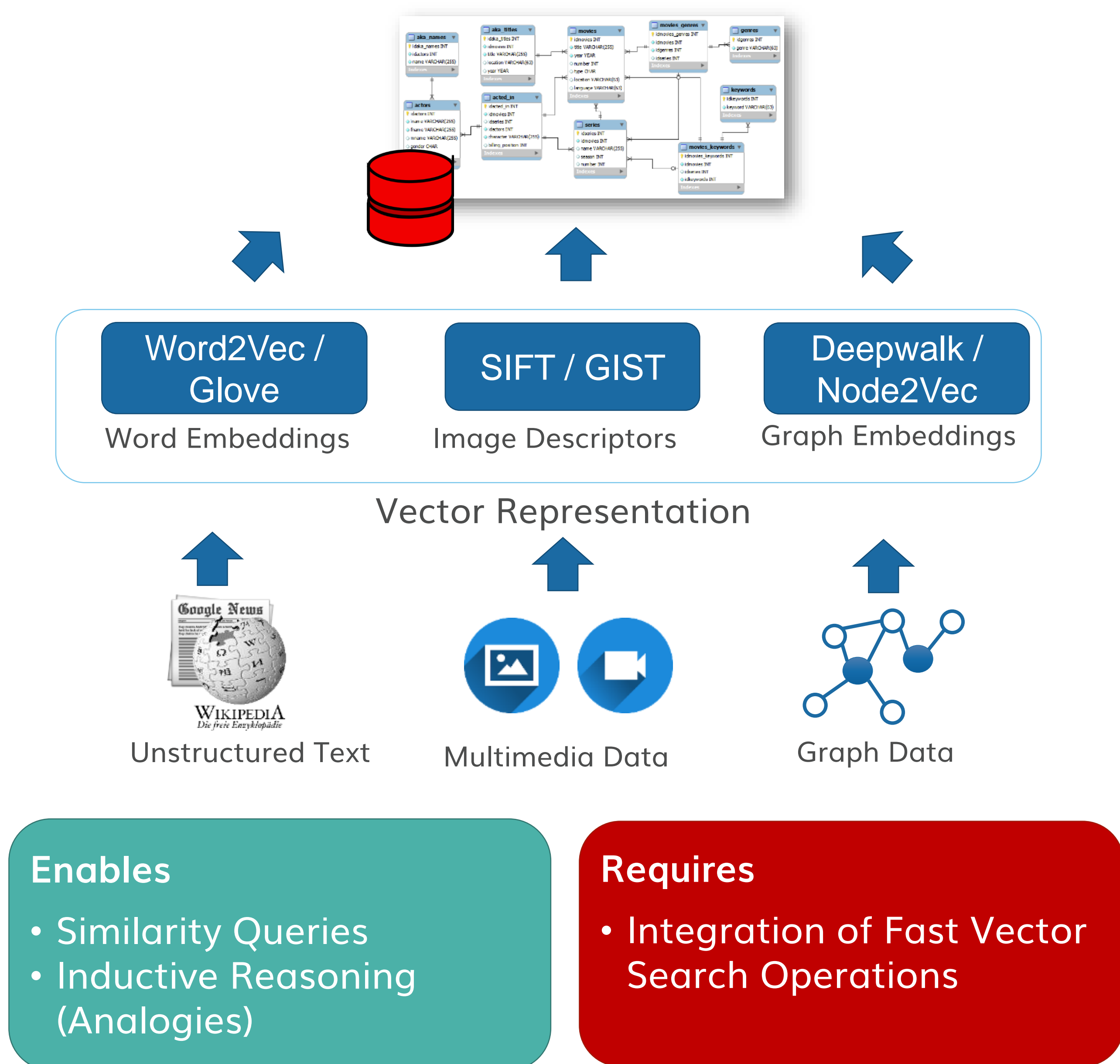
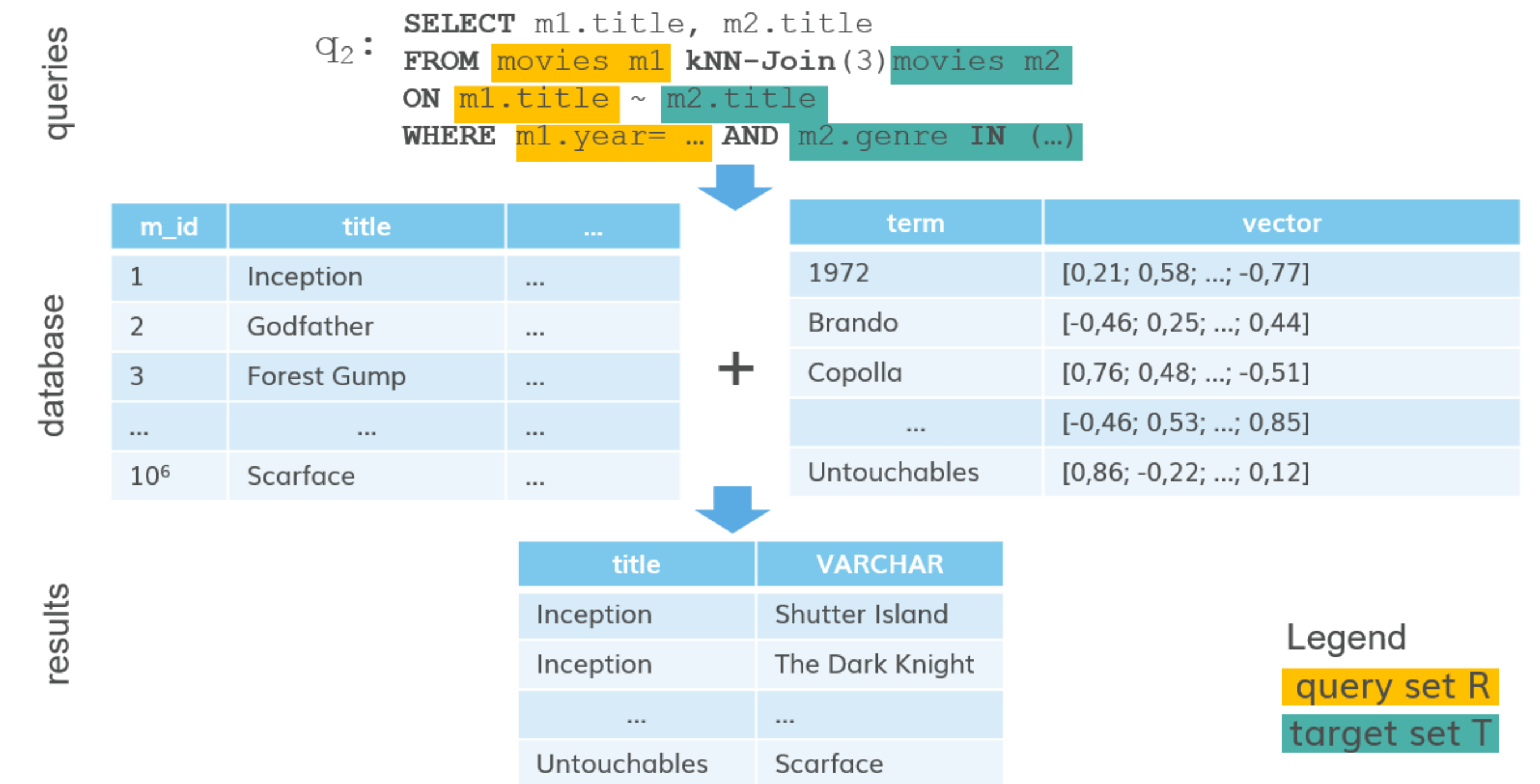


Fast Approximated Nearest Neighbor Joins For Relational Database Systems

Vector Data in Database Systems



Integration of kNN-Join Operation



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

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. geographical data)
- **Adaptive kNN-Join Algorithm**
→ Different cardinalities of join operands
→ Only one index for all vectors
- **Different Demands on Precision and Response Time**

Fast Inverted Search Algorithm

(1) Preprocessing

- Precompute distance values
- Adaptive to target set size

(2) Query Construction

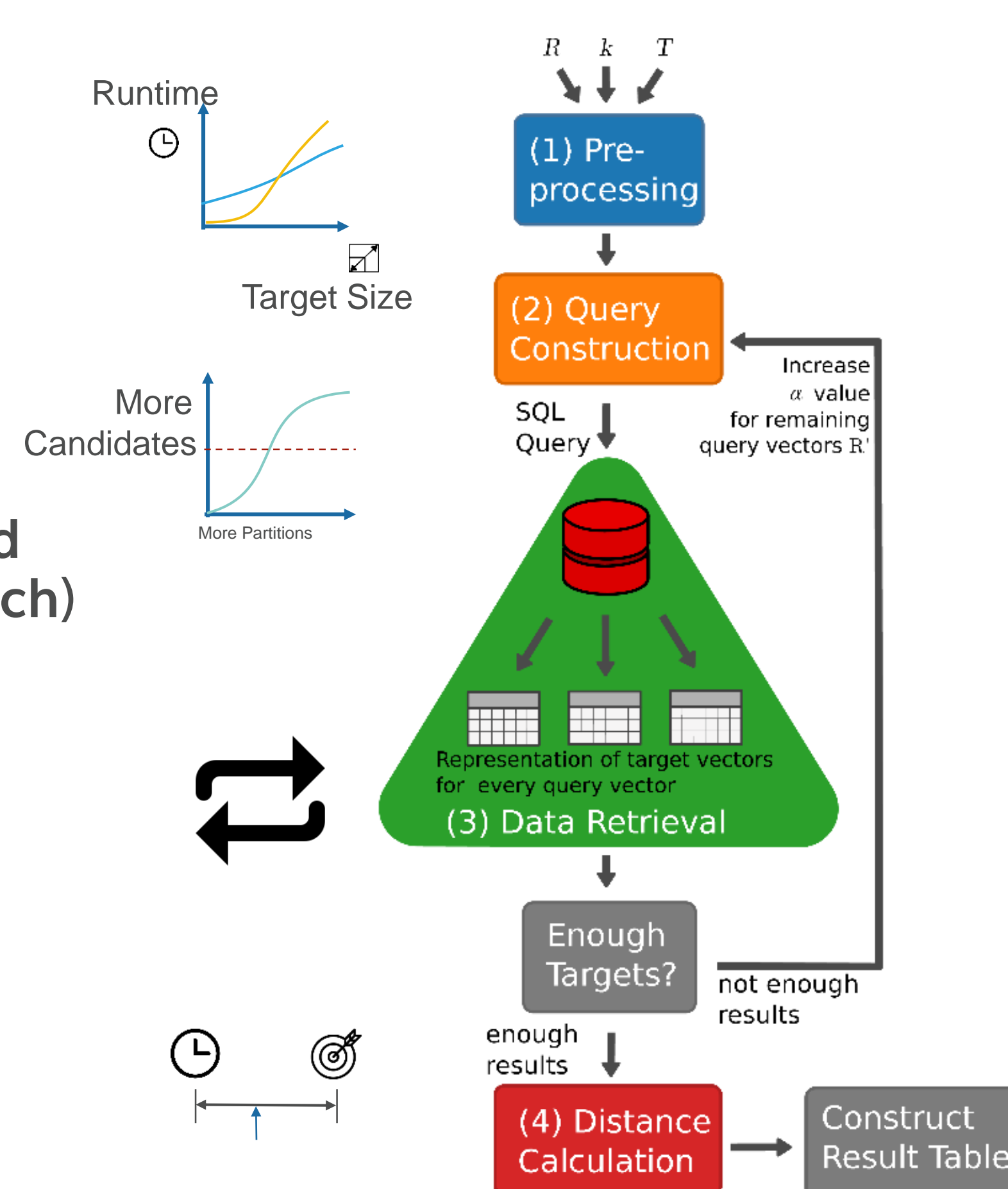
- Determine ω (Number of Part.)
- Estimate number of retrieved entries (Probabilistic Approach)

(3) Data Retrieval

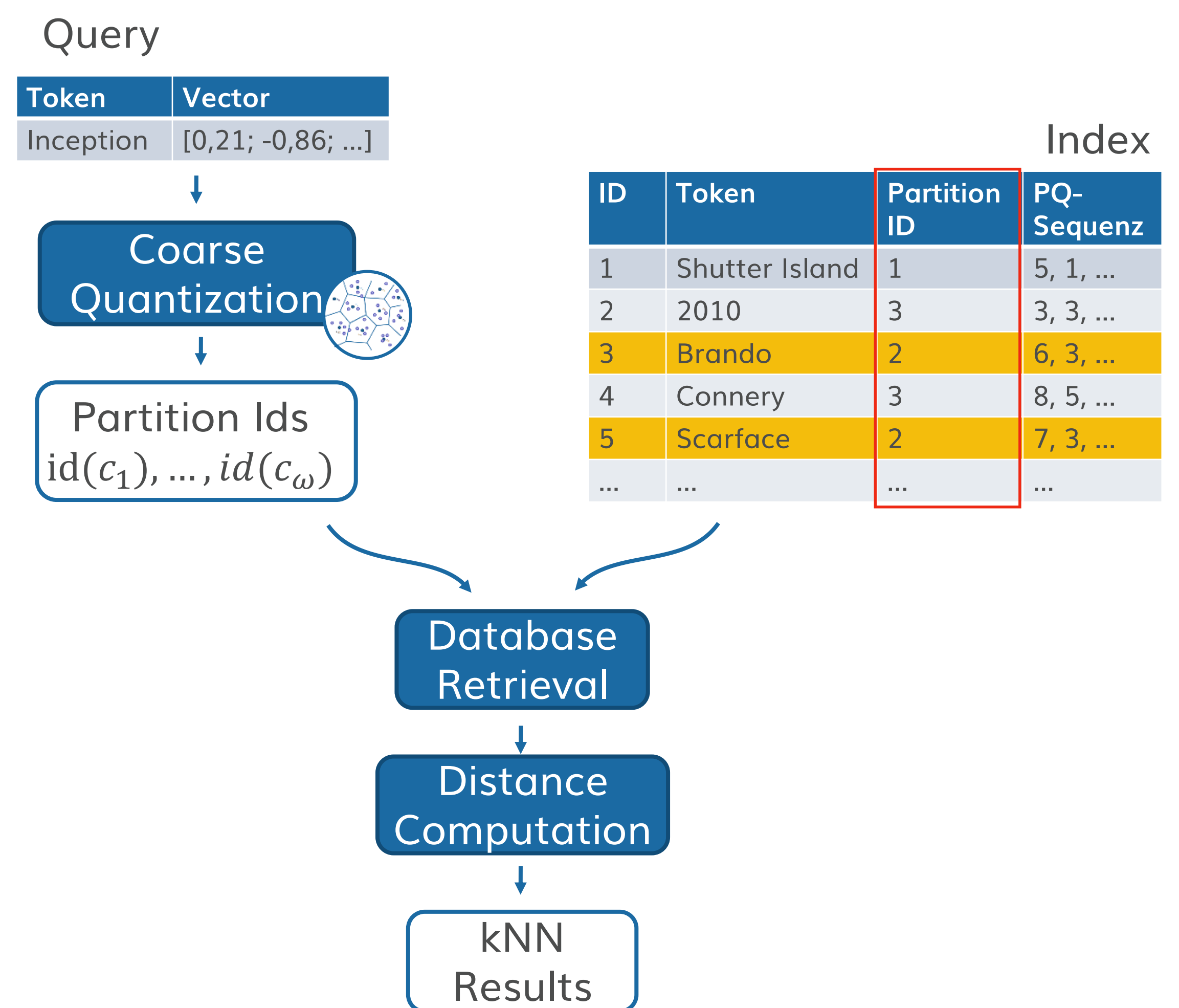
- Retrieve kNN candidates

(4) Distance Calculation

- User can adjust trade-off between precision and response time



Inverted Indexing



Estimation of Number of Partitions ω :

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

Parameter:

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

Evaluation of Algorithm with Different Distance Calculation Methods

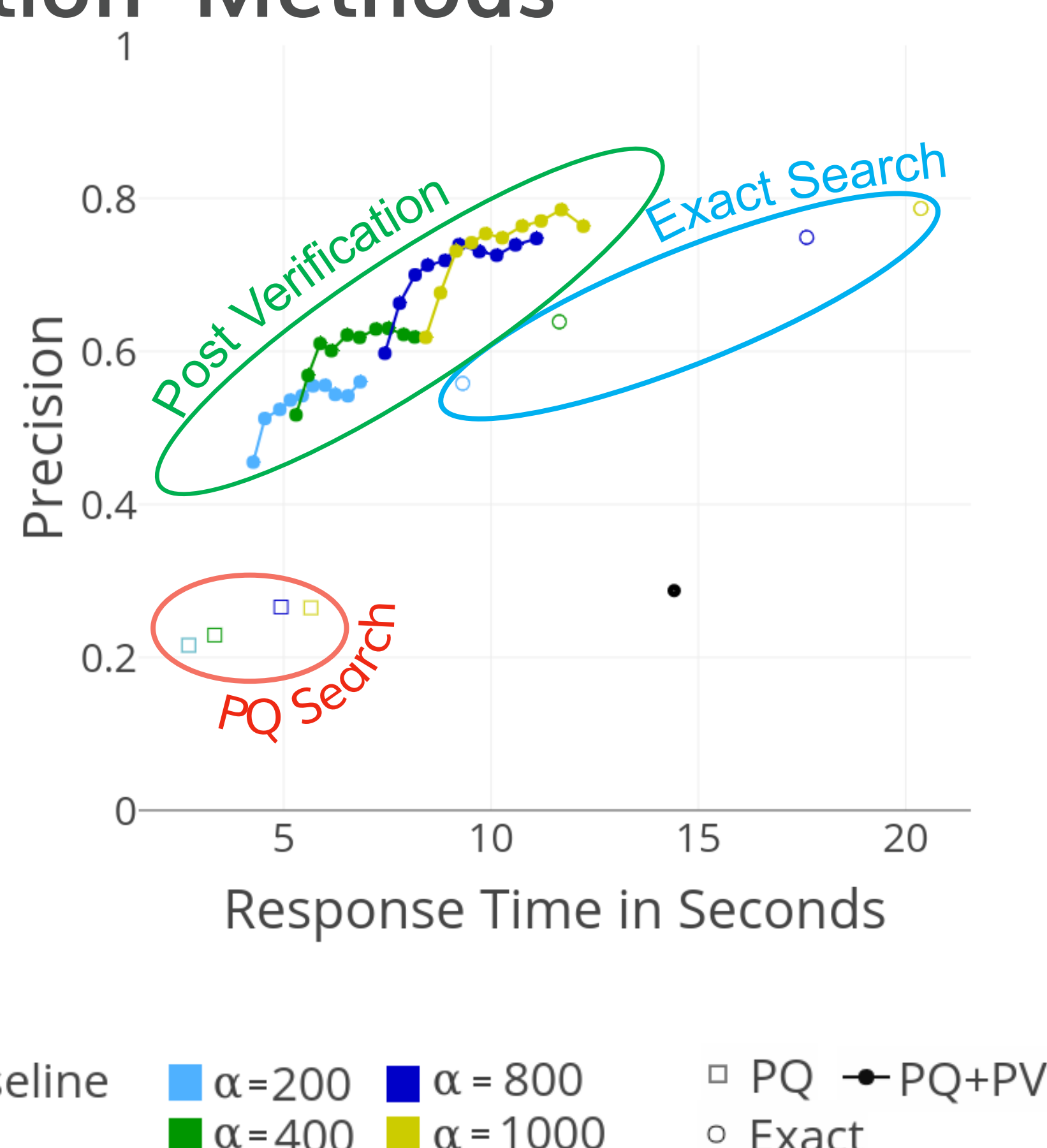
Product Quantization
Fast Approximated Search Method

Exact Distance Calculation
Fast Approximated Search Method

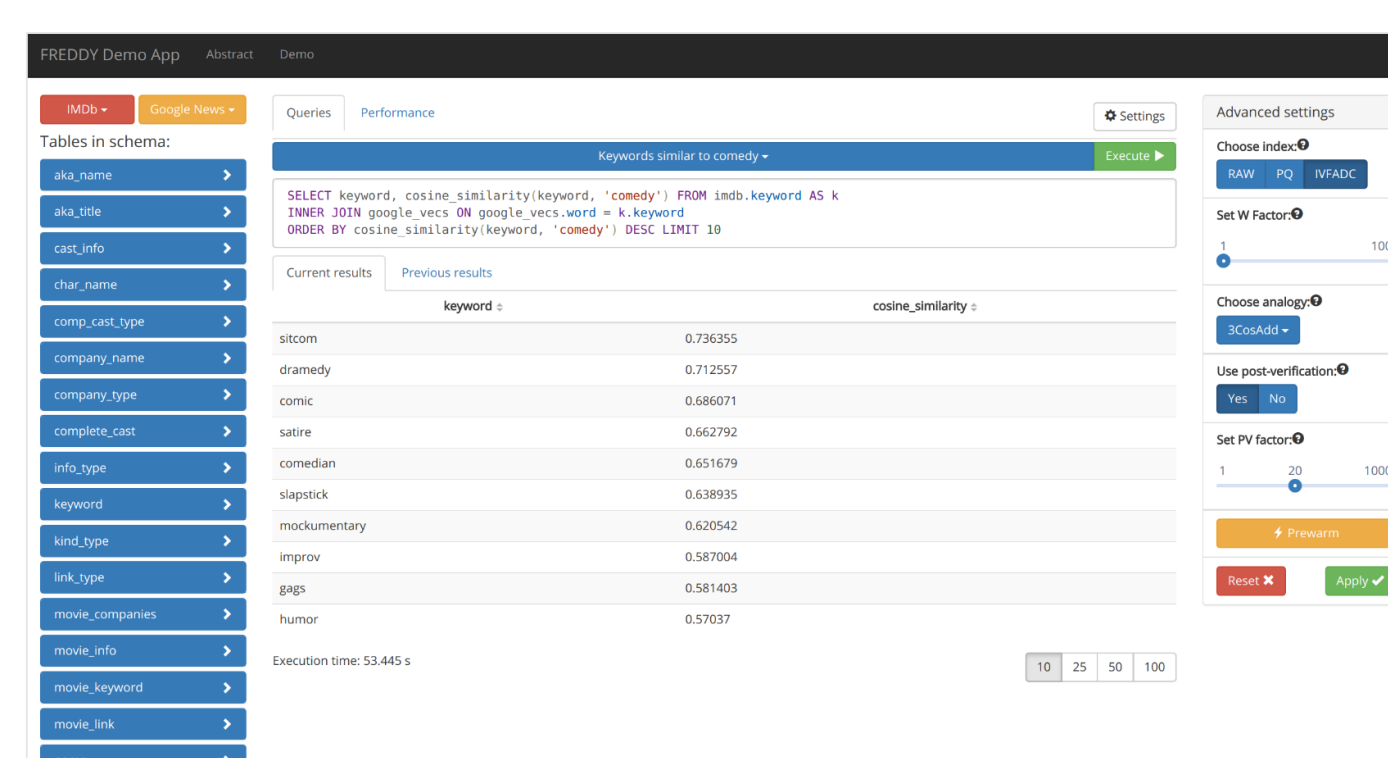
Post Verification
Post verify candidate solutions obtained by product quantization

Setup

5,000 query vectors
50,000 target vectors



Web Demo



Word Embedding operations that use the kNN-Join Index can be explored in our web demo

