

# Nearest Neighbor Search

## Brief Survey

# Caveat (setting the bar low)

- None of the projects I've worked on has nothing to do with this topic
- I just googled for fun => Do not expect too much

I. Introduction

II. Prerequisites

III. Main Topic

# What's Nearest Neighbor Search?

- Given  $N$   $D$ -dim vectors (dataset):  
 $X := \{x_i \in \mathbb{R}^D\}_{i=1,\dots,N}$
- For  $Q$  queries:
  - Given a query vector:  $q \in \mathbb{R}^D$
  - Find:  $k$ -argmin $_i \text{dist}(q, x_i)$ 
    - $\text{dist}: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}_{\geq 0}$
    - e.g. Euclidean distance, inner product, cosine similarity (argmax in this case)
- Typically,  $D = \mathcal{O}(100)$ ,  $N$  is large (thousands, millions, billions, ...)
- Brute Force linear search:
  - Space Complexity:  $\mathcal{O}(ND)$
  - Time Complexity per query:  $\mathcal{O}(ND)$
- Approximate solutions (= less accuracy) will do for better time & space complexity
- Tradeoffs

# Application

- Clustering
- Information Retrieval
  - Images
  - Texts
  - QA
    - REALM
- Recommendation System

# Glossaries

- NN: Nearest Neighbor Search ~ Similarity Search
- MIPS: Maximum Inner Product Search
- ANN: Approximate Nearest Neighbor Search
- Exact Nearest Neighbor Search

# It's been studied a lot

- Algorithms
    - Hashing
    - Partitioning
    - Graph traversal
    - Compression/Quantization
  - Tools
    - nmslib
    - FLANN
    - Falconn
    - Annoy
    - **Faiss**
    - **scaNN**

[ANN Benchmarks](#)
[Home](#)
[Datasets](#)
[Algorithms](#)
[Contact](#)

# Info

ANN-Benchmarks is a benchmarking environment for approximate nearest neighbor algorithms search. This website contains the current benchmarking results. Please visit <http://github.com/erikbern/ann-benchmarks/> to get an overview over evaluated data sets and algorithms. Make a pull request on [Github](#) to add your own code or improvements to the benchmarking system.

## Benchmarking Results

Results are split by distance measure and dataset. In the bottom, you can find an overview of an algorithm's performance on all datasets. Each dataset is annotated by ( $k = \dots$ ), the number of nearest neighbors an algorithm was supposed to return. The plot shown depicts *Recall* (the fraction of true nearest neighbors found, on average over all queries) against *Queries per second*. Clicking on a plot reveals detailed interactive plots, including approximate recall, index size, and build time.

### Benchmarks for Single Queries

### Results by Dataset

#### Distance: Angular

glove-100-angular ( $k = 10$ )

The graph displays the performance of several ANN algorithms on the glove-100-angular dataset with k=10. The y-axis represents 'queries per second (1/s)' on a logarithmic scale from 10<sup>1</sup> to 10<sup>4</sup>. The x-axis represents 'Recall' from 0.0 to 1.0. A title above the graph reads 'Recall-Queries per second (1/s) tradeoff - up and to the right is better'. The legend includes: annoy (green circles), BallTree(nmslib) (black squares), bruteforce-brais (red triangles), basic-of (blue diamonds), faiss (cyan stars), ivfnn/faiss (light blue crosses), ivfnn/nmslib (dark blue pluses), ivfnslb (yellow asterisks), milvus (orange dots), nyst (purple crosses), and NDT-paging (pink triangles). Annoy shows the highest performance at low recall, while BallTree(nmslib) and bruteforce-brais show higher performance at higher recall values.

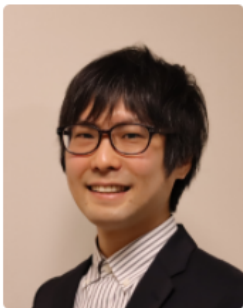
<http://ann-benchmarks.com>

# References & Today's topic

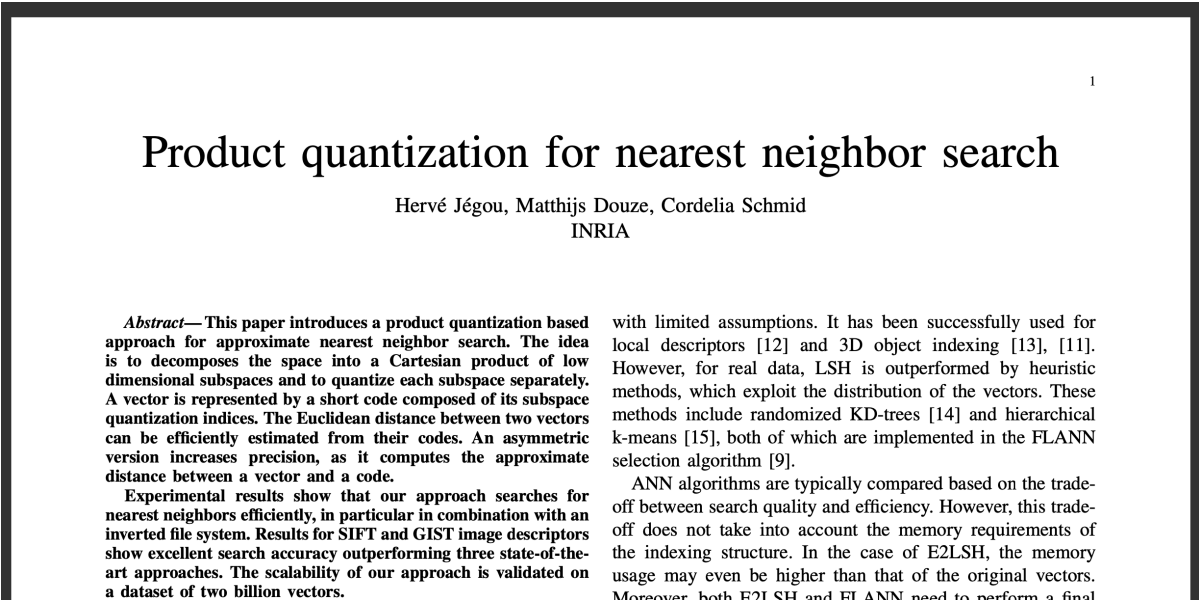


<- AWESOME Tutorial  
By Matsui-san

Highly recommend to watch this  
For comprehensive survey  
As of June 2020

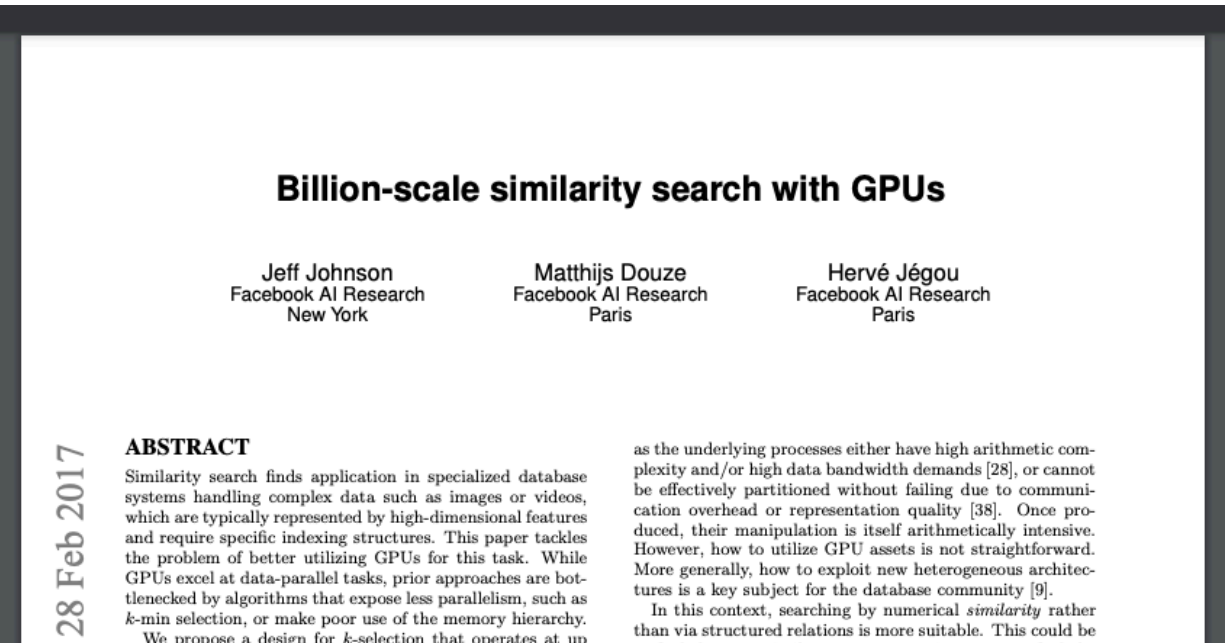


**Yusuke Matsui**  
Lecturer (Assistant Professor)  
Department of Information and Communication Engineering, Graduate School of Information Science and Technology, The University of Tokyo, Japan  
I am working as a Lecturer (Assistant Professor) at the University of Tokyo, Japan. My research is in the fields of computer vision and multimedia processing, with particular interests in large-scale indexing.



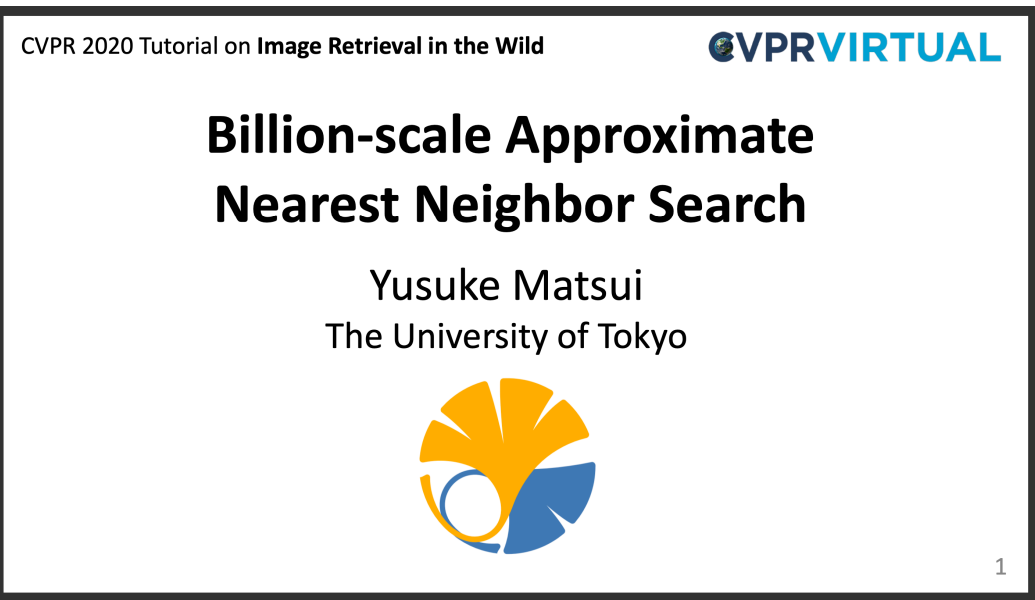
<- PQ (Product Quantization)  
Theory & experiments  
2013

Talk about some of the contents  
As prerequisites



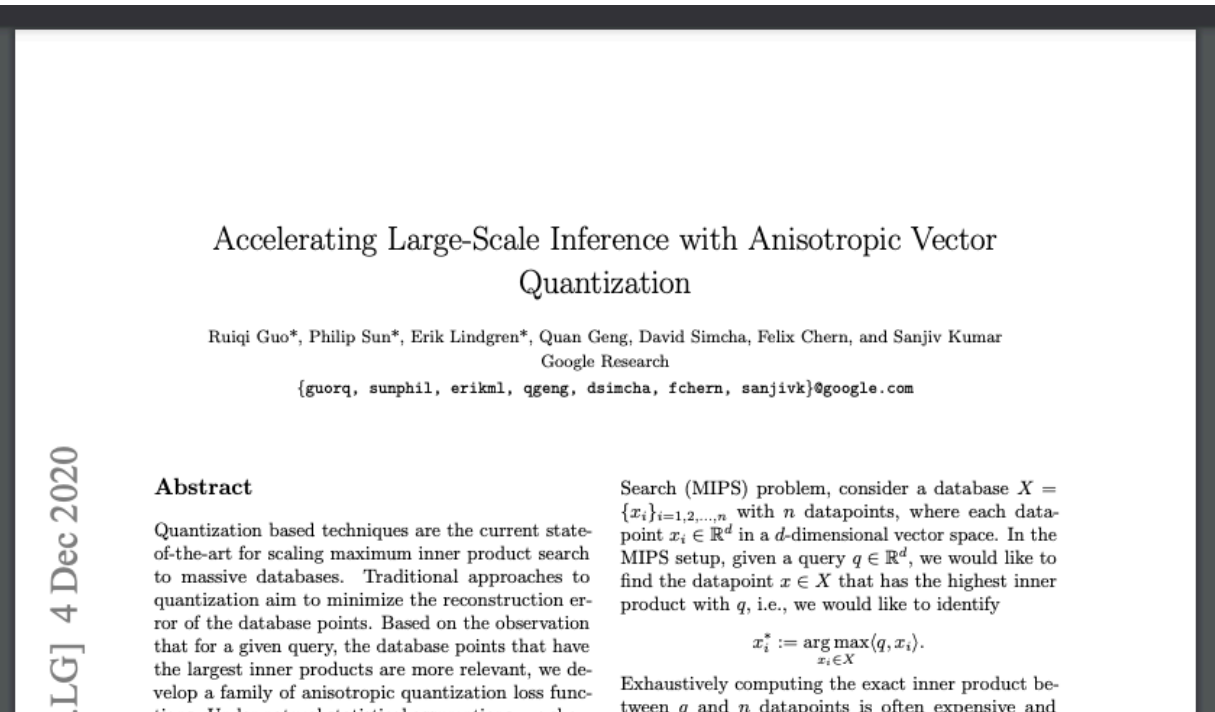
<- paper for Faiss at 2017  
PQ for GPU  
By PQ author & GPU expert

I couldn't understand the GPU  
part... :pien:



<- Slide

Some of the pages are  
Copy-and-pasted in prerequisites



<- paper for scaNN  
July 2020

Today's topic



I. Introduction

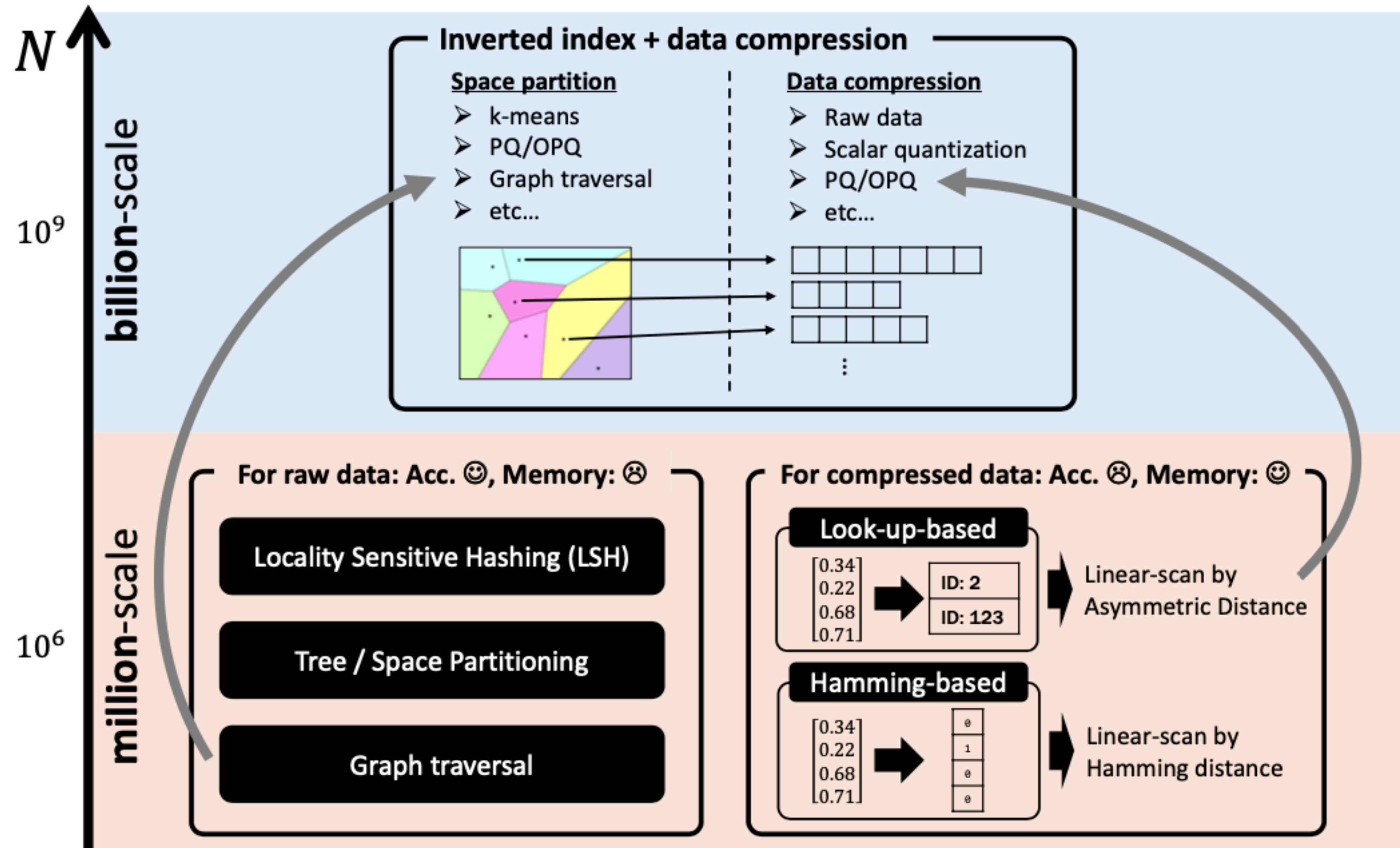
II. Prerequisites

III. Main Topic

# Prerequisites for scaNN

1. Overview of NN Algorithms
2. Product Quantization (PQ) as the best compression method
3. Inverted Index + PQ

# 1. Overview of NN Algorithms



# Prerequisites for scaNN

1. Overview of NN Algorithms
2. Product Quantization (PQ) as the best compression method
3. Inverted Index + PQ

# Compression (Quantization)

- Motivation: Cannot store all vectors  $X = \{x_i \in \mathbb{R}^D\}_{i=1..N}$  in RAM
- Basic idea: Convert each vector into a short-code (cordeword) so it takes up less RAM
- Goal: Define  $c : \mathbb{R}^D \rightarrow C$  s.t.  $\text{dist}(q, x) \sim \text{dist}(q, c(x))$  has the following properties
  - Less space & time cpx
- The error  $\mathbb{E}_q \left[ \frac{1}{N} \sum_i \text{dist}(q, x_i) - \text{dist}(q, c(x_i)) \right]$  should be small
- (Dimensional reduction satisfies above but is often thought of as a mere preprocess)

# E.g. Naive compression by k-means

- $c : \mathbb{R}^D \rightarrow C = \left\{ c_j \text{ (centroid)} \in \mathbb{R}^D \right\}_{j=1, \dots, K}$   
by just k-means clustering
  - $\mathcal{O}(K)$  sample from  $X$  is enough for construction
- $\text{dist}(q, x) \sim \text{dist}(q, c_j)$ 
  - Time cpx is still  $\mathcal{O}(D)$
- Space cpx:  $\mathcal{O}(KD + N \ln K)$ 
  - $KD$  for the codebook (set of centroids)
  - $N \ln K$  for lookup table
- Tradeoff
  - $K$  is small  $\Rightarrow$  better cpx, worse acc
  - $K$  is large  $\Rightarrow$  worse cpx, better acc
- For reasonable acc,  $K$  should be large enough
  - K-means for such  $K$  is unrealistic

# Product Quantization

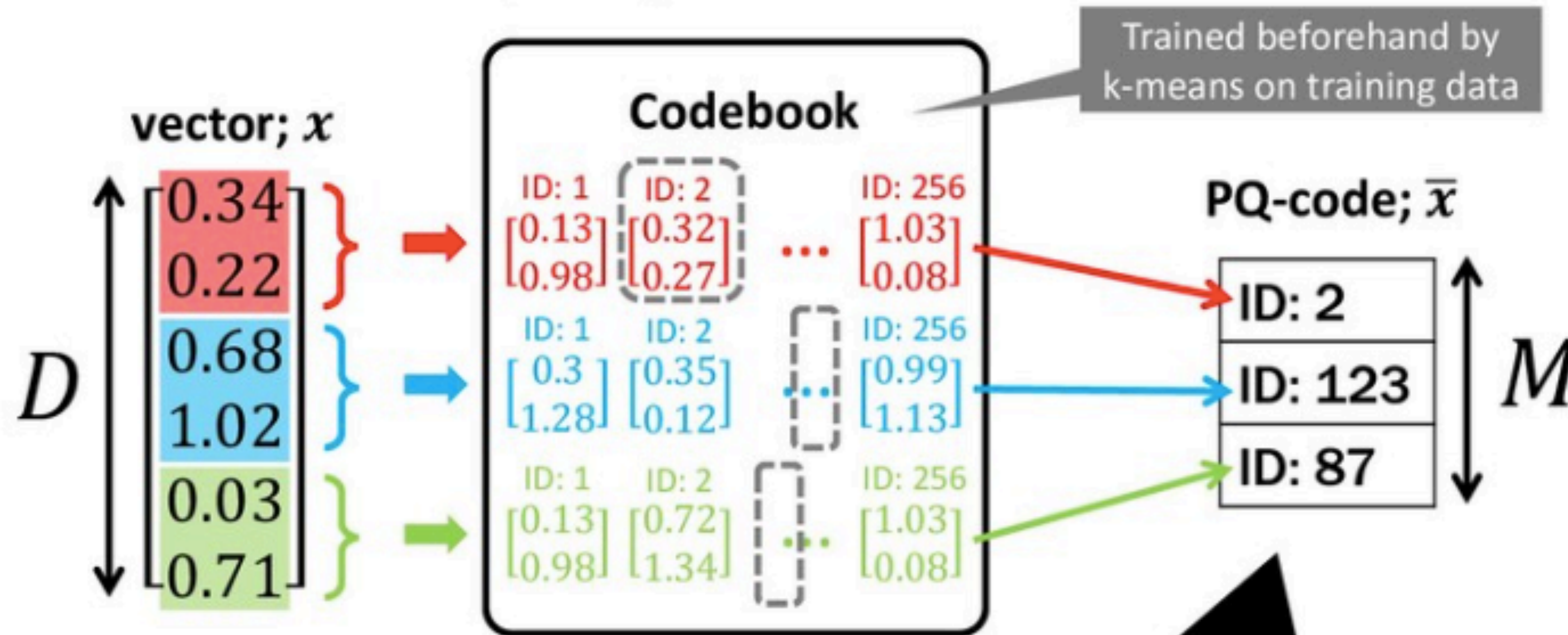
- Divide  $D$  dim into  $M$  sectors
  - Each sector has  $D^* = D/M$  dim
- Sub-quantize each sector
  - $c^* : \mathbb{R}^{D^*} \rightarrow C^* = \left\{ c_j \in \mathbb{R}^{D^*} \right\}_{j=1, \dots, K^*}$
- Total
  - $C = C_1 \times \dots \times C_M$
  - $|C| = K = (K^*)^M$
  - $M = 8, K^* = 256$  has enough acc
- $\text{dist}(q, x) \sim \text{dist}(q, c)$ 
  - Time cpx is  $\mathcal{O}(MD^* = D)$
  - ADC is better than SDC
    - Asymmetric Distance Computation: use  $q$  directly
    - Symmetric Distance computation: convert  $q$  to nearest centroids
- Space cpx:  $\mathcal{O}(K^*D + NM \ln K^*)$ 
  - $MK^*D^* = K^*D$  for the codebook ( $K^* D^*$ -dim centroids for  $M$  sectors)
  - $NM \ln K^*$  for lookup table



# \*Sanko- Kopipe

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

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



- Simple
- Memory efficient
- Distance can be estimated

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

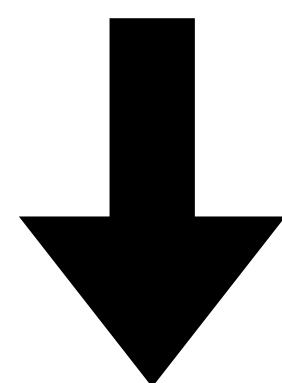


# Prerequisites for scaNN

1. Overview of NN Algorithms
2. Product Quantization (PQ) as the best compression method
3. Inverted Index + PQ to find the closest quickly

# Inverted Index

ID	w1	w2	w3	w4	...
1	I	love	you	so	...
2	I	am	a	cat	...
N	You	are	a	dog	...



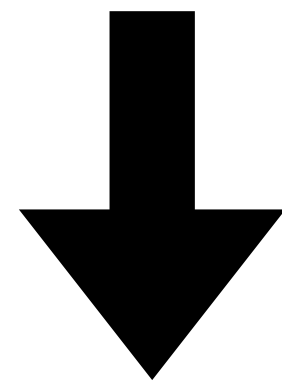
*K*

w	ID1	ID2	...
I	1	2	...
love	1	15	...
you	1	...	N
am	2	...	
a	2	...	
cat	2	...	

- For given query (= search word)
  - Find the word from the index  $\mathcal{O}(\ln K \text{ or } 1)$
  - Return the document IDs

# Inverted Index for NN

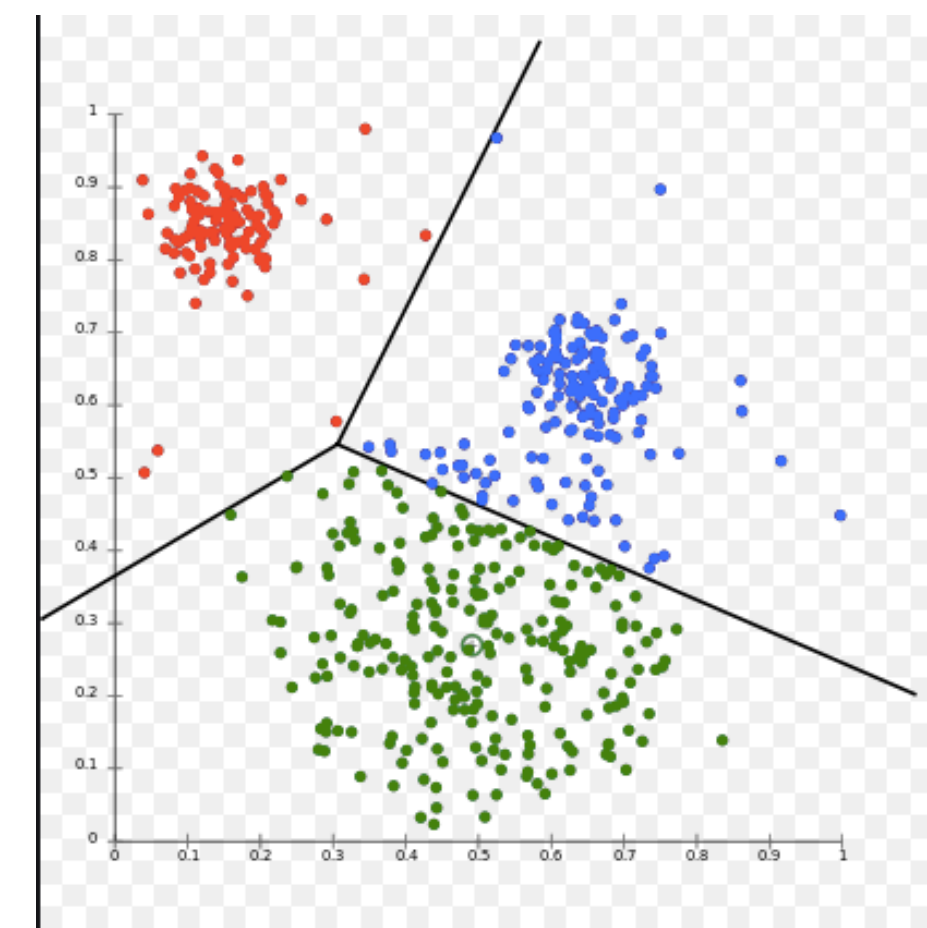
$$X = \{x_i \in \mathbb{R}^D\}_{i=1, \dots, N}$$



	ID1	ID2	...
<b>c1</b>	1	3	...
<b>c2</b>	4	15	...
<b>cK</b>	9	28	...

$$\mathcal{V}_j = \{i = 1, \dots, N \mid c(x_i) = c_j\}$$

- For given query vector
- Find the  $\tau$ -nearest centroids
  - Linear search:  $\mathcal{O}(KD)$
  - Better (e.g. HNSW)
- Find  $\operatorname{argmin}_{i \in \mathcal{V}_{j=1, \dots, \tau}} \operatorname{dist}(q, x_i)$



# Inverted Index + PQ

- Called IVFADC  
(InVerged File system with Asymmetric Distance Computation)
- $c : x \mapsto c(x) = c_c(x) + c_f(r := x - c_c(x))$ 
  - $c_c$  : coarse quantizer (K'-means)
  - $c_f$  : fine quantizer (PQ)
- Typically  
 $K' \sim \sqrt{N}, M = 8, K^* = 2^8$
- Space cpx:  
 $\mathcal{O}(K'D + K^*D + N(\ln K + \ln K^*))$
- Time cpx for search:  
 $\mathcal{O}(K'D + \tau N/K'D)$

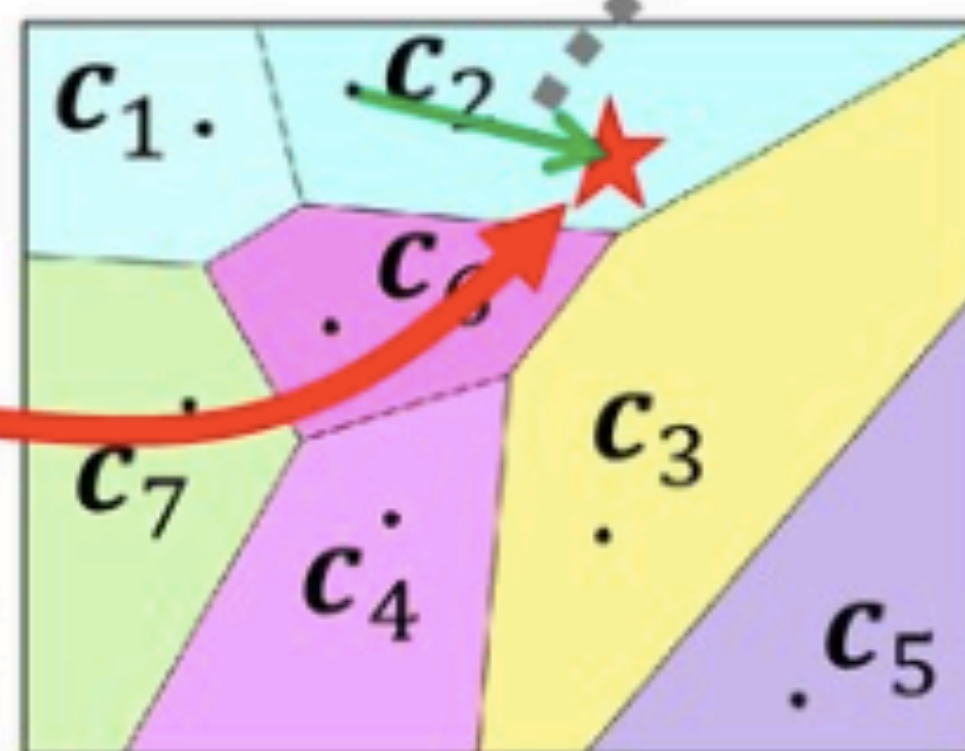
	ID1	ID2	...
<b>c1</b>	(1, cf(r1))	(3, cf(r3))	...
<b>c2</b>	(4, cf(r4))	(15, cf(r15))	...
<b>cK'</b>	(9, cf(r9))	(28, cf(r28))	...

## Inverted index + PQ: Record

Record  $x_1$

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

$x_1$



Coarse quantizer

$k = 1$

$k = 2$

$\vdots$

$k = K$

1
ID: 42
ID: 37
ID: 9

$i$

$\bar{r}_i$

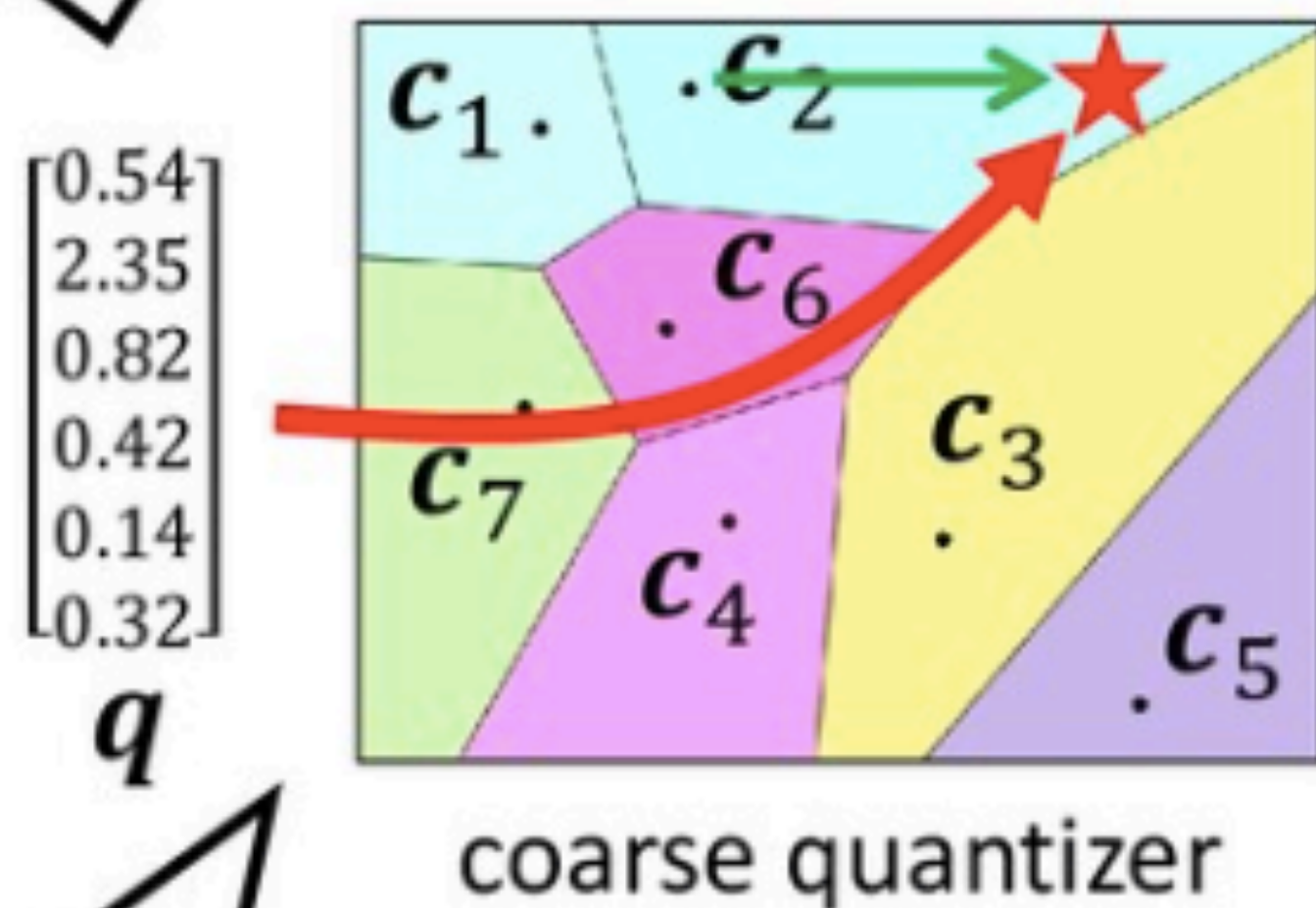
- $c_2$  is closest to  $x_1$
- Compute a residual  $r_1$  between  $x_1$  and  $c_2$ :  
 $r_1 = x_1 - c_2$  (  $\rightarrow$  )

- Quantize  $r_1$  to  $\bar{r}_1$  by PQ
- Record it with the ID, "1"
- i.e., record  $(i, \bar{r}_i)$



# Inverted index + PQ: Search

Find the nearest vector to  $q$



- $c_2$  is the closest to  $q$
- Compute the residual:  $r_q = q - c_2$

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

$i$  points to the first row (ID: 18).  
 $\bar{r}_i$  points to the second row (ID: 4).

- For all  $(i, \bar{r}_i)$  in  $k = 2$ , compare  $\bar{r}_i$  with  $r_q$ :  

$$d(q, x_i)^2 = d(q - c_2, x_i - c_2)^2$$

$$= d(r_q, r_i)^2 \sim d_A(r_q, \bar{r}_i)^2$$
- Find the smallest one (several strategies)

I. Introduction

II. Prerequisites

III. Main Topic

