

From Distance to Intelligence: Evolving Your Vector Search

Bilge Ince & Boriss Mejias



Bilge INCE MLE @ EDB Muay Thai, Running 🧡





@abilgegunduz

@bilge.bsky.social





Boriss Mejías

Solutions architect @ EDB Organizer of PGDay Lowlands Air guitar \m/

@tchorix@mastodon.world

@tchorix.bsky.social





Find 11 rock music outfits



















LLMs







Search: Rock music outfits Filter: women Limit to top 11





Rock music outfit for women



CREATE TABLE products (

- id BIGINT PRIMARY KEY
- , img_id BIGINT REFERENCES products_images(id)
- , description TEXT
- , gender TEXT

);



CREATE TABLE products (

- id BIGINT PRIMARY KEY
- , img_id BIGINT REFERENCES products_images(id)
- , description TEXT
- , gender TEXT

CREATE TABLE products_embeddings (product_id BIGINT REFERENCES products(id) , embedding vector(384)

);



);





























Why Getting Search Right Matters

Finding exactly the right product vs. irrelevant suggestions; Slack, Healthcare sector, e-commerce, game





Floral Dress for the summer party



The Promise of Understanding Semantic Meaning

- Finds **related concepts** even without keyword overlap.
- Excellent for **discovery** and **recommendation**.



The Promise of Understanding Semantic Meaning

- Finds related concepts even without keyword overlap.
- Excellent for discovery and recommendation.
- Handles **synonyms** and related concepts naturally.



The Promise of Understanding Semantic Meaning

- Finds related concepts even without keyword overlap.
- Excellent for discovery and recommendation.
- Handles synonyms and related concepts naturally.
- Enables **natural language querying** and Question Answering.
- Can search across modalities (text-image, text-speech).



The Reality Behind The Hype

"Nearest" ≠ "Best"

The "Black Box"

Computational Cost

Keyword Blindness

Context Ambiguity



What can we do?





©EDB 2025 - ALL RIGHTS RESERVED.

What can we do?





CREATE TABLE products (

- id BIGINT PRIMARY KEY
- , img_id BIGINT REFERENCES products_images(id)
- , description TEXT
- , gender TEXT

CREATE TABLE products_embeddings (product_id BIGINT REFERENCES products(id) , embedding vector(384)

);



);

CREATE TABLE products (

- id BIGINT PRIMARY KEY
- , img id BIGINT REFERENCES products images(id)
- , description TEXT
- , gender TEXT
- , embedding vector(384)

);



CREATE TABLE products (

- id BIGINT PRIMARY KEY
- , img_id BIGINT REFERENCES products_images(id)
- , description TEXT
- , gender TEXT
- , embedding vector(384)

);

CREATE INDEX ON products USING hnsw (embedding vector_cosine_ops)
WHERE gender = 'women';



What can we do?





What can we do?

















The Solution: Hybrid Search - Union

- Run Full Text Search and Vector Search queries independently.
- Get two ranked lists of results.
- Combine these lists using an algorithm. Popular option: Reciprocal Rank Fusion







The Solution: Hybrid Search - Intersect

- Leverage the best of both worlds
- Find semantically relevant items that also satisfy keyword constraints or other filters.

The Method:

Use full text search for essential keywords to get an initial set of candidates (e.g., top 100-500 results). Then run semantic search on the candidate data for user query that fits with semantic need.

For example; "summer outfit for Thailand trip under €50"
Summer outfit for Thailand - requires semantic search because...
<€50 - is a filtering.











Hybrid Search: using pgvector

"summer outfit for Thailand trip under €50"

Indexes: HNSW & IVFFlat are Approximate Nearest Neighbor (ANN) indexes.

Standard ANN index scans (especially HNSW in pgvector) often struggle to efficiently apply WHERE clause filters during the nearest neighbor search.

Common Issue: Retrieve Top K vector neighbours -> Then apply metadata filters. **Result:** The desired k value didn't meet as a result, it may not return an output at all.



Hierarchical Navigable Small Worlds

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

SET hnsw.ef_search = 100;





Hierarchical Navigable Small Worlds

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

SET hnsw.ef_search = 100;





Inverted File Flat

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

SET ivfflat.probes = 32





Ideal Scenario

```
-- Select top 5 similar image IDs
SELECT
    img id,
    (embedding <=> '{text embeddings}') AS score
    -- Calculate similarity score
FROM
    products embeddings
WHERE
    gender = '{gender}' -- Filter by gender
ORDER BY
    score
LIMIT 11;
```



```
The Scenario I can Use?
```

```
-- 1. Filter products by gender
WITH filtered products AS (
    SELECT img id, productdisplayname
    FROM products pgconf
    WHERE gender = '{selected gender}'
-- 2. Find top 100 similar embeddings for filtered products
SELECT
    r.img id, r.score
FROM
    filtered products fp
CROSS JOIN LATERAL (
    SELECT img id, (embedding <=> '{text embeddings}') AS score
    FROM products embeddings pgvector
    ORDER BY score
    LIMIT 100
) AS r
-- 3. Join and get the top 5 most similar
WHERE r.img id = fp.img id
ORDER BY r.score LIMIT 11;
```

```
The Scenario I can Use?
```

```
-- 1. Filter products by gender
WITH filtered products AS (
    SELECT img id, productdisplayname
    FROM products pgconf
    WHERE gender = '{selected gender}'
-- 2. Find top 100 similar embeddings for filtered products
SELECT
    r.img id, r.score
FROM
    filtered products fp
CROSS JOIN LATERAL (
    SELECT img id, (embedding <=> '{text embeddings}') AS score
    FROM products embeddings pgvector
    ORDER BY score
   LIMIT 100
) AS r
-- 3. Join and get the top 5 most similar
WHERE r.img id = fp.img id
ORDER BY r.score LIMIT 11;
```

Hybrid Search with Intersection: using pgvector





Hybrid Search: using pgvector v0.8.0

Iterative Index Scans

With approximate indexes, queries with filtering can return less results since filtering is applied *after* the index is scanned. Starting with 0.8.0, you can enable iterative index scans, which will automatically scan more of the index until enough results are found (or it reaches hnsw.max_scan_tuples or ivfflat.max_probes).

Iterative scans can use strict or relaxed ordering.

Strict ensures results are in the exact order by distance

SET hnsw.iterative_scan = strict_order;

Relaxed allows results to be slightly out of order by distance, but provides better recall

SET hnsw.iterative_scan = relaxed_order; # or SET ivfflat.iterative_scan = relaxed_order;

With relaxed ordering, you can use a materialized CTE to get strict ordering

WITH relaxed_results AS MATERIALIZED (
 SELECT id, embedding <-> '[1,2,3]' AS distance FROM items WHERE category_id = 123 ORDER BY
) SELECT * FROM relaxed_results ORDER BY distance;

For queries that filter by distance, use a materialized CTE and place the distance filter outside of it for best performance (due to the <u>current behavior</u> of the Postgres executor)

WITH nearest_results AS MATERIALIZED (
 SELECT id, embedding <-> '[1,2,3]' AS distance FROM items ORDER BY distance LIMIT 5
) SELECT * FROM nearest_results WHERE distance < 5 ORDER BY distance;</pre>



VectorChord Vchord+Vchord_BM25

- VectorChord-BM25 is a PostgreSQL extension for keyword search. It not only implements BM25 ranking but also includes a tokenizer and a Block-WeakAnd index to improve speed.
- Semantic Search is utilizing RabitQ algorithm.
 - RabitQ is a quantization algorithm for high-dimensional spaces, designed to improve the storage and retrieval efficiency of high-dimensional data, such as embedded vectors.



















Closing words

- Full text is not enough
- Vector search is not enough
- Existing tools only solves part of the problem or too expensive (and complicated)
- Iterative scans looks promising yet not cost-efficient





- 1. Hands on LLMs Jay Alammar & Maarten Grootendorst
- 2. <u>https://github.com/pgvector/pgvector/issues/259</u>
- 3. <u>https://github.com/pgvector/pgvector?tab=readme-ov-file#iterative-index-scans</u>
- 4. <u>https://github.com/pgvector/pgvector/issues/301</u>
- 5. <u>https://blog.vectorchord.ai/hybrid-search-with-postgres-native-bm25-and-vectorchord</u>



Any Questions?

Bilge Ince @bilge.bsky.social bilge.ince@enterprisedb.com Boriss Mejías @tchorix.bsky.social boriss.mejias@enterprisedb.com

