



**Developing Open LLM
applications with**



Apache OpenServerless

Lesson 6
Vector Database

Vector Database

- Vector Database
- Embedding and Search
- Importing PDF



Vector Database

Milvus: concepts

- Milvus is a No-SQL database optimized for **vector searches**
 - finding **similarities** in a dataset looking at their **numeric representation**
 - text is transformed in a numeric representatino using an **embedding model** then stored
- Multiple **databases**, each database has multiple **collections**
- Each collection has a **schema** and **indexes**

Access Milvus

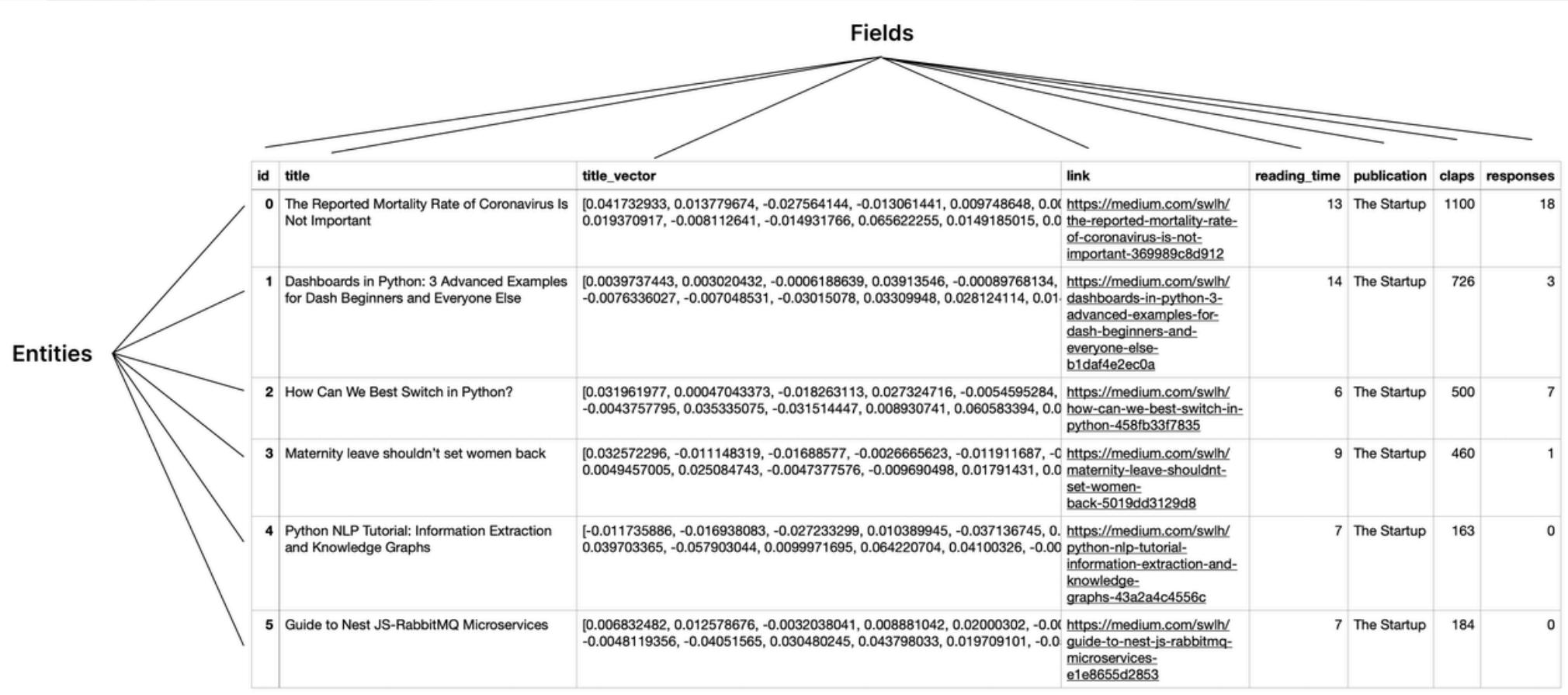
```
import os
from pymilvus import MilvusClient

uri = f"http://{{os.getenv("MILVUS_HOST")}}"
token = os.getenv("MILVUS_TOKEN")
db_name = os.getenv("MILVUS_DB_NAME")
client = MilvusClient(uri=uri, token=token, db_name=db_name)
```

- Using Milvus

```
client.list_collections()
client.drop_collection("test")
```

Collections



Create Schema

- Parameters

```
from pymilvus import DataType
COLLECTION = "test"
DIMENSION=1024
```

- Define Schema

```
schema = client.create_schema()
schema.add_field(field_name="id", datatype=DataType.INT64, is_primary=True, auto_id=True)
schema.add_field(field_name="text", datatype=DataType.VARCHAR, max_length=DIMENSION)
schema.add_field(field_name="embeddings", datatype=DataType.FLOAT_VECTOR, dim=DIMENSION)
```

Create Index and Collection

- Define Index

```
index_params = client.prepare_index_params()  
index_params.add_index("embeddings", index_type="AUTOINDEX", metric_type="IP")
```

- Create Collection with Index and Schema

```
client.create_collection(  
    collection_name=COLLECTION,  
    schema=schema, index_params=index_params)
```

Insert

```
text = "Hello World"
vec = [float(i) for i in range(0,DIMENSION)] # TO BE REPLACED with embedding
client.insert(COLLECTION, {"text":text, "embeddings": vec})
```

Retrieve

```
qit = client.query_iterator(collection_name=COLLECTION, batchSize=2, output_fields=[ "text"])
res = qit.next()
print(res[0].get("text"))
```

Hello World

Embedding

Embedding

- Use an embedding model

```
import sys, requests as req
MODEL="mxbai-embed-large:latest"
DIMENSION=1024
```

- Invoke the embedding API

```
inp = "Hello World"
url = f"https://{{os.getenv("AUTH")}}@{{os.getenv("OLLAMA_HOST")}}/api/embeddings"
msg = { "model": MODEL, "prompt": inp, "stream": False }
res = req.post(url, json=msg).json()
out = res.get('embedding', [])
```

VectorDB with embedding

- Using the `VectorDB` class

```
!code packages/vdb/load/vdb.py
import sys ; sys.path.append("packages/vdb/load")
import vdb
db = vdb.VectorDB({})
```

- Insert text with embedding

```
db.insert("Hello World")
db.insert("This is a test")
db.insert("This is another test")
db.insert("Testing")
```

Vector Search

- Prepare

```
text = "Test"  
vec = db.embed(text)
```

- Execute the actual search

```
cur = client.search(collection_name=COLLECTION, # collection  
    search_params={"metric_type": "IP"},           # how to measure distance  
    anns_field="embeddings",                      # where to search  
    data=[vec],                                    # what to search  
    output_fields=["text"]                         # field to return  
)
```

Vector Search results

```
for item in cur[0]:  
    dist = item.get('distance', 0)  
    text = item.get("entity", {}).get("text", "")  
    print(dist, text)
```

271.28466796875 Testing

252.88241577148438 This is a test

231.84872436523438 This is another test

181.67494201660156 Hello World

Note output is ordered by **distance**

The screenshot shows the Pinocchio web interface with a purple header bar. The title is "Pinocchio - Nuvolaris". The URL in the address bar is "msciab.openserverless.dev". The main content area has a teal header "Load". On the left, there's a sidebar with sections like "Puzzle", "Hello", "Stream", "Redis", "Milvus", "S3", "Ollama", "RAG", and "Test". Under "RAG", "Load" is selected. The main area shows a conversation between a "BOT" and "YOU". The BOT messages are in grey boxes, and the "YOU" messages are in blue boxes. The messages are:

- BOT: Inserted 456142061011228407 (21:26)
- YOU: This is another test (21:27)
- BOT: Inserted 456142061011228409 (21:27)
- YOU: *test (21:27)
- BOT: Found:
(271.28) Testing
(252.88) This is a test
(231.85) This is another test
(203.04) Hello (21:27)

At the bottom, there's an input field "Enter your message..." and a "Send" button.

vdb/load action

Using command line

```
!ops invoke vdb/load input="openserless cli import"  
!ops invoke vdb/load input="*openserless" | jq .body.output
```

Looking at code

```
!code packages/vdb/load/vdb.py  
!code packages/vdb/load/load.py
```

PDF Import

Extract content from PDF files

Pages from a PDF

```
import pymupdf
doc = pymupdf.open("lessons/bitcoin.pdf") # load a file
text = doc[0].get_text()                  # extract text from page
```

Split in sentences

```
import nltk.data
nltk.download('punkt')                   # load tokenizer model
from nltk.tokenize import sent_tokenize
sentences = sent_tokenize(text)          # extract sentences from text
enum = enumerate(sentences, 1)           # enumerate
sent = next(enum)                      # extract one sentence
```

ops ai loader

```
!ops ai loader
```

```
ai loader [--action=<action>] <file>...
```

```
!ops ai loader lessons/bitcoin.pdf
```

```
>>> converting lessons/bitcoin.pdf
```

```
saved lessons/bitcoin.pdf.txt
```

Excercise: import from the web

- Add to the `vdb/load` the ability to import content from http pages providing an `https://` url

Hints:

- process `input` starting with `https://`
- use `requests` to read the url content
- use `BeautifulSoup (bs4)` to process the content and extract text
- Tokenize text with regular expressions:
<https://stackoverflow.com/questions/75253187/tokenizing-the-text-without-the-use-of-libraries>