

The background features a dark purple gradient on the left and a vibrant, multi-colored geometric design on the right. The design consists of overlapping triangles and quadrilaterals in shades of magenta, blue, and orange, separated by thin white lines.

AWS re:Invent

DECEMBER 1 - 5, 2025 | LAS VEGAS, NV



Session: ANT213

Build GPU-boostered, auto-optimized billion-scale VectorDBs in hours

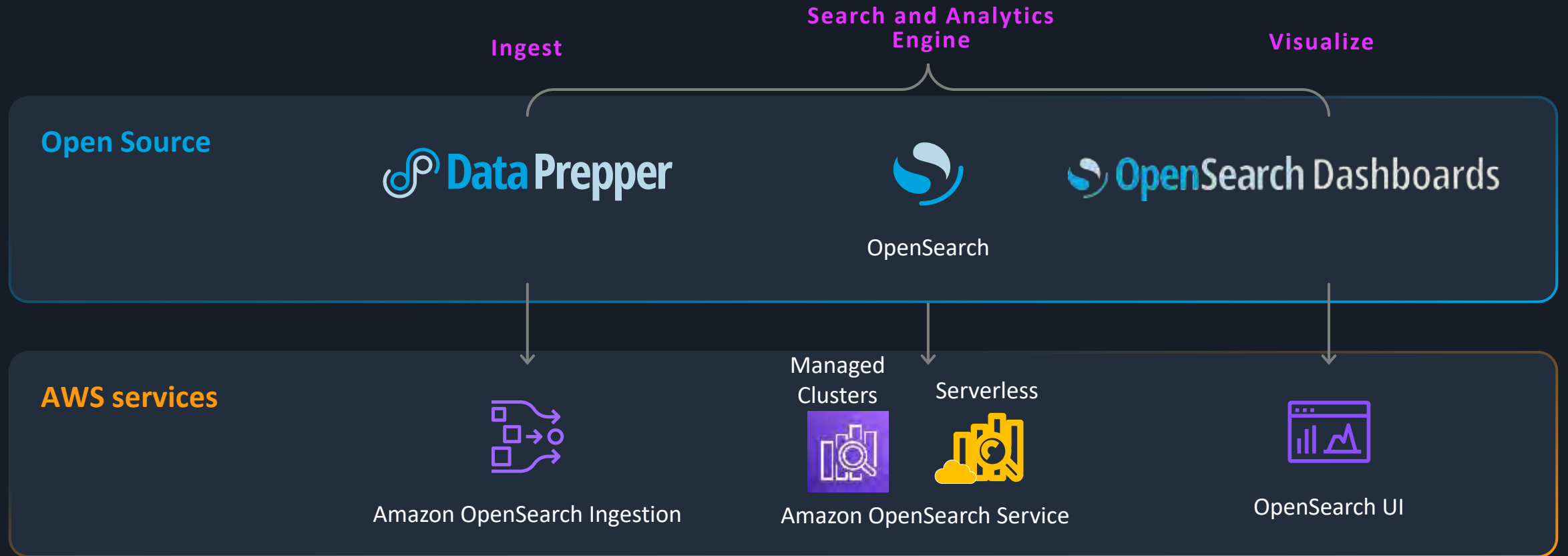
Dylan Tong

Product Lead, AI and vectors
Amazon OpenSearch Service
AWS

Vamshi Nakkirtha

Senior Manager, Software Dev, Vector Search
Amazon OpenSearch Service
AWS

OpenSearch: AWS and Open Source



Use Cases

**Log Analytics,
Observability**

**Search:
Keyword, Vector, Hybrid**

Why Vector Search?

1

Improves search quality
(relevance)

2

Versatile support for
content types

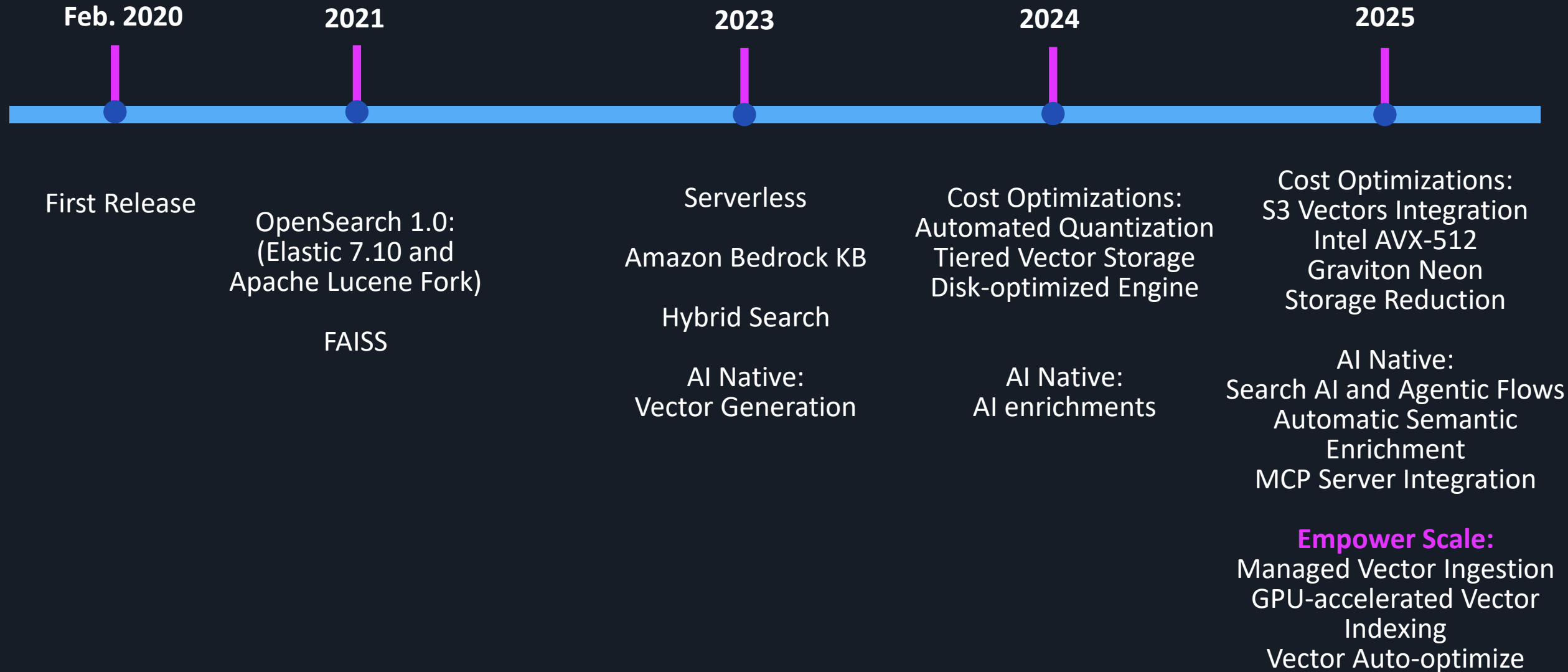
3

Diverse
Use Cases

4

Agentic (AI) App
Enabler

Vector Search Journey



Customers trending to billion-scale and beyond



IP Infringement Detection

8 billion attempted changes to product detail pages for signs of potential abuse (2022)

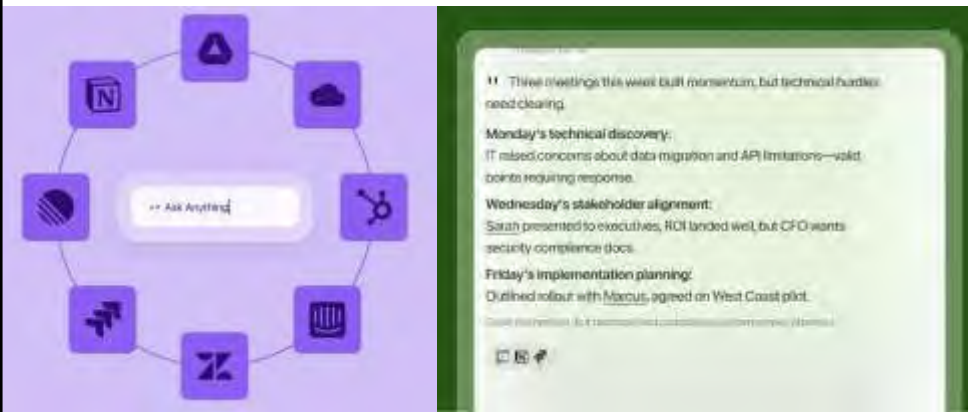
68 Billion vectors encoded from product information indexed into OpenSearch to power vector search.

99% of discovered infringements were automatically found or blocked through proactive controls



Source <https://aws.amazon.com/blogs/big-data/amazon-opensearch-services-vector-database-capabilities-explained/>

Agentic Application



AI teammate that **powers human and AI collaboration.**

Problem:

Silo data and tools across support, product and sales teams

Solution:

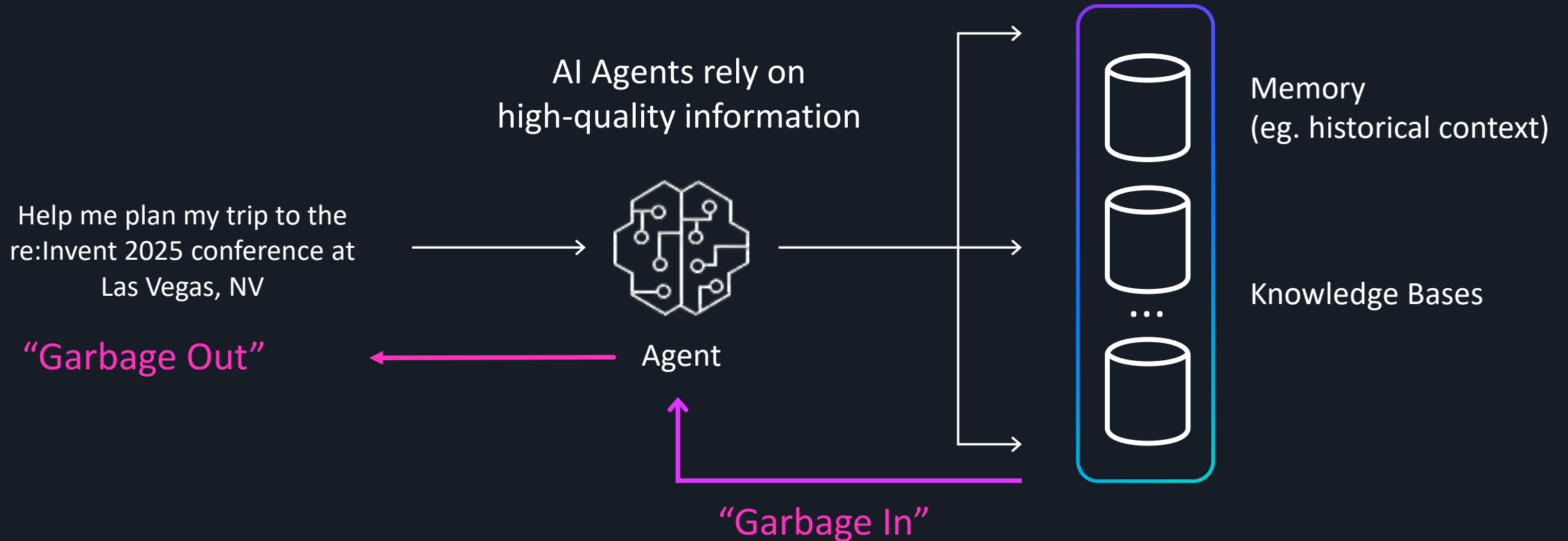
Unify, automate, search and insights via AI agents.

Hundreds millions of vectors with ~1M vector updates daily.

Results:

85% of tickets resolved with no human intervention,
50% cost reduction in customer support
10-hours saved per employee every week.

Agentic systems rely on high-quality search



Vector search and hybrid variations are state-of-the-art.

Keyword vs Semantic (Vector)

Experimental Feature

Compare results using the same search text with different queries. For more information, see the [Compare Search Results Documentation](#). To leave feedback, visit [forums.opensearch.com](#).

Search

Query 1

Index: flickr2_demo

Query

```
1- {
2-   "query": {
3-     "match": {
4-       "image_description_text": "%SearchText%"
5-     }
6-   },
7-   "_source": ["image_id", "image_description_text"]
8- }
9
```

Enter a query in OpenSearch Query DSL. Use %SearchText% to refer to the text in the search bar.

Query 2

Index: flickr2_demo

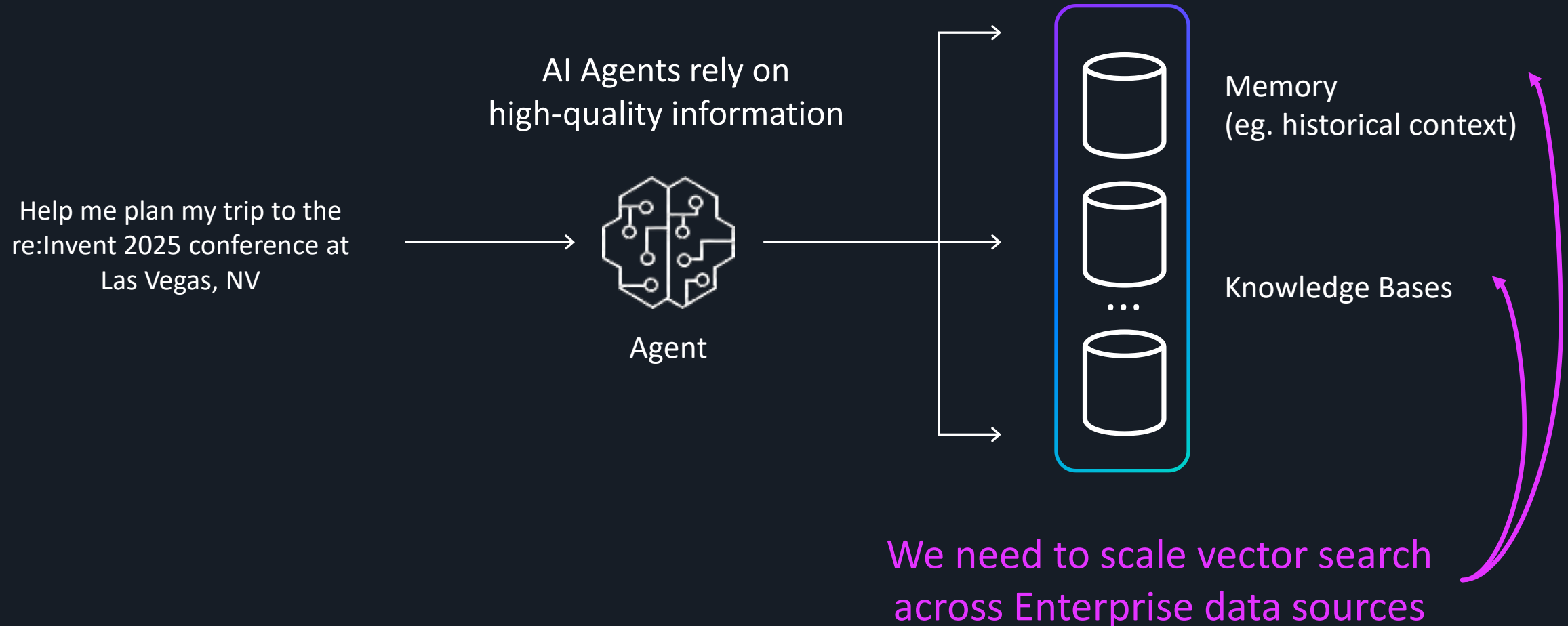
Query

```
1- {
2-   "script": {
3-     "source": "ctx._score = 1"
4-   },
5-   "image_description_embedding_custom": {
6-     "query_text": "%SearchText%",
7-     "model_id": "f0NbeYQ8J_M6WKhH8gc"
8-   }
9- },
10- "_source": ["image_description_text", "image_id"]
11- }
```

Enter a query in OpenSearch Query DSL. Use %SearchText% to refer to the text in the search bar.

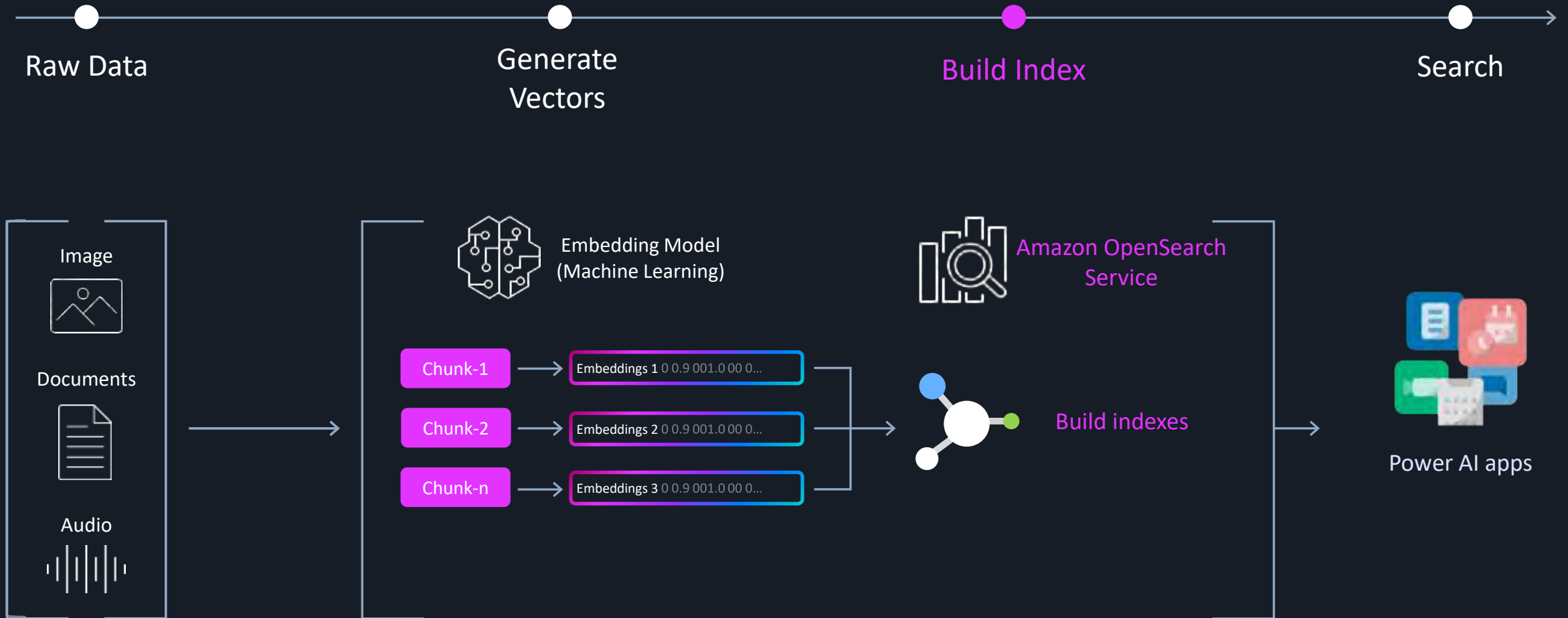
Add queries to compare search results.

Agentic systems need vector search across vast knowledge bases...

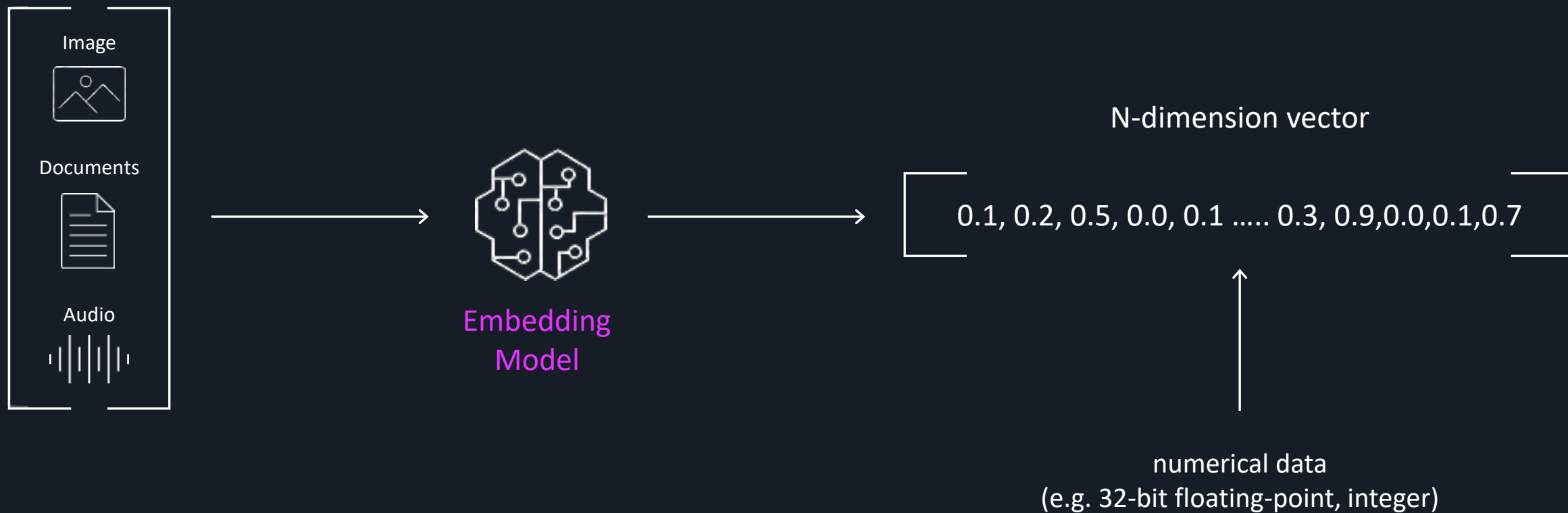


Big Vector Databases, Big Challenges

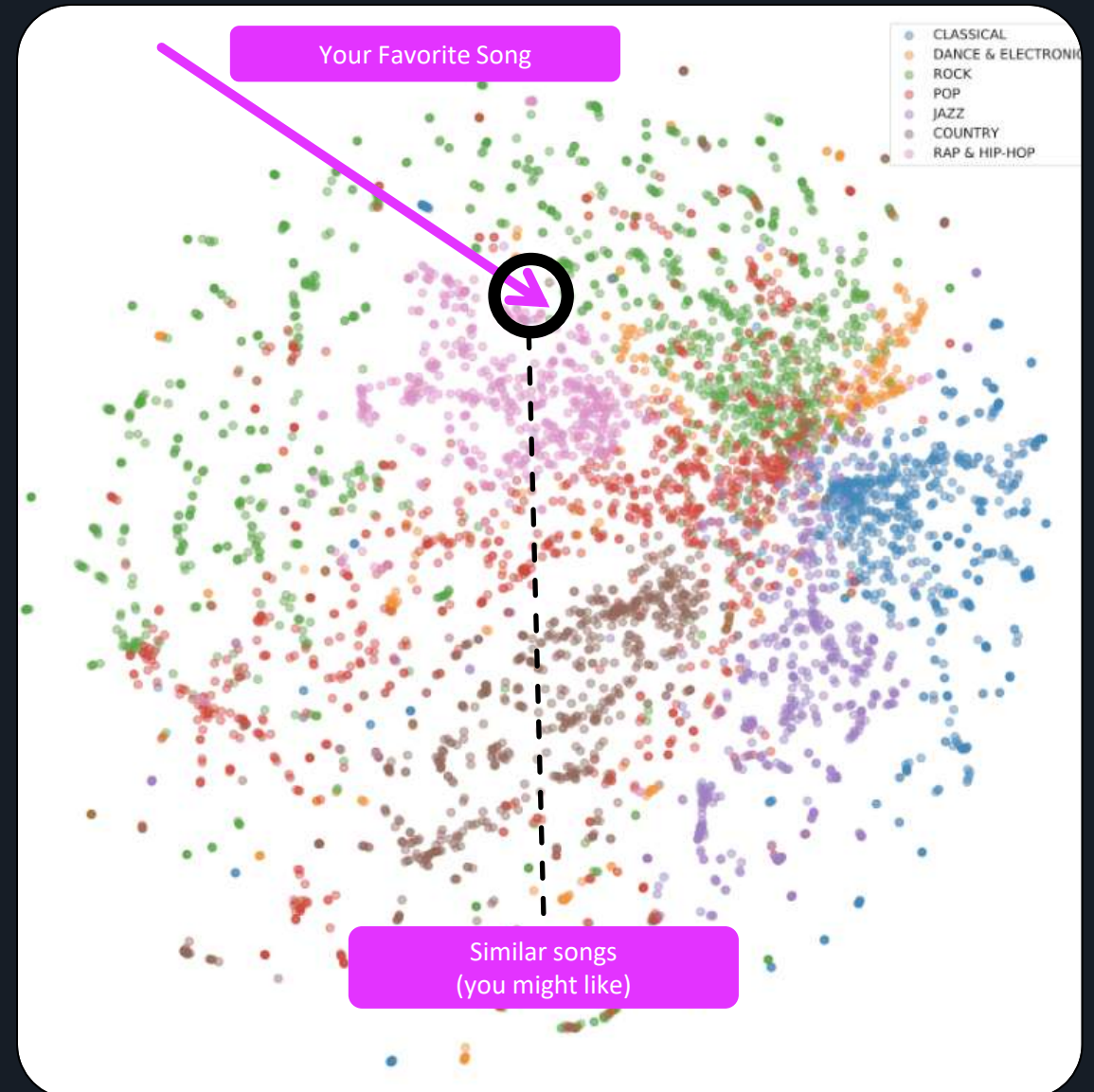
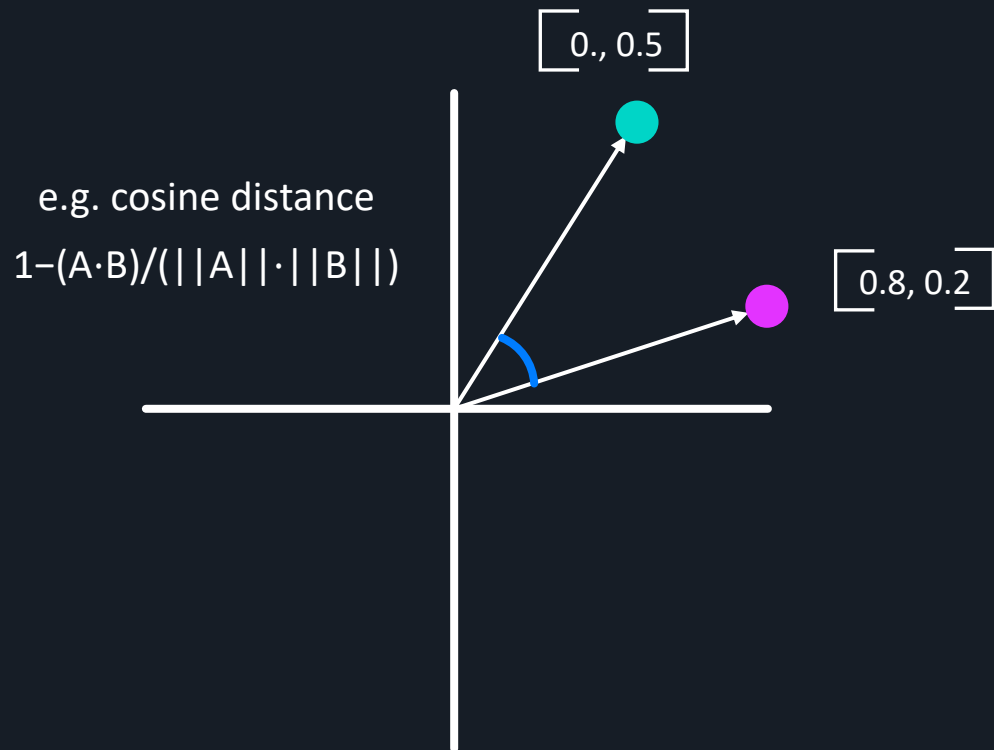
Building a Vector Database



Vectors



Content (Vector) Similarity



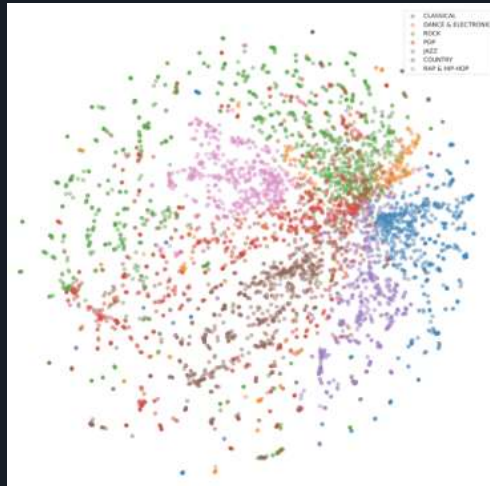
k-Nearest-Neighbors: Find the “Top K” most similar...

Exact (Brute-force) k-NN

query vector: $\begin{bmatrix} 0.1, 0.2, 0.6 \dots 0.3, 0.7 \end{bmatrix}$
(e.g. your favorite song)

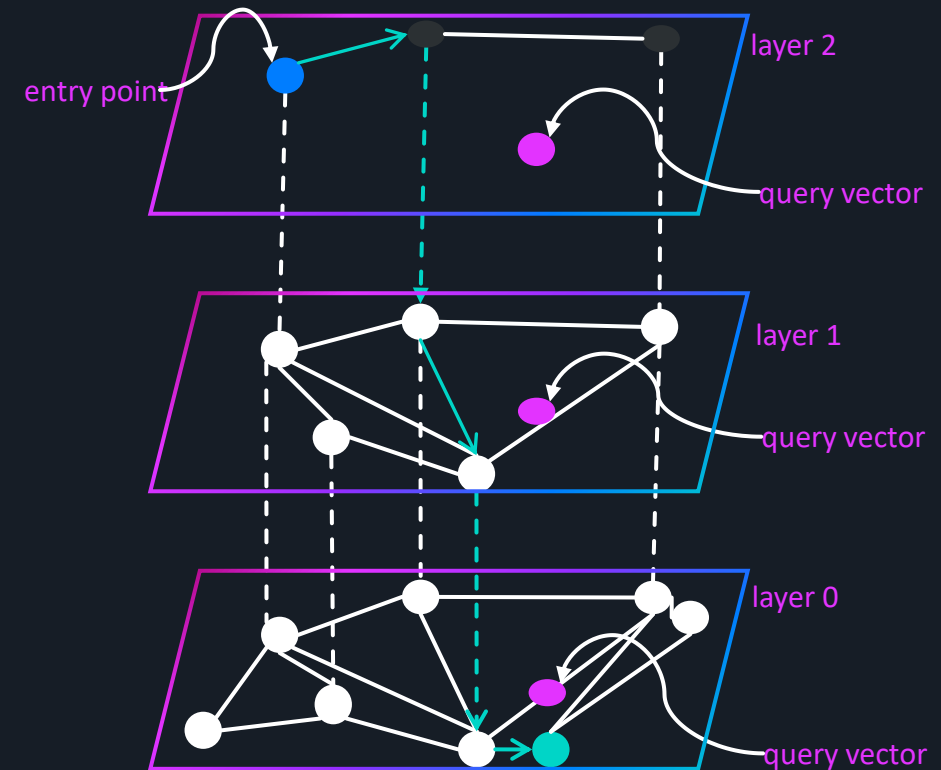
Calculate distance and rank

Song
vector
corpus



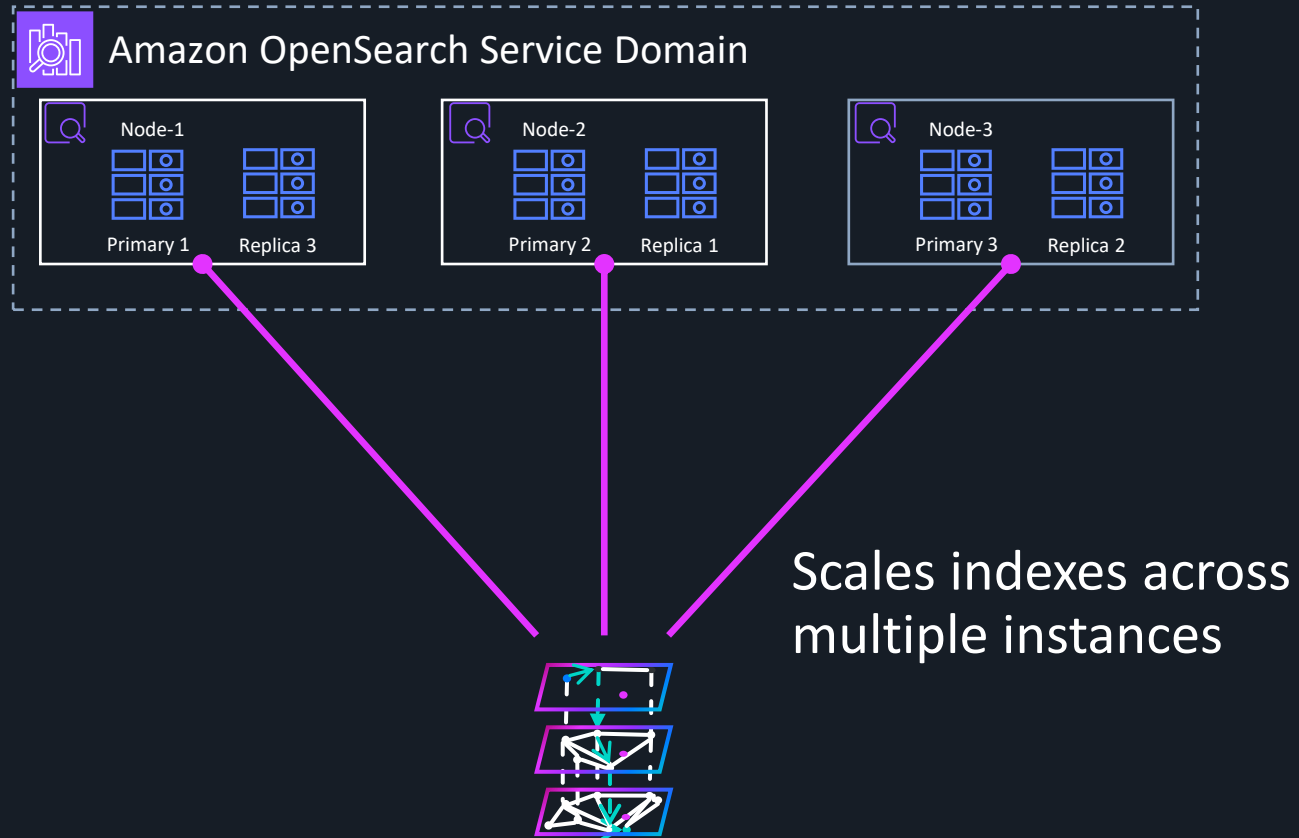
Approximated k-NN

Hierarchical Navigable Small Worlds (HNSW)

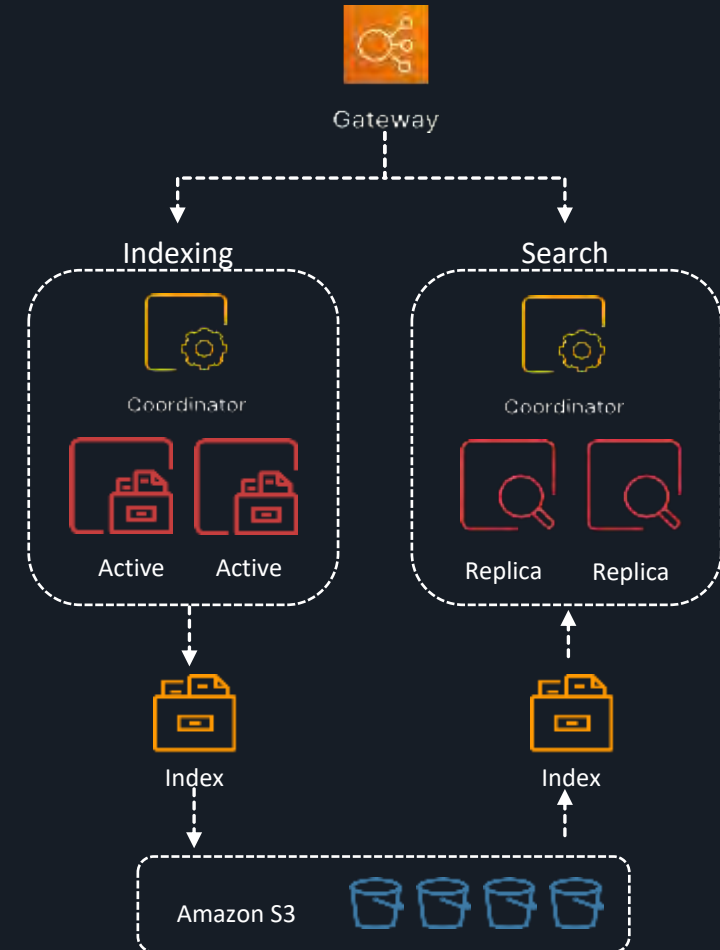


Scaling Vector Search

Managed Clusters



Serverless Collections



How long does it take to index 1-Billion?

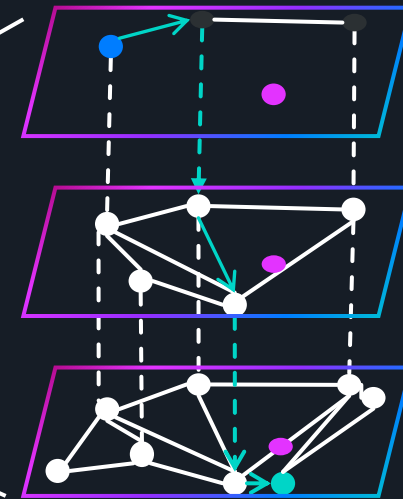
Typically, **Days**



1-Billion Vectors



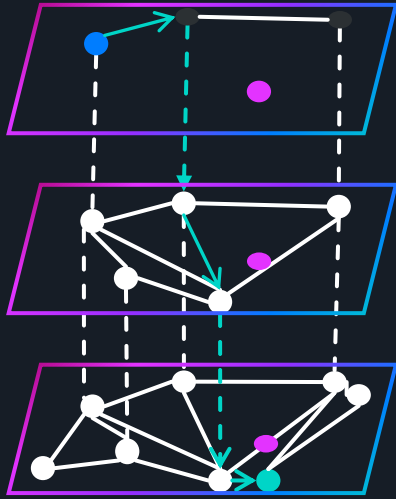
OpenSearch Service
Vector Database



Build indexes

Build Index

Life-cycle of a vector index...



Data Changes

**Add, modify and delete
content**

(index (HNSW) recall
degradation)



Model Changes

New provider or version



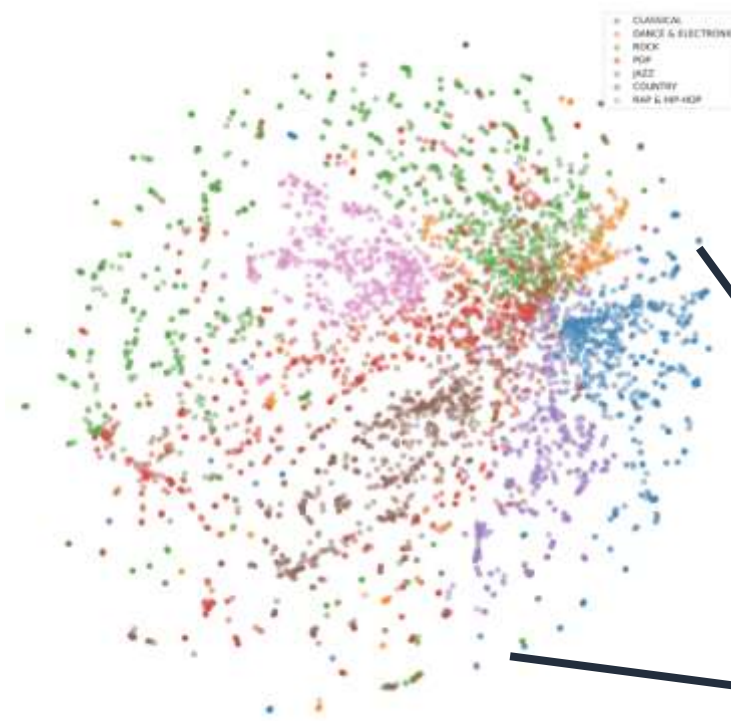
Model Changes

Fine-tuning

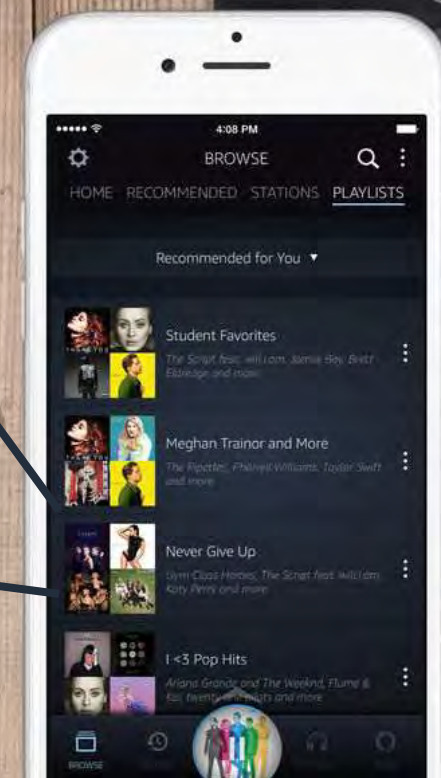
Personalization

amazon music

- **100 Million** music tracks with recommendations based on user listening history.
- **1.05 Billion** vectors indexed into OpenSearch to power item-item collaborative filtering.
- **Daily** model retraining to deliver high-quality recommendations



Get personalized recommendations based on your listening history



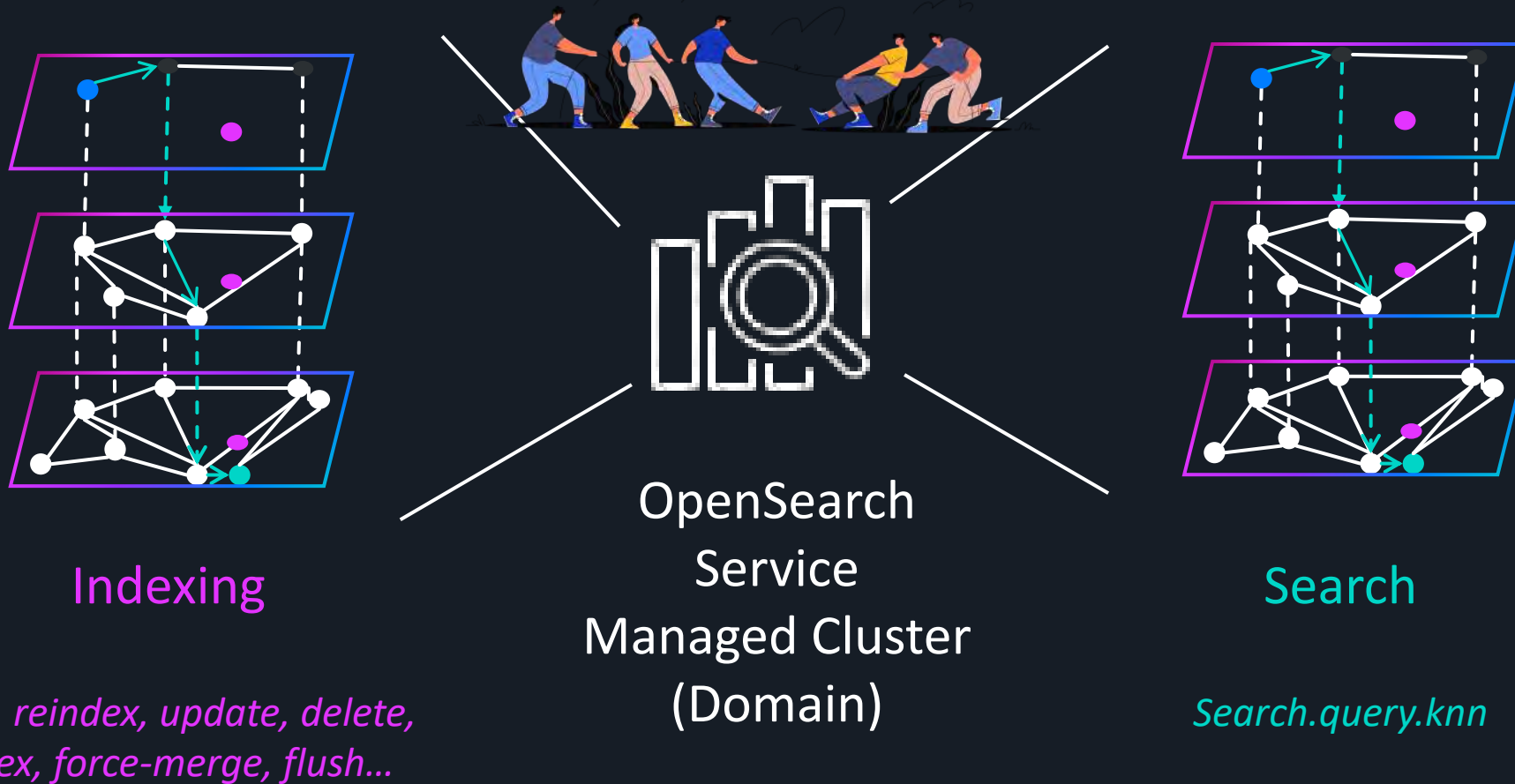
Source <https://aws.amazon.com/blogs/big-data/amazon-opensearch-services-vector-database-capabilities-explained/>



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

How to maintain fast search on dynamic applications?

Compete for significant compute and RAM



Challenges at scale...



amazon music

- Index build and maintenance takes **days**.
- Vector ingestion **can impact search times**.

Help customers...



- Maintain **innovation velocity** and **productivity**
- Build **responsive, dynamic**, AI applications

GPU-accelerated Vector Indexing



NVIDIA cuVS



Best Performance

20X faster index build time, 11X lower latency



Advanced Algorithms

Performance-tuned approximate nearest neighbor search



Flexible Integration

Supports multiple languages including C, C++, Python, and Rust



Interoperable

interoperable between CPU and GPU



Scalable

Enables massive-scale vector search and clustering



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Vector Search Integrations



Vector Databases



Open Source Libraries



Applications



Offline workflows

cuVS

C++

C

Python

Java

Rust

Go

Nearest Neighbors

Exact and Approximate Nearest Neighbors, Quantization, Pre-filtering, Dynamic Batching, GPU/CPU Interoperability, Sparse Nearest Neighbors, Epsilon Nearest Neighbors, k-NN Graph Construction

Distance

Pairwise Distance, 1-Nearest Neighbors, Kernel Gram Construction, Sparse Distances

Clustering

K-means, Hierarchical K-Means, Hierarchical Agglomerative Clustering, Spectral clustering

RAFT

High Performance Machine Learning Primitives

NCCL

CUDA Math Libraries

RMM

CCCL

CUDA

<https://developer.nvidia.com/cuvs>

GPU-Acceleration on Managed Clusters (Domains)

Vectors	CPU-only (Index + Merge)	Add GPU	Speed Gain	Compute Cost
1M 768-dim	1.4 hr.	9.9 min.	8X	8X Less
10M 768-dim	8.5 hr.	36.8 min.	14X	12X Less
113M 1024-dim	28.7 hr.	4.5 hr.	6X	6X Less
1B 128-dim	31.9 hr.	2.8 hr. (Index: 35 min.)	11X	10X Less

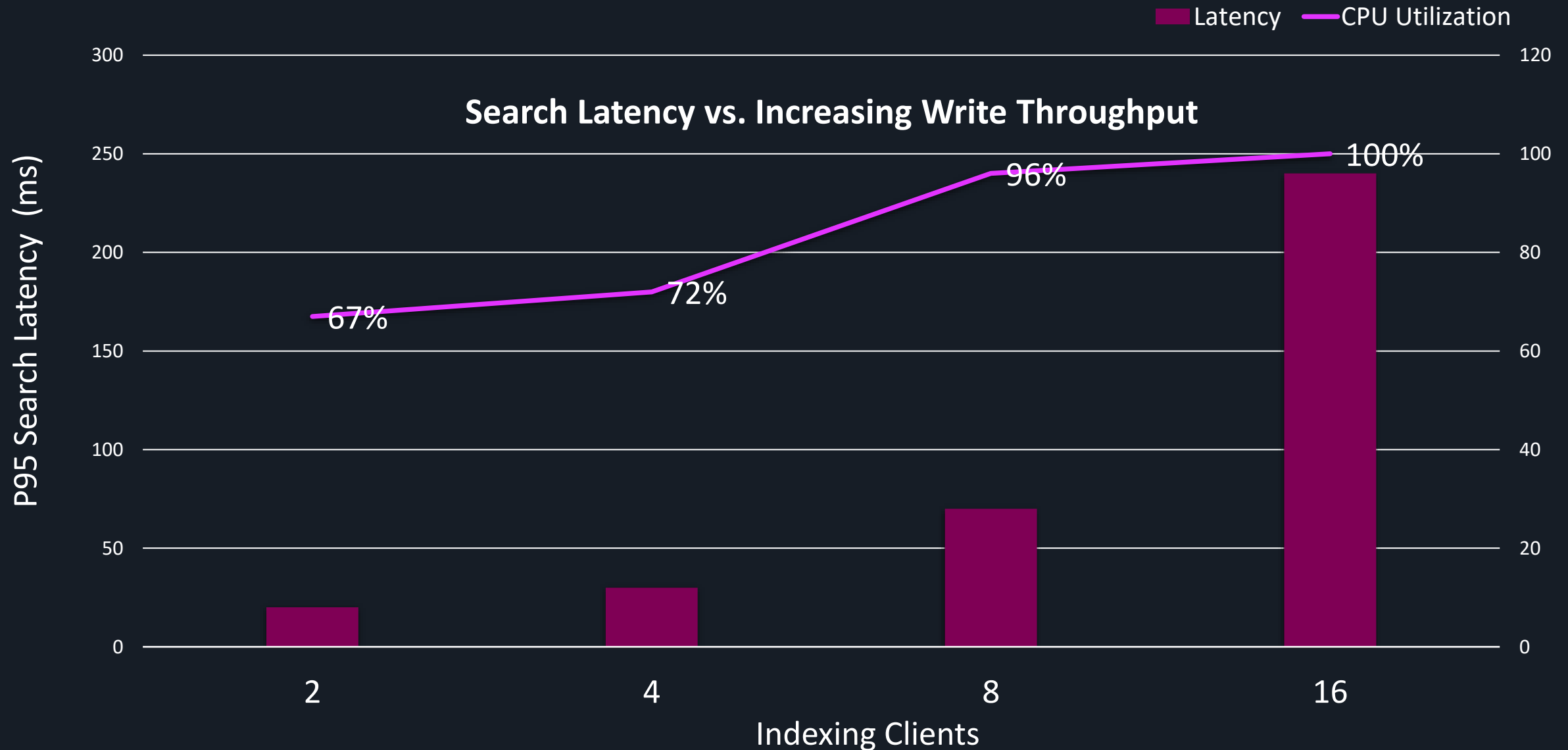


GPU-Acceleration on Serverless (Collections)

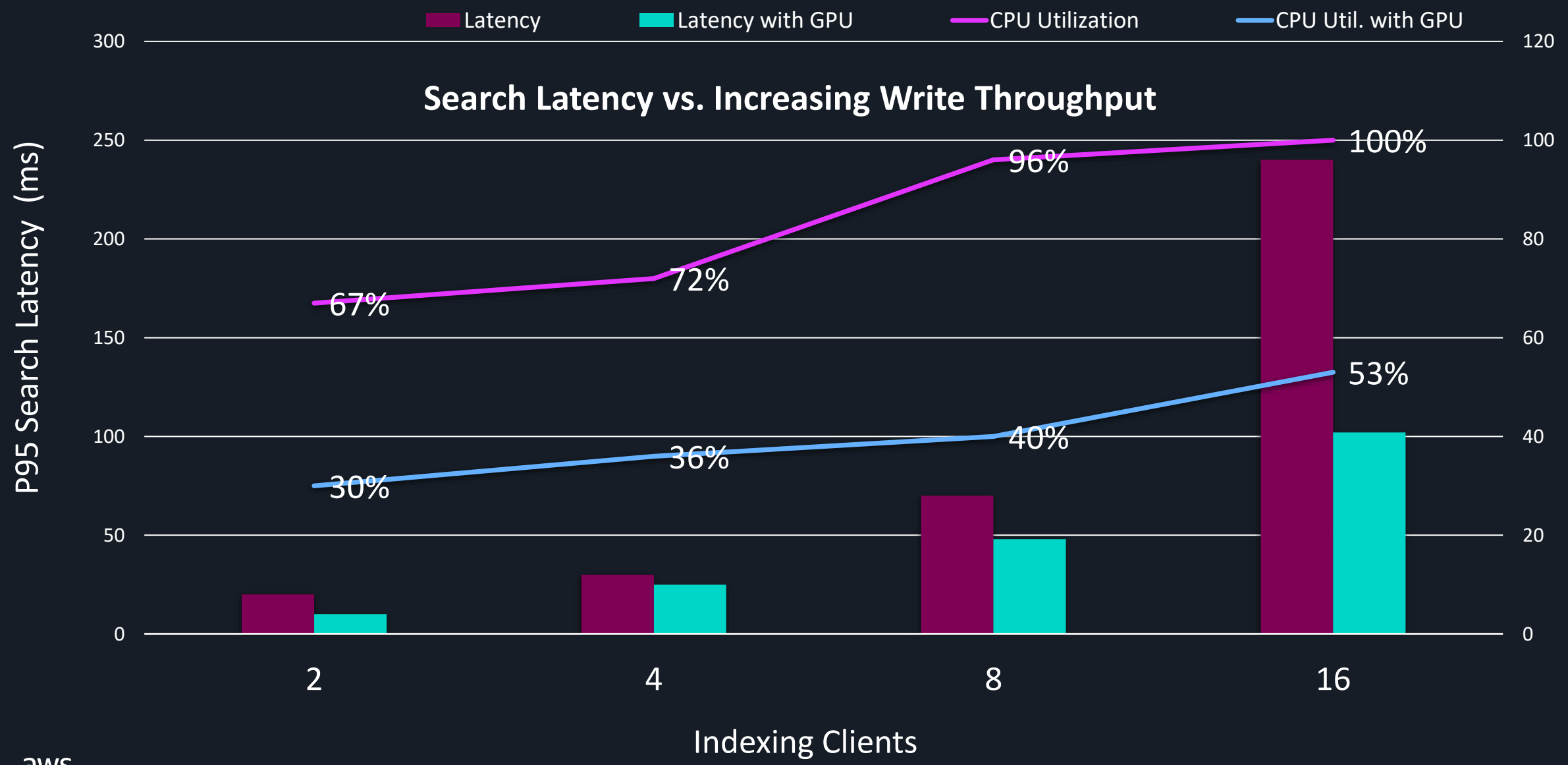
Vectors	Index OCU/hrs.	With GPU	Cost Reduction
1M 768-dim	8	1.5	5.3X
10M 768-dim (min. 32 OCUs)	78	20.3	3.8X
113M 1024-dim (min. 48 OCUs)	2721	304.5	8.9X \$653 vs. \$73
1B 128-dim (min. 48 OCUs)	1562	201	7.8X \$375 vs. \$48



Indexing Impacts Search Speed on Managed Clusters



Offloading Indexing to GPU Improves Search Speed



How can we deliver practical economics?

CPU Cluster

3 X 384 RAM

r8g.12xlarge.search:
48 vCPU, **384 RAM**



1B 1024 dim Vectors
32X compression

Indexing < 30%
of uptime

GPU Cluster

6 X 192 RAM

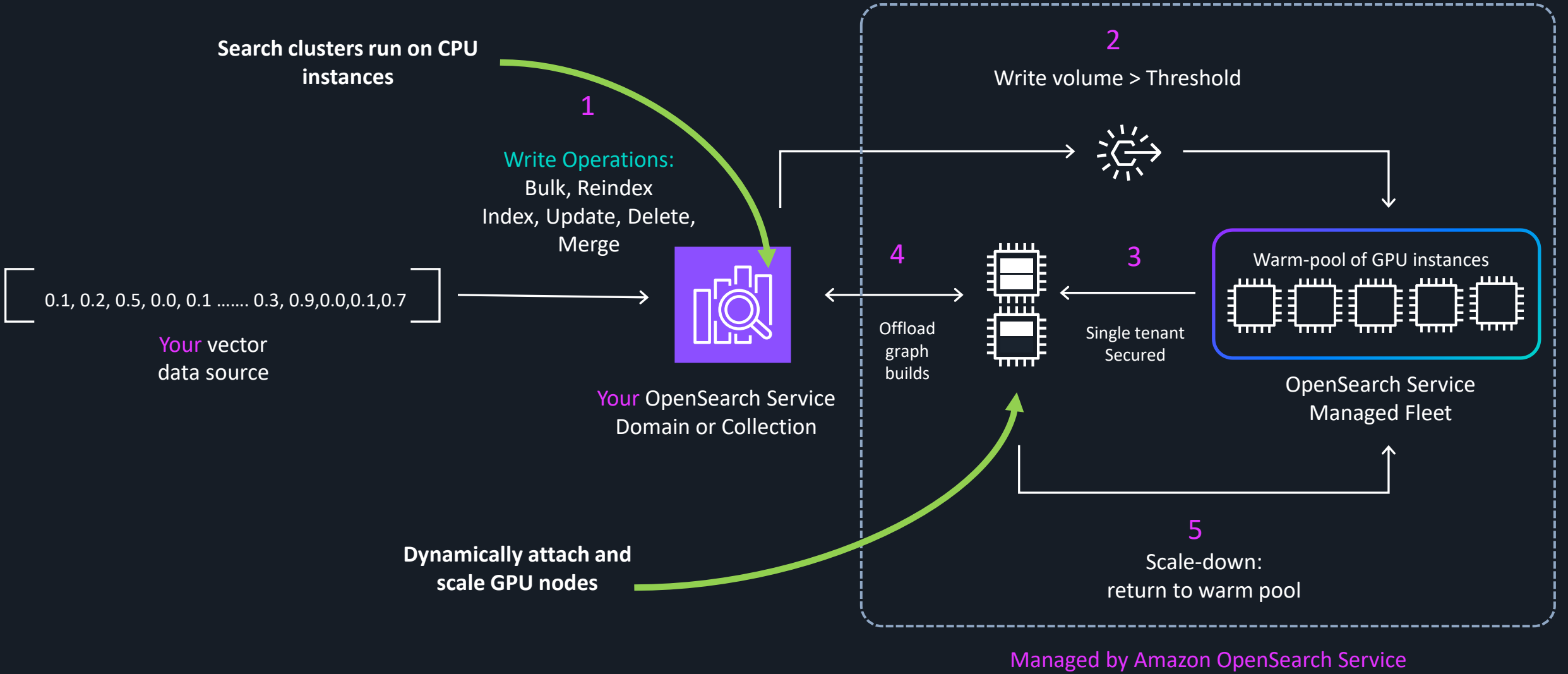
g6.12xlarge.search:
48 vCPU, **192 RAM**

2.4X Cost

Overprovisioned GPU

Wasted \$ from poor
utilization

On-demand GPU acceleration, pay-for-value



Serverless Acceleration for Domains and Collections

1

Enable on
Domain or Collection



2

Pay-on-use

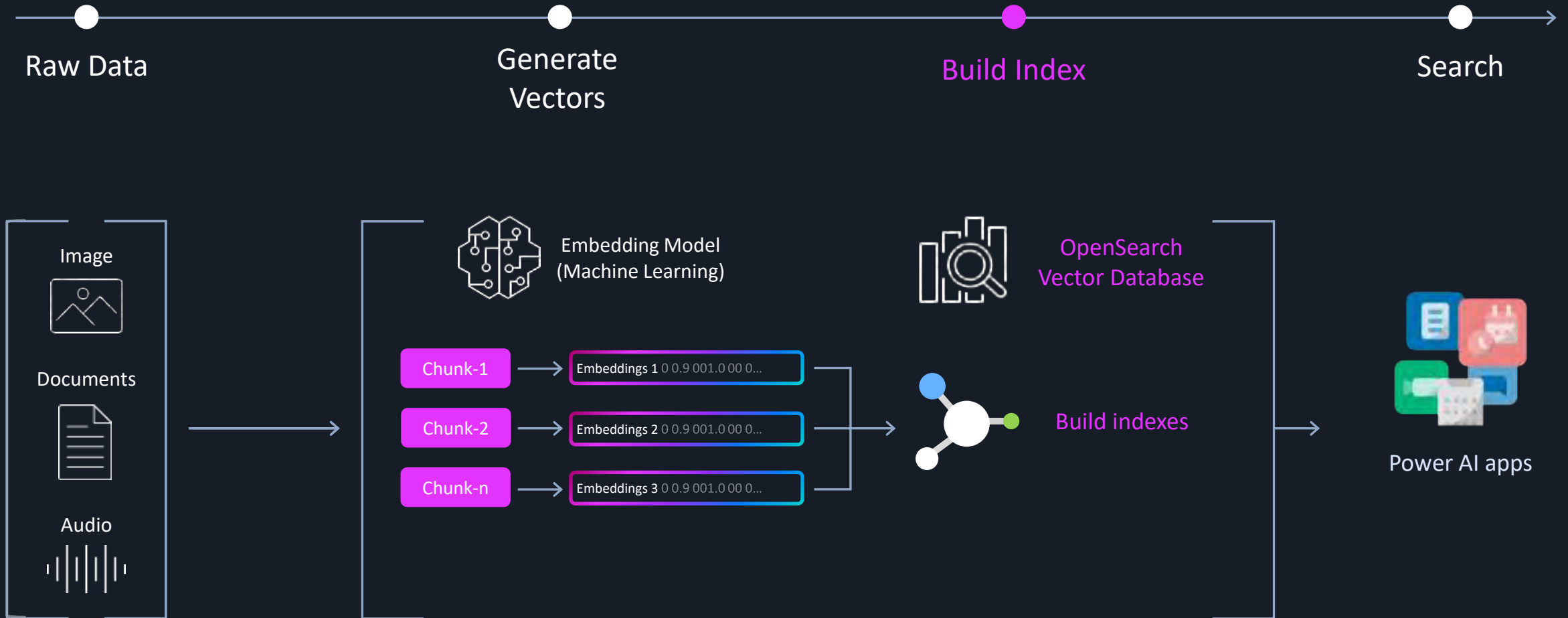
OpenSearch Compute Unit
(OCU)—Vector Acceleration

\$0.24 OCU/hr. (N. Virginia)

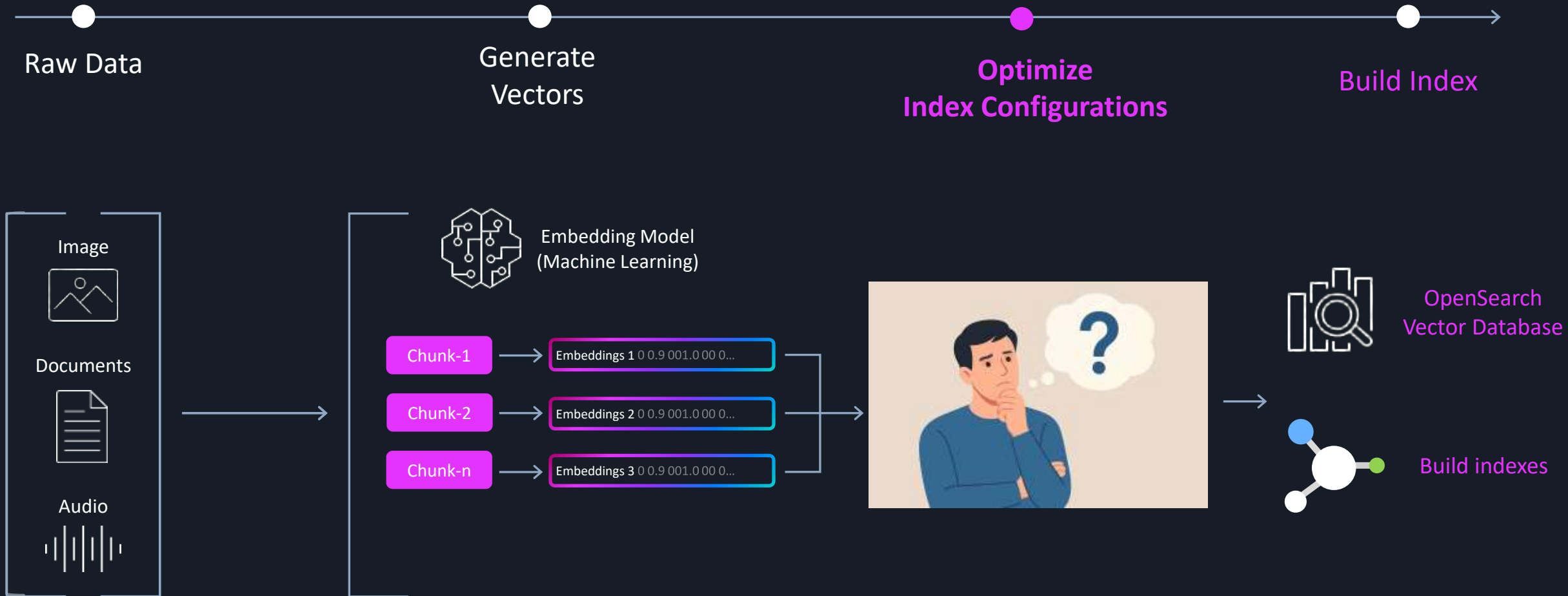
Optimizing Vector Indexes



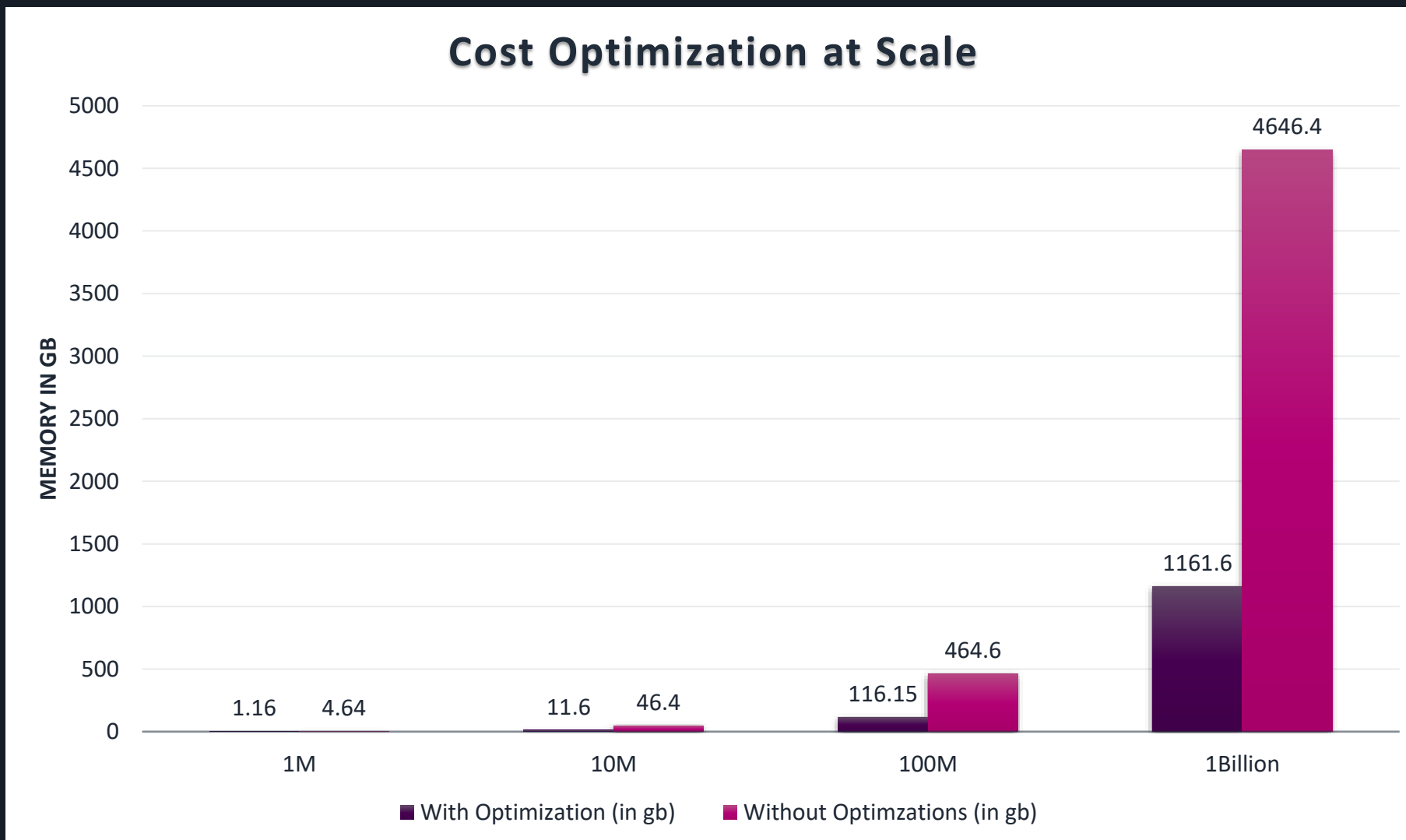
Building a Vector Database (Recap...)



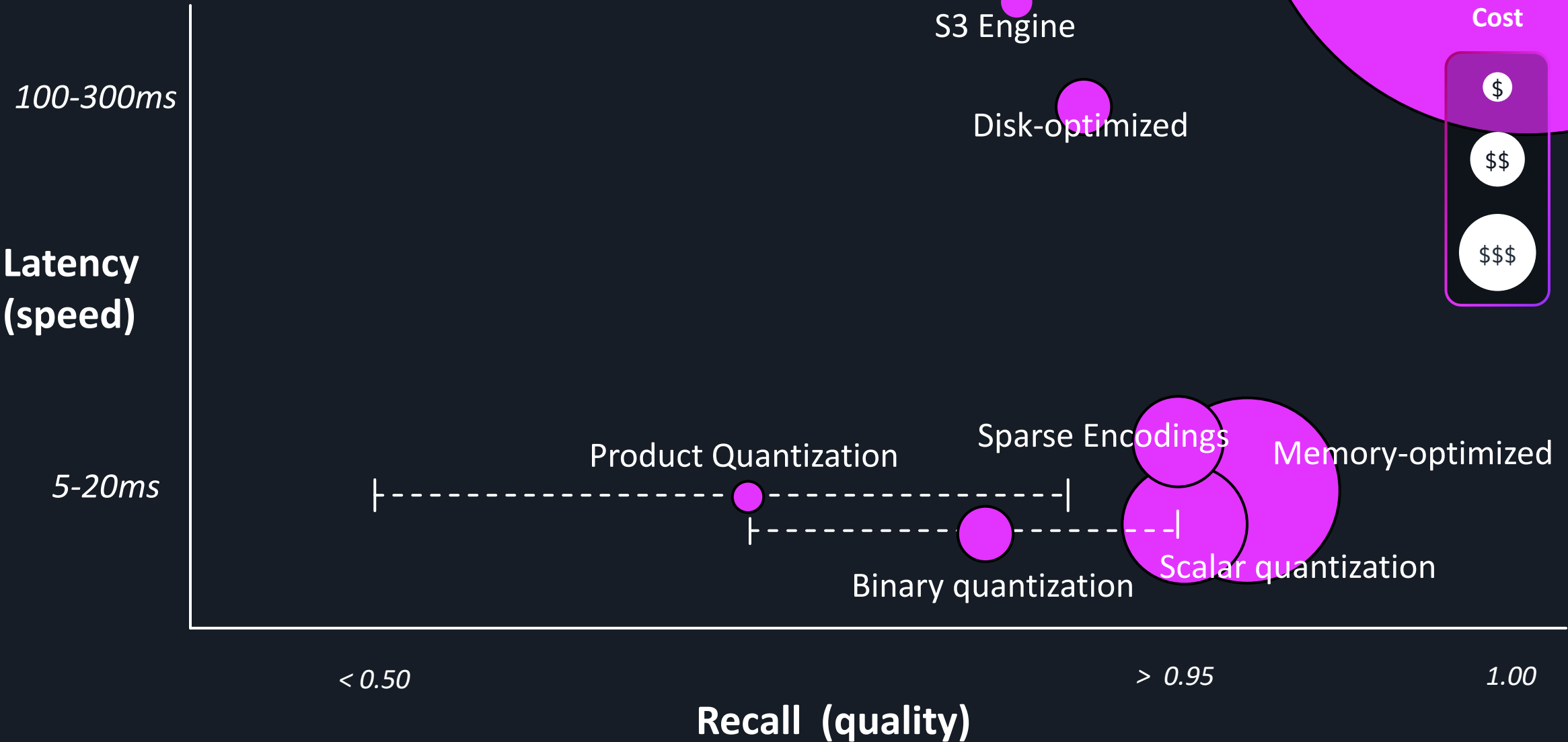
Building a Vector Database



Cost to host a 1-Billion vector index



Configure for favorable trade-offs



Time-consuming, expert-driven process

Select Index Parameters

Algorithms:

HNSW, ef_construction,
m...

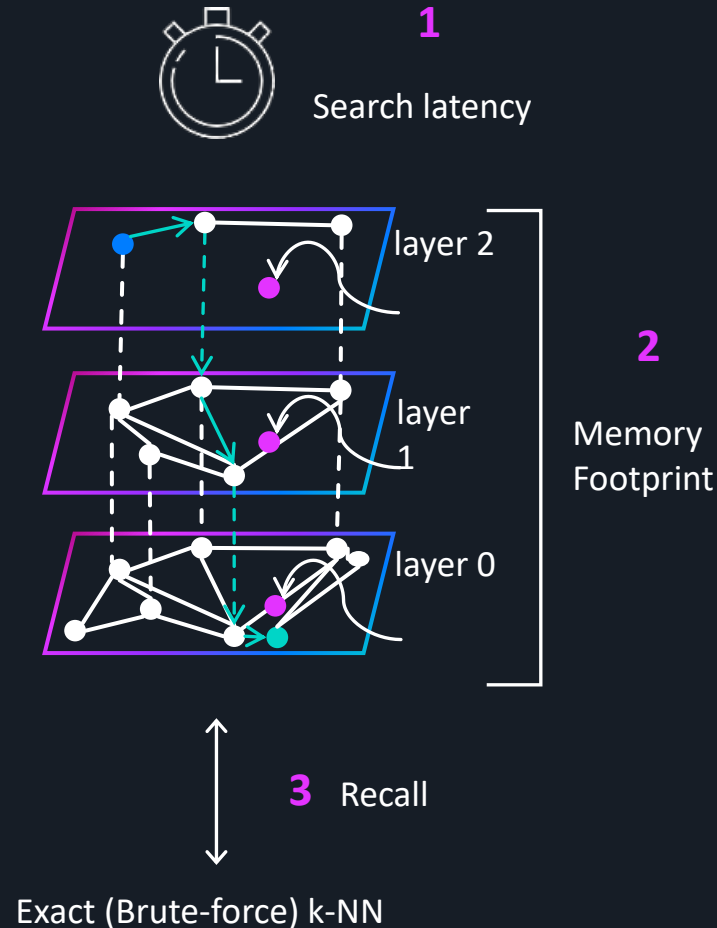
Quantization:

Scalar, Binary,
Product

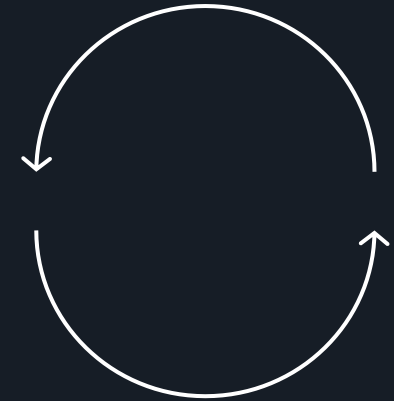
Engine Settings:

Disk-optimized,
In-memory,
Infrequent Queries

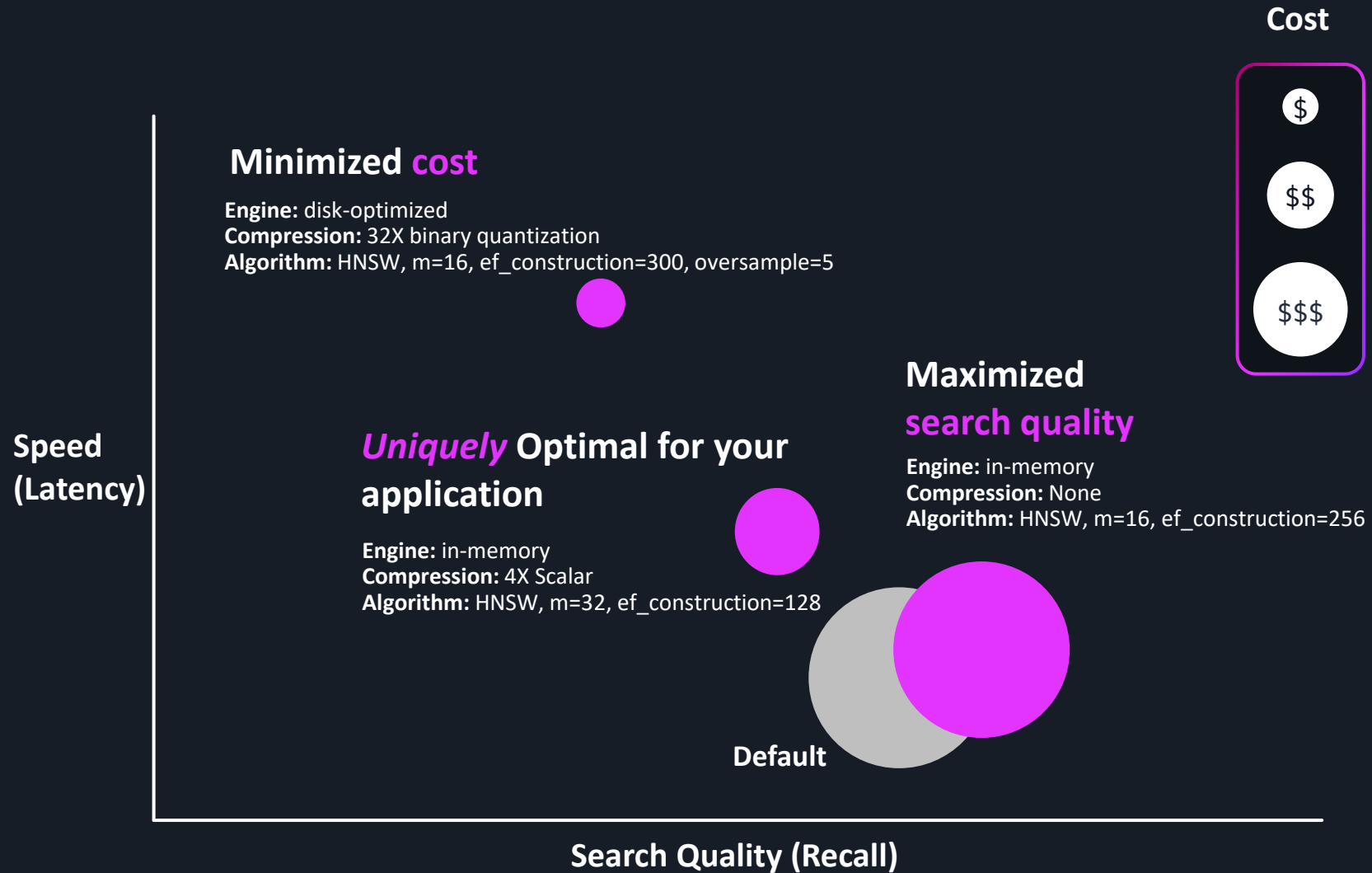
Build and Evaluate Index



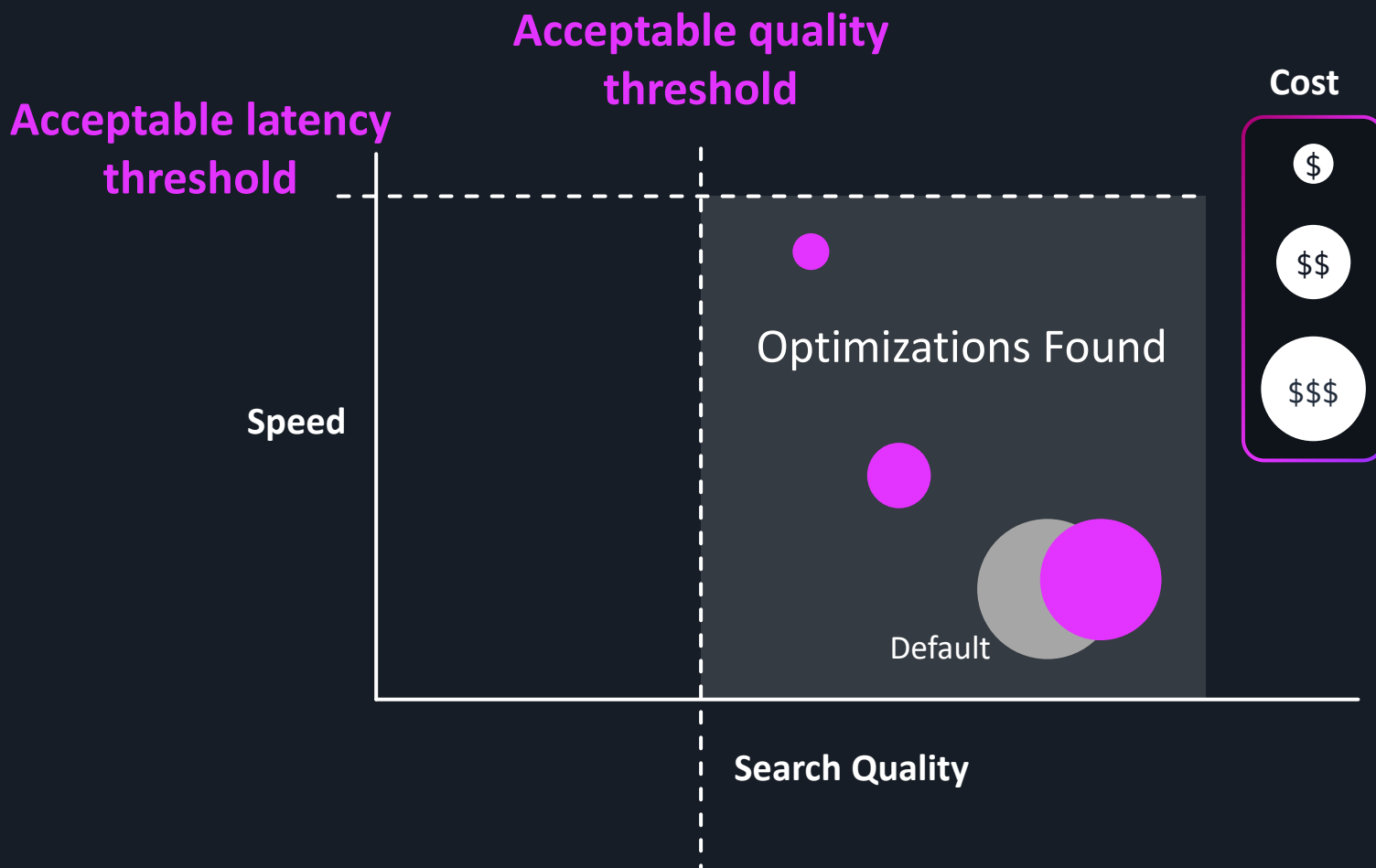
Adjust Parameters and Repeat



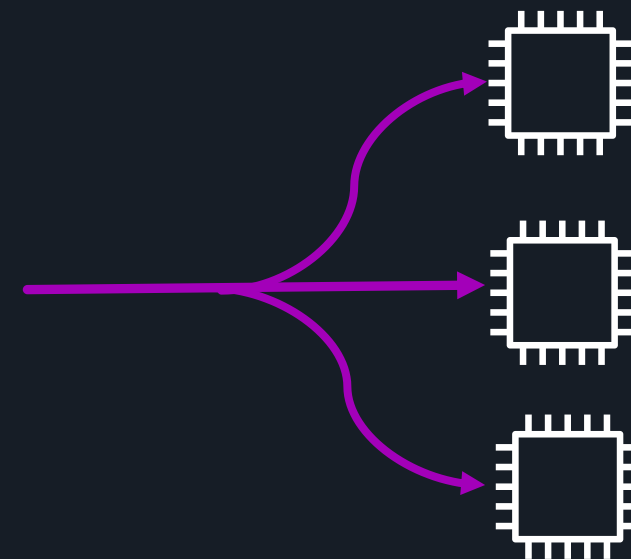
Best optimizations depend on data and use case



Let's simplify! Auto-optimize Vectors



Parallelize Index Builds and Evaluations



Serverless Auto-optimize jobs with a predictable flat rate

Demo: Build an auto-optimized, GPU-accelerated vector database

Level up your skills on AWS Skill Builder

Build AWS Cloud and AI skills, your way, with our online learning center.

LEARN

1,000+ free learning resources

PRACTICE

Hands-on labs and immersive real-world simulations powered by AI

PREPARE

Official AWS Certification exam prep

VALIDATE

Microcredentials to validate practical skills

Scan to start learning



skillbuilder.aws



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.



Thank you

Please complete the session
survey in the mobile app