

A Year-in-Review of vectors in PostgreSQL

Jonathan Katz

Principal Product Manager – Technical Amazon RDS

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Generative AI is powered by foundation models

Pretrained on vast amounts of unstructured data

Contain a large number of parameters that make them capable of learning complex concepts

Can be applied in a wide range of contexts

Customize FMs using your data for domainspecific tasks





Retrieval-augmented generation (RAG)

Configure FM to interact with your company data





How vector embeddings are used





Challenges with vectors

- Time to generate embeddings
- Embedding size
- Compression
- Query time



1536-dim 4-byte floats 6152B => 6KiB 1,000,000 => 5.7GB



Approximate nearest neighbor (ANN)

- Find similar vectors without searching all of them
- Methods include hash, tree, clustering, graph
- Faster than exact nearest neighbor

• "Recall" – % of expected results



Recall: 80%



Why use PostgreSQL for vector searches?

- Existing client libraries work without modification
- Convenient to co-locate app + AI/ML data in same database

• Interfacing with PostgreSQL storage gives ACID transactional storage



PostgreSQL support for vectors

Native

- ARRAY
- cube

aws

Extensions

- pgvector
 <u>pg_embedding</u>
- pgvecto.rs
- Lantern
- Timescale Vector
- pgvector-remote

What is pgvector?

An open-source extension that:

Adds support for storage, indexing, searching, metadata with choice of distance

vector data types

Co-locate with embeddings

Exact nearest neighbor (K-NN) Approximate nearest neighbor (ANN)

Supports IVFFlat/HNSW indexing

Distance operators (<->, <=>, <#>)

github.com/pgvector/pgvector

pgvector popularity



Does pgvector scale?

BIGANN 10M (128-dim) on RDS PostgreSQL r7g.12xlarge (pgvector 0.6.2) k=10 - x-axis is # concurrent clients



Querying an IVFFlat index



SET ivfflat.probes TO 2

SELECT id FROM products ORDER BY \$1 <-> embedding LIMIT 3



Querying an HNSW index



pgvector: Year-in-review timeline

- <u>v0.4.x</u> (1/2023 6/2023)
 - Improved IVFFlat plan costs
 - Increasing dimension of vectors stored in table + index
- <u>v0.5.x</u> (8/2023 10/2023)
 - Add HNSW index + distance function performance improvements
 - Parallel IVFFlat builds
- <u>v0.6.x</u> (1/2024 3/2024)
 - Parallel HNSW index builds + in-memory build optimizations
- <u>v0.7.x</u> (Unreleased, planner 4/2024)
 - halfvec (2-byte float), bit(n) index support, sparsevec (up to 1B dim)
 - Quantization (scalar/binary), Jaccard/hamming distance, explicit SIMD

dbpedia-openai-1000k-angular (1MM 1536-dim) - Index Build Time





dbpedia-openai-1000k-angular (1MM 1536-dim) – p99 Latency @ 99% recall

dbpedia-openai-1000k-angular (1MM 1536-dim)



Impact of parallelism on HNSW build time

HNSW index build (1,000,000 128-dim vectors)



Choosing m and ef_construction (serial)



Choosing m and ef_construction (parallel)



pgvector roadmap

- Enhanced index-based filtering ("pre-filtering") (in progress)
- Additional quantization techniques (integer, product)
- Parallel query
- Improvements to PostgreSQL async query pushdown for federated sorts



Thank you!

Jonathan Katz

jkatz@amazon.com @jkatz05



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.