# Hacking pgvector for performance



# Daniel Krefl Sednai

\$1.\$2 @ sedn.ai















#### Project funded by



Federal Department of Economic Affairs, Education and Research EAER State Secretariat for Education, Research and Innovation SERI

Swiss Confederation

Funded by the European Union. Views and opinions are however those of the author(s) only and do not necessarily reflect those of the European Union or the HaDEA. Neither the European Union nor the granting authority can be held responsible for them. Project number: 101092850.

AERO has also received funding from UKRI under grants no. 10048318 and 10048915, and the Swiss State Secretariat for Education, Research, and Innovation.





























#### Target platform:



HPC and Al processor

European high-performance energy-efficient processor (ARM based), dedicated to high performance computing, and designed to work with third-party accelerators, see [https://sipearl.com/]

RHEA images kindly provided by SIPEARL



























#### Target platform:



#### [ https://sipearl.com/ ]

RHEA images kindly provided by SIPEARL





























#### Objectives:

[ https://aero-project.eu/about/ ]

- Managed Programming Languages
- Native Programming Languages & Runtimes
- OS, drivers & virtualization support
- State-of-the-art cloud deployments
- Hardware acceleration for performance & security
- Adoption of the EU cloud ecosystem



























## Use cases / pilots:

- Automotive Digital Twins with IoT-Cloud Interoperability
- High Performance Algorithms for Space Exploration (Gaia)
- HPC/Cloud Database Acceleration for Scientific Computing



























## Use cases / pilots:

- Automotive Digital Twins with IoT-Cloud Interoperability
- High Performance Algorithms for Space Exploration (Gaia)
- HPC/Cloud Database Acceleration for Scientific Computing

Leadership / Coordination of pilots



























## Use cases / pilots:

- Automotive Digital Twins with IoT-Cloud Interoperability
- High Performance Algorithms for Space Exploration (Gaia)
- HPC/Cloud Database Acceleration for Scientific Computing

























... GAIA ...

# Gaia





Gaia was a space based astronomy telescope of ESA operational 2014-2025.

Gaia has made more than three trillion observations of two billion stars and other objects throughout our Milky Way galaxy and beyond, mapping their motions, luminosity, temperature and composition.

#### Scientific objectives:

- First 3d map of our galaxy
- Insides on the origin and formation of our galaxy
- Detection of diverse variable phenomena
- Many more ... see [ https://www.cosmos.esa.int/web/gaia/science ]

# Gaia





Gaia was a space based astronomy telescope of ESA operational 2014-2025.

Gaia has made more than three trillion observations of two billion stars and other objects throughout our Milky Way galaxy and beyond, mapping their motions, luminosity, temperature and composition.

#### Scientific objectives:

- First 3d map of our galaxy
- Insides on the origin and formation of our galaxy
- Detection of diverse variable phenomena
- Many more ... see [ https://www.cosmos.esa.int/web/gaia/science ]

# Gaia





Gaia was a space based astronomy telescope of ESA operational 2014-2025.

Gaia has made more than three trillion observations of two billion stars and other objects throughout our Milky Way galaxy and beyond, mapping their motions, luminosity, temperature and composition.

#### Data and compute challenge:

- Petabyte scale
- ~ 10 Billion photometric time series
- ~ 5 Billion spectra time series

We solve this with Postgres!





## Intertwined Gaia+SED pilots:

- Process efficiently constantly increasing volumes of data.
  - Enable GPU computations directly within the database
  - Optimize data and compute pipeline for RHEA



























## Intertwined Gaia+SED pilots:

- Process efficiently constantly increasing volumes of data.
  - Enable GPU computations directly within the database

Example: Hack to GPU accelerate a Postgres vector index

(WARNING: This will be technical ...)



























... VECTOR SEARCH ...



#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.



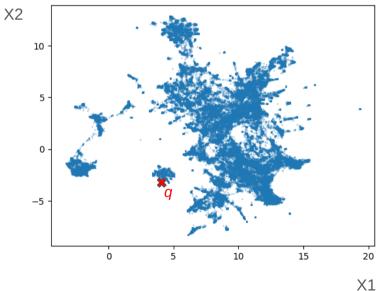
#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.

For a query point *q*:

What are its *k* nearest neighbors?

Wikipedia - OpenAI embeddings 1536d -> 2d projection (UMAP for 100k subset)



Dataset generated by S. Sturges, 2023 ] (https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings)



#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.

For a query point *q*:

What are its *k* nearest neighbors?

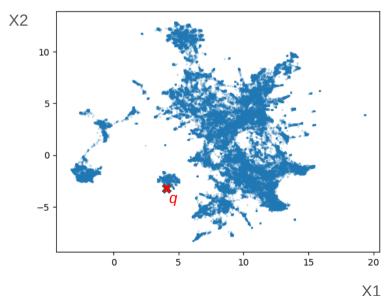
#### Note:

Growing interest due to ML / Al generated vector embeddings.

For instance:

Retrieval-Augmented Generation

Wikipedia - OpenAI embeddings 1536d -> 2d projection (UMAP for 100k subset)



[ Dataset generated by S. Sturges, 2023 ] (https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings



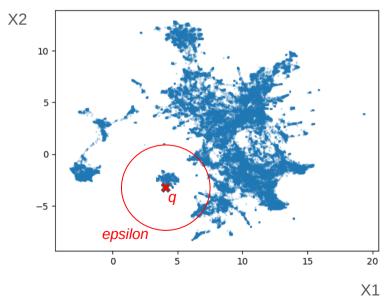
#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.

For a query point *q*:

What points are close by?

Wikipedia - OpenAI embeddings 1536d -> 2d projection (UMAP for 100k subset)



[ Dataset generated by S. Sturges, 2023 ] (https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings)



#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.

For a query point *q*:

What points are close by?

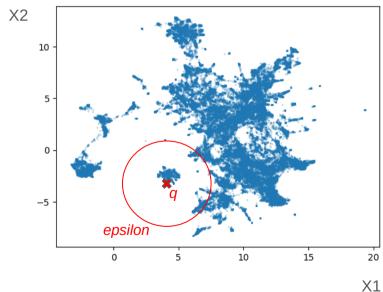
#### Note:

Of interest for classical unsupervised learning algorithms.

kNN, DBSCAN, ...







Dataset generated by S. Sturges, 2023 ] https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings)



#### **Motivation**

Many data analysis algorithms require a nearest neighbour search in a D-dim space. Often just referred to as *Vector Search*.

#### BUT:

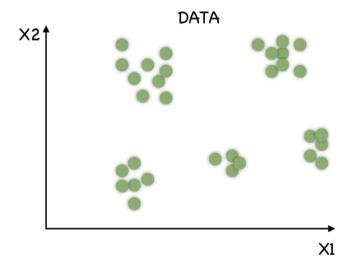
Requires for each query point N distance calculations + ranking. (With N the number of datapoints in the dataset )

How to scale to large datasets?



## **Approximate Nearest Neighbour search**

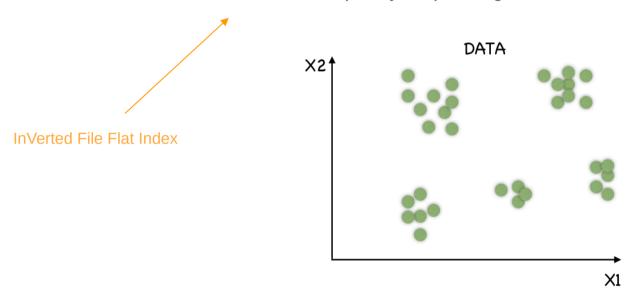
Several flavours exist, but *IVFFLAT* is the conceptually simplest algorithm.





# **Approximate Nearest Neighbour search**

Several flavours exist, but *IVFFLAT* is the conceptually simplest algorithm.

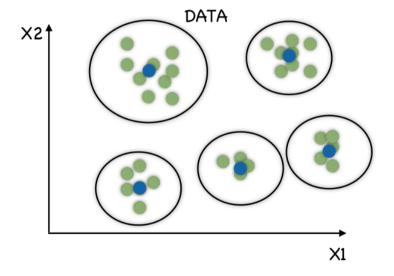




#### **IVFFLAT**

Several flavours exist, but *IVFFLAT* is the conceptually simplest algorithm.

- Build vector index on dataset via K-means clustering
- Index each datapoint to the closest cluster (centroid)



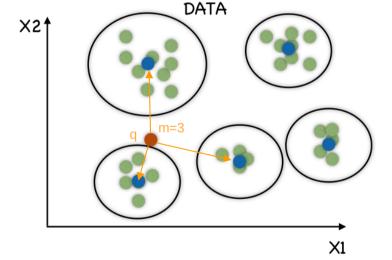


#### **IVFFLAT**

Several flavours exist, but *IVFFLAT* is the conceptually simplest algorithm.

- Build vector index on dataset via K-means clustering
- Index each datapoint to the closest cluster (centroid)
- Evaluate query point q only against cluster members of the m nearest centroids.

Recall / Performance tradeoff



#### BUT:

- Strongly depends on distribution of data
- High recall may effectively result in brute-force scan



#### **IVFFLAT**

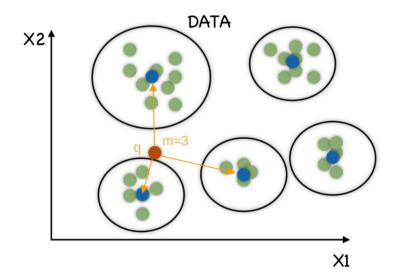
Several flavours exist, but *IVFFLAT* is the conceptually simplest algorithm.

- Build vector index on dataset via K-means clustering
- Index each datapoint to the closest cluster (centroid)
- Evaluate query point *q* only against cluster members of the *m* nearest centroids.

Recall / Performance tradeoff



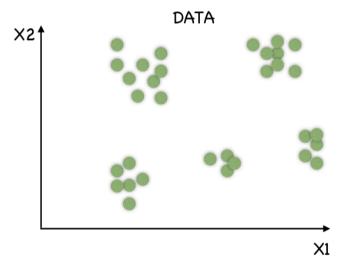
- Relatively fast and cheap to build index
- Very suitable for parallelization





## **Approximate Nearest Neighbour search**

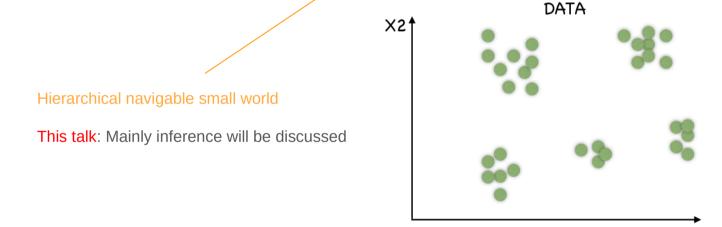
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





## **Approximate Nearest Neighbour search**

A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.

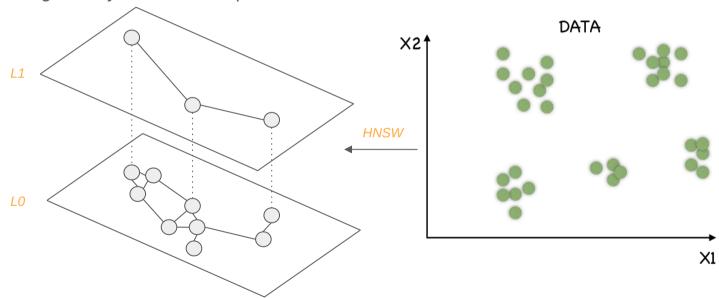


X1



#### **HNSW** inference

A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.

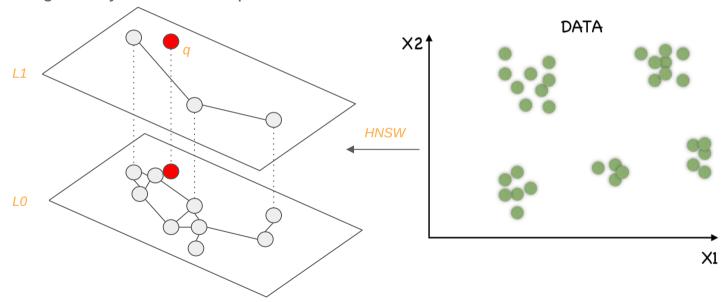


Multi-layered graph
(figure is simplified, usually many more layers)



#### **HNSW** inference

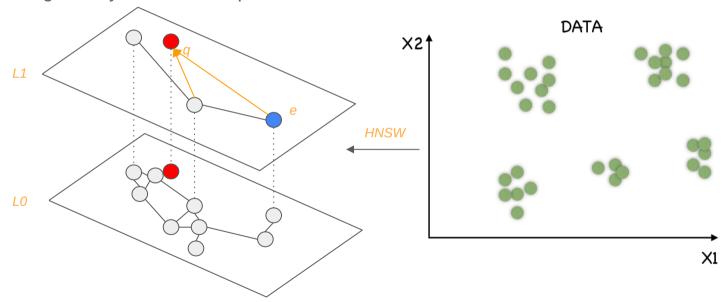
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

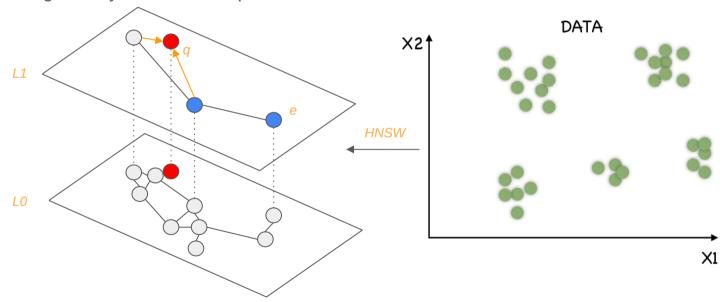
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

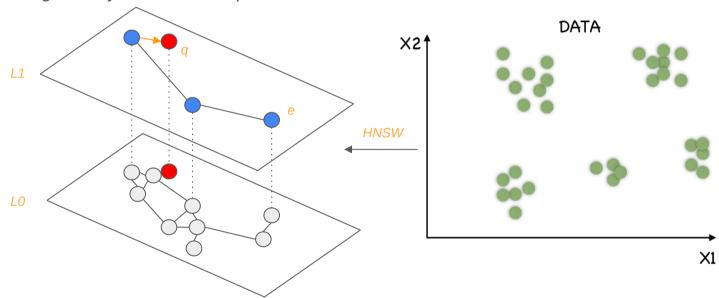
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

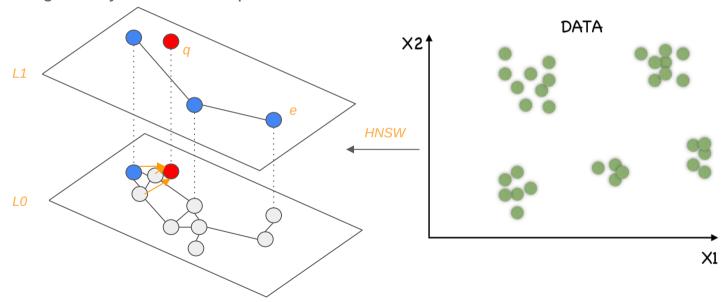
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

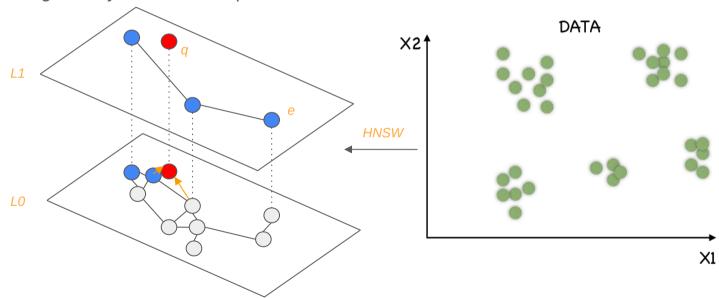
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

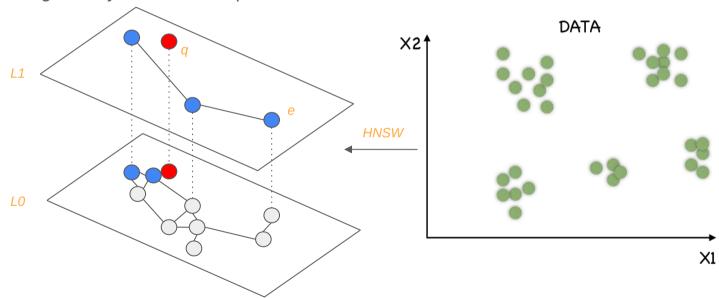
A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.





#### **HNSW** inference

A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.

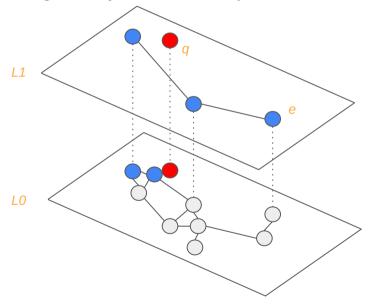


Search via greedy graph traversal (keeping closest k visited nodes along the way)



#### **HNSW** inference

A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.



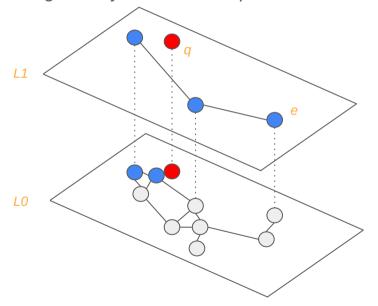
#### BUT:

- Graph construction can be very compute and memory intensive
- Difficult to parallelize and distribute
- Large top *k* with filters tricky



#### **HNSW** inference

A bit more conceptually challenging is *HNSW*. But nowadays often preferred algorithm due to generally better recall / performance behavior.



#### BUT:

- Graph construction can be very compute and memory intensive
- Difficult to parallelize and distribute
- Large top *k* with filters tricky



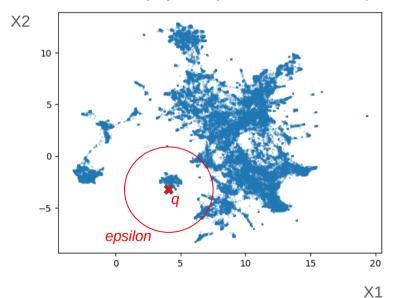
#### **HNSW** inference

"Large top k with filters tricky"

For a query point *q*:

What points are close by?

Wikipedia - OpenAI embeddings 1536d -> 2d projection (UMAP for 100k subset)



[ Dataset generated by S. Sturges, 2023 ] (https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings)



#### **HNSW** inference

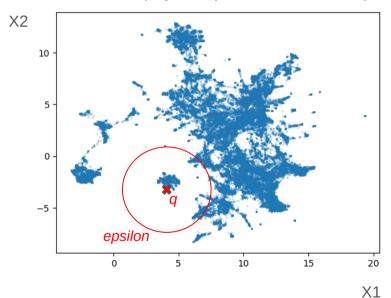
"Large top k with filters tricky"

For a query point *q*:

What points are close by?

IVFFLAT of main interest for this talk.

Wikipedia - OpenAI embeddings 1536d -> 2d projection (UMAP for 100k subset)



[ Dataset generated by S. Sturges, 2023 ] (https://www.kaggle.com/datasets/stephanst/wikipedia-simple-openai-embeddings)



### **Approximate Nearest Neighbour search**

#### Remark:

There are also newer, more specialized, algorithms which gain in popularity.

### For example:

DiskANN: Applicable to large scale problems (beyond RAM). [Subramanya et al., 2019]

CAGRA: HNSW variant optimized for GPU execution. [Ootomo et al., 2024]



### **Postgres implementation**

Most well-known and popular: *pgvector* ( <a href="https://github.com/pgvector/pgvector">https://github.com/pgvector/pgvector</a>)

Introduces new PG type (vector), distance operators acting on vectors, and can build indices for vector columns.

*Exact*, and *HNSW* or *IVFFLAT* based approximate vector search. Many different metrics, but here we are mainly interested in euclidean distance (<->) and cosine distance (<=>).

### Query example:

select \* from table where embedding <-> '[...]' < 10 order by embedding <-> '[...]' limit 10000



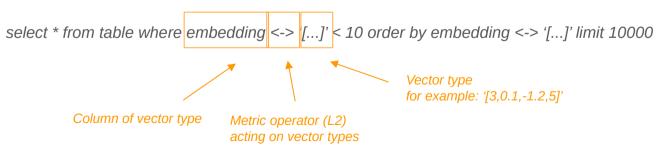
#### **Postgres implementation**

Most well-known and popular: *pgvector* ( <a href="https://github.com/pgvector/pgvector">https://github.com/pgvector/pgvector</a>)

Introduces new PG type (vector), distance operators acting on vectors, and can build indices for vector columns.

*Exact*, and *HNSW* or *IVFFLAT* based approximate vector search. Many different metrics, but here we are mainly interested in euclidean distance (<->) and cosine distance (<=>).

### Query example:



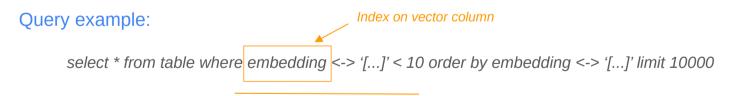


#### **Postgres implementation**

Most well-known and popular: *pgvector* ( <a href="https://github.com/pgvector/pgvector">https://github.com/pgvector/pgvector</a>)

Introduces new PG type (vector), distance operators acting on vectors, and can build indices for vector columns.

*Exact*, and *HNSW* or *IVFFLAT* based approximate vector search. Many different metrics, but here we are mainly interested in euclidean distance (<->) and cosine distance (<=>).



Reduce number of required distance calculations via "index" / approx vector search. (Worst case: For all rows)



### **Postgres implementation**

Newcomer: *pgvectorscale* (timescale)

( https://github.com/timescale/pgvectorscale/ )

DiskANN based approximate vector search.

Offers euclidean distance, cosine distance and inner product.

Interesting feature: Label filter push down into index scan.



### **Postgres implementation**

Newcomer: pgvectorscale (timescale)

( https://github.com/timescale/pavectorscale/ )

DiskANN based approximate vector search.

Offers euclidean distance, cosine distance and inner product.

Interesting feature: Label filter push down into index scan.

Note:

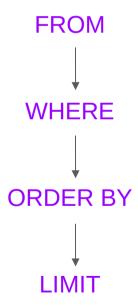
Mentioned only for completeness. In this talk I will not dive further into *pgvectorscale*.



... PGVECTOR ...



# pgvector



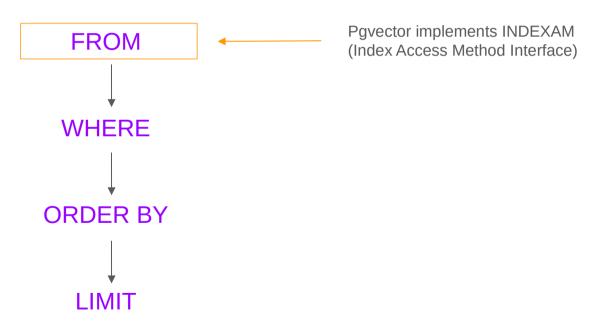


### pgvector



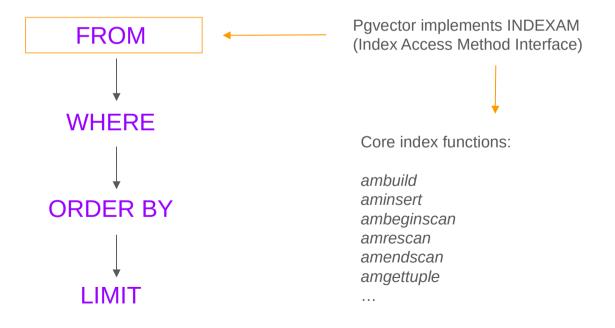


### pgvector



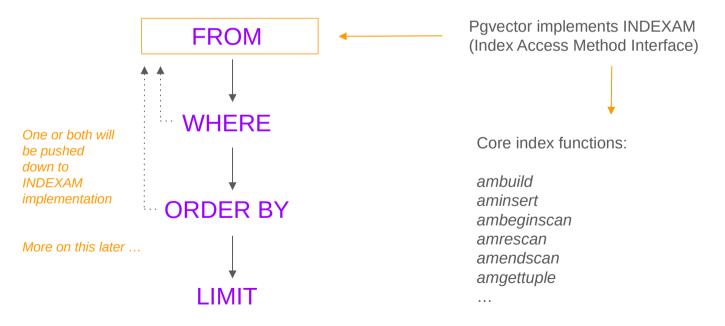


#### pgvector





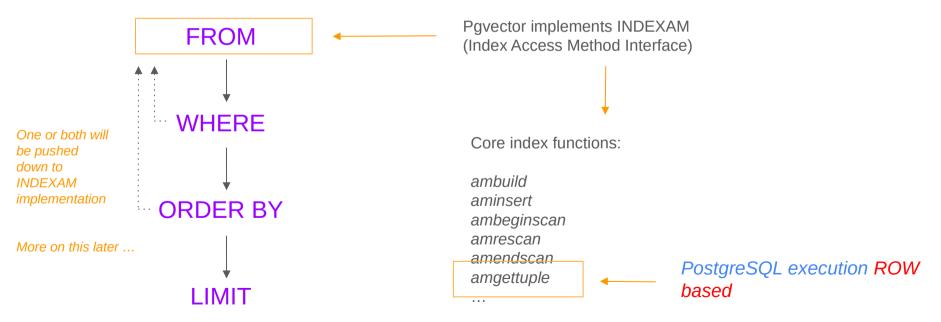
#### pgvector





#### pgvector

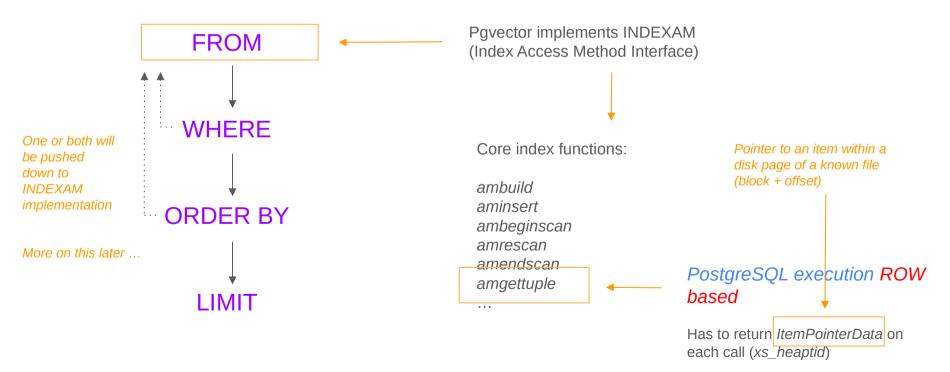
### PostgreSQL plan generation (simplified)



Has to return *ItemPointerData* on each call (xs\_heaptid)

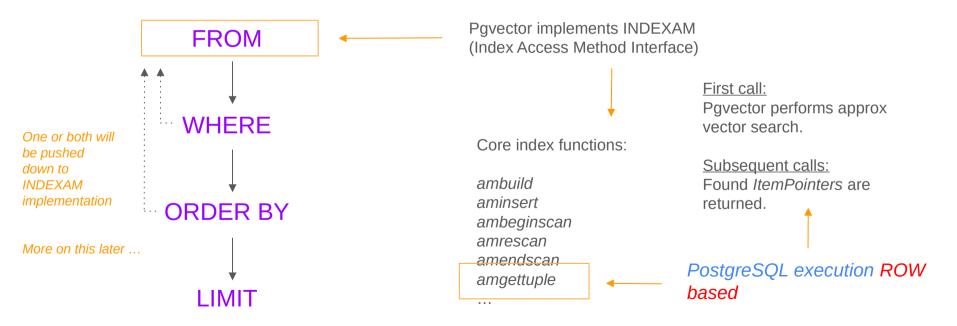


#### pgvector





#### pgvector





### pgvector

#### What about performance?

#### Hardware:

- i9-13900H + RTX 4070 [35W]

#### Software:

- Ubuntu 22.04
- gcc/g++-14; icpx
- Python 3.10.12 + psycopg2 2.9.6
- Standard Postgres v15 (no special compile flags besides -g)
- PGvector master on Mar 24, 2025 (commit 05182479a2a62e04300386b4da18be02fcb819b5) (compiled with -O3 -march=native -g)



#### pgvector

#### What about performance?

### Settings:

- Queried locally via python+psycopg2 (on persistent connection)
- IVFFLAT: 200 clusters
- Recall computation for IVFFLAT via variation of # probes. (Smallest # probes for fixed recall + median at fixed recall)
- Recall computation for HNSW via variation of *ef\_search* parameter. (Smallest parameter for fixed recall + median at fixed recall)



#### pgvector

#### What about performance?

# Settings:

- Queried locally via python+psycopg2
   (on persistent connection)
- IVFFLAT: 200 clusters
- Recall computation for IVFFLAT via variation of # probes. (Smallest # probes for fixed recall + median at fixed recall)
- Recall computation for HNSW via variation of *ef\_search* parameter. (Smallest parameter for fixed recall + median at fixed recall)

#### Variation in discrete steps

IVFFLAT: 10 probes given by cutting points of [1,100] into equal size ranges HNSW: [100, 200,..., 1000]



### pgvector

#### What about performance?

GIST-960 dataset [Jégou, Douze and Schmid, 2011] (http://corpus-texmex.irisa.fr/)

- 1M vectors
- 960d
- 1k test vectors
- Pre-computed 100 nearest neighbors (euclidean metric)

(Vectors given by global GIST descriptors of image dataset. GIST summarizes the gradient information for different parts of an image.)

### DEEP1B dataset [Babenko and Lempitsky, 2016] (https://www.tensorflow.org/datasets/catalog/deep1b)

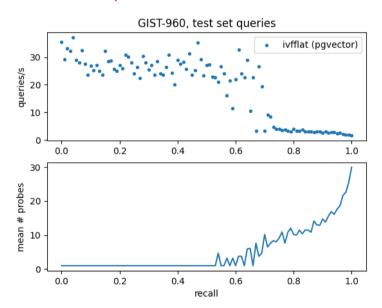
- 10M vectors (subset of 1B)
- 96d
- 10k test vectors (we use first 1k)
- Pre-computed 100 nearest neighbors (cosine metric)

(Vectors given by PCA of 1B image embeddings produced as outputs from the last fully-connected layer of a GoogLeNet CNN model)

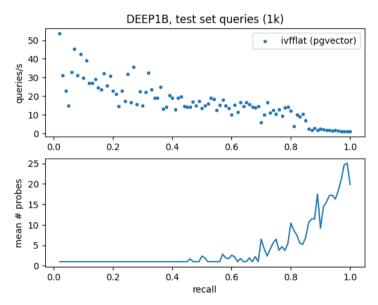


### pgvector

#### What about performance?



select id from gist order by embedding <-> "+q+" limit 100

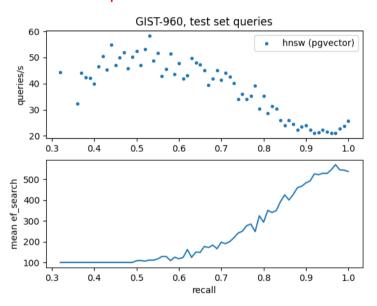


select id from deep1b order by embedding <=> "+q+" limit 100

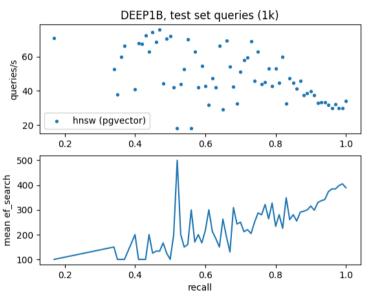


### pgvector

#### What about performance?



select id from gist order by embedding <-> "+q+" limit 100

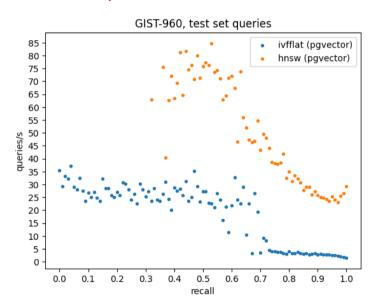


select id from deep1b order by embedding <=> "+q+" limit 100

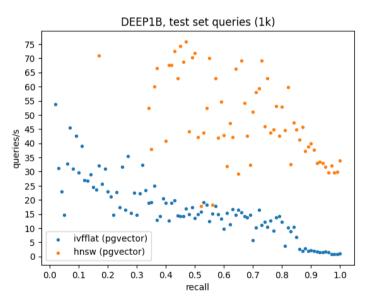


### pgvector

### What about performance?



select id from gist order by embedding <-> "+q+" limit 100



select id from deep1b order by embedding <=> "+q+" limit 100

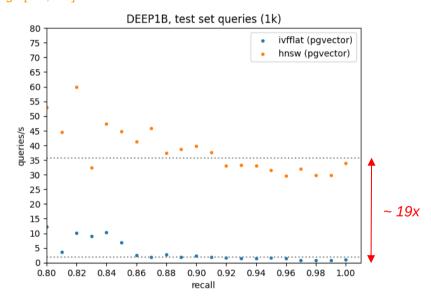


#### pgvector

#### What about performance?

#### GIST-960, test set queries ivfflat (pgvector) 85 hnsw (pgvector) 80 -75 70 65 60 55 45 35 30 20 15 $\sim 9x$ 10 0.84 0.86 0.88 0.90 0.92 recall

#### Median in range [0.8, 1.0]



select id from gist order by embedding <-> "+q+" limit 100

select id from deep1b order by embedding <=> "+q+" limit 100



### pgvector (ivfflat)

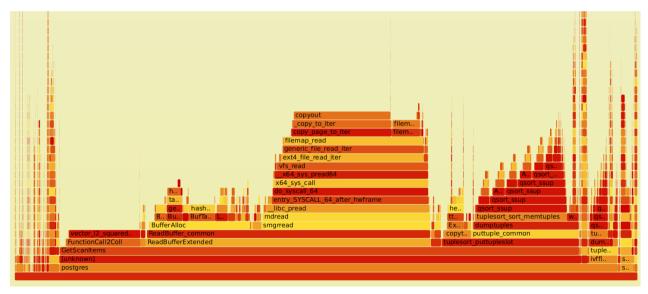
What about performance?

GIST-960 dataset (1M, 960d)

#### Time thieves:

- ReadBuffers
- Tuplesort
- Distance calculation

## Ivfflat inference perf analysis



[Generated with FlameGraph] ( https://github.com/brendangregg/FlameGraph )



### pgvector (ivfflat)

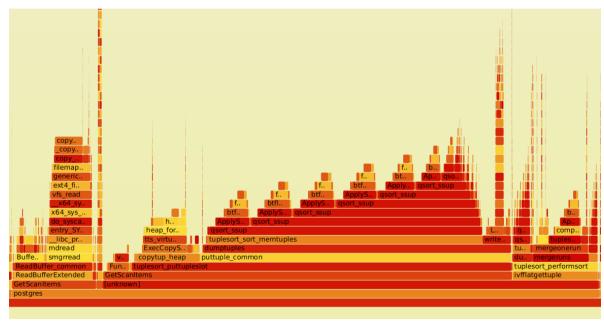
What about performance?

DEEP-1B dataset (10M, 96d)

#### Time thieves:

- Tuplesort
- ReadBuffers
- Distance calculation

### Ivfflat inference perf analysis



[Generated with FlameGraph] ( https://github.com/brendangregg/FlameGraph )



#### pgvector (ivfflat)

#### What about performance?

#### Time thieves:

- ReadBuffers
- Tuplesort
- Distance calculation

Portions Copyright (c) 1996-2024, PostgreSQL Global Development Group

Portions Copyright (c) 1994, The Regents of the University of California

Permission to use, copy, modify, and distribute this software and its documentation for any purpose, without fee, and without a written agreement is hereby granted, provided that the above copyright notice and this paragraph and the following two paragraphs appear in all copies.

IN NO EVENT SHALL THE UNIVERSITY OF CALIFORNIA BE LIABLE TO ANY PARTY FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, OR CONSEQUENTIAL OWANGES, INCLUDING LOST PROFITS, ARISING OUT OF THE USE OF THIS SOFTWARE AND ITS DOCUMENTATION, EVEN IF THE UNIVERSITY OF CALIFORNIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAWAGE.

THE UNIVERSITY OF CALIFORNIA SPECIFICALLY DISCLAIRS ANY WAGRANTIES,
INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY
AND FITNESS FOR A PARTICULAR PURPOSE. THE SOFTWARE PROVIDED HEREUNDER IS
ON AN "AS IS" BASIS, AND THE UNIVERSITY OF CALIFORNIA HAS NO OBLIGATIONS TO
PROVIDE MAINTENANCE, SUPPORT, UPDATES, ENHANCEMENTS, OR MODIFICATIONS.

License of pgvector, 2025

```
buf = ReadBufferExtended(scan->indexRelation, MAIN FORKNUM, searchPage, RBM NORMAL, so->bas);
LockBuffer(buf, BUFFER LOCK SHARE);
page = BufferGetPage(buf);
maxoffno = PageGetMaxOffsetNumber(page);
for (OffsetNumber offno = FirstOffsetNumber: offno <= maxoffno: offno = OffsetNumberNext(offno))</pre>
    IndexTuple itup;
    Datum
                datum:
    bool
                isnull:
                itemid = PageGetItemId(page, offno);
    ItemId
                                                         All points in cluster(s) are read
                                                         again from buffers on each call!
   itup = (IndexTuple) PageGetItem(page. itemid);
    datum = index getattr(itup. 1. tupdesc. &isnull);
     * Add virtual tuple
     * Use procinfo from the index instead of scan key for
     * performance
     */
   ExecClearTuple(slot);
   slot->tts values[0] = so->distfunc(so->procinfo, so->collation, datum, value);
   slot->tts_isnull[0] = false;
   slot->tts_values[1] = PointerGetDatum(&itup->t_tid);
   slot->tts_isnull[1] = false;
   ExecStoreVirtualTuple(slot);
    tuplesort puttupleslot(so->sortstate, slot);
```



... HACKING PGVECTOR ...



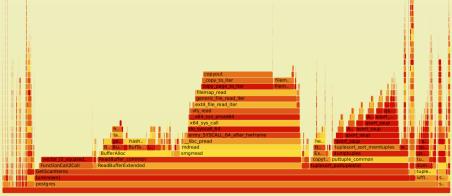
### pgvector hack (ivfflat)

Can we do better? ( https://github.com/Sednai/pgvector/tree/AERO v2 )

Keep index data persistent in Non-Postgres controlled memory:

- No Buffers but continuous arrays
- No TupleSort
- Possibility for better optimization for hardware







### pgvector hack (ivfflat)

Can we do better ? ( <a href="https://github.com/Sednai/pgvector/tree/AERO\_v2">https://github.com/Sednai/pgvector/tree/AERO\_v2</a> )

Keep index data persistent in Non-Postgres controlled memory:

- No Buffers but continuous arrays
- No TupleSort
- Possibility for better optimization for hardware

"In-memory vector database" (instead of in-memory buffer db)

Proof-of-concept for a next generation pgvector



### pgvector hack (ivfflat)

### *Implementation*

Own Postgres background process with task queue in shared memory.

```
BackgroundWorker worker:
BackgroundWorkerHandle *handle;
BgwHandleStatus status:
pid t
            pid;
memset(&worker, 0, sizeof(worker));
worker.bgw_flags = BGWORKER_SHMEM_ACCESS | BGWORKER_BACKEND_DATABASE_CONNECTION;
worker.bgw start time = BgWorkerStart RecoveryFinished;
worker.bgw restart time = BGW NEVER RESTART; // Time in s to restart if crash. Use BGW NEVER RESTART for no restart;
char* WORKER LIB = GetConfigOption("ivfflat.lib",true,true);
sprintf(worker.bgw_library_name, WORKER_LIB);
sprintf(worker.bgw function name, "pgv gpuworker main");
snprintf(worker.bgw_name, BGW_MAXLEN, "%s",buf);
worker.bgw notify pid = MyProcPid;
if (!RegisterDynamicBackgroundWorker(&worker, &handle))
    elog(ERROR, "Could not register background worker");
status = WaitForBackgroundWorkerStartup(handle, &pid);
```

Added bonus: Resources can be managed over all user sessions (via queuing system).



### pgvector hack (ivfflat)

### *Implementation*

Own Postgres background process with task queue in shared memory.

Re-routing the index tuple scan as task to background worker during scan:

```
* Fetch the next tuple in the given scan
bool
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                                Exec task
                            dlist node* dnode = dlist pop head node(&worker head->exec list);
         own process
                            worker exec entry* entry = dlist container(worker exec entry, node, dnode);
                            SpinLockRelease(&worker_head->lock);
                            load index(entry->nodeid, entry->tupdesc, entry->usetriangle);
                            // Compute
                            if(!entry->usegpu) {
                                entry->returns = exec_query_cpu(entry, worker_head);
                            else
```



#### pgvector hack (ivfflat)

### *Implementation*

For an indexed table, we store the index vectors and corresponding location info (*ItemPointerData*) as raw native arrays in Non-Postgres memory.

(The data will be persistent over the lifetime of the background process. Have not implemented active memory management yet. Pgvector background process will crash if you run out of memory!)

```
* Fetch the next tuple in the given scan
bool
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                                Exec task
                            dlist_node* dnode = dlist_pop_head_node(&worker_head->exec_list);
         own process
                            worker exec entry* entry = dlist container(worker exec entry, node, dnode);
                            SpinLockRelease(&worker_head->lock);
                            load index(entry->nodeid, entry->tupdesc, entry->usetriangle);
                            // Compute
                            if(!entry->useqpu) {
                                entry->returns = exec_query_cpu(entry, worker_head);
                            else
```



#### pgvector hack (ivfflat)

### *Implementation*

For a query point (vector), we compute its *distance* against all index vectors and distance *sort*. Location infos are returned to the user process.

(More precisely, the corresponding page number and ItemPointerData are returned.)

```
* Fetch the next tuple in the given scan
bool
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                                Exec task
                            dlist node* dnode = dlist pop head node(&worker head->exec list);
         own process
                            worker exec entry* entry = dlist container(worker exec entry, node, dnode);
                            SpinLockRelease(&worker_head->lock);
                            load index(entry->nodeid, entry->tupdesc, entry->usetriangle);
                            // Compute
                            if(!entry->usegpu) {
                                entry->returns = exec_query_cpu(entry, worker_head);
                            else
```



### pgvector hack (ivfflat)

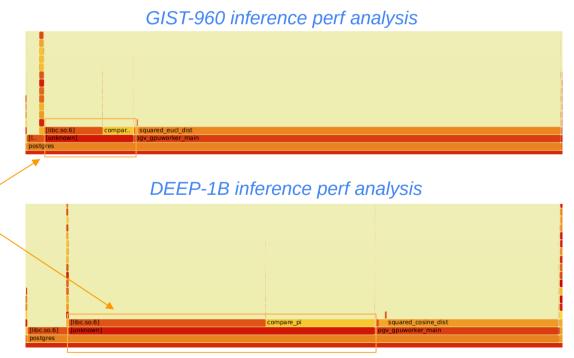
What about performance?

#### Note:

- No special tricks
- No multi-threading
- No manual vector instructions

std::qsort

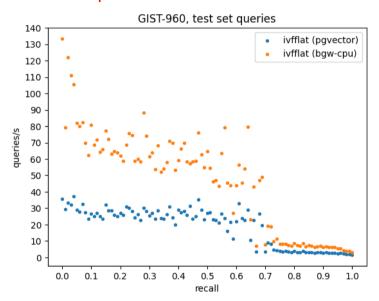
- No HBM memory



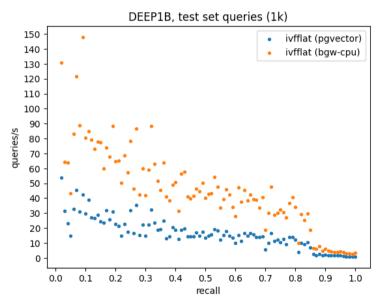


#### pgvector

#### What about performance?



select id from gist order by embedding <-> "+q+" limit 100

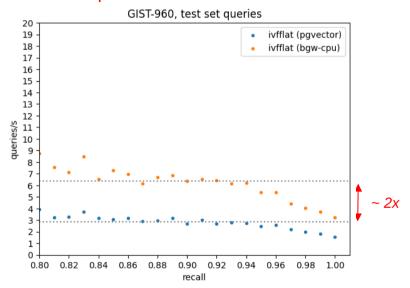


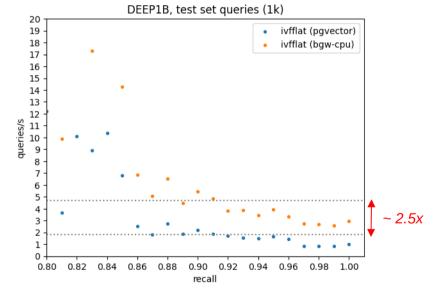
select id from deep1b order by embedding <=> "+q+" limit 100



#### pgvector

#### What about performance?





select id from gist order by embedding <-> "+q+" limit 100

select id from deep1b order by embedding <=> "+q+" limit 100



#### pgvector hack (ivfflat)

#### Can we do better?

Since we have already a setup to keep index persistent in Non-Postgres memory, we can easily go one step further and offload the index and compute to a GPU!

Compute of distances and sort on device. Return only sorted index ids from device (Mapping to location info on CPU)



#### pgvector hack (ivfflat)

#### Can we do better?

Since we have already a setup to keep index persistent in Non-Postgres memory, we can easily go one step further and offload the index and compute to a GPU!

Compute of distances and sort on device. Return only sorted index ids from device (Mapping to location info on CPU)

#### Example: Nvidia A100

80 GB in additional memory! (for FP32 vectors of 100d that is enough to keep >200M index points persistent)

~20 TFLOPS in FP32 compute power!



... GPU ACCELERATED VECTOR SEARCH ...



#### pgvector hack (ivfflat)

### *Implementation*

For an indexed table, we store now the index vectors on a GPU device.

(The data will be persistent over the lifetime of the background process. Have not implemented active memory management yet. Background process will crash if you run out of memory!)

```
* Fetch the next tuple in the given scan
bool
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                              Exec task
                          dlist node* dnode = dlist pop head node(&worker head->exec list);
        own process
                          worker exec entry* entry = dlist container(worker exec entry, node, dnode);
                          SpinLockRelease(&worker_head->lock);
                          load index(entry->nodeid, entry->tupdesc, entry->usetriangle);
                                                                                                      C++ / CUDA
                                                                                                      (essentially one big cudaMemcpy)
                          // Compute
                           if(!entry->useqpu) {
                              entry->returns = exec_query_cpu(entry, worker_head);
                           else
```



#### pgvector hack (ivfflat)

### *Implementation*

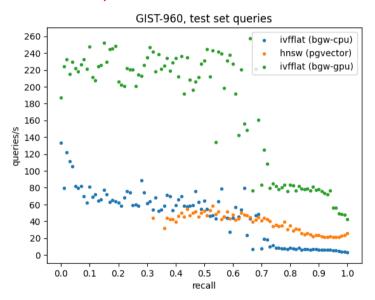
For a query point (vector), we compute its *distance* against all index vectors and distance *sort* on device. Only ordered index ids are returned from GPU.

```
* Fetch the next tuple in the given scan
bool
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                                load_index(entry->nodeid, entry->tupdesc, entry->usetriangle);
        own process
                                // Compute
                                if(!entry->usegpu) {
                                    entry->returns = exec_query_cpu(entry, worker_head);
                                else
                        #ifdef GPU
                                    entry->returns = exec_query_gpu(entry, worker_head);
                                                                                                 C++ / CUDA
                         #else
                                                                                                 (custom kernels)
                                    entry->returns = exec_query_cpu(entry, worker_head);
                         #endif
```

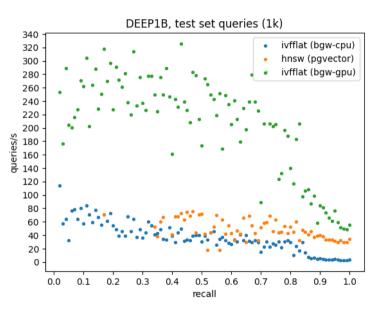


### pgvector hack (ivfflat)

#### What about performance?



select id from gist order by embedding <-> "+q+" limit 100

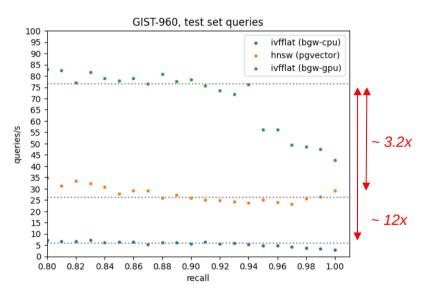


select id from deep1b order by embedding <=> "+q+" limit 100

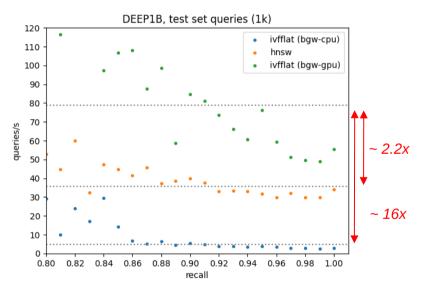


#### pgvector hack (ivfflat)

#### What about performance?



select id from gist order by embedding <-> "+q+" limit 100



select id from deep1b order by embedding <=> "+q+" limit 100



#### pgvector hack

Since we have now a basic vector search PG background worker:

## Hack in Nvidia RAFT and cuVS library?

( https://github.com/rapidsai/raft ) ( https://github.com/rapidsai/cuvs )

- Fullstack (HNSW, IFFLAT, CAGRA, DiskANN; Index build + inference)
- As GPU backend available for FAISS, Milvus, Redis and others ...
   ( according to <a href="https://big-ann-benchmarks.com/neurips23\_slides/NVIDIA\_Corey.pdf">https://big-ann-benchmarks.com/neurips23\_slides/NVIDIA\_Corey.pdf</a>)



#### pgvector hack

Since we have now a basic vector search PG background worker:

# Hack in Nvidia RAFT and cuVS library?

( https://github.com/rapidsai/raft ) ( https://github.com/rapidsai/cuvs )

- Fullstack (HNSW, IFFLAT, CAGRA, DiskANN; Index build + inference)
- As GPU backend available for FAISS, Milvus, Redis and others ...
   ( according to <a href="https://big-ann-benchmarks.com/neurips23\_slides/NVIDIA\_Corey.pdf">https://big-ann-benchmarks.com/neurips23\_slides/NVIDIA\_Corey.pdf</a>)

May be a low effort way to get something new and complete into PG ...

**BUT:** Do we really want to become vendor dependent?

### oneAPI





Aero aims to complement the efforts of the *EU Processor Initiative (EPI)* project by developing the open-source software ecosystem required to not only improve the efficiency of the EPI hardware but also accelerate and ease the processor's integration into the cloud.

Heterogeneous hardware

PROBLEM: Divers set of accelerators from different vendors (NVIDIA, AMD, INTEL,...)























#### oneAPI





Aero aims to complement the efforts of the *EU Processor Initiative (EPI)* project by developing the open-source software ecosystem required to not only improve the efficiency of the EPI hardware but also accelerate and ease the processor's integration into the cloud.

### Heterogeneous hardware



Open, cross-industry, standards-based, unified, multi-architecture, multi-vendor programming model, adopted by Intel.





























Aero aims to complement the efforts of the *EU Processor Initiative (EPI)* project by developing the open-source software ecosystem required to not only improve the efficiency of the EPI hardware but also accelerate and ease the processor's integration into the cloud.

### Heterogeneous hardware



Open, cross-industry, standards-based, unified, multi-architecture, multi-vendor programming model, adopted by Intel.

Intel oneAPI base toolkit plugins for NVIDIA and AMD



























#### pqvector hack (ivfflat)

## *Implementation*

bool

For a guery point (vector), we compute its *distance* against all index vectors and distance *sort* on device. Only ordered index ids are returned from GPU to CPU process.

```
* Fetch the next tuple in the given scan
ivfflatgettuple(IndexScanDesc scan, ScanDirection dir)
                                load index(entry->nodeid, entry->tupdesc, entry->usetriangle);
        own process
                                // Compute
                                if(!entrv->usegpu) {
                                    entry->returns = exec_query_cpu(entry, worker_head);
                                else
                         #ifdef GPU
                                    entry->returns = exec_query_gpu(entry, worker_head);
                                                                                                 C++ / oneAPI
                         #else
                                                                                                 (custom kernels)
                                    entry->returns = exec_query_cpu(entry, worker_head);
                         #endif
```



#### pgvector hack (ivfflat)

## *Implementation*

For a query point (vector), we compute its *distance* against all index vectors and distance *sort* on device. Only ordered index ids are returned from GPU to CPU process.



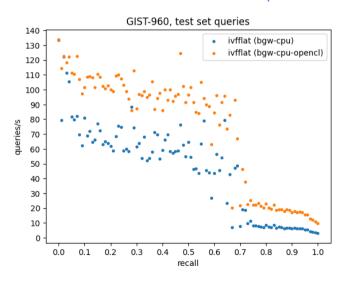
```
void calc squared euclidean distances(float* M, float* V, sort item* C, int* p, int N, int L, int probe) ∤
    Q->parallel for(range<1>(N),
    [=](id<1> k){}
        int pos = *p;
        float tmp = 0;
        for(int i = 0; i < L; i++) {</pre>
            tmp += (M[L*k+i] - V[i])*(M[L*k+i] - V[i]);
        C[pos+k].distance = tmp;
        C[pos+k].probe = probe;
        C[pos+k].pos = k;
    });
    0->wait():
```

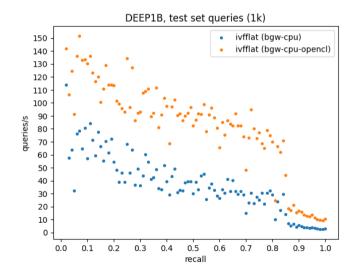


## pgvector hack (ivfflat)

#### What about performance?

#### With OpenCL CPU oneAPI backend (all cores)



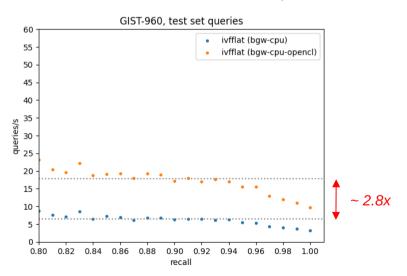


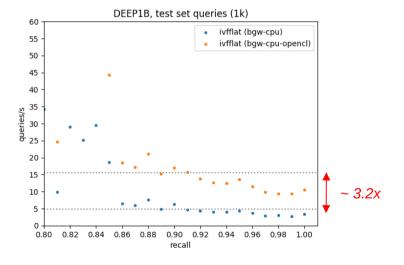


### pgvector hack (ivfflat)

#### What about performance?

#### With OpenCL CPU oneAPI backend (all cores)



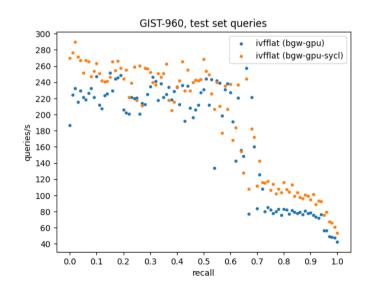


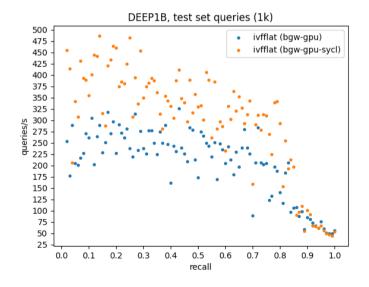


### pgvector hack (ivfflat)

#### What about performance?

#### With Nvidia GPU oneAPI backend





80

60

0.1 0.2

0.0

0.3

0.4 0.5

recall

0.6 0.7

0.8

0.9



Seems the automatic kernels are not had ....

### pgvector hack (ivfflat)

What about performance? Artefact of using unified memory? With Nvidia GPU oneAPI backend DEEP1B, test set queries (1k) GIST-960, test set queries 300 ivfflat (bgw-gpu) ivfflat (bgw-gpu) 280 ivfflat (bgw-gpu-sycl) ivfflat (bgw-gpu-sycl) 450 260 425 240 400 375 220 350 200 325 s/s 275 275 250 225 dneries/s 160 140 200 120 175 150 100

~ 1.25x

125

100

75 50 25

0.0 0.1

0.2

0.3 0.4

0.5

recall

0.6 0.7

0.8

0.9



... A SPECIAL CASE ...



## pgvector

What we want:

Fast way to retrieve (most) points up to a max distance from a query point.



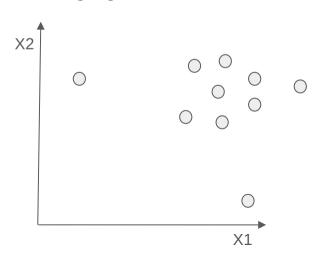
### pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





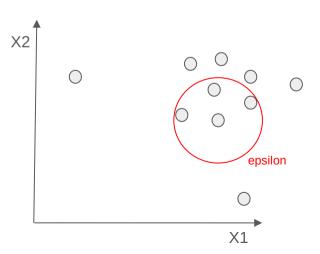
### pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

## Why?

Core ingredient to density based clustering algorithms.





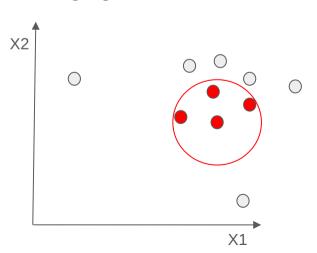
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





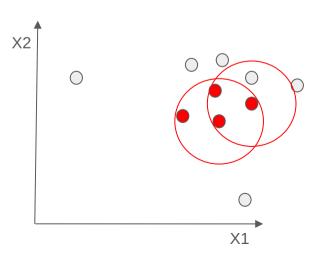
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





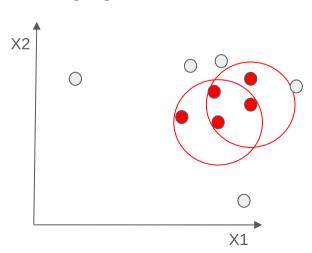
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





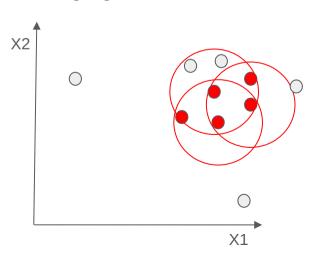
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

## Why?

Core ingredient to density based clustering algorithms.





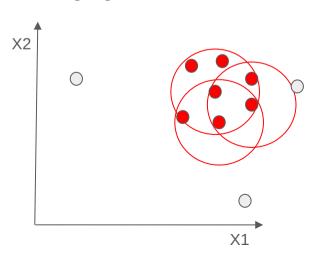
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

## Why?

Core ingredient to density based clustering algorithms.





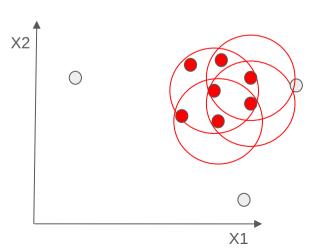
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





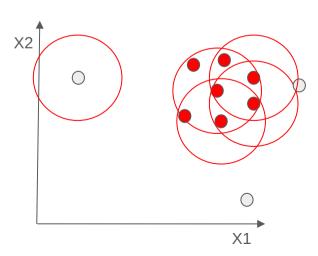
## pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

# Why?

Core ingredient to density based clustering algorithms.





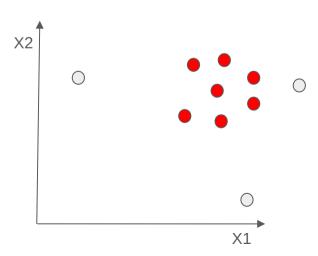
### pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

## Why?

Core ingredient to density based clustering algorithms.





#### pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

Get rows within a certain distance

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

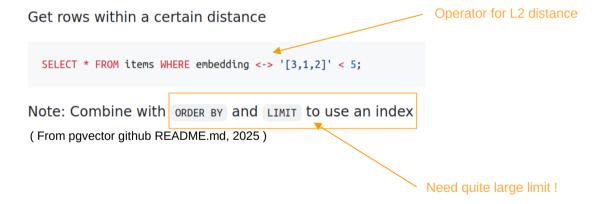
Note: Combine with ORDER BY and LIMIT to use an index (From pgyector github README.md, 2025)



# pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.





# pgvector

#### What we want:

Fast way to retrieve (most) points up to a max distance from a query point.

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 (From pgvector github README.md, 2025)

Need quite large limit!
```

PROBLEM: Very very slow ...



### pgvector

### What about performance?

- 2M row Gaia dataset, 79 float8 features
- 40 ivfflat clusters, ivfflat.probes = 20
- Retrieve points up to a max distance from a query point (sparse return)
- After warmup

## Original pgvector

```
Limit (cost=0.00..9416.48 rows=10000 width=16) (actual time=113.688..563.274 rows=2 loops=1)
                        -> Index Scan using lorenzo attrs idx on lorenzo (cost=0.00..569379.89 rows=604663 width=16) (actual time=113.685..563.269 rows=2 loops=1)
                                                                      Order By: (attrs <-> '\[ -0.36719987, 0.8524608, \cdot 0.5427666, \cdot 1.6615063, 1.1010165, 0.06038815, \cdot -0.8972259, \cdot -1.1181297, 0.23769811, 1.6662519, \cdot -1.58769811, 
473. - 0.042955805. 0.17839357. 0.08050123. 0.27791676. - 0.425645. 0.11280374. 0.84778684. 0.08167486. 1.8496727. 0.8007245. 0.5793525. - 0.5038844. 0.22990316. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486. 0.08167486
 0.12975255, -1.9553635, 0.9473572, -2.2433414, -0.30360684, 0.33857238, -0.21312521, 2.3237233, -0.060708717, 0.339421, -2.0196183, -0.35616732, 1.8636712, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616732, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.363616742, 0.366616742, 0.366616742, 0.366616742, 0.366616742, 0.366616742, 0.366616742, 0.366616742, 0.3
749731,-1.5635711,-1.2368783,-0.96010184,-0.6722383,0.8274576,-0.7714504,-0.16363333,-0.96023947,-0.16326201,-1.0754527,-0.6974341,-2.3611114,-1
03672114, 1.2685906, -0.22232309, -0.17129573, -0.30436236, -1.1221358, 0.6857615, -0.60302067, 0.22385728, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.6519747, -0.9824103, -1.5251932, -1.0727525, -1.671974, -1.0727525, -1.671974, -1.0727525, -1.671974, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727525, -1.0727
1835048,-2.293227,1.9016955,-2.8030064,-0.045054823,-0.14567287]'::vector)
                                                                        Filter: ((attrs <-> '[-0.36719987,0.8524608,-0.5427666,-1.6615063,1.1010165,0.06038815,-0.8972259,-1.1181297,0.23769811,1.6662519,-1.51
73, -0.042955805, 0.17839357, 0.08050123, 0.27791676, -0.425645, 0.11280374, 0.84778684, 0.08167486, 1.8496727, 0.8007245, 0.5793525, -0.5038844, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.22990316, 0.
0.12975255, -1.9553635, 0.9473572, -2.2433414, -0.30360684, 0.33857238, -0.21312521, 2.3237233, -0.060708717, 0.339421, -2.0196183, -0.35616732, 1.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.8636712, 7.86367
49731.-1.5635711.-1.2368783.-0.96010184.-0.6722383.0.8274576.-0.7714504.-0.16363333.-0.96023947.-0.16326201.-1.0754527.-0.6974341.-2.3611114.-1.
3672114.1.2685906. -0.22232309. -0.17129573. -0.30436236. -1.1221358.0.6857615. -0.60302067.0.22385728. -1.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -1.5251932. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -1.6519747. -0.9824103. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.0727525. -7.07275
835048,-2.293227,1.9016955,-2.8030064,-0.045054823,-0.14567287]'::vector) < '6'::double precision)
                                                                        Rows Removed by Filter: 578577
       Planning Time: 0.204 ms
       Execution Time: 563.326 ms
   (7 rows)
pgv2=#
```



### pgvector

# Why?

Let us look with GDB what actually happens for a query of type

select \* from table where embedding <-> '[...]' < 10 order by embedding <-> '[...]' limit 10000

by setting a breakpoint at ivfscan.c:ivfflatgettuple:

```
(gdb) break ivfscan.c:353

Breakpoint 1 at 0x7c106d843e7b: file src/ivfscan.c, line 358.

(gdb) c

Continuing.

Breakpoint 1, ivfflatgettuple (scan=0x60700e9406f8, dir=ForwardScanDirection) at src/ivfscan.c:360

360

if (so->first)

(gdb) p* scan

$1 = {heapRelation = 0x7c106d7925e8, indexRelation = 0x7c106d796818, xs_snapshot = 0x60700e8a8a68, numberOfKeys = 0, numberOfOrderBys = 1, keyData = 0x0, orderByData = 0x60700e940808, xs_want_itup = false, xs_temp_snap = false, kill_prior_tuple = false, ignore_killed_tuples = true, xactStartedInRecovery = false, opaque = 0x60700e940898, xs_itup = 0x0, xs_hitup = 0x0, xs_hitupdesc = 0x0, xs_heap_continue = false, xs_heapfetch = 0x60700e9409a8, xs_recheck = false, xs_orderbyvals = 0x0, xs_orderbynulls = 0x0, xs_recheckorderby = false, parallel_scan = 0x0}

(gdb) 

(gdb) 

(gdb) 

(gdb) break ivfscan.c:353

Breakpoint 1 at 0x7c106d843e7b: file src/ivfscan.c:360

ivf (so->first)

(so->first)

(so->first)

(gdb) p* scan

$1 = {heapRelation = 0x7c106d7925e8, indexRelation = 0x7c106d796818, xs_snapshot = 0x60700e8a8a68, numberOfKeys = 0, numberOfOrderBys = 1, keyData = 0x0, orderByData = 0x60700e940898, xs_want_itup = false, xs_temp_snap = false, kill_prior_tuple = false, ignore_killed_tuples = true, xactStartedInRecovery = false, opaque = 0x60700e940898, xs_itup = 0x0, xs_ituplesc = 0x0, xs_orderbynulls = 0x0, xs_orderbynulls = 0x0, xs_recheckorderby = false, xs_orderbynulls = 0x0, xs_orderbynulls = 0x0, xs_recheckorderby = false, ys_orderbynulls = 0x0, xs_orderbynulls = 0x0, x
```



### pgvector

# Why?

Let us look with GDB what actually happens for a query of type

select \* from table where embedding <-> '[...]' < 10 order by embedding <-> '[...]' limit 10000

by setting a breakpoint at ivfscan.c:ivfflatgettuple:

```
(gdb) break ivfscan.c:353
Breakpoint 1 at 0x7c106d843e7b: file src/ivfscan.c, line 358.

(gdb) c
Continuing.

Breakpoint 1, ivfflatgettuple (scan=0x60700e9406f8, dir=ForwardScanDirection) at src/ivfscan.c:360

(gdb) p* scan
$1 = {heapRelation = 0x7c106d7925e8, indexRelation = 0x7c106d796818, xs_snapshot = 0x60700e8a8a68, numberOfKeys = 0, numberOfOrderBys = 1, keyData = 0x0, orderByData = 0x60700e940808, xs_want_itup = false, xs_temp_snap = false, kill_prior_tuple = false, ignore_killed_tuples = true, xactStartedInRecovery = false, opaque = 0x60700e940898, xs_itup = 0x0, xs_hitup = 0x0, xs_hitupdesc = 0x0, xs_heaptid = {ip_blkid = {bi_hi = 0, bi_lo = 8680}, ip_posid = 4}, xs_heap_continue = false, xs_heapfetch = 0x60700e9409a8, xs_recheck = false, xs_orderbyvals = 0x0, xs_orderbynulls = 0x0, xs_recheckorderby = false, parallel_scan = 0x0}

(gdb) I
```



### pgvector

Why no scan keys?

We have to dig deeper ...

ExecInitBuildScanKeys: quals are Null

iss NumScanKevs = 0 already in IndexScanState

```
(qdb) bt
#0 ivfflatgettuple (scan=0x60700e9406f8, dir=ForwardScanDirection) at src/ivfscan.c:360
#1 0x000060700c44b319 in index getnext tid (scan=0x60700e9406f8, direction=ForwardScanDirection) at indexam.c:575
#2 0x000060700c44b529 in index getnext slot (scan=0x60700e9406f8 direction=ForwardScanDirection, slot=0x60700e9586d0) at indexam.c:667
#3 0x000060700c6b234d in IndexNextWithReorder (node=0x60700e958430) at nodeIndexscan.c:264
#4 0x000060700c68ad18 in ExecScanFetch (node=0x60700e958430, accessMtd=0x60700c6b2158 <IndexNextWithReorder>, recheckMtd=0x60700c6b2675 <IndexRecheck>)
    at execScan.c:132
#5 0x000060700c68adbd in ExecScan (node=0x60700e958430. accessMtd=0x60700c6b2158 <IndexNextWithReorder>. recheckMtd=0x60700c6b2675 <IndexRecheck>) at execScan.c:198
#6 0x000060700c6b2b73 in ExecIndexScan (pstate=0x60700e958430) at nodeIndexscan.c:533
#7 0x000060700c686cce in ExecProcNodeFirst (node=0x60700e958430) at execProcnode.c:464
#8 0x000060700c6b5207 in ExecProcNode (node=0x60700e958430) at ../../src/include/executor/executor.h:262
#9 0x000060700c6b53f7 in Exectimit (pstate=0x60700e958140) at nodeLimit.c:96
#10 0x000060700c686cce in ExecProcNodeFirst (node=0x60700e958140) at execProcnode.c:464
#11 0x000060700c67b20a in ExecProcNode (node=0x60700e958140) at ../../../src/include/executor/executor.h:262
#12 0x000060700c67dac5 in ExecutePlan (queryDesc=0x60700e965c68, operation=CMD SELECT, sendTuples=true, numberTuples=0, direction=ForwardScanDirection.
    dest=0x60700e9531e0) at execMain.c:1640
#13 0x000060700c67b71c in standard ExecutorRun (queryDesc=0x60700e965c68, direction=ForwardScanDirection, count=0, execute once=false) at execMain.c:362
#14 0x000060700c67b621 in ExecutorRun (queryDesc=0x60700e965c68, direction=ForwardScanDirection, count=0, execute_once=false) at execMain.c:311
#15 0x000060700c8b9b66 in PortalRunSelect (portal=0x60700e8f4328, forward=true, count=0, dest=0x60700e9531e0) at pguerv.c:922
#16 0x000060700c8b97d1 in PortalRun (portal=0x60700e8f4328, count=9223372036854775807, isTopLevel=true, run once=true, dest=0x60700e9531e0. altdest=0x60700e9531e0.
```

=> Already before execution level no scan keys!



### pgvector

### Why no scan keys?

We have to dig deeper ...

#### Let us look into indxpath.c:

```
2213
2214
        * match clause to indexcol()
2215
             Determine whether a restriction clause matches a column of an index.
             and if so, build an IndexClause node describing the details.
2216
2217
             To match an index normally, an operator clause:
2218
2219
2220
             (1) must be in the form (indexkey op const) or (const op indexkey);
2221
                  and
             (2) must contain an operator which is in the index's operator family
2222
                  for this column; and
2223
             (3) must match the collation of the index, if collation is relevant.
2224
2225
```

```
*
* indxpath.c

* Routines to determine which indexes are usable for scanning a
* given relation, and create Paths accordingly.

*
* Portions Copyright (c) 1996-2022, PostgreSQL Global Development Group
* Portions Copyright (c) 1994, Regents of the University of California
*
* IDENTIFICATION
* src/backend/optimizer/path/indxpath.c
*
```



### pgvector

## Why no scan keys?

We have to dig deeper ... Query:

Let us look into indxpath.c:

select \* from table where embedding <-> '[...]' < 10 ...

```
2213
2214
        * match clause to indexcol()
             Determine whether a restriction clause matches a column of an index.
2215
             and if so, build an IndexClause node describing the details.
2216
2217
             To match an index normally, an operator clause:
2218
2219
2220
             (1) must be in the form (indexkey op const) or (const op indexkey);
2221
                  and
             (2) must contain an operator which is in the index's operator family
2222
                  for this column; and
2223
             (3) must match the collation of the index, if collation is relevant.
2224
2225
```

Indexkey can not be matched!

op(indexkey, '[...']) op const



# pgvector

Why no scan keys?

We have to dig deeper ...

Query:

Let us look into indxpath.c:

select \* from table where embedding <-> '[...]' < 10 ...

=> Looks like Postgres enhancement required!

**BUT**: May take ages to get upstream ...

Indexkey can not be matched!

op(indexkey, '[...']) op const



# pgvector hack

# Quicker to production:

Let us introduce a new operator, thereby hacking the *WHERE* into the *ORDER BY*:

New operator for WHERE clause directly in index scan (only for euclidean metric so far):

vector <!> vector\_adv

with

vector\_adv = (vector,int,float,int)

int specifies the filter operator, float the condition value

2: >=
1: >
0: ==
-1: <
-2: <=
-100: no filter</pre>



# pgvector hack

### Quicker to production:

Let us introduce a new operator, thereby hacking the *WHERE* into the *ORDER BY*:

New operator for WHERE clause directly in index scan (only for euclidean metric so far):

vector <!> vector\_adv

Will come back to the last int later in this talk

vector\_adv = (vector,int,float,int)

int specifies the filter operator, float the condition value

2: >=
1: >
0: ==
-1: <
-2: <=
-100: no filter



# pgvector hack

Quicker to production:

Let us introduce a new operator, thereby hacking the *WHERE* into the *ORDER BY*:

Query:

select \* from table order by embedding <!> ('[...]', -1, 10.0, 0) limit 10000



# pgvector hack

# Quicker to production:

Let us introduce a new operator, thereby hacking the *WHERE* into the *ORDER BY*:

# Query:



Will be evaluated inside of pgvector!

-> Can also be executed on GPU!



## pgvector hack

### What about performance?

- 2M row Gaia dataset, 79 float8 features
- 40 ivfflat clusters, ivfflat.probes = 20
- Retrieve points up to a max distance from a query point (sparse return)
- After warmup

#### BGW with filter on CPU

```
Limit (cost=0.00..3105.50 rows=10000 width=16) (actual time=18.928..18.932 rows=1 loops=1)

-> Index Scan using lorenzo_attrs_idx on lorenzo (cost=0.00..563333.26 rows=1813988 width=16) (actual time=18.925..18.927 rows=1 loops=1)

Order By: (attrs <!> '("[-0.36719987,0.8524608,-0.5427666,-1.6615063,1.1010165,0.06038815,-0.8972259,-1.1181297,0.23769811,1.6662519,-1
18473,-0.042955805,0.17839357,0.08050123,0.27791676,-0.425645,0.11280374,0.84778684,0.08167486,1.8496727,0.8007245,0.5793525,-0.5038844,0.229905
75,0.12975255,-1.9553635,0.9473572,-2.2433414,-0.30360684,0.33857238,-0.21312521,2.3237233,-0.060708717,0.339421,-2.0196183,-0.35616732,1.863671
.0749731,-1.5635711,-1.2368783,-0.96010184,-0.6722383,0.8274576,-0.7714504,-0.16363333,-0.96023947,-0.16326201,-1.0754527,-0.6974341,-2.3611114,0.03672114,1.2685906,-0.22232309,-0.17129573,-0.30436236,-1.1221358,0.6857615,-0.60302067,0.22385728,-1.0727525,-1.6519747,-0.9824103,-1.5251932
1.1835048,-2.293227,1.9016955,-2.8030064,-0.045054823,-0.14567287]",-1,4)'::vector_adv)
Planning Time: 0.188 ms
Execution Time: 18.980 ms
(5 rows)
```



## pgvector hack

### What about performance?

- 2M row Gaia dataset, 79 float8 features
- 40 ivfflat clusters, ivfflat.probes = 20
- Retrieve points up to a max distance from a query point (sparse return)
- After warmup

#### BGW with filter on GPU

```
Limit (cost=0.00.3105.50 rows=10000 width=16) (actual time=2.482..2.488 rows=1 loops=1)

-> Index Scan using lorenzo_attrs_idx on lorenzo (cost=0.00.563333.26 rows=1813988 width=16) (actual time=2.478..2.482 rows=1 loops=1)

Order By: (attrs <!> '("[-0.36719987,0.8524608,-0.5427666,-1.6615063,1.1010165,0.06038815,-0.8972259,-1.1181297,0.23769811,1.6662519,-1
18473,-0.042955805,0.17839357,0.08050123,0.27791676,-0.425645,0.11280374,0.84778684,0.08167486,1.8496727,0.8007245,0.5793525,-0.5038844,0.229903
75,0.12975255,-1.9553635,0.9473572,-2.2433414,-0.30360684,0.33857238,-0.21312521,2.3237233,-0.060708717,0.339421,-2.0196183,-0.35616732,1.863671
.0749731,-1.5635711,-1.2368783,-0.96010184,-0.6722383,0.8274576,-0.7714504,-0.16363333,-0.96023947,-0.16326201,-1.0754527,-0.6974341,-2.3611114,
0.03672114,1.2685906,-0.22232309,-0.17129573,-0.30436236,-1.1221358,0.6857615,-0.60302067,0.22385728,-1.0727525,-1.6519747,-0.9824103,-1.5251932
1.1835048,-2.293227,1.9016955,-2.8030064,-0.045054823,-0.14567287]",-1,4)'::vector_adv)
Planning Time: 0.199 ms
Execution Time: 2.548 ms
(5 rows)
```



# pgvector hack

### What about performance?

- 2M row Gaia dataset, 79 float8 features
- 40 ivfflat clusters, ivfflat.probes = 20
- Retrieve points up to a max distance from a query point (sparse return)
- After warmup

#### BGW with filter on GPU

```
Limit (cost=0.00.3105.50 rows=10000 width=16) (actual time=2.482..2.488 rows=1 loops=1)

-> Index Scan using lorenzo_attrs_idx on lorenzo (cost=0.00.563333.26 rows=1813988 width=16) (actual time=2.478..2.482 rows=1 loops=1)

Order By: (attrs <!> '["[-0.36719987,0.8524608,-0.5427666,-1.6615063,1.1010165,0.06038815,-0.8972259,-1.1181297,0.23769811,1.6662519,-1
18473,-0.042955805,0.17839357,0.08050123,0.27791676,-0.425645,0.11280374,0.84778684,0.08167486,1.8496727,0.8007245,0.5793525,-0.5038844,0.229903
75,0.12975255,-1.9553635,0.9473572,-2.2433414,-0.30360684,0.33857238,-0.21312521,2.3237233,-0.060708717,0.339421,-2.0196183,-0.35616732,1.863671
.0749731,-1.5635711,-1.2368783,-0.96010184,-0.6722383,0.8274576,-0.7714504,-0.16363333,-0.96023947,-0.16326201,-1.0754527,-0.6974341,-2.3611114,
0.03672114,1.2685906,-0.22232309,-0.17129573,-0.30436236,-1.1221358,0.6857615,-0.60302067,0.22385728,-1.0727525,-1.6519747,-0.9824103,-1.5251932
1.1835048,-2.293227,1.9016955,-2.8030064,-0.045054823,-0.14567287]",-1,4)'::vector_adv)
Planning Time: 0.199 ms
Execution Time: 2.548 ms
(5 rows)
```

## => 200x speedup!

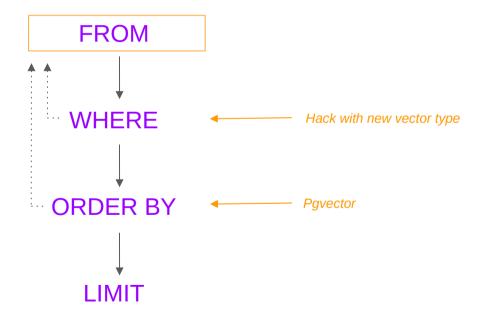


... LIMIT TRICKS ...



# pgvector hack

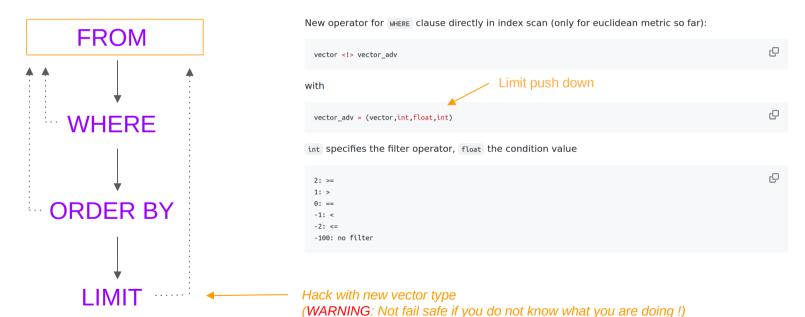
# PostgreSQL plan generation (simplified)





# pgvector hack

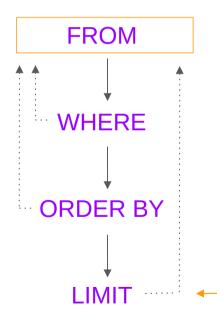
## PostgreSQL plan generation (simplified)





### pgvector hack





We can do instead of a full sort a partial sort (kth-element) + a small sort!

$$\rightarrow$$
 O(N) + O(k log k) scaling instead of O(N log N)

CPU: std::nth\_element

OneAPI: oneapi::dpl::nth\_element

Nvidia: Missing. Instead: raft::matrix::select\_k

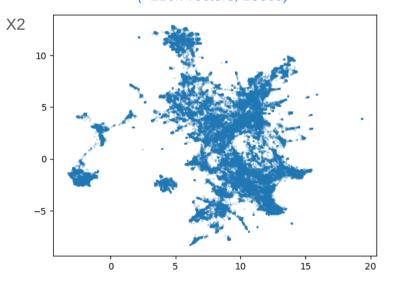
Hack with new vector type (WARNING: Not fail safe if you do not know what you are doing!)



### pgvector hack

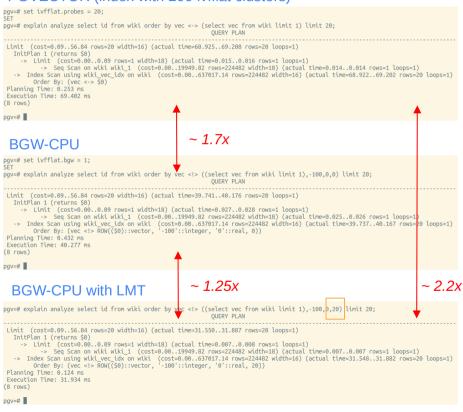
## Example:

Wikipedia - OpenAI embeddings (~220k vectors; 1536d)



X1

#### PGVECTOR (index with 200 ivfflat clusters)

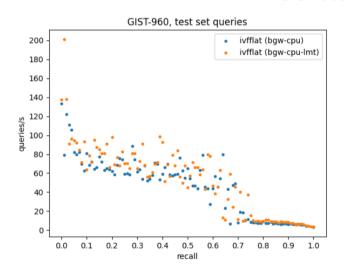


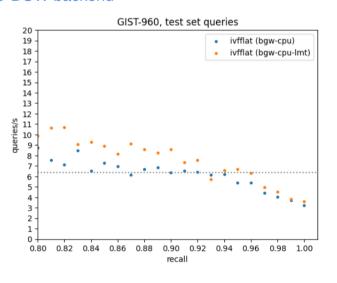


# pgvector hack

# Example:

#### GIST-960 with CPU-BGW backend







... REMARKS & OUTLOOK ...



# pgvector hack

#### **General remarks:**

- Proof-of-concept
- No active memory management (memory freed only upon killing the worker)
- Enough shared memory needs to be reserved for number of expected returns
- As more sparse the return, as better will be the speedup



## pgvector hack

#### Can we do more?

- Improvements of code (GPU kernels) are possible.
- Faster initial loading via Nvidia GPUDirect ( NVMe <-> GPU DMA )
- Product quantization
- Multi-threading on the worker level
- For significant performance improvement, more *vectorization* ... (for instance, to query for several points at once)
- Port other algorithms (HNSW, DiskANN, ...)
- Use Nvidia RAFT and cuVS instead ?!



#### ... THANK YOU ...







#### Project funded by



Federal Department of Economic Affairs, Education and Research EAER State Secretariat for Education, Research and Innovation SERI

Swiss Confederation

Funded by the European Union. Views and opinions are however those of the author(s) only and do not necessarily reflect those of the European Union or the HaDEA. Neither the European Union nor the granting authority can be held responsible for them. Project number: 101092850.

AERO has also received funding from UKRI under grants no. 10048318 and 10048915, and the Swiss State Secretariat for Education, Research, and Innovation.