# Try it Today for a Chance to Win!

Anyone who tries out the SingleStore tutorial today will be entered for a chance to win new Meta Ray Bans Smart Glasses or branded AirPods Pro.

Simply click this link, sign up for a SingleStore Free Trial, try out the features, and we'll announce a winner by the end of today's session!

SingleStore™

bit.ly/devday-boston-raffle

Google Cloud | SingleStore™

# Join us for our exclusive Happy Hour!

After our workshop, unwind and network with fellow developers over cold beverages and appetizers atShy Bird, just a short walk from the Google office:

**Shy Bird Kendall Square**
390 Third St
Cambridge, MA 02142

# RAG to Production

**1**

## Rich context for RAG

Data freshness to provide recency and knowledge graphs for better accuracy

**Accuracy + richness of context**

**2**

## Multi-modal RAG

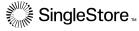Ability to store audio, video, images in addition to structured data in one place

**Real-time AI**

**3**

## RAG with analytics

Ability to provide richness of data back to the LLM and also for better evals

**Accuracy + richness of response**

Google Cloud | SingleStore ™

# Advanced RAG

**Evaluations**
(Scoring RAG and LLM responses)

**Guardrails**
(Vertex AI)

**Explainability**
(LangSmith)

**Re-training on the fly**
(LangSmith)

# SingleStore + Google Vertex AI

Google Cloud | SingleStore ™

# What Is SingleStore?
## Data in Real Time

**MODERN APPLICATIONS**
High volume transactions

**IN-APP ANALYTICS APPLICATIONS**
Rideshare and recommendations

**GEN AI APPLICATIONS**
Use of Retrieval Augmented Generation
(RAG) across enterprise data

# Who Is SingleStore?

Pro Max version with
ANN and compute

Rebranded to
SingleStore

$558M in
total funding

Patent awarded

Three-tiered architecture

$30M funding

Vectors + semantic
search support

Raj Verma
joins as CEO

$100M+ ARR •

Founded
as MemSQL

MemSQL 5.0

350+ customers •

2011    2012    2013    2014    2015    2016    2017    2018    2019    2020    2021    2022    2023    2024

**Investors**

Accel    GLYNN CAPITAL.    khosla ventures    سنابل للاستثمار SANABIL INVESTMENTS    IBM    GV    Prosperity7 VENTURES    DELL Technologies CAPITAL

iqt IN-Q-TEL™    Goldman Sachs    INSIGHT PARTNERS    R|E|V Funded by RELX    Hewlett Packard Enterprise

Google Cloud | SingleStore™

# SingleStore Tenets

Speed

Scale

**Semantic + full-text search**
Fast hybrid search

**Transactions + analytics**
Petabytes of enterprise data

**SQL + JSON**
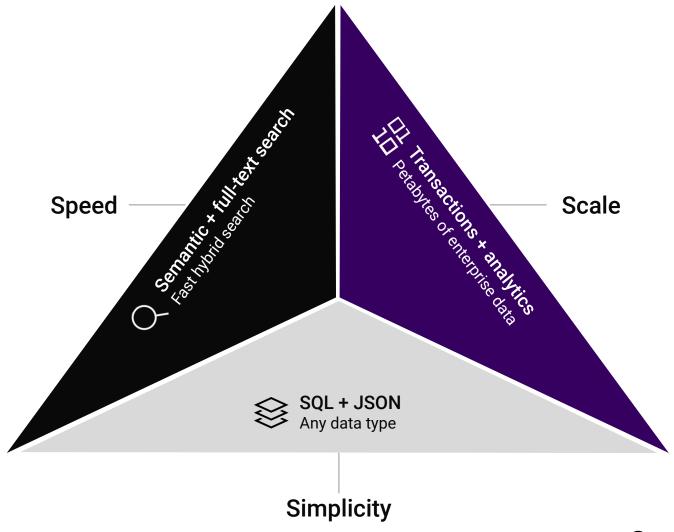Any data type

Simplicity

9

Google Cloud | SingleStore ™

# SingleStore
## Key product capabilities

**Universal Storage**
[Patented] single table type
for transactions and analytics.
Compatibility with ANSI SQL,
MySQL + MariaDB ecosystem;
NoSQL

**Pure speed**
Fast transactions, analytics,
vectorization and query
compilation

**Fast ingestion**
SingleStore Pipelines — load
data with updates

**Bottomless storage**
Separation of storage +
compute in a unified database

**MongoDB® API**
SingleStore Kai™ powers
up to 1,000x faster analytics
on JSON

**Multi-model**
Relational foundation with
support for various data types,
including vectors

**Run anywhere**
Hybrid, multi-cloud, SaaS,
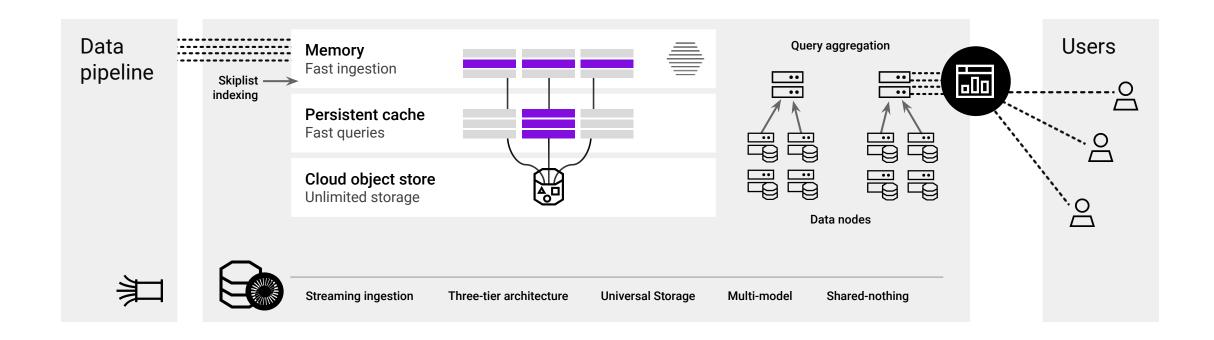on-premises, Kubernetes
operator

**Intelligence**
Build generative AI and other
apps with **Notebooks** on top
of SingleStore

Google Cloud | SingleStore™

# SingleStore: Under the Hood

**Data pipeline**

Skiplist indexing

**Memory**
Fast ingestion

**Persistent cache**
Fast queries

**Cloud object store**
Unlimited storage

Query aggregation

Data nodes

**Users**

Streaming ingestion    Three-tier architecture    Universal Storage    Multi-model    Shared-nothing
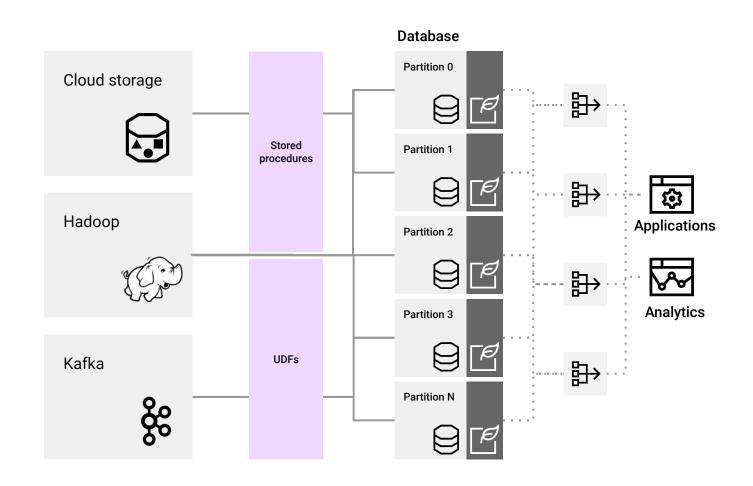
SingleStore ™

# Pipelines: Real-Time Ingest

## Product overview

- Millions of records / sec with immediate availability

- Built-in component of database

- Transactional consistency

- Exactly-once semantics

- Native integrations with Kafka, cloud storage, HDFS

- Ingest to stored procedures to shape/modify data

- Pipelines support JSON, Avro, Parquet and CSV data formats.
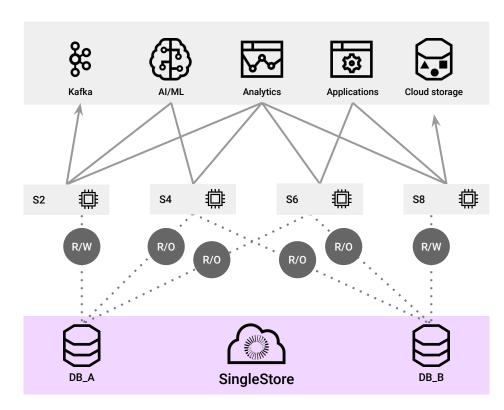
SingleStore ™

# Workspaces Overview

## Isolated workloads on shared data
Compute operates on shared data, without one workload affecting the performance of others

- Allow granular scalability and isolation of compute resources

- Eliminate cost of moving and maintaining data between multiple workloads

- Scale ingest and compute workloads independently

- Isolate internal and customer facing real-time applications simultaneously on shared data

## Separation of compute + storage

Separate write transactions from read-only workloads (analytics) each with its own dedicated compute resources — without data duplication.



Each workspace is isolated and can scale compute independently

SingleStore ™

# Searching + Contextualizing

## Simple semantic search

```
stmt = f"""
        SELECT
            text,
            DOT_PRODUCT_F64(JSON_ARRAY_PACK_F64(%s),
embedding) AS score
        FROM {table_name}
        ORDER BY score DESC
        LIMIT %s
    """.format(table_name=table_name)
```

## Hybrid search

```
SET @promptEmbedding = '[...]' :> vector(1536) :>
blob;
SELECT
  products.id,
  MATCH(products.description) AGAINST ('blue') AS
ft_score_color,
  title_v <*> @promptEmbedding AS v_score_title,
  description_v <*> @promptEmbedding AS
v_score_description,
  v_score_title + v_score_description AS score
FROM products
JOIN product_sku ON products.id =
product_sku.product_id
JOIN product_sizes ON product_sku.product_size_id =
product_sizes.id AND product_sizes.value = 'xs'
WHERE ft_score_color AND (v_score_title >= 0.75 OR
v_score_description >= 0.75)
AND price BETWEEN 100 AND 1000
GROUP BY products.id
ORDER BY score DESC
LIMIT 5;
```

SingleStore ™

# SingleStore AI Capabilities

## Core engine

MySQL wire compatibility

Scale out (distributed SQL)

SQL compiled to machine code

Universal storage table type

Three-layer storage architecture

Enterprise-grade availability, security

Point-in-Time Recovery (PITR)

Projections and branching

## Hybrid search

ANN with PQ, IVF, HNSW

Multiple indices on one column

Exact KNN search

Vector datatype

JLucene under-the-hood for exact keyword search

## Ecosystem

CDC in from MySQL, MongoDB®

Pipeline with Kafka

Connector in LangChain, Llamaindex, Flowise, etc.

SDKs for CRUD in different languages

Google Cloud | SingleStore ™

# Optimize Tradeoffs on Your Own Terms

## IVF-PQ (ANN)

+ Low index build time

+ Cost (smaller index size)

- Recall

## HNSW-FLAT (ANN)

+ Recall

+ Speed (performant at high dimensionality)

- Large index size

## kNN

+ Recall

+ Cost (no added size)

- Slow at scale (vs ANN)

SingleStore ™

# Vector Range Search using ANN index

```
select id, txt, vec <*> @query_vec as score
from t
where score > 0.70;
```

What it's good for:

- Answer, is there anything even close to this query vector?
    - "even close" might mean "vector_col <*> @query_vec > 0.60"
        - can be faster than a "top K" query


- Find "all the really close stuff"
    - "vector_col <*> @query_vec > 0.95"
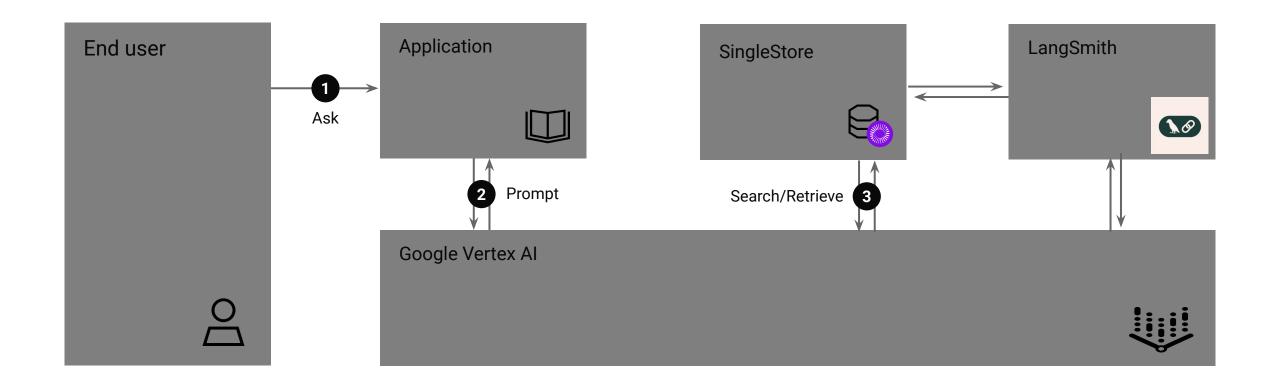    - works in one shot, easier to program that running top 5, top 10, top 20 in a loop

SingleStore ™

# 15-min break

SingleStore ™

# Demo: Production Grade RAG with SingleStore + Vertex AI

Google Cloud | SingleStore ™

# Production-ready RAG Architecture

**End user** → ①  Ask → **Application**

**Application** ② Prompt ↔ **Google Vertex AI**

**SingleStore** ↔ **LangSmith**

**SingleStore** — Search/Retrieve ③ ↔ **Google Vertex AI**

Google Cloud | SingleStore™

# Production Grade RAG

## Scenario

- **Input:** User question regarding contract data
- **Output:** Response to user question
- **Goal 1:** Provide the user with fast responses to their questions
- **Goal 2:** Log all relevant information

## Why SingleStore + GCP?

ML Model Deployment on Vertex AI

Streaming ingest through Pipelines

Fast Vector Search in SingleStore

Large context window with Gemini

Google Cloud | SingleStore ™

Google Cloud | SingleStore™

# Thank You