# High-dimensional Approximate Nearest Neighbor Search
## Applications, Algorithms, Current Challenges

Martin Aumüller

IT University of Copenhagen

[maau@itu.dk](mailto:maau@itu.dk)

DAMVISEH Workshop, November 2025

🌐 http://itu.dk/people/maau

 @maumueller

**Slides will be available on my website.**

IT UNIVERSITY OF CPH

# Martin Aumüller

## Associate Professor, IT University of Copenhagen, Denmark
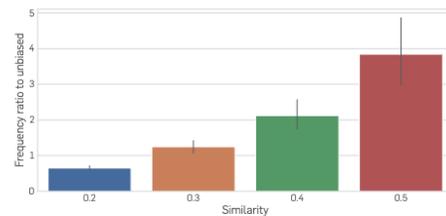
🌐 http://itu.dk/people/maau        @maumueller
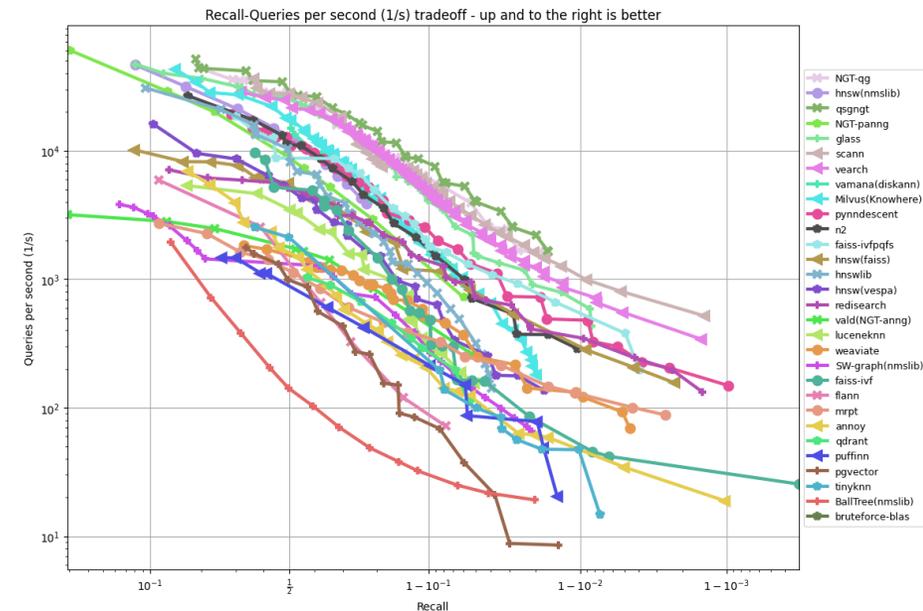
✓ Theory (LSH) and Practice of ANN

✓ Benchmarks & Challenges

**Figure 1. Bias introduced by uniform sampling from LSH buckets on the Last.FM dataset. The task is to (repeatedly) retrieve a uniform user among all users with similarity at least 0.2 to a fixed user. The result is split up into four buckets by rounding down the similarity to the first decimal. Error bars show the standard deviation. Compared to an unbiased sample, user vectors with small similarity are underrepresented, and users with high similarity are, by a factor of approximately 4 on average, overrepresented.**

LEMMA 9. *Given a set S of n points and a parameter r, we can preprocess it such that given a query* $\mathbf{q}$, *one can report a point* $\mathbf{p} \in S$ *with probability* $1/n(\mathbf{q}, r)$. *The algorithm uses space* $O(L \log n)$ *and has expected query time* $O\left(\left(L + \frac{n(\mathbf{q}, cr)}{n(\mathbf{q}, r)}\right) \log^5 n\right)$.

PUFFINN: k-NN with strong guarantees [A, Christiani, Pagh, Vesterli, ESA 2019]

Fair Near Neighbor Search [A, Har-Peled, Mahabadi, Pagh, Silvestri, Comm. ACM 2022]

Recall-Queries per second (1/s) tradeoff - up and to the right is better

ANN-Benchmarks [A, Bernhardsson, Faithfull, 2020]
Billion-Scale ANN Challenge
[Simhadri, A, et al., NeurIPS 21/23, Competition]

# Credits

**Yusuke Matsui**



https://yusukematsui.me/

Provided examples for RAG, PQ, BeamSearch.

**Matteo Ceccarello**



https://www.dei.unipd.it/~ceccarello/

Provided LSH visualization and workload generation slides.

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?
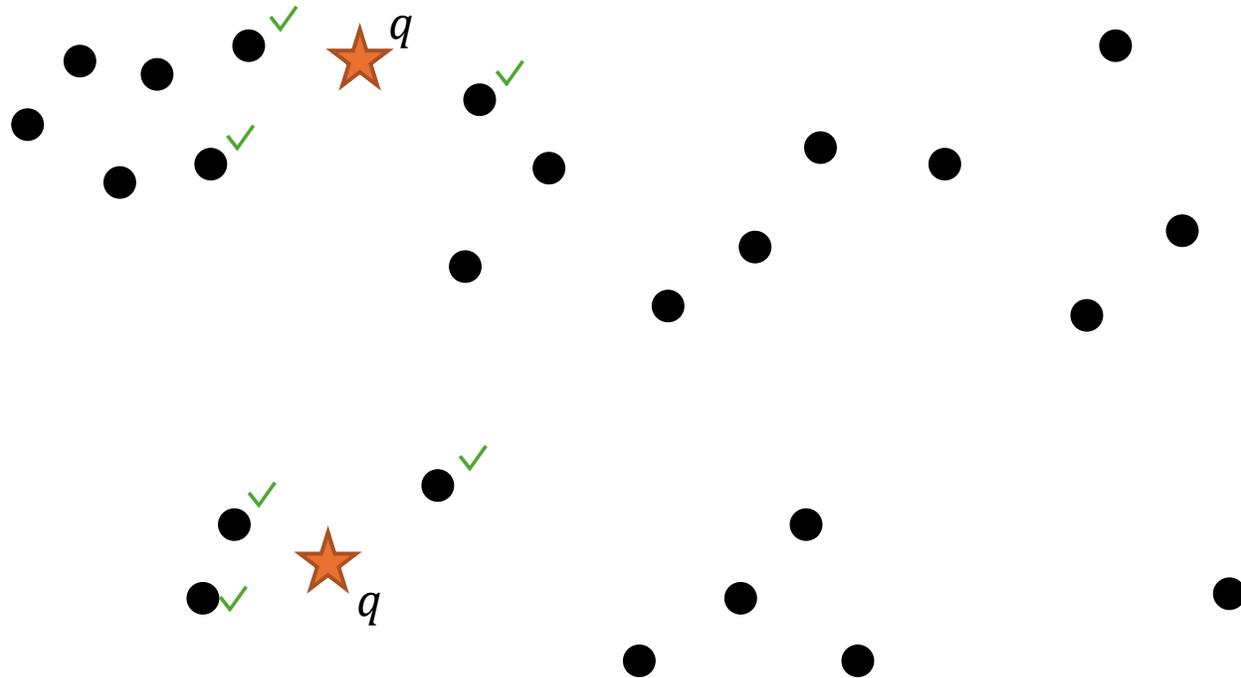
Part 4: How to use search to speed up data mining?

**Focus:**
- Explain the basics
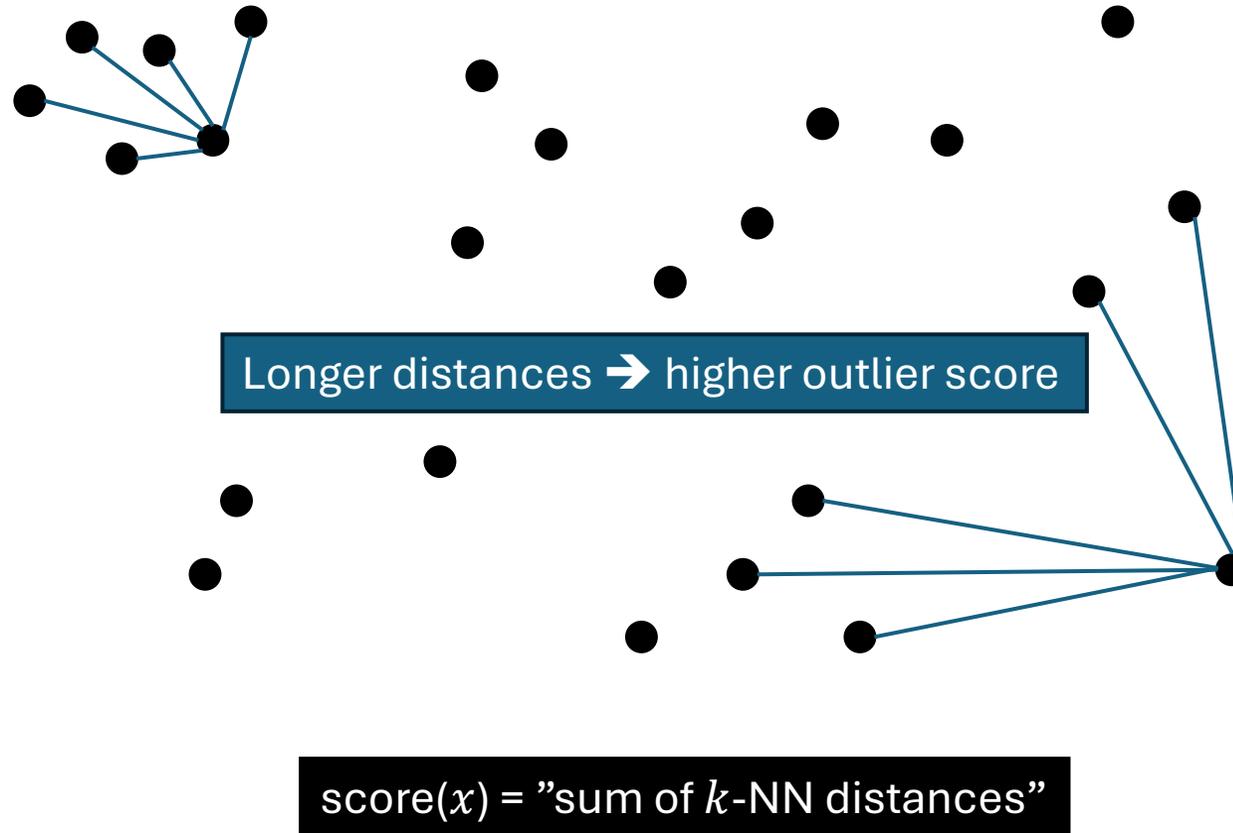- Point to interesting recent work

# Search: $k$-Nearest Neighbor Problem

- **Preprocessing**: Build DS for set $S \subseteq \mathbb{R}^d$ of $n$ data points

- **Task**: For query $q \in \mathbb{R}^d$ and $k \geq 1$, return $k$ closest points to $q$ in $S$
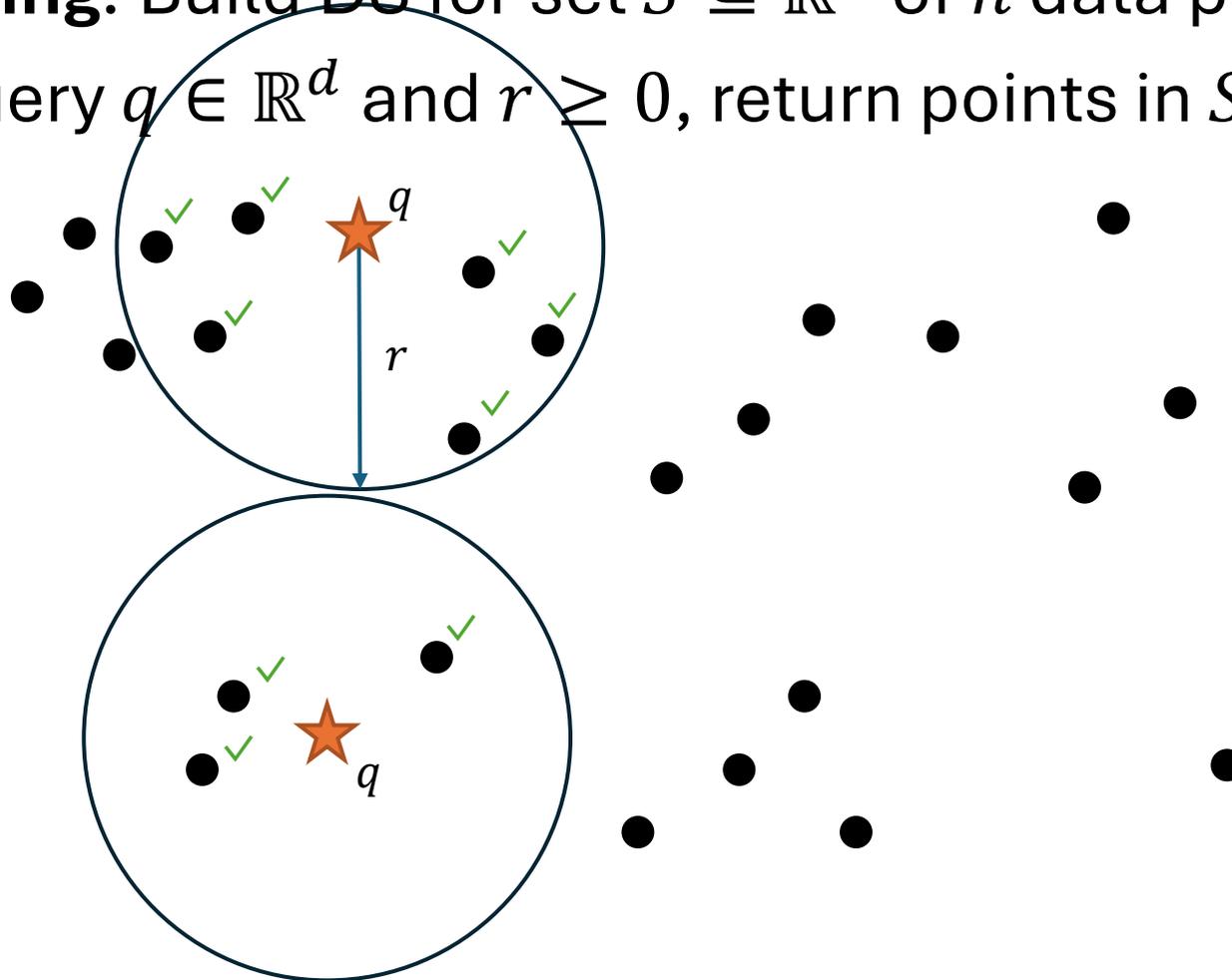
Naïve Solution: Linear Scan through $S$, Time $O(dn)$

# Application: Outlier detection ($k$-NN)

**Task**: Given $S \subseteq \mathbb{R}^d$, assign each point an outlier score.

Longer distances ➜ higher outlier score

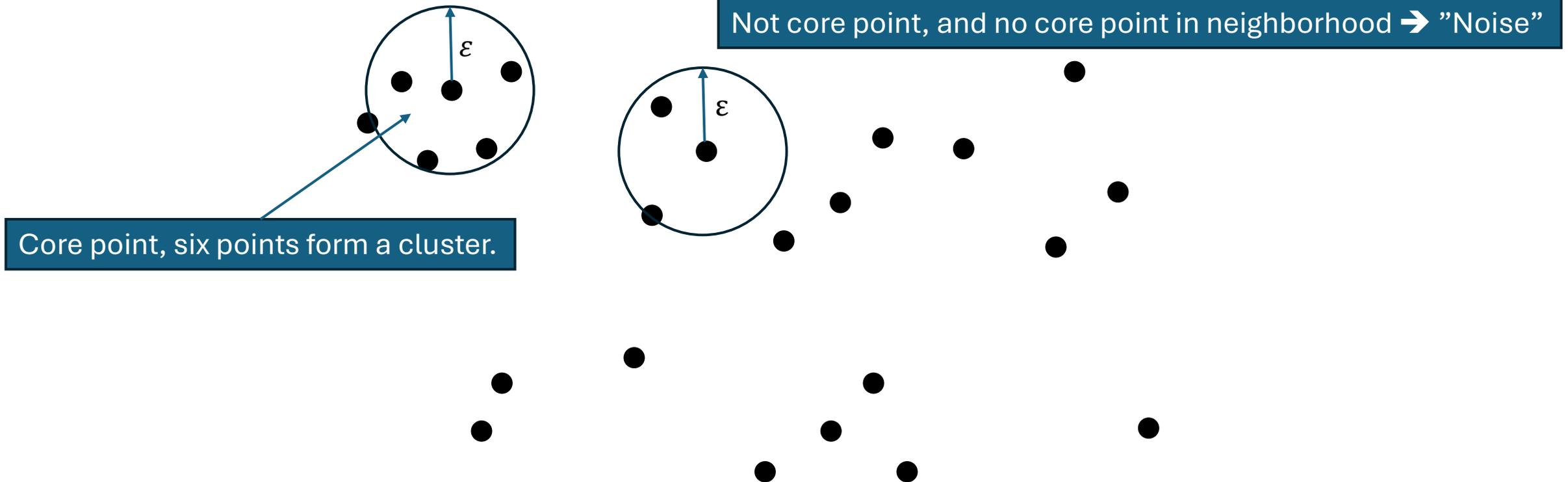score($x$) = "sum of $k$-NN distances"

# Search: Radius search

- **Preprocessing**: Build DS for set $S \subseteq \mathbb{R}^d$ of $n$ data points
- **Task**: For query $q \in \mathbb{R}^d$ and $r \geq 0$, return points in $S$ at dist $\leq r$

# Application: DBSCAN [Ester et al., 1996]

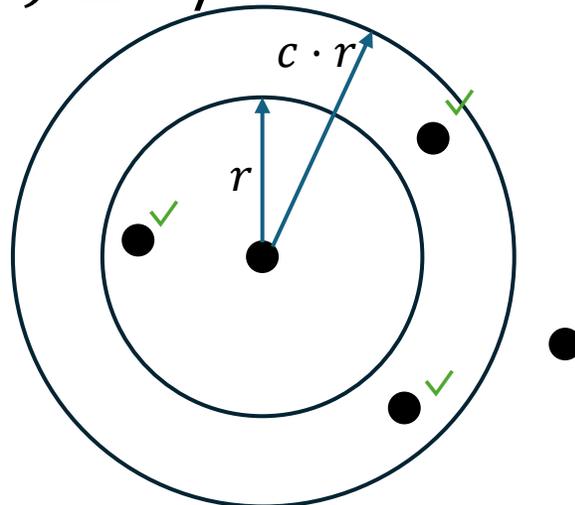**Input**: Dataset S, distance threshold $\varepsilon$, minPts (e.g., 3).

**Core point**: at least **minPts** neighbors at distance $\leq \varepsilon$.

Not core point, and no core point in neighborhood ➔ "Noise"

Core point, six points form a cluster.

# **Approximate** Near(est) Neighbor Search

**Theoretician's view**

- (Probabilistic) guarantee to return $c$-approximate near neighbor

- Guaranteed running time: $O\left(n^{\rho(c)}\right), \rho(c) \leq 1/c$

**Practitioner's view**

- *Approximate = "inexact"*
  - "no guarantees are given"
  - Quality has to be measured

- **Quality measures**
  - **Recall** = Fraction of true nearest neighbors identified
  - Often: **Average recall** over set of queries.

- **Running time** vs. **Quality**

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?

Part 4: How to use search to speed up data mining?

# Real-world use cases: LLM + embedding

"Who won curling
gold at the 2022
Winter Olympics?"

ChatGPT 3.5
(trained in 2021)

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"

Ask →



ChatGPT 3.5
(trained in 2021)

→

"I'm sorry, but as an AI language model, I don't have information about the future events."

☹

# Real-world use cases: LLM + embedding

"Who won curling
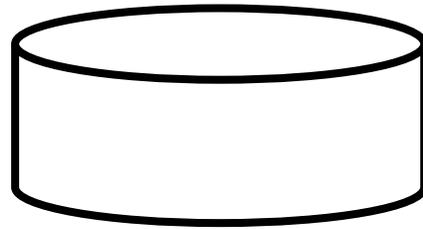gold at the 2022
Winter Olympics?"



ChatGPT 3.5
(trained in 2021)

Texts are from: https://github.com/openai/openaicookbook/blob/main/examples/Question_answering_using_embeddings.ipynb     Icon credit: https://ja.wikipedia.org/wiki/ChatGPT

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"

ChatGPT 3.5
(trained in 2021)

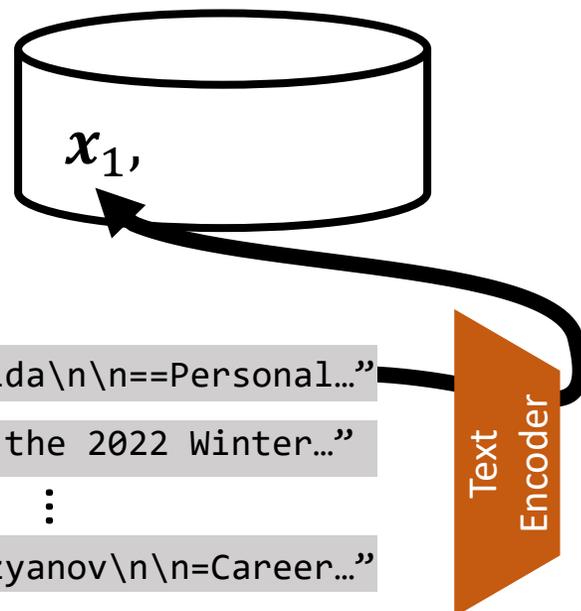"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

Texts are from: https://github.com/openai/openaicookbook/blob/main/examples/Question_answering_using_embeddings.ipynb     Icon credit: https://ja.wikipedia.org/wiki/ChatGPT

# Real-world use cases: LLM + embedding

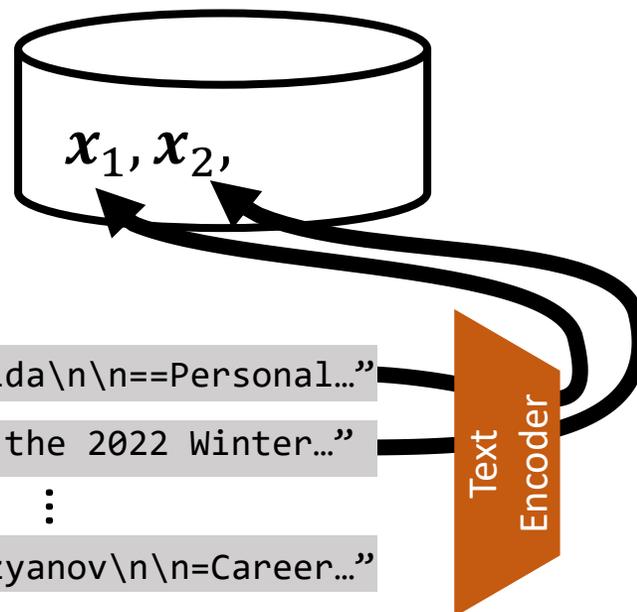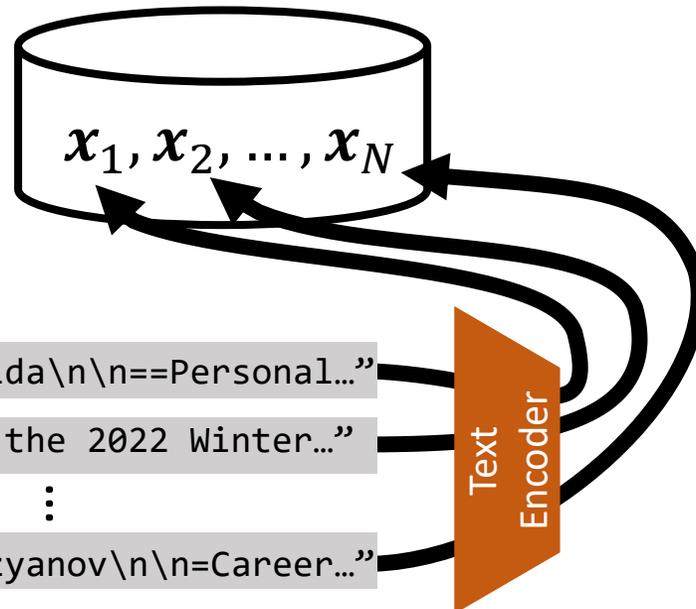"Who won curling gold at the 2022 Winter Olympics?"

ChatGPT 3.5
(trained in 2021)

$x_1,$

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

Text Encoder

15

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"



ChatGPT 3.5
(trained in 2021)

$$x_1, x_2,$$

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

Text Encoder

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"

ChatGPT 3.5
(trained in 2021)

$$x_1, x_2, \ldots, x_N$$

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

Text Encoder

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"

Text Encoder

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

ChatGPT 3.5
(trained in 2021)

$x_1, x_2, \ldots, x_N$

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

Text Encoder

# Real-world use cases: LLM + embedding



"Who won curling gold at the 2022 Winter Olympics?"

Text Encoder

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

Search

$x_1, x_2, \ldots, x_N$

ChatGPT 3.5
(trained in 2021)

$$\mathrm{argmin} \| q - x_n \|_2^2$$

Result
$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

"List of 2022 Winter Olympics medal winners…"

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

Text Encoder

"Damir Sharipzyanov\n\n=Career…"

Texts are from: https://github.com/openai/openaicookbook/blob/main/examples/Question_answering_using_embeddings.ipynb    Icon credit: https://ja.wikipedia.org/wiki/ChatGPT

# Real-world use cases: LLM + embedding



"Who won curling gold at the 2022 Winter Olympics?"

**Update** →

"Who won curling gold at the 2022 Winter Olympics? Use the below articles: List of 2022 Winter Olympics medal winners…"

ChatGPT 3.5
(trained in 2021)

**Text Encoder**

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

**Search** →

$$x_1, x_2, \ldots, x_N$$

$$\mathrm{argmin}\|q - x_n\|_2^2$$

Result

$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

"List of 2022 Winter Olympics medal winners…"

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

**Text Encoder**

Texts are from: https://github.com/openai/openaicookbook/blob/main/examples/Question_answering_using_embeddings.ipynb    Icon credit: https://ja.wikipedia.org/wiki/ChatGPT

20

# Real-world use cases: LLM + embedding



"Who won curling gold at the 2022 Winter Olympics?"

**Update**

"Who won curling gold at the 2022 Winter Olympics?
Use the below articles: List of 2022 Winter Olympics medal winners…"

ChatGPT 3.5
(trained in 2021)

"**Niklas Edin, Oskar Eriksson, …**"

**Text Encoder**

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

**Search**

$$x_1, x_2, \ldots, x_N$$

$$\mathrm{argmin} \| \boldsymbol{q} - \boldsymbol{x}_n \|_2^2$$

Result

$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

"List of 2022 Winter Olympics medal winners…"

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

⋮

"Damir Sharipzyanov\n\n=Career…"

**Text Encoder**

# Real-world use cases: LLM + embedding

"Who won curling gold at the 2022 Winter Olympics?"

**Update** →

"Who won curling gold at the 2022 Winter Olympics? Use the below articles: List of 2022 Winter Olympics medal winners…"

→ ChatGPT 3.5 (trained in 2021) →

"Niklas Edin, Oskar Eriksson, …"

Text Encoder

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

**Search** →

$x_1, x_2, \ldots, x_N$

$$\mathrm{argmin}\|\boldsymbol{q} - \boldsymbol{x}_n\|_2^2$$

Result

$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

"List of 2022 Winter Olympics medal winners…"

"Chinami Yoshida\n\n==Personal…"

"Lviv bid for the 2022 Winter…"

"Damir Sharipzyanov\n\n=Career…"

Text Encoder

**Embedding+ANN is the current easiest way to provide knowledge to LLM**

22

# Where does the data come from?

**Semantic embeddings**

- "Handcrafted vectors", e.g.,
  - SIFT descriptors for object detection
  - BM25/Bag of Words descriptors for text
  - Pixel-by-Pixel encoding of image

**Deep-Learning-based embeddings**

- Neural-network-based embeddings, often trained on classification tasks
- Vector = "activation values on *some* layer"
- Cross-Modal embeddings possible, e.g., CLIP

# Example of Handcrafted Embedding
## Fashion-MNIST



- 70,000 images
- Each image:
  - 28x28=784 dimensions
  - grayscale 0,...,255 value per coordinate

https://mlops.systems/posts/2022-10-24-foundations-mnist-basics.html

# Example of Neural Network Embedding Vision Transformer (ViT)



**Source**: AN IMAGE IS WORTH 16X16 WORDS, Dosovitskiy et al., ICLR 2021

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?

Part 4: How to use search to speed up data mining?

# The ANN search pipeline



Data vectors

Index structure (Graph, Hashing, Clustering, Tree)

**BUILD**

$$x_1, x_2, \ldots, x_N$$

$$x_n \in \mathbb{R}^D$$

Index building

*~several hours*

$x_{13}$

$x_{91}$

**SEARCH**

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

$$q \in \mathbb{R}^D$$

Candidate selection

*~milliseconds*

$x_{13}$

$x_{91}$

Scan candidates

$x'_1, x'_2, \ldots, x'_L$

$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

$$x_{74}$$

# How a Theoretician Searches in a High-dimensional Space

… she doesn't, she already proved the existence of the algorithm!

# How a Theoretician Searches in a High-dimensional Space

PUFFINN [Aumüller, Christiani, Pagh, Vesterli, ESA 2019]

https://github.com/puffinn/puffinn

# How does it work?
Locality-Sensitive Hashing (LSH) [Indyk-Motwani, 1998]

$$h(p) = h_1(p) \circ h_2(p) \circ h_3(p) \in \{0,1\}^3$$

A family $\mathcal{H}$ of hash functions is **locality-sensitive**, if the collision probability of two points is decreasing with their distance to each other.
[Charikar, 2002]

**Visualization**: cecca.github.io/attimo/VLDB-supplemental/

# Standard LSH for Reporting Points at Distance $\leq r$



Dataset $S$

$h_1 = {h'}_{1,1} \circ \ldots \circ {h'}_{K,1} \in R^K$

$h_2$

$h_3$

$h_4$

$h_5$

$h_L$

$\ldots$

$k$-NN?

**Query Algorithm**: Collect all points that collide with $q$ under $h_1, \ldots, h_L$. Return all points at distance $\leq r$.

# Our Approach: Solving $k$-NN using LSH

- Check buckets $j \in \{1, \dots, L\}$, one-by-one

- keep track of **current** closest $k$ points

- **Goal**:
  Report with prob. $\geq 1 - \delta$



- What if there is no such $j$?
  - Try again with smaller $K$

**Termination**: If $(1 - p)^j \leq \delta$, report **current** top-$k$.

probability of the **current** $k$-th nearest neighbor to collide.

Why does that work? Monotonicity of the LSH collision prob.

# A Bound on the Expected Running Time

- **Assumption:** Given query $q$, know best stopping point in data structure

$$OPT(L, K, k, \delta) = \min\left\{\frac{\ln(1/\delta)}{p(q, x_k)^i}\left(i + \sum_{x \in P} p(q, x)^i\right) \mid 0 \leq i \leq K, \frac{\ln(1/\delta)}{p(q, x_k)^i} \geq L\right\}$$

#repetitions

Hash function eval

Expected bucket size

- **Lemma**: In expectation, proposed algorithm takes time

$$O(OPT(L, K, k, \delta/k) + L(K + k))$$

Additional results using this technique:
- closest pair [A., Ceccarello, SISAP 2023]
- time series motifs [Ceccarello+, VLDB 2023]
- single-linkage & hdbscan [A.+, submitted]

Recall-Queries per second (1/s) tradeoff - up and to the right is better

1.2M vectors, 100d, GloVe word embeddings

better

Better Throughput

Better Quality

Graph-based

Clustering-based

Tree-based

LSH-based

Legend:
- NGT-qg
- hnsw(nmslib)
- qsgngt
- NGT-panng
- glass
- scann
- vearch
- vamana(diskann)
- Milvus(Knowhere)
- pynndescent
- n2
- faiss-ivfpqfs
- hnsw(faiss)
- hnswlib
- hnsw(vespa)
- redisearch
- vald(NGT-anng)
- luceneknn
- weaviate
- SW-graph(nmslib)
- faiss-ivf
- flann
- mrpt
- annoy
- qdrant
- puffinn
- pgvector
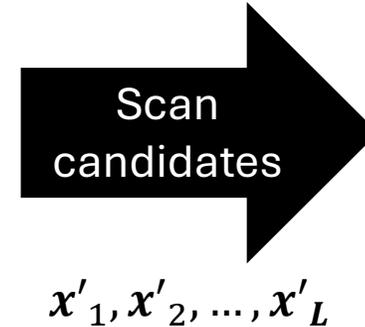- tinyknn
- BallTree(nmslib)
- bruteforce-blas

[A., Bernhardsson, Faithfull,

https://github.com/erikbern/ann-benchmarks

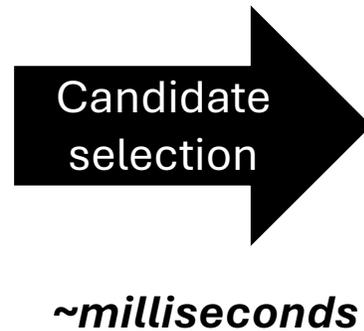| | LSH | Clustering-based | Graph-based |
|---|---|---|---|
| Supports | • range search<br>• k-NN | | |
| Pros | • strong guarantees on running time/quality<br>• data independent<br>• adaptive | | |
| Cons | • many points need to be inspected to get decent quality<br>• typically large space requirements<br>• space/distance must be "lshable" | | |

# How a Practitioner Searches in a High-dimensional Space

# Clustering-based approaches

SCANN (Google) [Guo+, 2020], FAISS IVF (Meta) [Johnson+, 2021],
LoRANN [Jääsaari+, 2024]

# IVF-based solutions ("inverted file index")

2 steps:
(1) Train partition
(2) Add vectors



$$x_1, x_2, \ldots, x_N$$

$$x_n \in \mathbb{R}^D$$

**Finding a space partition: <u>Clustering-based (k-means)</u>, LSH-based, ...**

# IVF: insert a vector

Record $x_1$

$$x_1 = \begin{bmatrix} 1.02 \\ 0.73 \\ 0.56 \\ 1.37 \\ 1.37 \\ 0.72 \end{bmatrix}$$

$x_1$



**Cells: all points closest to given centroid ("Voronoi cells")**
**Build parameter: #clusters**

# IVF: search

Find the nearest vector to $q$

$$\begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \\ 0.14 \\ 0.32 \end{bmatrix}$$

$q$

**Search parameter: #clusters to inspect**
**Candidates: #clusters inspected * avg. cluster size**

Recall-Queries per second (1/s) tradeoff - up and to the right is better

1.2M vectors, 100d, GloVe word embeddings

better

**Better Throughput**

**Better Quality**

Queries per second (1/s)

Recall

Graph-based

Clustering-based

Tree-based

LSH-based

NGT-qg
hnsw(nmslib)
qsgngt
NGT-panng
glass
scann
vearch
vamana(diskann)
Milvus(Knowhere)
pynndescent
n2
faiss-ivfpqfs
hnsw(faiss)
hnswlib
hnsw(vespa)
redisearch
vald(NGT-anng)
luceneknn
weaviate
SW-graph(nmslib)
faiss-ivf
flann
mrpt
annoy
qdrant
puffinn
pgvector
tinyknn
BallTree(nmslib)
bruteforce-blas

[A., Bernhardsson, Faithfull, 2
https://github.com/erikbern/ann-benchmarks

41

# The ANN search pipeline

Data vectors                                          Index structure (Graph, Hashing, Clustering, Tree)

**BUILD**

$$x_1, x_2, ..., x_N$$

Index building

$x_{13}$

$x_{91}$

$x_n \in \mathbb{R}^D$

*~several hours*

**SEARCH**

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

Candidate selection

$x_{13}$

$x_{91}$

Scan candidates

$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$

$q \in \mathbb{R}^D$     *~milliseconds*

$x'_1, x'_2, ..., x'_L$

$x_{74}$

42

# Basic idea

① ② Ⓝ

$$D \begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \end{bmatrix} \begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \end{bmatrix} \dots \begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \end{bmatrix}$$

**Convert**

① ② Ⓝ

| code | | code | | | code |

# Basic idea

$D$ $\begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \end{bmatrix}$ (1) $\begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \end{bmatrix}$ (2) ... $\begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \end{bmatrix}$ (N)

**Convert**

(1) code (2) code ... (N) code

What kind of conversion is preferred?

1. The "distance" between two codes can be calculated

2. The distance can be computed quickly

3. That distance approximates the distance between the original vectors (e.g., $L_2$)

4. Sufficiently small length of codes can achieve the above three criteria

# Product Quantization; PQ [Jégou+, TPAMI 2011]

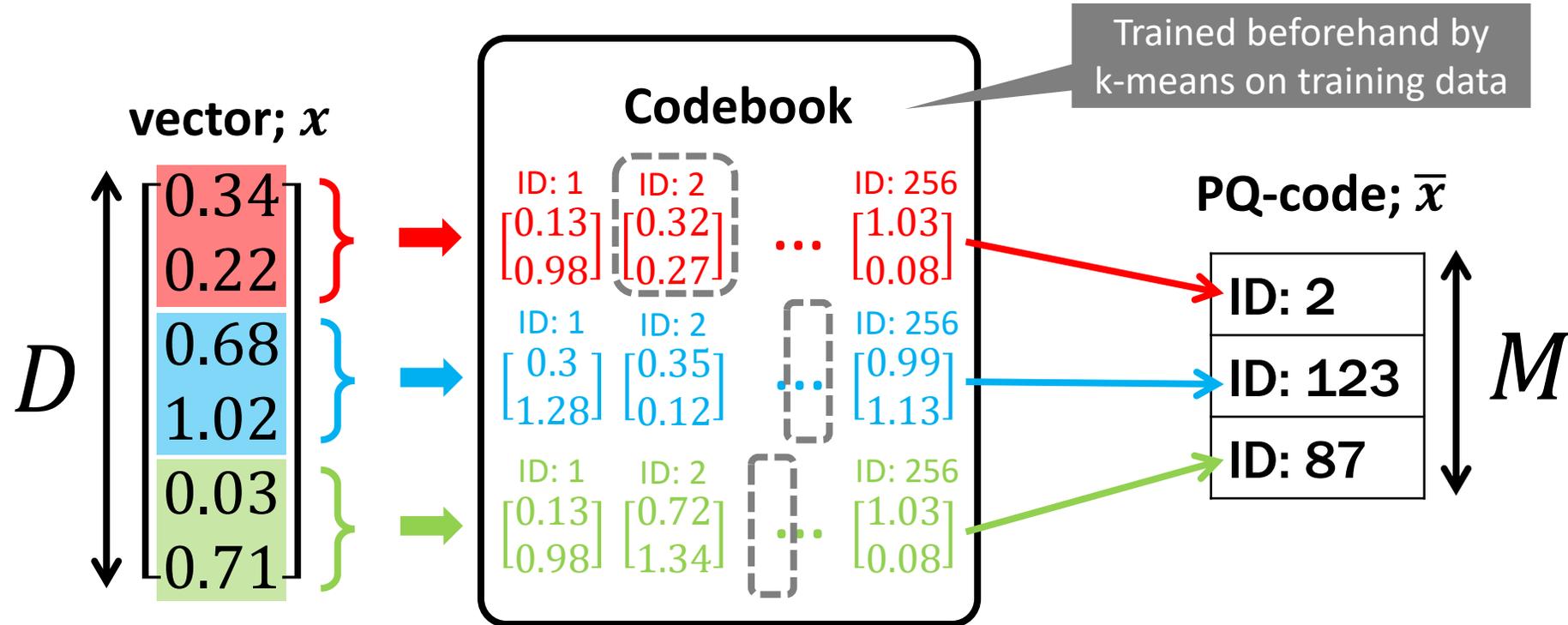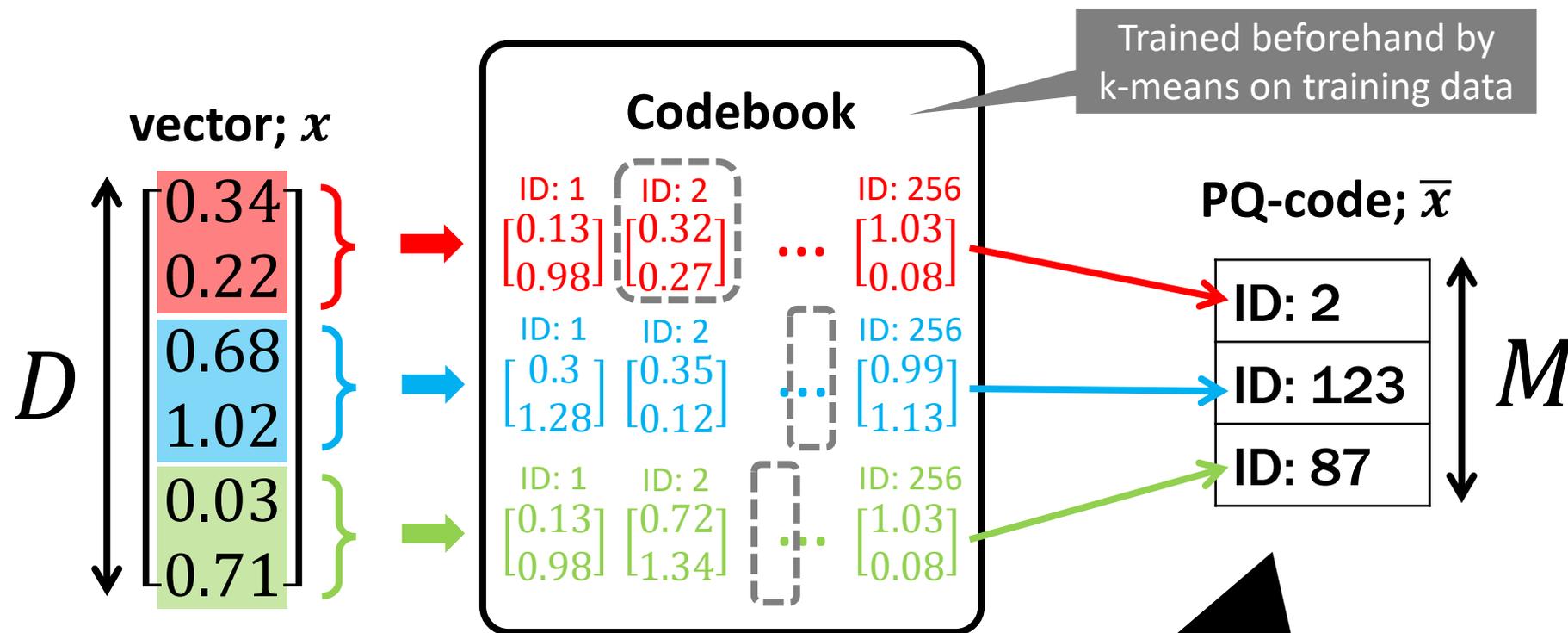➢ Split a vector into sub-vectors, and quantize each sub-vector

**vector;** $x$

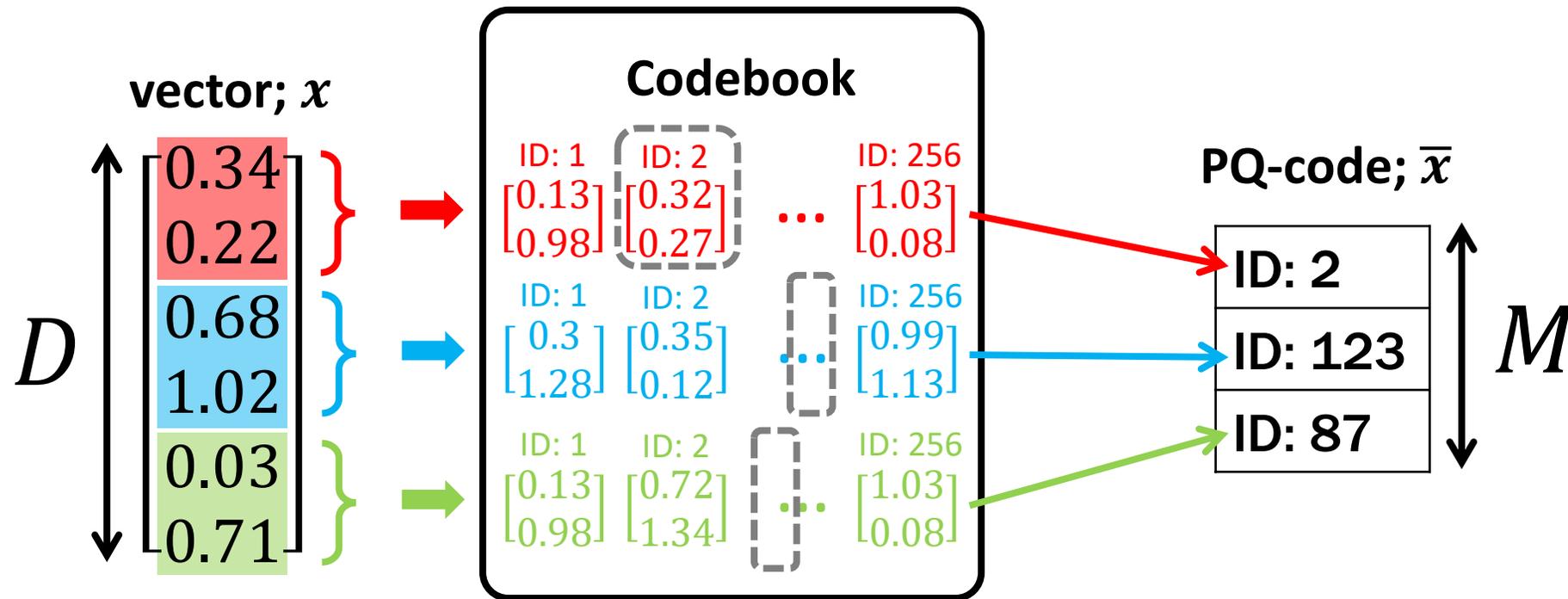$$D \begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

**Codebook**

Trained beforehand by k-means on training data

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$   ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$   $\cdots$   ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$   ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$   $\cdots$   ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$   ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$   $\cdots$   ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code;** $\bar{x}$

$M$

# Product Quantization; PQ [Jégou+, TPAMI 2011]

➤ Split a vector into sub-vectors, and quantize each sub-vector

**vector; $x$**

$$D \begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

**Codebook**

Trained beforehand by k-means on training data

ID: 1          ID: 2                    ID: 256
$\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$  $\cdots$  $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1          ID: 2                    ID: 256
$\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$  $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$  $\cdots$  $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1          ID: 2                    ID: 256
$\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$  $\cdots$  $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code; $\overline{x}$**

$M$

# Product Quantization; PQ [Jégou+, TPAMI 2011]

➢ Split a vector into sub-vectors, and quantize each sub-vector



**vector;** $x$

$D$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

**Codebook**

Trained beforehand by k-means on training data

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code;** $\overline{x}$

ID: 2

$M$

# Product Quantization; PQ [Jégou+, TPAMI 2011]

➢ Split a vector into sub-vectors, and quantize each sub-vector

**vector;** $x$

$$
D \left\{
\begin{array}{c}
0.34 \\
0.22 \\
0.68 \\
1.02 \\
0.03 \\
0.71
\end{array}
\right.
$$

**Codebook**

Trained beforehand by k-means on training data

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$  ...  ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code;** $\overline{x}$

| |
|---|
| ID: 2 |
| ID: 123 |
| |

$M$

# **Product Quantization; PQ** [Jégou+, TPAMI 2011]

➢ Split a vector into sub-vectors, and quantize each sub-vector

# Product Quantization; PQ [Jégou+, TPAMI 2011]

➢ Split a vector into sub-vectors, and quantize each sub-vector

**vector;** $x$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix} D$$

**Codebook**

Trained beforehand by k-means on training data

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code;** $\bar{x}$

| ID: 2 |
| ID: 123 |
| ID: 87 |

$M$

Bar notation for PQ-code:
$$x \in \mathbb{R}^D \mapsto \bar{x} \in \{1, \ldots, 256\}^M$$

➢ Simple
➢ Memory efficient
➢ Distance can be estimated

# Product Quantization: **Memory efficient**

# Product Quantization: **Memory efficient**



**vector;** $x$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

$D$

float: 32bit

**Codebook**

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$ $\cdots$ ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$ $\cdots$ ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$ ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$ $\cdots$ ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code;** $\bar{x}$

| |
|---|
| **ID: 2** |
| **ID: 123** |
| **ID: 87** |

$M$

e.g., $D = 128$
$128 \times 32 = 4096$ [bit]

# Product Quantization: **Memory efficient**

# Product Quantization: **Memory efficient**



**vector; $x$**

$$D \begin{cases} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{cases}$$

float: 32bit

**Codebook**

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.32 \\ 0.27 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.3 \\ 1.28 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.35 \\ 0.12 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 0.99 \\ 1.13 \end{bmatrix}$

ID: 1 $\begin{bmatrix} 0.13 \\ 0.98 \end{bmatrix}$  ID: 2 $\begin{bmatrix} 0.72 \\ 1.34 \end{bmatrix}$ ... ID: 256 $\begin{bmatrix} 1.03 \\ 0.08 \end{bmatrix}$

**PQ-code; $\bar{x}$**

$$M \begin{cases} \text{ID: 2} \\ \text{ID: 123} \\ \text{ID: 87} \end{cases}$$

uchar: 8bit

e.g., $D = 128$
$128 \times 32 = 4096$ [bit]

e.g., $M = 8$
$8 \times 8 = 64$ [bit]

**1/64**

**Product Quantization: Distance estimation**

Database vectors

Query; $q \in \mathbb{R}^D$

$x_1 \quad x_2 \qquad x_N$

$$q = \begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix} \qquad \begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \\ 0.14 \\ 0.32 \end{bmatrix} \begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \\ 1.43 \\ 0.74 \end{bmatrix} \cdots \begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \\ 0.12 \\ 2.32 \end{bmatrix}$$

**Product Quantization: <span style="color:red">Distance estimation</span>**

Database vectors

Query; $\boldsymbol{q} \in \mathbb{R}^D$  $\boldsymbol{x}_1$  $\boldsymbol{x}_2$  $\boldsymbol{x}_N$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix} \qquad \begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \\ 0.14 \\ 0.32 \end{bmatrix} \begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \\ 1.43 \\ 0.74 \end{bmatrix} \quad \cdots \quad \begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \\ 0.12 \\ 2.32 \end{bmatrix}$$

Product quantization

# Product Quantization: Distance estimation

Query; $\boldsymbol{q} \in \mathbb{R}^D$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

$\overline{x_1} \in \{1, \dots, 256\}^M$

$\overline{\boldsymbol{x}_1}$

| ID: 42 |
| ID: 67 |
| ID: 92 |

$\overline{\boldsymbol{x}_2}$

| ID: 221 |
| ID: 143 |
| ID: 34 |

$\cdots$

$\overline{\boldsymbol{x}_N}$

| ID: 99 |
| ID: 234 |
| ID: 3 |

# Product Quantization: Distance estimation

Query; $\boldsymbol{q} \in \mathbb{R}^D$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

Linear Scan Through Candidates

$\overline{x_1} \in \{1, \dots, 256\}^M$

| $\overline{\boldsymbol{x}_1}$ |
|---|
| ID: 42 |
| ID: 67 |
| ID: 92 |

| $\overline{\boldsymbol{x}_2}$ |
|---|
| ID: 221 |
| ID: 143 |
| ID: 34 |

$\cdots$

| $\overline{\boldsymbol{x}_N}$ |
|---|
| ID: 99 |
| ID: 234 |
| ID: 3 |

Asymmetric distance

➤ $d(\boldsymbol{q}, \boldsymbol{x})^2$ can be efficiently approximated by $d_A(\boldsymbol{q}, \overline{\boldsymbol{x}})^2$

➤ Lookup-trick: Looking up pre-computed distance-tables

➤ Candidate selection by $d_A$

```python
import numpy as np
from scipy.cluster.vq import vq, kmeans2
from scipy.spatial.distance import cdist


def train(vec, M):
    Ds = int(vec.shape[1] / M) # Ds = D / M
    # codeword[m][k] = c_k^m
    codeword = np.empty((M, 256, Ds), np.float32)

    for m in range(M):
        vec_sub = vec[:, m * Ds : (m + 1) * Ds]
        codeword[m], label = kmeans2(vec_sub, 256)

    return codeword

def encode(codeword, vec): # vec = {x_n}_{n=1}^N
    M, _K, Ds = codeword.shape
    # pqcode[n] = i(x_n),   pqcode[n][m] = i^m(x_n)
    pqcode = np.empty((vec.shape[0], M), np.uint8)

    for m in range(M): # Eq. (2) and Eq. (3)
        vec_sub = vec[:, m * Ds: (m + 1) * Ds]
        pqcode[:, m], dist = vq(vec_sub, codeword[m])

    return pqcode
```

```python
def search(codeword, pqcode, query):
    M, _K, Ds = codeword.shape
    # dist_table = D(m, k)
    dist_table = np.empty((M, 256), np.float32)

    for m in range(M):
        query_sub = query[m * Ds: (m + 1) * Ds]
        dist_table[m, :] = cdist([query_sub],
        ↪    codeword[m], 'sqeuclidean')[0] # Eq. (5)

    # Eq. (6)
    dist = np.sum(dist_table[range(M), pqcode], axis=1)

    return dist

if __name__ == "__main__":
    # Read vec_train, vec ({x_n}_{n=1}^N), and query (y)
    M = 4
    codeword = train(vec_train, M)
    pqcode = encode(codeword, vec)
    dist = search(codeword, pqcode, query)
    print(dist)
```

➢ Only tens of lines in Python
➢ Pure Python library: nanopq
   https://github.com/matsui528/nanopq
➢ `pip install nanopq`

# The ANN search pipeline (with quantization)

Data vectors

Index structure (Graph, IVF, Tree)

**BUILD**

$x_1, x_2, \ldots, x_N$

$x_n \in \mathbb{R}^D$

Index building

*~several hours*

$x_{13}$

$x_{91}$

**SEARCH**

$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$

$q \in \mathbb{R}^D$

*~milliseconds*

Candidate selection

$x_{13}$

$x_{91}$

Scan candidates by code

$\bar{x}'_1, , \ldots, \bar{x}'_K$

Rerank with true vectors

$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$

$x_{74}$

**Typically 10-100x more quantized vectors than target**

Recall-Queries per second (1/s) tradeoff - up and to the right is better

1.2M vectors, 100d, GloVe word embeddings

Legend:
- NGT-qg
- hnsw(nmslib)
- qsgngt
- NGT-panng
- glass
- scann
- vearch
- vamana(diskann)
- Milvus(Knowhere)
- pynndescent
- n2
- faiss-ivfpqfs
- hnsw(faiss)
- hnswlib
- hnsw(vespa)
- redisearch
- vald(NGT-anng)
- luceneknn
- weaviate
- SW-graph(nmslib)
- faiss-ivf
- flann
- mrpt
- annoy
- qdrant
- puffinn
- pgvector
- tinyknn
- BallTree(nmslib)
- bruteforce-blas

Graph-based

Clustering-based

Tree-based

LSH-based

Better Throughput

Better Quality

better

[A., Bernhardsson, Faithfull, 2
https://github.com/erikbern/ann-benchmarks

61

# Current Frontier in Quantization

- **SCANN** [Guo et al., ICML'20]
  - Quantization for Inner Product Spaces
- **RaBitQ for quantizing high-dimensional v** [Gao, Long, SIGMOD 2024]
  - D dimensional vectors → D bit strings
  - Asymmetric distance estimation
  - The RaBitQ Library, Gao et al., VecDB 2025
- **LoRANN: Low-Rank Matrix Factorization for Approximate Nearest Neighbor Search**, [*Jääsaari,Hyvönen,Roos,*NeurIPS 2024]
  - reduced-rank regression as quantizer



**Source**: Accelerating Large-Scale Inference with Anisotropic Vector Quantization, Guo et al.

|  | LSH | Clustering-based | Graph-based |
|---|---|---|---|
| Supports | • range search<br>• k-NN | • k-NN |  |
| Pros | • strong guarantees on running time/quality<br>• data independent<br>• adaptive | • small space requirements<br>• fast search through quantization<br>• fast index building<br>• parameters somewhat easy to tune |  |
| Cons | • many points need to be inspected to get decent quality<br>• typically large space requirements<br>• space/distance must be "lshable" | • many points inspected to reach decent quality |  |

# When is quantization used in data mining?

# Graph-based solutions

HNSW [Malkov,Yashunin, 2020], DiskANN [Subramanya+, 2019],
pyNNDescent [McInnes], [Dong, Wei, Charikar, Li, 2011]....

# Graph search

➢ De facto standard for million-scale data (i.e., if all data can be loaded on memory)

➢ Fast and accurate for real-world data

➢ Useful for billion-scale situation as well

  ✓ Graph-search is a building block for billion-scale systems

query

entry point

➢ Traverse graph towards the query
➢ Seems intuitive, but not so much easy to understand
➢ Review the algorithm carefully

Images are from [Malkov+, Information Systems, 2013]

# Incremental approach

➢ Add a new item to the current graph incrementally

# Refinement approach

➢ Iteratively refine an initial graph

# Incremental approach

➤ Add a new item to the current graph incrementally

# Refinement approach

➤ Iteratively refine an initial graph

Graph of
$x_1, \dots, x_{90}$

$x_{13}$



➤Each node is a database vector

Images are from [Malkov+, Information Systems, 2013]

Graph of
$x_1, \dots, x_{90}$

$x_{13}$

$x_{91}$

➢Each node is a database vector
➢Given a new database vector,

70

Images are from [Malkov+, Information Systems, 2013]

Graph of
$x_1, \dots, x_{90}$

$x_{13}$

$x_{91}$

➢Each node is a database vector
➢Given a new database vector, create new edges to neighbors

71

Images are from [Malkov+, Information Systems, 2013]

Graph of
$x_1, \dots, x_{90}$

$x_{13}$

$x_{91}$

> Each node is a database vector
> Given a new database vector, create new edges to neighbors

Images are from [Malkov+, Information Systems, 2013]

> ➤ Prune edges if some node have too many edges
> ➤ Several strategies (e.g., RNG-rule)

Graph of
$x_1, \ldots, x_{90}$

$x_{13}$

$x_{91}$

➤ Each node is a database vector
➤ Given a new database vector, create new edges to neighbors

73

# Incremental approach

➤ Add a new item to the current graph incrementally

# Refinement approach

➤ Iteratively refine an initial graph

Images are from [Subramanya+, NeurIPS 2019]



> Create an initial graph (e.g., random k-regular graph or approx. kNN graph)
> Refine it iteratively (pruning/adding edges)

Images are from [Subramanya+, NeruIPS 2019]



> Need to be moderately sparse (otherwise the graph traverse is slow)
> Some "long" edges are required for shortcut

> Create an initial graph (e.g., random graph or approx. kNN graph)
> Refine it iteratively (pruning/adding edges)

Images are from [Malkov+, Information Systems, 2013]

Name each node for explanation

Close to the query

Candidates (size = 3)

query

➢ **Given a query vector**

Images are from [Malkov+, Information Systems, 2013]



query

Close to the query

Candidates
(size = 3)

➢ Given a query vector

Close to the query

Candidates
(size = 3)

entry point

query

> Given a query vector
> Start from an entry point (e.g., M )

Close to the query

M   23.1

Candidates
(size = 3)

query

23.1

entry point

> Given a query vector
> Start from an entry point (e.g., M ). Record the distance to q.

**Search** Images are from [Malkov+, Information Systems, 2013]

Close to the query

M 23.1

Candidates (size = 3)

query

entry point

81

Images are from [Malkov+, Information Systems, 2013]



1st iteration

Close to the query

query

entry point

M  23.1

Candidates
(size = 3)

Images are from [Malkov+, Information Systems, 2013]



Close to the query

query

M  23.1

**Best**

Candidates
(size = 3)

entry point

Best

➢ Pick up the unchecked best candidate ( M )

Images are from [Malkov+, Information Systems, 2013]

Close to the query

M 23.1

Best

Candidates
(size = 3)

check!

Best

➢ Pick up the unchecked best candidate ( M ). Check it.

Close to the query

M 23.1

Best

Candidates (size = 3)

check!

Best

➤ Pick up the unchecked best candidate ( M ). Check it.
➤ Find the connected points.

Close to the query

Candidates
(size = 3)

M 23.1

check!

Best

query

➢ Pick up the unchecked best candidate ( M ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

86

**Search** Images are from [Malkov+, Information Systems, 2013]

Close to the query

Candidates (size = 3)

| | |
|---|---|
| M | 23.1 |
| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

check!

Best

query

➢ Pick up the unchecked best candidate ( M ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

| | |
|---|---|
| M | 23.1 |
| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

Candidates (size = 3)

entry point

query

Best

➢ Pick up the unchecked best candidate ( M ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

Candidates
(size = 3)

| | |
|---|---|
| M | 3.1 |
| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

entry point

Best

query

➢ Pick up the unchecked best candidate ( M ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

➢ Maintain the candidates (size=3)

**Search** Images are from [Malkov+, Information Systems, 2013]

Close to the query

Candidates (size = 3)

K 19.4

N 15.3

J 11.1

query

entry point

Best

➤ Pick up the unchecked best candidate ( M ). Check it.
➤ Find the connected points.
➤ Record the distances to q.
➤ Maintain the candidates (size=3)
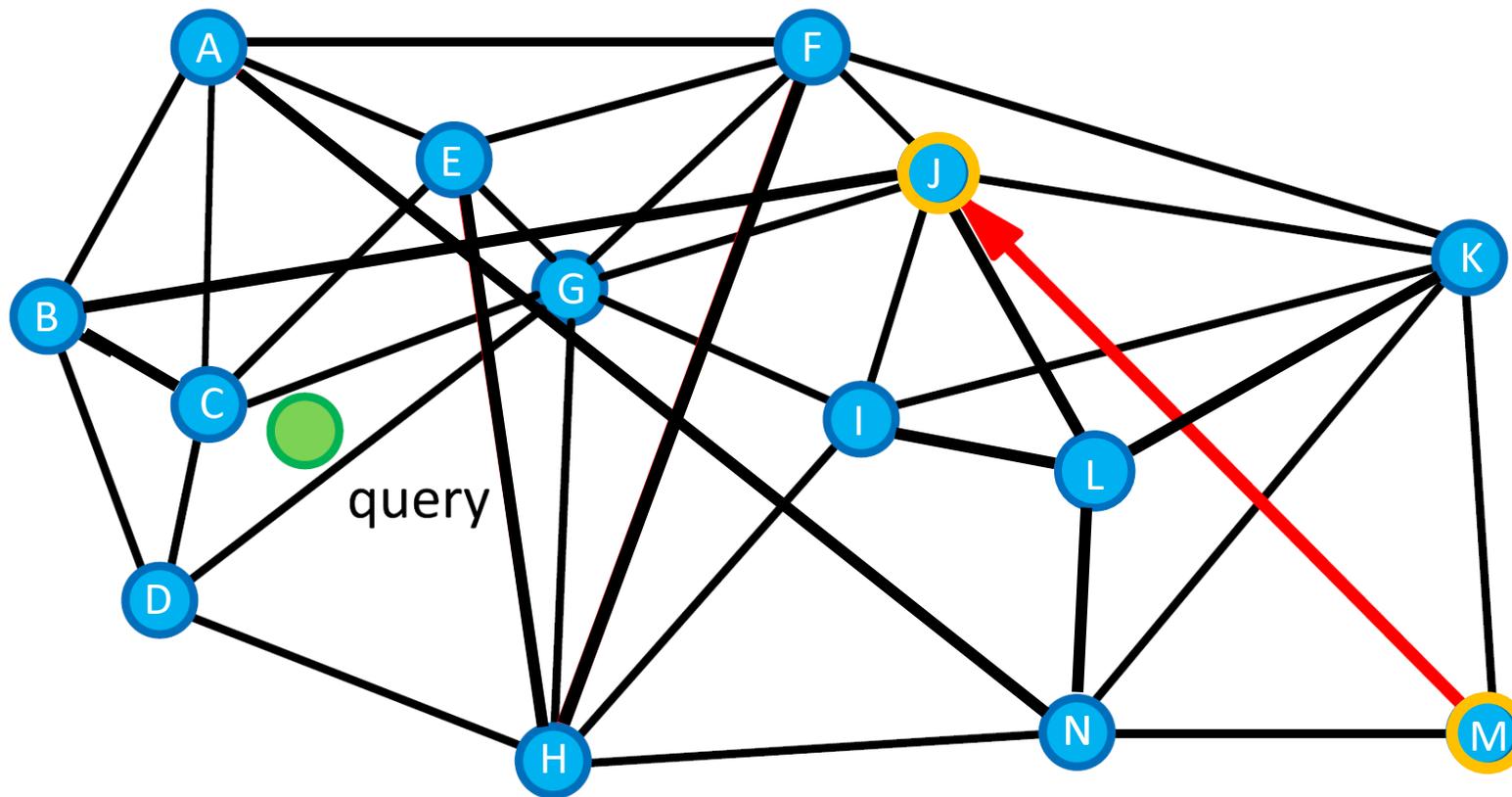
90

Close to the query

| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

Candidates (size = 3)

query

entry point

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Best

| | |
|---|---|
| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

Best

Candidates (size = 3)

query

entry point

➤ Pick up the unchecked best candidate ( J )

Images are from [Malkov+, Information Systems, 2013]

➢ Pick up the unchecked best candidate ( J ). Check it.

Images are from [Malkov+, Information Systems, 2013]



Best

check!

Close to the query

| | |
|---|---|
| K | 19.4 |
| N | 15.3 |
| J | 11.1 |

query

Best

entry point

Candidates
(size = 3)

➢ Pick up the unchecked best candidate ( J ). Check it.

➢ Find the connected points.

➤ Pick up the unchecked best candidate ( J ). Check it.

➤ Find the connected points.

➤ Record the distances to q.

Images are from [Malkov+, Information Systems, 2013]

➤ Pick up the unchecked best candidate ( J ). Check it.

➤ Find the connected points.

➤ Record the distances to q.

Images are from [Malkov+, Information Systems, 2013]

Best

Close to the query

| | |
|---|---|
| N | 15.3 |
| L | 13.2 |
| J | 11.1 |
| F | 10.2 |
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

Candidates (size = 3)

query

entry point

➢ Pick up the unchecked best candidate ( J ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

Best

query

entry point

Candidates
(size = 3)

| | |
|---|---|
| | 5.3 |
| | 3.2 |
| | 1.1 |
| | 1.2 |
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

➢ Pick up the unchecked best candidate ( J ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

➢ Maintain the candidates (size=3)

99

Close to the query

Best

| | |
|---|---|
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

Candidates (size = 3)

entry point

query

➢ Pick up the unchecked best candidate ( J ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

➢ Maintain the candidates (size=3)

100

Images are from [Malkov+, Information Systems, 2013]

Close to the query

| | |
|---|---|
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

Candidates (size = 3)

3rd iteration

query

entry point

102

Close to the query

Best

Best

I 9.7

G 3.5

B 2.3

Candidates (size = 3)

entry point

query

➢ Pick up the unchecked best candidate ( B )

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Best

check!

query

Best

entry point

Candidates (size = 3)

| | |
|---|---|
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

➢ Pick up the unchecked best candidate ( B ). Check it.

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Best

check!

query

Best

entry point

Candidates
(size = 3)

| | |
|---|---|
| I | 9.7 |
| G | 3.5 |
| B | 2.3 |

➢ Pick up the unchecked best candidate ( B ). Check it.

➢ Find the connected points.

105

Images are from [Malkov+, Information Systems, 2013]

Close to the query

Best

check!

3.6

0.5

2.1

query

I   9.7

G   3.5

B   2.3

Candidates
(size = 3)

entry point

➤ Pick up the unchecked best candidate ( B ). Check it.

➤ Find the connected points.

➤ Record the distances to q.

Close to the query

I 9.7
A 3.6
G 3.5
B 2.3
D 2.1
C 0.5

Candidates (size = 3)

A F E J K G B I L C D H N M

Best

check!

3.6

0.5

query

2.1

entry point

➢ Pick up the unchecked best candidate ( B ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

| | |
|---|---|
| I | 9.7 |
| A | 3.6 |
| G | 3.5 |
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

Candidates
(size = 3)

entry point

Best

query

➤ Pick up the unchecked best candidate ( B ). Check it.

➤ Find the connected points.

➤ Record the distances to q.

108

Close to the query

Best

query

entry point

Candidates (size = 3)

| | |
|---|---|
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

➢ Pick up the unchecked best candidate ( B ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

➢ Maintain the candidates (size=3)

Close to the query

Best

query

entry point

Candidates (size = 3)

| | |
|---|---|
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

➢ Pick up the unchecked best candidate ( B ). Check it.
➢ Find the connected points.
➢ Record the distances to q.
➢ Maintain the candidates (size=3)

110

Images are from [Malkov+, Information Systems, 2013]

Close to the query

| | |
|---|---|
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

Candidates (size = 3)

query

entry point

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Best

B 2.3

D 2.1

C 0.5

Best

entry point

Candidates
(size = 3)

query

➤ Pick up the unchecked best candidate ( C ).

113

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Best

B  2.3

D  2.1

C  0.5

Best

Candidates
(size = 3)

entry point

Best

check!

query

➤ Pick up the unchecked best candidate ( C ). Check it.

Close to the query

B 2.3
D 2.1
C 0.5

Best

Candidates
(size = 3)

entry point

Best

check!

query

➢ Pick up the unchecked best candidate ( C ). Check it.
➢ Find the connected points.

**Search** — Images are from [Malkov+, Information Systems, 2013]

Already visited

Already visited

Already visited

Already visited

Best

check!

query

entry point

Close to the query

B 2.3

D 2.1

C 0.5

Candidates (size = 3)

➢ Pick up the unchecked best candidate ( C ). Check it.
➢ Find the connected points.
➢ Record the distances to q.

Images are from [Malkov+, Information Systems, 2013]



Close to the query

Candidates
(size = 3)

B 2.3

D 2.1

C 0.5

Best

query

entry point

➤ Pick up the unchecked best candidate ( C ). Check it.
➤ Find the connected points.
➤ Record the distances to q.
➤ Maintain the candidates (size=3)

117

Close to the query

B 2.3

D 2.1

C 0.5

Candidates
(size = 3)

query

entry point

Images are from [Malkov+, Information Systems, 2013]

Close to the query

A

F

E

J

K

B

G

B 2.3

D 2.1

C 0.5

C

query

I

L

Candidates
(size = 3)

D

entry point

H

N

M

# 5ᵗʰ iteration

Close to the query

Best

B  2.3
D  2.1
C  0.5

Candidates (size = 3)

Best

query

entry point

➤ Pick up the unchecked best candidate ( D ).

Images are from [Malkov+, Information Systems, 2013]



Close to the query

B 2.3

D 2.1

C 0.5

**Best**

Candidates (size = 3)

query

Best

check!

entry point

> Pick up the unchecked best candidate ( D ). Check it.

121

Close to the query

Best

Candidates (size = 3)

| | |
|---|---|
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

query

Best

check!

entry point

➢ Pick up the unchecked best candidate ( D ). Check it.

➢ Find the connected points.

122

Images are from [Malkov+, Information Systems, 2013]

Already visited

Already visited

Already visited

Close to the query

B 2.3

D 2.1

C 0.5

Candidates (size = 3)

query

3.9

Best

check!

entry point

➢ Pick up the unchecked best candidate ( D ). Check it.

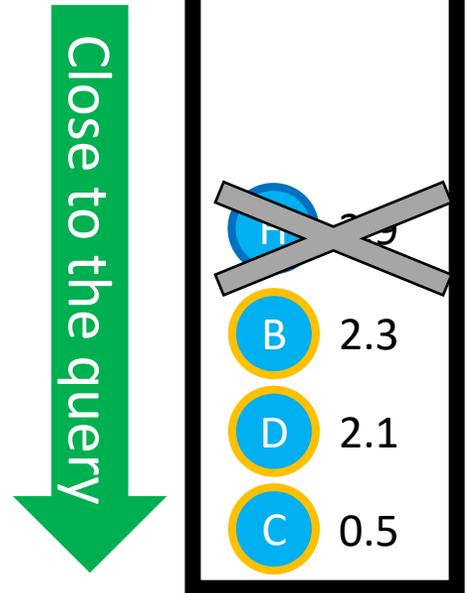➢ Find the connected points.

➢ Record the distances to q.

Close to the query

Already visited
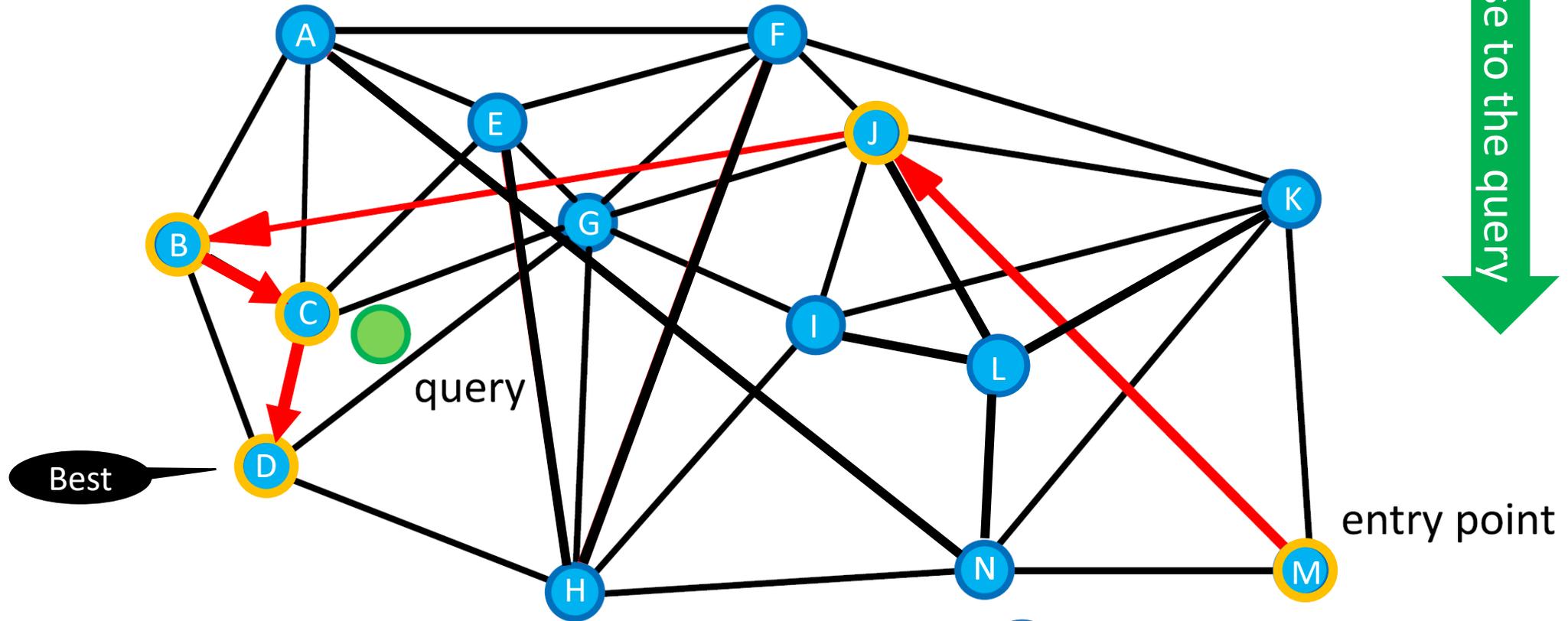
Already visited

Already visited

query

3.9

Best

check!

entry point

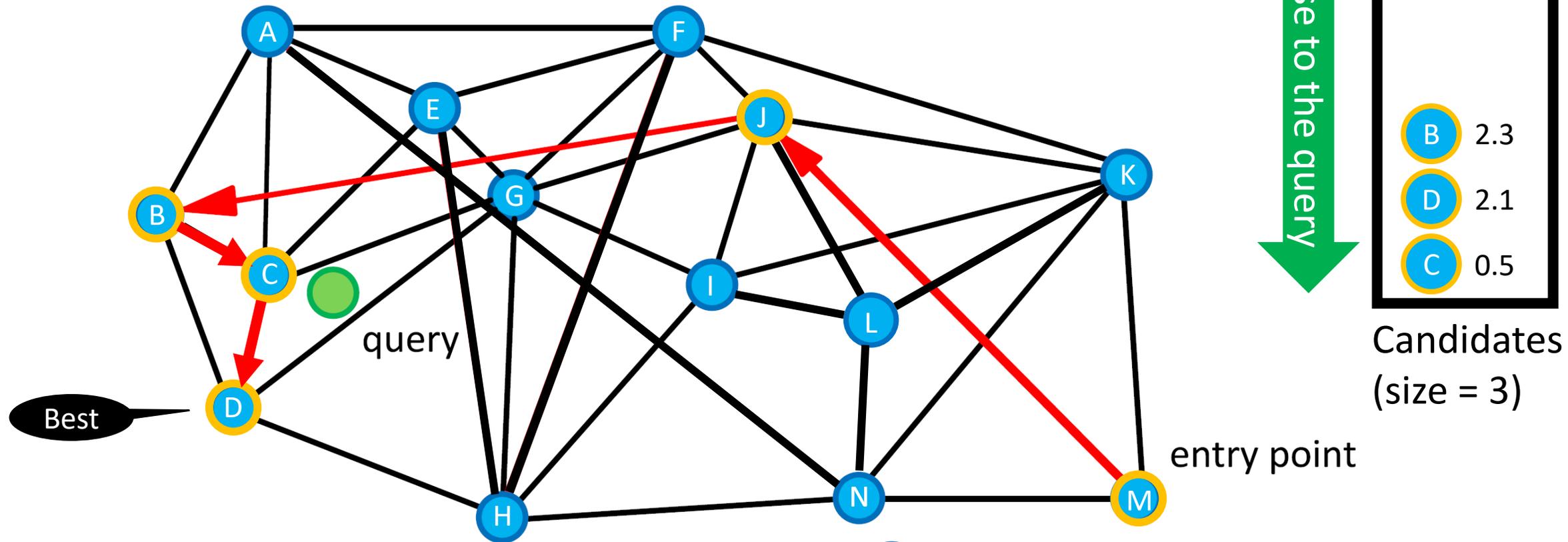| | |
|---|---|
| H | 3.9 |
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

Candidates (size = 3)

➢ Pick up the unchecked best candidate ( D ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

| | |
|---|---|
| H | 3.9 |
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

Candidates (size = 3)

query

Best

entry point

➢ Pick up the unchecked best candidate ( D ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

Close to the query

A  F
E  J
B  K
G
C
query
I
L
D
Best
entry point
H  N  M

Candidates
(size = 3)

H
B  2.3
D  2.1
C  0.5

➢ Pick up the unchecked best candidate ( D ). Check it.
➢ Find the connected points.
➢ Record the distances to q.
➢ Maintain the candidates (size=3)

126

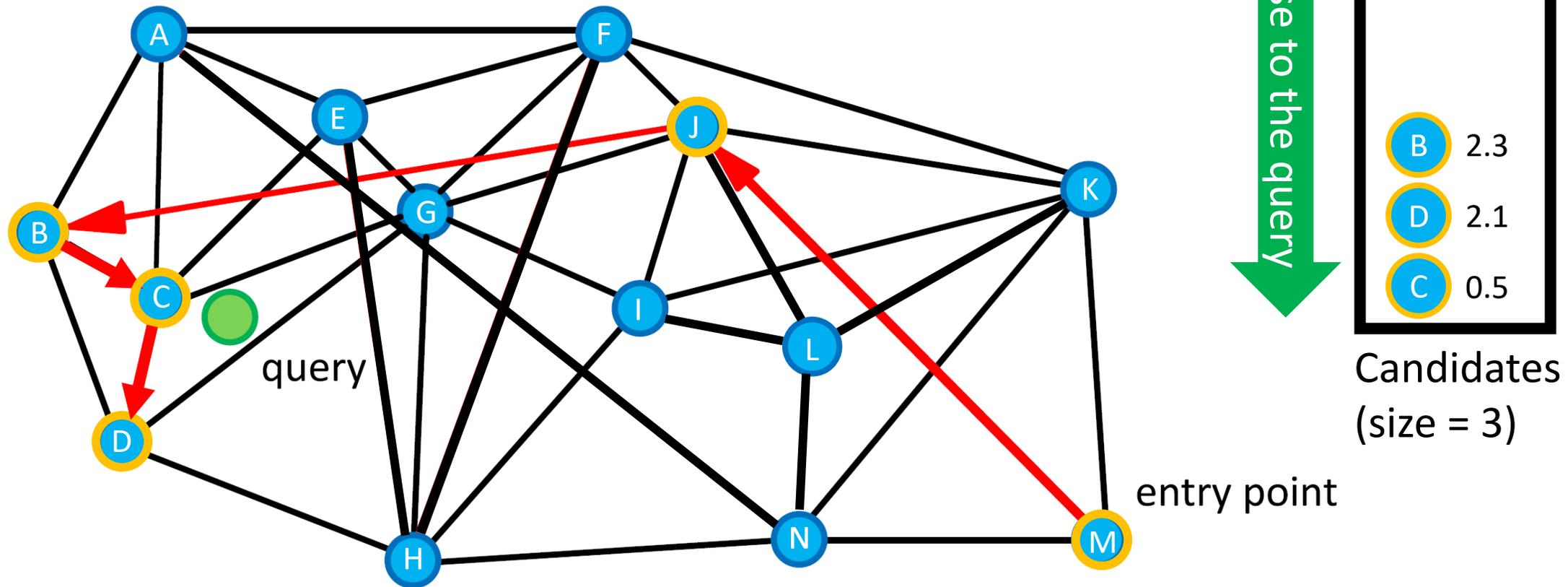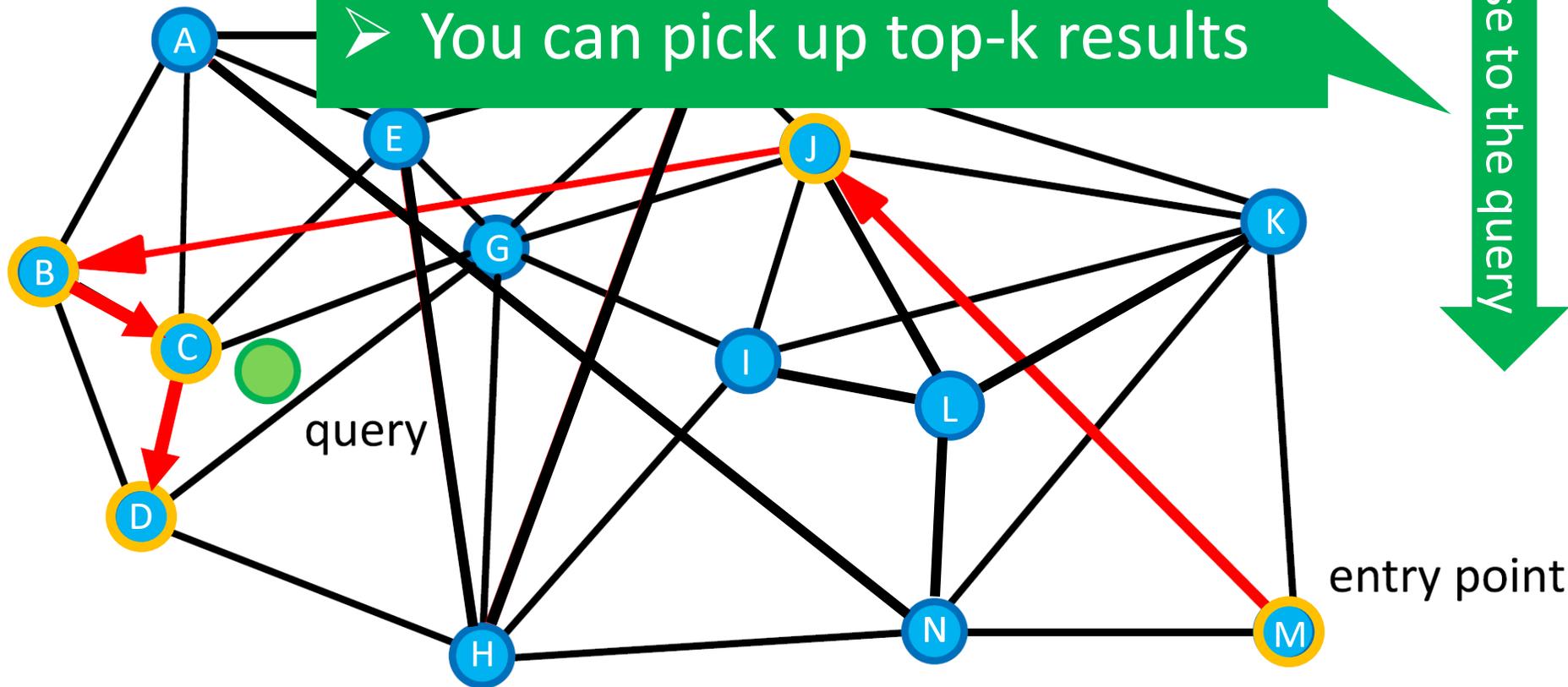Close to the query

B 2.3

D 2.1

C 0.5

Candidates (size = 3)

query

Best

entry point

➢ Pick up the unchecked best candidate ( D ). Check it.

➢ Find the connected points.

➢ Record the distances to q.

➢ Maintain the candidates (size=3)

127

Close to the query

| | |
|---|---|
| B | 2.3 |
| D | 2.1 |
| C | 0.5 |

Candidates (size = 3)

entry point

query

➢ All candidates are checked. Finish.
➢ Here, C is the closet to the query ( 🟢 )

**Final output 1: Candidates**
➢ You can pick up top-k results

Close to the query

B 2.3
D 2.1
C 0.5

Candidates (size = 3)

query

entry point

➢ All candidates are checked. Finish.
➢ Here, C is the closet to the query ( )

Final output 1: **Candidates**
➢ You can pick up top-k results

Close to the query

A

E

J

K

B 2.3

G

D 2.1

C 0.5

B

C

I

query

L

D

Candidates
(size = 3)

entry point

N

M

H

➢ All candidates are checked. Finish.

Final output 2: **Checked items**
➢ i.e., search path

query ( 🟢 )

Final output 1: **Candidates**
➢ You can pick up top-k results

Close to the query

A
E
J
K
B   2.3
G
B
C
D   2.1
I
C   0.5
L
query
Candidates
(size = 3)
D
entry point
H
N
M

➢ All candidates are checked. Finish

Final output 2: **Checked items**
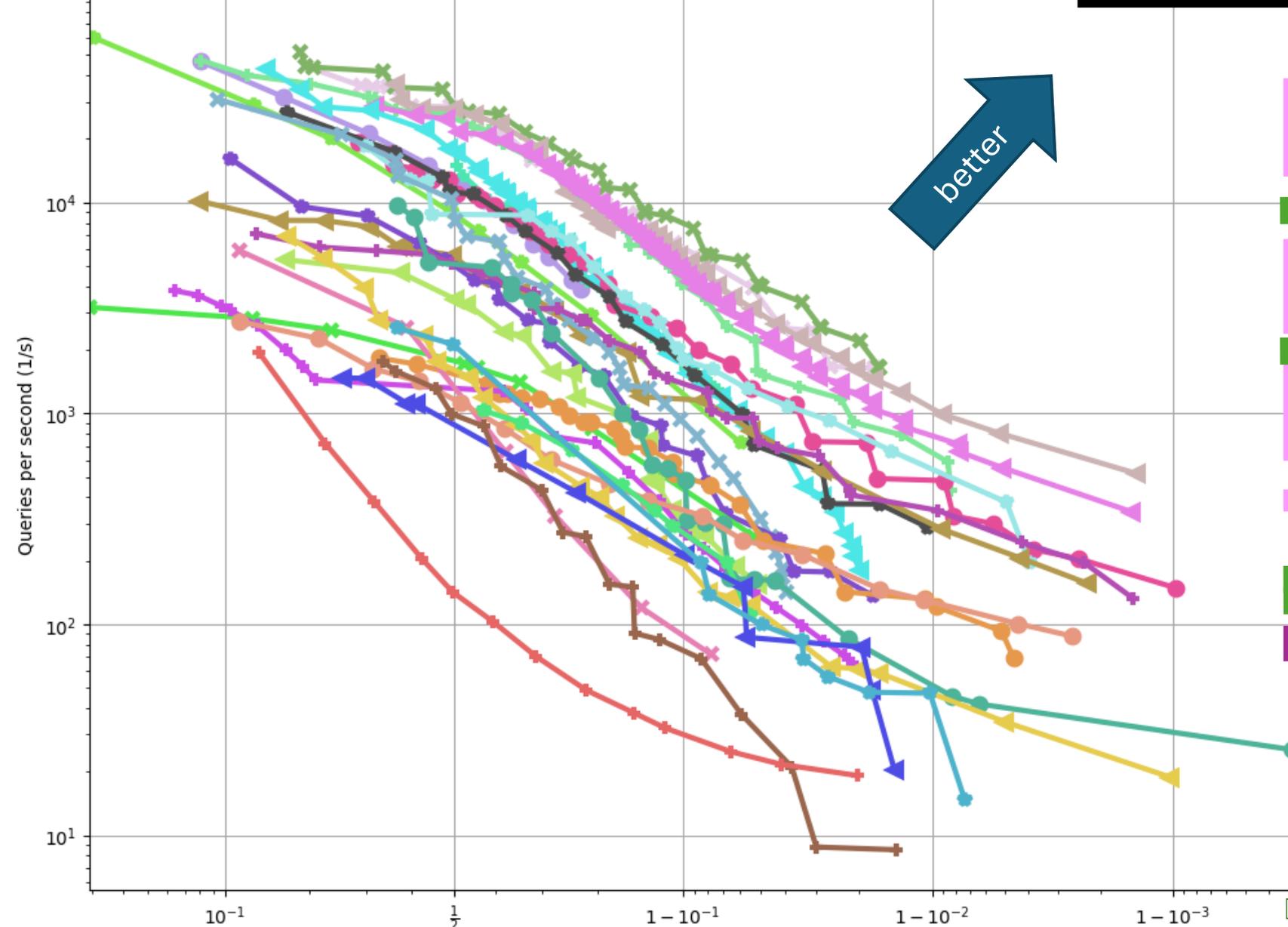➢ i.e., search path

Final output 3: **Visit flag**
➢ For each item, visited or not

Recall-Queries per second (1/s) tradeoff - up and to the right is better

1.2M vectors, 100d, GloVe word embeddings

better

Better Throughput

Queries per second (1/s)

$10^4$

$10^3$

$10^2$

$10^1$

Graph-based

Clustering-based

Tree-based

LSH-based

Legend:
- NGT-qg
- hnsw(nmslib)
- qsgngt
- NGT-panng
- glass
- scann
- vearch
- vamana(diskann)
- Milvus(Knowhere)
- pynndescent
- n2
- faiss-ivfpqfs
- hnsw(faiss)
- hnswlib
- hnsw(vespa)
- redisearch
- vald(NGT-anng)
- luceneknn
- weaviate
- SW-graph(nmslib)
- faiss-ivf
- flann
- mrpt
- annoy
- qdrant
- puffinn
- pgvector
- tinyknn
- BallTree(nmslib)
- bruteforce-blas

Better Quality

$10^{-1}$    $\frac{1}{2}$    $1-10^{-1}$    $1-10^{-2}$    $1-10^{-3}$

Recall

[A., Bernhardsson, Faithfull, 2

132

https://github.com/erikbern/ann-benchmarks

# Design choices for Graph-based ANN

- **Starting point**
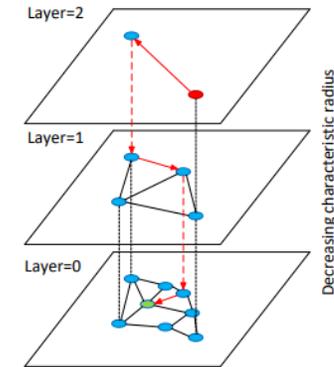  - Low-Quality Index (NGT), Hierarchy (HNSW), Medoid (DiskANN), Random

- **Neighborhood diversification**
  - Which edges to add?
  - How to keep "navigability"?

- **Degree bound**
  - Fix degree of a node
  - Directed/undirected?

- **Search algorithm** = Beam Search



**Source**: Malkov+, HNSW, 2016.



**Algorithm 2:** RobustPrune($p, \mathcal{V}, \alpha, R$)

**Data:** Graph $G$, point $p \in P$, candidate set $\mathcal{V}$, distance threshold $\alpha \geq 1$, degree bound $R$

**Result:** $G$ is modified by setting at most $R$ new out-neighbors for $p$

**begin**
$\quad \mathcal{V} \leftarrow (\mathcal{V} \cup N_{\text{out}}(p)) \setminus \{p\}$
$\quad N_{\text{out}}(p) \leftarrow \emptyset$
$\quad$**while** $\mathcal{V} \neq \emptyset$ **do**
$\quad\quad p^* \leftarrow \arg\min_{p' \in \mathcal{V}} d(p, p')$
$\quad\quad N_{\text{out}}(p) \leftarrow N_{\text{out}}(p) \cup \{p^*\}$
$\quad\quad$**if** $|N_{\text{out}}(p)| = R$ **then**
$\quad\quad\quad$ break
$\quad\quad$**for** $p' \in \mathcal{V}$ **do**
$\quad\quad\quad$**if** $\alpha \cdot d(p^*, p') \leq d(p, p')$ **then**
$\quad\quad\quad\quad$ remove $p'$ from $\mathcal{V}$

**Source**: Subramanya+, NeurIPS 2019

# What matters?



## Graph-Based Vector Search: An Experimental Evaluation of the State-of-the-Art

ILIAS AZIZI, UM6P, Université Paris Cité, Morocco - France
KARIMA ECHIHABI, UM6P, Morocco
THEMIS PALPANAS, Université Paris Cité, France

Vector data is prevalent across business and scientific applications, and its popularity is growing with the proliferation of learned embeddings. Vector data collections often reach billions of vectors with thousands of dimensions, thus, increasing the complexity of their analysis. Vector search is the backbone of many critical analytical tasks, and graph-based methods have become the best choice for analytical tasks that do not require guarantees on the quality of the answers. We briefly survey in-memory graph-based vector search, outline the chronology of the different methods and classify them according to five main design paradigms: seed selection, incremental insertion, neighborhood propagation, neighborhood diversification, and divide-and-conquer. We conduct an exhaustive experimental evaluation of twelve state-of-the-art methods on seven real data collections, with sizes up to 1 billion vectors. We share key insights about the strengths and limitations of these methods; e.g., the best approaches are typically based on incremental insertion and neighborhood diversification, and the choice of the base graph can hurt scalability. Finally, we discuss open research directions, such as the importance of devising more sophisticated data-adaptive seed selection and diversification strategies.
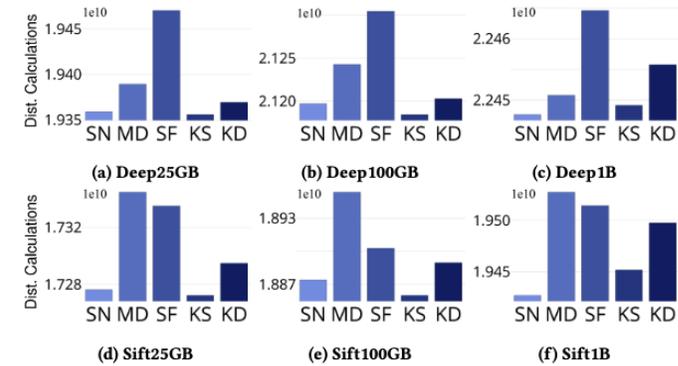
https://arxiv.org/pdf/2502.05575

Fig. 6. The impact of SS Methods on query answering

| Method | Query Answering | | | Index Building | | |
|---|---|---|---|---|---|---|
| | Efficiency | Accuracy | Tuning | Efficiency | Footprint | Tuning |
| HNSW | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| ELPIS | ✓ | ✓ | ~ | ✓ | ✓ | ~ |
| VAMANA | ✓ | ✓ | ✓ | ✓ | ✓ | ~ |
| NSG | ✓ | ✓ | ✓ | ~ | ~ | ~ |
| SSG | ✓ | ✓ | ✓ | ~ | ~ | ~ |
| EFANNA | × | ~ | × | × | × | × |
| KGRAPH | × | × | × | × | × | × |
| DPG | × | ~ | ~ | ~ | ~ | ~ |
| SPTAG | ~ | ✓ | × | × | ✓ | × |
| HCNNG | ✓ | ✓ | ✓ | ✓ | ✓ | ~ |
| LSHAPG | × | ~ | × | ~ | ✓ | ✓ |
| NGT | ~ | ~ | × | × | ✓ | × |
| SPTAG | ~ | ~ | ~ | × | ✓ | × |

✓ Good  ~ Medium  × Bad

Table 3. Comparative Analysis

# 4 recent must-read papers

## ParlayANN: Scalable and Deterministic Parallel Graph-Based Approximate Nearest Neighbor Search Algorithms

Magdalen Dobson Manohar
Carnegie Mellon University
mrdobson@cs.cmu.edu

Zheqi Shen
UC Riverside
zshen055@ucr.edu

Guy E. Blelloch
Carnegie Mellon University
guyb@cs.cmu.edu

Laxman Dhulipala
University of Maryland
laxman@umd.edu

Yan Gu
UC Riverside
ygu@cs.ucr.edu

Harsha Vardhan Simhadri
Microsoft Research
harshasi@microsoft.com

Yihan Sun
UC Riverside
yihans@cs.ucr.edu

## SymphonyQG: Towards Symphonious Integration of Quantization and Graph for Approximate Nearest Neighbor Search

Yutong Gou[†]
Nanyang Technological University
Singapore
yutong003@e.ntu.edu.sg

Jianyang Gao[†]
Nanyang Technological University
Singapore
jianyang.gao@ntu.edu.sg

Yuexuan Xu
Nanyang Technological University
Singapore
yuexuan001@e.ntu.edu.sg

Cheng Long[*]
Nanyang Technological University
Singapore
c.long@ntu.edu.sg

## RoarGraph: A Projected Bipartite Graph for Efficient Cross-Modal Approximate Nearest Neighbor Search

Meng Chen
Fudan University
mengchen22@m.fudan.edu.cn

Kai Zhang
Fudan University
zhangk@fudan.edu.cn

Zhenying He
Fudan University
zhenying@fudan.edu.cn

Yinan Jing
Fudan University
jingyn@fudan.edu.cn

X.Sean Wang
Fudan University
xywangcs@fudan.edu.cn

## Worst-case Performance of Popular Approximate Nearest Neighbor Search Implementations: Guarantees and Limitations

Piotr Indyk
MIT
indyk@mit.edu

Haike Xu
MIT
haikexu@mit.edu

| | LSH | Clustering-based | Graph-based |
|---|---|---|---|
| Supports | • range search<br>• k-NN | • k-NN | - k-NN<br>- range search [Manohar, Kim, Blelloch, 2025] |
| Pros | • strong guarantees on running time/quality<br>• data independent<br>• adaptive | • small space requirements<br>• fast search through quantization<br>• fast index building<br>• parameters somewhat easy to tune | • unmatched search performance on real-world datasets |
| Cons | • many points need to be inspected to get decent quality<br>• typically large space requirements<br>• space/distance must be "lshable" | • many points inspected to reach decent quality | • long index build times<br>• Space requirements can be large<br>• parameters can be obscure |

# Short break

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?

Part 4: How to use search to speed up data mining?

# I came up with a great ANN method!! Now what?

- Make a good implementation of your method
  - kANNolo
    https://github.com/TusKANNy/kannolo
  - PANNA
    https://github.com/Cecca/panna
- Expose a Python API
  - build, search
- Run through standardized benchmark





kANNolo

paper ECIR 2025 arXiv 2501.06121

PANNA: Playground for Approximate Nearest Neighbor Algorithms

This library aims at providing useful building blocks to implement algorithms for approximate nearest neighbor search.

# Overview over Benchmarks

**"Standard ANN" (my focus)**
- ann-benchmarks
- big-ann-benchmarks
- vibe

**Vector Databases**
- https://github.com/zilliztech/VectorDBBench
- https://github.com/qdrant/vector-db-benchmark

- annbench (light-weight!)

# "My Overview"



Star History

# Overview over Architecture

# Design Principles
## Evaluation != Implementation

- We only provide evaluation infrastructure

- We provide datasets and workloads

- We are in charge of running evaluation

- System tests via Github actions
  - All implementations run on small dummy datasets



```
ann-benchmarks / .github / workflows / benchmarks.yml

Code   Blame

23
24      run-benchmarks:
25        runs-on: ubuntu-22.04
26        timeout-minutes: 30
27        strategy:
28          fail-fast: false
29          matrix:
30            dataset: [random-xs-20-angular]
31            library:
32              - annoy
33              - balltree
34              - bruteforce
35              - ckdtree
36              - descartes
37              - diskann
38              - dolphinnpy
39              - elasticsearch
40              - elastiknn
41              - expann
42              - faiss
43              - flann
44              - glass
```

# Design Principles
## Authors are in full control

- Authors contribute Docker install file, implementation and hyperparameters

- Managed through Pull Request system on GitHub

```
float:
  any:
   annoy:
    constructor: Annoy
    base-args: ["@metric"]
    run-groups:
     one-or-two-hundred-trees:
      args: [[100, 200], [100, 200, 400,
1000]]
      four-hundred-trees:
       args: [400, [1000, 2000, 4000,
10000]]
```

**config.yml**

Annoy("euclidean", 100, 100)
Annoy("euclidean", 100, 200)
Annoy("euclidean", 100, 400)          Annoy("euclidean", 400, 1000)
Annoy("euclidean", 100, 1000)         Annoy("euclidean", 400, 2000)
Annoy("euclidean", 200, 100)          Annoy("euclidean", 400, 4000)
Annoy("euclidean", 200, 200)          Annoy("euclidean", 400, 10000)
Annoy("euclidean", 200, 400)
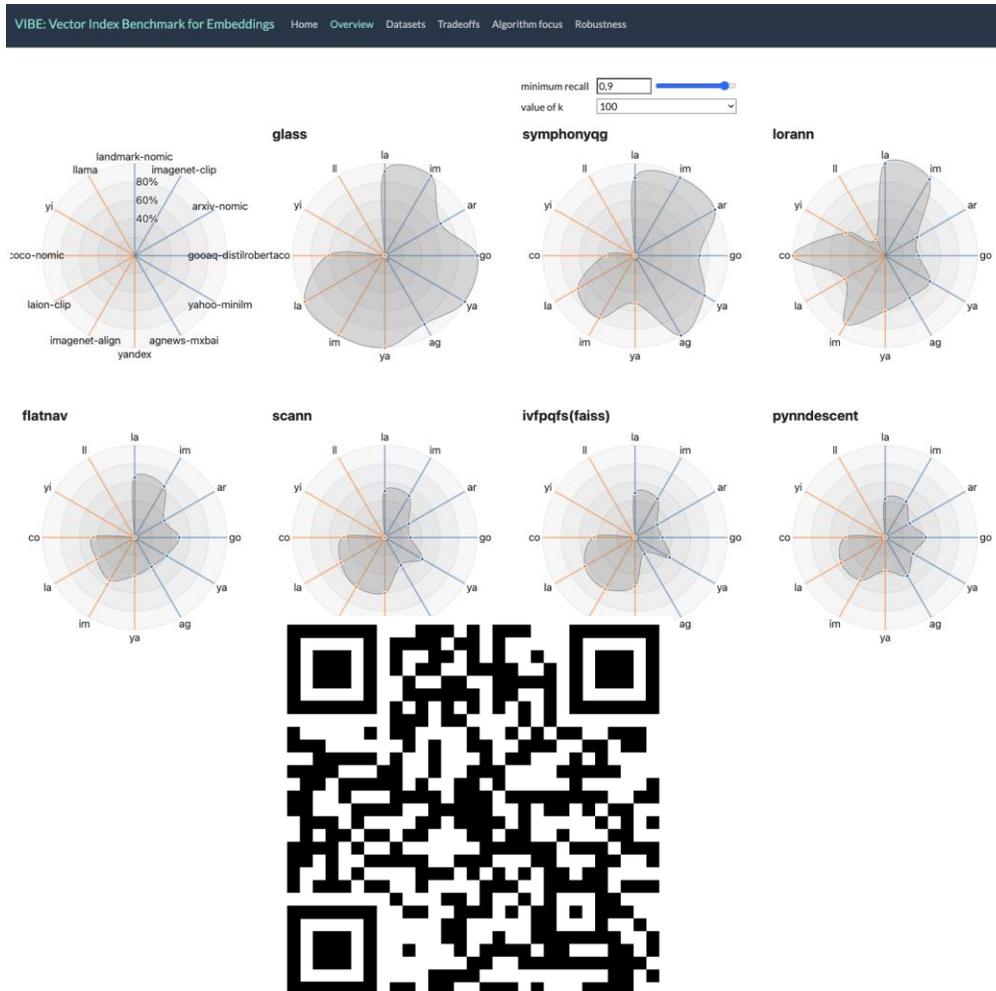Annoy("euclidean", 200, 1000)

# Design Principles
Modularity and Usability

- Implementation **isolation** through **containers**

- **Infrastructure usable outside** of main evaluation pipeline:
  - **Datasets**: Available as single download, contain groundtruth
  - **Evaluation**: Write raw results, everything else is postprocessing

```
{
  "name": "Annoy(n_trees=100, search_k=100)",
  "results": [
    [
      0.00017213821411132812, [
        [240, 0.24573669837836787],
        [1250, 0.24768519849506598],
        [341, 0.32758953153834214],
        [729, 0.3286839883166234],
        [1627, 0.3457532474865024],
        [1100, 0.34957032952284117],
        [1631, 0.361215601983618],
        [1310, 0.4161217668585876],
        [1672, 0.43032373474308594],
        [281, 0.4417039809892499]
      ]
    ], ...,
    [

        9.989738464355469e-05, [
          [1286, 0.15623190857718705], ...,
          [1610, 0.22010540054623218]
        ]
    ]
  ],
  "build_time": 0.8827729225158691, "index_size": 696,
  "candidates": 10,                   "library": "annoy",
  "run_alone": true,                  "run_count": 1,
  "best_search_time": 0.00018098115921020509
}
```

# Recent: Vector Index Benchmark (VIBE)
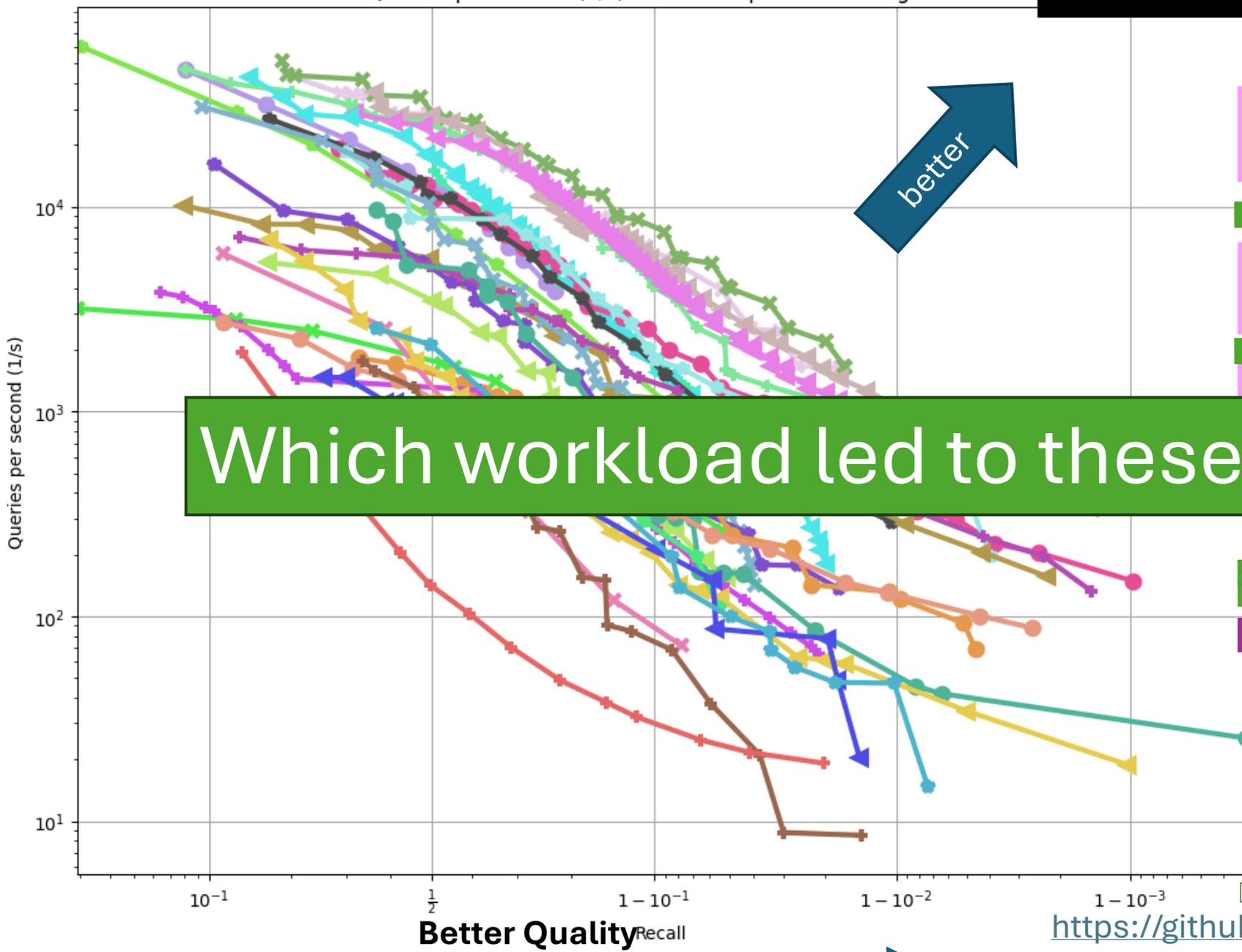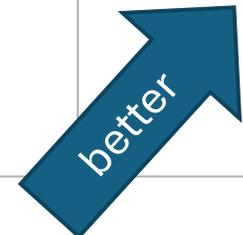## [ Jääsaari, Hyvönen, Ceccarello, Roos, A., submitted]



- Research-oriented variant of ann-benchmarks
- transparent dataset creation
  - Modern neural embedding
  - In-distribution vs. out of distribution
- Supports HPC infrastructure
- Better evaluation/inspection

https://vector-index-bench.github.io/

Recall-Queries per second (1/s) tradeoff - up and to the right is better

1.2M vectors, 100d, GloVe word embeddings

**Better Throughput**

Queries per second (1/s)

better

Which workload led to these results?

Graph-based

Clustering-based

Tree-based

LSH-based

Legend:
- NGT-qg
- hnsw(nmslib)
- qsgngt
- NGT-panng
- glass
- scann
- vearch
- vamana(diskann)
- Milvus(Knowhere)
- pynndescent
- n2
- faiss-ivfpqfs
- hnsw(faiss)
- weaviate
- SW-graph(nmslib)
- faiss-ivf
- flann
- mrpt
- annoy
- qdrant
- puffinn
- pgvector
- tinyknn
- BallTree(nmslib)
- bruteforce-blas

**Better Quality** Recall

[A., Bernhardsson, Faithfull, 2
https://github.com/erikbern/ann-benchmarks

# Easy and hard queries

**Cost measure:** the number distance computations required to answer a query with a given recall (say 0.9)

## Easy

queries require few distance computations, hence are fast to answer

## Hard

queries require computing many distances, hence are slower

The overall performance depends on the distribution of difficulties of queries, that should not be left to random chance

# Assessing Workload Difficulty [Aumüller, Ceccarello 2021]



$$\text{LID}_k(\vec{q}) = -\left(\frac{1}{k}\sum_{i=1}^{k}\log\frac{r_i}{r_k}\right)^{-1}$$

$$\text{RC}_k(\vec{q}) = \frac{\frac{1}{n}\sum_{i=1}^{n}r_i}{r_k}$$

$$\text{Expansion}_k(\vec{q}) = \frac{r_{2k}}{r_k}$$

**Takeaways**
- Can pick query points adversarially from dataset to manipulate difficulty.
- Can make "diverse query workloads" to challenge algorithms in terms of robustness and adaptiveness.

# What makes a query intrinsically difficult?

- A good candidate is the position of the query relative to the dataset
- We can look at the distribution of distances, condensing it to a single number

# Generating workloads [Ceccarelo, Levchenko, Ileana, Palpanas, KDD 2025]

Input:

Data points

Hephaestus

Output:



Target hardness: RC=1.4
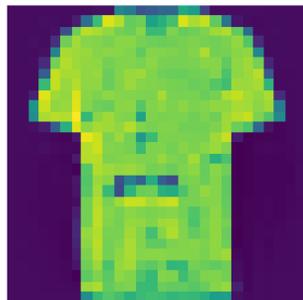
Query placed to achieve the desired hardness

```
hephaestus --dataset fashion-mnist-784-euclidean.hdf5 --output queries.hdf5 -k 10 -q 1:1.4
```

Command line

https://github.com/Cecca/hephaestus

# Easy query

Relative contrast 4

IVF index inspects 4.6% of the dataset to answer it



All answers are shirt items, very similar to the query

# Hard query

Relative contrast 1.09

IVF index inspects 45.1% of the dataset to answer it!



Some answers are shirt items, while others are bags

How does all of this differ from data mining best practices?

Where are opportunities?

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?

Part 4: How to use search to speed up data mining?

# Credits

## On the Design of Scalable Outlier Detection Methods using Approximate Nearest Neighbor Graphs

Camilla Birch Okkels[1],
Martin Aumüller[1][0000−0002−7212−6476], and Arthur Zimek[2][0000−0001−7713−4208]

[1] IT University of Copenhagen, Denmark, {cabi, maau}@itu.dk
[2] University of Southern Denmark, Denmark, zimek@imada.sdu.dk

## Approximate Single-Linkage Clustering Using Graph-based Indexes: MST-based Approaches and Incremental Searchers

Camilla Birch Okkels[1],
Erik Thordsen[3], Martin Aumüller[1], Arthur Zimek[2], and Erich Schubert[3]

[1] IT University of Copenhagen, Denmark
{cabi, maau}@itu.dk
[2] University of Southern Denmark, Denmark
zimek@imada.sdu.dk
[3] TU Dortmund, Germany
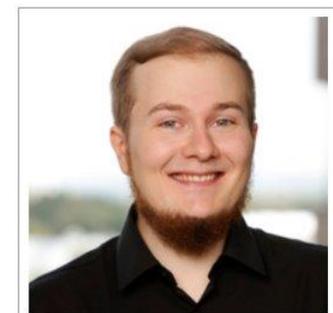{erik.thordsen, erich.schubert}@tu-dortmund.de

Camilla Okkels

Arthur Zimek

Erich Schubert

Erik Thordsen

# Using Approximate Search to Accelerate Data Mining Algorithms: Two Stories...

## Outlier Detection

- **Input**: Dataset $S \subseteq R^d$

- **Task**: Compute *outlier score* for each vector in $S$

- **Unsupervised methods:**
  - **K-NN**: score = sum of $k$-NN distances
  - **LOF**: Contrast in neighborhood distances

**"Easy" using ANN Search**

## Single-Linkage Clustering

- **Input**: Dataset $S \subseteq R^d$

- **Task**: Compute Dendrogram of single-linkage clustering

- = Computing MST in $(S, S \times S)$.

**How to solve using ANN search?**

# Using Approximate Search to Accelerate Data Mining Algorithms: Two Stories... One idea!

**Black-box**

- Build ANN index on dataset
- Use search to produce approximate neighborhoods
- Use these to solve the downstream task (clustering/outlier detection/...)

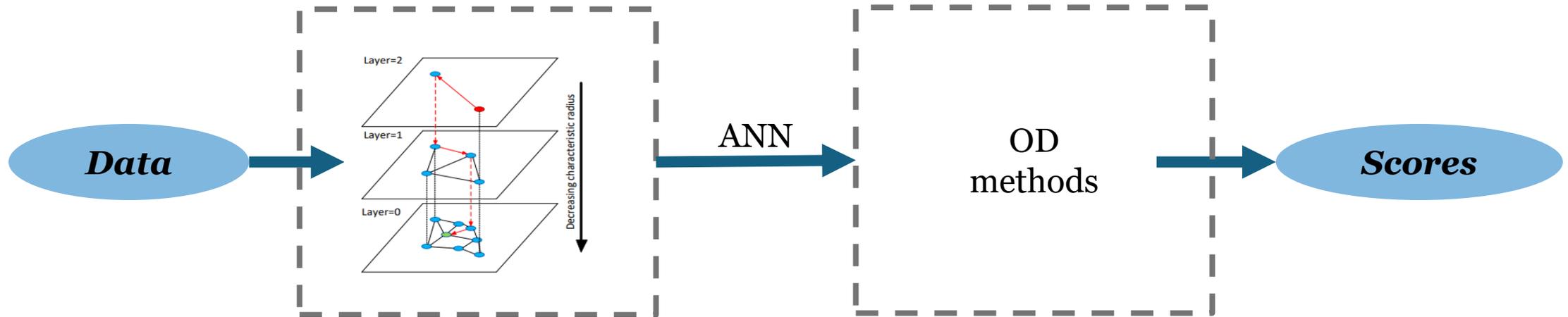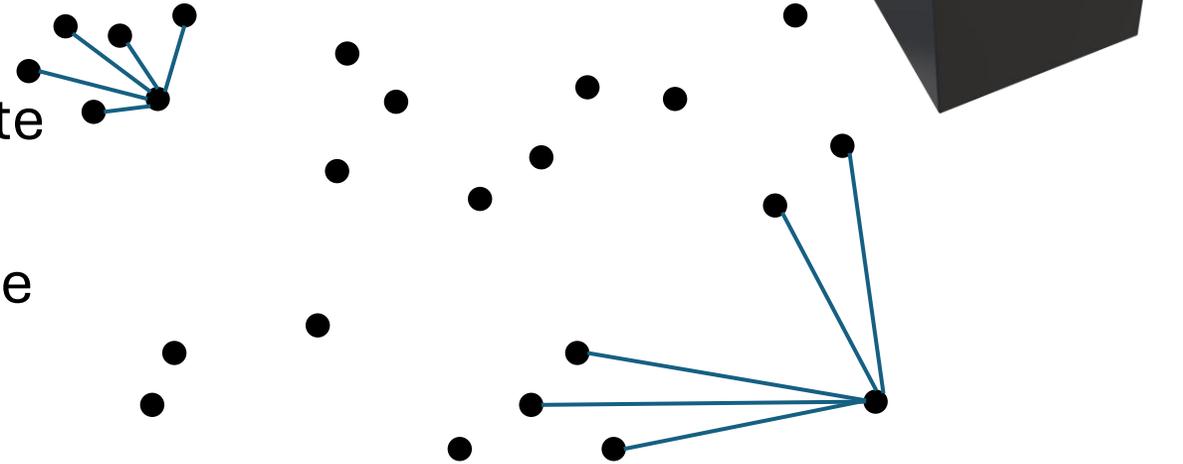**Game is different: Index building is part of execution time!**

**White-box**

- Build ANN index on dataset
- Directly solve downstream task by "clever inspection" of the data structure

- **Hope**:
    - More efficient than "blackbox" search
    - Applicable when search is incompatible (range vs. top-$k$)
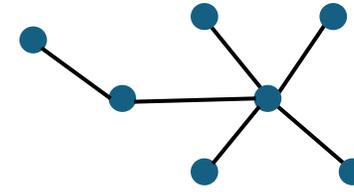
# Story 1:Outlier detection
## Blackbox method

- Query on the graph to find approximate nearest neighbours.

- Calculate k-NN/LOF from approximate nearest neighbours.
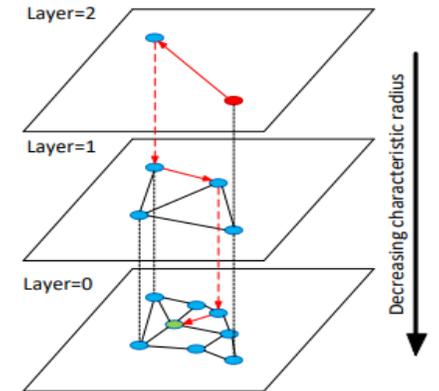
# Story 1: Outlier detection
## Whitebox method



- Store ground layer of HNSW graph.

- Calculate k-NN/LOF-like scores from the graph directly.
  - For a point, consider the neighbours in the graph
  - Sum the distances and divide by the degree of the point.
  - For the LOF-like score, sum the scores of the neighbouring points, and divide by the degree and score of the point itself.



```
4  foreach v ∈ V do
5  │   score[v] ← 1/deg(v) Σ(v,w)∈E dist(v,w).
6  foreach v ∈ V do
7  │   contrast_score[v] ← 1/(score[v]deg(v)) Σ(v,w)∈E score[w].
8  return (score, contrast_score)
```

$$\text{score}[v] \leftarrow \frac{1}{\deg(v)} \sum_{(v,w)\in E} \text{dist}(v,w).$$

$$\text{contrast\_score}[v] \leftarrow \frac{1}{\text{score}[v]\deg(v)} \sum_{(v,w)\in E} \text{score}[w].$$
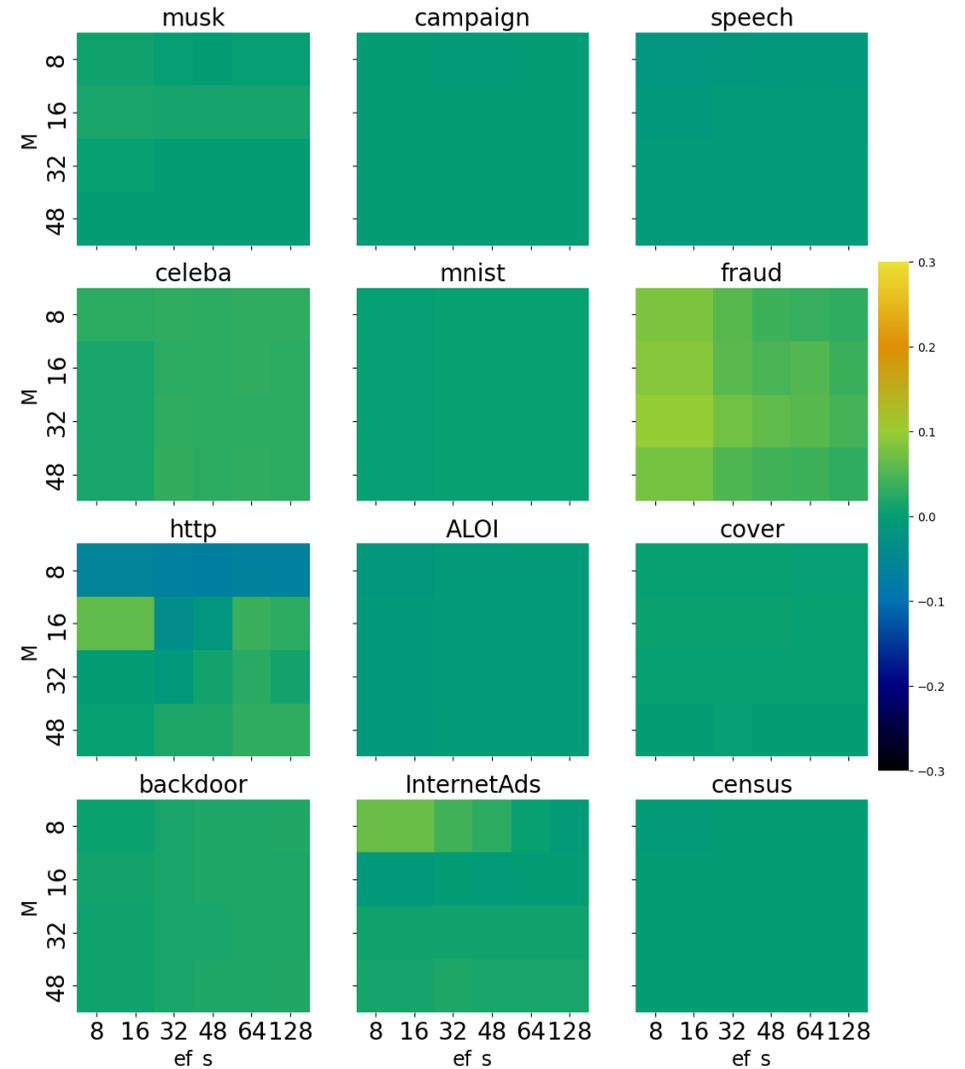
**Linear time!**

# Quality

- AUCROC score comparison between exact LOF and Blackbox LOF.

- *Difference* between exact and new methods shown – positive means the new method achieved higher AUCROC score.

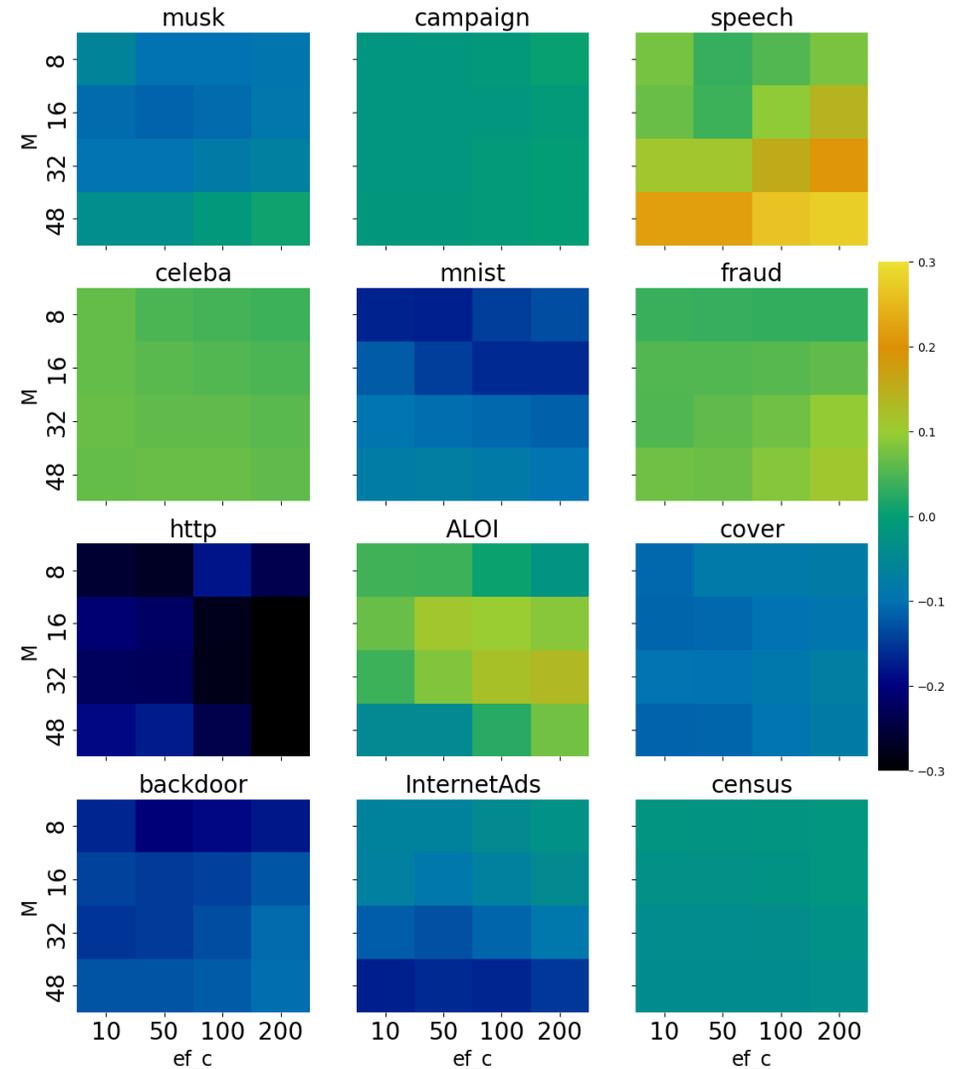- Blackbox got similar if not higher AUCROC in majority of cases.
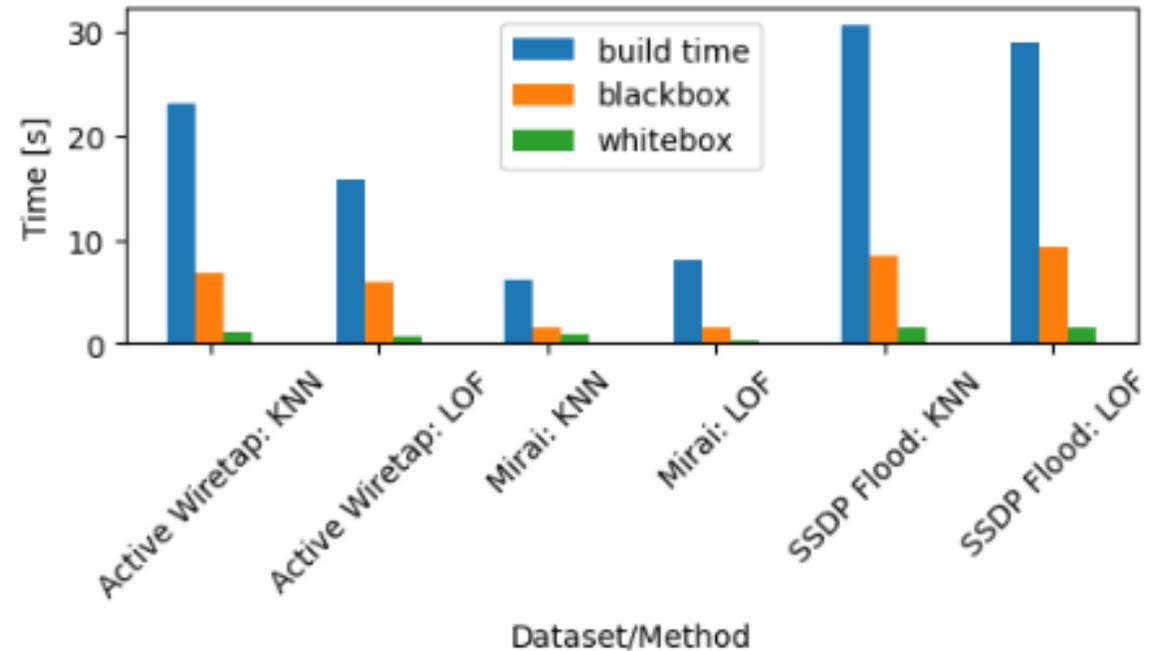


**Blackbox - LOF**

# Quality

- AUCROC score comparison between exact LOF and Whitebox LOF.

- *Difference* between exact and new methods shown – positive means the new method achieved higher AUCROC score.

- Whitebox showed more variance, but most still didn't give an AUCROC less than 0.1 smaller than the exact – for some datasets it even got a higher AUCROC!

# Running times

- Running time distribution for Kitsune datasets.

- Reported *smallest* observed difference between blackbox and whitebox.

- Build time had the most significant contribution.

- Blackbox is about 4x faster than build time, whitebox is 10-30x.
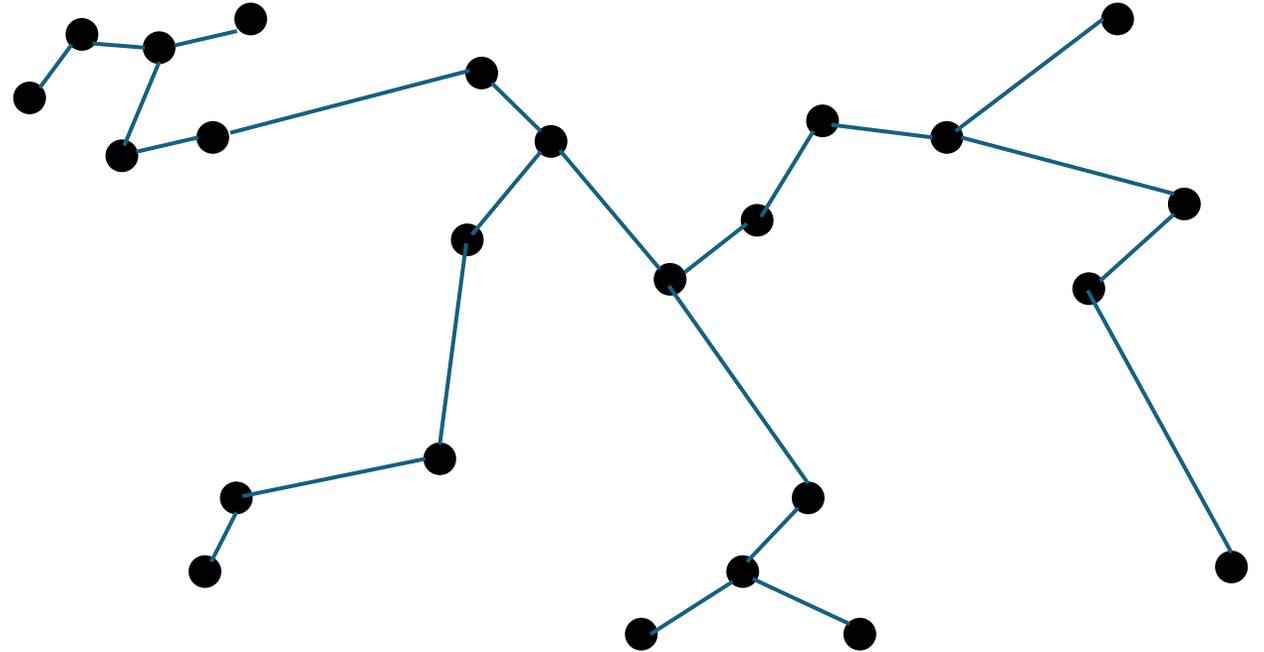


*- Baseline methods cannot finish within several hours*

# Story 2: Single-Linkage Clustering
Computing the MST

- **Want:** Enumerate all edges sorted by distance (Kruskal)

- How to generate these pairs?

- **Problem**:
  - ANN-Search good at returning close neighbors, but we might need far neighbors.

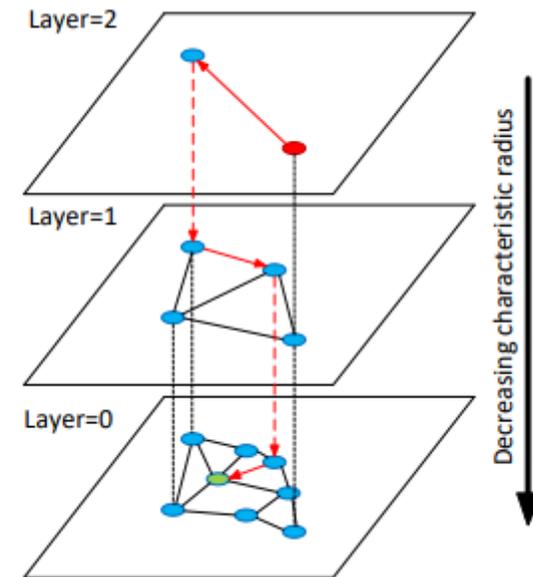# Approaches (on 'flattened out' HNSW graph)

- **MST**
  - MST of dataset = MST of HNSW graph
- **Extend**
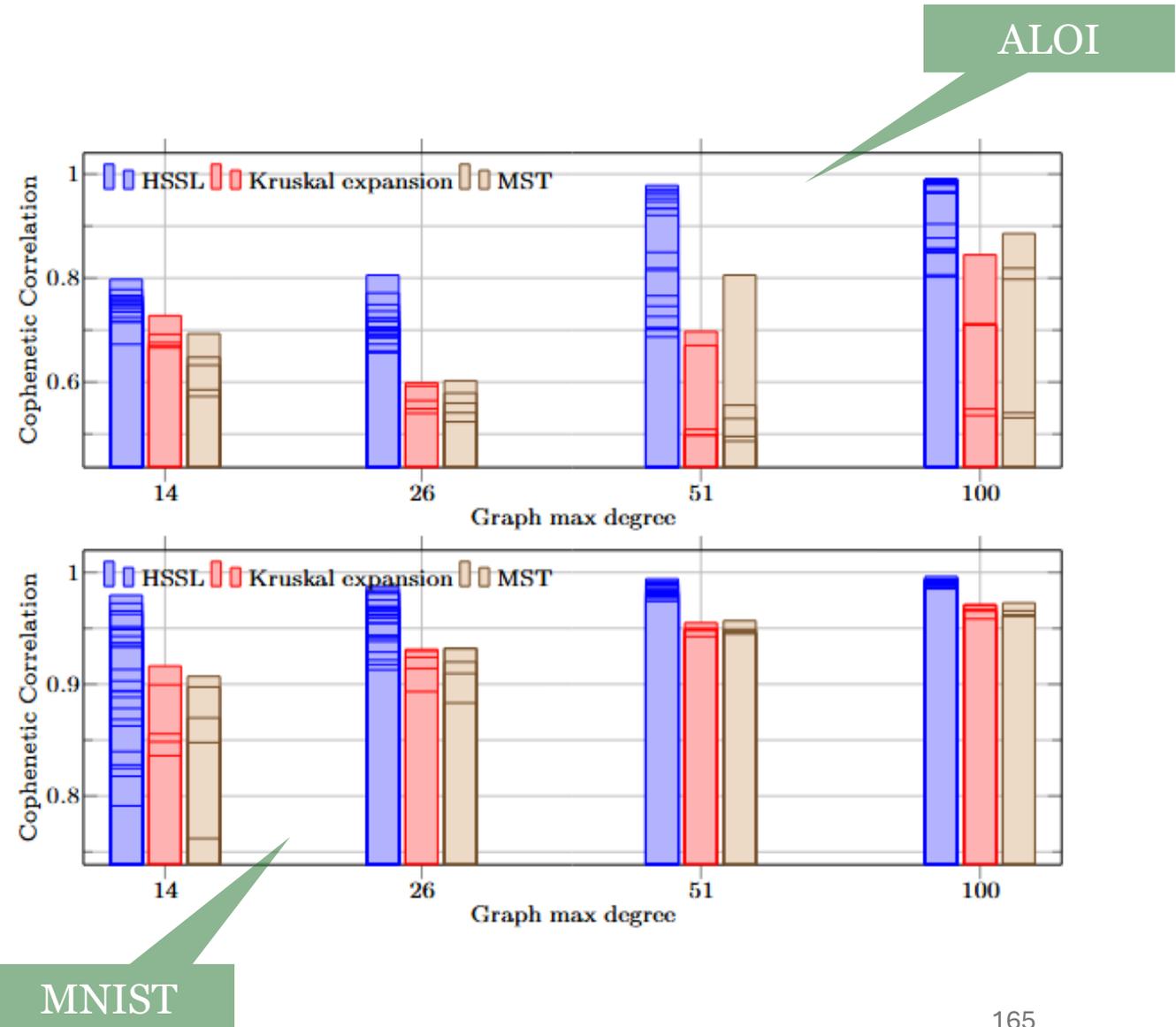  - When edge $(v, w)$ is traversed, add edges between $v$ and neighbors of $w$.
- **HSSL: Incremental Searchers [Schubert, SISAP 2024]**
  - Each point keeps local queue of close neighbors
  - Always make progress with the point that has best candidate
  - Fill up queue if empty (incremental search)



Layer=2

Layer=1

Layer=0

Decreasing characteristic radius

# Quality

- Horizontal lines →
  different parameter setting

- Two datasets: ALOI/MNIST

- ALOI more difficult.

- Mostly ⬆ M ⟶ ⬆ Quality

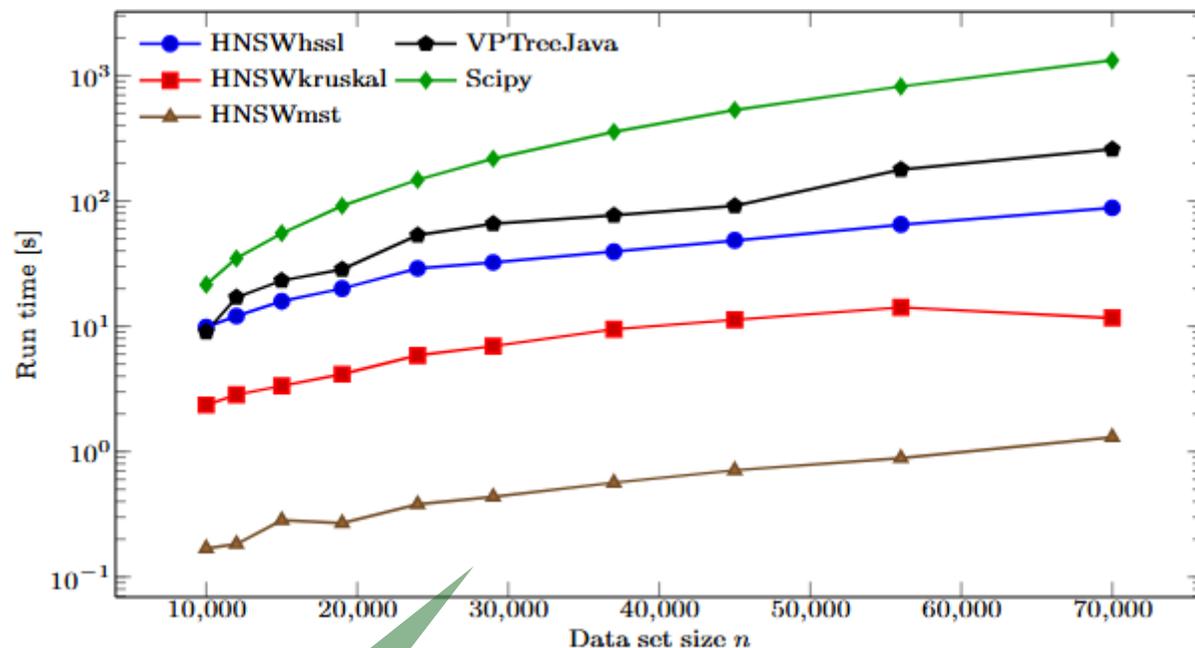- HSSL > MST/Kruskal
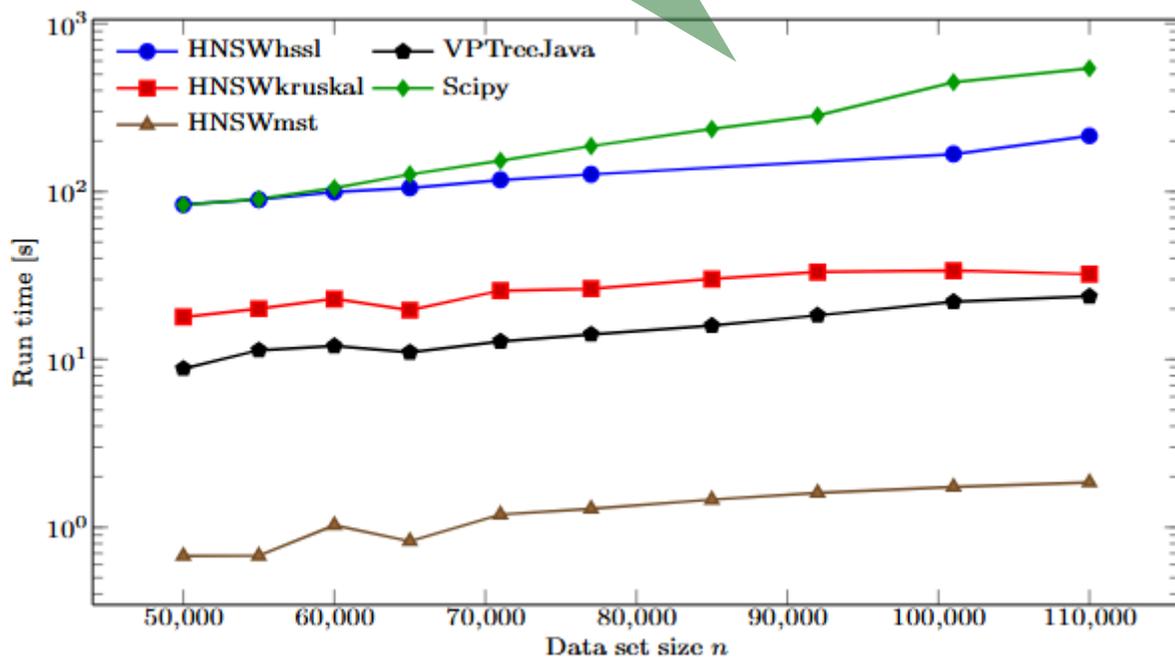
- MST performed comparatively to
  Kruskal.



ALOI

MNIST

165

# Running time scaling

Fastest parameter setting **above a 0.8** cophenetic correlation chosen

**Speedup comparison:**
MST → (~x10) → Kruskal → (~x5) → HSSL→ (often x2 or more) → SciPy

# Improving Data Mining Through ANN

- **Blackbox vs. Whitebox**
  - **Design novel ways** to make creative use of the data structure
  - **New game**: Whole pipeline has to be tuned to downstream task
  - In particular: **Expensive index building** is **infeasible**

- **Whitebox standard** in "hashing-based data mining"
  - Scalable Density-based Clustering with Random Projections, [Xu, Pham, NeurIPS 2024]
  - High-dimensional Density-based Clustering Using Locality-Sensitive Hashing, [Okkels, A., Thomsen, Zimek, EDBT'25]

# What I have missed when starting this topic

- **How to evaluate the quality of approximate methods?**
  - What do we compare to?
    - faithfullness to original method
    - clustering quality

- **How important is it to get accurate results?**

- **Do embeddings present meaningful mining tasks?**
  - Make for good baselines
  - Where do we get the data from otherwise?

# International Workshop on Data Mining, Visualization, and Search in Very High-Dimensional Spaces

Part 1: Why search high-dimensional spaces?

Part 2: How to search high-dimensional spaces?

Part 3: How to asses high-dimensional search?

Part 4: How to use search to speed up data mining?

Thanks!