
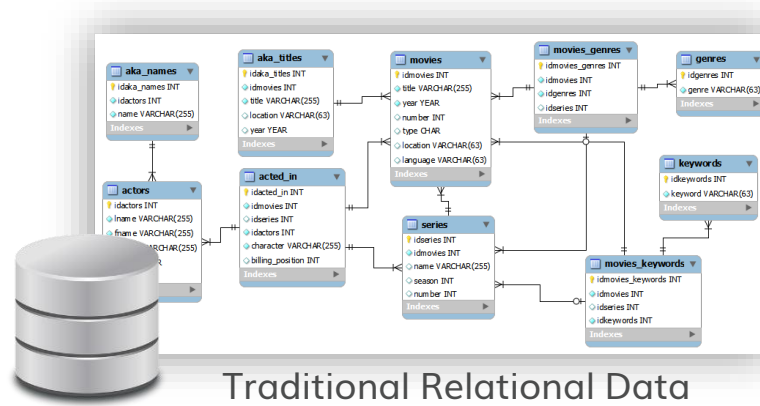





Fast Approximated Nearest Neighbor Joins For Relational Database Systems

Michael Günther, Maik Thiele, and Wolfgang Lehner BTW 2019, 07.03.2019

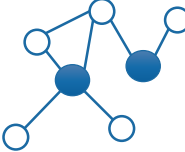
Motivation



Unstructured Text

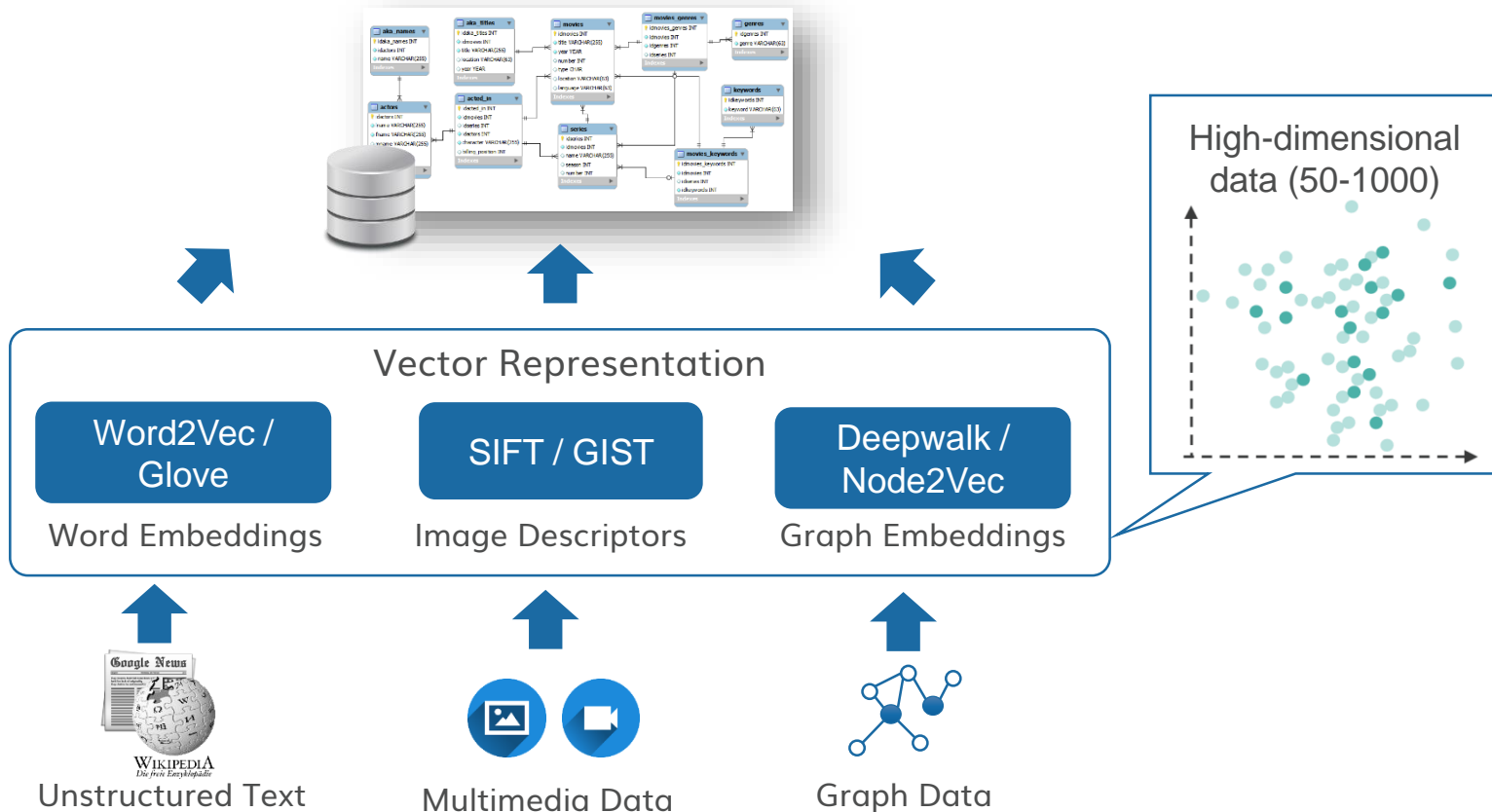


Multimedia Data



Graph Data

Motivation



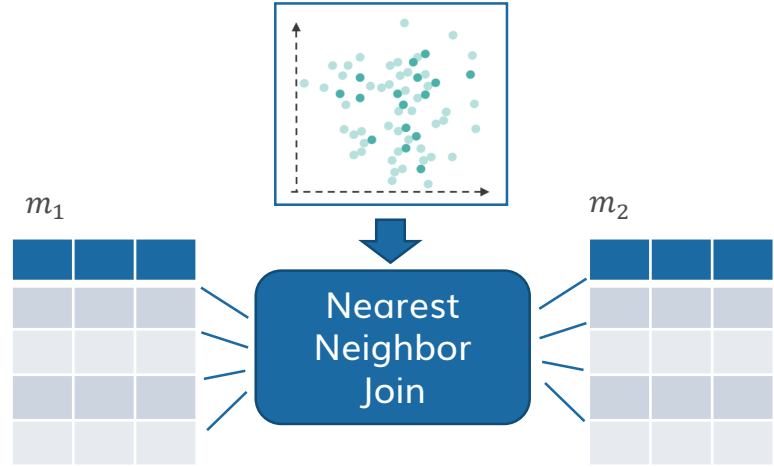
Motivation (2)

Word Embedding Operations within SQL

Movies Relation

m_id	title	...
1	Inception	...
2	Godfather	...
3	Forrest Gump	...
...		
10 ⁶	Scarface	

$kNN (movies.title AS m1 \bowtie movies.title AS m2)$



m1	m2
Inception	Shutter Island
...	...

Bulk operation

→ millions of vector distance calculations

Key Operation: kNN

queries

```
q1: SELECT kNN (movies.title, 3)  
FROM movies
```

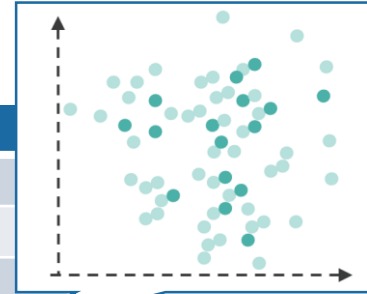
database

m_id	title	...
1	Inception	...
2	Godfather	...
3	Forrest Gump	...
...
10 ⁶	Scarface	...



+

term	vector
2010	[0,21; 0,58; ...; -0,77]
Brando	[-0,46; 0,25; ...; 0,44]
Connery	[0,76; 0,48; ...; -0,51]
...	...
Untouchables	[0,86; -0,22; ...; 0,12]



results

title	VARCHAR
Inception	2010
Inception	Shutter Island
...	...
Godfather	Brando

3 Million Vectors

query set R

Key Operation: kNN-Join

queries

```

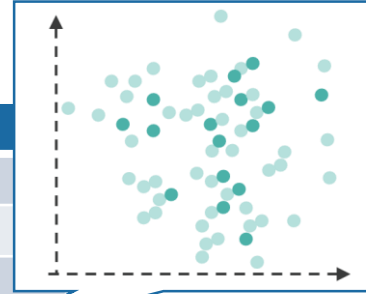
q2: SELECT m1.title, m2.title
      FROM movies m1 kNN-Join(3) movies m2
      ON m1.title ~ m2.title
      WHERE m1.year= ... AND m2.genre IN (...)
  
```

database

m_id	title	...
1	Inception	...
2	Godfather	...
3	Forrest Gump	...
...
10 ⁵	Scarface	...

+

term	vector
Shutter Island	[-0,46; 0,12; ...; 0,87]
Forrest Gump	[0,16; 0,62; ...; -0,19]
The Dark Knight	[0,76; 0,48; ...; -0,51]
...	...
Untouchables	[0,86; -0,22; ...; 0,12]



results

title	VARCHAR
Inception	Shutter Island
Inception	The Dark Knight
...	...
Scarface	Untouchables

100.000 Vectors

query set R

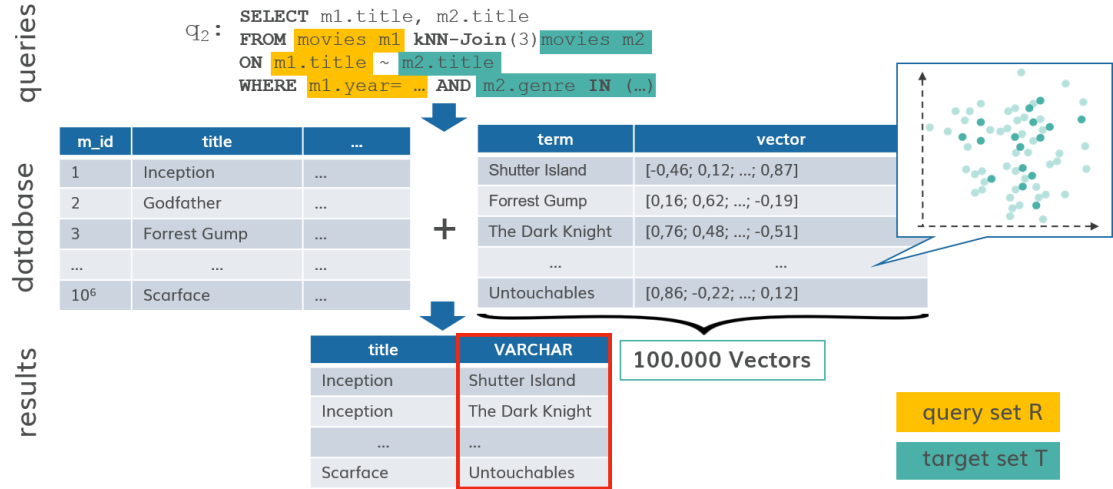
target set T

$$kNN(R \bowtie T) = \{(r, t) \mid t \in kNN(r, T), r \in R\}$$

kNN-Joins in RDBMS

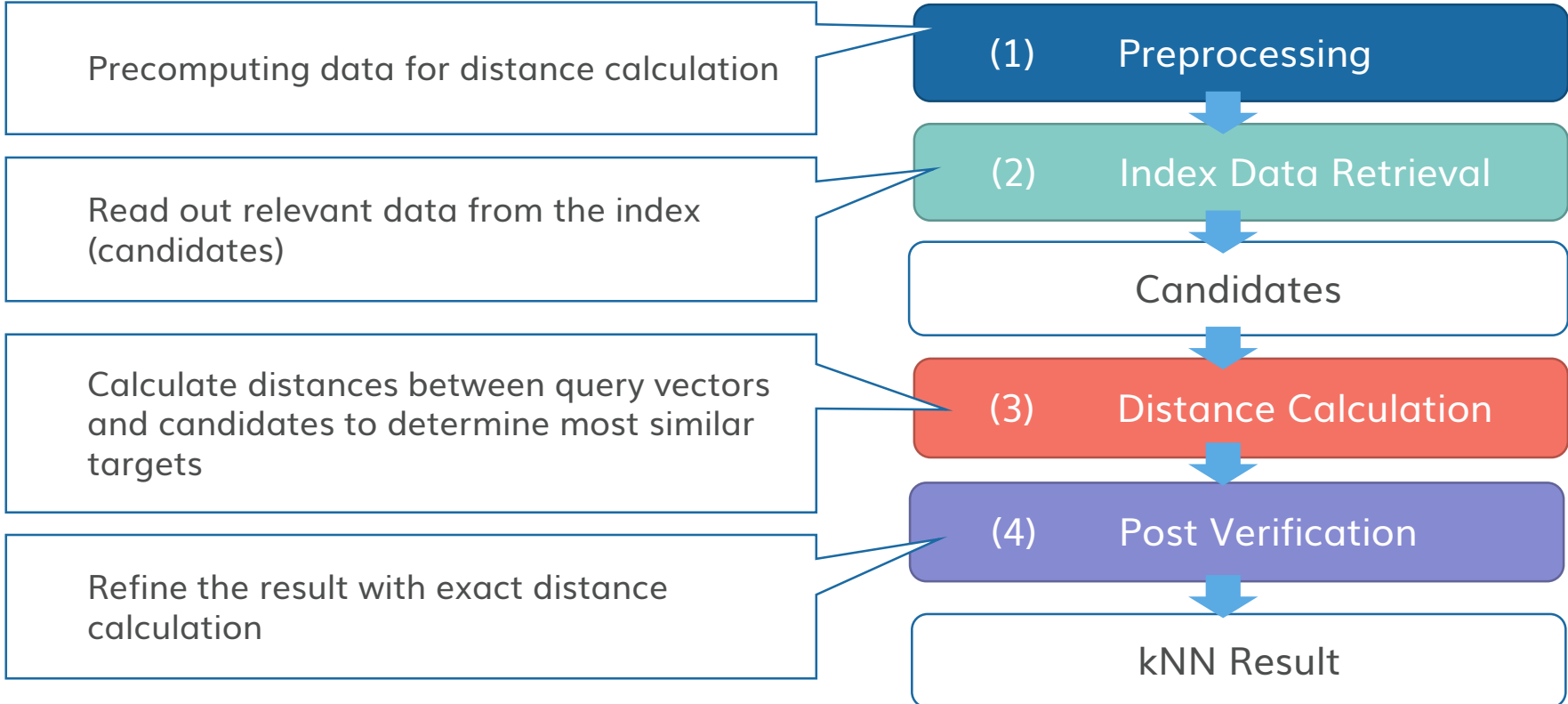
Challenges

- **Batch-wise kNN Search** for large query sets
→ reduce interface and retrieval times
- **High Dimensional Data**
→ Previous Work mainly focus on low dimensional data (e.g. spatial data)
- **Adaptive kNN-Join Algorithm**
→ Different cardinalities of join operands
→ Only one index for all vectors
- **Different Demands on Precision and Response Time**



$$kNN(R \bowtie T) = \{(r, t) \mid t \in kNN(r, T), r \in R\}$$

kNN-Join Algorithm



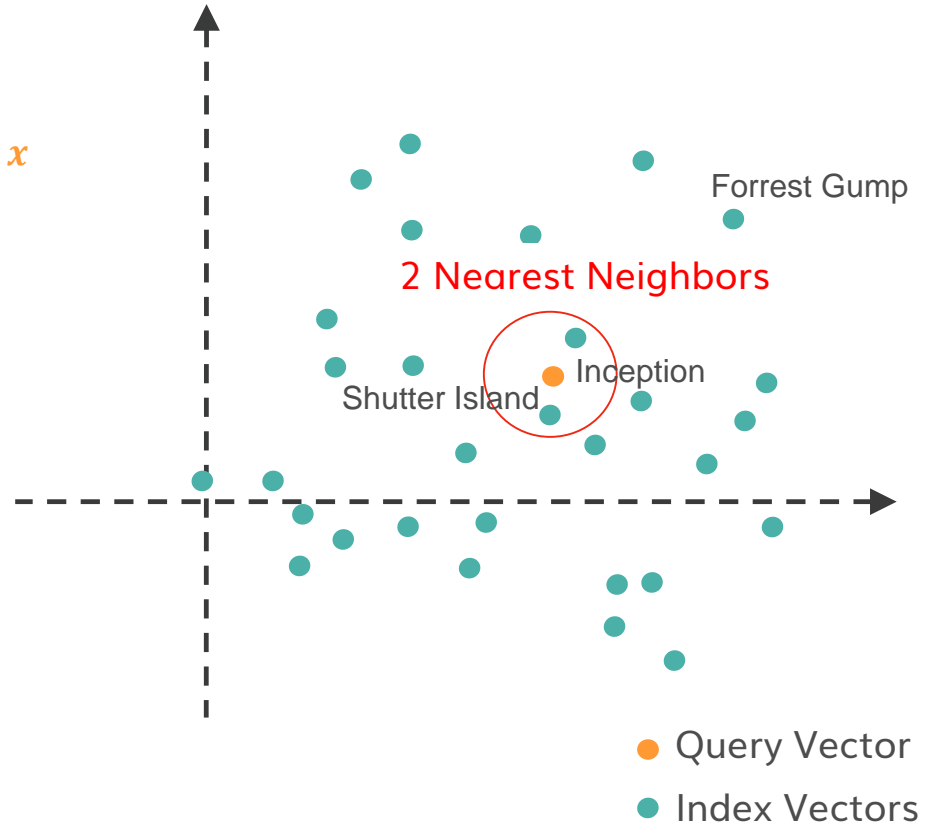


Approximated kNN Search

K Nearest Neighbors

Naïve Algorithm

- Determine all distances between **query vector x** and **vectors $p \in P$** in the target set
 - Select vectors with lowest distances (kNN)
- Complexity: $O(|P| \cdot D)$



Vector Quantization

Objective

- Transform vector data into **compact representation**

Quantization

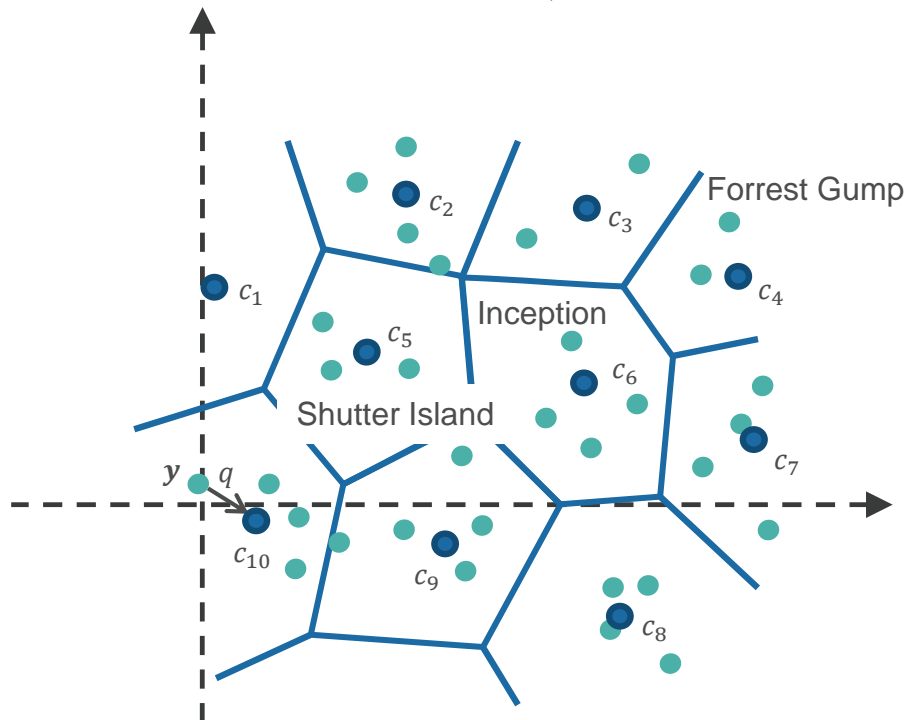
- Cluster Data: kMeans
- Objective:** Obtain a set of centroids which minimize the distortion:
→ Minimize $d(\mathbf{y}, q(\mathbf{y}))$

Quantization Function

- Def.: $q: \mathbf{y} \in \mathbb{R}^D \rightarrow \mathcal{C}, \mathcal{C} \subset \mathbb{R}^D$

Assign every vector to its nearest neighbor out of a set of centroids

Represent vector by centroid id



$$q(\mathbf{y}) = c_{10} \rightarrow \text{id: 10}$$

Product Quantization Search

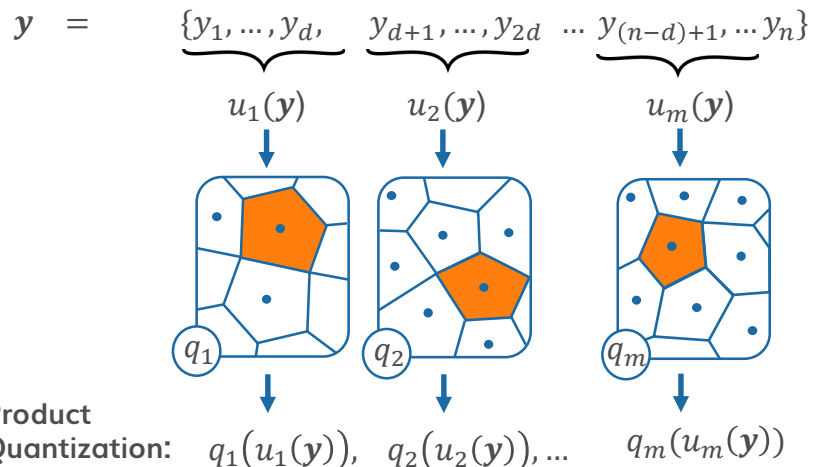
Fast Computation of kNN via Pre-Processing

- Approximated Distances

$$\hat{d}(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_j d(u_j(\mathbf{x}), q_j(u_j(\mathbf{y})))^2}$$

- Precompute

$$d(\underbrace{u_j(\mathbf{x})}_{\text{Query Sub Vector}}, \underbrace{q_j(u_j(\mathbf{y}))}_{\text{Centroid of Quantizer}})^2$$



→ Computation of square distances resolves to a sum of m precomputed distances

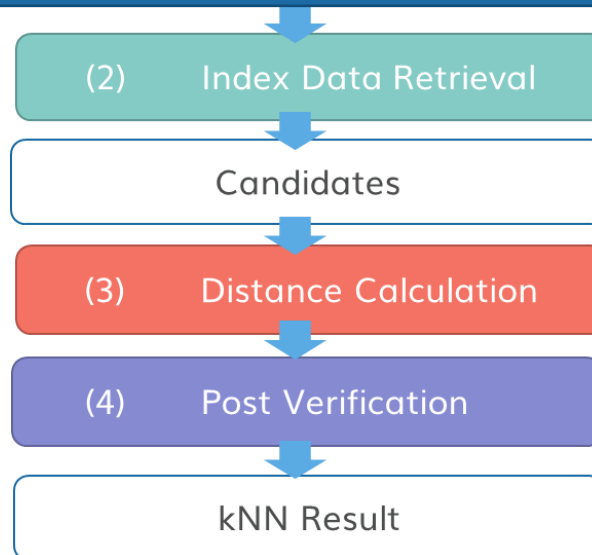
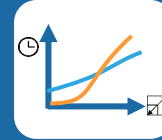


Optimized kNN-Join

Product Quantization Search

(1) Preprocessing can be time consuming

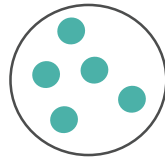
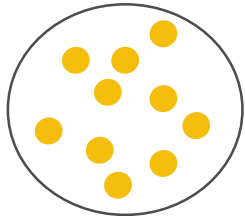
- Should be adaptive to cardinalities of query and target set



Adaptive Preprocessing

QUERY:

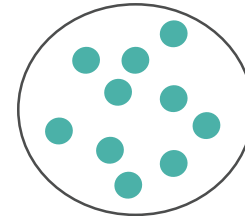
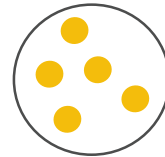
```
SELECT m1.title, m2.title  
FROM movies m1 kNN-Join(3) movies m2  
ON m1.title ~ m2.title  
WHERE m1.year= ... AND m2.genre IN (...)
```



Large Query Set

Small Target Set

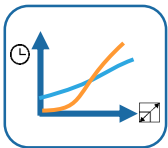
→ Objective: Safe preprocessing time



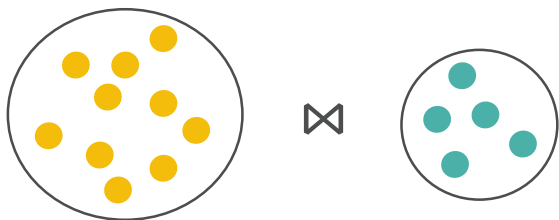
Small Query Set

Large Target Set

→ Objective: Safe time during search



Adaptive Preprocessing



Large Query Set Small Target Set

→ **Objective:** Safe preprocessing time

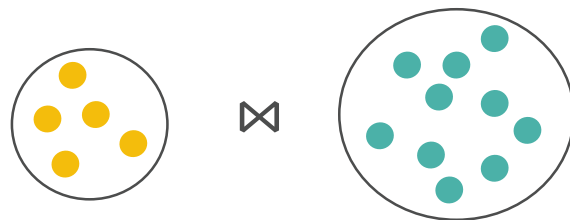


few centroids **but** many quantizer functions

LONG Codes

Number of centroids: 8×10
PQ-Sequences e.g.: 1, 7, 0, 5, 6, 2, 6, 1, 1, 4

Fast preprocessing
(80 distance calculations per query)
Slower search
(10 additions per query-target-pair)



Small Query Set Large Target Set

→ **Objective:** Safe time during search

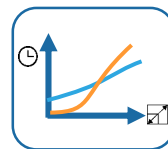


many centroids **but** only few quantizer functions

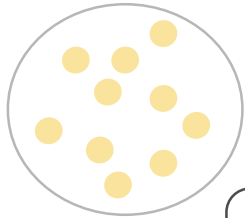
SHORT Codes

Number of centroids: 64×5
PQ-Sequences e.g.: 15, 5, 50, 49, 12

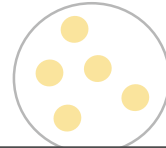
Slow preprocessing
(320 distance calculations per query)
Faster search
(5 additions per query-target-pair)



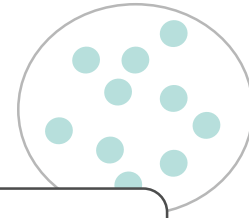
Adaptive Preprocessing



×



×



Large Query Set

Target Set

→ Objective:

search

It is only a single index necessary!

few centroids but

antizer functions

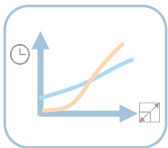
LONG Codes

Number of centroids:
PQ-Sequences e.g.: 1, 7, 0, 5, 6, 2, 6, 1, 1, 4

Fast preprocessing
(80 distance calculations per query)
Slower search
(10 additions per query-target-pair)

PQ-Sequences e.g.: 15, 5, 50, 49, 12

Slow preprocessing
(320 distance calculations per query)
Faster search
(5 additions per query-target-pair)



Adaptive Preprocessing (2)

Evaluation Setup

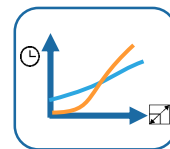
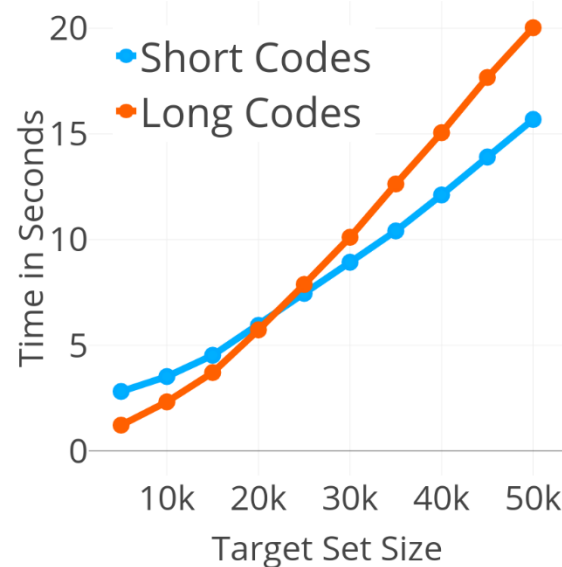
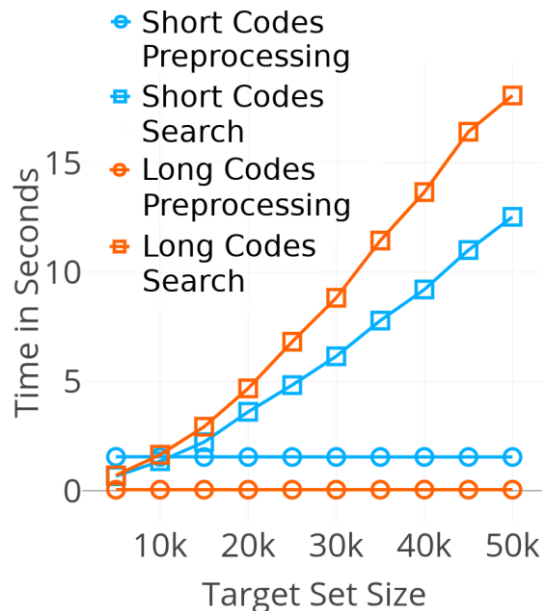
- kNN-Joins with different target set sizes
- Vectors are sampled randomly
- Measurements are done 10 times for each point with different vectors

Dataset

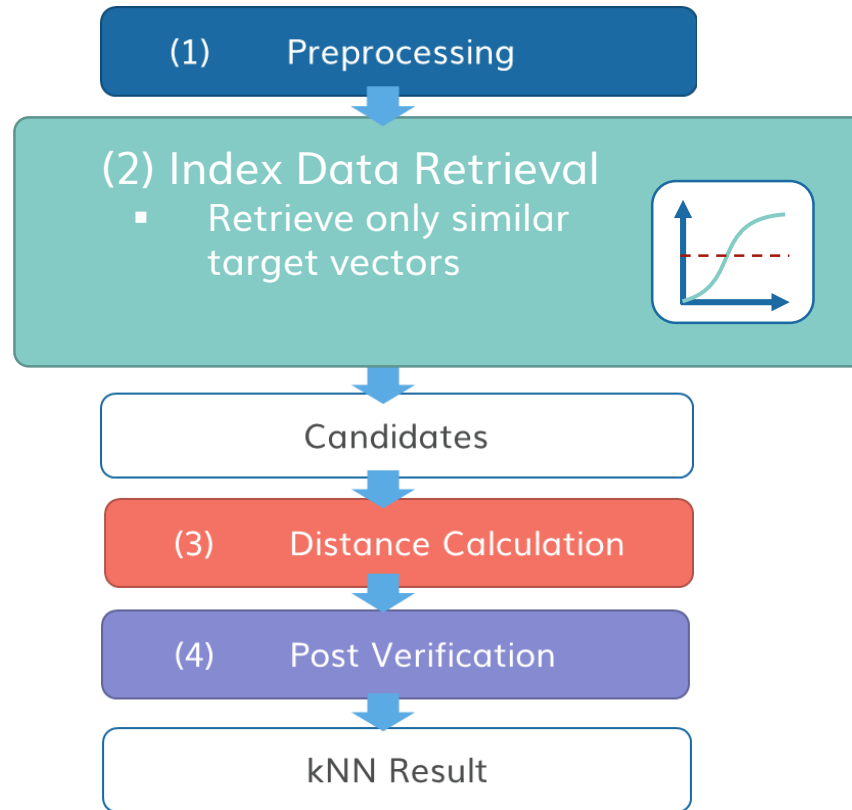
- 3M 300-dimensional word vectors

Parameters

- Query Set Size: 5,000
- K: 5



kNN-Join Algorithmus



Inverted Indexing

Idea

Accelerate kNN-Search by **non-exhaustive** search schema

Preprocessing

Additional *coarse quantizer* q_c is applied to the complete vectors

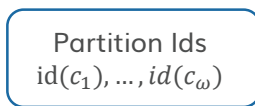
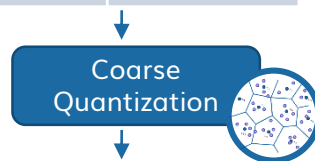
→ Assigns vectors to partitions

kNN Calculation

- Coarse quantization of query vector is calculated → $id(c_1), \dots, id(c_\omega)$
..... ω determine number of candidates
- Partitions with $PID \in \{id(c_1), \dots, id(c_\omega)\}$ are retrieved
- Distances Calculation → kNN

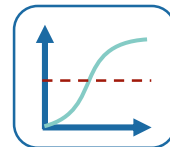
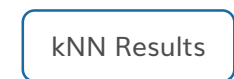
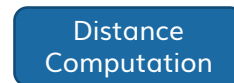
Query

Token	Vector
Inception	[0,21; -0,86; ...]

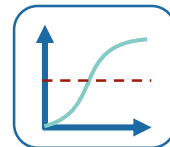
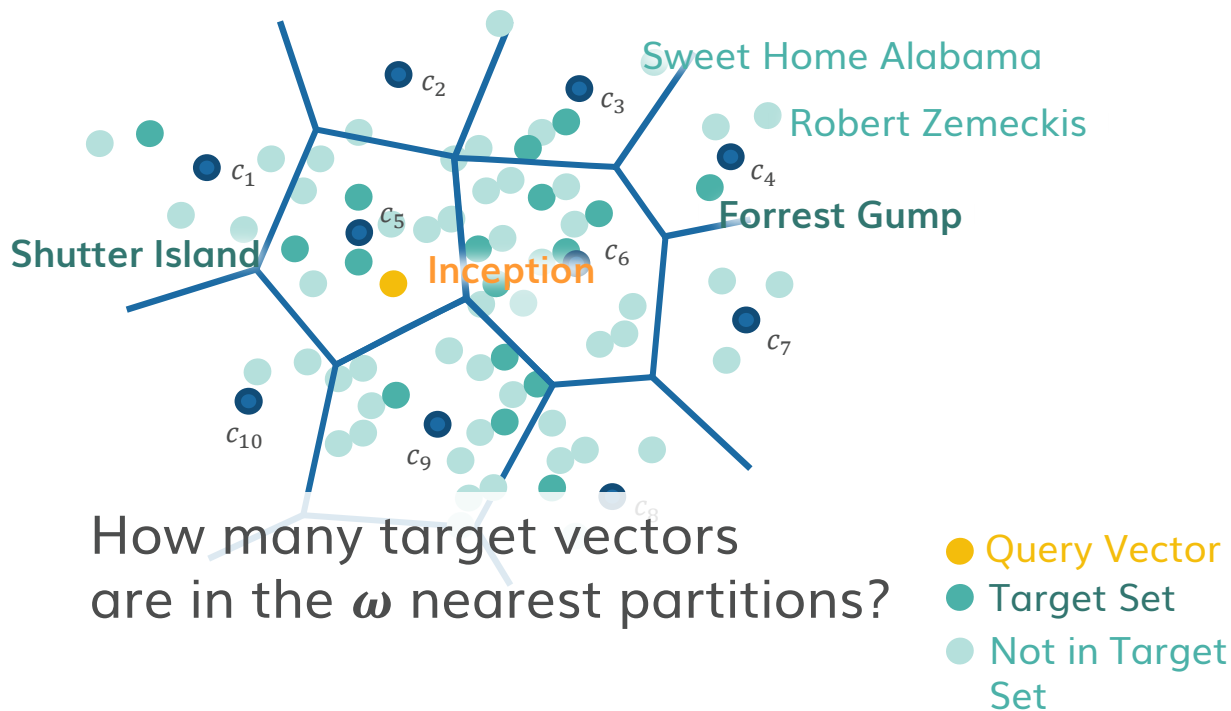


Index

ID	Token	Partition ID	PQ-Sequenz
1	Shutter Island	1	5, 1, ...
2	2010	3	3, 3, ...
3	Brando	2	6, 3, ...
4	Connery	3	8, 5, ...
5	Scarface	2	7, 3, ...
...



Query Construction



Query Construction

Problem

- Number of partitions to retrieve depends on target set size

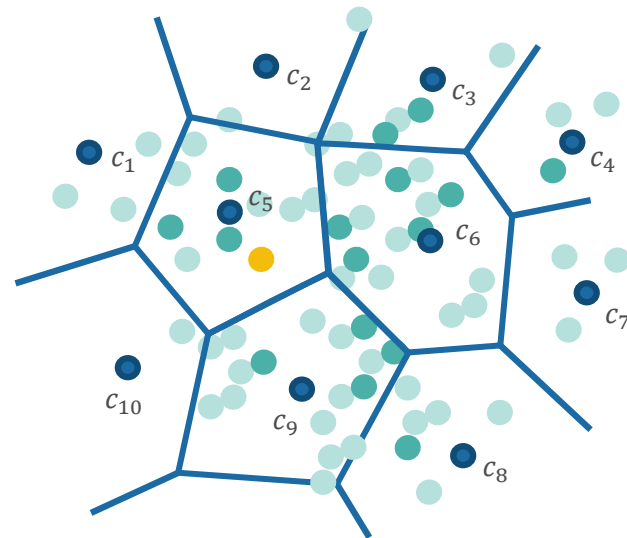
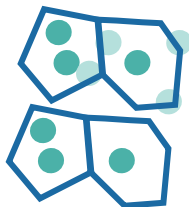
Estimate Probability of 'Getting at least n targets'

→ Hypergeometric Distr. (draw without replacement)

$$P(m \geq n) = \sum_{n \leq m \leq d} \frac{\binom{K}{m} \binom{N-K}{d-m}}{\binom{N}{d}}$$

Parameter

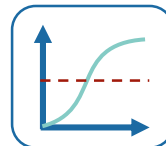
- Population Size N :
Number of all index vectors
- Number of success states in the population K :
Number of target vectors
- Number of draws d :
Number of vectors in selected partitions
- Number of Successes m :
Number of targets in selected partitions



● Query Vector

● Target Set

● Not in Target Set



Query Construction

Problem

- Number of partitions to retrieve depends on target set size

Estimate Probability of 'Getting at least n targets'

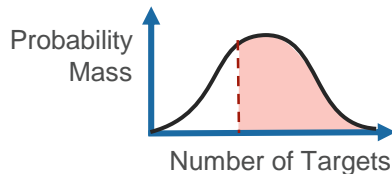
→ Hypergeometric Distr. (draw without replacement)

Calculation is very expensive!

→ Approximation via **normal distribution**

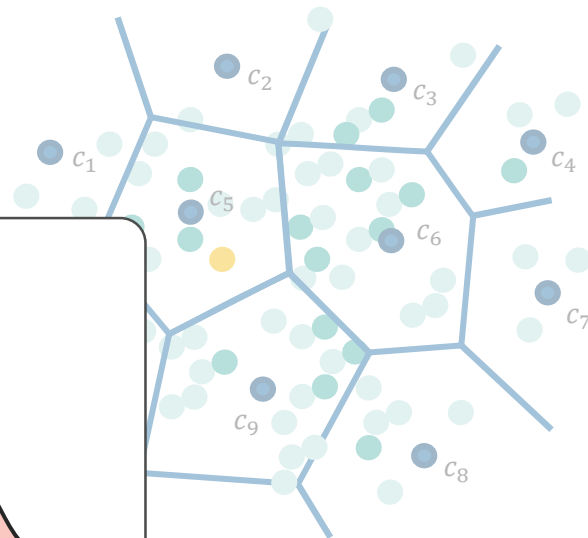
→ Integral of the normal distribution describes the confidence of 'getting at least n targets'

→ **Statistics of partition sizes** for accurate estimation



Parameter

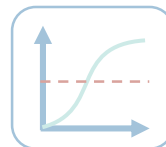
- Population
- Number of
- Number of
- Number of
- Number of
- Number of vectors in selected partitions
- Number of Successes m :
Number of targets in selected partitions



● Query Vector

● Target Set

● Not in Target Set



Query Construction (2)

Setup

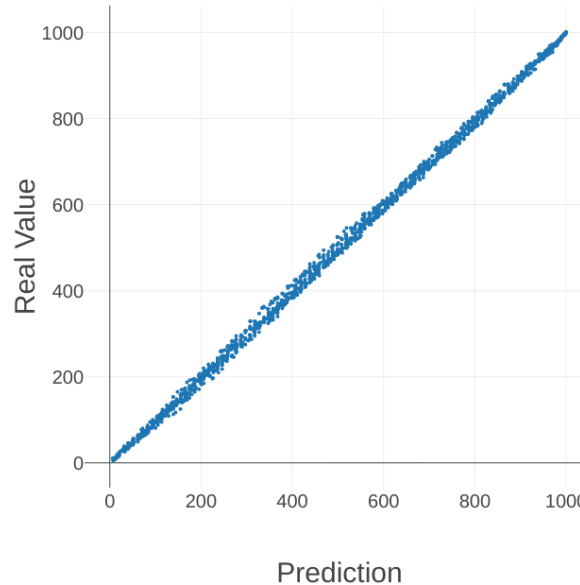
- kNN-Joins with a **single** randomly sampled **query vector** and 1000 target vectors are executed

Estimation of Target Numbers

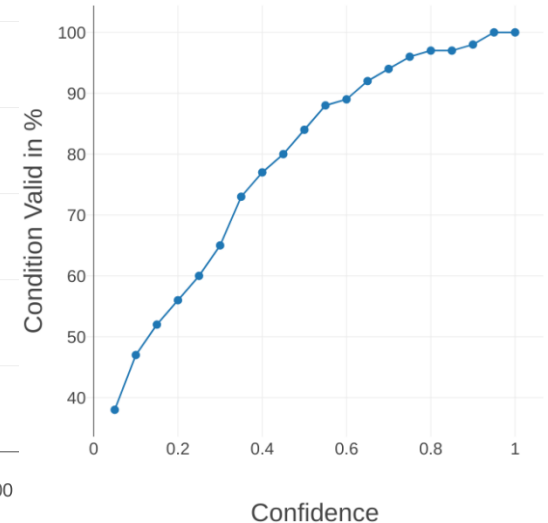
- Different selectivity values
- Prediction of the normal distribution displayed

Estimation of Confidence

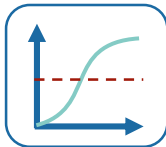
- Number of partitions is adjusted to a probability of getting at least n targets $>$ *confidence*
- Plot shows validity of condition (each rate based on 1000 queries)



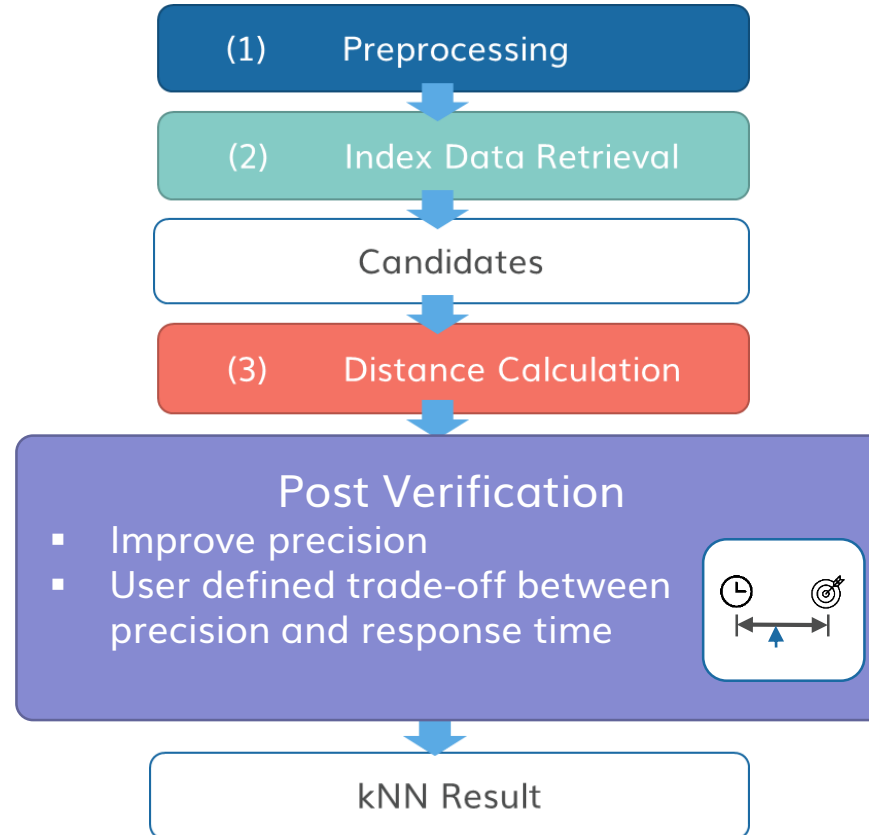
Estimated Number of Targets
(Mean of the Distribution)



Probability Estimation

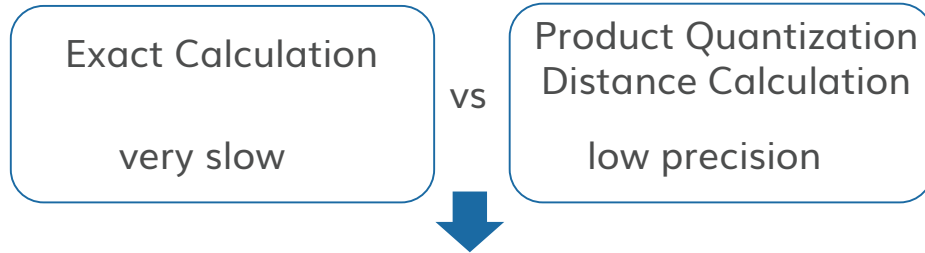


kNN-Join Algorithm



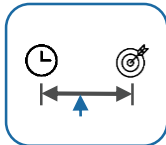
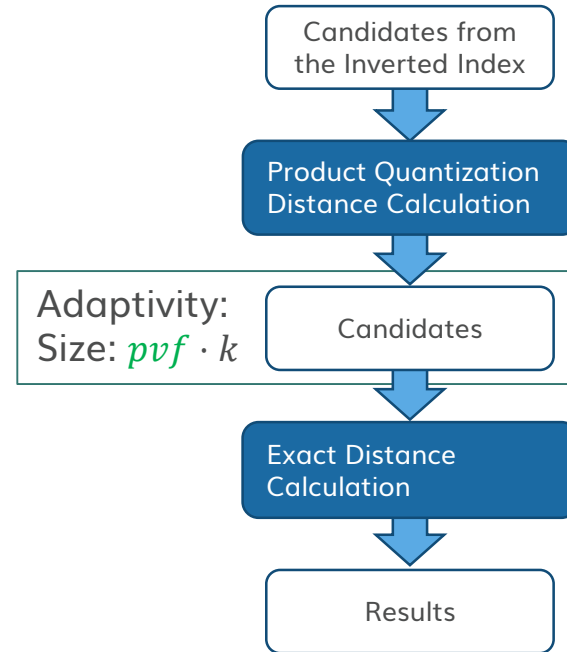
Distance Calculation with Post Verification

Combined Distance Calculation – Post Verification



Post Verification Method

- Product Quantization Methods can be used to obtain **candidate set** – size determined by pvf
 - Refine candidate set with exact calculation
 - pvf is set based on precision and response time demands
- Post verification converges fast to results of exact calculation



Evaluation

Setup

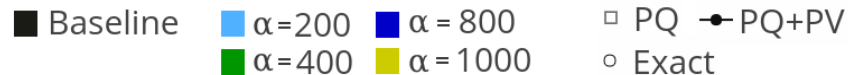
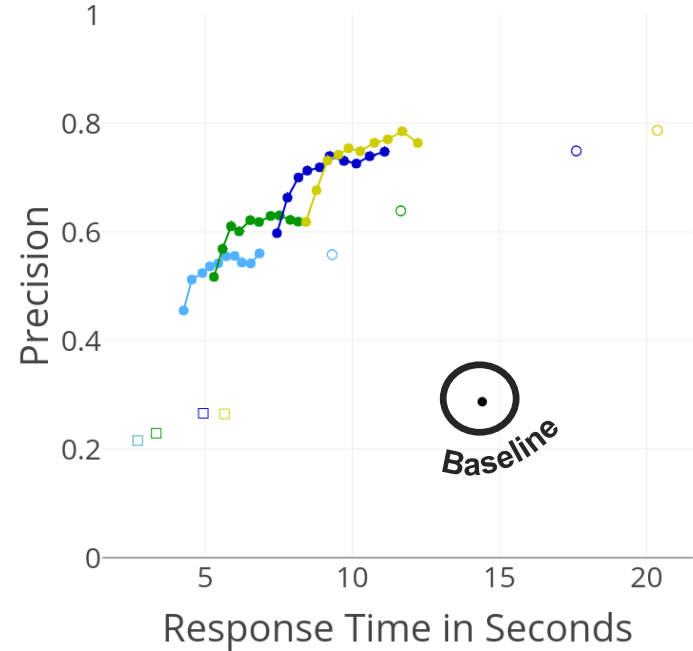
- Time and precisions measurements for randomly sampled word vectors

Dataset

- 3M 300-dimensional word vectors

Parameters

- Number of query vectors: 5,000
- Number of target vectors: 50,000
- K: 5
- Selectivity of the filter α (high values are less selective)
- Different sizes of candidate sets for post verification



Evaluation

Setup

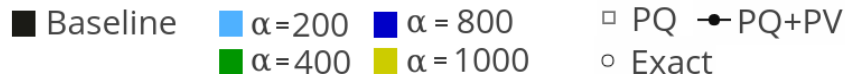
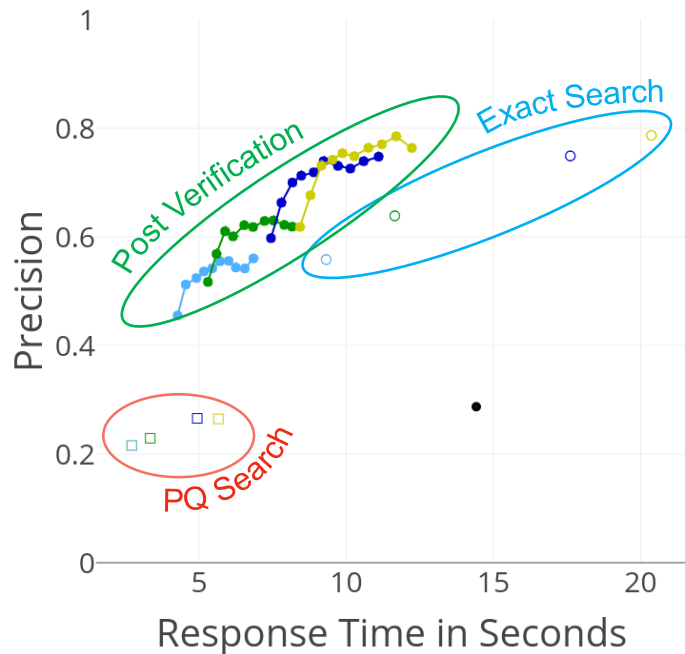
- Time and precisions measurements for randomly sampled word vectors

Dataset

- 3M 300-dimensional word vectors

Parameters

- Number of query vectors: 5,000
- Number of target vectors: 50,000
- K: 5
- Selectivity of the filter α (high values are less selective)
- Different sizes of candidate sets for post verification



Conclusion

Contributions

- Adaptive product quantization search
- **kNN-Join** based **inverted** product quantization search and post verification for **flexible target sets** → faster than the product quantization baseline
- Implementation into a relational database system

Outlook

- Automatic parameter configuration for maximal response time or minimal precision requirements

