



# Couchbase

## AI-powered Adaptive Applications and Vector Search

JUNE 2024

# Agenda

New Opportunities with AI

Vectors and GenAI

Why Couchbase for Vectors and GenAI

# AI Creates an Opportunity for a new Application Design



## Everywhere

Interactions are mobile



## Responsive

Latency is unacceptable



## Contextually Proactive

Inputs: who, what, where,  
why, when, real-time data



## Hyper Personalized

Accurate responses that fit  
the situation perfectly

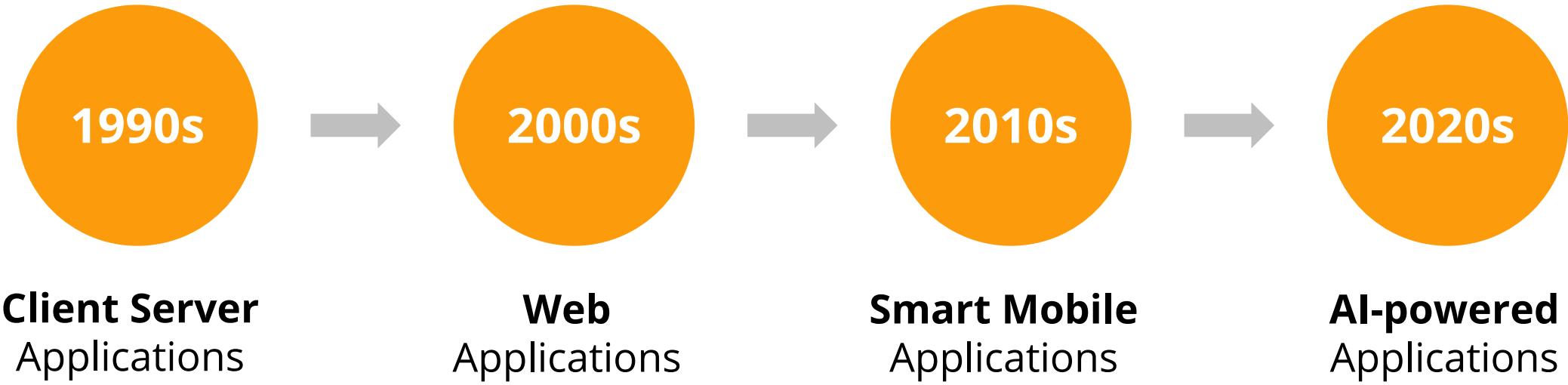


## Innovative

The best adaptations are  
fueled by trustworthy data



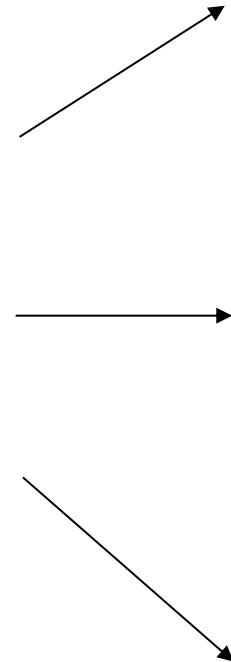
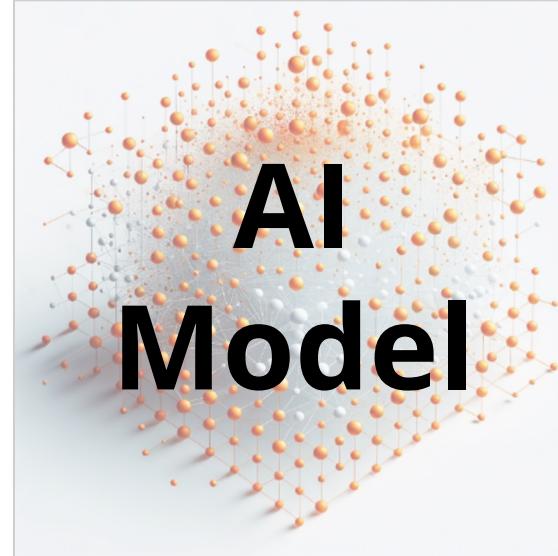
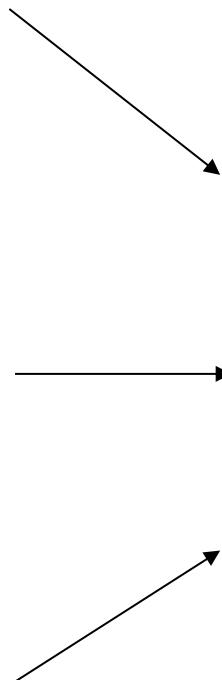
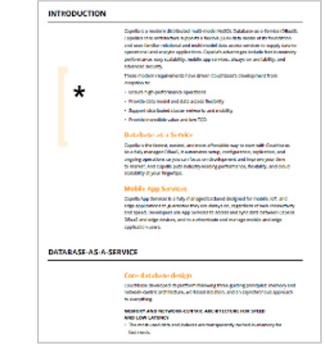
# The AI Shift is Happening



# Vectors and GenAI



# What is a Vector, Vector Search and Why is It Important?



**Embeddings =  
Arrays of Floats**

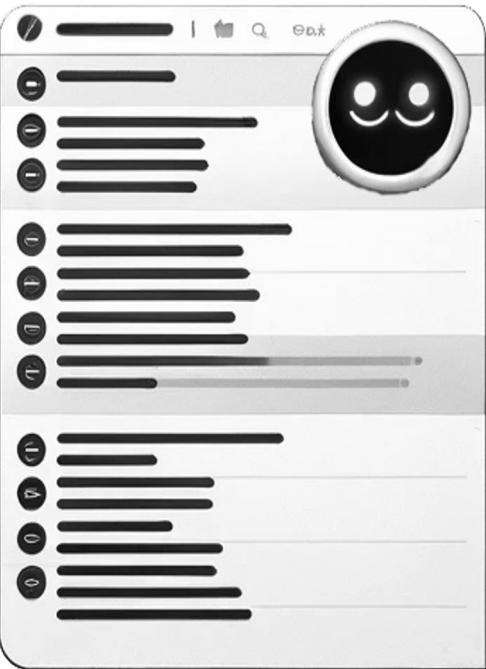
[0.9229, 0.7824, 0.194, 0.0402, 0.206,  
0.4446, 0.5875, 0.4214, 0.2371,  
0.4038, 0.9229, 0.7824, 0.194, 0.0402,  
0.206, 0.4446, 0.5875, 0.4214, 0.2371,  
0.4038]

[0.0509, 0.8862, 0.179, 0.4116,  
0.5995, 0.769, 0.5282, 0.8079, 0.0468,  
0.272, 0.0509, 0.8862, 0.179, 0.4116,  
0.5995, 0.769, 0.5282, 0.8079, 0.0468,  
0.272]

[0.3819, 0.1375, 0.0828, 0.0856,  
0.9081, 0.1556, 0.6251, 0.7356,  
0.2964, 0.2198, 0.3819, 0.1375,  
0.0828, 0.0856, 0.9081, 0.1556,  
0.6251, 0.7356, 0.2964, 0.2198]

# Use Cases

## AI-powered Chatbots and Applications



## Content Generation

BLACK+DECKER 12-Cup Digital Coffee Maker, CM1160B, Programmable, Washable Basket Filter, Sneak-A-Cup, Auto Brew, Water Window, Keep Hot Plate, Black

**Customers say**

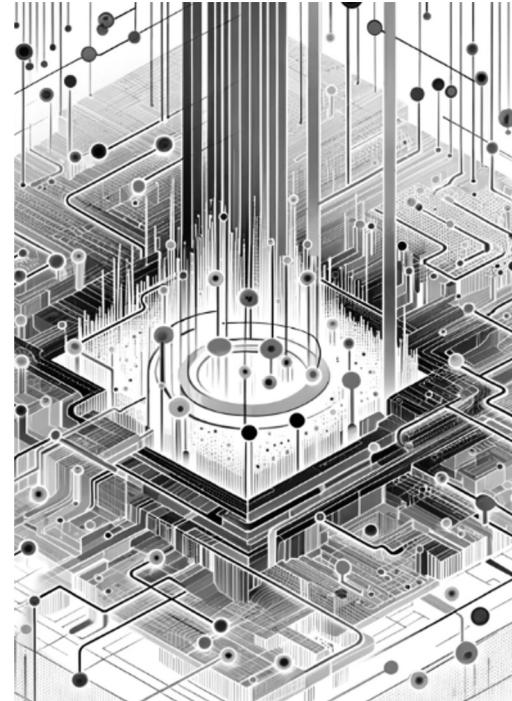
Customers like the ease of use of the coffee maker. They say it's very simple to set and use. Customers are also satisfied with ease of cleaning, value, and speed. However, some customers have reported issues with drips. They mention that the inside will flood over with coffee grounds. Customers disagree on performance, quality, and temperature.

AI-generated from the text of customer reviews

## Advanced Semantic / Hybrid Search



## Data Analysis: Classification / Anomalies



# Vector Similarity Search Example

Searching



Search for "French Open"



Object Loading



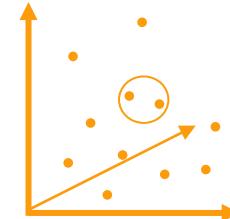
Modern Application

Search for "French open" with a distance of 5.

Model



Nearest neighbors' results  
[.315, .236, 1.43,...]  
[.250, .457, .435,...]



# Key Risks with Apps that Share Data with AI Models

These are C-level showstoppers if they are not addressed

## Sharing proprietary and sensitive data



## Sharing data that induces hallucinations



Retrieval-augmented generation (RAG) enhances the accuracy and reliability of GenAI-powered applications

# Why Couchbase for Vectors and GenAI

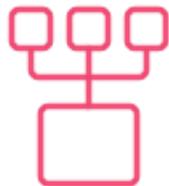


# Extending Couchbase Search for Vectors

No Separate Database Needed

## Full-Text Search

```
my: Doc 1, Doc 2, Doc 3  
dog: Doc 1, Doc 2, Doc 81  
has: Doc 1, Doc 2, Doc 3  
fleas: Doc 1, Doc 81  
...
```



- Documents scanned for words
- Inverted indexes created** linking word to document id and text positions.

### Pre-query work

- User** wants docs contains "Big Cat"
- Text analysis performed, to find same or similar misspelled words, term frequency, other criteria
- Documents are scored

### Query Process

- Documents returned to **User**
- Sorted by highest ranking score

### Results

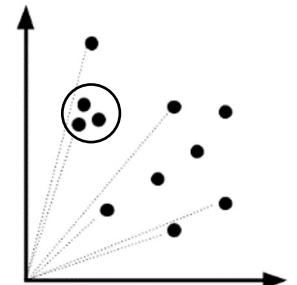
**Results:** Docs w/ "Cat in a Hat, Big Hat"



## Vector Search

- Text, image, audio, other scanned into vectors for "nearest-neighbor" similarities
- These multidimensional **embeddings are indexed** as a special "vector" type.

.2758, .2637, .4315,  
.1372, .6194, .2164,  
.1654, .2856, .2785,  
.4615, .7516, .3194,  
.5456, .5229, .9567,  
.2342, .5234, .6345,  
.3634, .3153, .9502,  
.9124, .2856, .2785,  
.9634, .7516, .3194,  
.5456, .5229, .9567



- User** sends query sentence "Big Cat" that is converted to its own vector embedding
- Analysis performed, to find top N nearest matches with other
- Results returned to AI-based Application
- App sends top N result (context) plus the query (prompt) to AI model

Answers

---

---

---

---

---

- Retrieval Augmented Generation**
- AI model generates output
- AI app delivers result to **User**

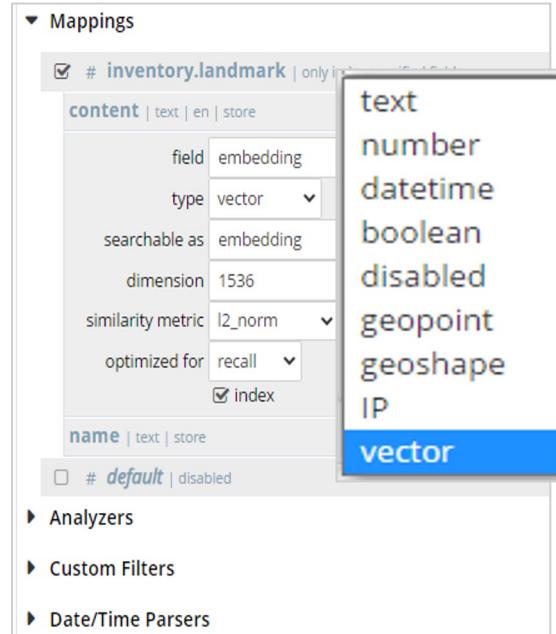
**Results:** Text discussing tigers, leopards, lions and pumas

# Vector Storage, Indexing and Query

## JSON Storage

```
{ "type": "shoes",
  "productId": "CP123456",
  "category": "Gym Shoes",
  "name": "Beach Sneakers",
  "brand": "Ultimate Surf",
  },
  "description": "The ultimate companion for beach adventurers, designed to seamlessly transition from sandy shores to urban landscapes. This innovative sneaker features a water-resistant, quick-drying mesh upper, allowing your feet to breathe while keeping them dry.",
  "descriptionVector": [0.131, 0.339, -0.611, 0.981,...]
```

## Indexes

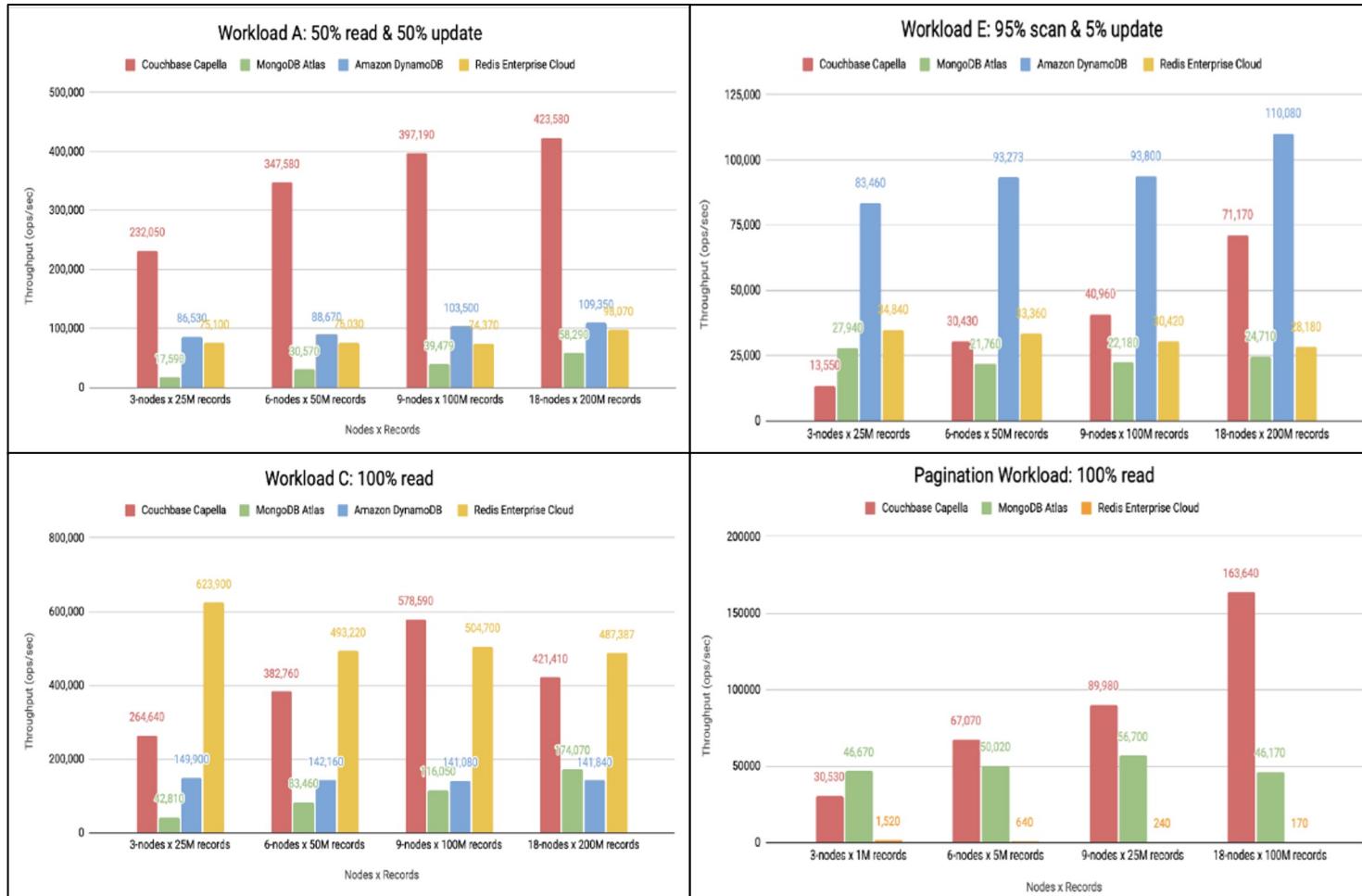


## SQL++

```
SELECT *
FROM product
WHERE LOWER(product.type) = 'shoes'
AND product.size = 11
AND product.price between 50 and 80
/* desc SIMILAR TO 'blue running
shoes' */
ORDER BY
GSI_VECTOR_ORDER(desc_embedding,
{
  "knn": [
    "field": "desc_embedding",
    "vector": [0.1, 0.334, -0.604, 0.985]
  ]
})
LIMIT 4
```

# Performance Proven via YCSB Tests

## Proven Speed and Flexibility



## Capella

- Excels at Read and Write
- Scales effectively for multiple workloads
- Best price-performance

## MongoDB

- Struggles to scale, strongest at 3 nodes
- Price performance makes it worse
- Weakest performer

## DynamoDB

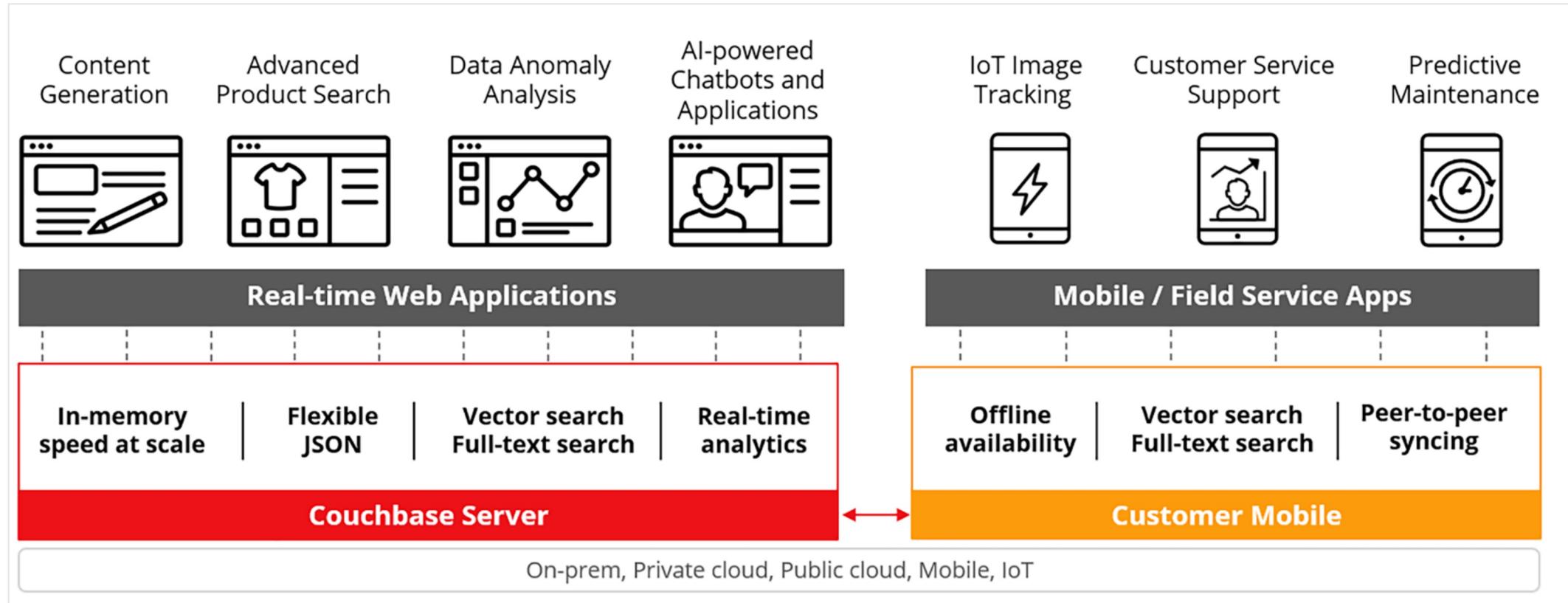
- Excels at scans, but throws excessive errors
- Challenged across other workloads

## Redis

- Excels at read-only, Capella meets at scale
- Regularly used as cache for Atlas or DynamoDB
- “Fails” Pagination workload

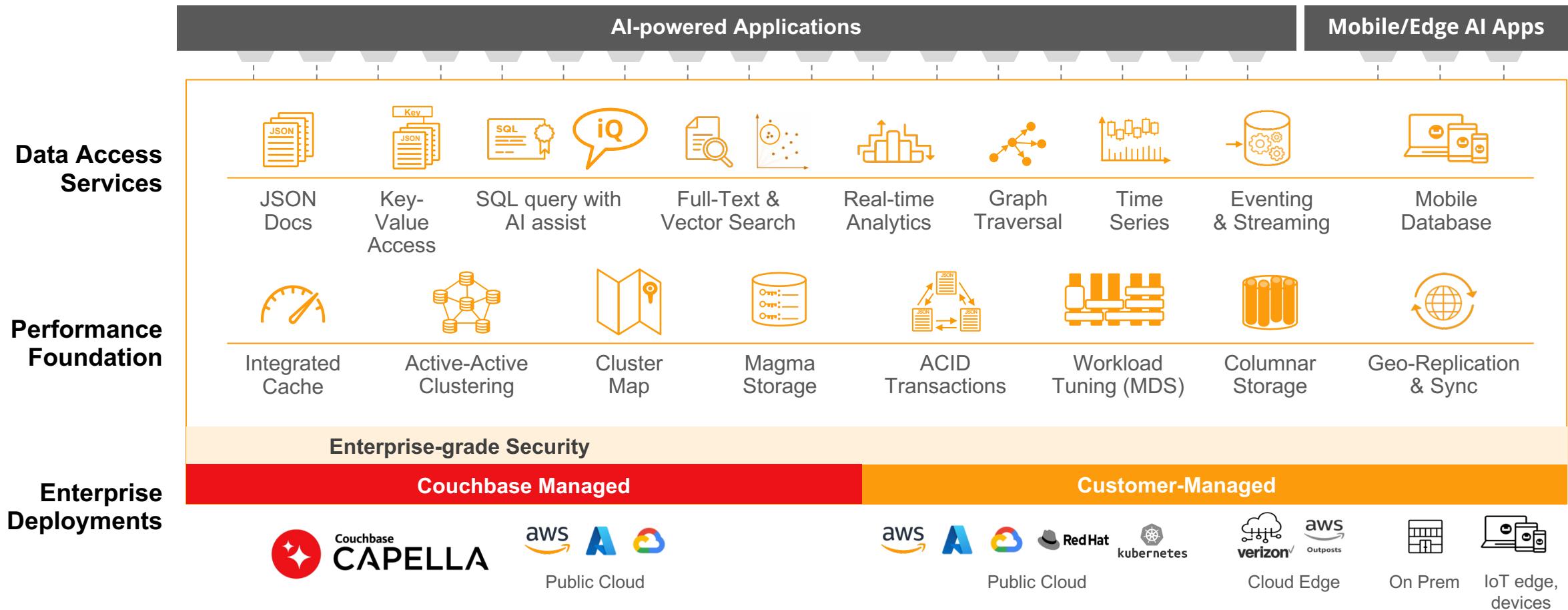
# Vector Across our Products

First in the industry to announce support for all 3 deployments: Cloud, on-prem, mobile



# Couchbase's Differentiated Architecture

Innovate Faster with Couchbase



# Beyond Vector Search: Hybrid Search

## Building real-world use cases

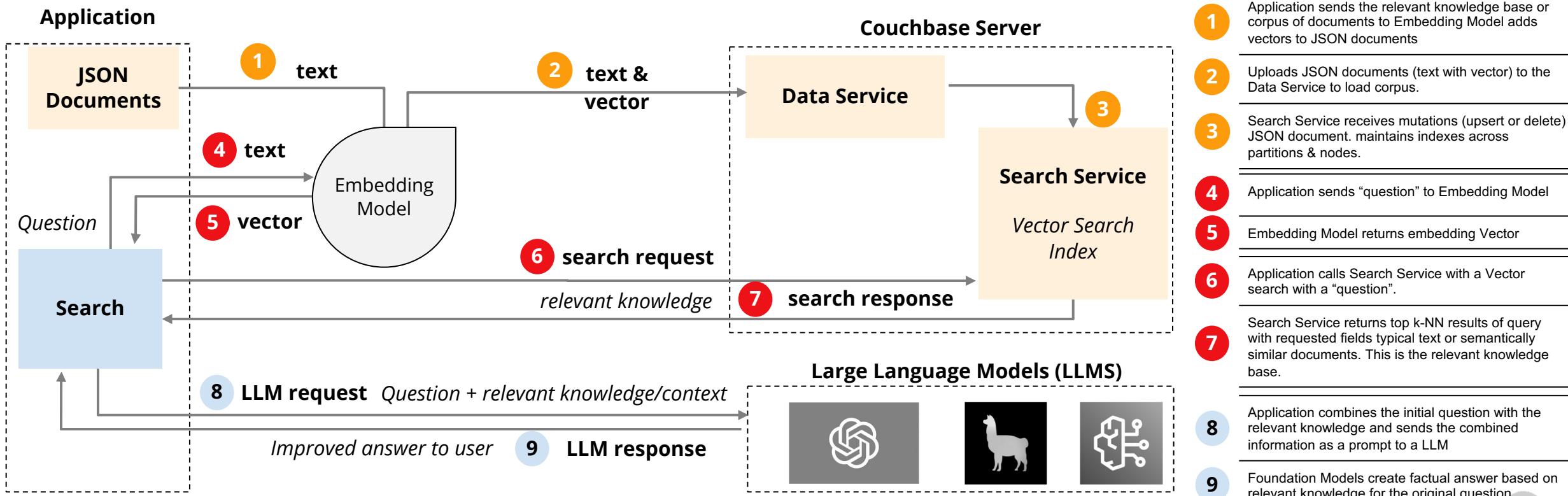
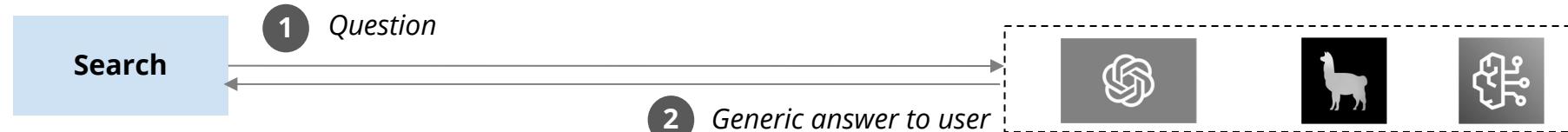
### Customer wants new shoes

- Match color and style of an object (semantic/vector)
- Description to mention “casual and fun” (text/fuzzy)
- Price between \$100 and \$200 (range)
- With rating over 4.5 stars (range)
- Within 15 miles (geospatial)
- Available today in stores (inventory)



# RAG Application: Chatbot for Auto Parts Supplier Chain

**Question:** What is the best way to reduce muffler noise on my 1967 Ford Bronco?



# Ecosystem Integration

1

Drive developer productivity

2

Optimize AI processing

3

Enable AI-driven apps anywhere, including the Edge

4

Part of a vibrant AI partner ecosystem

Today

Roadmap



Amazon  
SageMaker



Amazon Bedrock

# Couchbase: Architected for AI-based Applications



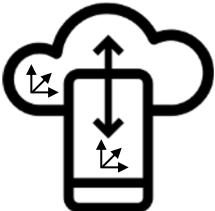
## No Separate Vector Database Needed

Simplify | Lower latency | Save



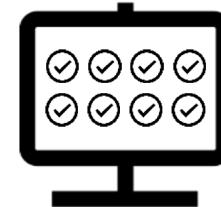
## Proven Speed and Flexibility

In-memory architecture | Flexible JSON | Powering Indexing



## Vector Across our Products

On-prem | Cloud | Offline mobile



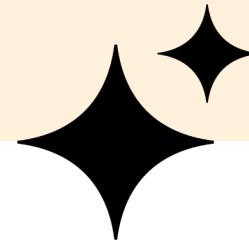
## Broad Platform Capabilities

Profiles & Catalogs | Vector Search | Text Search | Graph | Real-time Analytics | Powerful SQL



## Ecosystem Integrations

LangChain | LlamaIndex



# Thank you!



Couchbase