

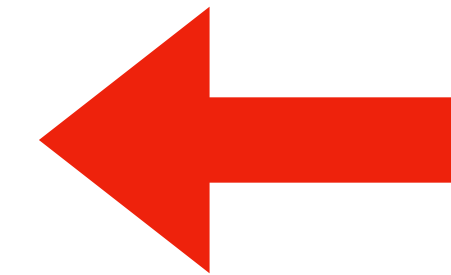
Nearest Neighbor Search

Brief Survey

Caveat (setting the bar low)

- None of the projects I've worked on has nothing to do with this topic
- I just studied for fun
- The purpose of this talk is to understand a part of the concepts
- I've never used any of the tools presented
- Some of the statements may be wrong => Please correct me

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} dist(q, x_i)
 - Focus on $k = 1$ since generalization is obvious using heap
 - 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 Cpx per query: $\mathcal{O}(ND)$
- Approximate solutions (= less accuracy) will do for better time & space complexity
- They are 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 (2017)**
- **ScaNN (July 2020)**

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 ($k = 10$)

<http://ann-benchmarks.com>

References & Today's topic



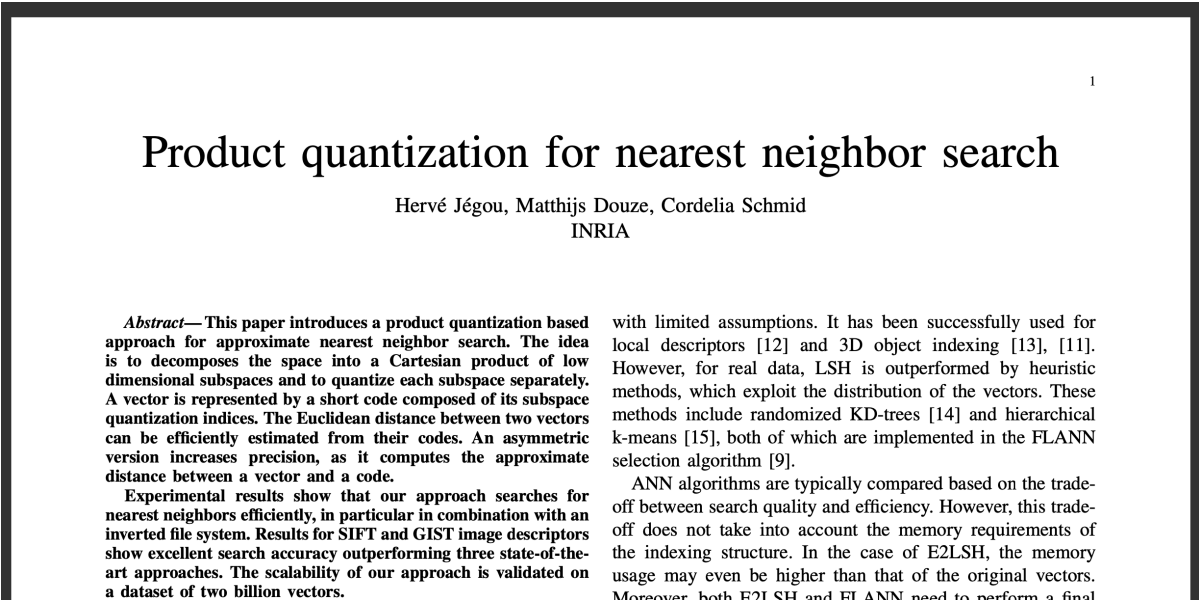
<- AWESOME Tutorial
By Matsui-san

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



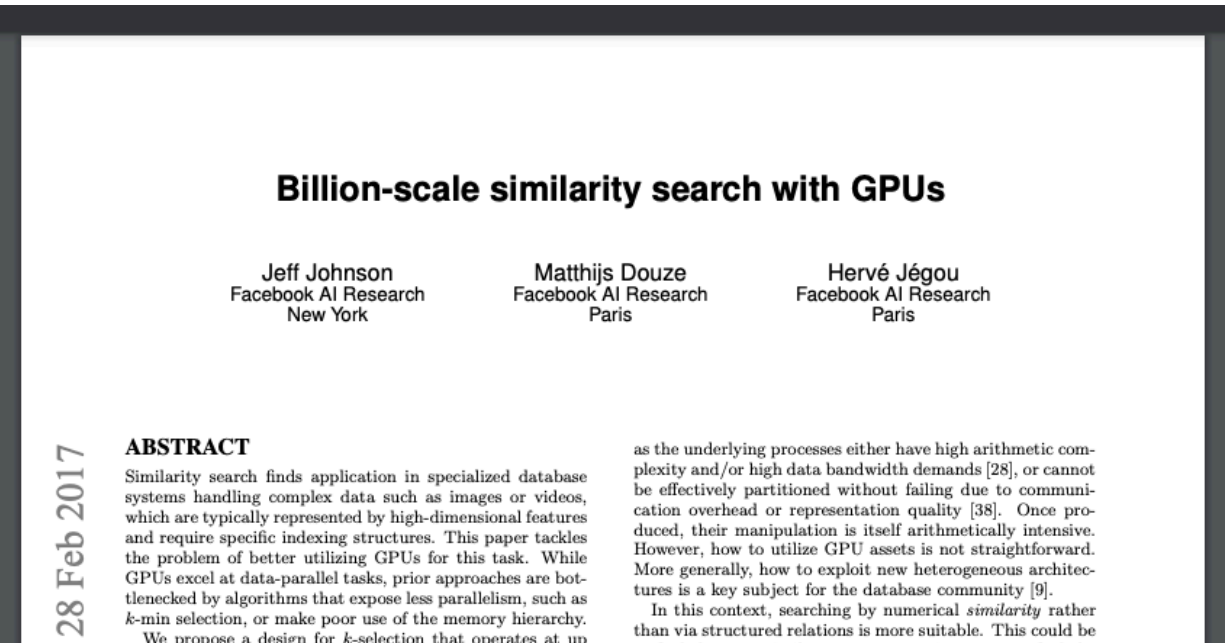
Yusuke Matsui
Lecturer (Assistant Professor)
Department of Information and Communication Engineering, Graduate School
of Information Science and Technology, The University of Tokyo, Japan

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.



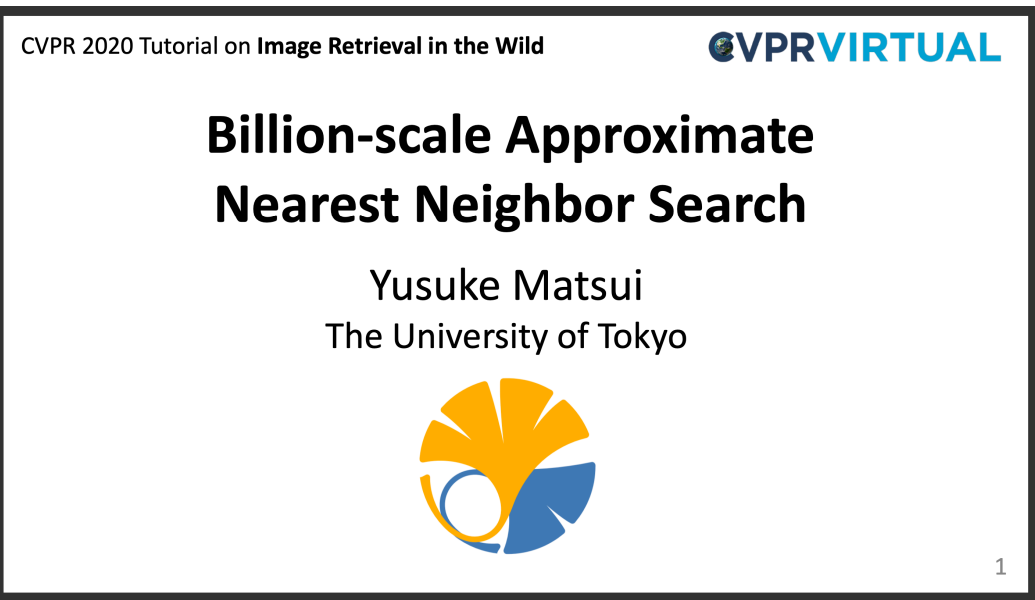
<- PQ (Product Quantization)
Theory & experiments
2013

Talk about some of the contents
As prerequisites



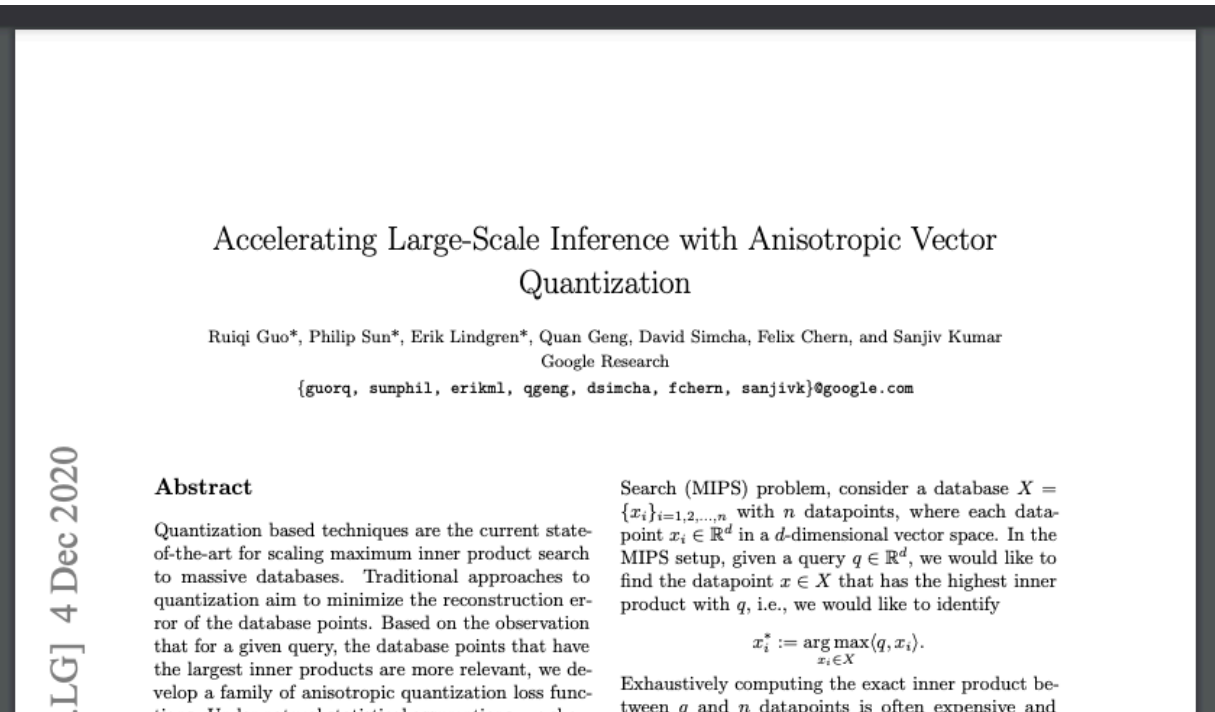
<- paper for Faiss at 2017
PQ for GPU
By PQ author & GPU expert

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



<- Slide

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

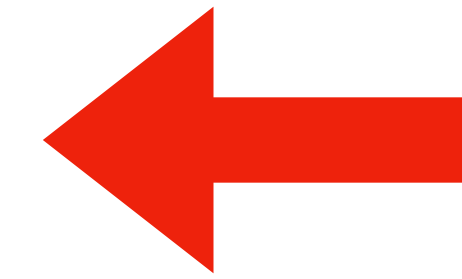


<- paper for ScaNN
July 2020

Today's topic

I. Introduction

II. Prerequisites

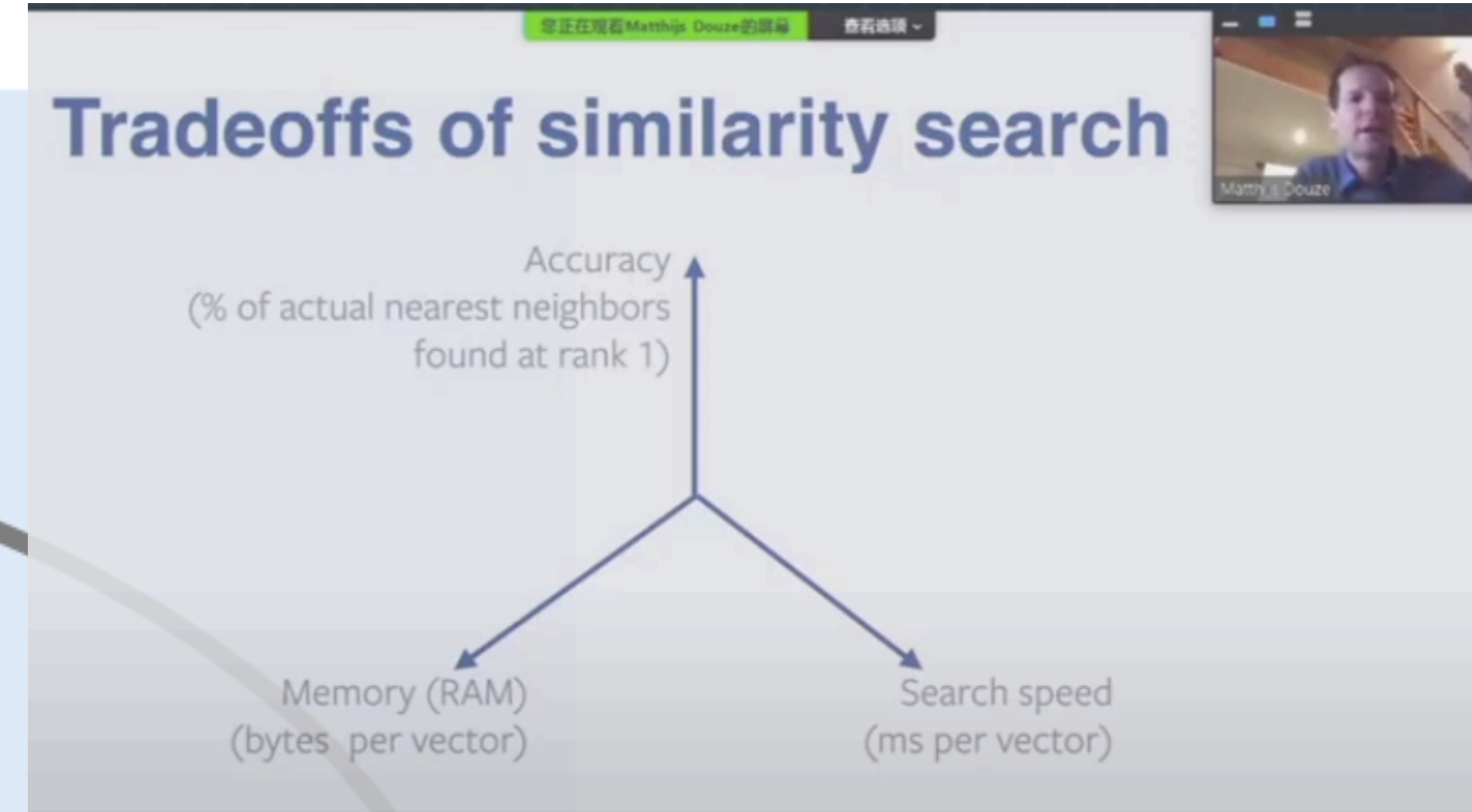
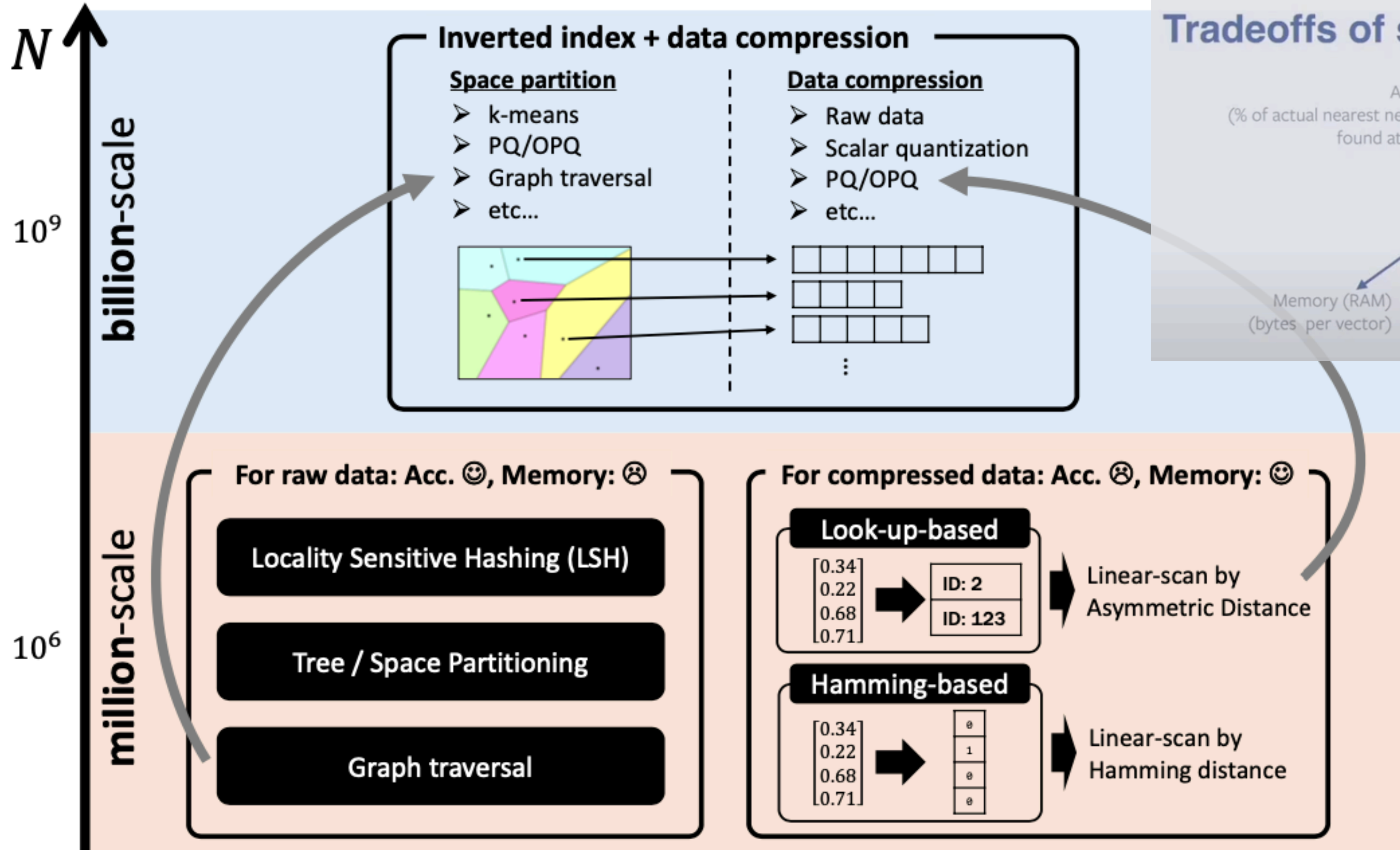


III. Main Topic

Prerequisites for ScaNN

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

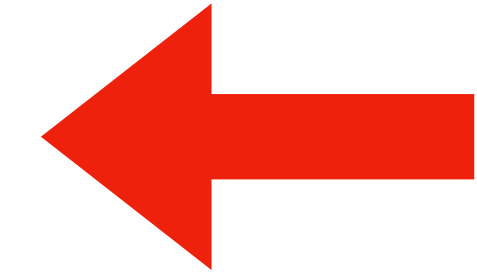
1. Overview of NN Algorithms



<https://www.youtube.com/watch?v=Un1Q92IfhPM>

Prerequisites for ScaNN

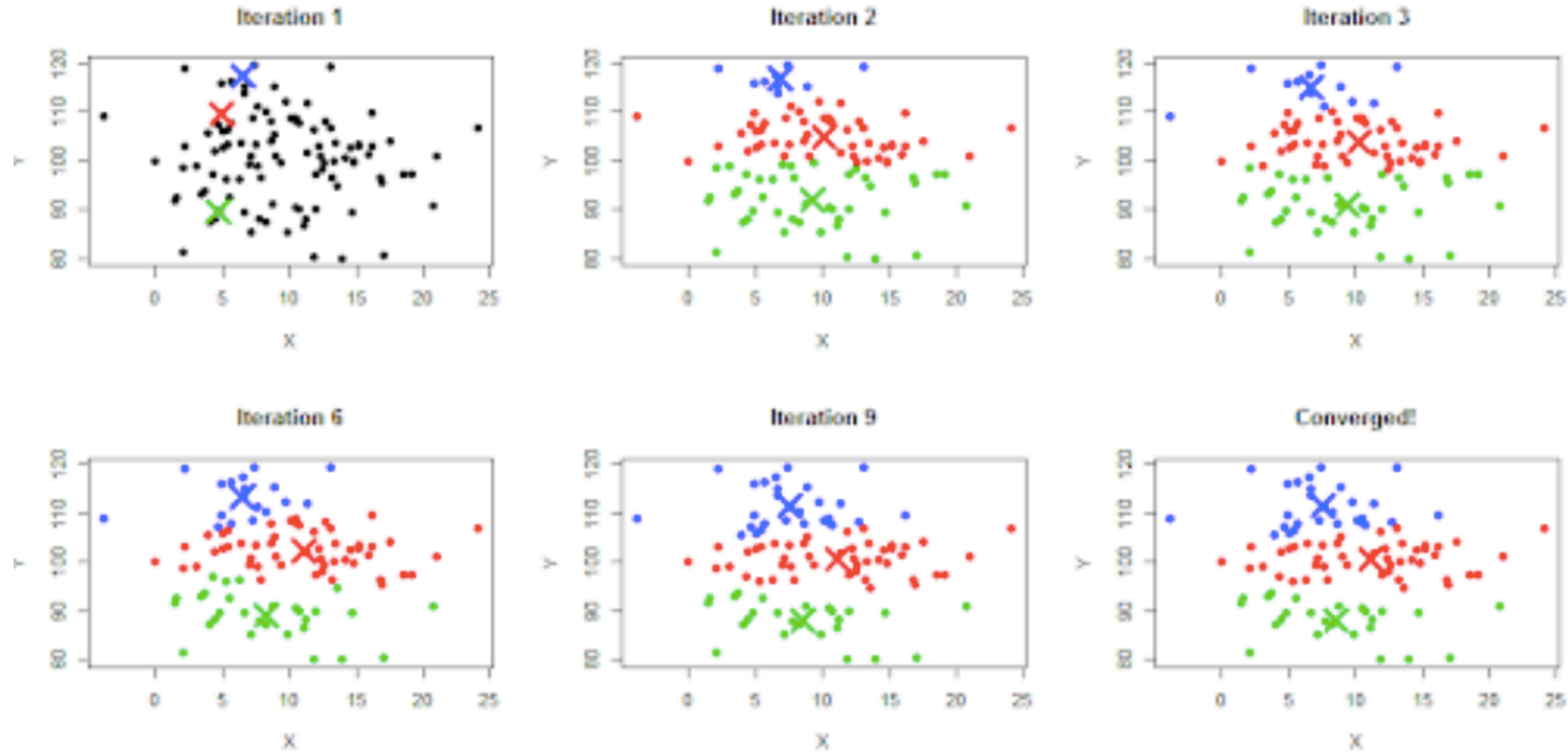
1. Overview of NN Algorithms
2. Product Quantization (PQ) as the best compression method
3. Inverted Index + PQ to find the closest quickly



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 (codeword) 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 $\mathbb{E}_q \left[\frac{1}{N} \sum_{i=1}^N \left(\text{dist}(q, x_i) - \text{dist}(q, c(x_i)) \right)^2 \right]$ should be small
- (Dimensional reduction will do but is often thought of as mere a preprocess)

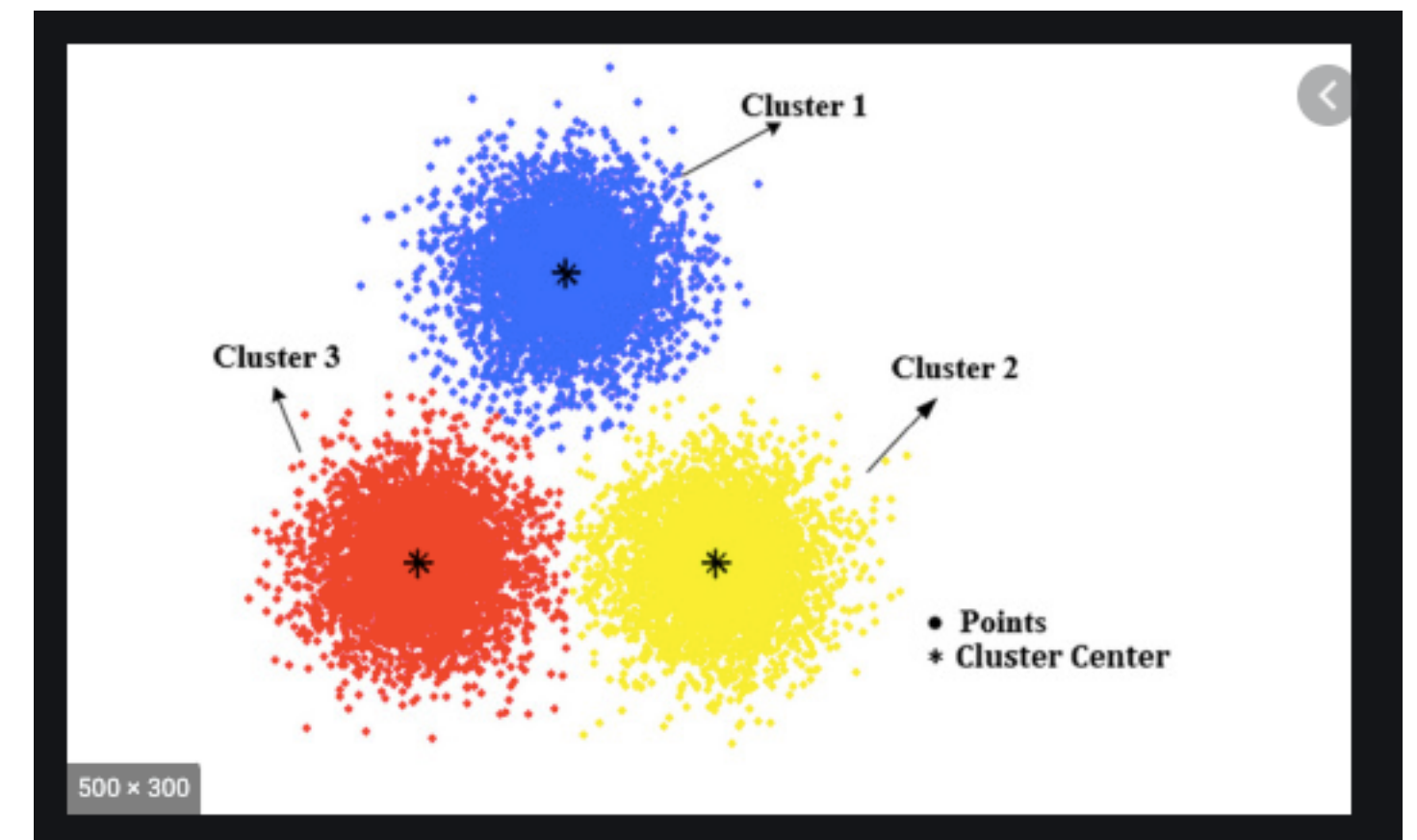
K-means clustering



- Time cpx: $\mathcal{O}(\text{iters} \times NKD)$
- Number of samples from N should be at least K

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
- $\text{dist}(q, x) \sim \text{dist}(q, c_j)$
- Space cpx: $\mathcal{O}(KD + N)$
 - KD for the codebook (set of centroids)
 - N for lookup table (set of data id for each centroid)
- Time cpx to calc $\left\{ \text{dist}(q, c(x_i)) \right\}_{i=1, \dots, N}$
: $\mathcal{O}(KD + N)$
- 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

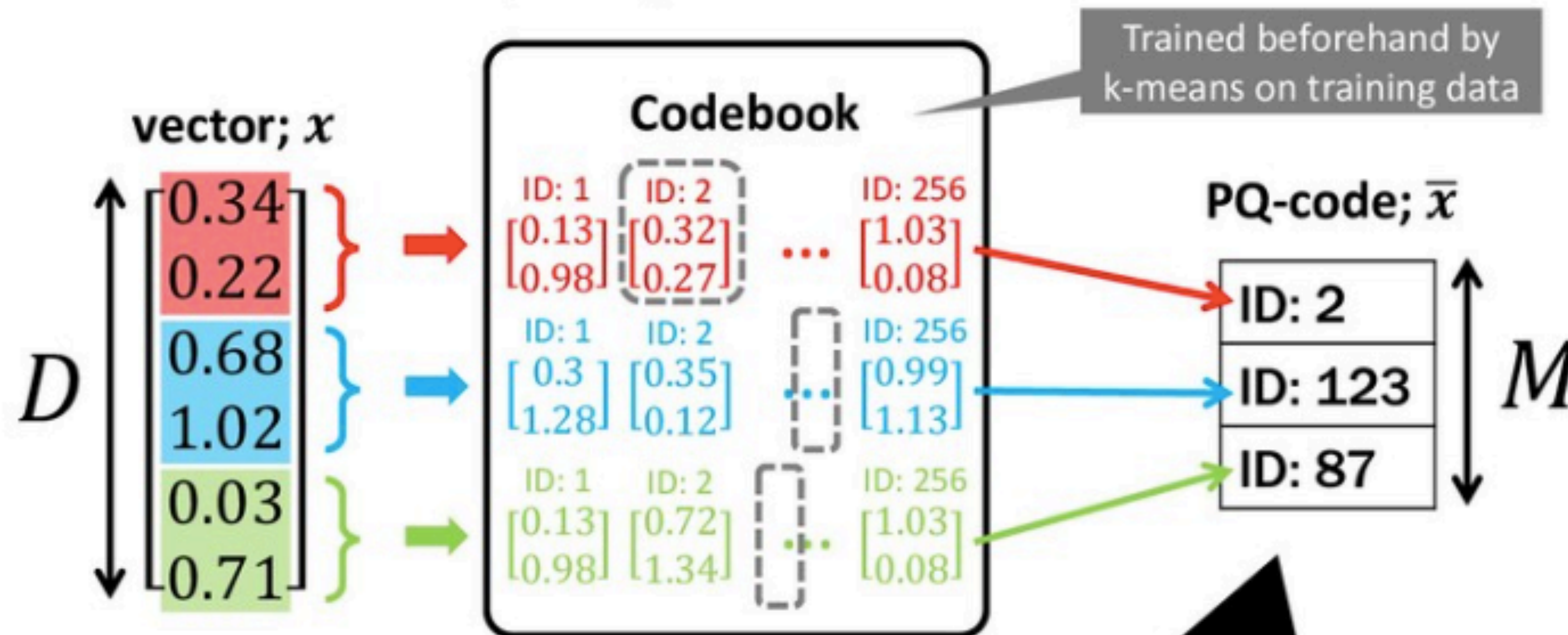


Product Quantization

- 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
 - $C = C_1 \times \dots \times C_M$
 - $|C| = K = (K^*)^M$
- $M = 8, K^* = 256$ gives enough acc
- $\text{dist}(q, x) \sim \text{dist}(q, c(x))$
- Space cpx: $\mathcal{O}(K^*D + NM)$
 - $MK^*D^* = K^*D$ for the codebook ($K^* D^*$ -dim centroids for M sectors)
 - NM for lookup table
- Time cpx to calc $\left\{ \text{dist}(q, c(x_i)) \right\}_{i=1, \dots, N}$
: $\mathcal{O}(K^*D + NM)$

Product Quantization; PQ [Jégou, TPAMI 2011]

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

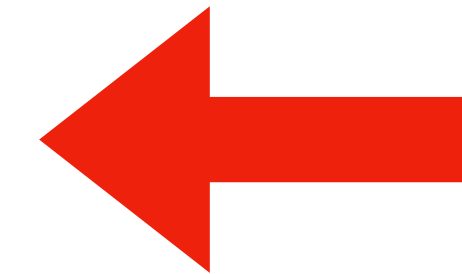


- Simple
- Memory efficient
- Distance can be estimated

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

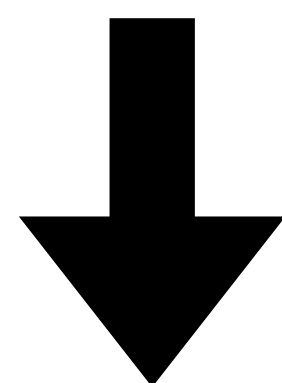
Prerequisites for ScaNN

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	...



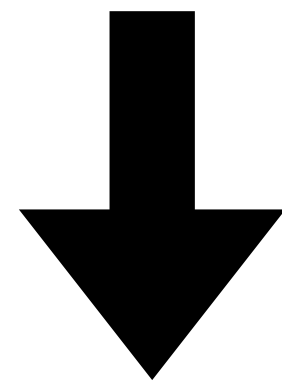
K

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

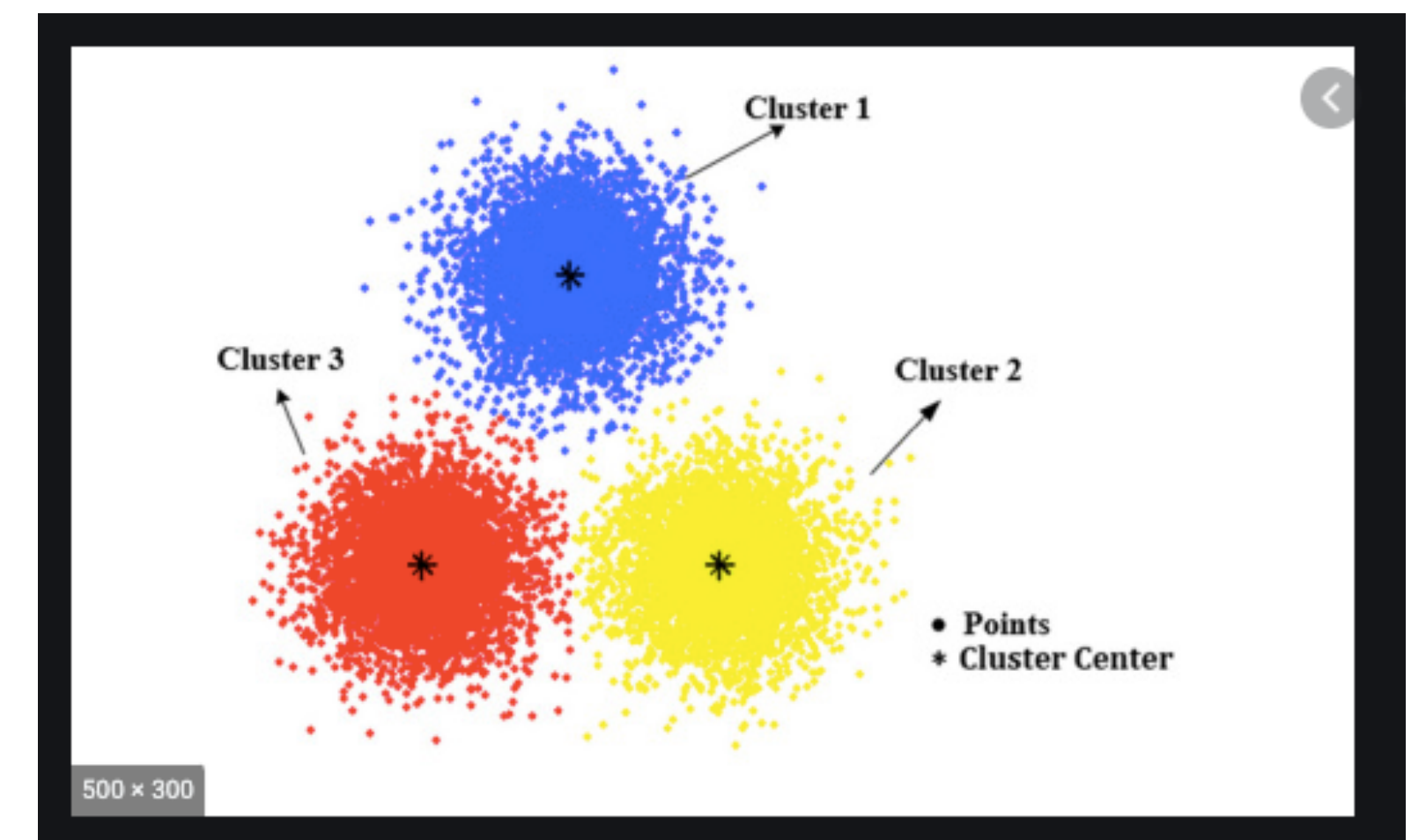
$$X = \{x_i \in \mathbb{R}^D\}_{i=1, \dots, N}$$



	ID1	ID2	...
c1	1	3	...
c2	4	15	...
cK	9	28	...

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

- For given query vector
 - Find the τ -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 or IVFPQ (InVerged File system with Asymmetric Distance Computation)
- Two stage quantize
 $c : x \mapsto c(x) = c_c(x) + c_f(r := x - c_c(x))$
 - c_c : coarse quantizer (K'-means)
 - For Time cpx
 - c_f : fine quantizer (PQ)
 - For Space cpx with acc

- Typically $K' \sim \sqrt{N}$, $M = 8$, $K^* = 2^8$
- Space cpx:
 $\mathcal{O}(K'D + K^*D + N(1 + M \ln K^*))$
- Time cpx for search:
 $\mathcal{O}(K'D + \tau N/K'D)$

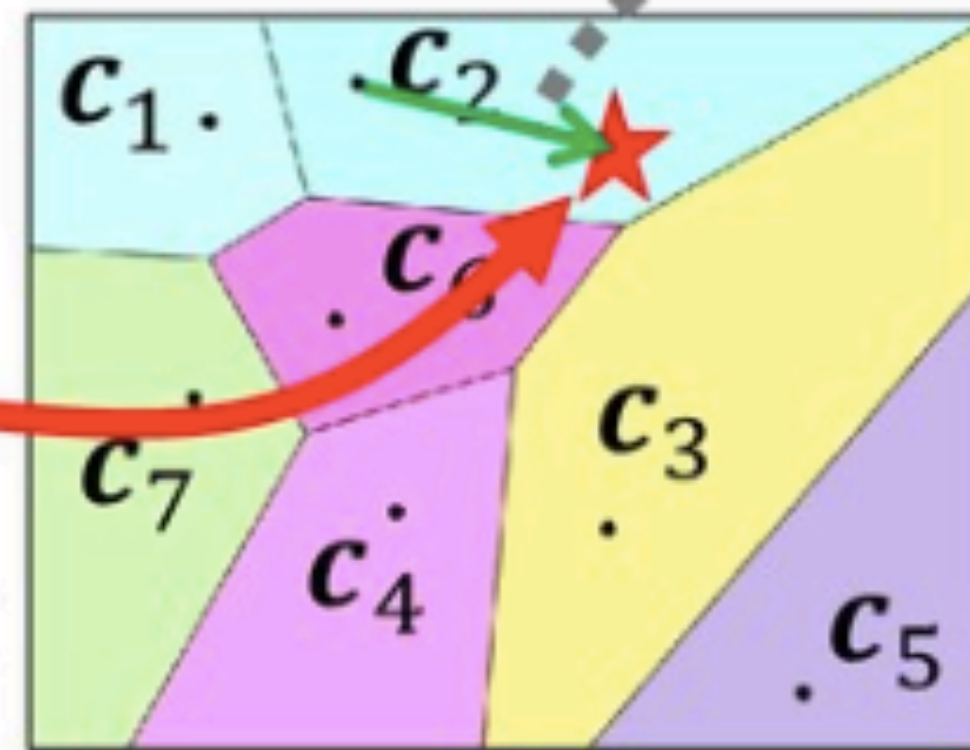
	ID1	ID2	...
c1	(1, cf(r1))	(3, cf(r3))	...
c2	(4, cf(r4))	(15, cf(r15))	...
cK'	(9, cf(r9))	(28, cf(r28))	...

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

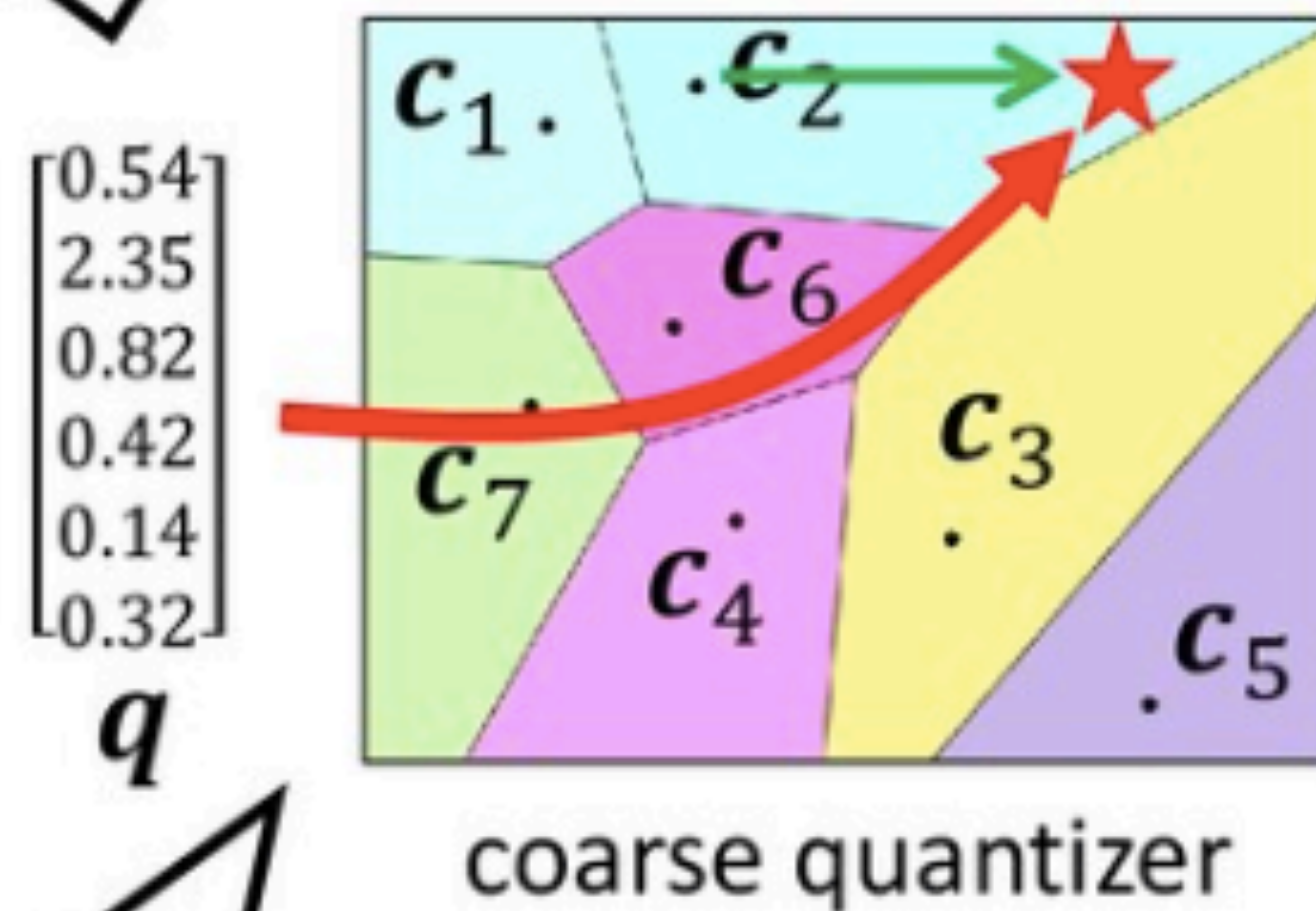
\bar{r}_i

- 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

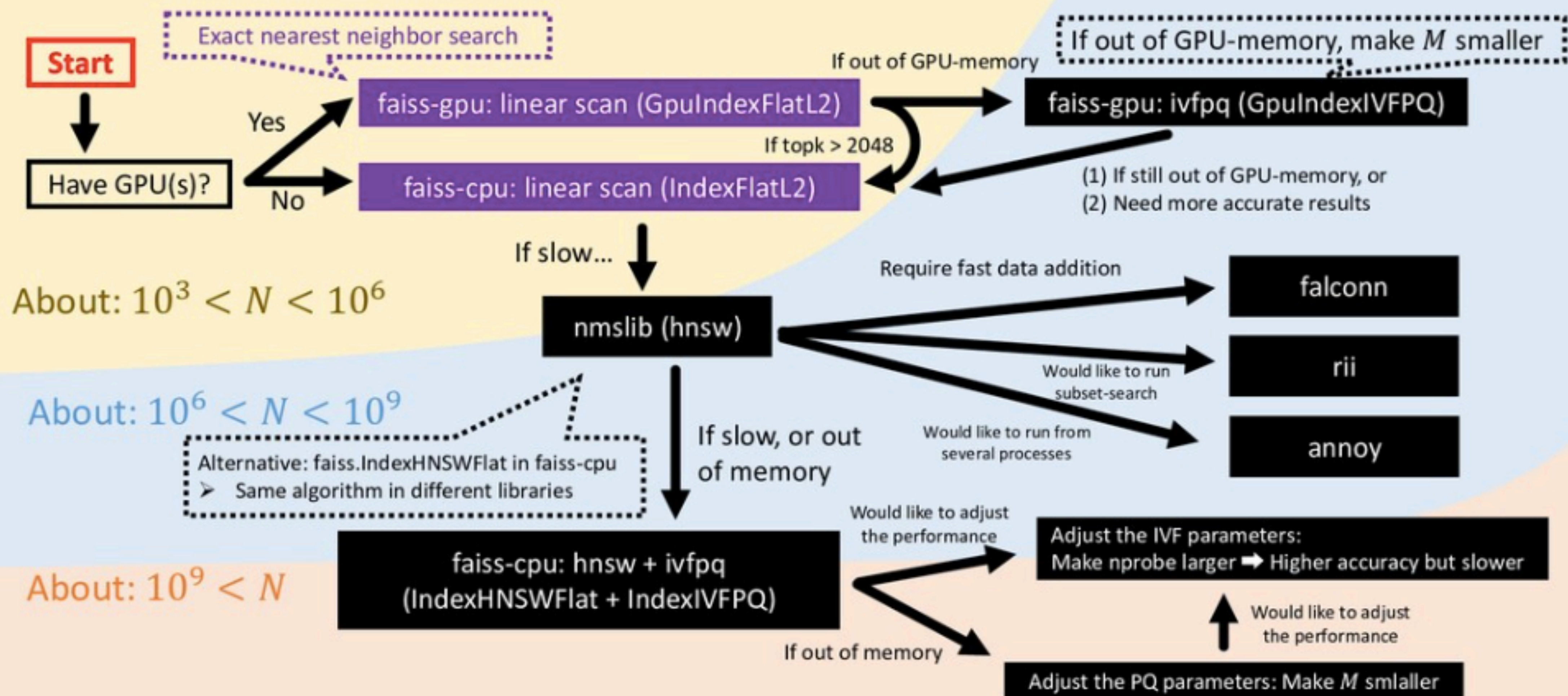
i points to 3721
 \bar{r}_i points to ID: 4

- 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)

cheat-sheet for ANN in Python (as of 2020. Can be installed by conda or pip)

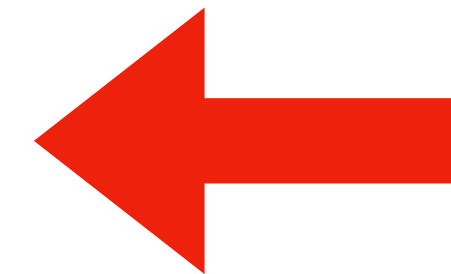


Note: Assuming $D \cong 100$. The size of the problem is determined by DN . If $100 \ll D$, run PCA to reduce D to 100

I. Introduction

II. Prerequisites

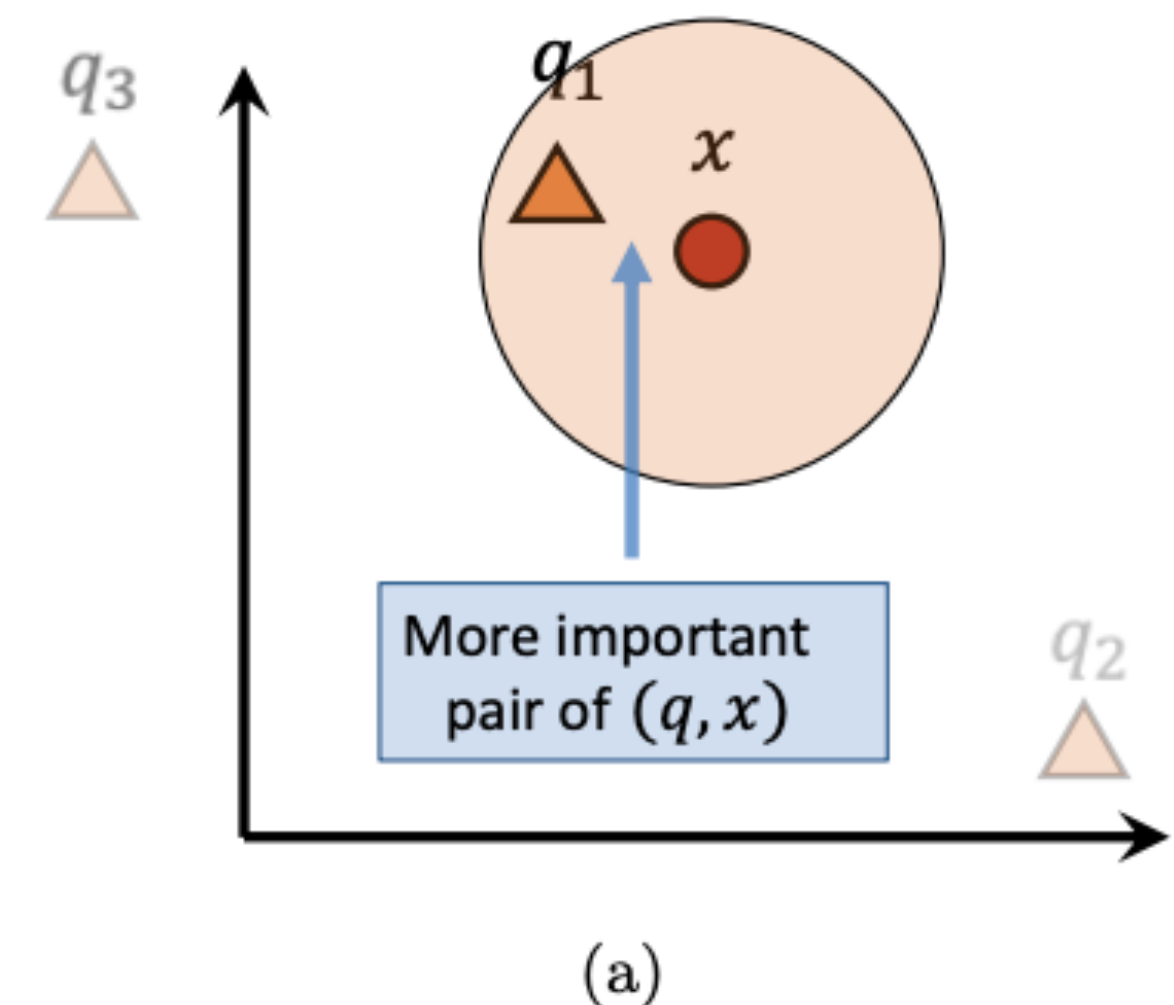
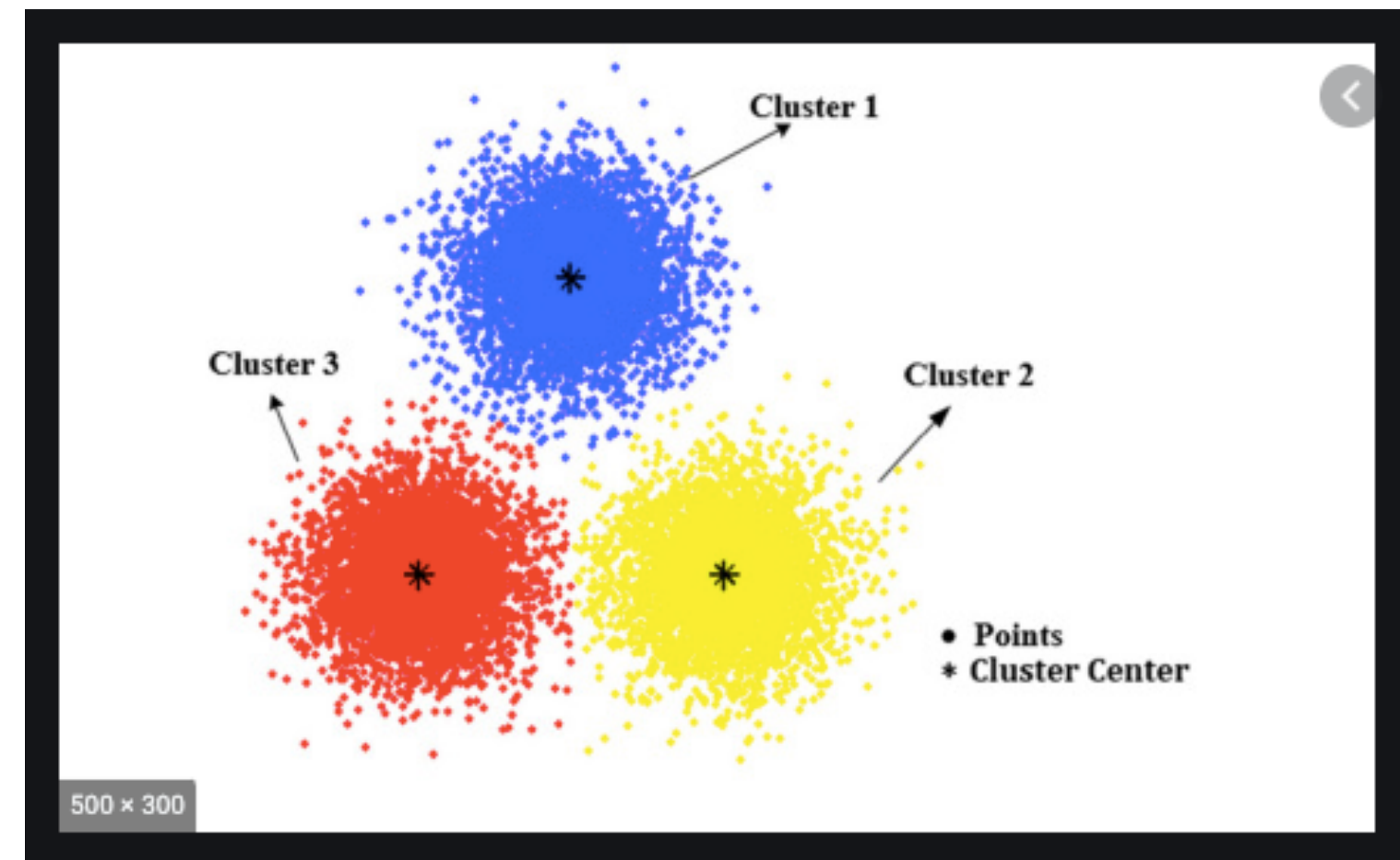
III. Main Topic



Looking back on the objective

- Quantization is the key (e.g. IVFPQ)
- Consider MIPS: $\text{dist}(x, y) = \langle x, y \rangle$
- Quantization Error: $E(c) = \frac{1}{N} \sum_i l(x_i, c(x_i))$
- $$l_{\text{naive}}(x_i, c(x_i)) = \mathbb{E}_q \left[\left(\langle q, x_i \rangle - \langle q, c(x_i) \rangle \right)^2 \right]$$

$$= \mathbb{E}_q \left[\langle q, x_i - c(x_i) \rangle^2 \right] \propto \sum_i |x_i - c(x_i)|^2$$
- Assumed q is isotropic, i.e., $\mathbb{E}_q [qq^T] \propto I$
- K-means is valid
- $$l_{\text{new}}(x_i, c(x_i), w) = \mathbb{E}_q \left[w \left(\langle q, x_i \rangle \right) \langle q, x_i - c(x_i) \rangle^2 \right]$$
 - Score-aware quantization loss**
 - Weight function $w : \mathbb{R} \rightarrow \mathbb{R}^+$
 - Sigmoid
 - Step function: $\mathbf{I}(t \geq T)$



Anisotropic vector quantization problem (Parallel is more important than orthogonal)

Theorem 3.2. Suppose we are given a datapoint x_i , its quantization \tilde{x}_i , and a weight function w . Assuming that query q is uniformly distributed in the d -dimensional unit sphere, the score-aware quantization loss equals

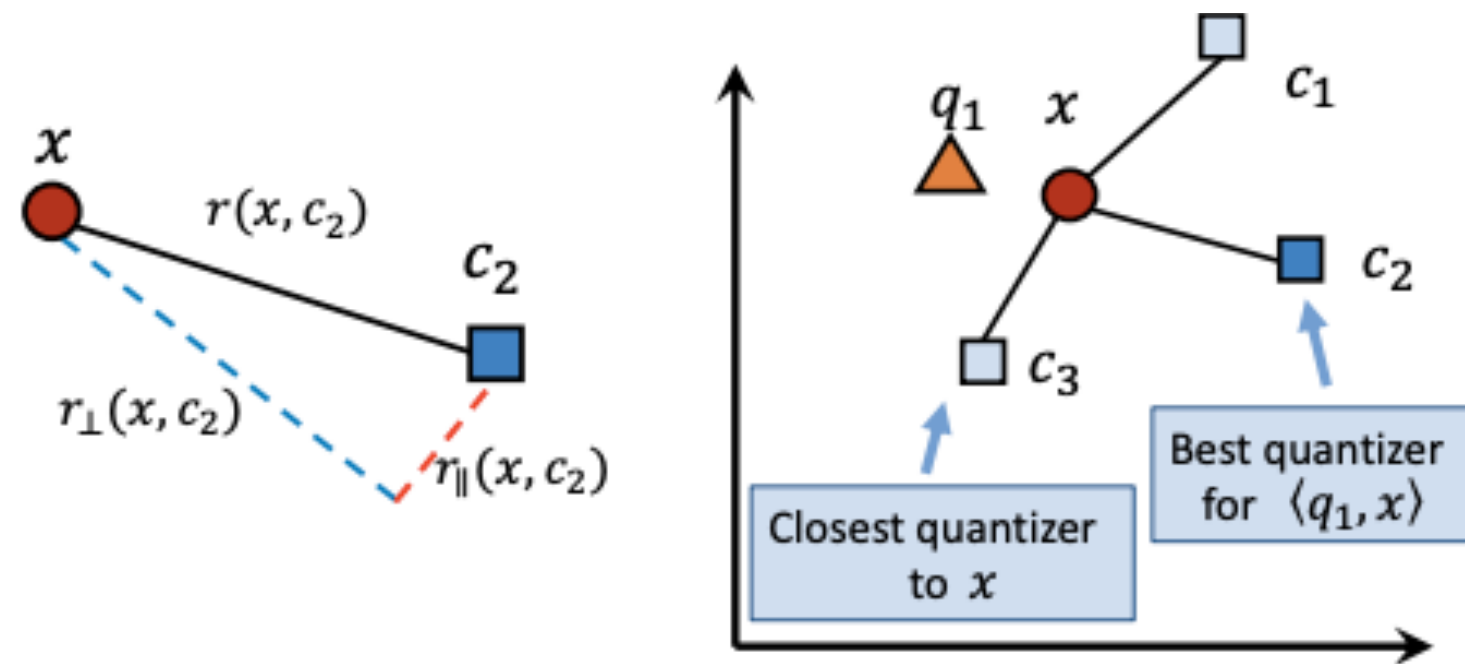
$$\ell(x_i, \tilde{x}_i, w) = h_{\parallel}(w, \|x_i\|) \|r_{\parallel}(x_i, \tilde{x}_i)\|^2 + h_{\perp}(w, \|x_i\|) \|r_{\perp}(x_i, \tilde{x}_i)\|^2$$

with h_{\parallel} and h_{\perp} defined as follows:

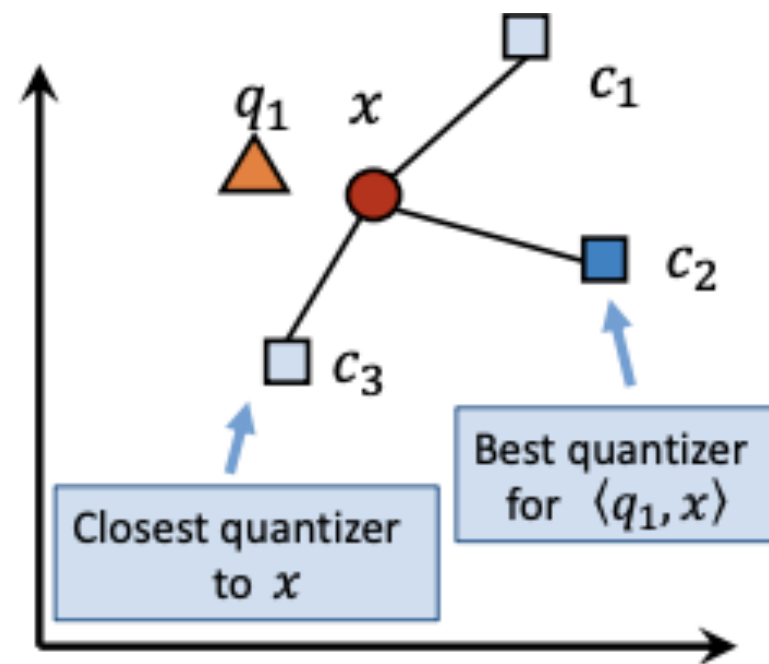
$$h_{\parallel} := (d-1) \int_0^{\pi} w(\|x_i\| \cos \theta) (\sin^{d-2} \theta - \sin^d \theta) d\theta$$

$$h_{\perp} := \int_0^{\pi} w(\|x_i\| \cos \theta) \sin^d \theta d\theta.$$

Proof. See Appendix Section 7.1. \square



(b)



(c)

Theorem 3.3. For any w such that $w(t) = 0$ for $t < 0$ and $w(t)$ is monotonically non-decreasing for $t \geq 0$,

$$h_{\parallel}(w, \|x_i\|) \geq h_{\perp}(w, \|x_i\|)$$

with equality if and only if $w(t)$ is constant for $t \in [-\|x_i\|, \|x_i\|]$.

Proof. See Appendix Section 7.2. \square

$$\begin{aligned} \ell(x_i, \tilde{x}_i, w) &= h_{\parallel}(w, \|x_i\|) \|r_{\parallel}(x_i, \tilde{x}_i)\|^2 + h_{\perp}(w, \|x_i\|) \|r_{\perp}(x_i, \tilde{x}_i)\|^2 \\ &\propto \eta(w, \|x_i\|) \|r_{\parallel}(x_i, \tilde{x}_i)\|^2 + \|r_{\perp}(x_i, \tilde{x}_i)\|^2 \end{aligned}$$

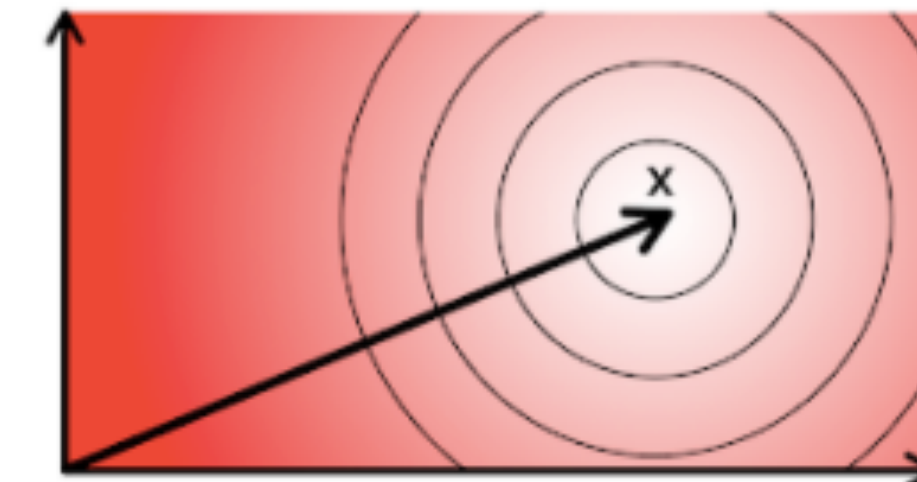
$$\text{where } \eta(w, \|x_i\|) := \frac{h_{\parallel}(w, \|x_i\|)}{h_{\perp}(w, \|x_i\|)}.$$

Theorem 3.4.

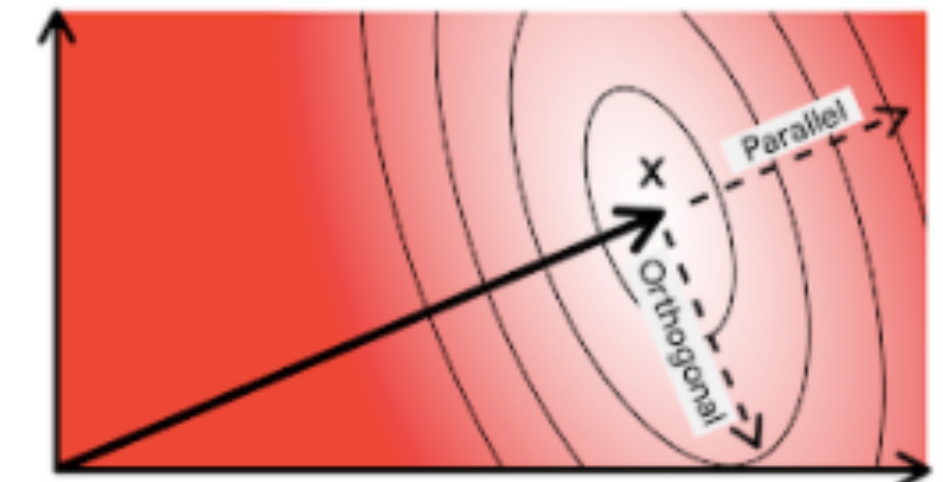
$$\lim_{d \rightarrow \infty} \frac{\eta(\mathbf{I}(t \geq T), \|x_i\|)}{d-1} = \frac{(T/\|x_i\|)^2}{1 - (T/\|x_i\|)^2} \quad (3)$$

Proof. See Appendix Section 7.3. \square

As special cases, when $T = 0$, $\eta(\mathbf{I}(t \geq 0), \|x_i\|) = 1$ which implies parallel and orthogonal errors are weighted equally. When $T = \|x_i\|$, we have $\eta(\mathbf{I}(t \geq \|x_i\|), \|x_i\|) = \infty$ which indicates we should only consider parallel error.



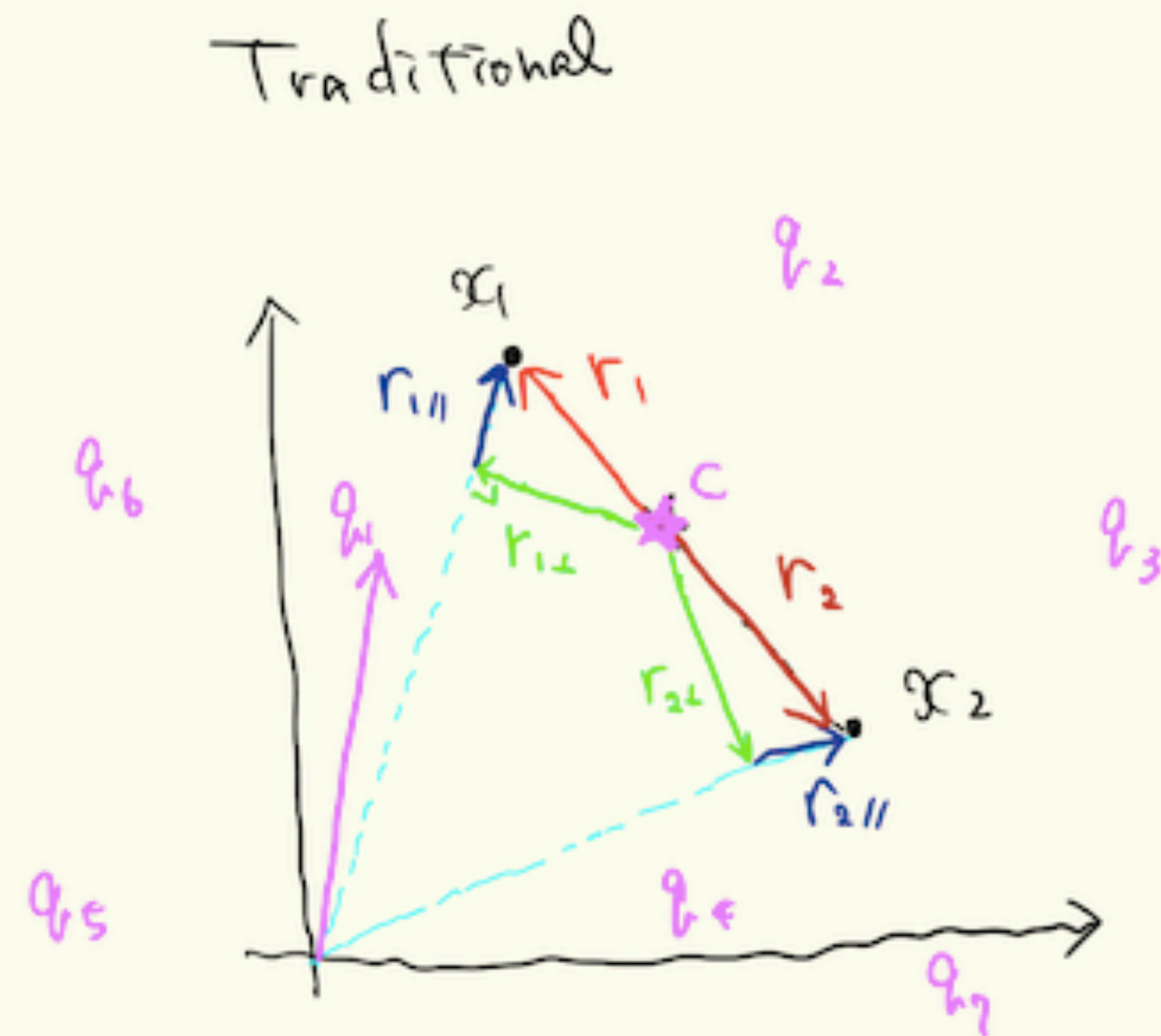
Traditional Loss



Anisotropic Loss

Why score-aware leads to anisotropic loss?

Intuitive explanation

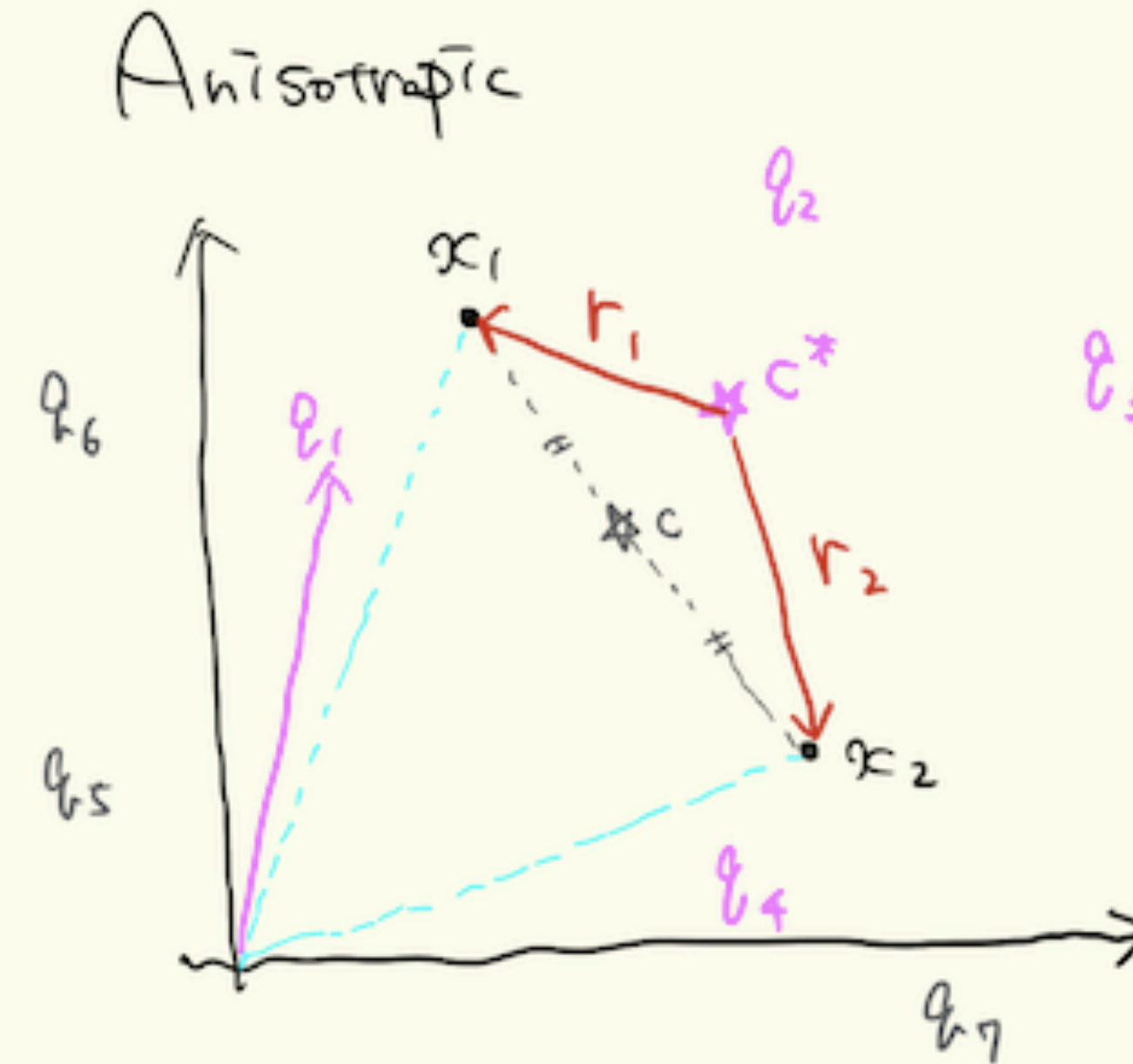


q is everywhere



minimize $\|r\|_2^2 = \|x - c\|_2^2$
 $= \|r_{||}\|_2^2 + \|r_{\perp}\|_2^2$

to minimize $\mathbb{E}_q[\langle q, r \rangle]$



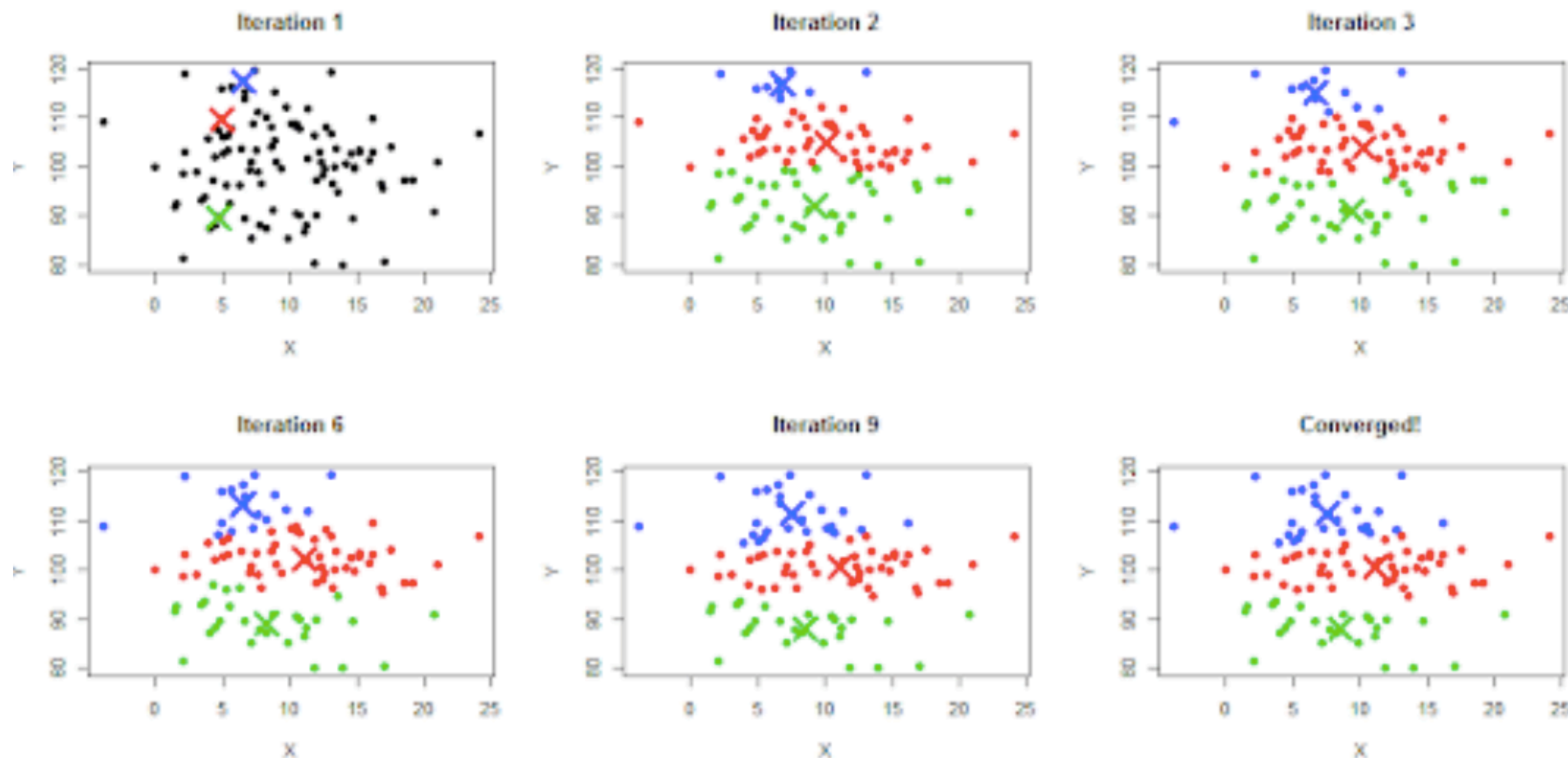
only q close to x matters



smaller $r_{||}$ is better

to minimize $\mathbb{E}_q[\omega(\langle q, x \rangle) \langle q, r \rangle]$

How to update centroids with anisotropic loss



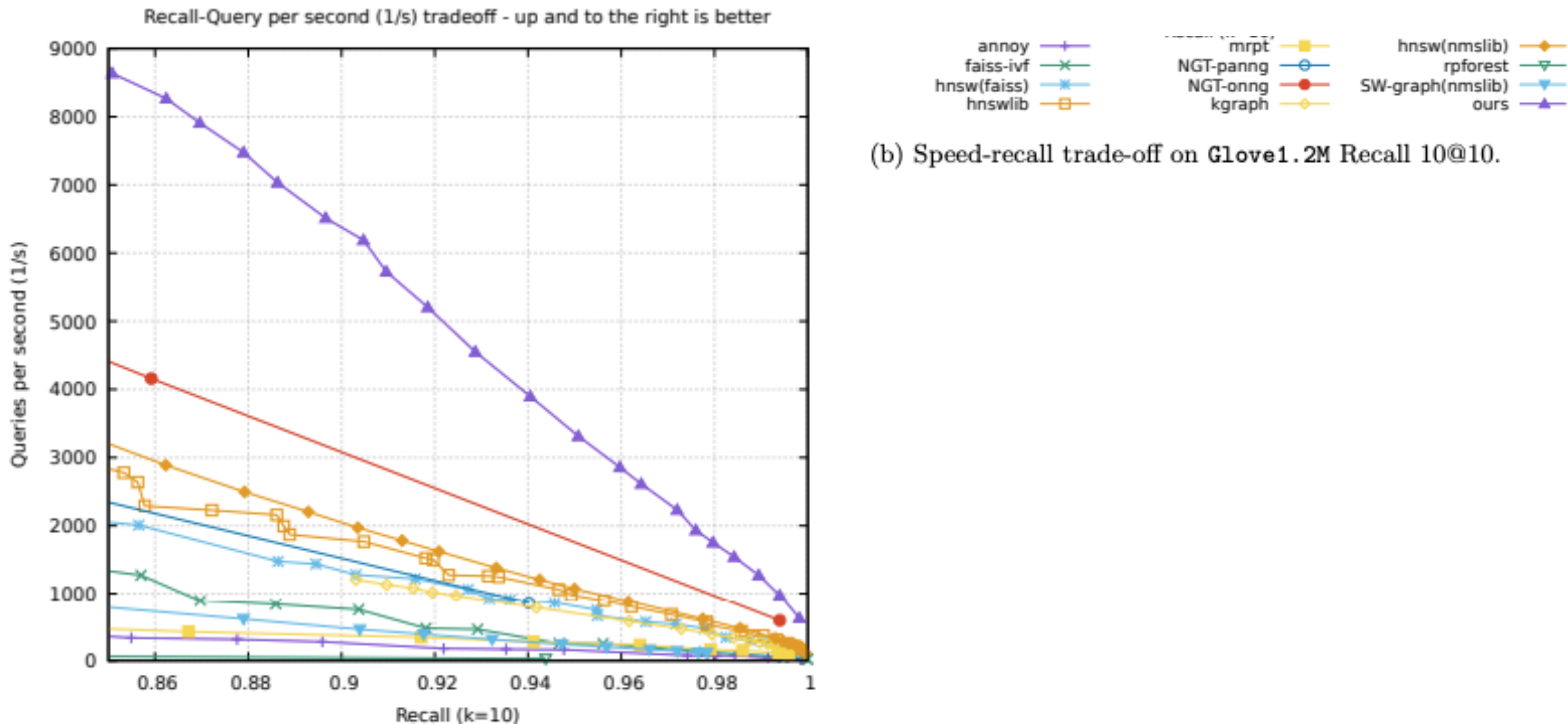
3. (Codebook Update Step) For every codeword c_j , let X_j be all datapoints x_i such that $\tilde{x}_i = c_j$. Update c_j by

$$c_j \leftarrow \arg \min_{c \in \mathbb{R}^d} \sum_{x_i \in X_j} \ell(x_i, c).$$

Theorem 4.2. Optimal codeword c_j can be obtained in closed form by solving the optimization problem in Equation (4) for a partition X_j . The update rule for the codebook is

$$c_j^* = \left(I \sum_{x_i \in X_j} h_{i,\perp} + \sum_{x_i \in X_j} \frac{h_{i,\parallel} - h_{i,\perp}}{\|x_i\|^2} x_i x_i^T \right)^{-1} \sum_{x_i \in X_j} h_{i,\parallel} x_i$$

We win!



Thank you

Memo

- N D-dim vectors look small enough for RAM?
Why compress?
 - $N \cdot D \cdot 4\text{Bytes} = 10^9 \cdot 10^3 \cdot 4 \text{ Bytes} = 4\text{TB}$
 - Just to reduce the server cost?
 - I may be missing something
 - My reasoning
 - Billion-scale app has high QPS
 - Multiple TB servers are too expensive