

## pgvector

Open-source vector similarity search for Postgres

Store your vectors with the rest of your data. Supports:

- exact and approximate nearest neighbor search
- L2 distance, inner product, and cosine distance
- any language with a Postgres client

Plus ACID compliance, point-in-time recovery, JOINS, and all of the other great features of Postgres

## Installation

### Linux and Mac

Compile and install the extension (supports Postgres 12+)

```
cd /tmp
git clone --branch v0.6.2 https://github.com/pgvector/pgvector.git
cd pgvector
make
make install # may need sudo
```

See the installation notes if you run into issues

You can also install it with Docker, Homebrew, PGXN, APT, Yum, pkg, or conda-forge, and it comes preinstalled with Postgres.app and many hosted providers. There are also instructions for GitHub Actions.

### Windows

Ensure C++ support in Visual Studio is installed, and run:

```
call "C:\Program Files\Microsoft Visual Studio\2022\Community\VC\Auxiliary\Build\vcvars64.bat"
```

Note: The exact path will vary depending on your Visual Studio version and edition

Then use `nmake` to build:

```
set "PGRROOT=C:\Program Files\PostgreSQL\16"
cd %TEMP%
git clone --branch v0.6.2 https://github.com/pgvector/pgvector.git
cd pgvector
nmake /F Makefile.win
nmake /F Makefile.win install
```

See the installation notes if you run into issues

You can also install it with Docker or conda-forge.

## Getting Started

Enable the extension (do this once in each database where you want to use it)

```
CREATE EXTENSION vector;
```

Create a vector column with 3 dimensions

```
CREATE TABLE items (id bigserial PRIMARY KEY, embedding vector(3));
```

Insert vectors

```
INSERT INTO items (embedding) VALUES ('[1,2,3]'), ('[4,5,6]');
```

Get the nearest neighbors by L2 distance

```
SELECT * FROM items ORDER BY embedding <-> '[3,1,2]' LIMIT 5;
```

Also supports inner product (<#>) and cosine distance (<=>)

Note: <#> returns the negative inner product since Postgres only supports ASC order index scans on operators

## Storing

Create a new table with a vector column

```
CREATE TABLE items (id bigserial PRIMARY KEY, embedding vector(3));
```

Or add a vector column to an existing table

```
ALTER TABLE items ADD COLUMN embedding vector(3);
```

Insert vectors

```
INSERT INTO items (embedding) VALUES ('[1,2,3]'), ('[4,5,6]');
```

Or load vectors in bulk using COPY (example)

```
COPY items (embedding) FROM STDIN WITH (FORMAT BINARY);
```

Upsert vectors

```
INSERT INTO items (id, embedding) VALUES (1, '[1,2,3]'), (2, '[4,5,6]')  
ON CONFLICT (id) DO UPDATE SET embedding = EXCLUDED.embedding;
```

Update vectors

```
UPDATE items SET embedding = '[1,2,3]' WHERE id = 1;
```

Delete vectors

```
DELETE FROM items WHERE id = 1;
```

## Querying

Get the nearest neighbors to a vector

```
SELECT * FROM items ORDER BY embedding <-> '[3,1,2]' LIMIT 5;
```

Get the nearest neighbors to a row

```
SELECT * FROM items WHERE id != 1 ORDER BY embedding <-> (SELECT embedding FROM items WHERE
```

Get rows within a certain distance

```
SELECT * FROM items WHERE embedding <-> '[3,1,2]' < 5;
```

Note: Combine with ORDER BY and LIMIT to use an index

**Distances** Get the distance

```
SELECT embedding <-> '[3,1,2]' AS distance FROM items;
```

For inner product, multiply by -1 (since <#> returns the negative inner product)

```
SELECT (embedding <#> '[3,1,2]') * -1 AS inner_product FROM items;
```

For cosine similarity, use 1 - cosine distance

```
SELECT 1 - (embedding <=> '[3,1,2]') AS cosine_similarity FROM items;
```

**Aggregates** Average vectors

```
SELECT AVG(embedding) FROM items;
```

Average groups of vectors

```
SELECT category_id, AVG(embedding) FROM items GROUP BY category_id;
```

## Indexing

By default, pgvector performs exact nearest neighbor search, which provides perfect recall.

You can add an index to use approximate nearest neighbor search, which trades some recall for speed. Unlike typical indexes, you will see different results for queries after adding an approximate index.

Supported index types are:

- HNSW - added in 0.5.0
- IVFFlat

## HNSW

An HNSW index creates a multilayer graph. It has better query performance than IVFFlat (in terms of speed-recall tradeoff), but has slower build times and uses more memory. Also, an index can be created without any data in the table since there isn't a training step like IVFFlat.

Add an index for each distance function you want to use.

L2 distance

```
CREATE INDEX ON items USING hnsw (embedding vector_l2_ops);
```

Inner product

```
CREATE INDEX ON items USING hnsw (embedding vector_ip_ops);
```

Cosine distance

```
CREATE INDEX ON items USING hnsw (embedding vector_cosine_ops);
```

Vectors with up to 2,000 dimensions can be indexed.

### Index Options

Specify HNSW parameters

- `m` - the max number of connections per layer (16 by default)
- `ef_construction` - the size of the dynamic candidate list for constructing the graph (64 by default)

```
CREATE INDEX ON items USING hnsw (embedding vector_l2_ops) WITH (m = 16, ef_construction = 64);
```

A higher value of `ef_construction` provides better recall at the cost of index build time / insert speed.

### Query Options

Specify the size of the dynamic candidate list for search (40 by default)

```
SET hnsw.ef_search = 100;
```

A higher value provides better recall at the cost of speed.

Use `SET LOCAL` inside a transaction to set it for a single query

```
BEGIN;  
SET LOCAL hnsw.ef_search = 100;  
SELECT ...  
COMMIT;
```

## Index Build Time

Indexes build significantly faster when the graph fits into `maintenance_work_mem`

```
SET maintenance_work_mem = '8GB';
```

A notice is shown when the graph no longer fits

```
NOTICE: hnsw graph no longer fits into maintenance_work_mem after 100000 tuples
DETAIL: Building will take significantly more time.
HINT: Increase maintenance_work_mem to speed up builds.
```

Note: Do not set `maintenance_work_mem` so high that it exhausts the memory on the server

Like other index types, it's faster to create an index after loading your initial data

Starting with 0.6.0, you can also speed up index creation by increasing the number of parallel workers (2 by default)

```
SET max_parallel_maintenance_workers = 7; -- plus leader
```

For a large number of workers, you may also need to increase `max_parallel_workers` (8 by default)

## Indexing Progress

Check indexing progress with Postgres 12+

```
SELECT phase, round(100.0 * blocks_done / nullif(blocks_total, 0), 1) AS "%" FROM pg_stat_p
```

The phases for HNSW are:

1. initializing
2. loading tuples

## IVFFlat

An IVFFlat index divides vectors into lists, and then searches a subset of those lists that are closest to the query vector. It has faster build times and uses less memory than HNSW, but has lower query performance (in terms of speed-recall tradeoff).

Three keys to achieving good recall are:

1. Create the index *after* the table has some data
2. Choose an appropriate number of lists - a good place to start is `rows / 1000` for up to 1M rows and `sqrt(rows)` for over 1M rows
3. When querying, specify an appropriate number of probes (higher is better for recall, lower is better for speed) - a good place to start is `sqrt(lists)`

Add an index for each distance function you want to use.

L2 distance

```
CREATE INDEX ON items USING ivfflat (embedding vector_l2_ops) WITH (lists = 100);
```

Inner product

```
CREATE INDEX ON items USING ivfflat (embedding vector_ip_ops) WITH (lists = 100);
```

Cosine distance

```
CREATE INDEX ON items USING ivfflat (embedding vector_cosine_ops) WITH (lists = 100);
```

Vectors with up to 2,000 dimensions can be indexed.

### Query Options

Specify the number of probes (1 by default)

```
SET ivfflat.probes = 10;
```

A higher value provides better recall at the cost of speed, and it can be set to the number of lists for exact nearest neighbor search (at which point the planner won't use the index)

Use SET LOCAL inside a transaction to set it for a single query

```
BEGIN;  
SET LOCAL ivfflat.probes = 10;  
SELECT ...  
COMMIT;
```

### Index Build Time

Speed up index creation on large tables by increasing the number of parallel workers (2 by default)

```
SET max_parallel_maintenance_workers = 7; -- plus leader
```

For a large number of workers, you may also need to increase max\_parallel\_workers (8 by default)

### Indexing Progress

Check indexing progress with Postgres 12+

```
SELECT phase, round(100.0 * tuples_done / nullif(tuples_total, 0), 1) AS "%" FROM pg_stat_p
```

The phases for IVFFlat are:

1. initializing
2. performing k-means
3. assigning tuples

#### 4. loading tuples

Note: % is only populated during the `loading tuples` phase

### Filtering

There are a few ways to index nearest neighbor queries with a `WHERE` clause

```
SELECT * FROM items WHERE category_id = 123 ORDER BY embedding <-> '[3,1,2]' LIMIT 5;
```

Create an index on one or more of the `WHERE` columns for exact search

```
CREATE INDEX ON items (category_id);
```

Or a partial index on the vector column for approximate search

```
CREATE INDEX ON items USING hnsw (embedding vector_l2_ops) WHERE (category_id = 123);
```

Use partitioning for approximate search on many different values of the `WHERE` columns

```
CREATE TABLE items (embedding vector(3), category_id int) PARTITION BY LIST(category_id);
```

### Hybrid Search

Use together with Postgres full-text search for hybrid search.

```
SELECT id, content FROM items, plainto_tsquery('hello search') query
WHERE textsearch @@ query ORDER BY ts_rank_cd(textsearch, query) DESC LIMIT 5;
```

You can use Reciprocal Rank Fusion or a cross-encoder to combine results.

### Performance

#### Tuning

Use a tool like PgTune to set initial values for Postgres server parameters.

#### Loading

Use `COPY` for bulk loading data (example).

```
COPY items (embedding) FROM STDIN WITH (FORMAT BINARY);
```

Add any indexes *after* loading the initial data for best performance.

#### Indexing

See index build time for HNSW and IVFFlat.

In production environments, create indexes concurrently to avoid blocking writes.

```
CREATE INDEX CONCURRENTLY ...
```

## Querying

Use EXPLAIN ANALYZE to debug performance.

```
EXPLAIN ANALYZE SELECT * FROM items ORDER BY embedding <-> '[3,1,2]' LIMIT 5;
```

**Exact Search** To speed up queries without an index, increase `max_parallel_workers_per_gather`.

```
SET max_parallel_workers_per_gather = 4;
```

If vectors are normalized to length 1 (like OpenAI embeddings), use inner product for best performance.

```
SELECT * FROM items ORDER BY embedding <#> '[3,1,2]' LIMIT 5;
```

**Approximate Search** To speed up queries with an IVFFlat index, increase the number of inverted lists (at the expense of recall).

```
CREATE INDEX ON items USING ivfflat (embedding vector_l2_ops) WITH (lists = 1000);
```

## Vacuuuming

Vacuuuming can take a while for HNSW indexes. Speed it up by reindexing first.

```
REINDEX INDEX CONCURRENTLY index_name;  
VACUUM table_name;
```

## Monitoring

Monitor performance with `pg_stat_statements` (be sure to add it to `shared_preload_libraries`).

```
CREATE EXTENSION pg_stat_statements;
```

Get the most time-consuming queries with:

```
SELECT query, calls, ROUND((total_plan_time + total_exec_time) / calls) AS avg_time_ms,  
       ROUND((total_plan_time + total_exec_time) / 60000) AS total_time_min  
FROM pg_stat_statements ORDER BY total_plan_time + total_exec_time DESC LIMIT 20;
```

Note: Replace `total_plan_time + total_exec_time` with `total_time` for Postgres < 13

Monitor recall by comparing results from approximate search with exact search.

```
BEGIN;  
SET LOCAL enable_indexscan = off; -- use exact search  
SELECT ...  
COMMIT;
```



## Scaling

Scale pgvector the same way you scale Postgres.

Scale vertically by increasing memory, CPU, and storage on a single instance. Use existing tools to tune parameters and monitor performance.

Scale horizontally with replicas, or use Citus or another approach for sharding (example).

## Languages

Use pgvector from any language with a Postgres client. You can even generate and store vectors in one language and query them in another.

Language	Libraries / Examples
C	pgvector-c
C++	pgvector-cpp
C#, F#, Visual Basic	pgvector-dotnet
Crystal	pgvector-crystal
Dart	pgvector-dart
Elixir	pgvector-elixir
Go	pgvector-go
Haskell	pgvector-haskell
Java, Kotlin, Groovy, Scala	pgvector-java
JavaScript, TypeScript	pgvector-node
Julia	pgvector-julia
Lisp	pgvector-lisp
Lua	pgvector-lua
Nim	pgvector-nim
OCaml	pgvector-ocaml
Perl	pgvector-perl
PHP	pgvector-php
Python	pgvector-python
R	pgvector-r
Ruby	pgvector-ruby, Neighbor
Rust	pgvector-rust
Swift	pgvector-swift
Zig	pgvector-zig

## Frequently Asked Questions

**How many vectors can be stored in a single table?** A non-partitioned table has a limit of 32 TB by default in Postgres. A partitioned table can have thousands of partitions of that size.

**Is replication supported?** Yes, pgvector uses the write-ahead log (WAL), which allows for replication and point-in-time recovery.

**What if I want to index vectors with more than 2,000 dimensions?** You'll need to use dimensionality reduction at the moment.

**Can I store vectors with different dimensions in the same column?** You can use `vector` as the type (instead of `vector(3)`).

```
CREATE TABLE embeddings (model_id bigint, item_id bigint, embedding vector, PRIMARY KEY (model_id, item_id));
```

However, you can only create indexes on rows with the same number of dimensions (using expression and partial indexing):

```
CREATE INDEX ON embeddings USING hnswn ((embedding::vector(3)) vector_l2_ops) WHERE (model_id = 123);
```

and query with:

```
SELECT * FROM embeddings WHERE model_id = 123 ORDER BY embedding::vector(3) <-> '[3,1,2]' LIMIT 5;
```

**Can I store vectors with more precision?** You can use the `double precision[]` or `numeric[]` type to store vectors with more precision.

```
CREATE TABLE items (id bigserial PRIMARY KEY, embedding double precision[]);
```

```
-- use {} instead of [] for Postgres arrays
```

```
INSERT INTO items (embedding) VALUES ('{1,2,3}'), ('{4,5,6}');
```

Optionally, add a check constraint to ensure data can be converted to the `vector` type and has the expected dimensions.

```
ALTER TABLE items ADD CHECK (vector_dims(embedding::vector) = 3);
```

Use expression indexing to index (at a lower precision):

```
CREATE INDEX ON items USING hnswn ((embedding::vector(3)) vector_l2_ops);
```

and query with:

```
SELECT * FROM items ORDER BY embedding::vector(3) <-> '[3,1,2]' LIMIT 5;
```

**Do indexes need to fit into memory?** No, but like other index types, you'll likely see better performance if they do. You can get the size of an index with:

```
SELECT pg_size_pretty(pg_relation_size('index_name'));
```

## Troubleshooting

**Why isn't a query using an index?** The query needs to have an `ORDER BY` and `LIMIT`, and the `ORDER BY` must be the result of a distance operator, not an expression.

```
-- index
ORDER BY embedding <=> '[3,1,2]' LIMIT 5;
```

```
-- no index
ORDER BY 1 - (embedding <=> '[3,1,2]') DESC LIMIT 5;
```

You can encourage the planner to use an index for a query with:

```
BEGIN;
SET LOCAL enable_seqscan = off;
SELECT ...
COMMIT;
```

Also, if the table is small, a table scan may be faster.

**Why isn't a query using a parallel table scan?** The planner doesn't consider out-of-line storage in cost estimates, which can make a serial scan look cheaper. You can reduce the cost of a parallel scan for a query with:

```
BEGIN;
SET LOCAL min_parallel_table_scan_size = 1;
SET LOCAL parallel_setup_cost = 1;
SELECT ...
COMMIT;
```

or choose to store vectors inline:

```
ALTER TABLE items ALTER COLUMN embedding SET STORAGE PLAIN;
```

**Why are there less results for a query after adding an HNSW index?** Results are limited by the size of the dynamic candidate list (`hnsw.ef_search`). There may be even less results due to dead tuples or filtering conditions in the query. We recommend setting `hnsw.ef_search` to at least twice the `LIMIT` of the query. If you need more than 500 results, use an IVFFlat index instead.

**Why are there less results for a query after adding an IVFFlat index?** The index was likely created with too little data for the number of lists. Drop the index until the table has more data.

```
DROP INDEX index_name;
```

Results can also be limited by the number of probes (`ivfflat.probes`).

## Reference

### Vector Type

Each vector takes  $4 * \text{dimensions} + 8$  bytes of storage. Each element is a single precision floating-point number (like the `real` type in Postgres), and all

elements must be finite (no NaN, Infinity or -Infinity). Vectors can have up to 16,000 dimensions.

### Vector Operators

Operator	Description	Added
+	element-wise addition	
-	element-wise subtraction	
*	element-wise multiplication	0.5.0
<->	Euclidean distance	
<#>	negative inner product	
<=>	cosine distance	

### Vector Functions

Function	Description	Added
cosine_distance(vector, vector) → double precision	cosine distance	
inner_product(vector, vector) → double precision	inner product	
l2_distance(vector, vector) → double precision	Euclidean distance	
l1_distance(vector, vector) → double precision	taxicab distance	0.5.0
vector_dims(vector) → integer	number of dimensions	
vector_norm(vector) → double precision	Euclidean norm	

### Aggregate Functions

Function	Description	Added
avg(vector) → vector	average	
sum(vector) → vector	sum	0.5.0

## Installation Notes - Linux and Mac

### Postgres Location

If your machine has multiple Postgres installations, specify the path to pg\_config with:

```
export PG_CONFIG=/Library/PostgreSQL/16/bin/pg_config
```

Then re-run the installation instructions (run `make clean` before `make` if needed). If `sudo` is needed for `make install`, use:

```
sudo --preserve-env=PG_CONFIG make install
```

A few common paths on Mac are:

- EDB installer - `/Library/PostgreSQL/16/bin/pg_config`
- Homebrew (arm64) - `/opt/homebrew/opt/postgresql@16/bin/pg_config`
- Homebrew (x86-64) - `/usr/local/opt/postgresql@16/bin/pg_config`

Note: Replace 16 with your Postgres server version

### Missing Header

If compilation fails with `fatal error: postgres.h: No such file or directory`, make sure Postgres development files are installed on the server.

For Ubuntu and Debian, use:

```
sudo apt install postgresql-server-dev-16
```

Note: Replace 16 with your Postgres server version

### Missing SDK

If compilation fails and the output includes `warning: no such sysroot directory on Mac`, reinstall Xcode Command Line Tools.

### Portability

By default, pgvector compiles with `-march=native` on some platforms for best performance. However, this can lead to `Illegal instruction` errors if trying to run the compiled extension on a different machine.

To compile for portability, use:

```
make OPTFLAGS=""
```

## Installation Notes - Windows

### Missing Header

If compilation fails with `Cannot open include file: 'postgres.h': No such file or directory`, make sure `PGROOT` is correct.

### Permissions

If installation fails with `Access is denied`, re-run the installation instructions as an administrator.

## Additional Installation Methods

### Docker

Get the Docker image with:

```
docker pull pgvector/pgvector:pg16
```

This adds pgvector to the Postgres image (replace 16 with your Postgres server version, and run it the same way).

You can also build the image manually:

```
git clone --branch v0.6.2 https://github.com/pgvector/pgvector.git
cd pgvector
docker build --build-arg PG_MAJOR=16 -t myuser/pgvector .
```

### Homebrew

With Homebrew Postgres, you can use:

```
brew install pgvector
```

Note: This only adds it to the postgresql@14 formula

### PGXN

Install from the PostgreSQL Extension Network with:

```
pgxn install vector
```

### APT

Debian and Ubuntu packages are available from the PostgreSQL APT Repository. Follow the setup instructions and run:

```
sudo apt install postgresql-16-pgvector
```

Note: Replace 16 with your Postgres server version

### Yum

RPM packages are available from the PostgreSQL Yum Repository. Follow the setup instructions for your distribution and run:

```
sudo yum install pgvector_16
```

```
# or
```

```
sudo dnf install pgvector_16
```

Note: Replace 16 with your Postgres server version

### pkg

Install the FreeBSD package with:

```
pkg install postgresql15-pg_vector
```

or the port with:

```
cd /usr/ports/databases/pgvector
make install
```

## conda-forge

With Conda Postgres, install from conda-forge with:

```
conda install -c conda-forge pgvector
```

This method is community-maintained by [@mmcauliffe](https://github.com/mmcauliffe)

## Postgres.app

Download the latest release with Postgres 15+.

## Hosted Postgres

pgvector is available on these providers.

## Upgrading

Install the latest version (use the same method as the original installation). Then in each database you want to upgrade, run:

```
ALTER EXTENSION vector UPDATE;
```

You can check the version in the current database with:

```
SELECT extversion FROM pg_extension WHERE extname = 'vector';
```

## Upgrade Notes

### 0.6.0

**Postgres 12** If upgrading with Postgres 12, remove this line from `sql/vector--0.5.1--0.6.0.sql`:

```
ALTER TYPE vector SET (STORAGE = external);
```

Then run `make install` and `ALTER EXTENSION vector UPDATE;`.

**Docker** The Docker image is now published in the `pgvector` org, and there are tags for each supported version of Postgres (rather than a `latest` tag).

```
docker pull pgvector/pgvector:pg16
```

```
# or
```

```
docker pull pgvector/pgvector:0.6.0-pg16
```

Also, if you've increased `maintenance_work_mem`, make sure `--shm-size` is at least that size to avoid an error with parallel HNSW index builds.

```
docker run --shm-size=1g ...
```

## Thanks

Thanks to:

- PASE: PostgreSQL Ultra-High-Dimensional Approximate Nearest Neighbor Search Extension
- Faiss: A Library for Efficient Similarity Search and Clustering of Dense Vectors
- Using the Triangle Inequality to Accelerate k-means
- k-means++: The Advantage of Careful Seeding
- Concept Decompositions for Large Sparse Text Data using Clustering
- Efficient and Robust Approximate Nearest Neighbor Search using Hierarchical Navigable Small World Graphs

## History

View the changelog

## Contributing

Everyone is encouraged to help improve this project. Here are a few ways you can help:

- Report bugs
- Fix bugs and submit pull requests
- Write, clarify, or fix documentation
- Suggest or add new features

To get started with development:

```
git clone https://github.com/pgvector/pgvector.git
cd pgvector
make
make install
```

To run all tests:

```
make installcheck          # regression tests
make prove_installcheck   # TAP tests
```

To run single tests:

```
make installcheck REGRESS=functions          # regression test
make prove_installcheck PROVE_TESTS=test/t/001_ivfflat_wal.pl # TAP test
```

To enable assertions:

```
make clean && PG_CFLAGS="-DUSE_ASSERT_CHECKING" make && make install
```

To enable benchmarking:

```
make clean && PG_CFLAGS="-DIVFFLAT_BENCH" make && make install
```



To show memory usage:

```
make clean && PG_CFLAGS="-DHNSW_MEMORY -DIVFFLAT_MEMORY" make && make install
```

To get k-means metrics:

```
make clean && PG_CFLAGS="-DIVFFLAT_KMEANS_DEBUG" make && make install
```

Resources for contributors

- Extension Building Infrastructure
- Index Access Method Interface Definition
- Generic WAL Records