

Advances in Product Quantization

RaBitQ & TurboQuant as examples

大数据算法课程助教
陈磊



Contents

- Discussions
- Introduction for PQ
- Background and notations
- RabbitQ
- TurboQuant



01 Discussions



中国科学技术大学
University of Science and Technology of China



Discussions on OpenReview

- Proceedings of the ACM on Management of Data, 2024

RaBitQ: Quantizing High-Dimensional Vectors with a Theoretical Error Bound for Approximate Nearest Neighbor Search

Jiayang Gao
jiayang.gao@ntu.edu.sg
Nanyang Technological University
Singapore

Cheng Long*
c.long@ntu.edu.sg
Nanyang Technological University
Singapore

- ICLR 2026 Poster

TurboQuant: Online Vector Quantization with Near-optimal Distortion Rate

Amir Zandieh
Google Research
zandieh@google.com

Majid Daliri
New York University
daliri.majid@nyu.edu

Majid Hadian
Google DeepMind
majidh@google.com

Vahab Mirrokni
Google Research
mirrokni@google.com

Discussions on OpenReview

- Comparisons between **different strategies** for (1) codebook construction and (2) unbiased estimations.
- **Imcomplete discussion** for RabbitQ?
- **Ungrounded theoretical discussion?**
- **Transparent empirical comparisons** required?
- Dive into these two papers, learn some new algorithm designs, and something to take care about.

TurboQuant

<https://openreview.net/forum?id=tO3ASKZlok>

Introduction for Product Quantization (PQ)

- Why product? Dimension reduction and clustering.
- Mutual transformation between **2-norms and cosine similarity**:

$$\begin{aligned}\|\mathbf{o}_r - \mathbf{q}_r\|^2 &= \|(\mathbf{o}_r - \mathbf{c}) - (\mathbf{q}_r - \mathbf{c})\|^2 \\ &= \|\mathbf{o}_r - \mathbf{c}\|^2 + \|\mathbf{q}_r - \mathbf{c}\|^2 - 2 \cdot \|\mathbf{o}_r - \mathbf{c}\| \cdot \|\mathbf{q}_r - \mathbf{c}\| \cdot \langle \mathbf{q}, \mathbf{o} \rangle\end{aligned}$$

- Application:
 - attention scores, recommendations.
 - Kv cache quantization, information retrieval (RAG)
 - Approximation nearest neighbor search (ANN)

Introduction for PQ

- Product quantization (**PQ**) is a family of popular methods for approximate nearest neighbors (**ANN**), while ANN is a relaxed counterpart of **NN**.
- The quality of ANN (PQ algorithms) affects the recall rate after re-ranking (in recommendation system).
- Given data vectors and query vector(s).
- Two phases in PQ:
 - **Index phase**: construct quantization codebook, map each data vector to its quantized vector.
 - **Query phase**: pre-compute $\langle \text{query vector}, \text{quantized vectors} \rangle$, assign as the estimated distance of each data vectors.

Challenges of PQ and proposal of RabbitQ

- The **same** estimated distance for data vectors within same cluster. (Does RabbitQ changes that?)
- Hard to quantify the **approximation error** during codebook construction and distance estimation
- It implies that PQ can unpredictably fail anytime, moderately or severely.
- RabbitQ:
 - Provide geometric **relationship** between $\langle \text{query}, \text{data} \rangle$ and $\langle \text{query}, \text{quantified data} \rangle$.
 - Provide **unbiased estimator** for $\langle \text{query}, \text{data} \rangle$ based on that geometric relationship.
 - Provide **error bound** analysis.
 - Practical implementations.

Basic idea of RabbitQ

- **Normalize**: the normalized data vectors are expected to spread evenly on the unit hypersphere
- Intuitively construct codebook: also spread evenly on the unit hypersphere (and think about the benefit in theory)
- Inject some **random rotations** to remove inductive bias.
- **Assign** Quantized Codes of Data Vectors.
 - Index phase \uparrow | query phase \downarrow
- **Quantizing** query vector (uniform scalar quantization with $\log\log D$ resolution and randomness)
- Efficiently compute as a weighted sum of the inner product of the binary vectors

RabitQ preliminary

Normalize	$\mathbf{o} := \frac{\mathbf{o}_r - \mathbf{c}}{\ \mathbf{o}_r - \mathbf{c}\ }$	$\mathbf{o}_r, \mathbf{q}_r$	The raw data and query vectors.
		\mathbf{o}, \mathbf{q}	The normalized data and query vectors.
Construct	$C := \left\{ +\frac{1}{\sqrt{D}}, -\frac{1}{\sqrt{D}} \right\}^D$	C, C_{rand}	The quantization codebook, its randomized version.
	$C_{rand} := \{P\mathbf{x} \mathbf{x} \in C\}$	P	A random orthogonal transformation matrix.
Assign	$\bar{\mathbf{x}} = \arg \min_{\mathbf{x} \in C} \ \mathbf{o} - P\mathbf{x}\ ^2$	$\bar{\mathbf{x}}$	The code in C s.t. $P\bar{\mathbf{x}}$ is the quantized vector of \mathbf{o} .
by		$\bar{\mathbf{o}}$	The quantized vector of \mathbf{o} in C_{rand} , i.e., $\bar{\mathbf{o}} = P\bar{\mathbf{x}}$.
Picking	$\text{Sign}(\bar{\mathbf{x}}) = \text{Sign}(P^{-1}\mathbf{o})$	$\bar{\mathbf{x}}_b$	The quantization code of \mathbf{o} as a D -bit string.
		\mathbf{q}'	The inversely transformed query vector, i.e., $P^{-1}\mathbf{q}$.
		$\bar{\mathbf{q}}$	The quantized query vector of \mathbf{q}' .
Compute	$\langle \bar{\mathbf{x}}_b, \bar{\mathbf{q}}_u \rangle$	$\bar{\mathbf{q}}_u$	The unsigned integer representation of $\bar{\mathbf{q}}$.

Approximation error of RabbitQ

$$\mathbb{E} \left[\frac{\langle \bar{\mathbf{o}}, \mathbf{q} \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle} \right] = \langle \mathbf{o}, \mathbf{q} \rangle$$

$$\frac{\langle \bar{\mathbf{o}}, \mathbf{q} \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle} = \langle \mathbf{o}, \mathbf{q} \rangle + \sqrt{1 - \langle \mathbf{o}, \mathbf{q} \rangle^2} \cdot \frac{\langle \bar{\mathbf{o}}, \mathbf{e}_1 \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle}$$

$$\left| \frac{\langle \bar{\mathbf{o}}, \mathbf{q} \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle} - \langle \mathbf{o}, \mathbf{q} \rangle \right| = O\left(\frac{1}{\sqrt{D}}\right) \text{ with high probability}$$

$$\mathbb{P} \left\{ \left| \frac{\langle \bar{\mathbf{o}}, \mathbf{q} \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle} - \langle \mathbf{o}, \mathbf{q} \rangle \right| > \sqrt{\frac{1 - \langle \bar{\mathbf{o}}, \mathbf{o} \rangle^2}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle^2}} \cdot \frac{\epsilon_0}{\sqrt{D} - 1} \right\} \leq 2e^{-c_0 \epsilon_0^2}$$

Approximation error of RabbitQ

$$\langle \bar{\mathbf{o}}, \mathbf{q} \rangle = \langle P\bar{\mathbf{x}}, \mathbf{q} \rangle = \langle P^{-1}P\bar{\mathbf{x}}, P^{-1}\mathbf{q} \rangle = \langle \bar{\mathbf{x}}, \mathbf{q}' \rangle$$

THEOREM 3.3. $B_q = \Theta(\log \log D)$ suffices to guarantee that $|\langle \bar{\mathbf{x}}, \mathbf{q}' \rangle - \langle \bar{\mathbf{x}}, \bar{\mathbf{q}} \rangle| = O(1/\sqrt{D})$ with high probability.

Experiments of RabbitQ

Time in the **Indexing Phase**

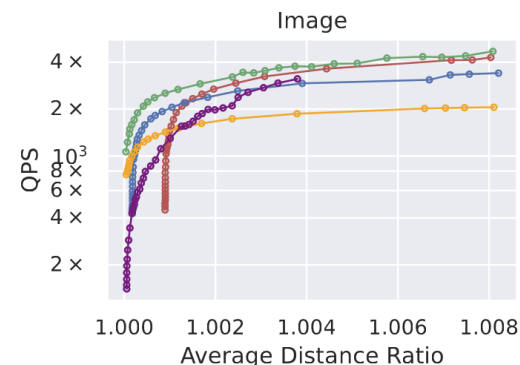
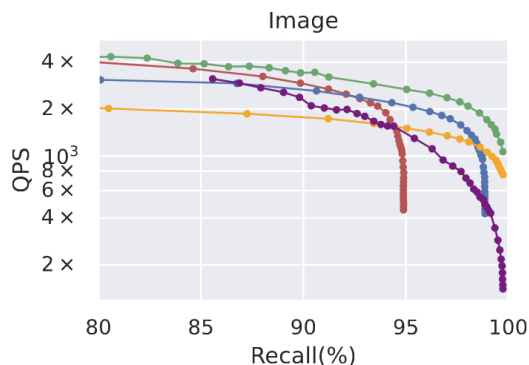
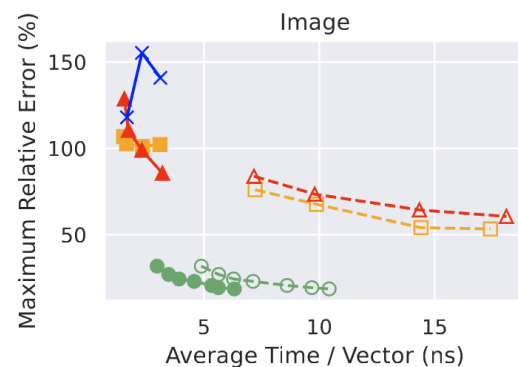
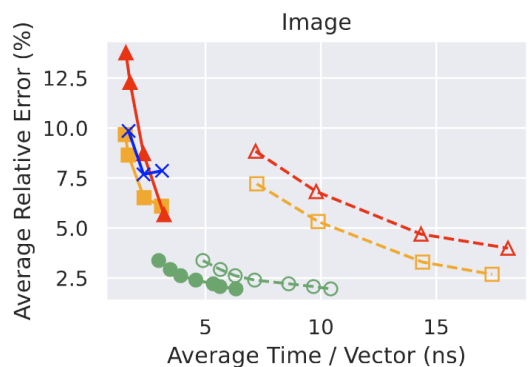
Table 4: The Indexing Time for the GIST Dataset

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

Time-Accuracy Trade-Off per Vector for

Distance Estimation
(first row)

ANN Search (second
row)



Experiments of RabbitQ

“the theoretical analysis provides explicit suggestions on the parameters. Thus, our method requires **no tuning**.”

However...

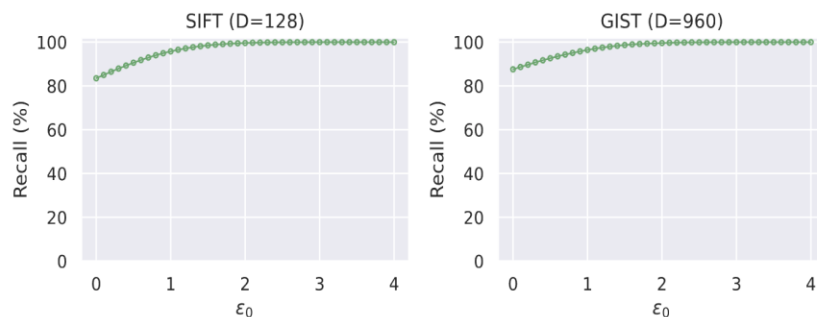


Figure 5: Verification Study on ϵ_0 .

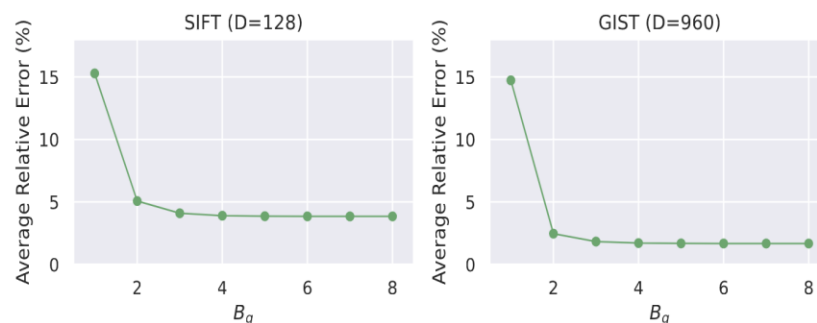


Figure 6: Verification Study on B_q .

Because, they empirically use a **confidence interval** to decide which samples to **re-rank or drop out** in practice.

$$\frac{\langle \bar{\mathbf{o}}, \mathbf{q} \rangle}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle} \pm \sqrt{\frac{1 - \langle \bar{\mathbf{o}}, \mathbf{o} \rangle^2}{\langle \bar{\mathbf{o}}, \mathbf{o} \rangle^2}} \cdot \frac{\epsilon_0}{\sqrt{D-1}}$$

TurboQuant

- Define **different** objectives for cosine similarity and 2-norm, quantizers optimized for MSE do not produce unbiased estimators for inner products in their cases due to its **definitions**.

$$\text{(MSE)} \quad D_{\text{mse}} := \mathbb{E}_Q \left[\|\mathbf{x} - Q^{-1}(Q(\mathbf{x}))\|_2^2 \right]$$

$$\text{(inner-prod error)} \quad D_{\text{prod}} := \mathbb{E}_Q \left[\left| \langle \mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, Q^{-1}(Q(\mathbf{x})) \rangle \right|^2 \right].$$

- Reweighting -> quantize the **residual term**

providing the following unbiased inner product estimator:

$$\langle \mathbf{y}, Q_{\text{mse}}^{-1}(Q_{\text{mse}}(\mathbf{x})) \rangle + \|\mathbf{r}\|_2 \cdot \left\langle \mathbf{y}, Q_{\text{qj1}}^{-1}(Q_{\text{qj1}}(\mathbf{r})) \right\rangle.$$

More formally, the quantization map $Q_{\text{prod}} : \mathbb{S}^{d-1} \rightarrow [2^{b-1}]^d \times \{-1, 1\}^d \times \mathbb{R}$ is defined as:

$$Q_{\text{prod}}(\mathbf{x}) = [Q_{\text{mse}}(\mathbf{x}), Q_{\text{qj1}}(\mathbf{x} - Q_{\text{mse}}^{-1}(Q_{\text{mse}}(\mathbf{x}))), \|\mathbf{x} - Q_{\text{mse}}^{-1}(Q_{\text{mse}}(\mathbf{x}))\|_2].$$

TurboQuant (MSE)

Algorithm 1 TURBOQUANT_{mse}: optimized for MSE

- 1: **input:** dimension d and bit-width b
 // Global Parameters for Setting up TURBOQUANT_{mse}
 - 2: Generate a **random rotation matrix** $\mathbf{\Pi} \in \mathbb{R}^{d \times d}$
 - 3: Construct **codebook** by finding centroids $c_1, c_2, \dots, c_{2^b} \in [-1, 1]$ that minimize MSE cost in Eq. (4)
-
- 4: **Procedure** QUANT_{mse}(\mathbf{x})
 - 5: $\mathbf{y} \leftarrow \mathbf{\Pi} \cdot \mathbf{x}$
 - 6: $\text{idx}_j \leftarrow \arg \min_{k \in [2^b]} |\mathbf{y}_j - c_k|$ for every $j \in [d]$ $\{\text{idx}_j\text{'s are } b\text{-bit integers}\}$
 - 7: **output:** idx
-
- 8: **Procedure** DEQUANT_{mse}(idx)
 - 9: $\tilde{\mathbf{y}}_j \leftarrow c_{\text{idx}_j}$ for every $j \in [d]$
 - 10: $\tilde{\mathbf{x}} \leftarrow \mathbf{\Pi}^\top \cdot \tilde{\mathbf{y}}$
 - 11: **output:** $\tilde{\mathbf{x}}$
-

TurboQuant (Inner Product)

Algorithm 2 $\text{TURBOQUANT}_{\text{prod}}$: optimized for inner product

- 1: **input:** dimension d and bit-width b
 // Global Parameters for Setting up $\text{TURBOQUANT}_{\text{prod}}$
 - 2: Instantiate a $\text{TURBOQUANT}_{\text{mse}}$ with bit-width $b - 1$ as per Algorithm 1
 - 3: Generate a **random projection matrix** $\mathbf{S} \in \mathbb{R}^{d \times d}$ with i.i.d. entries $\mathbf{S}_{i,j} \sim \mathcal{N}(0, 1)$
-
- 4: **Procedure** $\text{QUANT}_{\text{prod}}(\mathbf{x})$
 - 5: $\text{idx} \leftarrow \text{QUANT}_{\text{mse}}(\mathbf{x})$
 - 6: $\mathbf{r} \leftarrow \mathbf{x} - \text{DEQUANT}_{\text{mse}}(\text{idx})$ {residual vector}
 - 7: $\text{qjl} \leftarrow \text{sign}(\mathbf{S} \cdot \mathbf{r})$ {QJL on residual vector}
 - 8: **output:** $(\text{idx}, \text{qjl}, \|\mathbf{r}\|_2)$
-
- 9: **Procedure** $\text{DEQUANT}_{\text{prod}}(\text{idx}, \text{qjl}, \gamma)$
 - 10: $\tilde{\mathbf{x}}_{\text{mse}} \leftarrow \text{DEQUANT}_{\text{mse}}(\text{idx})$
 - 11: $\tilde{\mathbf{x}}_{\text{qjl}} \leftarrow \frac{\sqrt{\pi/2}}{d} \cdot \gamma \cdot \mathbf{S}^\top \cdot \text{qjl}$
 - 12: **output:** $\tilde{\mathbf{x}}_{\text{mse}} + \tilde{\mathbf{x}}_{\text{qjl}}$
-

TurboQuant

- Provide error bounds for different **bit-width budgets** of quantizations

In Theorem 1 we prove that the b -bit **MSE** optimized TURBOQUANT $Q_{\text{mse}} : \mathbb{R}^d \rightarrow \{0, 1\}^{b \cdot d}$ achieves the following distortion for any worst-case vector $\mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\| = 1$:

- $D_{\text{mse}}(Q_{\text{mse}}) := \mathbb{E} \left[\|\mathbf{x} - Q_{\text{mse}}^{-1}(Q_{\text{mse}}(\mathbf{x}))\|_2^2 \right] \leq \frac{\sqrt{3}\pi}{2} \cdot \frac{1}{4^b}$ for any $b \geq 0$.
- For small bit-widths the above distortion upper bound can be further refined. Specifically, for $b = 1, 2, 3, 4$ we have $D_{\text{mse}}(Q_{\text{mse}}) \approx \mathbf{0.36}, \mathbf{0.117}, \mathbf{0.03}, \mathbf{0.009}$, respectively.

In Theorem 2 we prove that the b -bit **inner product** optimized TURBOQUANT $Q_{\text{prod}} : \mathbb{R}^d \rightarrow \{0, 1\}^{b \cdot d}$ achieves the following distortion for any worst-case vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\|\mathbf{x}\| = 1$:

- $\mathbb{E} \left[\left\langle \mathbf{y}, Q_{\text{prod}}^{-1}(Q_{\text{prod}}(\mathbf{x})) \right\rangle \right] = \langle \mathbf{y}, \mathbf{x} \rangle$
- $D_{\text{prod}}(Q_{\text{prod}}) := \mathbb{E} \left[\left| \langle \mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, Q_{\text{prod}}^{-1}(Q_{\text{prod}}(\mathbf{x})) \rangle \right|^2 \right] \leq \frac{\sqrt{3}\pi^2 \cdot \|\mathbf{y}\|_2^2}{d} \cdot \frac{1}{4^b}$ for any $b \geq 0$.
- For small bit-widths the above distortion upper bound can be further refined. Specifically, for $b = 1, 2, 3, 4$ we have $D_{\text{prod}}(Q_{\text{prod}}) \approx \frac{\mathbf{1.57}}{d}, \frac{\mathbf{0.56}}{d}, \frac{\mathbf{0.18}}{d}, \frac{\mathbf{0.047}}{d}$, respectively.

TurboQuant

- Also provide lower bounds

Lower Bound. In Theorem 3, we leverage Shannon's lower bound and Yao's minimax principle to prove that for any randomized quantization algorithm $Q : \mathbb{R}^d \rightarrow \{0, 1\}^{b \cdot d}$ with bit-width b , there exist hard input instances $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\|\mathbf{x}\| = 1$ such that the following lower bounds hold:

- $D_{\text{mse}}(Q) := \mathbb{E} \left[\|\mathbf{x} - Q^{-1}(Q(\mathbf{x}))\|_2^2 \right] \geq \frac{1}{4^b}$
- $D_{\text{prod}}(Q) = \mathbb{E} \left[|\langle \mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, Q^{-1}(Q(\mathbf{x})) \rangle|^2 \right] \geq \frac{\|\mathbf{y}\|_2^2}{d} \cdot \frac{1}{4^b}$

谢谢!

请提宝贵意见



中国科学技术大学
University of Science and Technology of China



