

PostgreSQL as a Vector database

Gülçin Yildirim Jelinek Staff Engineer 4th of June, 2024 Boriss Mejías Solutions Architect

The Muffin Chihuahua Challenge













































©EDB 2024 – ALL RIGHTS RESERVED

pgvector

Vectors in your PostgreSQL database (pgvector)

- Specialized vector data type
- Indexing
- Search capabilities
- Distance operators



pgvector is an extension

CREATE EXTENSION vector;



vector is a data type



```
vector is a data type
```

CREATE SCHEMA p2d2; **CREATE TABLE** p2d2.pictures (id BIGSERIAL PRIMARY KEY filename text NOT NULL embedding vector(768) NOT NULL ,);





p2d2.pictures id BIGSERIAL PLAIN filename text EXTENDED embedding vector(768) EXTERNAL





ALTER TABLE p2d2.pictures ALTER COLUMN embedding SET STORAGE PLAIN;





p2d2.pictures id BIGSERIAL PLAIN filename text EXTENDED embedding vector(768) PLAIN



Lab – Part I



Each Participant



- Ip address: https://42.0.0.101 Self-Signed Cert
- Credentials: user42 secrettowel
- Database: p2d2db42

psql -d p2d2db42 CREATE EXTENSION vector; \dx



CREATE TABLE

- psql -d p2d2db42
- **CREATE SCHEMA** p2d2;

CREATE TABLE p2d2.pictures (

- id BIGSERIAL PRIMARY KEY
- , filename text NOT NULL
- , embedding vector(768) NOT NULL

);

ALTER TABLE p2d2.pictures ALTER COLUMN embedding SET STORAGE PLAIN;



Lab – Part II












https://www.kaggle.com/datasets/samuelcortinhas/muffin-vs-chihuahua-image-classification/

python3 with imgbeddings



https://www.kaggle.com/datasets/samuelcortinhas/muffin-vs-chihuahua-image-classification/

Check /home/user42/.load_embeddings

load_embeddings.py



Check /home/user42/.load_embeddings

load_embeddings.py

ibed = imgbeddings()
ibed.to_embeddings(img)



Check /home/user42/.load_embeddings

load_embeddings.py

ibed = imgbeddings()
ibed.to_embeddings(img)

psycopg2.connect(...)
conn.

Check /home/user42/.load_embeddings

load_embeddings.py

ibed = imgbeddings()
ibed.to_embeddings(img)

psycopg2.connect(...)
conn.commit()



\\
./load_embeddings.py
/home/dataset/chihuahua_muffin/chihuahua \
/var/run/postgresql 5432 p2d2db42 user42

./load_embeddings.py
 /home/dataset/chihuahua_muffin/muffin \
 /var/run/postgresql 5432 p2d2db42 user42



©EDB 2024 - ALL RIGHTS RESERVED.

Lab – Part III



©EDB 2024 – ALL RIGHTS RESERVED.



querying



















Querying



©EDB 2024 - ALL RIGHTS RESERVED.



SELECT embedding <-> '[0.2, 0.1, 0.0]' AS distance FROM p2d2.pictures;





- SELECT embedding <-> '[0.2, 0.1, 0.0]' AS distance
 FROM p2d2.pictures;
- SELECT embedding <+> '[0.2, 0.1, 0.0]' AS l1_distance
 FROM p2d2.pictures;





- SELECT embedding <-> '[0.2, 0.1, 0.0]' AS distance
 FROM p2d2.pictures;
- SELECT embedding <+> '[0.2, 0.1, 0.0]' AS l1_distance
 FROM p2d2.pictures;
- SELECT embedding <=> '[0.2, 0.1, 0.0]' AS cosine_simil
 FROM p2d2.pictures;



- SELECT embedding <-> '[0.2, 0.1, 0.0]' AS distance
 FROM p2d2.pictures;
- SELECT embedding <+> '[0.2, 0.1, 0.0]' AS l1_distance
 FROM p2d2.pictures;
- SELECT embedding <=> '[0.2, 0.1, 0.0]' AS cosine_simil
 FROM p2d2.pictures;
- SELECT embedding <#> '[0.2, 0.1, 0.0]' AS inner_prod
 FROM p2d2.pictures;



SELECT id, filename FROM p2d2.pictures
WHERE filename ~ '/chihuahua/' LIMIT 10;
SELECT id, filename FROM p2d2.pictures
WHERE filename ~ '/muffin/' LIMIT 10;

SELECT embedding FROM p2d2.pictures WHERE id=1) FROM p2d2.pictures WHERE id = 2600;



Lab – Part IV



©EDB 2024 – ALL RIGHTS RESERVED.

indexing



©EDB 2024 – ALL RIGHTS RESERVED.



©EDB 2024 - ALL RIGHTS RESERVED.

• HNSW

• IVFFlat



• HNSW

Hierarchical Navigable Small Worlds

• IVFFlat



• HNSW

Hierarchical Navigable Small Worlds

• IVFFlat

InVerted File with Flat compression



• HNSW

Hierarchical Navigable Small Worlds

• IVFFlat

InVerted File with Flat compression

Inverted lists



Creating indexes for vectors



Creating HNSW indexes for vectors

CREATE INDEX ON p2d2.pictures USING hnsw (embedding vector_l2_ops);



Creating HNSW indexes for vectors

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l2_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l1_ops);



Creating HNSW indexes for vectors

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l2_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l1_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_cosine_ops);


Creating HNSW indexes for vectors

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l2_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_l1_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_cosine_ops);

CREATE INDEX ON p2d2.pictures
USING hnsw (embedding vector_ip_ops);



Creating IVFFlat indexes for vectors

CREATE INDEX ON p2d2.pictures
USING ivfflat (embedding vector_l2_ops);

CREATE INDEX ON p2d2.pictures
USING ivfflat (embedding vector_l1_ops);

CREATE INDEX ON p2d2.pictures
USING ivfflat (embedding vector_cosine_ops);

CREATE INDEX ON p2d2.pictures
USING ivfflat (embedding vector_ip_ops);



Querying



©EDB 2024 - ALL RIGHTS RESERVED.

- SELECT embedding <-> '[0.2, 0.1, 0.0]' AS distance
 FROM p2d2.pictures;
- SELECT embedding <+> '[0.2, 0.1, 0.0]' AS l1_distance
 FROM p2d2.pictures;
- SELECT embedding <=> '[0.2, 0.1, 0.0]' AS cosine_simil
 FROM p2d2.pictures;
- SELECT embedding <#> '[0.2, 0.1, 0.0]' AS inner_prod
 FROM p2d2.pictures;



Lab – Part V



©EDB 2024 – ALL RIGHTS RESERVED.

The Challenge



©EDB 2024 – ALL RIGHTS RESERVED.













Check /home/user42/.the_challenge

chihuahua_or_muffin.py



The Challenge Check /home/user42/.the_challenge chihuahua_or_muffin.py

connect

close

cursor

fetchall

Vector operator



The Challenge

./chihuahua_or_muffin.py \

/home/dataset/chihuahua_muffin/challenge/img1a.jpg \
/var/run/postgresql 5432 p2d2db42 user42



The Challenge

./chihuahua_or_muffin.py \

/home/dataset/chihuahua_muffin/challenge/img1a.jpg \
/var/run/postgresql 5432 p2d2db42 user42

python3 -m http.server 8042



Closing Words



©EDB 2024 — ALL RIGHTS RESERVED.

Why not PostgreSQL?



©EDB 2024 – ALL RIGHTS RESERVED.

Why not PostgreSQL?

- Relational database designed for OLTP/OLAP
- Multi-Version Concurrency Control



Why not PostgreSQL?

- Relational database designed for OLTP/OLAP
- Multi-Version Concurrency Control
- Vertical Scalability





©EDB 2024 – ALL RIGHTS RESERVED.

- Shortcomings of specialized vector databases
 - Lack of transactional support
 - Difficult to combine queries with other filters



- Shortcomings of specialized vector databases
 - Lack of transactional support
 - Difficult to combine queries with other filters
- Integration with existing software



- Shortcomings of specialized vector databases
 - Lack of transactional support
 - Difficult to combine queries with other filters
- Integration with existing software
- Integration with existing data



- Shortcomings of specialized vector databases
 - Lack of transactional support
 - Difficult to combine queries with other filters
- Integration with existing software
- Integration with existing data
- Established community



- Shortcomings of specialized vector databases
 - Lack of transactional support
 - Difficult to combine queries with other filters
- Integration with existing software
- Integration with existing data
- Established community to be improved



What's the buzz?



©EDB 2024 — ALL RIGHTS RESERVED.



• Multi-disciplinary work





- Multi-disciplinary work
- pgvector already turns PostgreSQL into a vector database
- pgvector keeps improving





- Multi-disciplinary work
- pgvector already turns PostgreSQL into a vector database
- pgvector keeps improving
- Look at the larger picture and collaborate







©EDB 2024 — ALL RIGHTS RESERVED.