The Wonderful World of WAL

BRUCE MOMJIAN



The write-ahead log, or WAL, provides key features of Postgres, which are covered by this presentation.

https://momjian.us/presentations





Last updated: November 2025

Outline

- 1. Inside the write-ahead log (WAL)
- 2. Crash recovery
- 3. Point-in-time recovery
- 4. Streaming binary replication
- 5. Logical replication
- 6. Replication slots

1. Inside the Write-Ahead Log (WAL): History

- 2001 (PG 7.1): Write-ahead log
- 2005 (PG 8.0): Point-in-time recovery
- 2010 (PG 9.0): Streaming binary replication
- 2017 (PG 10): Logical replication and replication slots

Install pg_walinspect

CREATE EXTENSION IF NOT EXISTS pg_walinspect;

Create Function to Remap Resource_Manager

```
-- adjust resource_manager for historical names

CREATE FUNCTION remap_resource_manager(resource_manager TEXT)

RETURNS TEXT AS

$

SELECT CASE resource_manager

WHEN 'CLOG' THEN 'Pg_xact'

WHEN 'XLOG' THEN 'Pg_wal'

ELSE resource_manager

END

$$ LANGUAGE SQL;
```

Pg_walinspect with resource_manager

```
-- adjust resource manager for historical names
WITH wal records AS (
       SELECT remap resource manager(resource manager) AS res manager,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info('0/01000000', 'FFFFFFFFFFFFFF')
        GROUP BY 1
SELECT
       res manager,
        count.
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager;
```

Output of Pg_walinspect with resource_manager

| res_manager | count | % | size% |
|-------------|-------|------------|-------|
| Btree | 15118 | 57.7 | 12.9 |
| Database | 3 | 0.0 | 0.0 |
| Неар | 5638 | 21.5 | 15.9 |
| Heap2 | 2686 | 10.3 | 17.1 |
| Pg_wal | 1263 | 4.8 | 52.0 |
| Pg_xact | 1 | 0.0 | 0.0 |
| Re1Map | 1 | 0.0 | 0.0 |
| Standby | 414 | 1.6 | 0.2 |
| Storage | 310 | 1.2 | 0.1 |
| Transaction | 756 | 2.9 | 1.7 |

Pg_walinspect with added record_type

```
-- add record type
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type.
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info('0/01000000', 'FFFFFFFFFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count.
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

Pg_walinspect: Btree

| res_manager | record_type | count | % | size% |
|-------------|--------------|-------|------|-------|
| Dtuca | DEDUC | | | |
| Btree | DEDUP | 50 | 0.2 | 0.0 |
| Btree | DELETE | 9 | 0.0 | 0.0 |
| Btree | INSERT_LEAF | 14888 | 56.8 | 11.9 |
| Btree | INSERT_UPPER | 43 | 0.2 | 0.0 |
| Btree | NEWROOT | 30 | 0.1 | 0.0 |
| Btree | SPLIT_L | 15 | 0.1 | 0.4 |
| Btree | SPLIT_R | 38 | 0.1 | 0.5 |
| Btree | VACUUM | 45 | 0.2 | 0.0 |

Pg_walinspect: Database

| res_manager | record_type | ٠. | count | • | ' | ٠ | |
|-------------|---------------------------------|-----------|-------|---|-----|---|--|
| Database | CREATE_FILE_COPY CREATE_WAL_LOG | | 2 | | 0.0 | | |

Pg_walinspect: Heap

| res_manager | record_type | count | % - | size% |
|--------------|---------------------|-------------|---------------|--------------|
| Heap Heap | DELETE HOT UPDATE | 37 393 | 0.1 | 0.0 0.4 |
| Неар | INPLACE | 211 | 0.8 | 0.8 |
| Неар | INSERT | 4326 | 16.5 | 9.1 |
| Неар | INSERT+INIT | 156 | 0.6 | 1.5 |
| Неар | LOCK | 258 | 1.0 | 2.8 |
| Неар | UPDATE | 239 | 0.9 | 1.1 |
| Неар | UPDATE+INIT | 18 | 0.1 | 0.2 |

Pg_walinspect: Heap2

| res_manager | record_type | count | % + | | size% |
|-------------|----------------------|-------|--------|----|-------|
| Heap2 | MULTI_INSERT | 1665 | 6 | .4 | 2.6 |
| Heap2 | MULTI_INSERT+INIT | 52 | 0 | .2 | 0.6 |
| Heap2 | PRUNE_ON_ACCESS | 269 | 1 | .0 | 0.1 |
| Heap2 | PRUNE_VACUUM_CLEANUP | 49 | 0 | .2 | 0.0 |
| Heap2 | PRUNE_VACUUM_SCAN | 256 | 1 | .0 | 0.3 |
| Heap2 | VISIBLE | 395 | 1 | .5 | 13.4 |

[&]quot;Heap2" is the same as "Heap" and is used to support additional heap record types.

Pg_walinspect: Pg_wal

| res_manager | record_type | count | % size% |
|-------------|---------------------------------------------------------------------------------------------------------------|-----------------------------------------------------|-----------------------------------------------------------------------------|
| <u> </u> | CHECKPOINT_ONLINE CHECKPOINT_REDO CHECKPOINT_SHUTDOWN FPI FPI_FOR_HINT NEXTOID | ++- 4 4 3 1106 142 | 0.0 0.0 0.0 0.0 0.0 0.0 4.2 42.5 0.5 9.5 0.0 0.0 |

Pg_walinspect: Pg_xact

| _ | record_type | count | • | ' |
|---------|-----------------|-------|-----|---|
| Pg_xact | | 1 | 0.0 | |

Pg_walinspect: RelMap

Pg_walinspect: Standby

| res_manager | record_type + | count + | • | |
|--------------------|--------------------------------------------|--------------|---|---------|
| Standby Standby | INVALIDATIONS LOCK RUNNING_XACTS | 139 271 | | 0.1 |

Pg_walinspect: Storage

```
res_manager | record_type | count | % | size% | storage | CREATE | 310 | 1.2 | 0.1
```

Pg_walinspect: Transaction

CREATE TABLE

```
-- lsn (log sequence numbers) values represent locations in the WAL
SELECT pg_current_wal_lsn() AS start_lsn
\gset
```

```
-- PRIMARY KEY creates an index
CREATE TABLE wal_test (x INTEGER PRIMARY KEY);
```

WAL Records for CREATE TABLE

```
-- uses psql variable start lsn
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type.
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count.
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for CREATE TABLE

| res_manager reco | ord_type count | % | size% + |
|--------------------------------------------------|------------------------------------|----------------------------------------------------------------------|------------|
| Heap INPLAC Heap INSERT Heap2 MULTI_ | T 7 _INSERT 9 _ON_ACCESS 5 _ 1 _ 2 | 63.3 2.5 8.9 11.4 6.3 1.3 2.5 2.5 | + |

INSERT

```
SELECT pg_current_wal_lsn() AS start_lsn
\gset
INSERT INTO wal_test VALUES (1);
```

WAL Records for INSERT

```
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFF/FFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count,
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for INSERT

| Btree INSERT_LEAF 1 25.0 25 | _ | er record_type | • | • | • |
|-------------------------------------|------------------------|-------------------------------------------|-------|--------------------------|--------------------------|
| ! ! ! ! | Btree Btree Heap | INSERT_LEAF NEWROOT INSERT+INIT | 1 1 1 | 25.0 25.0 25.0 | 25.9 36.4 23.9 |

UPDATE

```
SELECT pg_current_wal_lsn() AS start_lsn
\gset
UPDATE wal_test SET x = 2;
```

WAL Records for UPDATE

```
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFF/FFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count,
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for UPDATE

| _ | record_type | • | • | • |
|-------|-------------------------|-------|--------------------------|--------------|
| Btree | INSERT_LEAF UPDATE | 1 1 | 33.3 33.3 33.3 | 38.3 41.3 |

DELETE

```
SELECT pg_current_wal_lsn() AS start_lsn
\gset

DELETE FROM wal_test;
```

WAL Records for DELETE

```
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFF/FFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count,
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for DELETE

| _ | record_type + | | • | | • | |
|---------------------|------------------|---|---|------|---|--------------|
| Heap Transaction | DELETE | 1 | Ì | 50.0 | Ì | 61.4 38.6 |

DROP TABLE

```
SELECT pg_current_wal_lsn() AS start_lsn
\gset
DROP TABLE wal_test;
```

WAL Records for DROP TABLE

```
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFF/FFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count,
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for DROP TABLE

| res_manager | record_type + | count + | ' | ' |
|---------------|-----------------------------|------------------|---------------------------------|--------------------|
| Heap Heap2 | DELETE PRUNE_ON_ACCESS LOCK | 21 1 2 | 84.0 4.0 8.0 4.0 | 57.6 2.9 4.3 |

INSERT a Second Row

```
CREATE TABLE wal_test (x INTEGER PRIMARY KEY);
INSERT INTO wal_test VALUES (1);
SELECT pg_current_wal_lsn() AS start_lsn
\gset
INSERT INTO wal_test VALUES (2);
```

WAL Records for Second INSERT

```
WITH wal records AS (
        SELECT remap resource manager(resource manager) AS res manager,
                record type,
                COUNT(*) AS count,
                SUM(record length) as size
        FROM pg get wal records info(:'start lsn', 'FFFFFFFF/FFFFFFF')
        GROUP BY 1, 2
SELECT res manager,
        record type,
        count,
        round(count * 100.0 / (SUM(count) OVER ()), 1) AS "%",
        round(size * 100.0 / (SUM(size) OVER ()), 1) AS "size%"
FROM wal records
ORDER BY res manager, record type;
```

WAL Records for Second INSERT

| _ | record_type + | | | |
|-------|-------------------------|-----|--------------------------|----------------|
| Btree | INSERT_LEAF INSERT | 1 1 | 33.3 33.3 33.3 | 40.8 37.6 |

Detailed WAL Records for Second INSERT

```
\x
SELECT *
FROM pg_get_wal_records_info(:'start_lsn', 'FFFFFFFFFFFFFF')
ORDER BY start_lsn;
```

WAL Record #1

```
-[ RECORD 1 ]---+
start 1sn
                   0/01BB8540
end 1sn
                   0/01BB8580
prev 1sn
                   0/01BB8518
xid
                   766
resource manager
                   Heap
record type
                   INSERT
record length
                   59
main data length
fpi length
description
                   off: 2, flags: 0x00
block ref
                   blkref #0: rel 1663/16384/16399 fork main blk 0
```

WAL Record #2

```
-[ RECORD 2 ]---+
start 1sn
                   0/01BB8580
end lsn
                   0/01BB85C0
prev lsn
                   0/01BB8540
xid
                   766
resource manager
                   Btree
record type
                   INSERT LEAF
record length
                   64
main data length
fpi_length
description
                   off: 2
block ref
                   blkref #0: rel 1663/16384/16403 fork main blk 1
```

WAL Record #3

```
-[ RECORD 3 ]----+
start 1sn
                   0/01BB85C0
end 1sn
                   0/01BB85E8
prev 1sn
                   0/01BB8580
xid
                   766
                   Transaction
resource manager
record type
                   COMMIT
record length
                   34
main_data length
fpi_length
description
                   2025-10-31 19:00:00.978944-04
block ref
```

More Detailed WAL Record #1

```
SELECT *
FROM pg_get_wal_block_info(:'start_lsn', 'FFFFFFFFFFFFF')
ORDER BY start_lsn
LIMIT 1;
```

More Detailed WAL Record #1

```
-[ RECORD 1 ]--
start 1sn
                    0/01BB8540
end 1sn
                    0/01BB8580
prev 1sn
                    0/01BB8518
block id
reltablespace
                     1663
reldatabase
                     16384
relfilenode
                     16399
relforknumber
relblocknumber
xid
                    766
resource manager
                    Heap
record type
                     INSERT
record length
                     59
main data length
block data length
block fpi length
block fpi info
                     (null)
description
                    off: 2, flags: 0x00
block data
                     \x01000008180002000000
block fpi data
```

Logical Replication Adds WAL Record #4

```
ALTER SYSTEM SET wal level = 'logical';
$ pg ctl restart
SELECT pg current wal lsn() AS start lsn
\qset
INSERT INTO wal test values (3);
\x
SELECT *
FROM pg_get_wal_records_info(:'start lsn', 'FFFFFFFFFFFFFF')
ORDER BY start 1sn;
```

Logical Replication Adds WAL Record #4

```
-[ RECORD 4 ]----+
start 1sn
                   0/01BC09D0
end 1sn
                   0/01BC0A08
prev 1sn
                   0/01BC09A0
xid
resource manager
                   Standby
record type
                   RUNNING XACTS
record length
                   50
main data length
                   24
fpi length
description
                   nextXid 770 latestCompletedXid 769 oldestRunningXid 770
block ref
```

Generative AI Use Cases

- Text
 - chatbots
 - semantic/vector search
 - summarization
 - language translation
- Image/audio/video creation and search
- Software
 - code creation
 - code analysis
 - data analysis
 - programming language conversion *
 - natural language interface for tasks previously requiring specialized software

```
https://www.turing.com/resources/generative-ai-applications
https://www.youtube.com/watch?v=zjkBMFhNj_g
* https://www.codeconvert.ai/
```

Discriminative and Generative AI Compared

| Discriminative | Generative |
|---------------------------------------|--------------------------|
| supervised learning requires training | outcomes determined from |
| data with known outcomes | massive training data |
| often problem-domain specific | usually general |
| output is a determination/prediction | output is new content |

Discriminative unsupervised learning is similar to generative because it does not include known outcomes.

Discriminative and Generative AI Compared

A generative algorithm models how the data was "generated", so you ask it "what's the likelihood this or that class generated this instance?" and pick the one with the better probability. A discriminative algorithm uses the data to create a decision boundary, so you ask it "what side of the decision boundary is this instance on?" So it doesn't create a model of how the data was generated, it makes a model of what it thinks the boundary between classes looks like. – Anthony, Pittsburgh

Discriminative Boundary

boundary.eps

They focus on boundaries between groups.

Generative Vectors

vector.eps

The color of the arrows is meant to represent an aspect of the word, or an aspect of how the word appeared in the training data. It focuses on relationships between items.

2. Hyper-Dimensional Vectors: Math

Mathematics allows the representation of concepts that can't be represented in the real world:

- Numbers greater than the number of atoms in the universe (10^{80})
- Infinity (∞)
- Function derivatives calculate slopes using infinitesimally small adjustments
- Function integration involves the infinite sum of infinitesimally small areas

Duality

Mathematics allows concepts to be represented in multiple forms. For example, vectors can be thought of as representing the

- Physical world
- Computer science arrays
- Abstract representations

Hyper-Dimensional Vectors

hyper-vector.eps

Mathematics allows vectors in dimensions beyond the three-dimensional world to be used for AI. This is a two-dimensional image represents 10 dimensions. ChatGPT3's text embedding vectors can use 1024 to 12288 dimensions, meaning the image above would require 100-1,200 times more nesting.

Hyper-Dimensional Vector Range

hyper-vector.eps

Each axis is typically a 4-byte floating point number. With 1k-12k dimensions, that yields a total range of 10^{9k} to 10^{118k} values, far more than the number of atoms in the universe. (Even if we use only 4-byte floats between ± 1 , and adjusting the the limited use of bits by floating point numbers, the minimum range is still 10^{8k} .)

Hyper-Dimensional Vector Magnitude

hyper-vector.eps

Notice that each vector has the same length or magnitude.

3. Training Text Embeddings

To train generative AI

- 1. Create vectors
 - one for each English word/token (30-50k)
 - words can be added automatically
 - tokenizers can also be made up by letter combinations (byte-pair based, ideal for rare words and proper nouns)
 - assign each vector the same fixed length/magnitude (normalized vectors)
 - assign each vector a random direction in the 1k-12k dimensional space
- 2. Adjust the direction of vectors using a massive number of training documents (e.g., word2vec) by either
 - for each word, adjust its vector to be closer to the vectors of surrounding words, e.g., bag of words
 - for each word, adjust the vectors of surrounding words to be closer its vector, e.g., skip-gram

https://neptune.ai/blog/vectorization-techniques-in-nlp-guide

https://jalammar.github.io/illustrated-word2vec/

Text Embeddings with Random Directions

embedding.eps

Training Example

For the training text:

The king is a tall man.

we move each red word closer to the blue words, or the blue words closer to the red word (depending on the training method):

- 1. The king is a tall man.
- 2. The king is a tall man.
- 3. The king is a tall man.
- 4. The king is a tall man.
- 5. The king is a tall man.
- 6. The king is a tall man.

Training Example

For the training text:

The king is a tall man. The queen is a beautiful woman. They sit together in the throne room of the castle.

we have

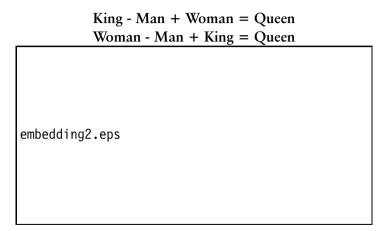
- "king" getting closer to "man"
- "tall" getting closer to "man" and "king"
- "queen" getting closer to "woman"
- "beautiful" getting closer to "woman" and "queen"
- "throne" getting closer to "castle"
- "king", "man" "queen", "woman" getting closer to "throne" and "castle" (spans sentences)

Training Example

- phrases and minor words also move closer to each other, e.g.,
 - "The" is closer to "king"
 - "throne" is closer to "room"
 - "throne" and "room" are closer to "castle"

Consider there are 1-12 thousand dimensions, so moving words closer in one dimension might not affect closeness in other dimensions.

Text Embeddings with Learned Directions



Text Embeddings

chatgpt_embedding.jpg

4. Semantic/Vector Search

Postgres has supported full text/phrase search since 2003, but that only finds base/stemmed words. It has no concept of synonyms (except those explicitly configured) or word relationships. As you can imagine from our previous slides, semantic/vector search promises much richer search capabilities. Full text/phrase search is better for precise queries and handles proper nouns better. It is possible to use a hybrid search which combines full text/phrase search and semantic/vectors search.

Semantic/Vector Search Setup

- 1. Download a pre-trained text embedding, or create your own
- 2. Choose a chunk size (sentence, paragraph, section, page, document)
- 3. Find text embedding vectors of all words (or byte pairs) in the chunk
- 4. Average the vectors
- 5. Store them in the database along with the chunk

```
https://qdrant.tech/articles/what-is-rag-in-ai/
```

Semantic/Vector Search

- 1. Find text embedding vectors of all words in the query
- 2. Average the vectors
- 3. Search for the closest (nearest neighbor) vectors in the database and return their chunks

Search Setup Example

Assume you wish to index this document using semantic/vector search:

The king is a tall man. The queen is a beautiful woman. They sit together in the throne room of the castle.

If you choose a chunk size of sentences, you will store and index:

- 1. The king is a tall man. (average 6 vectors)
- 2. The queen is a beautiful woman. (average 6 vectors)
- 3. They sit together in the throne room of the castle. (average 10 vectors)

Search Example

A search of:

Who is the king?

will find the text embedding vectors of the query words, average them, and find the closest (nearest neighbor) stored vector which is this chunk:

The king is a tall man.

Example: Install pgvector

CREATE EXTENSION IF NOT EXISTS vector;

Pgvector can be found at https://github.com/pgvector/pgvector. All queries in this presentation can be downloaded from https://momjian.us/main/writings/pgsql/trenches.sql.

Create Table to Store Documents and Vectors

```
CREATE TABLE document (
   id SERIAL PRIMARY KEY,
   content TEXT NOT NULL
);

CREATE TABLE document_embedding (
   id INTEGER PRIMARY KEY,
   -- ChatGPT4's text embedding vectors have 1536 dimensions
   embedding vector(1536) NOT NULL
);
```

Example based on the URL at the bottom right.

Populate Document Table

Populate Text Embedding

```
#! /usr/bin/env python
""" Add document embeddings to the vector search table """
import os
import sys
from openai import OpenAI
import psycopg
if len(sys.argv) != 1:
    print("Usage: " + os.path.basename( file ), file=sys.stderr)
    sys.exit(1)
# Get OpenAI API key
client = OpenAI(api key=os.getenv("OPENAI API KEY"))
# Specify the embedding model
MODEL ID = "text-embedding-ada-002"
```

Populate Text Embedding

```
# Connect to the database
conn = psycopq.connect("host=localhost port=5432 dbname=ai user=postgres")
# Fetch documents
cur = conn.cursor()
cur.execute("SELECT id, content FROM document")
# Create embeddings for each document and store in the database
for doc id, doc content in cur.fetchall():
    embedding = (
        client.embeddings.create(input=doc content, model=MODEL ID).data[0].embedding
    cur.execute(
INSERT INTO document embedding (id. embedding)
VALUES (%s, %s);""",
        (doc id, embedding).
# Commit and close the database connection
conn.commit()
conn.close()
```

Embeddings Stored

```
-- show embeddings for documents
SELECT content.
       substring(embedding::text, 1, 30) AS embedding
FROM document JOIN document embedding USING (id);
                       content
                                                                   embedding
The king is a tall man.
                                                        [0.004228712,-0.0132633485,0.0...
                                                        [-0.006298352.0.008259722.-0.0...
 The gueen is a beautiful woman.
 They sit together in the throne room of the castle. \lceil 0.01191413, -0.03098976, -0.010...
-- each dimension is four bytes (float4)
SELECT pg column size(embedding),
       pg column size(embedding) / 1536 AS bytes per dim
FROM document embedding;
pg column size | bytes per dim
           6148 I
           6148 I
           6148 I
```

Query Text Embedding

```
#! /usr/bin/env python
""" Perform vector search of documents """
import os
import sys
from openai import OpenAI
import psycopg
if len(sys.argv) != 2:
    print("Usage: " + os.path.basename( file ) + " search string", file=sys.stderr)
    sys.exit(1)
# Get the user query
search = sys.argv[1]
# Get OpenAI API kev
client = OpenAI(api key=os.getenv("OPENAI API KEY"))
# Specify the embedding model
MODEL ID = "text-embedding-ada-002"
```

Query Text Embedding

```
# Get embedding for the user query
embedding = client.embeddings.create(input=search, model=MODEL ID).data[0].embedding
# Connect to the database
conn = psycopg.connect("host=localhost port=5432 dbname=ai user=postgres")
# Fetch documents in order of their vector search similarity to the user query
cur = conn.cursor()
cur.execute(
    .....
SELECT content, embedding <-> %s::vector
FROM document JOIN document embedding USING (id)
ORDER BY 2:""".
    (embedding.).
# Output documents ordered by similarity
for doc content, distance in cur.fetchall():
    print(f"{doc content:<52} {distance}")</pre>
# Commit and close the database connection
conn.commit()
conn.close()
```

Query Examples

```
-- This is matching a word, so it could have been done with full text search.

Who is the king?
The king is a tall man.

The queen is a beautiful woman.

They sit together in the throne room of the castle.

-- same

Who is tall?
The king is a tall man.

They sit together in the throne room of the castle.

0.6273594902980766
They sit together in the throne room of the castle.

0.6273594902980766
They sit together in the throne room of the castle.

0.7227548068086962
```

Query Examples

```
-- The vector for "short" is near "tall".
Who is short?
The king is a tall man.
                                                    0.6878928016051950
The gueen is a beautiful woman.
                                              0.7552441994053801
They sit together in the throne room of the castle. 0.7593614930258540
-- The vector for "pretty" is near "beautiful".
Who is pretty?
The queen is a beautiful woman.
                                               0.6579895352151753
They sit together in the throne room of the castle. 0.7473554877225373
The king is a tall man.
                                                    0.7480856783292117
-- The vector for "palace" is near "castle".
Who is in the palace?
They sit together in the throne room of the castle. 0.5872522481044645
The king is a tall man.
                                                    0.6646927703591086
The queen is a beautiful woman.
                                                    0.6908769898680150
-- The vector for "chair" is near "sit".
Where is the chair?
They sit together in the throne room of the castle. 0.6865519576029434
The king is a tall man.
                                                    0.7417508497217221
The gueen is a beautiful woman.
                                                    0.7578119887517404
```

5. Transformers

Transformers allow the embedding vectors just shown and attention blocks to generate output. While images and videos can also be created, this presentation will focus on text generation. How is this done?

- 1. Load the first attention block with the text embedding vectors of the words used in the user query
- 2. Adjust vectors to be closer to previous words/vectors
- 3. Repeat step 2 several times
- 4. Find the word closest to the last vector, and use it as the first output word
- 5. Repeat step 2 and later to get successive words

Attention Block Details

- Shown are self-attention blocks:
 - only previous vectors affect current vectors, not later ones
 - used to generate new words
 - translation uses cross-attention blocks where later vectors can also affect earlier ones
- While the following slides show attention block 1 in seven steps, attention blocks do vector calculations in parallel as matrix multiplication
 - GPUs are very efficient at matrix multiplication
 - internally uses a query/key/value process

https://sebastianraschka.com/blog/2023/self-attention-from-scratch.html

Attention Block Details

- While text embedding vectors can be thousands of dimensions, attention blocks usually use one-tenth the number of dimensions:
 - GPT-3 uses 96 attention blocks with 128-dimension vectors
- Shown is zero-shot learning
 - some models use single or multi-shot learning, e.g., "The capital of Spain is Madrid."
 is prepended to the user query
- Not shown
 - word positions are numbered in the attention vectors
 - sentence endings are also encoded; this is shown as "?"
 - used for sentence construction

Populate Initial Attention Block

transformer1_1.eps

Current Attention Block Vectors

- 1. "What"
- 2. "is"
- 3. "the"
- 4. "capital"
- 5. "of"
- 6. "France"
- 7. "?"

Adjust the Second Vector

transformer1_2.eps

Current Attention Block Vectors

- 1. "What"
- 2. "What" "is"
- 3. "the"
- 4. "capital"
- 5. "of"
- 6. "France"
- 7. "?"

Matrix Mask Prevents Adjustments from Later Words/Vectors

| Words | What | is | the | capital | of | France | 3 |
|---------|------|----|-----|---------|----|--------|---|
| What | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| is | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| the | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| capital | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| of | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| France | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| 3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

transformer1_3.eps

Current Attention Block Vectors

- 1. "What"
- 2. "What" "is"
- 3. "What" "is" "the"
- 4. "capital"
- 5. "of"
- 6. "France"
- 7. "?"

transformer1_4.eps

Current Attention Block Vectors

- 1. "What"
- 2. "What is"
- 3. "What" "is" "the"
- 4. "What" "is" "the" "capital"
- 5. "of"
- 6. "France"
- 7. "?"

transformer1_5.eps

transformer1_6.eps

transformer1_7.eps

End of Attention Block 1 Vectors

- 1. "What"
- 2. "What" "is"
- 3. "What" "is" "the"
- 4. "What" "is" "the" "capital"
- 5. "What" "is" "the" "capital" "of"
- 6. "What" "is" "the" "capital" "of" "France"
- 7. "What" "is" "the" "capital" "of" "France" "?"

Attention Block Details

- Each vector of an attention block potentially can be moved closer to other vectors, except the first vector:
 - this is the "attention" aspect of attention blocks
- Vectors are normalized to maintain the same magnitude
- Vector movement is not uniform
 - movement favors dimensionally-close vectors and previous movement (backpropagation)
 - distances are computed via vector dot products

Continue Running Attention Blocks

transformer2+.eps

The next attention block uses vectors that were adjusted by the previous attention block. The headings are now labeled as "vectors" because they no longer point to the original words of the user query.

End of Attention Block 2 Vectors

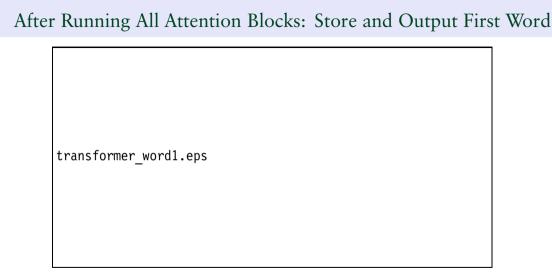
- 1. ("What")
- 2. ("What")×("What" "is")
- 3. ("What")×("What" "is")×("What" "is" "the")
- 4. ("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")
- 5. ("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital" "of")
- 6. ("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of" "France")
- 7. ("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of" "France")×("What" "is" "the" "capital" "of" "France" "?")

End of Attention Block 3 Vectors

- 1. ("What")
- 2. ("What")×(("What")×("What" "is"))
- 3. (("What")×("What" "is"))×(("What")×("What" "is")×("What" "is" "the"))
- 4. (("What")×("What" "is")×("What" "is" "the"))×(("What")×("What" "is")×("What" "is" "the")×("What" "is" "the")×("What" "is" "the")
- 5. (("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital"))×(("What")×("What" "is")×("What" "is" "the" "capital")×("What" "is" "the" "capital" "of"))
- 6. (("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital")×("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of" "France"))
- 7. (("What")×("What" "is")×("What" "is" "the")×("What" "is" "the" "capital")×("What" "is" "the" "capital")×("What" "is")×("What" "is")×("What" "is")×("What" "is") ("What" "is" "the" "capital") ("What" "is" "the" "capital" "of")×("What" "is" "the" "capital" "of" "France")×("What" "is" "the" "capital" "of" "France" "?"))

End of Attention Block 4 Vectors

- 1. ("What")
- 2. ("What")×("What")×("What" "is"))
- 3. ("What")×(("What")×("What" "is"))×(("What")×("What" "is"))×(("What")×("What" "is")×("What" "is")×
- 4. (("What")×("What" "is"))×(("What")×("What" "is")×("What" "is")×("What
- 5. (("What") × ("What" "is") × ("What" "is" "the")) × (("What") × ("What" "is") × ("What" "is" "the") × ("What
- 6. (("What") × ("What" "is") × ("What" "is" "the") × ("What" "is" "the") × ("What" "is" "the") × ("What" "is") × ("What" "is" "the") × ("What" "is" "the")
- 7. (("What") × ("What" "is" "the") × ("What" "is" "the") × ("What" "is" "the" "capital") × ("What" "is" "the" "capital" "of") × (("What") × ("What" "is") × ("What" "is" "the" "capital") × ("What" "is" "the" "capital" of") × ("What" "is" "the" "capital" of") × ("What" "is" "the" "capital" of" "France") × ("What" "is" "the" "capital" of ""france") × ("What" "is" "the" "capital") × ("What" "is" "the" "capital" "of" "France") ×



The last vector of the final attention block has dimensional values which represent the entire user query; its position among other word vectors is used to generate answer text. AI chat apps usually display results one word at a time because they generate results one word at a time.

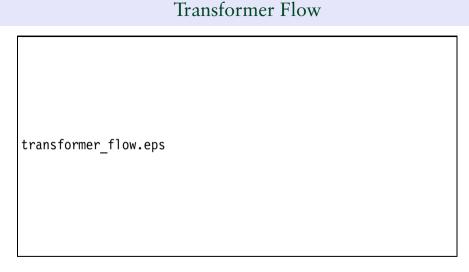
Adjust the First Word

transformer_word1 attention.eps

Vectors only adjust to previous vectors, so only vector8 needs adjustment.

After Running All Attention Blocks Again: Store and Output Second Word

transformer word2.eps



https://sebastianraschka.com/blog/2023/self-attention-from-scratch.html

6. Generative AI

Generated using large language model (LLM) ChatGPT4 o1-mini:

What is the capital of France? The capital of France is Paris.

The capital of France is? The capital of France is Paris.

Where is Paris?

Paris is the capital city of France, located in the north-central part of the country along the Seine River. It is known for its rich history, art, culture, and landmarks such as the Eiffel Tower, the Louvre Museum, and Notre-Dame Cathedral.

Generative AI

Columbus is in Ohio. Where is Paris?

Paris is in France. There is also a city named Paris in the United States, located in Texas. If you were referring to a different Paris, please specify!

I live in Texas. Where is Paris?

In Texas, there is a city named Paris, which is located in the northeastern part of the state, in Lamar County. It is approximately 100 miles northeast of Dallas. There is also a well-known city named Paris in France, which is the capital and largest city of the country. If you're referring to the Texas city, it's known for its local attractions, including the Eiffel Tower replica with a cowboy hat. If you meant the one in France, it's famous for its art, culture, and history.

7. Retrieval-Augmented Generation (RAG)

Retrieval-augmented generation (RAG) allows for data (below in red) to be passed to the generative transformer to affect the user query (transformers can accept thousands of words/tokens).

Where is paris?

Paris is the capital city of France, located in the northern part of the country along the Seine River. It is known for its rich history, culture, art, and landmarks such as the Eiffel Tower, the Louvre Museum, and Notre-Dame Cathedral.

reply briefly

Where is Paris?

Paris is the capital city of France, located in the north-central part of the country along the Seine River.

RAG Internals

```
$ rag -n 'reply briefly' 'Where is Paris?'
Payload:
  "model": "gpt-4o-mini",
  "stream": true,
  "temperature": 0.5.
  "max tokens": 8192,
  "messages": [
      "role": "system"
      "content": "reply briefly",
      "role": "user",
      "content": "Where is Paris?"
```

RAG

Columbus is in Ohio.

Where is Paris?

Paris is the capital city of France, located in the north-central part of the country along the Seine River. It is known for its art, fashion, gastronomy, and culture, as well as landmarks such as the Eiffel Tower, the Louvre Museum, and Notre-Dame Cathedral.

I live in Texas.

Where is Paris?

Paris is the capital city of France, located in the north-central part of the country along the Seine River. It is known for its rich history, culture, art, and landmarks such as the Eiffel Tower, the Louvre Museum, and Notre-Dame Cathedral. If you meant a different Paris, such as Paris, Texas, it is a city located in the northeastern part of the state of Texas.

Data Can Be Passed via RAG

Using the following RAG data for each query:

```
{"order_id": 735, "product": "gears", "amount": 21, "inv_date": 2024-04-12} 
{"order_id": 4423, "product": "pencils", "amount": 13, "inv_date": 2024-02-05} 
{"order_id": 1334, "product": "glue", "amount": 5, "inv_date": 2024-03-25}
```

A more sophisticated example would be to first create a list of questions or keywords and their relevant SQL queries. A semantic/vector search could then be performed on each user query to choose the relevant SQL query, run it, and supply its output as RAG data.

Data Can Be Passed via RAG

When were the orders made?
The orders were made on the following dates:

- 1. Order ID 735: April 12, 2024
- 2. Order ID 4423: February 5, 2024
- 3. Order ID 1334: March 25, 2024

What is the earliest invoice?

The earliest invoice is for the product "pencils" with order_id 4423, dated February 5, 2024.

Data Can Be Passed via RAG

What is the total of all invoices?

To calculate the total of all invoices, we need to sum the amounts for each order:

- 1. Order 735: Amount = 21
- 2. Order 4423: Amount = 13
- 3. Order 1334: Amount = 5

Now, adding them together:

$$Total = 21 + 13 + 5 = 39$$

The total of all invoices is 39.

Data Can Be Passed via RAG

What is the average invoice amount?

To calculate the average invoice amount, we first need to sum the amounts from the provided orders and then divide by the number of orders.

Here are the amounts from the orders:

- Order 735: 21
- Order 4423: 13
- Order 1334: 5

Now, let's calculate the total amount:

Total amount =
$$21 + 13 + 5 = 39$$

Average invoice amount = Total amount / Number of orders =
$$39 / 3 = 13$$
 So, the average invoice amount is 13.

Today is April 15, 2024. With 30-day payment expected, which invoices are overdue? To determine which invoices are overdue as of April 15, 2024, we need to check the invoice dates and see if they are more than 30 days old from the current date.

- 1. **Invoice for gears**:
- Inv Date: April 12, 2024
- Payment Due Date: May 12, 2024 (April 12 + 30 days)
- Status: Not overdue
- 2. **Invoice for pencils**:
- Inv Date: February 5, 2024
- Payment Due Date: March 6, 2024 (February 5 + 30 days)
- Status: Overdue

- 3. **Invoice for glue**:
- Inv Date: March 25, 2024
- Payment Due Date: April 24, 2024 (March 25 + 30 days)
- Status: Not overdue

Based on this analysis, the only overdue invoice is for **pencils** (order_id 4423).

What is the total of invoices that are not pencils?

To calculate the total of invoices that are not for pencils, we will look at the provided orders and sum the amounts for each product except for pencils.

Here are the relevant orders:

- 1. Gears: 21
- 2. Glue: 5

Now, we add these amounts together:

- 21 (gears) + 5 (glue) = 26
- Therefore, the total of invoices that are not for pencils is **26**.

Today is April 15, 2024. What is the next expected order date?

Since today's date is April 15, 2024, and the most recent order date in the provided data is April 12, 2024, we can estimate the next expected order date by looking at the pattern of previous orders. However, there is not enough information in the provided data to determine a specific frequency or pattern of orders. If we assume that orders are placed regularly, we can consider the time intervals between the previous orders:

- 1. Order on February 5, 2024
- 2. Order on March 25, 2024 (48 days after the previous order)
- 3. Order on April 12, 2024 (18 days after the previous order)

Given the decreasing intervals, it's challenging to predict a specific date. However, if we were to assume a similar timeframe, we could expect the next order to occur sometime in the following weeks. A rough estimate might suggest the next order could be expected around late April to early May 2024, but without a clear pattern, it remains uncertain.

SQL Queries Can Be Generated

What SQL query will generate the latest invoice date? To generate the latest invoice date from a table (let's assume the table is named 'invoices'), you can use the following SQL query:

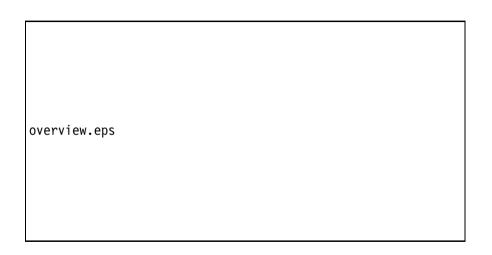
```
```sql
SELECT MAX(inv_date) AS latest_invoice_date
FROM invoices;
```

This query uses the 'MAX()' function to find the most recent date in the 'inv\_date' column, returning it as 'latest\_invoice\_date'. Make sure to replace 'invoices' with the actual name of your table if it's different.

### Data Can Be Taken from a Database

```
CREATE TABLE invoice (
 order id INTEGER PRIMARY KEY,
 product TEXT NOT NULL,
 amount NUMERIC(10,2),
 inv date DATE
);
INSERT INTO invoice VALUES
 (735, 'gears', 21, '2024-04-12'),
 (4423, 'pencils', 13, '2024-02-05'),
 (1334, 'glue', 5, '2024-03-25'):
$ # This invoice list can be customized by customer.
$ rag "$(psgl --tuples-only -c '\
 SELECT to json(invoice.*) \
 FROM invoice \
 WHERE cust id = 12; 'ai)" \
> 'How many invoices are there?'
There are three invoices.
```

# 8. Deployment



### Relational Databases for Semantic/Vector Search & Generative AI

| Section | AI Feature                           | Details                           | DB Appropriateness |
|---------|--------------------------------------|-----------------------------------|--------------------|
| 3       | text embedding vector training       | batched changes, mostly static    | poor               |
| 4       | semantic/vector search               | search existing database contents | good               |
| 5 & 6   | tranformers & generative AI          | billions of comparisons           | poor               |
| 7       | retrieval-augmented generation (RAG) | add details to user queries       | good               |
| 7       | language queries on retrieved data   | supplied as JSON                  | good               |
| 7       | generate SQL queries                 | natural language to SQL           | unknown*           |
|         | data analysis                        | regression and time series        | good               |

Relational databases continue to be appropriate for discriminative AI. AI tools can also help with database migrations, e.g., Oracle to Postgres.

# **Deployment Options**

As you have seen, there are several options for using generative AI:

- Cloud service model, e.g., ChatGPT
- Self-managed model
  - only publicly trained, e.g., Meta's Llama
    - publicly trained with private fine-tuningpre, e.g., domain-specific chat assistant
  - locally trained public and private data sets, e.g., structured-wikipedia
  - locally trained with only private data
  - private training requires an open data model
    - Hugging Face offers many pretrained models and public data sets
    - most OpenAI models are not open\*
- The above options can be augmented with local data, i.e., RAG; data derived from
  - personal preferences
  - relational data and supplied as JSON, e.g., PostgreSQL
  - text retrieved either via full text/phrase search or semantic/vector search

```
https://zapier.com/blog/hugging-face/
https://www.infoworld.com/article/2338922/5-easy-ways-to-run-an-llm-locally.html
https://www.nature.com/articles/d41586-024-02998-y
*https://analyticsindiamag.com/ai-mysteries/6-open-source-models-from-openai/
```

# Postgres AI Solutions

- pgvector, already covered
- Tembo's pg vectorize
- EDB's AI extension *aidb* (Pipelines)
- Timescale's AI extension pgai
- PostgresML's AI toolkit
- AI toolkits from cloud vendors
- Lantern Cloud
- Workik, PopSQL, and SQL AI SQL query generators
- Postgres.AI's PostgreSQL chatbot
- DBtune for server parameter tuning

### 9. Conclusion

- Pre-computer philosophy
- 1950's Turing test
- 1980's Expert systems
- 1970's 1990's AI winter
- 2010's Robotics
- 2013 word2vec by Google
- 2017 Attention blocks by Google
- 2022 ChatGPT for generating text, DALL-E for generating images

# How Did Google Miss the Boat?

- Groundbreaking research by Google in the 2010's to support language translation
- Products focused on revenue-generating activities like web search and advertising
- Worked on AI also to support device control
- Did not focus on aggregating knowledge across web pages like ChatGPT, or did not wish to risk existing revenue streams

# The Future: Revolutionary Vision

A year ago, if you had said to me in our lifetime will we have capabilities like we have today now with ChatGPT4 ... if you explain the kinds of things that ChatGPT4 can do I probably would have said to you a year ago I don't know if we will have those capabilities in our lifetime — and now we have it today — so the speed at which this is moving is staggering. — Jon Krohn, June 12, 2023

### The Future: Incremental Vision

We are used to the idea that people or entities that can express themselves, or manipulate language, are smart — but that's not true. You can manipulate language and not be smart, and that's basically what LLMs (large language models) are demonstrating. — Yann LeCun, October 11, 2024

```
https://www.linkedin.com/posts/yann-lecun_an-article-in-the-wall-street-journal-in-activity-7250915579228827648-WWA2/
https://www.msn.com/en-in/money/news/why-were-unlikely-to-get-artificial-general-intelligence-anytime-soon/ar-AA1EWy4y
https://www.theatlantic.com/culture/archive/2025/06/artificial-intelligence-illiteracy/683021/
```

### Conclusion

24017301626\_14e1fa2565\_c.jpg

