#### CS4221 Tutorial 4: Vector Databases - Milvus

Yao LU 2024 Semester 2

National University of Singapore School of Computing

# Objective

By the end of this tutorial, you will:

- Set up Milvus server using Docker.
- > Prepare the virtual environment and interact with Milvus server.
- > Develop a fundamental understanding of schema design principles in Milvus.
- > Perform basic operations such as create, insert, search and delete data.
- > Explore various index types, fine-tune parameters, and evaluate search performance.

#### Setup

#### **Step 1: Install Milvus**

- Download the provided docker-compose file.
- Start up the docker container and verify Milvus server is running.

#### **Step 2: Prepare the environment**

- Make sure you have installed the *anaconda*
- Download the provided environment file and set up the virtual environment, or follow the documentation to manually download required dependencies.
- Execute quickStart.ipynb to make sure everything is ok.

# Prerequisites

Before starting, ensure you have gone through :

- Milvus concept overview:
- https://milvus.io/docs/manage-collections.md
- https://milvus.io/docs/schema.md
- Index related documentations
- IVF\_FLAT: <a href="https://milvus.io/docs/ivf-flat.md">https://milvus.io/docs/ivf-flat.md</a>
- HNSW: https://milvus.io/docs/hnsw.md
- Blog: <u>https://zilliz.com/learn/how-to-pick-a-vector-index-in-milvus-visual-guide</u>

# Prerequisites

Before starting, ensure you have gone through :

- Milvus concept overview
- https://milvus.io/docs/manage-collections.md
- https://milvus.io/docs/schema.md
- Index related documentations
- Metric type: <u>https://milvus.io/docs/metric.md</u>
- IVF\_FLAT: <a href="https://milvus.io/docs/ivf-flat.md">https://milvus.io/docs/ivf-flat.md</a>
- HNSW: https://milvus.io/docs/hnsw.md
- Blog: <u>https://zilliz.com/learn/how-to-pick-a-vector-index-in-milvus-visual-guide</u>

# Hello-World

After you start the provided docker-compose, you will get

lingzo@worker 012. /cc/221/woot	or dh tutor	ialt dackar	compoco	un d	
<pre>lingze@worker-012:~/cs4221/vect [+] Running 23/23</pre>			compose	up -u	
✓ standalone 7 layers [########	] 0B/0I	B Pulle	be		
✓ d5fd17ec1767 Pull complete			su		
✓ 0f5a22c44678 Pull complete					
✓ 72ad4f350efb Pull complete					
<ul> <li>2f5ee08a99b8 Pull complete</li> </ul>					
✓ 3fe28e251347 Pull complete					
✓ 1f27396f6efc Pull complete					
✓ fe556ec02776 Pull complete					
<pre> v etcd 7 layers [[]]] </pre>	0B/0B	Pulled			
✓ dbba69284b27 Pull complete	,				
✓ 270b322b3c62 Pull complete					
✓ 7c21e2da1038 Pull complete					
✓ cb4f77bfee6c Pull complete					
✓ e5485096ca5d Pull complete					
✓ 3ea3736f61e1 Pull complete					
✓ 1e815a2c4f55 Pull complete					
✓ minio 6 layers [░░░░░:)	0B/0B	Pulled			
🖌 🗸 c7e856e03741 Pull complete					
<pre>✓ c1ff217ec952 Pull complete</pre>					
✓ b12cc8972a67 Pull complete					
✓ 4324e307ea00 Pull complete					
152089595ebc Pull complete					
Ø5f217fb8612 Pull complete					
[+] Running 4/4					
<ul> <li>Network milvus</li> </ul>	Created				
<ul> <li>Container milvus-etcd</li> </ul>	Started				
	Started				
🖌 Container milvus-standalone	Started				

# Hello-World

Follow the documentation to prepare you environment.

Before executing the scripts, something to note:

- Update the server address (HOST)
- > No restriction on the database name

HOST = '10.10.10.250'
PORT = 19530
DB\_NAME = 'testdb'
URL="http://"+HOST+':'+str(PORT)
# if you deployed the standalone milvus, and connect to the database server
# if you deployed in your local machine, use "http://localhost:19530"
client = MilvusClient(URL)

## Hello-World

Before executing the scripts, something to note:

- > For the embedding model example, ensure *pymilvus[model]* dependency installed.
- > Otherwise, refer to the following synthetic data example.

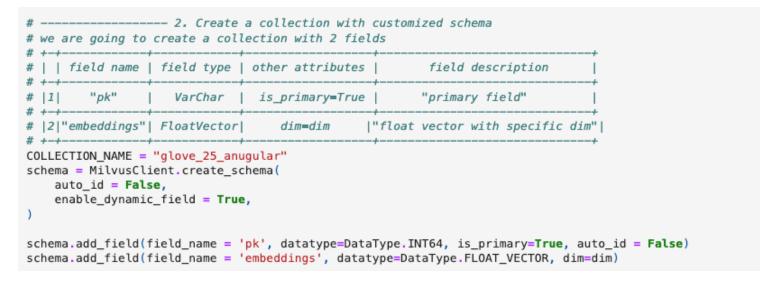
```
# generate vectors for the documents
# pre-requisite is to download pymilvus[model] first
embedding_fn = model.DefaultEmbeddingFunction()
vectors = embedding_fn.encode_documents(docs)
print("Dim:", embedding_fn.dim, vectors[0].shape)
# Each record is named "entity"m entity has id, vector representation, raw text, and a subject label that we use
# the usage is similar to the NoSQL database like mongodb.
data = [
        {"id":i, "vector": vectors[i], "text": docs[i], "subject":"history"}
        for i in range(len(docs))
]
print("Data has", len(data), "entities, each with fields: ", data[0].keys())
print("Vector dim:", len(data[0]["vector"]))
```

In this tutorial, we explore different search indexes and fine-tune them to evaluate retrieval performance on a real-world dataset.

For details, refer to *fine\_tune\_index.ipynb*.

We focus on two fundamental indexes: **IVF\_FLAT** and **HNSW**, using the **Glove-25-angular** dataset. This dataset has a dimension of **25**, with **1,183,514** training samples and **10,000** test samples

The schema design is as follows:



Take the Inverted File Flat index (**IVF\_FLAT**) as an example.

**IVF\_FLAT** offers two tunable hyperparameters:

- **nlist**: the number of partitions to create using the k-means algorithm.
- **nprobe**: the number of partitions to consider during the search for candidate

The **nlist** parameter is set when building the **IVF\_FLAT** index, while **nprobe** is adjusted dynamically for each query request.

First, we fix **nlist** and tune **nprobe** to evaluate query performance in terms of **latency** and **recall**. We build the index by specifying:

- The field for the index\_type (in this case, "embedding").
- The metric\_type (here, we use COSINE).
- The **index type** and **nlist** parameter.

```
# first, we fix nlist, tune nprobe to check the query performance (latency and recall)
nlist = 1024
index_params = MilvusClient.prepare_index_params()
index_params.add_index(
    field_name="embeddings",
    metric_type = "COSINE",
    # related distance metric to angular is CONSINE
    index_type = "IVF_FLAT",
    index_name = "vector_index",
    params = {
        "nlist":nlist
start_time = time.time()
client.create_index(collection_name = COLLECTION_NAME,
                   index_params = index_params,
                   sync = True)
end_time = time.time()
print(f"create index time: {end_time-start_time:.4f}s")
```

Next, we tune **nprobe** while executing the same search queries to analyze its impact on query latency and average recall.

```
=== ------ nprobe=1 search ------ ===
# ------ 4[IVF_FLAT]. Search with different nprobe, check the performance
for nprobe in [1,16, 64, 256, 1024]:
   start time = time.time()
                                                                                   search latency = 3.8174s
   res = client.search(
                                                                                   Mean Average Recall = 0.3681
      collection_name=COLLECTION_NAME,
      data=query_embedding,
                                                                                   === ----- nprobe=16 search ----- ===
      limit = 100,
      search_params={
                                                                                   search latency = 4.2593s
          "params" : {"nprobe":nprobe}
                                                                                   Mean Average Recall = 0.8873
   end time = time.time()
                                                                                   === ----- nprobe=64 search ----- ===
   print(fmt.format(f"------ nprobe={nprobe} search ------"))
   print(search_latency_fmt.format(end_time-start_time))
                                                                                   search latency = 6.7913s
                                                                                   Mean Average Recall = 0.9804
   # calculate the Mean Average Recall.
   # Recall@K = (# of true positive in top K) / (# of true positive)
                                                                                   === ----- nprobe=256 search ----- ===
   # MAR (Mean Average Recall) = 1/C * sum(Recall@K). C is the number of queries
   mar = []
                                                                                   search latency = 13.6901s
   for i, candidate_res in enumerate(res):
      y = neighbors[i]
                                                                                   Mean Average Recall = 0.9994
      y_ = [j['id'] for j in candidate_res]
      y, y_{-} = set(y), set(y_{-})
                                                                                   === ------ nprobe=1024 search ------ ===
      mar_.append(1.0 * len(y & y_) / len(y))
   mar = np.mean(mar_)
                                                                                   search latency = 43.0186s
   print(f"Mean Average Recall = {mar:.4f}")
                                                                                   Mean Average Recall = 1.0000
```

From this result, we can observe that as the **nprobe** increases, the **search latency** increase while the average recall degrades.

Next, we fix the **nprobe** and build the index with different **nlist** values to analyze its impact

# 4[IVF_FLAT]. Search with different nprobe, check the performance	=== nprobe=1 search ===
for nprobe in [1,16, 64, 256, 1024]:	
<pre>start_time = time.time()</pre>	search latency = 3.8174s
res = client.search(	Mean Average Recall = 0.3681
collection_name=COLLECTION_NAME,	-
data=query_embedding,	===
limit = 100,	
search_params={	course latency = 4 2502c
"params" : {"nprobe":nprobe}	search latency = 4.2593s
}	Mean Average Recall = 0.8873
)	
end_time = time.time()	=== nprobe=64 search
<pre>print(fmt.format(f""nprobe={nprobe} search"))</pre>	
<pre>print(search_latency_fmt.format(end_time-start_time))</pre>	search latency = $6.7913s$
	Mean Average Recall = 0.9804
# calculate the Mean Average Recall.	
# Recall@K = (# of true positive in top K) / (# of true positive)	=== nprobe=256 search
# MAR (Mean Average Recall) = 1/C * sum(Recall@K). C is the number of queries	
mar_ = []	
for i, candidate_res in enumerate(res):	search latency = 13.6901s
y = neighbors[i]	Mean Average Recall = 0.9994
<pre>y_ = [j['id'] for j in candidate_res]</pre>	
$y, y_{-} = set(y), set(y_{-})$	=== nprobe=1024 search
marappend(1.0 * len(y & y_) / len(y))	
<pre>mar = np.mean(mar_)</pre>	search latency = 43.0186s
print(f"Mean Average Recall = {mar:.4f}")	Mean Average Recall = 1.0000
	Hear Average Recatt - 110000

# Homework

In this homework, we try to tune the parameters of HNSW and analyze the impact on the query efficiency and accuracy.

**HNSW** offers three tunable hyperparameters:

- M: the maximum number of connections for each node in the graph
- **efConstruction**: the size of the dynamic candidate list which controls index search speed/build tradeoff.
- ef: the size of the dynamic candidate list during search

#### Task:

Follow the process of analyzing **IVF\_FLAT**, explore the impact and summarize the trends observed with varying values for these parameters.