# Reconfigurable Inverted Index

Yusuke Matsui[1]    Ryota Hinami[2]    Shin'ichi Satoh[1]
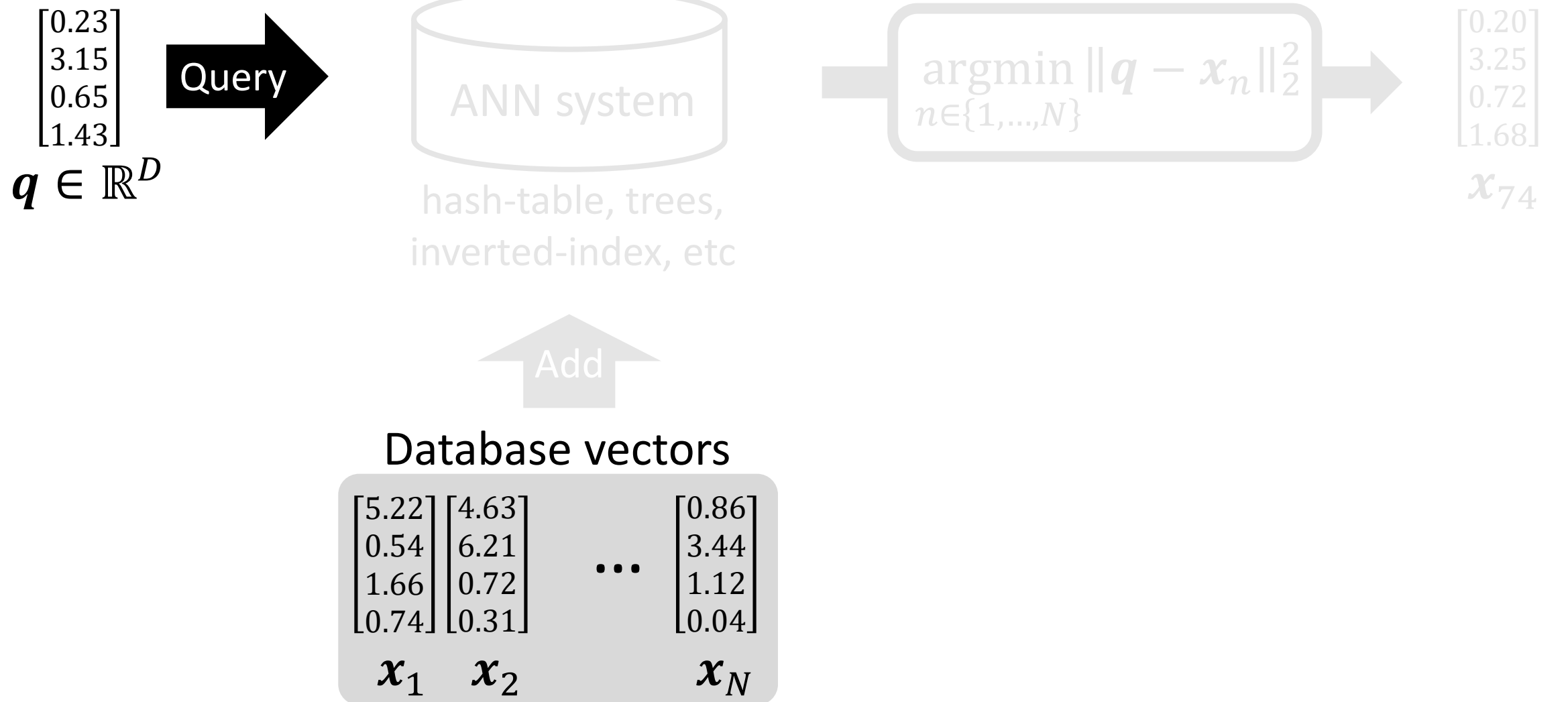
[1]National Institute of Informatics        [2]The University of Tokyo
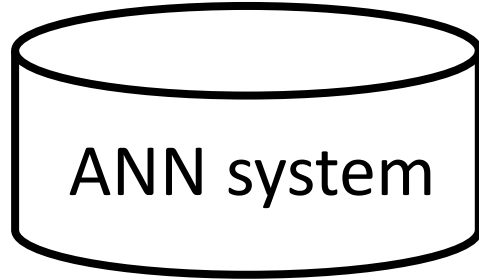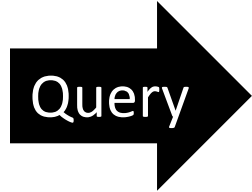
slides: https://bit.ly/2P0KuW1

# **Approximate** nestra neighbor search

# Approximate nearest neighbor search

# **Approximate** nearest neighbor search

$$\boldsymbol{q} = \begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

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

Query →

**ANN system**

hash-table, trees,
inverted-index, etc

↑ Add

## Database vectors

$$\boldsymbol{x}_1 = \begin{bmatrix} 5.22 \\ 0.54 \\ 1.66 \\ 0.74 \end{bmatrix} \quad \boldsymbol{x}_2 = \begin{bmatrix} 4.63 \\ 6.21 \\ 0.72 \\ 0.31 \end{bmatrix} \quad \cdots \quad \boldsymbol{x}_N = \begin{bmatrix} 0.86 \\ 3.44 \\ 1.12 \\ 0.04 \end{bmatrix}$$

Approximate NN search

$$\underset{n \in \{1,\ldots,N\}}{\operatorname{argmin}} \| \boldsymbol{q} - \boldsymbol{x}_n \|_2^2$$

Result

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

$\boldsymbol{x}_{74}$

# **Approximate** nearest neighbor search

Approximate NN search

Result

$$\underset{n \in \{1,\ldots,N\}}{\arg\min} \|\boldsymbol{q} - \boldsymbol{x}_n\|_2^2$$

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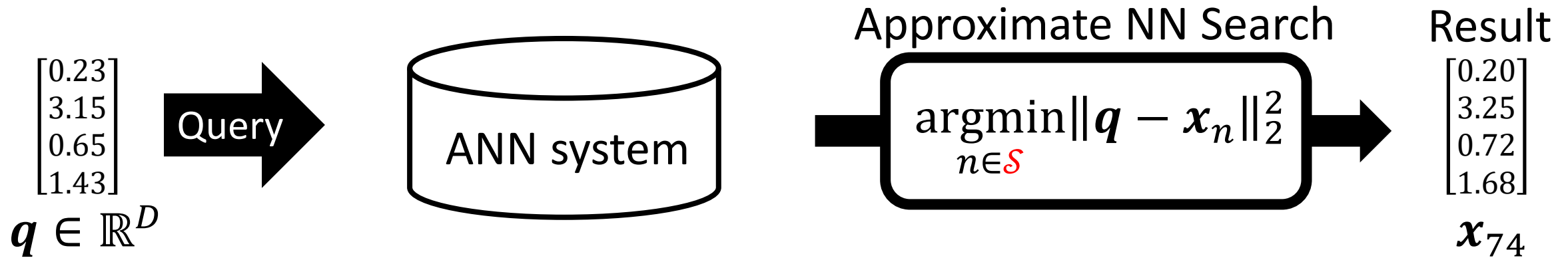
Query

ANN system

hash-table, trees,
inverted-index, etc

Add

Database vectors

# Related work

- Locality-sensitive-hashing (LSH)
  - FALCONN [Andoni+, 15] [Razenshteyn+, 18]
- Project/tree-based
  - FLANN [Muja+, 14]
  - Annoy [Bernhardsson, 18]
- Graph traversal
  - NSW/HNSW on NMSLIB [Malkov+, 16][Boytsov+, 13]
- Product quantization (PQ)
  - IVFPQ on Faiss [Jégou+, 11][Johnson+, 17] etc.
  - ***Our Reconfigurable Inverted Index***

$$\boldsymbol{q} = \begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

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

**Query** →

ANN system

Approximate NN Search

$$\underset{n \in \mathcal{S}}{\mathrm{argmin}} \|\boldsymbol{q} - \boldsymbol{x}_n\|_2^2$$

Result

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

$\boldsymbol{x}_{74}$

# Subset search problem

➢ Existing ANN systems are fast for the all vectors
- Search is over $\mathcal{S} = \{1, \dots, N\}$

➢ However, it is **hard** to run the search for a subset
- Search is over $\mathcal{S} \subseteq \{1, \dots, N\}$
- e.g., searching from $\{\boldsymbol{x}_{1000}, \dots, \boldsymbol{x}_{2000}\}$
- **Why?** Systems are usually optimized for $\mathcal{S} = \{1, \dots, N\}$

# There is a demand for subset search!

# There is a demand for subset search!

Propose: **Reconfigurable inverted index (Rii)**
- ✓ Subset search
- ✓ A comparative performance with IVFPQ (Faiss)
- ✓ 10 ms for billion-scale data

# Reconfigurable inverted index (Rii)

➢Preliminary
- PQ linear scan
- IVFPQ

➢Data structure

➢Search



Fast if $|\mathcal{S}|$ is small



Fast if $|\mathcal{S}|$ is large



Cherry pick!
Always fast

# Reconfigurable inverted index (Rii)

➤ Preliminary
  - PQ linear scan

Fast if $|\mathcal{S}|$ is small

  - IVFPQ

Fast if $|\mathcal{S}|$ is large

➤ Data structure

➤ Search

Cherry pick!
Always fast

# Preliminary: Product quantization (PQ) [Jégou+, TPAMI 11]

**PQ**: Compress a vector into a short code

$$\begin{bmatrix} 5.22 \\ 0.54 \\ 1.66 \\ 0.74 \end{bmatrix}$$

VQ →

VQ →

$$\mathbb{R}^4 \rightarrow \{ \ \blacksquare \ , \ \blacksquare \ , \dots \}^2$$

All database vectors are PQ-encoded beforehand

$x_1$  $x_2$  $x_N$

$$\begin{bmatrix} 5.22 \\ 0.54 \\ 1.66 \\ 0.74 \end{bmatrix} \begin{bmatrix} 4.63 \\ 6.21 \\ 0.72 \\ 0.31 \end{bmatrix} \dots \begin{bmatrix} 0.86 \\ 3.44 \\ 1.12 \\ 0.04 \end{bmatrix}$$

PQ  PQ  PQ

① ② Ⓝ

...

# **Preliminary: Product quantization (PQ)** [Jégou+, TPAMI 11]

➢ The subset search is possible with a linear cost of $|\mathcal{S}|$

$$\operatorname*{argmin}_{n \in \mathcal{S}} d\left(\boldsymbol{q}, \overset{\textstyle\textcircled{\scriptsize n}}{\boxed{\phantom{x}}}\right)$$

e.g., $\mathcal{S} = \{2, 4, 5, 8\}$

$$\boldsymbol{q} = \begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

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

Linearly compared

➢ The search is efficient
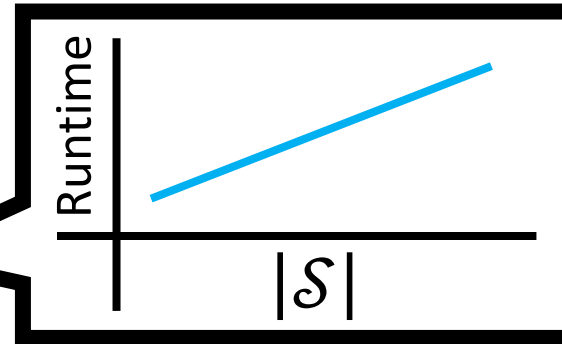  only if $|\mathcal{S}|$ is small

# Reconfigurable inverted index (Rii)
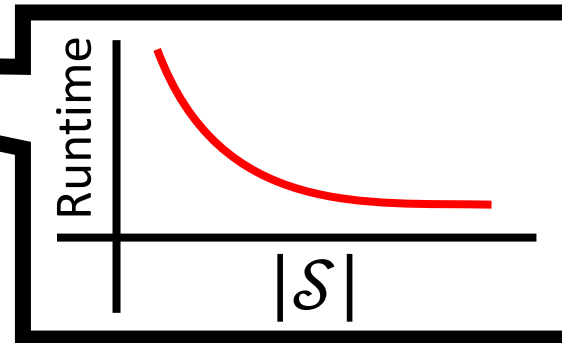
➢Preliminary
- PQ linear scan



Fast if $|\mathcal{S}|$ is small

- IVFPQ



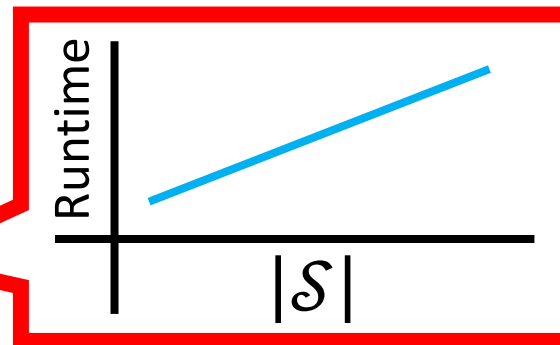Fast if $|\mathcal{S}|$ is large

➢Data structure

➢Search



Cherry pick!
Always fast

➢Evaluation

# Preliminary: Inverted Index + PQ (IVFPQ) [Jégou+, TPAMI 11]

➤ Current basic data structure for a large-scale search

➤ Subset-search is possible <span style="color:red">only if $|\mathcal{S}|$ is large</span>



Space partitioning

# Preliminary: Inverted Index + PQ (IVFPQ) [Jégou+, TPAMI 11]

➢ Current basic data structure for a large-scale search

➢ Subset-search is possible only if $|\mathcal{S}|$ is large

e.g., $\mathcal{S} = \{13, 92, 105, \dots\}$

$\widehat{n} \in \mathcal{S}$ or not



$$\boldsymbol{q} = \begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$

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

Space partitioning

Re-rank via PQ-linear scan

1. Find the closest space: $k^* = \operatorname{argmin}_k \|\boldsymbol{q} - \boldsymbol{c}_k\|_2^2$
2. Focus the $k^*$th space, accept items $\in \mathcal{S}$
3. Re-rank the items via PQ-linear scan

# Preliminary: Inverted Index + PQ (IVFPQ) [Jégou+, TPAMI 11]

➢ Current basic data structure for a large-scale search

➢ Subset-search is possible only if $|\mathcal{S}|$ is large

e.g., $\mathcal{S} = \{13, 92, 105, \dots\}$

$\text{\textcircled{n}} \in \mathcal{S}$ or not

$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$
$q \in \mathbb{R}^D$

126    225

13    92    188

✓    ✓

Re-rank via
PQ-linear scan

92

Space partitioning

1. Find the closest space: $k = \arg\min_k \|q - c_k\|_2^2$

2. Focus the $k$-th space, accept items $\in \mathcal{S}$

3. Re-rank the items via PQ-linear scan

**Why is it slow for small $|\mathcal{S}|$?**
e.g., if $|\mathcal{S}|$ is small and they are far away from the query, we might need to scan all items

Runtime

$|\mathcal{S}|$    $N$

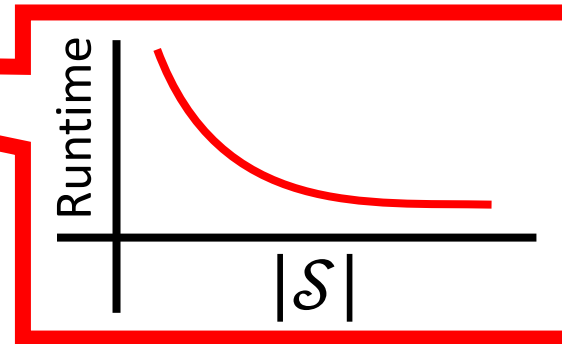# Reconfigurable inverted index (Rii)

➢Preliminary
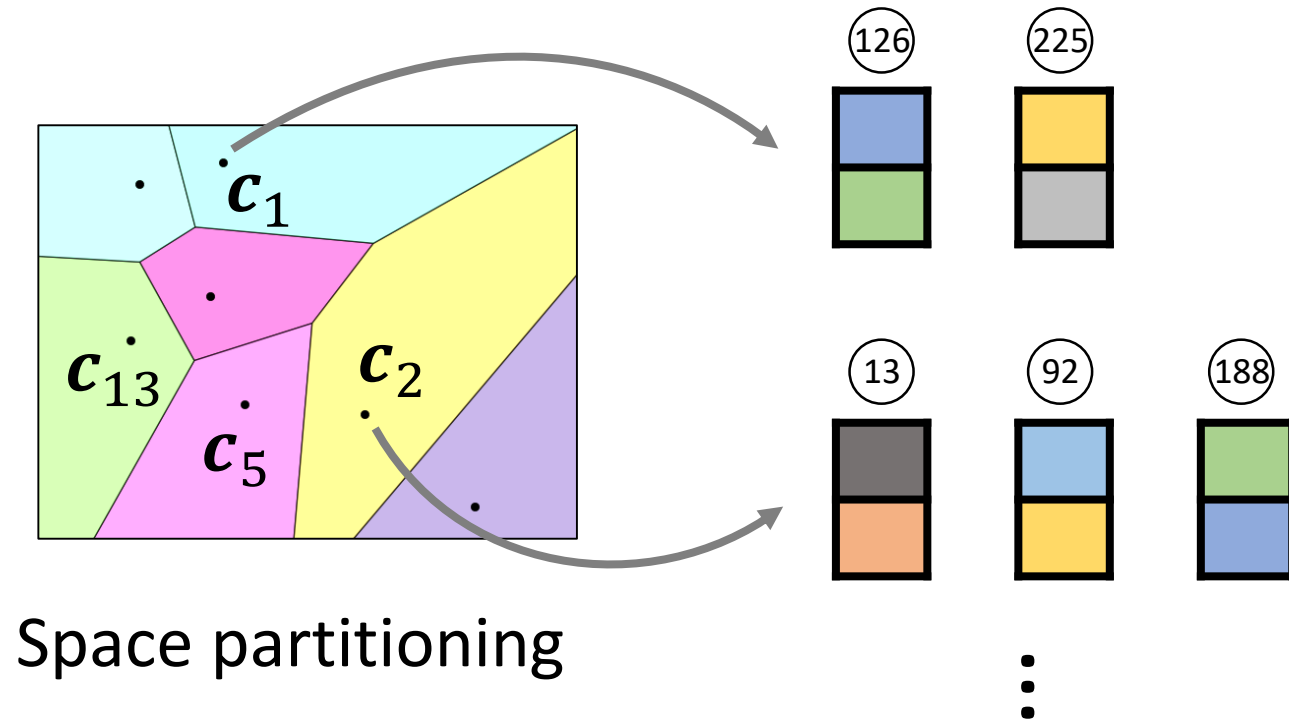  - PQ linear scan



Fast if $|\mathcal{S}|$ is small

  - IVFPQ



Fast if $|\mathcal{S}|$ is large

➢Data structure

➢Search



Cherry pick! Always fast

# Data structure

➢ Store (1) PQ-codes **linearly**, and (2) IDs as an inverted index
➢ Can run either PQ-linear-scan or IVFPQ with a **single data structure**

Key: store codes linearly

### cf. IVFPQ

➢ PQ-codes are also chunked. Natural
➢ Slight, but critical change

# Reconfigurable inverted index (Rii)

➢Preliminary
  - PQ linear scan
  - IVFPQ

➢Data structure

➢<span style="color:red">Search</span>



Fast if $|\mathcal{S}|$ is small

Fast if $|\mathcal{S}|$ is large

Cherry pick! Always fast

# Search

➢ If $|\mathcal{S}|$ is small, run PQ-linear scan

➢ If $|\mathcal{S}|$ is large, run IVFPQ

# Search

➢If $|\mathcal{S}|$ is small, run PQ-linear scan

➢If $|\mathcal{S}|$ is large, run IVFPQ

# Search

➢If $|\mathcal{S}|$ is small, run PQ-linear scan

➢If $|\mathcal{S}|$ is large, run IVFPQ

# Search

➢If $|\mathcal{S}|$ is small, run PQ-linear scan
➢If $|\mathcal{S}|$ is large, run IVFPQ

# Search

➢If $|\mathcal{S}|$ is small, run PQ-linear scan
➢If $|\mathcal{S}|$ is large, run IVFPQ

➢Set a threshold $\theta$
➢Key: Switch two methods based on $|\mathcal{S}| \lessgtr \theta$



Use PQ-linear-scan

Use IVFPQ

# Search

➢If $|\mathcal{S}|$ is small, run PQ-linear scan
➢If $|\mathcal{S}|$ is large, run IVFPQ

➢Set a threshold $\theta$
➢Key: Switch two methods based on $|\mathcal{S}| \lessgtr \theta$

# Evaluation

➢ SIFT1M ($N = 10^6, D = 128$). Results for top-R search

# Evaluation

> Existing system: Annoy
> Force to search a subset

> SIFT1M ($N = 10^6, D = 128$). Results for top-R search



The existing system is slow, especially when $|\mathcal{S}|$ is small

Proposed Rii is always fast regardless of $|\mathcal{S}|$ and $R$

Legend:
- Annoy+PC: $R = 1$
- Annoy+PC: $R = 10$
- Annoy+PC: $R = 100$
- Rii: $R = 1$
- Rii: $R = 10$
- Rii: $R = 100$

Y-axis: Runtime per query [msec]
X-axis: $|\mathcal{S}|$

```
$ pip install rii
```

https://github.com/matsui528/rii

```python
import rii
import nanopq

# Prepare a PQ/OPQ codec with M=32 sub spaces
codec = nanopq.PQ(M=32).fit(vecs=Xt)  # Trained using Xt

# Instantiate a Rii class with the codec
e = rii.Rii(fine_quantizer=codec)

# Add vectors
e.add_configure(vecs=X)

# Search
ids, dists = e.query(q=q, topk=3, target_ids=S)
print(ids, dists)  # e.g., [7484 8173 1556] [15.062 15.385 16.169]
```
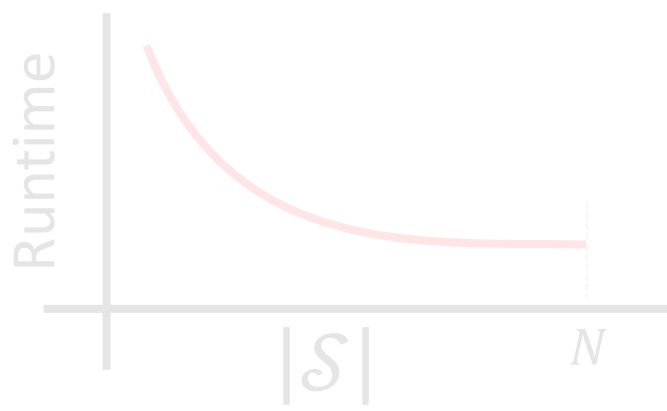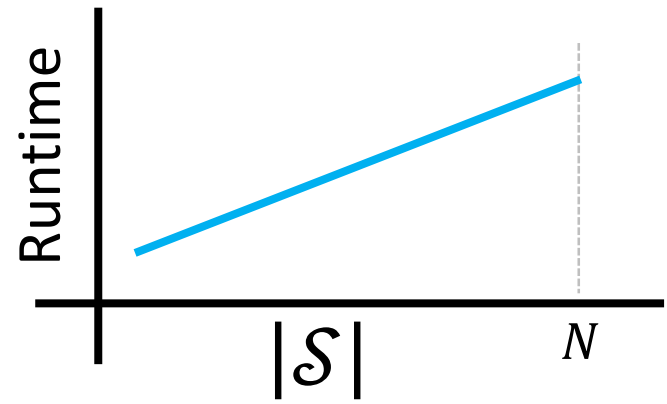
# Summary

Approximate NN Search

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$ **Query** → ANN system → $\underset{n \in \mathcal{S}}{\operatorname{argmin}} \|\boldsymbol{q} - \boldsymbol{x}_n\|_2^2$ → Result

Reconfigurable inverted index:
➤ Store PQ-codes linearly
➤ Switch method based on $|\mathcal{S}|$



Runtime

$\theta$    $N$

$|\mathcal{S}|$

Use PQ-linear-scan    Use IVFPQ

➤ PyPI:

```
$ pip install rii
```

➤ One thing I couldn't mention:

**Reconfiguration**: the system remains fast even after many new items are added

➤ See our paper, or come to our poster:
✓ Poster session 5 (15:30 – 16:30)
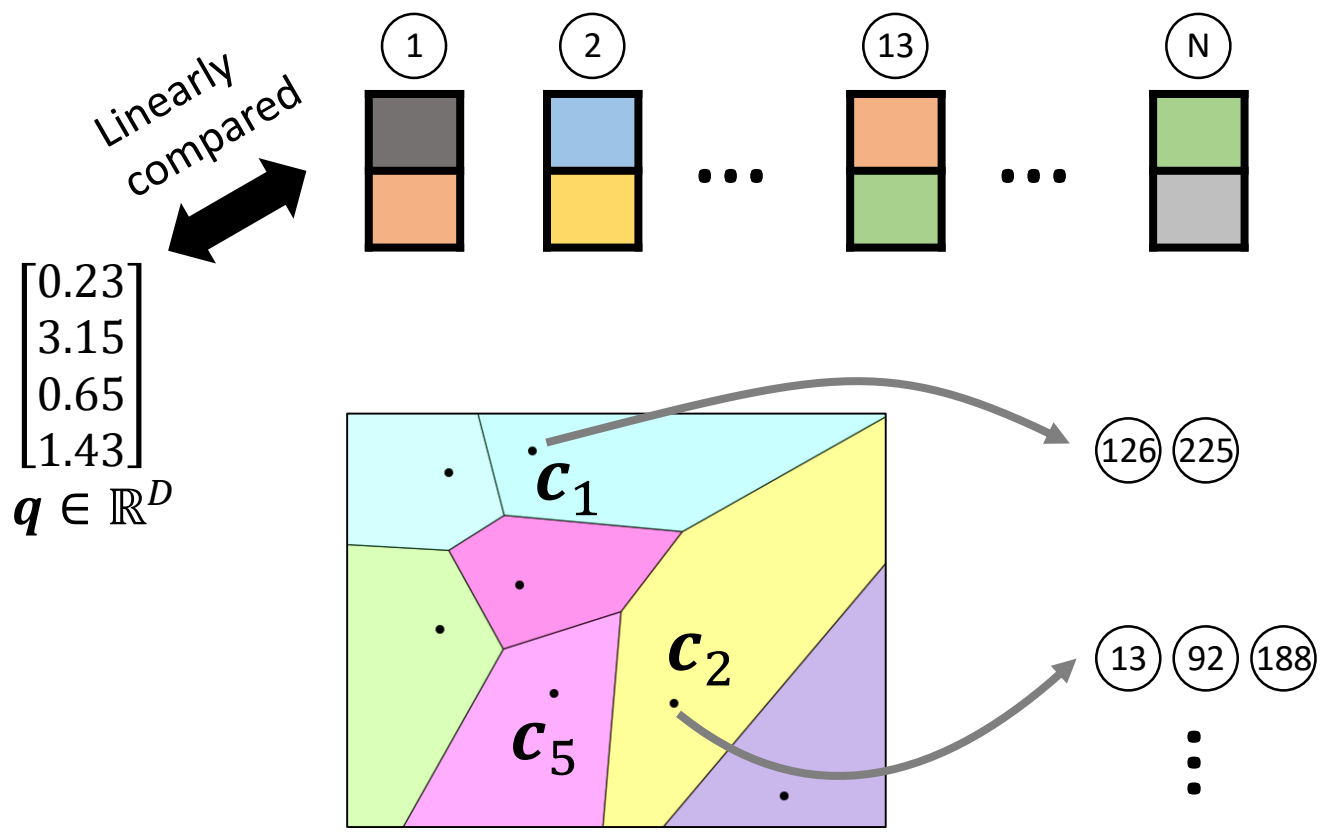
# Example: Image search

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

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

Query

ANN system

Approximate NN search

$$\operatorname*{argmin}_{n \in \{1, \dots, N\}} \|\boldsymbol{q} - \boldsymbol{x}_n\|_2^2$$

Result

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

$\boldsymbol{x}_{74}$

Extract a VGG feature

Query image

$N$ database images

74th image

Example of subset search for image

slides: https://bit.ly/2P0KuW1

# Evaluation

➢Extensive comparison against existing methods

➢For a fixed accuracy (Recall@1), check runtime and its disk space

| Dataset | Method | Parameters | Recall@1 (fixed) | Runtime/query | Disk space | Build time |
|---|---|---|---|---|---|---|
| SIFT1M | Annoy [11] | $n_{\text{trees}} = 2000, k_{\text{search}} = 400$ | 0.67 | 0.18 ms | 1703 MB | 899 sec |
| | FALCONN [1, 41] | $n_{\text{probes}} = 16$ | 0.63 | 0.87 ms | - | **1.8 sec** |
| | NMSLIB (HNSW) [14, 33, 39] | efS = 4 | 0.67 | **0.043 ms** | 669 MB | 436 sec |
| | Faiss (IVFADC) [25, 26] | $K = 10^3, M = 64, n_{\text{probe}} = 4$ | 0.67 | 0.61 ms | 73 MB | 30 sec |
| | Rii (proposed) | $K = 10^3, M = 64, L = 5000$ | 0.64 | 0.73 ms | **69 MB** | 82 sec |
| | Rii-OPQ (proposed) | $K = 10^3, M = 64, L = 5000$ | 0.65 | 0.82 ms | **69 MB** | 85 sec |
| GIST1M | Annoy [11] | $n_{\text{trees}} = 2000, k_{\text{search}} = 2000$ | 0.49 | 1.2 ms | 5023 MB | 2088 sec |
| | FALCONN [1, 41] | $n_{\text{probes}} = 512$ | 0.53 | 8.6 ms | - | **7.2 sec** |
| | NMSLIB (HNSW) [14, 33, 39] | efS = 8 | 0.49 | **0.19 ms** | 3997 MB | 1576 sec |
| | Faiss (IVFADC) [25, 26] | $K = 10^3, M = 240, n_{\text{probe}} = 8$ | 0.52 | 3.8 ms | 253 MB | 51 sec |
| | Rii (proposed) | $K = 10^3, M = 240, L = 8000$ | 0.45 | 3.2 ms | **246 MB** | 353 sec |
| | Rii-OPQ (proposed) | $K = 10^3, M = 240, L = 8000$ | 0.50 | 3.8 ms | 249 MB | 388 sec |

# Evaluation

➢Extensive comparison against existing methods
➢For a fixed accuracy (Recall@1), check runtime and its disk space

| Dataset | Method | Parameters | Recall@1 (fixed) | Runtime/query | Disk space | Build time |
|---------|--------|-----------|------------------|---------------|------------|------------|
| SIFT1M | Annoy [11] | $n_{\text{trees}} = 2000, k_{\text{search}} = 400$ | 0.67 | 0.18 ms | 1703 MB | 899 sec |
| | FALCONN [1, 41] | $n_{\text{probes}} = 16$ | 0.63 | 0.87 ms | - | **1.8 sec** |
| | NMSLIB (HNSW) [14, 33, 39] | efS = 4 | 0.67 | **0.043 ms** | 669 MB | 436 sec |
| | Faiss (IVFADC) [25, 26] | $K = 10^3, M = 64$, probe = 4 | 0.67 | 0.61 ms | 73 MB | 30 sec |
| | Rii (proposed) | $K = 10^3, M$ = 5000 | 0.64 | 0.73 ms | **69 MB** | 82 sec |
| | Rii-OPQ (proposed) | $K = 10^3$ = 5000 | 0.65 | 0.82 ms | **69 MB** | 85 sec |
| | Annoy [11] | = 2000 | 0.49 | 1.2 ms | 5023 MB | 2088 sec |
| | | | 0.53 | 8.6 ms | - | **7.2 sec** |
| | | | 0.49 | **0.19 ms** | 3997 MB | 1576 sec |
| | | | 0.52 | 3.8 ms | 253 MB | 51 sec |
| | | | 0.45 | 3.2 ms | **246 MB** | 353 sec |
| | Rii-OPQ (proposed) | $K = 10^3, M = 240, L = 8000$ | 0.50 | 3.8 ms | 249 MB | 388 sec |

NMSLIB is extremely fast, but consume relatively large disk space (~memory)

# Evaluation

➤ Extensive comparison against existing methods

➤ For a fixed accuracy (Recall@1), check runtime and its disk space

| Dataset | Method | Parameters | Recall@1 (fixed) | Runtime/query | Disk space | Build time |
|---------|--------|------------|------------------|---------------|------------|------------|
| SIFT1M | Annoy [11] | $n_{trees} = 2000, k_{search} = 400$ | 0.67 | 0.18 ms | 1703 MB | 899 sec |
| | FALCONN [1, 41] | $n_{probes} = 16$ | 0.63 | 0.87 ms | - | **1.8 sec** |
| | NMSLIB (HNSW) [14, 33, 39] | efS = 4 | 0.67 | **0.043 ms** | 669 MB | 436 sec |
| | Faiss (IVFADC) [25, 26] | $K = 10^3, M = 64$, probe = 4 | 0.67 | 0.61 ms | 73 MB | 30 sec |
| | Rii (proposed) | $K = 10^3, M$ = 5000 | 0.64 | 0.73 ms | **69 MB** | 82 sec |
| | Rii-OPQ (proposed) | $K = 10^3$ = 5000 | 0.65 | 0.82 ms | **69 MB** | 85 sec |
| | Annoy [11] | = 2000 | 0.49 | 1.2 ms | 023 MB | 2088 sec |
| | | | 0.53 | 8.6 ms | | **7.2 sec** |
| | | | 0.49 | **0.19 ms** | MB | 1576 sec |
| | | | 0.52 | 3.8 ms | | 51 sec |
| | | | 0.45 | 3.2 ms | | 353 sec |
| | Rii-OPQ (proposed) | $K = 10^3, M = 240, L =$ | | | | |

NMSLIB is extremely fast, but consume relatively large disk space (~memory)

Proposed Rii achieved a comparative performance with Faiss