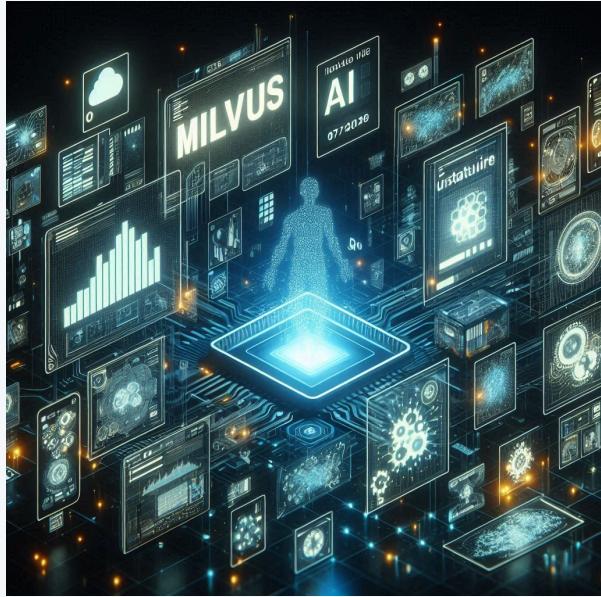


It's in the Air Tonight. Sensor Data in RAG

Tim Spann @ Zilliz



Slides



X

Overview

I will do a quick overview of the basics of Vector Databases and Milvus and then dive into a practical example of how to use one as part of an application. I will demonstrate how to consume air quality data and ingest it into Milvus as vectors and scalars. We will then use our vector database of Air Quality readings to feed our LLM and get proper answers to Air Quality questions. I will show you how to all the steps to build a RAG application with Milvus, LangChain, Ollama, Python and Air Quality Reports. Finally after demos I will answer questions.

A G E N D A

Introduction
Overview of Vector Databases
A Quick Introduction to Milvus
Consume and Ingest Air Quality Data
Building a local RAG application
Q&A

01

Introduction



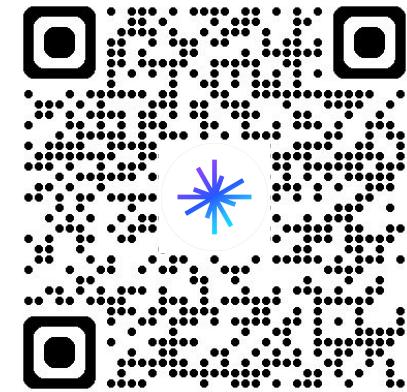
Tim Spann

Principal Developer
Advocate, Zilliz

tim.spann@zilliz.com

<https://www.linkedin.com/in/timothyspann/>

<https://x.com/PaaSDev>



Show Me A Demo



A screenshot of the Milvus UI interface. At the top, there's a navigation bar with "Database: default" and "airquality". Below it is a sidebar with icons for "Overview", "Vector Search", "Data", "Partitions", and "Segments". The main area shows an "Overview" of the "airquality" index, which has 34 partitions. It includes fields like "Created Time" (2024-04-24 03 PM) and "Consistency" (Bounded). A "Schema" table lists fields such as "id", "DateObserved", "HourObserved", "Latitude", "Longitude", "ParameterName", "ZipCode", "AQI", "vector", "details", and "location", each with its type (e.g., Int64, VarChar), index name, index type (e.g., id, vector), and parameters.



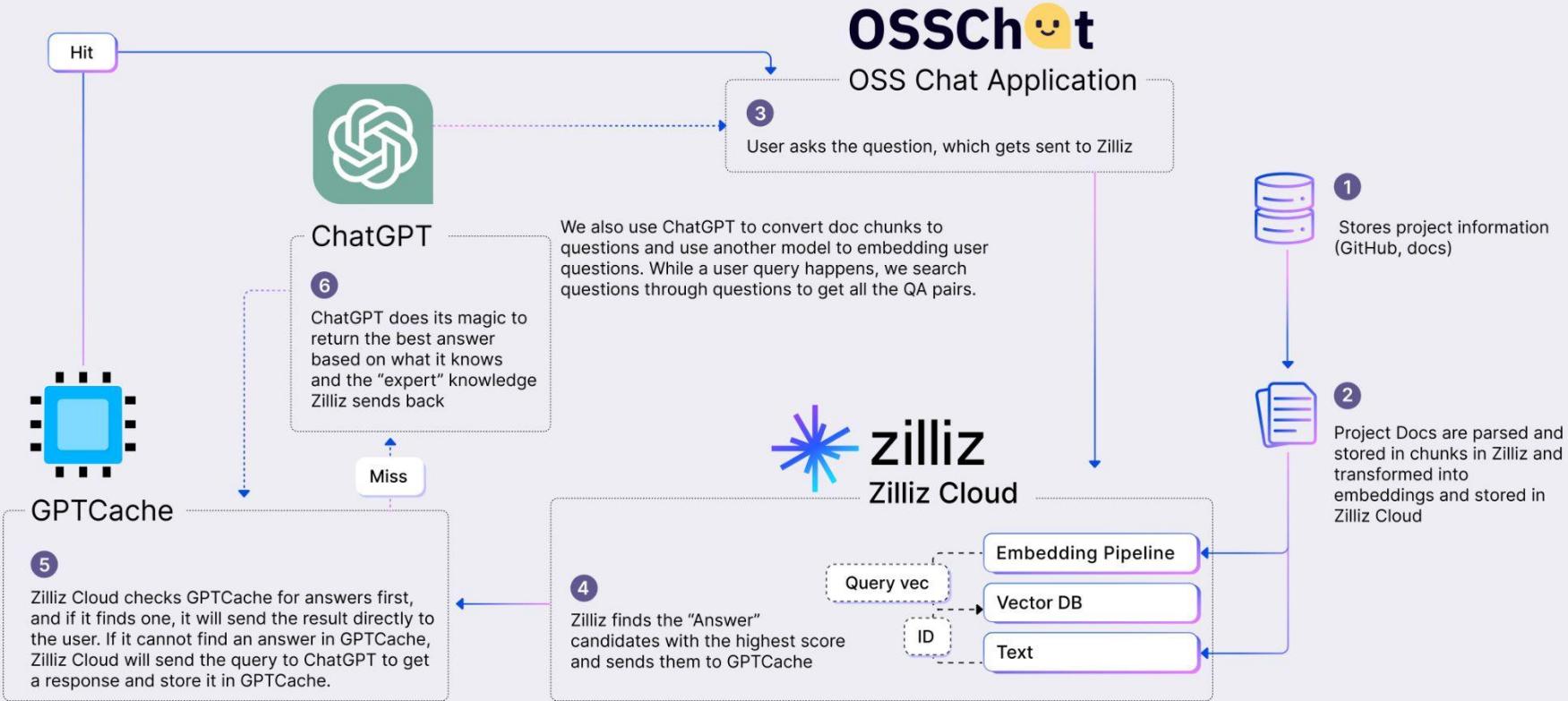
Lots
of Slides

Cool Demo



<https://multimodal-demo.milvus.io/>

<https://milvus.io/milvus-demos/reverse-image-search>



Unstructured Data is Everywhere

Unstructured data is any data that does not conform to a predefined data model.

Currently, 90% of unstructured data is never analyzed.



Text



Images



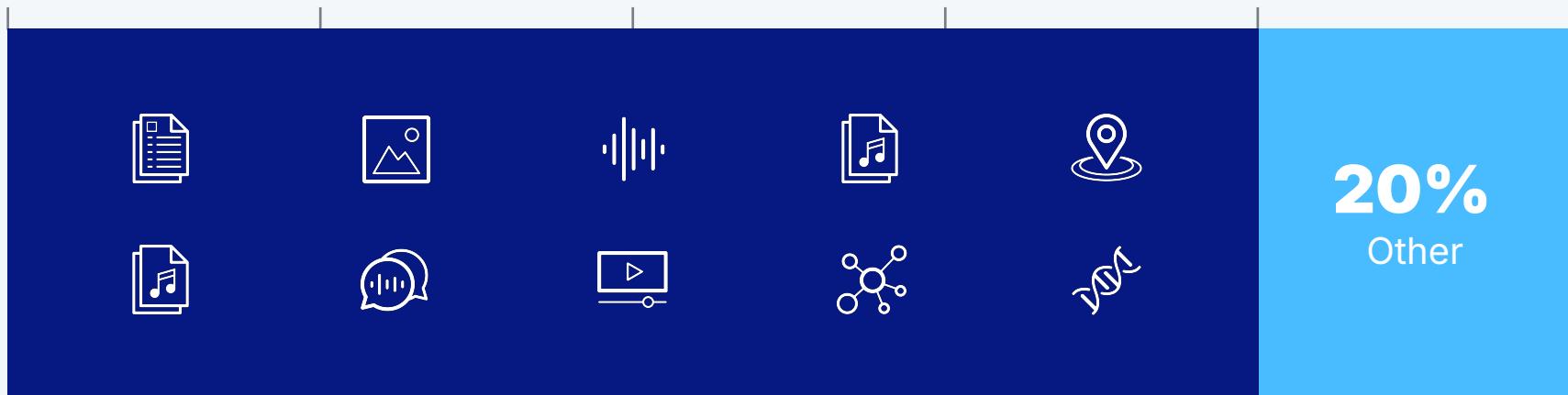
Videos



and more!

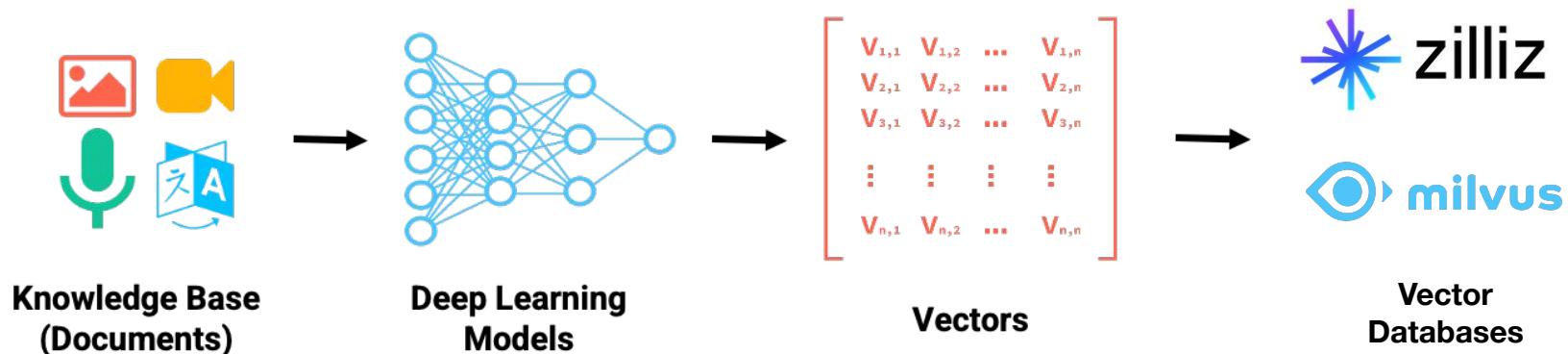
...and cannot process increasingly growing unstructured data

< 80% newly generated data in 2025
will be unstructured data



The challenge of unstructured data

- **Problem:** Unstructured data comes in lots of forms, no easy way to interact with it all
- **Solution:** Vector embeddings
- **How:** Neural networks e.g. embedding models



02

Overview of Vector Databases

Why a Vector Database?

Purpose-built to store, index and query vector embeddings from unstructured data.

- Vector database

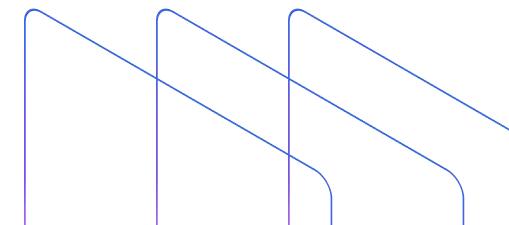
- Advanced filtering (filtered vector search, chained filters)
- Hybrid search (e.g. full text + dense vector)
- Durability (any write in a db is durable, a library typically only supports snapshotting)
- Replication / High Availability
- Sharding
- Aggregations or faceted search
- Backups
- Lifecycle management (CRUD, Batch delete, dropping whole indexes, reindexing)
- Multi-tenancy

- Vector search library

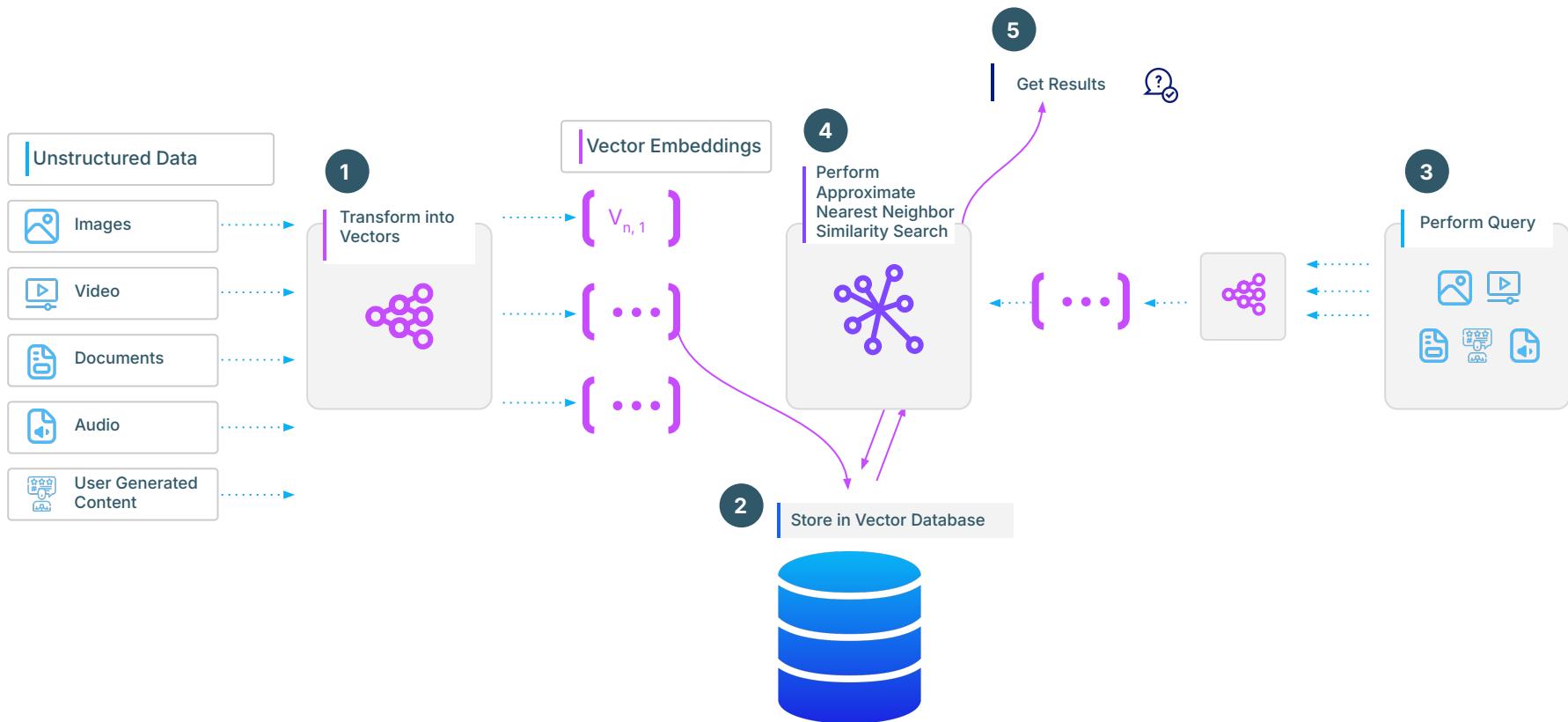
- High-performance vector search

- How do I support different applications?

- High query load
- High insertion/deletion
- Full precision/recall
- Accelerator support (GPU, FPGA)
- Billion-scale storage

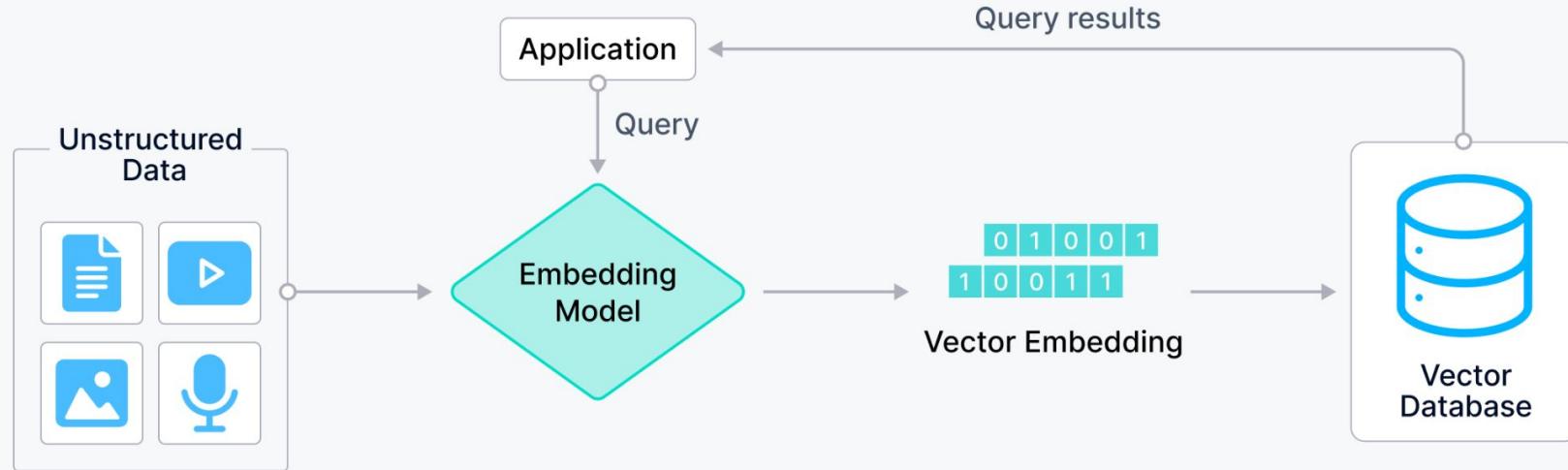


How Similarity Search Works



Vector Database: Making Sense of Unstructured Data

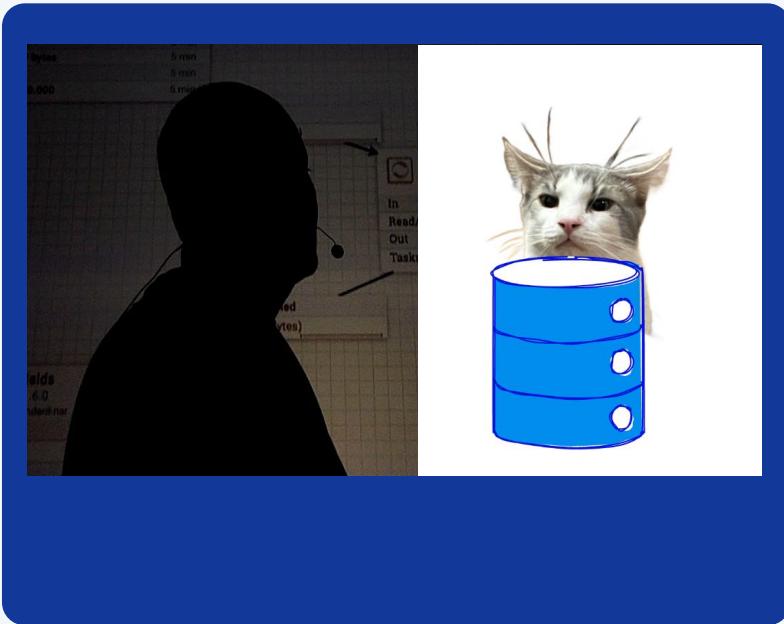
A vector database stores embedding vectors and allows for semantic retrieval of various types of unstructured data.



03

A Quick Introduction to Milvus

Milvus Features



**Scalable and Elastic
Architecture**

**Diverse Index
Support**

**Versatile Search
Capabilities**

**Tunable
Consistency**



Multi-Tenancy

**Hardware-
Accelerated
Compute Support**

**Python, Java,
Golang, NodeJS**

**Milvus Lite, K8,
Zilliz Cloud, Docker**



Technologies for various types of Use cases



Index Types

Offer a wide range of **15 indexes** support, including popular ones like Hierarchical Navigable Small Worlds (HNSW), PQ, Binary, Sparse, DiskANN and GPU index

Empower developers with tailored search optimizations, catering to performance, accuracy and cost needs



Search Types

Support multiple types such as **top-K ANN, Range ANN, sparse & dense, multi-vector, grouping, and metadata filtering**

Enable query flexibility and accuracy, allowing developers to tailor their information retrieval needs



Multi-tenancy

Enable **multi-tenancy** through collection and partition management

Allow for efficient resource utilization and customizable data segregation, ensuring secure and isolated data handling for each tenant

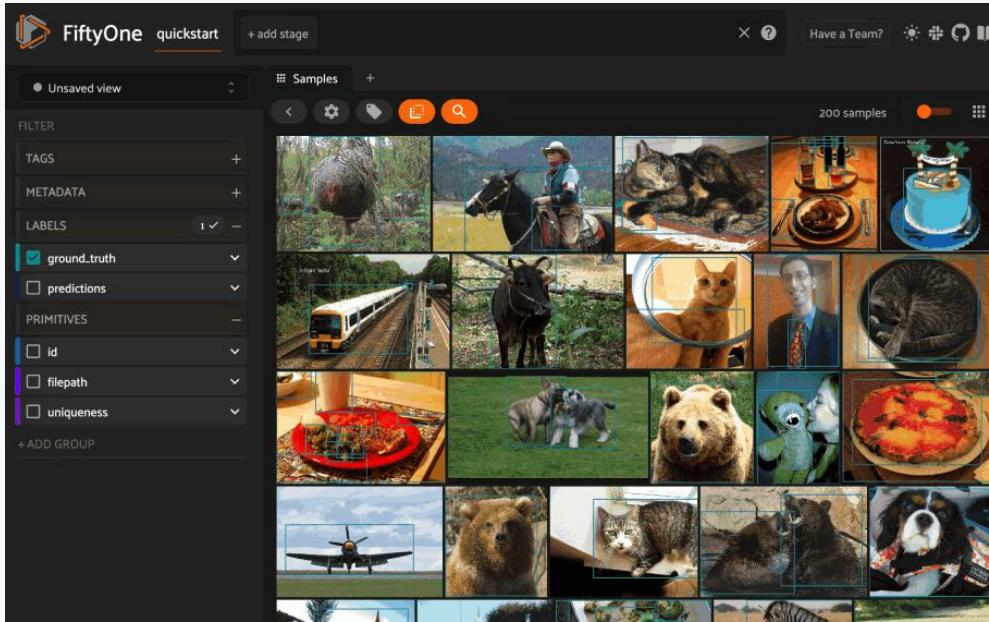


Compute Types

Designed for various compute powers, such as **AVX512, Neon for SIMD, quantization cache-aware optimization and GPU**

Leverage strengths of each hardware type, ensuring high-speed processing and cost-effective scalability for different application needs

FiveOne + Milvus



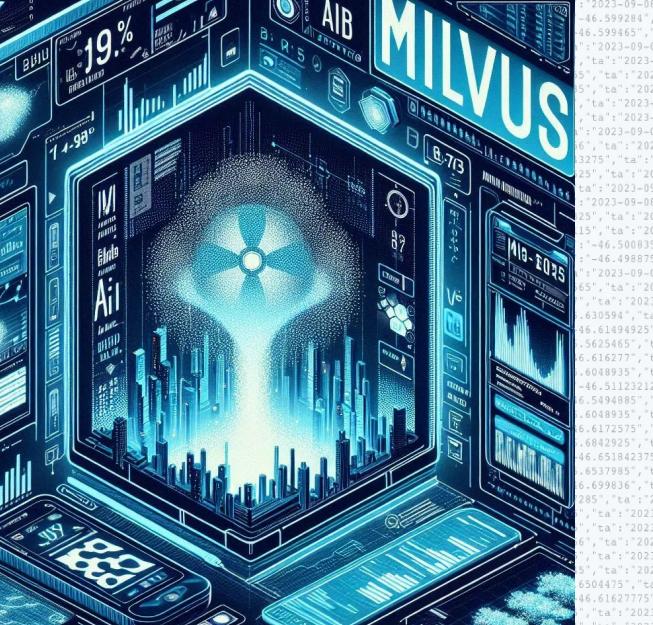
<https://docs.voxel51.com/integrations/milvus.html>

https://milvus.io/docs/integrate_with_voxel51.md

04

Consume and Ingest Air Quality Data

DATA!!!!



DATA



 zilliz

REST JSON Ingest

```
url = 'https://api.openaq.org/v2/measurements?country=US&date_from={0}&date_to={1}&limit=1000&page={2}&offset=0&sort=desc&radius=1000&order_by=datetime'.format(
    str(YESTERDAYS_DATE), str(TODAYS_DATE), str(aqpage) )
headers = {"accept": "application/json", "x-api-key": str(API_KEY)}

response = requests.get(url, headers=headers)

if ( len(response.text) > 0 ):
    openaq2 = json.loads(response.text)
else:
    openaq2['results'] = {}

try:
    for jsonitems in openaq2['results']:
        count = count + 1
        fullolocation = 'Location {0}: {1}, {2}, {3} @ {4},{5}'.format(jsonitems.get('locationId'),
                                                                    jsonitems.get('location'),jsonitems.get('city'),jsonitems.get('country'),
                                                                    jsonitems.get('coordinates')['latitude'],jsonitems.get('coordinates')['longitude'] )
        details = 'Current Air Quality Reading for {0} is {1} {2} for {3} at {4}. Is Mobile: {5} Is Analysis: {6} Entity: {7} Sensor Type: {8}'.format(jsonitems.get('parameter'),
jsonitems.get('unit'),fullolocation, jsonitems.get('date')['local'],
jsonitems.get('isMobile'), jsonitems.get('isAnalysis'), jsonitems.get('entity'), jsonitems.get('sensorType'))

        data.append({ "locationId": int(jsonitems.get('locationId')),
                      "location": str(jsonitems.get('location','')),
                      "parameter": str(jsonitems.get('parameter','')),
                      "value": float(jsonitems.get('value')),
                      "datelocal": str(jsonitems.get('date')['local']),
                      "unit": str(jsonitems.get('unit','')),
                      "latitude": float(jsonitems.get('coordinates')['latitude']),
                      "longitude": float(jsonitems.get('coordinates')['longitude']),
                      "country": str(jsonitems.get('country','')),
                      "city": str(jsonitems.get('city','')),
                      "isMobile": str(jsonitems.get('isMobile','')),
                      "isAnalysis": str(jsonitems.get('isAnalysis','')),
                      "entity": str(jsonitems.get('entity','')),
                      "sensorType": str(jsonitems.get('sensorType','')),
                      "vector": model(details), "details": str(details), "fulllocation": str(fullolocation))}

res = milvus_client.insert(collection_name=OAQ3_COLLECTION_NAME, data=data)
print(count)
except Exception as ex:
```

Scalars and Vectors in Milvus

The screenshot shows the Milvus UI interface. On the left, there's a sidebar with icons for Home, Databases, Collections, and a search bar. The main area shows the 'Database default' and 'Collection airquality' selected. The 'Data' tab is active, displaying a table with 34 entities. Each entity has a vector ID and a detailed description of its air quality reading. The table includes columns for 'vector', 'details', and 'location'. A query bar at the top allows for filtering by consistency (Bounded) and entering a data query. The bottom of the page shows navigation links and a footer.

vector	details	location
39957,-0.022822190076...	Current Air Quality Reading for O3 is 32 for Location Central NJ @ 40.401,-74.325 at Hour 15 on 2024-08-17.	Location Central NJ @ 40.401,-74.325
3004,-0.0165540408343...	Current Air Quality Reading for PM2.5 is 67 for Location Central NJ @ 40.401,-74.325 at Hour 15 on 2024-08-17.	Location Central NJ @ 40.401,-74.325
813995,-0.027215896174...	Current Air Quality Reading for O3 is 15 for Location Redwood City CA @ 37.48,-122.22 at Hour 12 on 2024-08-17.	Location Redwood City CA @ 37.48,-122.22
4523,-0.0204664301127...	Current Air Quality Reading for PM2.5 is 0 for Location Redwood City CA @ 37.48,-122.22 at Hour 12 on 2024-08-17.	Location Redwood City CA @ 37.48,-122.22
35587,-0.054002948105...	Current Air Quality Reading for O3 is 28 for Location Boston MA @ 42.351,-71.051 at Hour 15 on 2024-08-17.	Location Boston MA @ 42.351,-71.051
15228,-0.038178324699...	Current Air Quality Reading for PM2.5 is 76 for Location Boston MA @ 42.351,-71.051 at Hour 15 on 2024-08-17.	Location Boston MA @ 42.351,-71.051
'39334,-0.045572765171...	Current Air Quality Reading for PM10 is 29 for Location Boston MA @ 42.351,-71.051 at Hour 15 on 2024-08-17.	Location Boston MA @ 42.351,-71.051
30693,-0.055957559496...	Current Air Quality Reading for O3 is 40 for Location Madison CT @ 41.2583,-72.5506 at Hour 15 on 2024-08-17.	Location Madison CT @ 41.2583,-72.5506
125964,-0.04203481972...	Current Air Quality Reading for O3 is 14 for Location San Rafael CA @ 37.97,-122.52 at Hour 12 on 2024-08-17.	Location San Rafael CA @ 37.97,-122.52
86656,-0.043939355760...	Current Air Quality Reading for PM2.5 is 39 for Location San Rafael CA @ 37.97,-122.52 at Hour 12 on 2024-08-17.	Location San Rafael CA @ 37.97,-122.52
383163,-0.062174882739...	Current Air Quality Reading for O3 is 35 for Location NW Coastal LA CA @ 34.0505,-118.4566 at Hour 12 on 2024-08-17.	Location NW Coastal LA CA @ 34.0505,-118.4566
813995,-0.027215896174...	Current Air Quality Reading for O3 is 15 for Location Redwood City CA @ 37.48,-122.22 at Hour 12 on 2024-08-17.	Location Redwood City CA @ 37.48,-122.22
4523,-0.0204664301127...	Current Air Quality Reading for PM2.5 is 0 for Location Redwood City CA @ 37.48,-122.22 at Hour 12 on 2024-08-17.	Location Redwood City CA @ 37.48,-122.22
1687181,-0.03408430889...	Current Air Quality Reading for O3 is 46 for Location Northeast Urban NJ @ 40.692,-74.187 at Hour 15 on 2024-08-17.	Location Northeast Urban NJ @ 40.692,-74.187

1 - 14 of 34 Entities (15 ms)

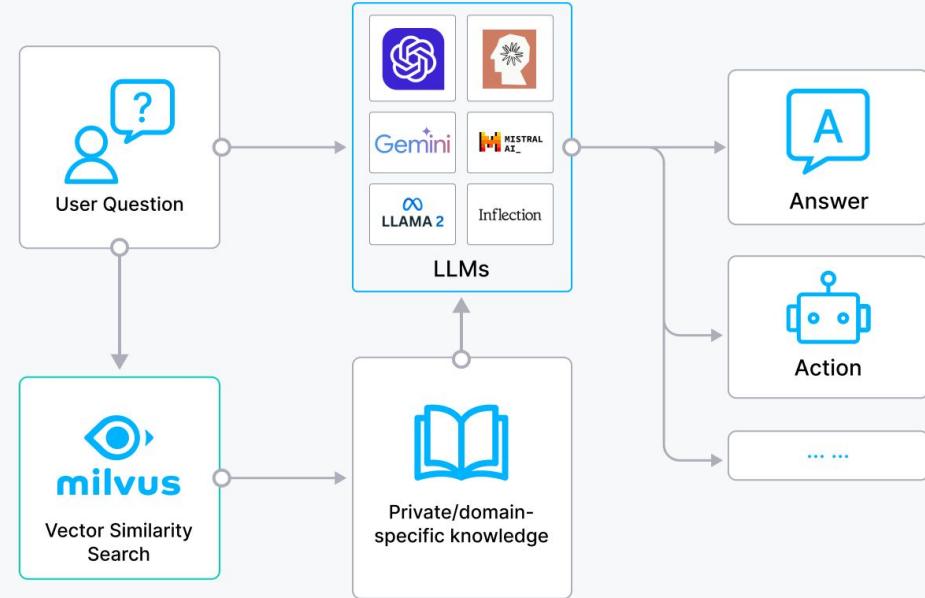
05

Building a local RAG application

Retrieval-Augmented Generation (RAG)

A technique that combines the strength of retrieval-based and generative models:

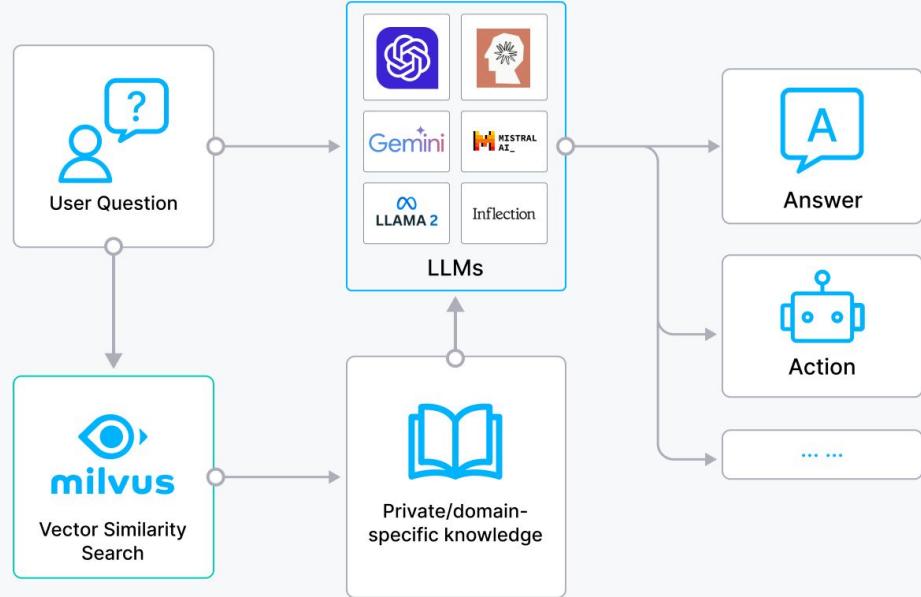
- Improve accuracy and relevance
- Eliminate hallucination
- Provide domain-specific knowledge



RAG : an economic perspective

A business model that bridges public data and private data

- Data sovereignty
- You can't and shouldn't give your private data to others



Ollama + Llama 3.1 + Milvus + LangChain = RAG

```
# -----
# Connect to Milvus
# https://zilliz.com/blog/a-beginners-guide-to-using-llama-3-with-ollama-milvus-langchain
# https://github.com/stephen37/ollama_local_rag/blob/main/rag_berlin_parliament.py
# https://python.langchain.com/v0.2/docs/integrations/vectorstores/milvus/
# https://api.python.langchain.com/en/latest/vectorstores/langchain_milvus.vectorstores.Milvus.html

vector_store = Milvus(
    embedding_function=embeddings,
    collection_name=QAQ3_COLLECTION_NAME,
    primary_field = "id",
    text_field="details",
    connection_args={"uri": MILVUS_URL},
)

def run_query() -> None:
    llm = Ollama(
        model="llama3.1",
        callback_manager=CallbackManager([StreamingStdOutCallbackHandler()]),
        stop=["||"],
    )

    query = input("\nQuery: ")
    qa_chain = RetrievalQA.from_chain_type(
        llm, retriever=vector_store.as_retriever(collection = QAQ3_COLLECTION_NAME, expr="parameter == 'o3'"))

    result = qa_chain.invoke({"query": query})

    resultforslack = str(result["result"])
    print(resultforslack)

    try:
        response = client.chat_postMessage(mrkdnw=True, channel="C06NE1FU6SE", text="",
                                             blocks=[{"type": "section", "text": {"type": "mrkdwn", "text": "*" + str(query) + "* \n\n" + str(resultforslack) + "\n\n"}}])

    except SlackApiError as e:
        # You will get a SlackApiError if "ok" is False
        print("Slack failed")

while True:
    run_query()
```

Ollama + Llama 3.1 + Milvus + LangChain = RAG

Query: Provide a detailed air quality report for Edison

Based on the provided context, here is a detailed air quality report for Edison:

Location: Edison, US (Latitude: 35.345607, Longitude: -118.851825)

Current Air Quality Reading for 03:

* **Highest Value:** 0.092 ppm recorded at 2024-08-30T15:00:00-07:00

* **Lowest Value:** 0.052 ppm recorded at 2024-08-29T19:00:00-07:00

* **Current Value:** 0.062 ppm recorded at 2024-08-30T11:00:00-07:00

Trend Analysis:

* The current value is slightly higher than the lowest value recorded, indicating a slight increase in ozone levels.

* However, it's still within a relatively safe range and not exceeding any health or safety thresholds.

Sensor Information:

* **Sensor Type:** Reference grade

* **Entity:** Governmental Organization

Mobility: The sensor is stationary (Mobile: False), indicating that the readings are taken from a fixed location. Based on the provided context, here is a detailed air quality report for Edison:

Location: Edison, US (Latitude: 35.345607, Longitude: -118.851825)

Current Air Quality Reading for 03:

* **Highest Value:** 0.092 ppm recorded at 2024-08-30T15:00:00-07:00

* **Lowest Value:** 0.052 ppm recorded at 2024-08-29T19:00:00-07:00

* **Current Value:** 0.062 ppm recorded at 2024-08-30T11:00:00-07:00

Trend Analysis:

* The current value is slightly higher than the lowest value recorded, indicating a slight increase in ozone levels.

* However, it's still within a relatively safe range and not exceeding any health or safety thresholds.

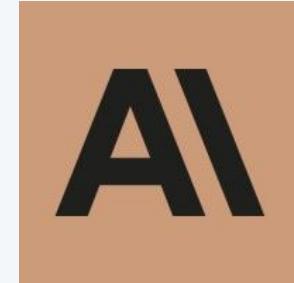
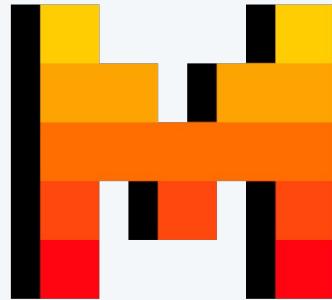
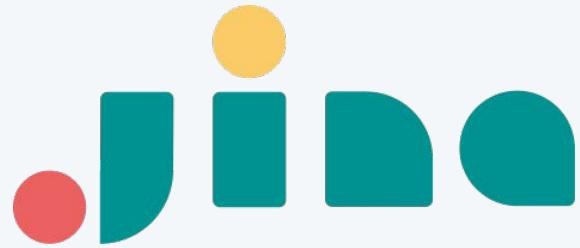
Sensor Information:

* **Sensor Type:** Reference grade

* **Entity:** Governmental Organization

Mobility: The sensor is stationary (Mobile: False), indicating that the readings are taken from a fixed location.

Embeddings Models



06

Q & A



<https://medium.com/@tspann/whats-in-the-air-tonight-mr-milvus-fbd42f06e482>

RESOURCES



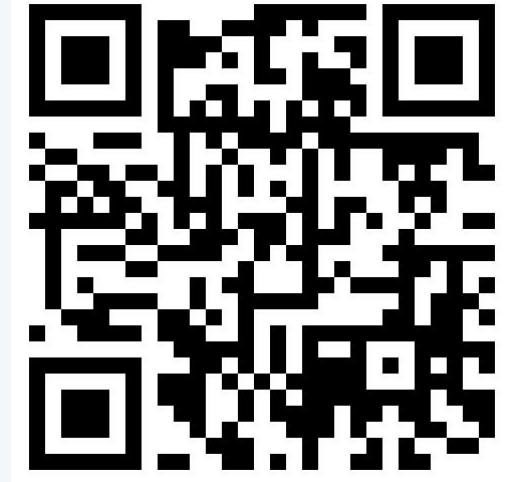
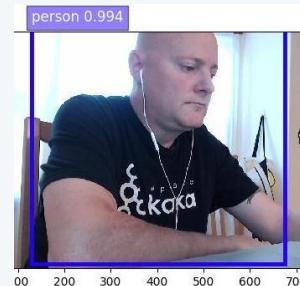
Vector Database Resources

Give Milvus a Star!

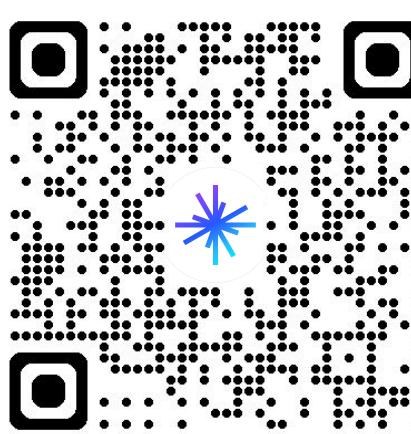


<https://github.com/milvus-io/milvus>

Chat with me on Discord!



Unstructured Data Meetup



<https://www.meetup.com/unstructured-data-meetup-new-york/>

This meetup is for people working in unstructured data. Speakers will come present about related topics such as vector databases, LLMs, and managing data at scale. The intended audience of this group includes roles like machine learning engineers, data scientists, data engineers, software engineers, and PMs.

This meetup was formerly Milvus Meetup, and is sponsored by [Zilliz](#) maintainers of [Milvus](#).

Generative AI Resource Hub

Tutorials, Code Examples, and Best Practices for Developing and Deploying GenAI Applications.



Learn



Build



Explore

<https://zilliz.com/learn/generative-ai>



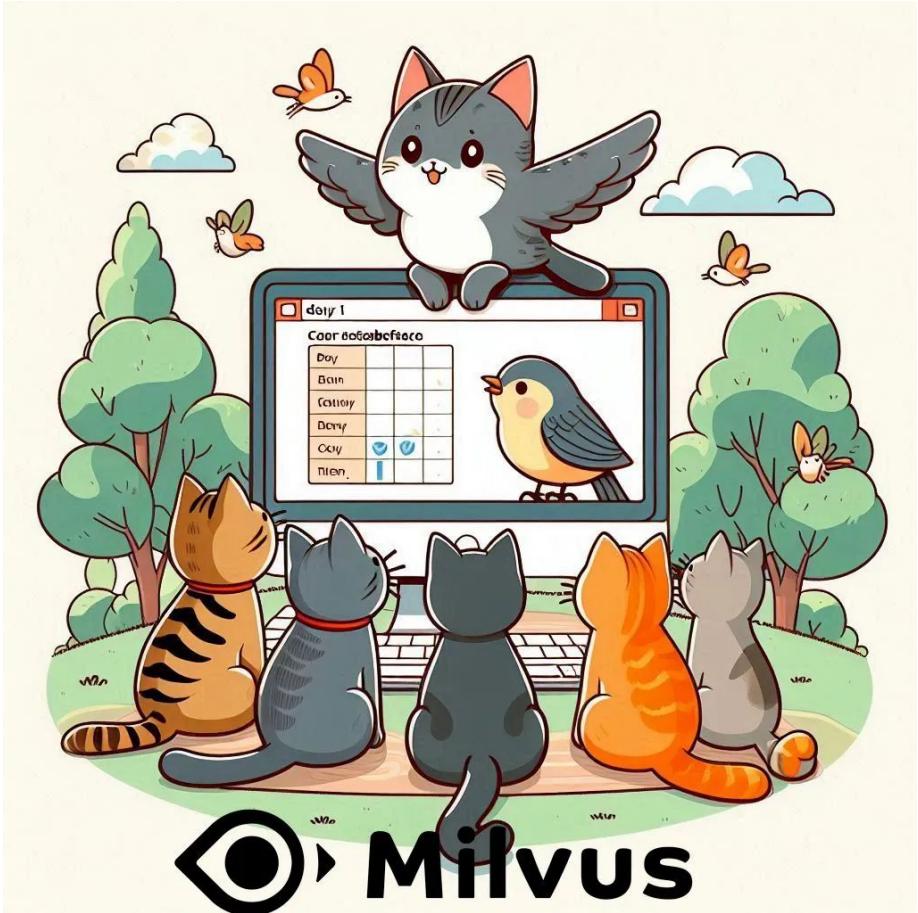




<https://medium.com/@tspann/unstructured-street-data-in-new-york-8d3cde0a1e5b>



<https://medium.com/@tspann/not-every-field-is-just-text-numbers-or-vectors-976231e90e4d>



<https://medium.com/@tspann/shining-some-light-on-the-new-milvus-lite-5a0565eb5dd9>

Raspberry Pi AI Kit - Hailo
Edge AI



Milvus



<https://medium.com/@tspann/unstructured-data-processing-with-a-raspberry-pi-ai-kit-c959dd7fff47>

AIM Weekly by Tim Spann



<https://bit.ly/32dAJft>

<https://github.com/milvus-io/milvus>

This week in Milvus, Towhee, Attu, GPT Cache, Gen AI, LLM, Apache NiFi, Apache Flink, Apache Kafka, ML, AI, Apache Spark, Apache Iceberg, Python, Java, Vector DB and Open Source friends.

Thank you!



milvus.io



github.com/milvus-io/



[@milvusio](https://twitter.com/milvusio)

Connect with me!



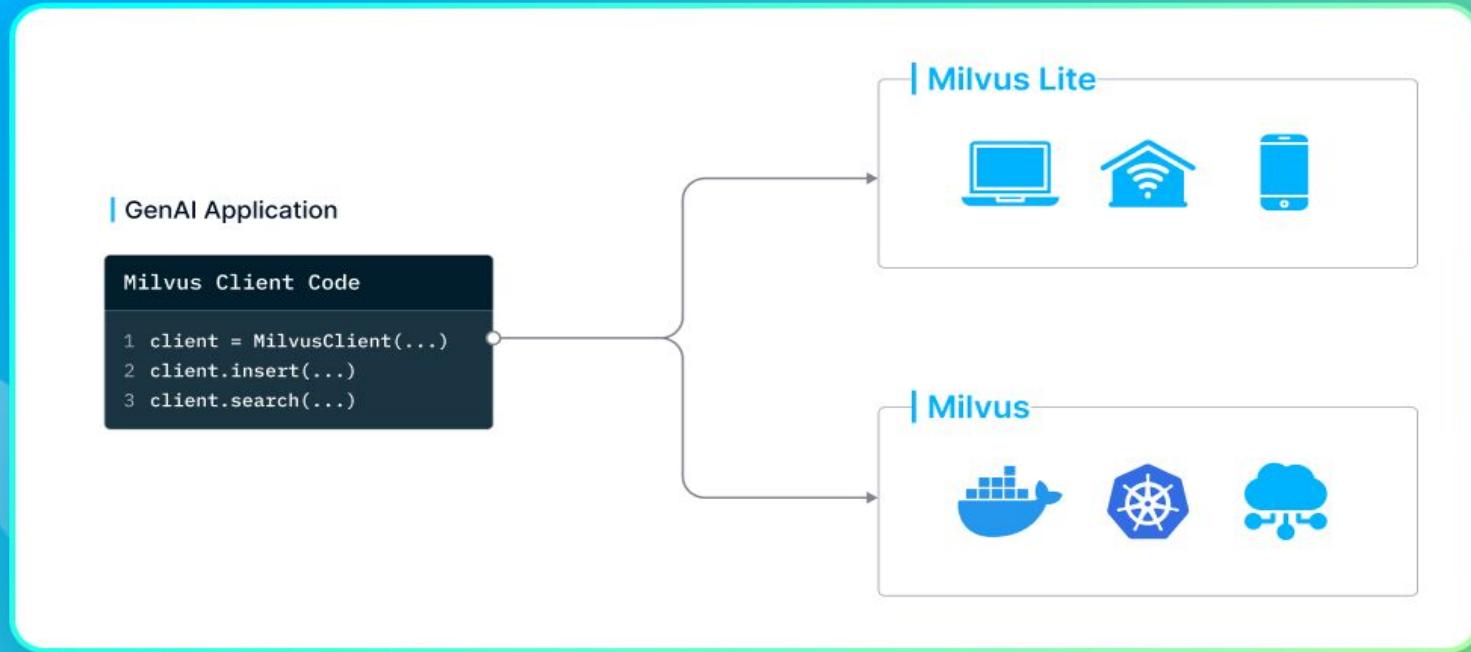
[@paasDev](https://twitter.com/paasDev)



[/in/timothyspann](https://www.linkedin.com/in/timothyspann/)



Build Once Deploy Anywhere



BEFORE MILVUS





Join us at our next meetup!
meetup.com/unstructured-data-meetup-new-york/

THANK YOU

