

Related topics

Hardware-based acceleration

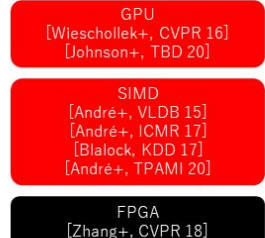
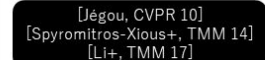
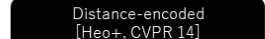


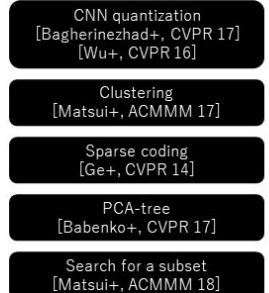
Image search with PQ



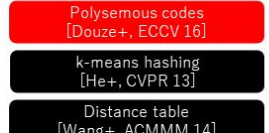
Additional bit management



Applications using PQ



Connection to binary hashing



credit: Yusuke Matsui and FAISS



This CVPR2013 paper is the Open Access version, provided by the Computer Vision Foundation.
The authoritative version of this paper is available in IEEE Xplore.

Optimized Product Quantization for Approximate Nearest Neighbor Search

Tiezheng Ge^{1*}

Kaiming He²

Qifa Ke³

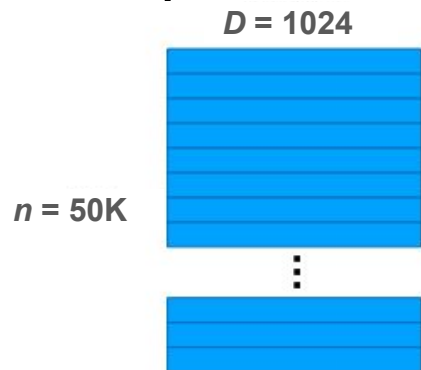
Jian Sun²

¹University of Science and Technology of China

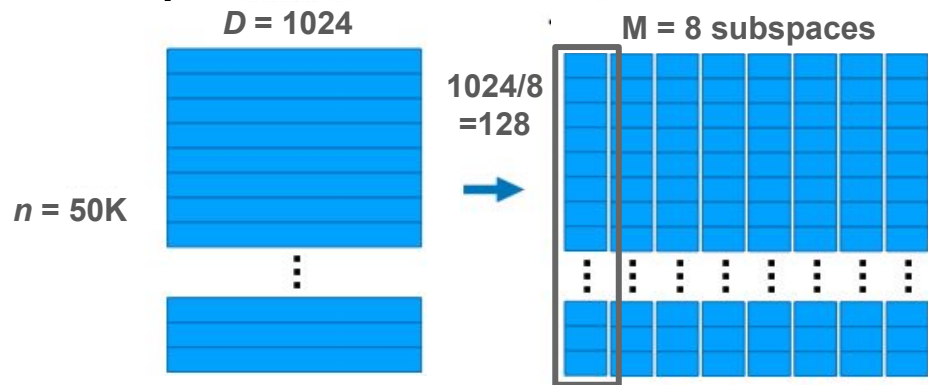
²Microsoft Research Asia

³Microsoft Research Silicon Valley

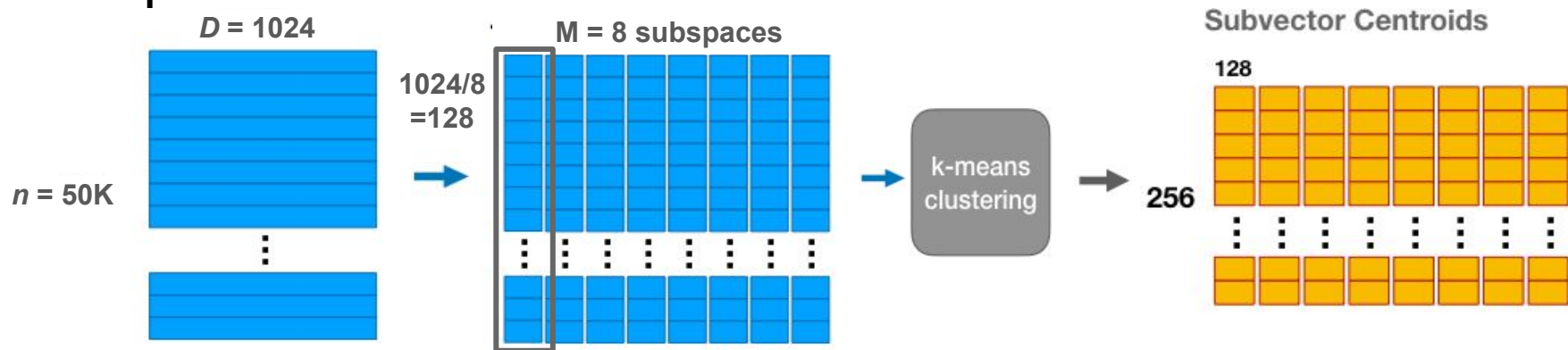
Recap: Product Quantization



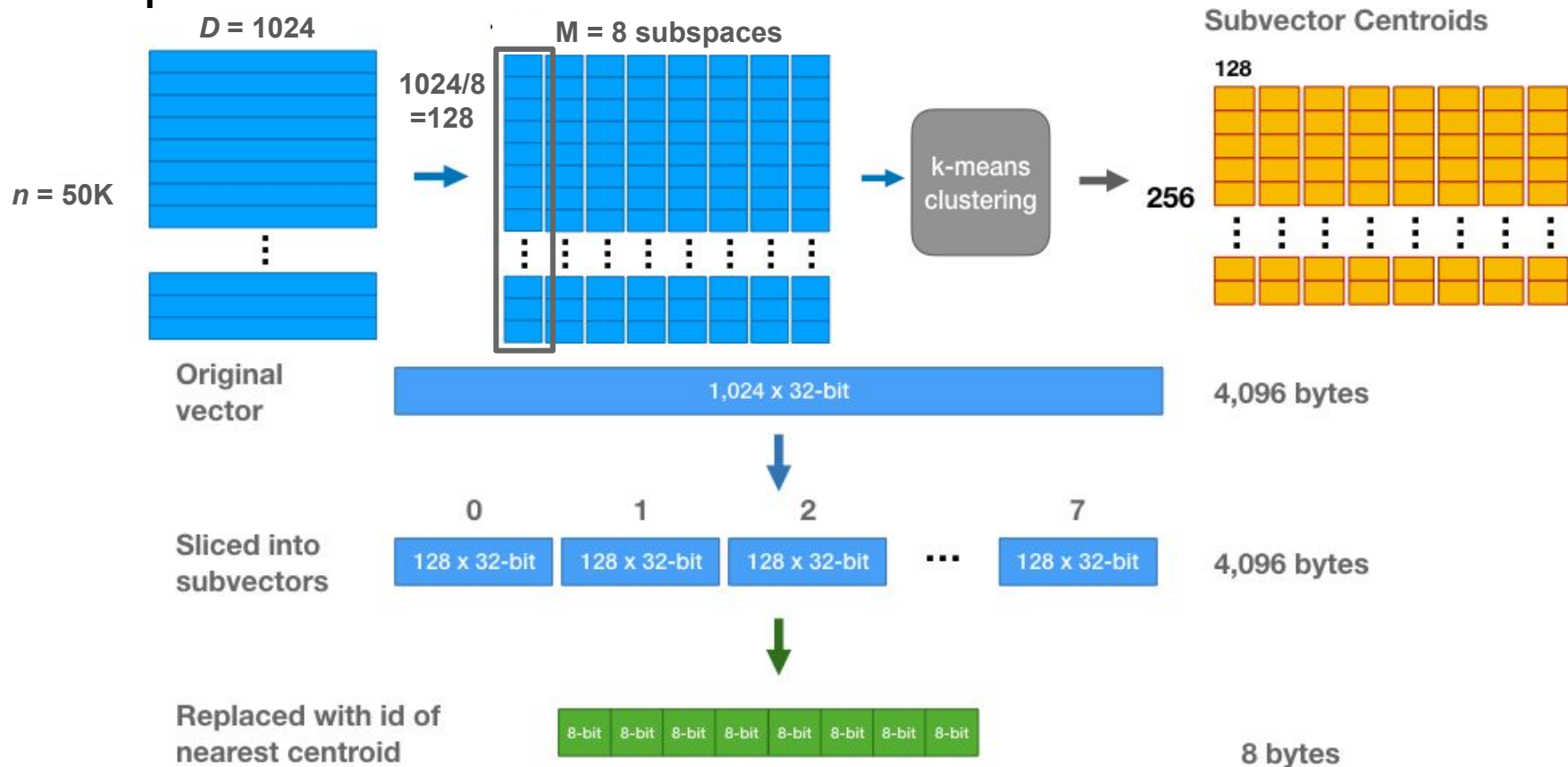
Recap: Product Quantization



Recap: Product Quantization



Recap: Product Quantization



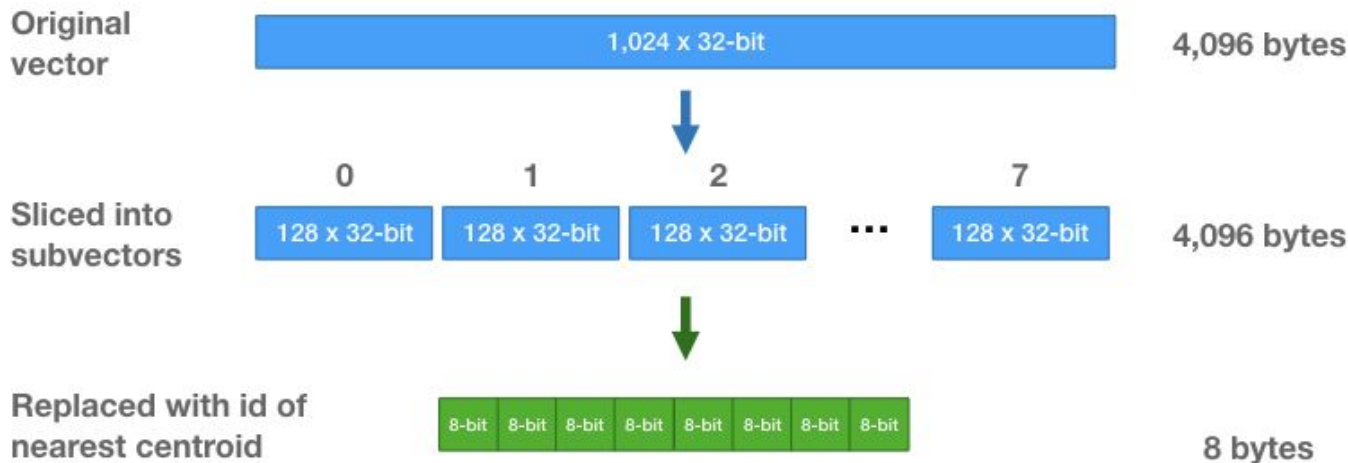
Recap: Product Quantization

- Formally:

- Quantizer: $\mathbf{x} \rightarrow \mathbf{c}(i(\mathbf{x}))$

M subvectors: $\mathbf{x} = [\mathbf{x}^1, \dots \mathbf{x}^m, \dots \mathbf{x}^M]$

M sub-codewords: $\mathbf{c} = [\mathbf{c}^1, \dots \mathbf{c}^m, \dots \mathbf{c}^M]$



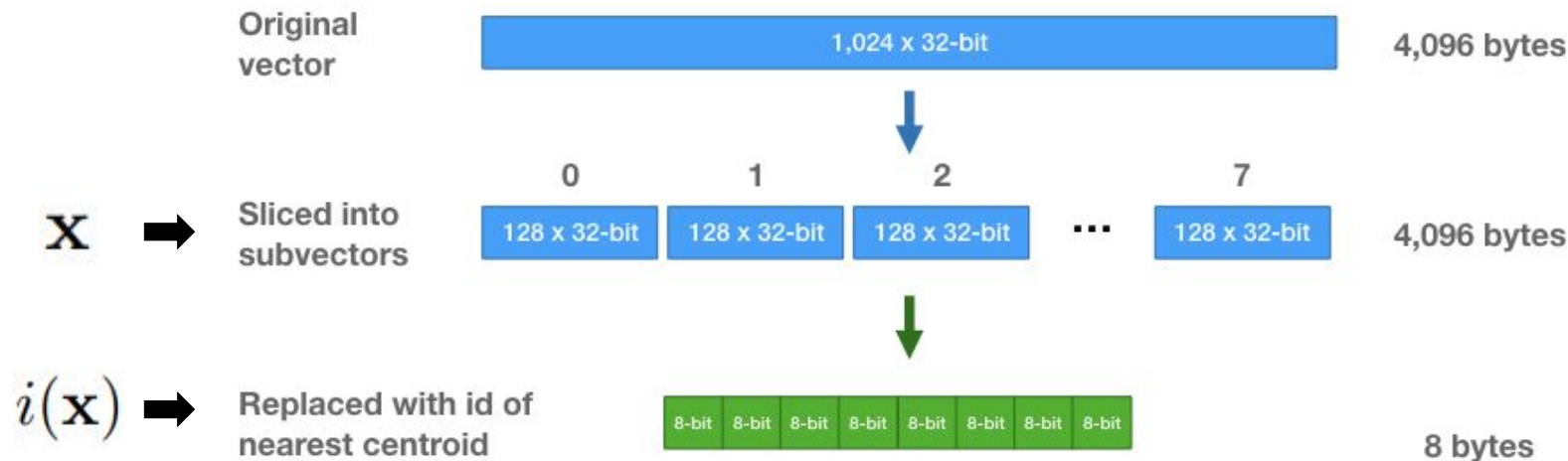
Recap: Product Quantization

- Formally:

- Quantizer: $\mathbf{x} \rightarrow \mathbf{c}(i(\mathbf{x}))$

M subvectors: $\mathbf{x} = [\mathbf{x}^1, \dots \mathbf{x}^m, \dots \mathbf{x}^M]$

M sub-codewords: $\mathbf{c} = [\mathbf{c}^1, \dots \mathbf{c}^m, \dots \mathbf{c}^M]$



Quantization Distortion

- Formally:

- Quantizer: $\mathbf{x} \rightarrow \mathbf{c}(i(\mathbf{x}))$

M subvectors: $\mathbf{x} = [\mathbf{x}^1, \dots, \mathbf{x}^m, \dots, \mathbf{x}^M]$

M sub-codewords: $\mathbf{c} = [\mathbf{c}^1, \dots, \mathbf{c}^m, \dots, \mathbf{c}^M]$

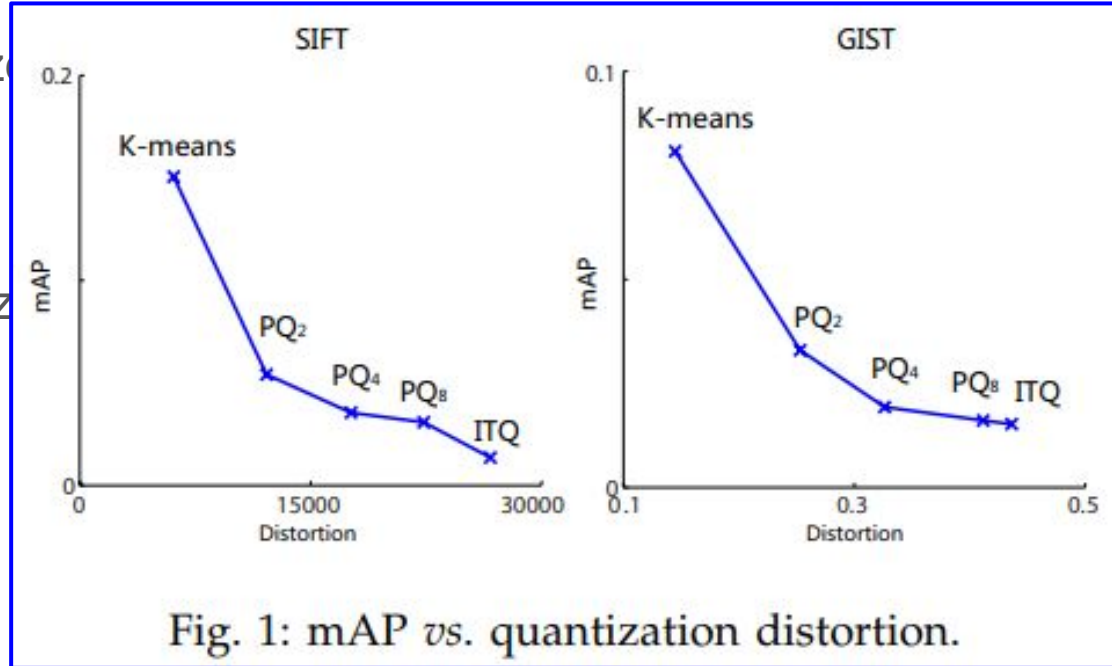
- PQ's quantization distortion (i.e. loss function)

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

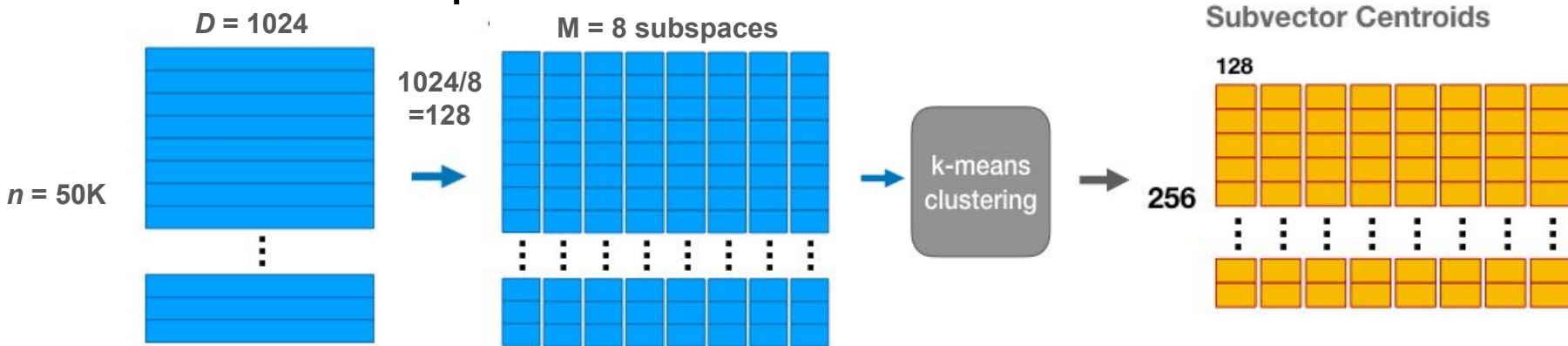
$$s.t. \quad \mathbf{c} \in \mathcal{C} = \mathcal{C}^1 \times \dots \times \mathcal{C}^M.$$

Quantization Distortion

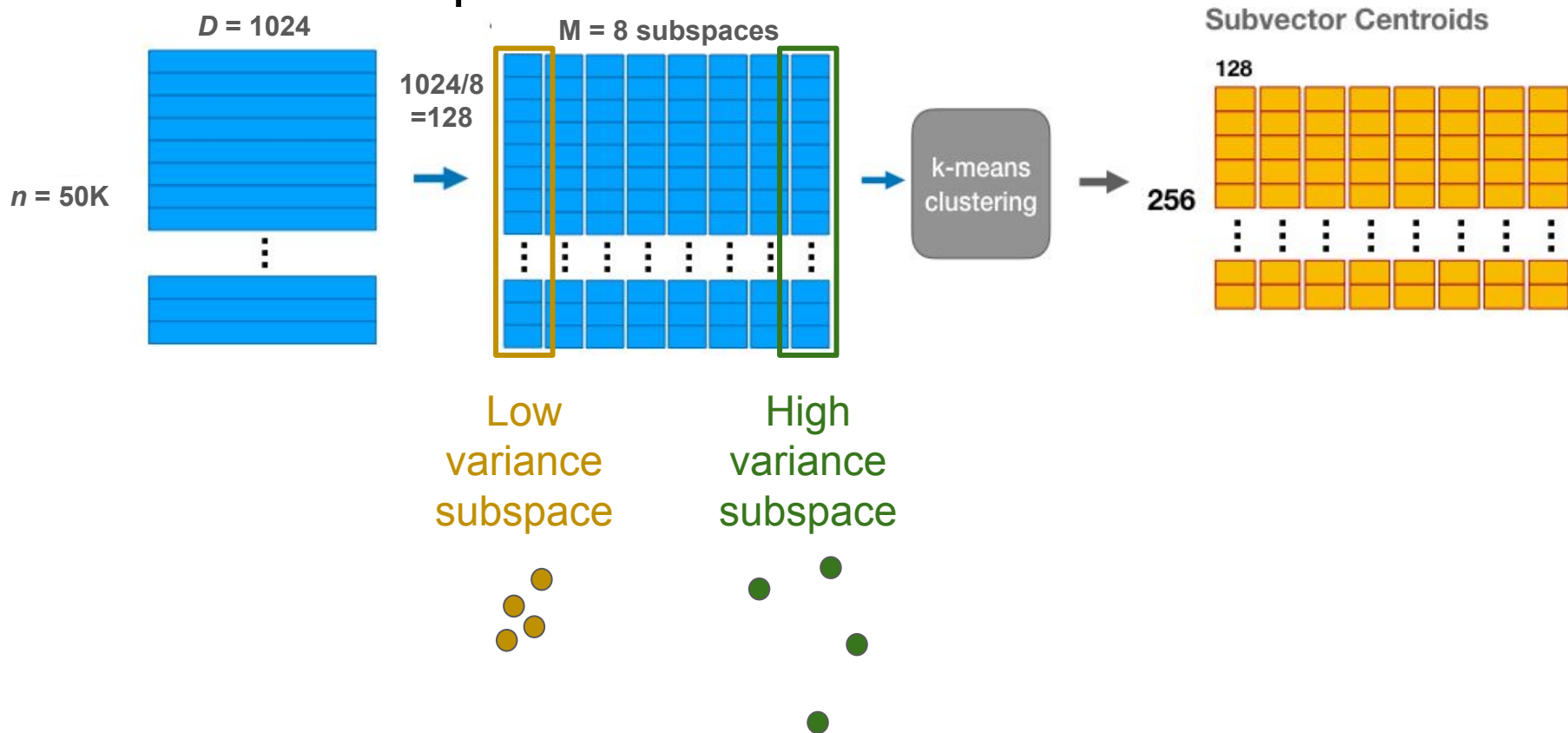
- Formally:
 - Quantization
- PQ's quantization



How can we improve PQ?



How can we improve PQ?



How can we improve PQ?

$D = 1024$

$M = 8$ subspaces

Subvector Centroids

128

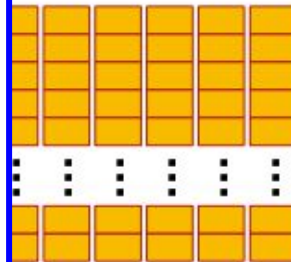
$n = 50K$

⋮

From PQ:

“...to provide good quantization ...each subvector should have...comparable energy”

→ multiply by a random rotation matrix”



How can we improve PQ?

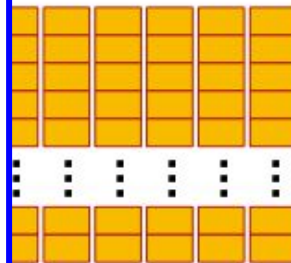
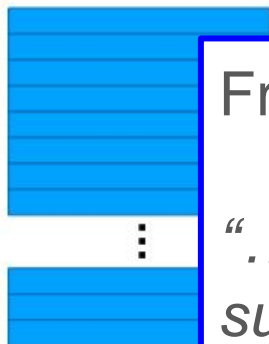
$D = 1024$

$M = 8$ subspaces

Subvector Centroids

128

$n = 50K$



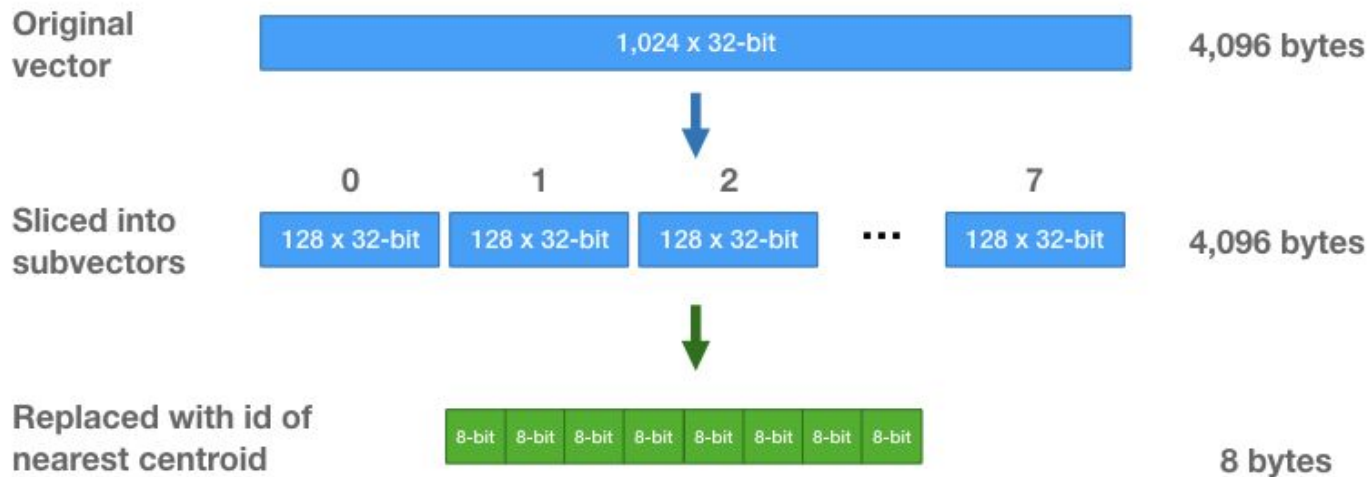
From PQ:

“...to provide good quantization ...each subvector should have...comparable energy”

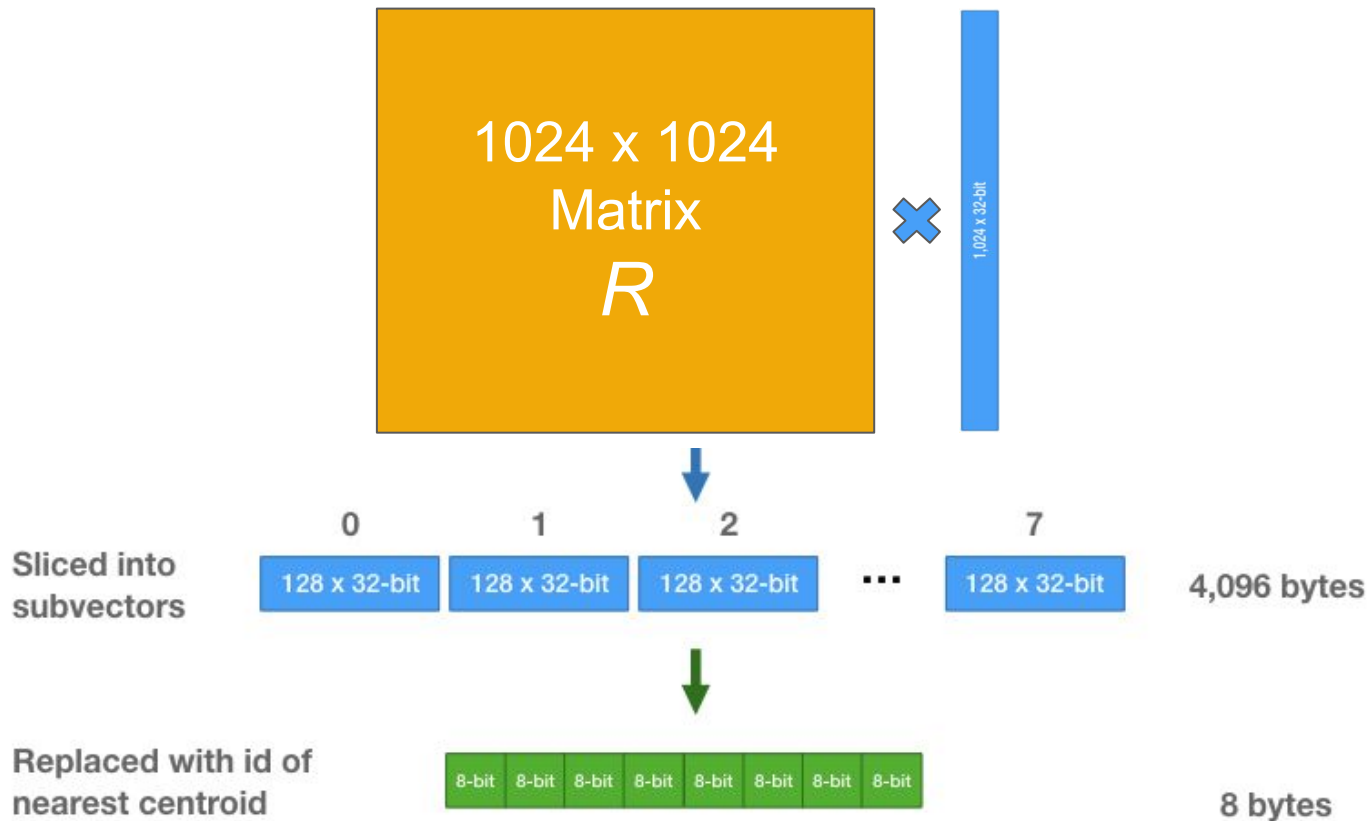
→ multiply by a random rotation matrix”

OPQ: Find optimal rotation matrix

Optimized Product Quantization



Optimized Product Quantization



Optimized Quantization Distortion

- PQ's quantization distortion (i.e. loss function)

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$
$$s.t. \quad \mathbf{c} \in \mathcal{C} = \mathcal{C}^1 \times \dots \times \mathcal{C}^M.$$

- Optimized PQ proposes to minimize:

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$
$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

How do we find matrix R ?

1. Non-parametric
 - Optimize two easier subproblems
2. Parametric
 - Gaussian assumption
 - Still works on non-Gaussian data

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

1. Fix R , optimize codebooks

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\hat{\mathbf{x}}} \|\hat{\mathbf{x}} - \hat{\mathbf{c}}(i(\hat{\mathbf{x}}))\|^2,$$

$$s.t. \quad \hat{\mathbf{c}} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M.$$

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

1. Fix R , optimize codebooks

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\hat{\mathbf{x}}} \|\hat{\mathbf{x}} - \hat{\mathbf{c}}(i(\hat{\mathbf{x}}))\|^2,$$

$$s.t. \quad \hat{\mathbf{c}} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M.$$

2. Fix codebooks, optimize R

$$\min_R \sum_{\mathbf{x}} \|R\mathbf{x} - \hat{\mathbf{c}}(i(\hat{\mathbf{x}}))\|^2,$$

$$s.t. \quad R^T R = I.$$

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

1. Fix R , optimize codebooks

Same as PQ

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

\rightarrow Solve with K-means

2. Fix codebooks, optimize R

$$\min_R \sum_{\mathbf{x}} \|R\mathbf{x} - \hat{\mathbf{c}}(i(\hat{\mathbf{x}}))\|^2,$$

$$s.t. \quad R^T R = I.$$

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, \quad R^T R = I\}$$

1. Fix R , optimize codebooks

Same as PQ

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

\rightarrow Solve with
K-means

2. Fix codebooks, optimize R

$$\min_R \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$XY^T = USV^T$$

$$R = VU^T$$

Non-parametric solution to minimize distortion

$$\min_{R, \mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$s.t. \quad \mathbf{c} \in \mathcal{C} = \{\mathbf{c} \mid R\mathbf{c} \in \mathcal{C}^1 \times \dots \times \mathcal{C}^M, R^T R = I\}$$

1. Fix R , optimize codebooks

Same as PQ

$$\min_{\mathcal{C}^1, \dots, \mathcal{C}^M} \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

\rightarrow Solve with
K-means

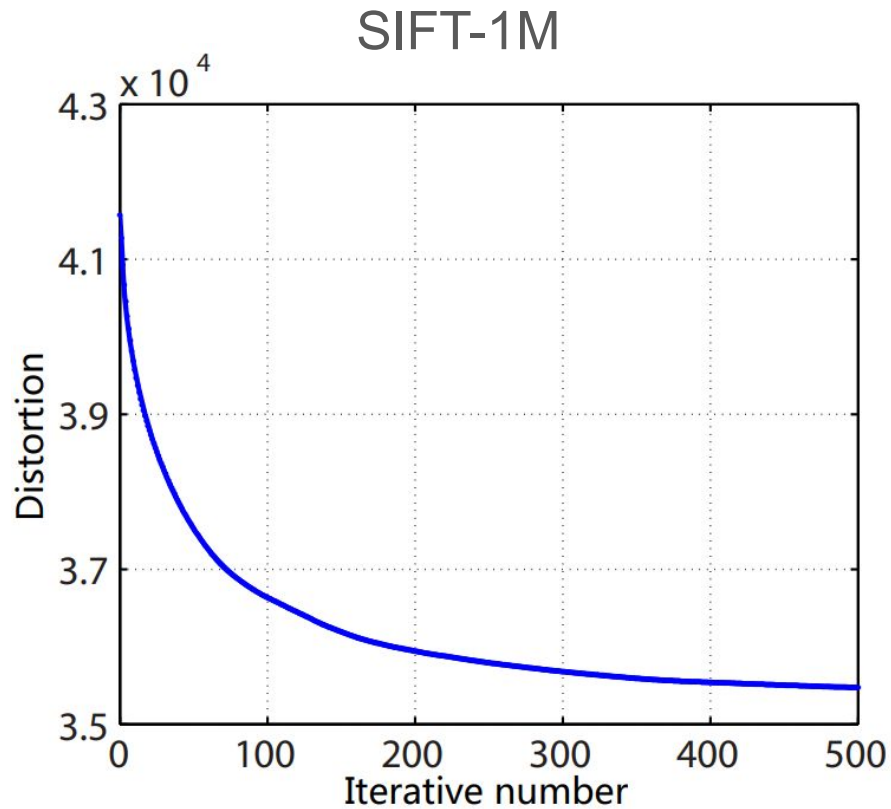
2. Fix codebooks, optimize R

$$\min_R \sum_{\mathbf{x}} \|\mathbf{x} - \mathbf{c}(i(\mathbf{x}))\|^2,$$

$$XY^T = USV^T$$

$$R = VU^T$$

Non-parametric solution to minimize distortion



Parametric Solution

- If data is Gaussian, distortion E of PQ is:

$$E_{\text{PQ}} = k^{-\frac{2M}{D}} \frac{D}{M} \sum_{m=1}^M |\Sigma_{mm}|^{\frac{M}{D}}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \cdots & \Sigma_{1M} \\ \vdots & \ddots & \vdots \\ \Sigma_{M1} & \cdots & \Sigma_{MM} \end{pmatrix}.$$

- Lower bound of distortion:

$$\sum_{m=1}^M |\hat{\Sigma}_{mm}|^{\frac{M}{D}} \geq M |\Sigma|^{\frac{1}{D}}.$$

Parametric Solution

- If data is Gaussian, distortion E of PQ is:

$$E_{PQ} = k^{-\frac{2M}{D}}$$

Minimal distortion with:

1. Vector dimension independence
2. Balanced subspaces variance

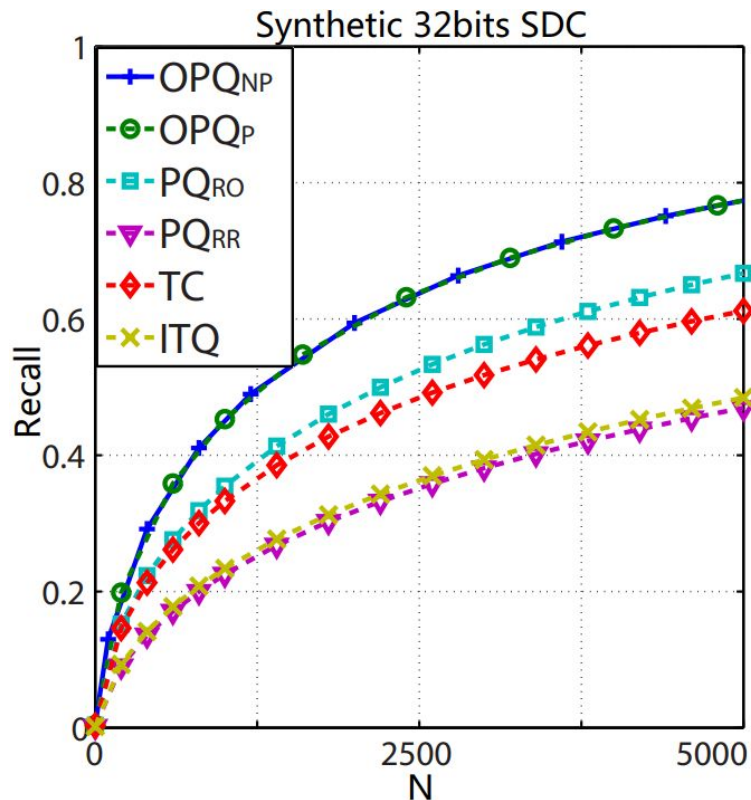
$m=1$

$$\begin{pmatrix} \cdots & \Sigma_{1M} \\ \ddots & \vdots \\ \cdots & \Sigma_{MM} \end{pmatrix}.$$

Evaluation

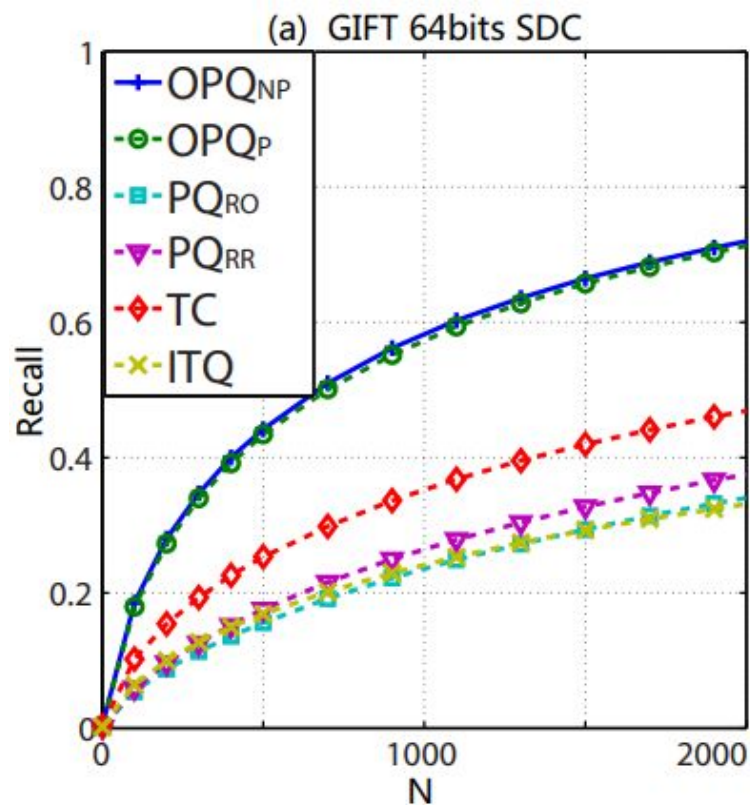
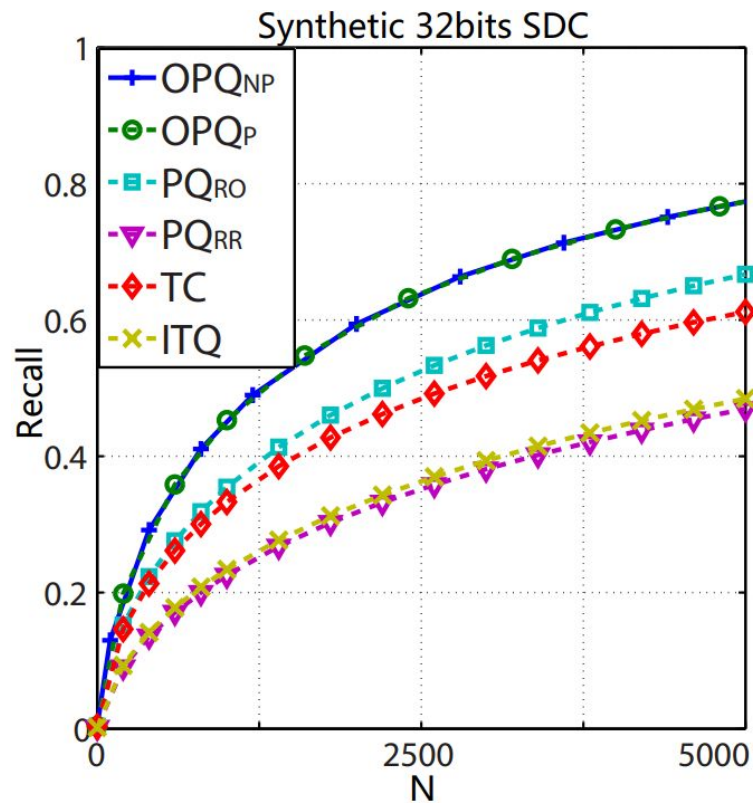
- Datasets: Synthetic Gaussian dataset, GIST1M, SIFT1M, MNIST,
- Compare OPQ_P and OPQ_NP with:
 - PQ_RO: randomly ordered dimensions
 - PQ_RR: PCA alignment then random rotation
 - TC: scalar quantizer
 - ITQ: vector quantizer

Evaluation

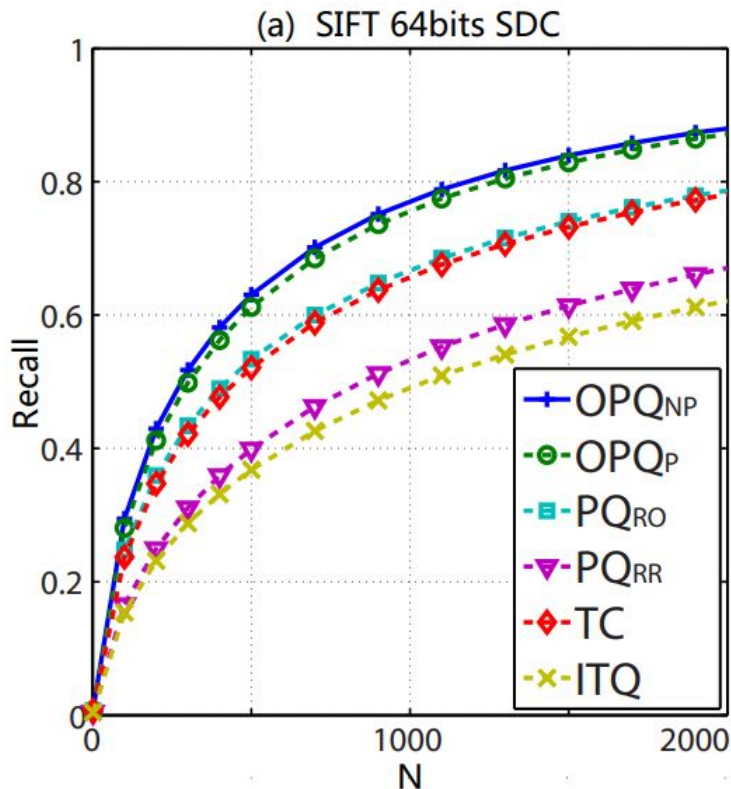


- Synthetic Gaussian data
 - 128D
 - 1M data points
- PQ_{RO} better than PQ_{RR}
 - Vector dimension independence

Evaluation

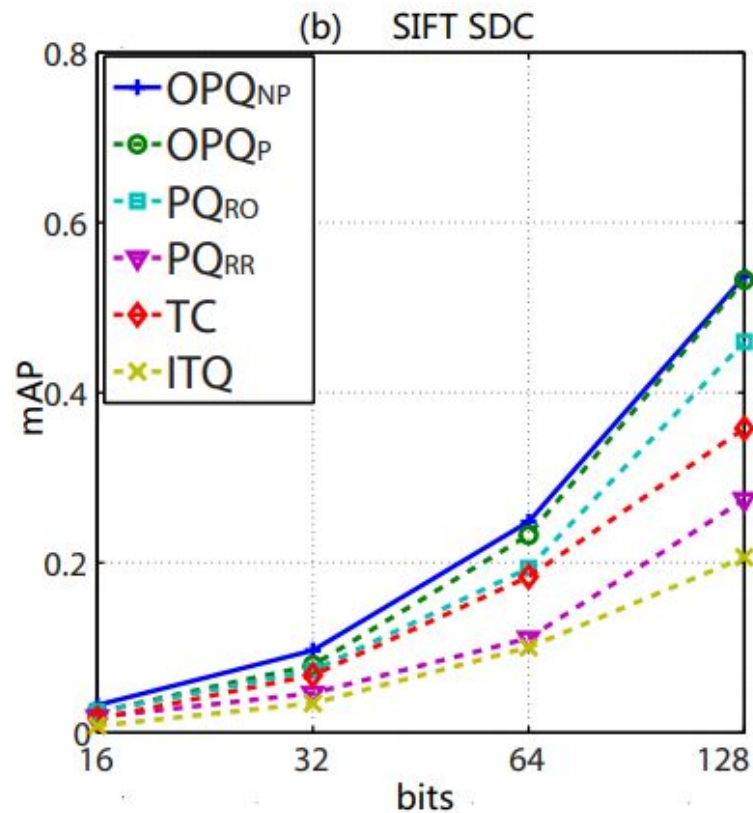
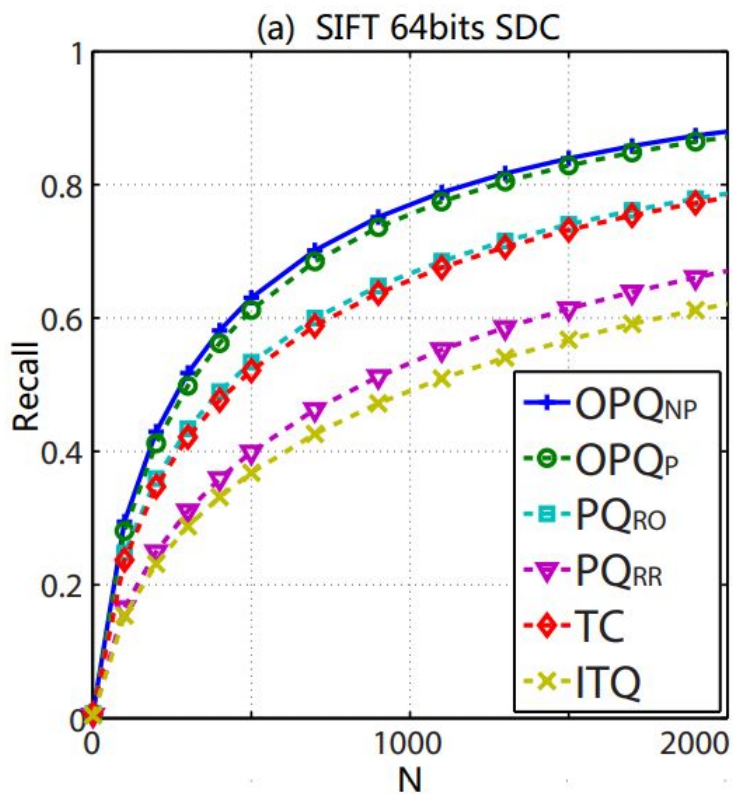


Evaluation

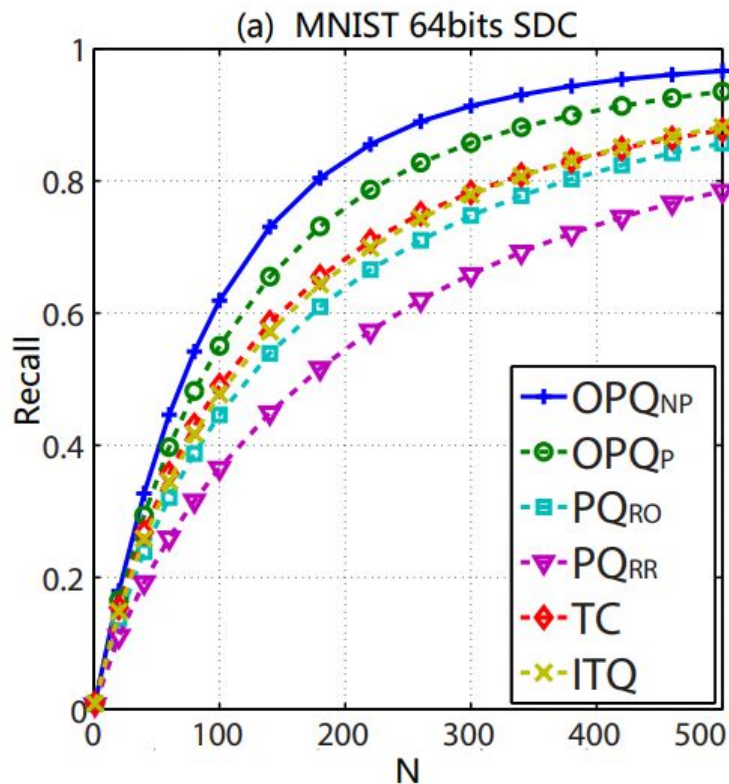


- SIFT1M: two distinct clusters
- OPQ_{NP} begins to outperform OPQ_P

Evaluation



Evaluation



- MNIST: ten distinct clusters
- Greater difference between OPQ_{NP} and OPQ_P
- More complicated data?

Takeaways

- PQ is sensitive to data distribution!
- Optimizing transformation matrix can improve PQ accuracy
- Gaussian solution:
 - Independent vectors dimensions
 - Balanced subspace variance
- Limitations:
 - No non-parametric convergence guarantee
 - No evaluation of overhead
 - Gaussian assumption

Extras: Overhead

Table 4. The Indexing Time for the GIST Dataset

	RaBitQ	PQ	OPQ	LSQ
Time	117s	105s	291s	time-out (>24 hours)

RaBitQ: Quantizing High-Dimensional Vectors with a Theoretical Error Bound for Approximate Nearest Neighbor Search, JIANYANG GAO, CHENG LONG