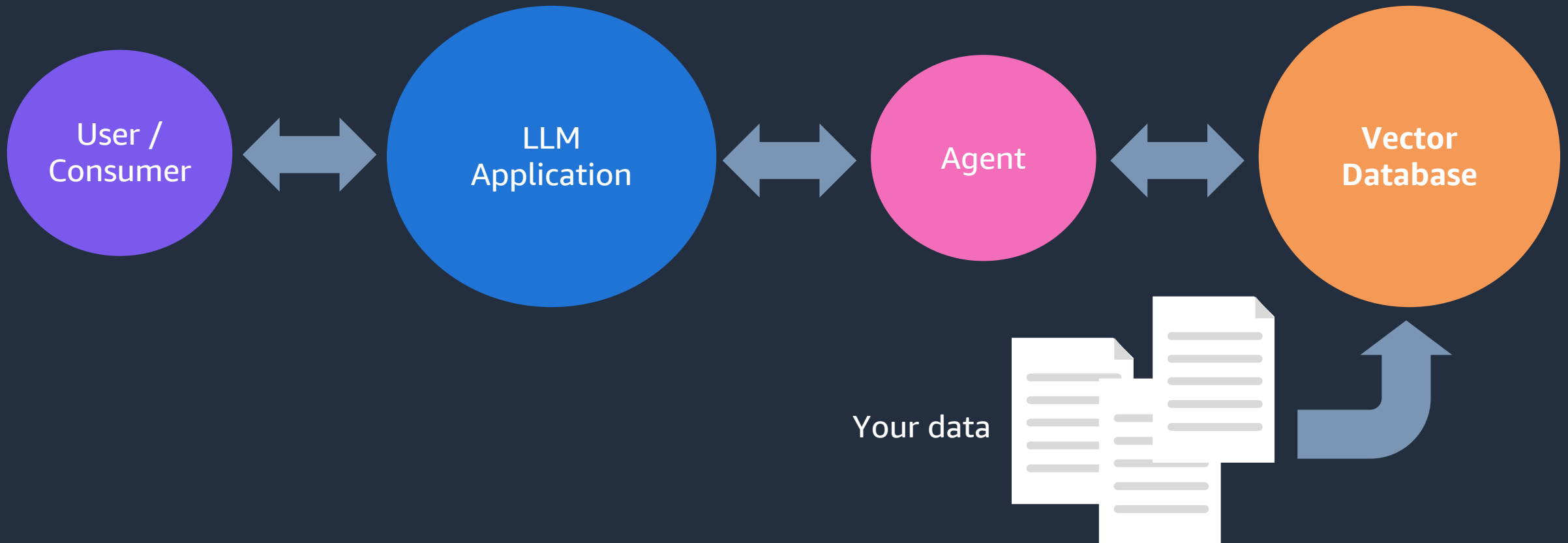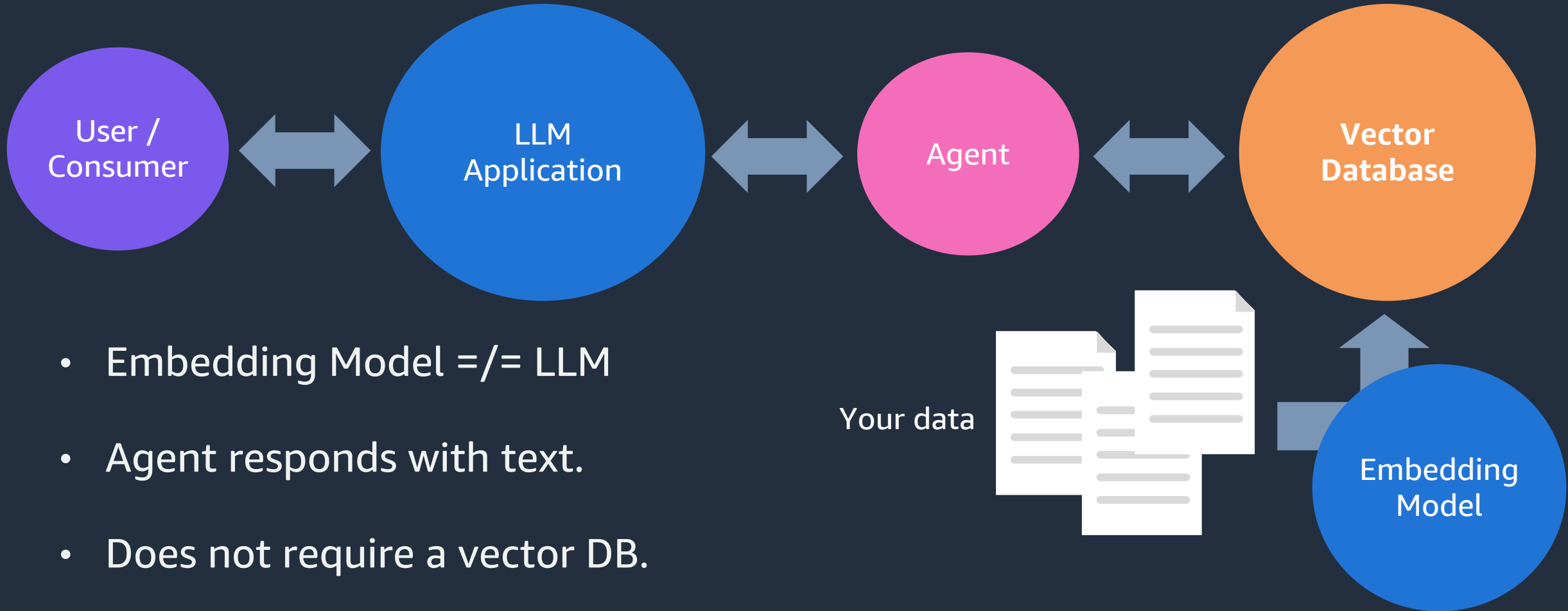**What *is* the weather in Oslo right now?**



(Preview)

# Retrieval Augmented Generation

# Retrieval Augmented Generation

# Retrieval Augmented Generation



- Embedding Model =/= LLM

- Agent responds with text.

- Does not require a vector DB.

# Retrieval Augmented Generation

# Fine-tuning and Continuation training

"I have experimented with different models and different prompts, and I am not getting the generation I want."

"Gather some examples and provide them to the model. You can do this in the prompt with in-context-learning… Or if you have enough, you can **fine-tune** the model."

# Fine-tuning and Continuation training

"I work in a specialized field, and the model does not seem to understand the terminology we use."

"Gather as much domain specific example text that you can, and perform **continuation training** on the model."
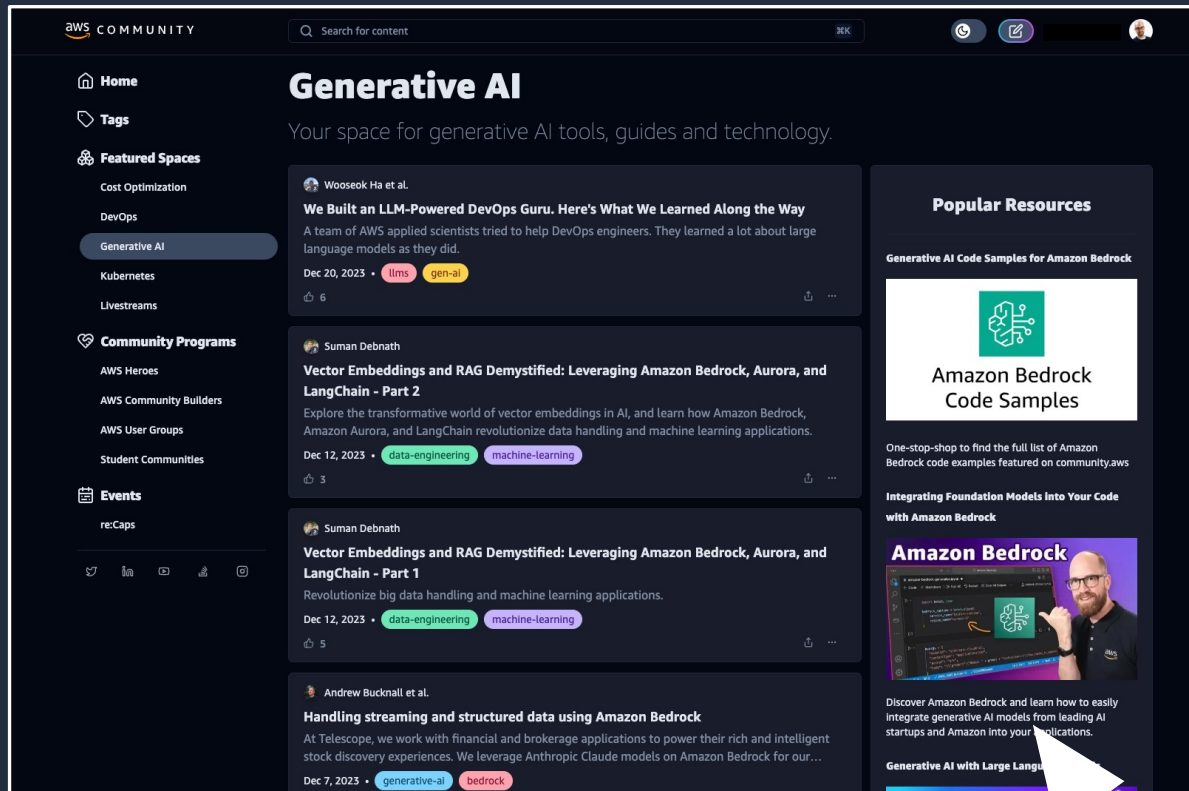
# Questions



- Please tell us what you think

- Your generative ai journey

- What should we do next?

# Find out more



**aws COMMUNITY**

# Generative AI

Your space for generative AI tools, guides and technology.

https://community.aws/generative-ai

👋 thanks

linkedin.com/in/mikegchambers