aws

ADC – D4

# Leveraging pgvector and Amazon Aurora PostgreSQL for natural language processing, chatbots, and sentiment analysis

Divya Sharma

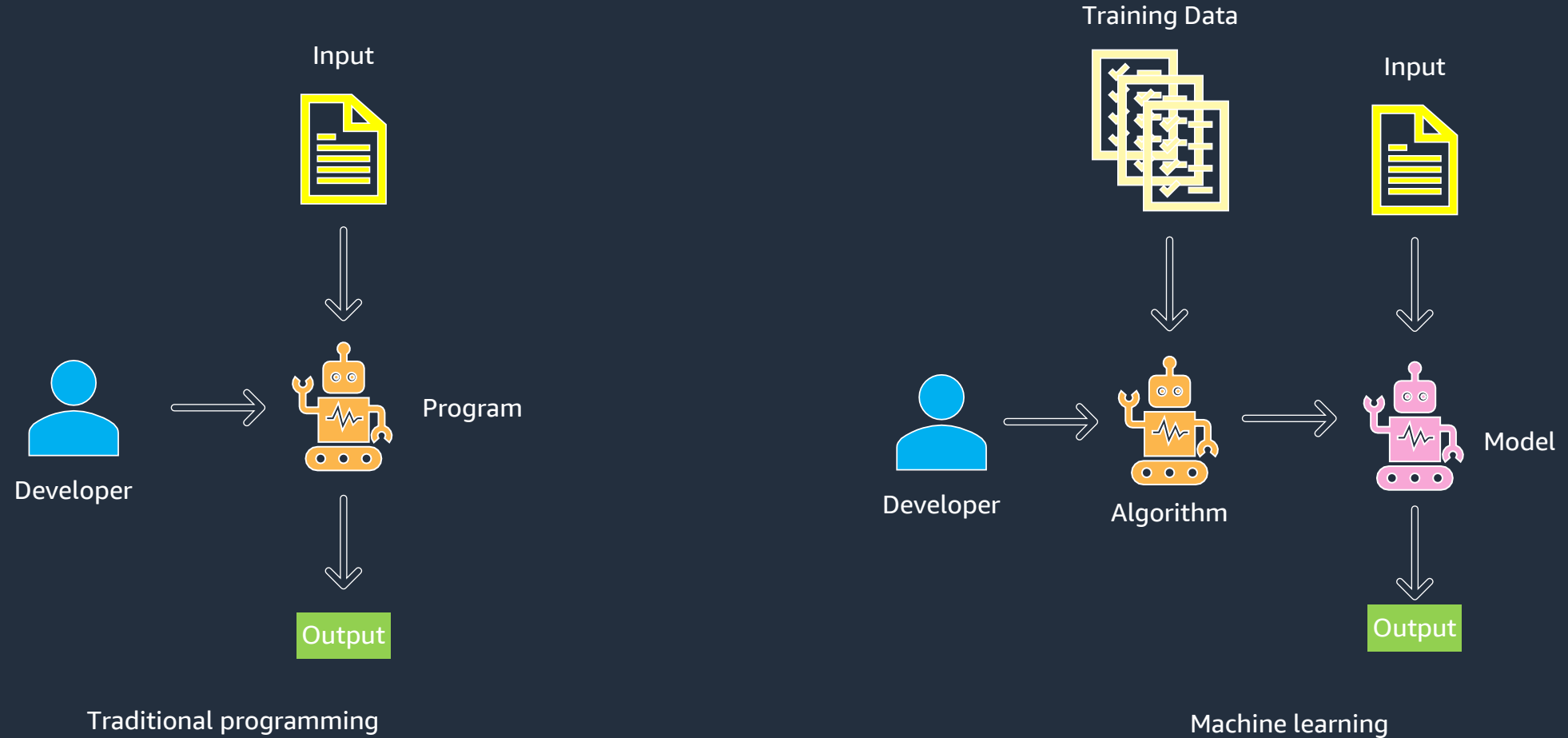Sr. RDS PostgreSQL Solutions Architect

# Agenda

➢ Machine Learning and Generative AI

➢ Vector Embeddings

➢ Vector Database

➢ Using pgvector with RDS and Aurora PostgreSQL

➢ Workshop - Generative AI Use Cases with Aurora PostgreSQL and pgvector
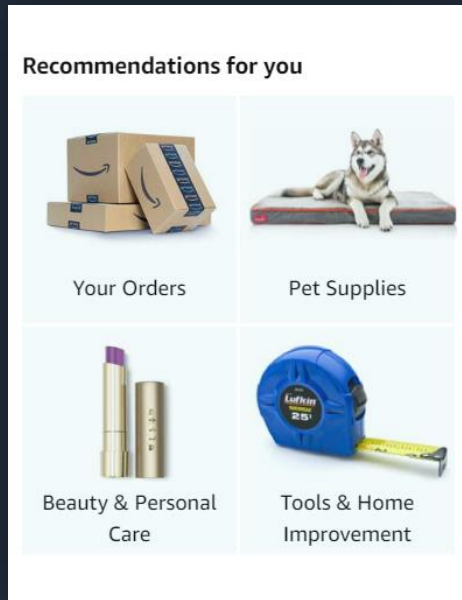
# Machine Learning and Generative AI

# Traditional programming vs. machine learning

Input

Program

Developer

Output

Traditional programming

Training Data

Input

Developer

Algorithm

Model

Output

Machine learning

# ML innovation is in the Amazon DNA



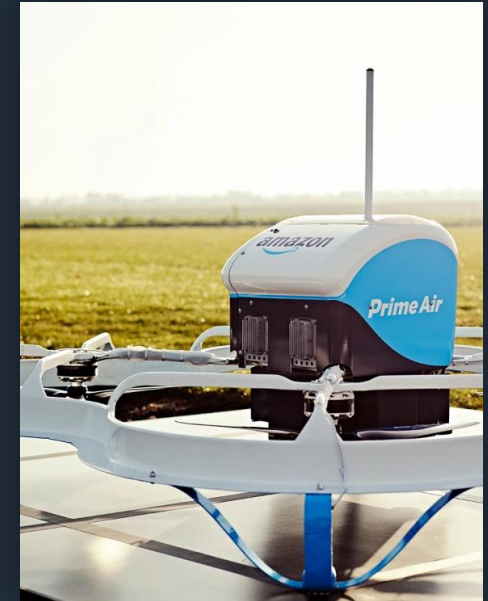**4,000 products per minute** sold on Amazon.com



**1.6M packages** every day



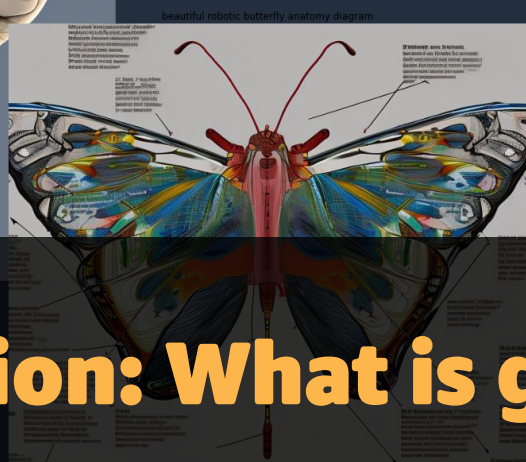**Billions** of Alexa interactions each week



First Prime Air delivery on **December 7, 2016**

# Question: What is generative artificial intelligence (AI)?

- Creates new content and ideas, including conversations, stories, images, videos, and music

- Powered by large models that are pretrained on vast corpora of data and commonly referred to as foundation models (FMs)
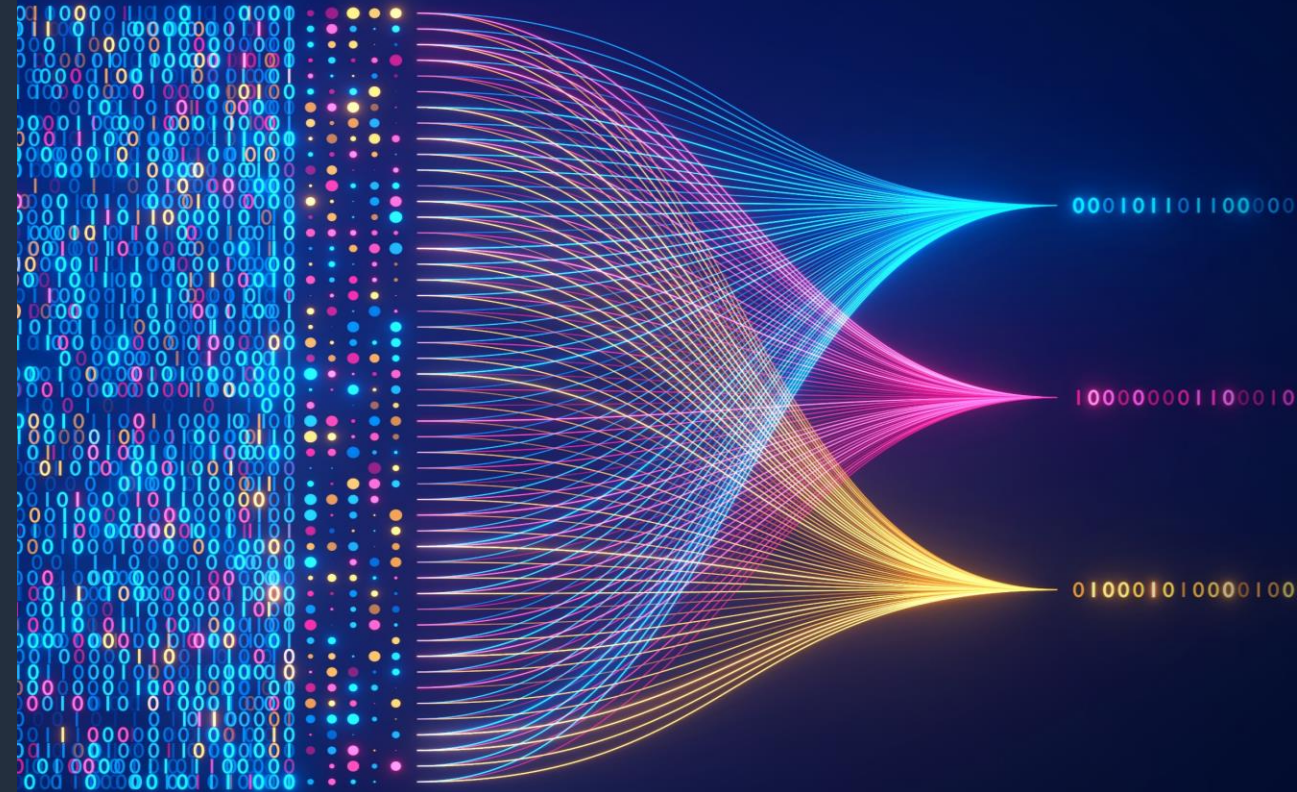
# Generative AI is powered by foundation models

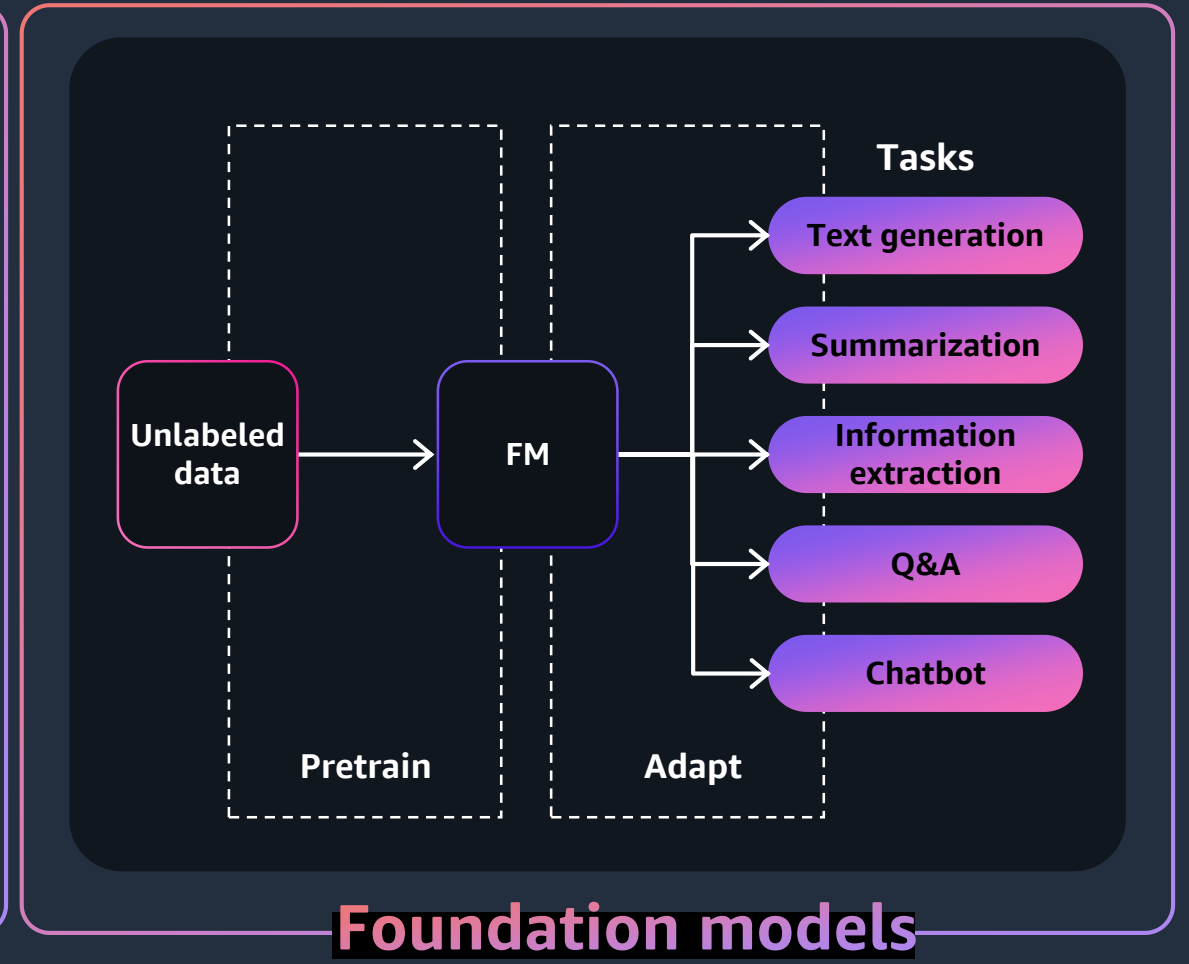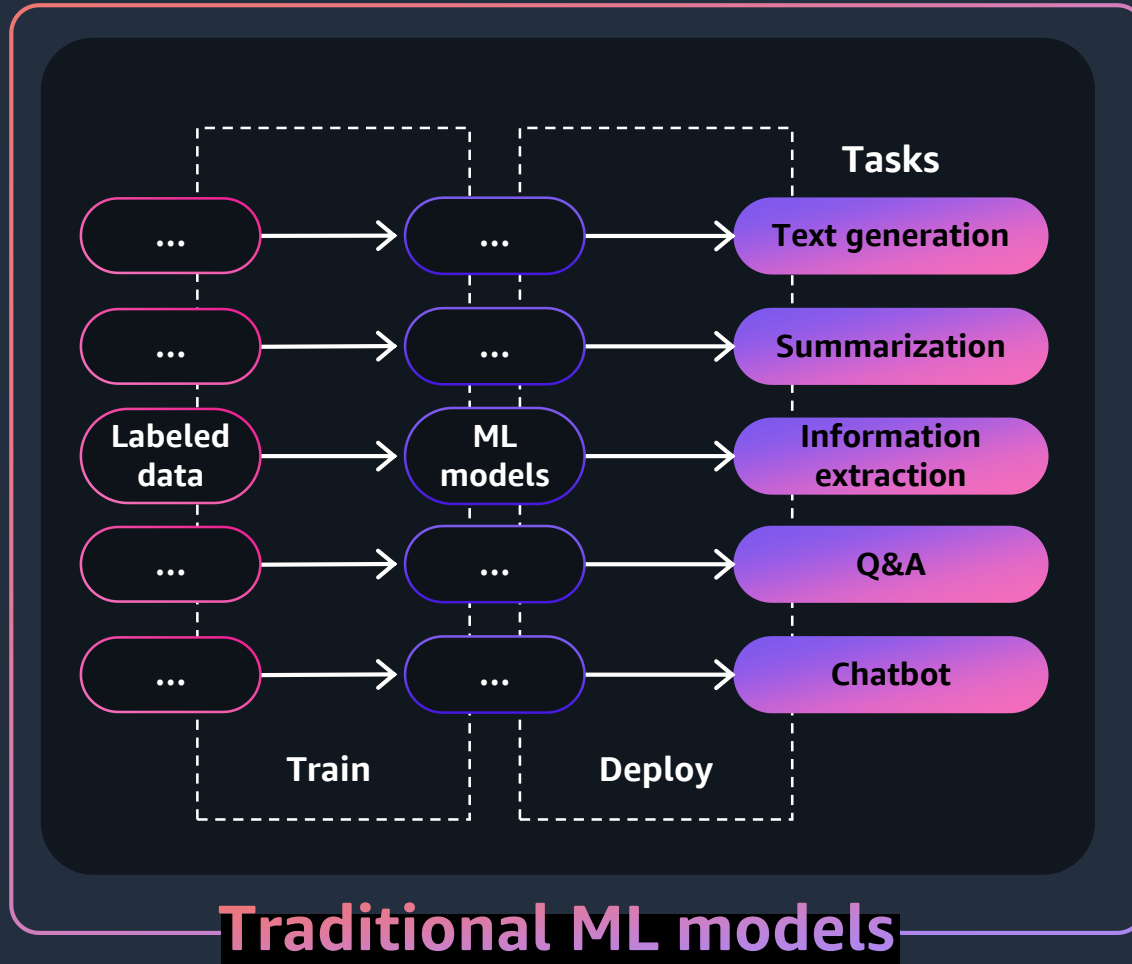Pretrained on vast amounts of unstructured data

---

Contain large number of parameters that make them capable of learning complex concepts

---

Can be applied in a wide range of contexts

---

Customize FMs using your data for domain specific tasks

# How foundation models differ from other ML models



**Traditional ML models**

Tasks
- Text generation
- Summarization
- Information extraction
- Q&A
- Chatbot

... → ... → Text generation
... → ... → Summarization
Labeled data → ML models → Information extraction
... → ... → Q&A
... → ... → Chatbot

Train | Deploy

**Foundation models**

Tasks
- Text generation
- Summarization
- Information extraction
- Q&A
- Chatbot

Unlabeled data → FM → Text generation / Summarization / Information extraction / Q&A / Chatbot

Pretrain | Adapt

# Generative AI can be used for a wide range of use cases

Chatbots &
Virtual assistants

Agent Assist

Contact Center
Analytics

Personalization

**Enhance
customer
experience**

Conversational search

Content Localization

Text, image,
video generation

Text summarization

Code generation

**Boost
employee
productivity**

Document processing

Content moderation

Synthetic data creation

Maintenance assistance

Anomaly detection

**Improve
business
operations**

Image generation
for web pages

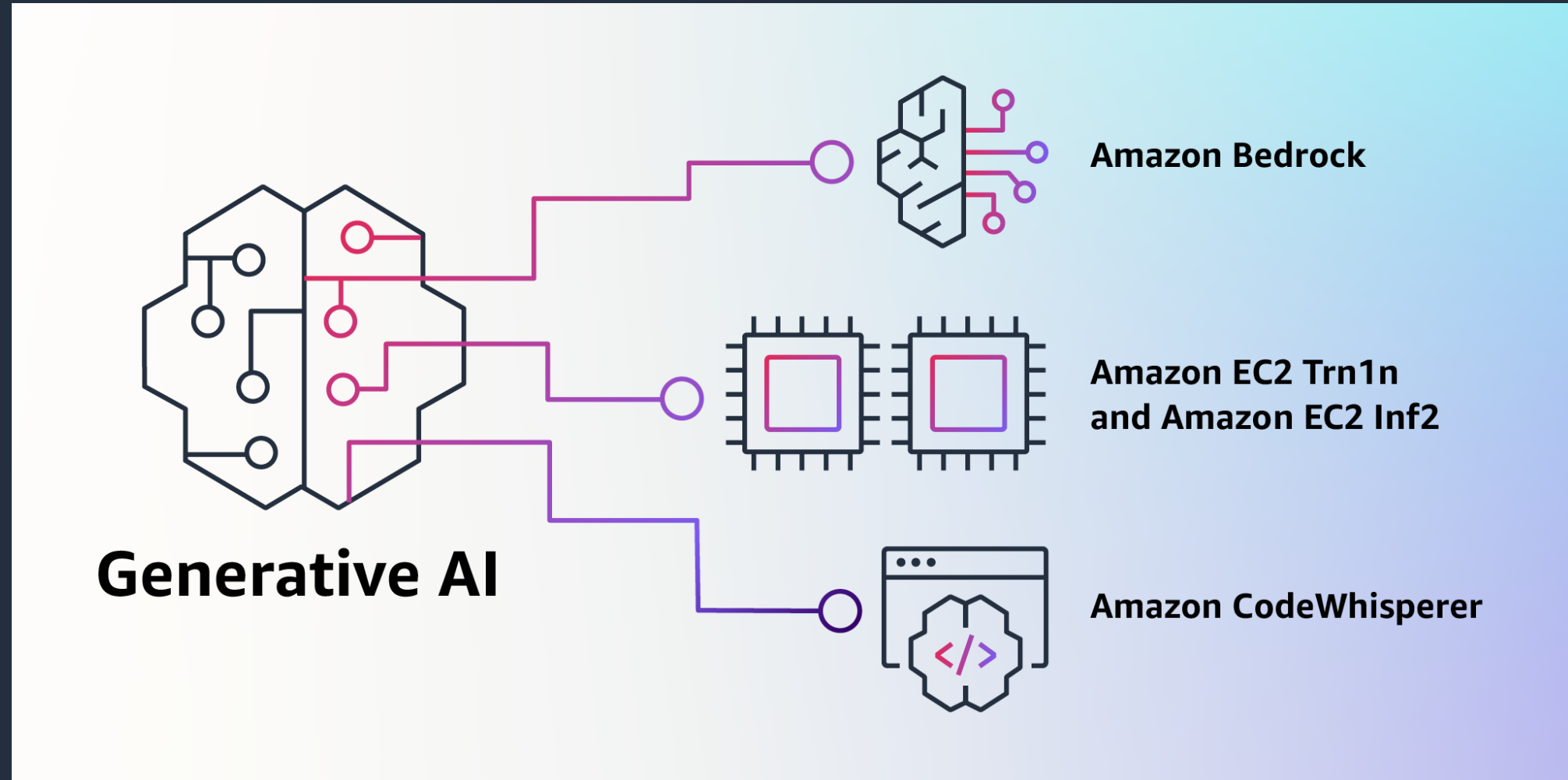Video enhancement

Music creation

Image enhancement

Creating animations

**Creativity**

# Building with generative AI on AWS



Generative AI

Amazon Bedrock

Amazon EC2 Trn1n
and Amazon EC2 Inf2

Amazon CodeWhisperer

# Vector Embeddings

Magnitude

Direction

# Notion of vectors in search

*A vector is a list of values describing some attributes of an item.*

$v_1$ = [5, 4, 854, $1.1M]

$v_3$ = [6, 4, 530, $2.1M]

1. How many bedrooms?

2. How many bathrooms?

3. Size of the house in sqm?

4. Price of the house?

$v_2$ = [4, 3, 335, $530K]

$v_4$ = [6, 4, 500, $2.5M]

# Similarity search using vectors

$v_1 = [5, 4, 854, \$1.1M]$

0.96

0.92

$v_2 = [4, 3, 335, \$530K]$

0.90

0.87

0.90

$v_3 = [6, 4, 530, \$2.1M]$

0.99

$v_4 = [6, 4, 500, \$2.5M]$

$$\cos(\theta) = \frac{v_a \cdot vb}{|v_a| \times |vb|}$$

# Notion of vectors in search

A *vector* is a list of values describing some attributes of an item.

$v_1$ = [5, 4, 854, $1.1M]

$v_3$ = [6, 4, 530, $2.1M]

1. How many bedrooms?

2. How many bathrooms?

3. Size of the house in sqm?

4. Price of the house?

$v_2$ = [4, 3, 335, $530K]

$v_4$ = [6, 4, 500, $2.5M]

# What is a vector embedding?

- Numerical representation of words or sentences, used in Natural Language Processing (NLP) to facilitate efficient analysis and manipulation of text

- By converting text into vector embeddings, NLP models can easily perform tasks such as querying, classification, and applying machine learning algorithms on textual data

- Mathematical vector generated to be used in machine-learning tasks

# What is a vector embedding?

New York →
Beijing →
Paris →

Embedding Model

Animal →
Horse →

Airplane →

| 0.027 | -0.011 | ... | -0.023 |
| 0.025 | -0.009 | ... | -0.025 |
| 0.024 | -0.012 | ... | -0.021 |

| -0.011 | 0.021 | ... | 0.013 |
| -0.009 | 0.019 | ... | 0.015 |

| -0.048 | 0.079 | ... | 0.076 |

Text

Text as vector embeddings

# Image Embeddings

Chicken, Wolf, Dog, and Cat

Banana, Apple (fruit) , Apple (corp)

# Vector Database

# What is a Vector database?

Raw
Data

Vector
Embedding
Space

Dev-ready and
Operationalized

Consumable

Image

Machine Learning Model
(Embedding)

Documents

Audio

*Retrieve content most similar to some content: question context, image, music clip...*

Dense Vector Encodings

0.35 0.1 0 0.9 001.0 00 0001.0 0 0...

0.35 0.1 0 0.8 001.0 00 0001.0 0 0...

0.15 0.1 0 0.7 001.0 00 0001.0 0 0...

Sparse Vector Encodings
(Automatic metadata extraction)

Content classification
Salient terms and topics
...

PostgreSQL

Vector Database

Build
AI-powered Application

*Retrieve most relevant content by key terms (metadata)...*

# Visual Search

# Retrieval Augmented Generation (RAG)

# Using pgvector with RDS and Aurora PostgreSQL

# pgvector: An open-source library for vector search

- An open-source extension for PostgreSQL to store/search vectors
- Provides IVF FLAT indexing for L2 distance, inner products, and cosine similarities
- Also provides HSNW indexing for vectors
- Developer focused: works with existing client libraries
- Currently available for Amazon RDS for PostgreSQL and Amazon Aurora PostgreSQL-Compatible edition

```sql
CREATE EXTENSION vector;

SELECT typname FROM pg_type WHERE
typname = 'vector';

typname
---------
vector
```

# Why use PostgreSQL for vector searches?

- ➤ Open source

- ➤ Integrated solution: keep your relational data and vector data in one place

- ➤ Enterprise-level robustness and operations

- ➤ Full-featured SQL

- ➤ Scalability and performance

- ➤ Existing client libraries work without modification

- ➤ Both PostgreSQL-native data types have current limitations for modern AI/ML workloads

  - ➤ ARRAY – does not support indexing for "nearest-neighbor" queries

  - ➤ cube – limited to 100 dimensions

# pgvector example: Querying nearest neighbor

- Supports exact and approximate nearest neighbor (ANN) search

    - L2 distance <->

    - Inner product <#>

    - Cosine distance <=>

```
CREATE TABLE test_embeddings(product_id bigint, embeddings vector(3) );
INSERT INTO test_embeddings VALUES
(1, '[1, 2, 3]'), (2, '[2, 3, 4]'), (3, '[7, 6, 8]'), (4, '[8, 6, 9]');

SELECT product_id, embeddings, embeddings <-> '[3,1,2]' AS distance
FROM test_embeddings ORDER BY embeddings <-> '[3,1,2]' limit 2;

product_id | embeddings |      distance
------------+------------+--------------------
         1 | [1,2,3]    | 2.449489742783178
         2 | [2,3,4]    |                  3
(2 rows)
```

# Indexing for vectors

- By default, pgvector performs exact nearest neighbor search, which provides perfect recall.

- You can add an index to use approximate nearest neighbor search, which trades some recall for speed. Unlike typical indexes, you will see different results for queries after adding an approximate index.

- Supported index types are:

  - IVFFlat – Inverted File Flat

  - HNSW - Hierarchical Navigable Small World - added in 0.5.0

# Workshop – Generative AI Use Cases with Aurora PostgreSQL and pgvector

# Workshop outline

1: Prerequisites - Attending an AWS event

2: Product Recommendations

3: Document QnA Chatbot using RAG

4: Similarity Search & Sentiment Analysis

5: Cleanup

# Vector Database Overview



Simply store, search, index, and query ML embeddings

Store

Search

Index

Query

Amazon Aurora PostgreSQL

# Step 1: Sign in via your preferred method

https://catalog.workshops.aws/join

# Step 2: Enter event access code

Enter 12-digit event access code: <mark>d8f3-03a520-5c</mark>

# Step 3: Review terms and join event

# Step 4: Access AWS account

Access the AWS Management Console, or generate AWS CLI credentials as needed

# Workshop Office Hours : 28th Sept & 29 Sept

- Objective : Ask questions and clarify your doubts

- Meeting Date & Time : 28th Sept – 12PM – 1PM

  - Join the Meeting : https://chime.aws/5150551774

- Meeting Date & Time : 29th Sept – 11AM – 12PM

  - Join the Meeting : https://chime.aws/8395367338

# Additional Resources

# Resources

- [Amazon Bedrock](#)

- [Build with Generative AI on AWS](#)

- [Building AI-powered search in PostgreSQL using Amazon SageMaker and pgvector](#)

- [Leverage pgvector and Amazon Aurora PostgreSQL for Natural Language Processing, Chatbots and Sentiment Analysis](#)