

Abstract

Quantization, which replaces vectors with codewords, is widely used to enable **fast, accurate inner-product-based similarity search over large-scale data**. ScaNN is a popular approach for quantization. ScaNN computes the quantization error between each vector and all possible codewords to select the codeword with the smallest error, achieving **high approximation accuracy**. However, since ScaNN requires error computation with all codewords, it incurs a high computational cost, making quantization extremely slow for large-scale datasets. The proposed approach uses upper bounds on quantization errors and efficiently evaluates them to **prune codeword candidates**. This significantly reduces the number of error computations required. As a result, the proposed approach can **substantially accelerate vector quantization while preserving ScaNN's search accuracy**. Consequently, it facilitates practical large-scale data processing in applications such as image retrieval and natural language processing.

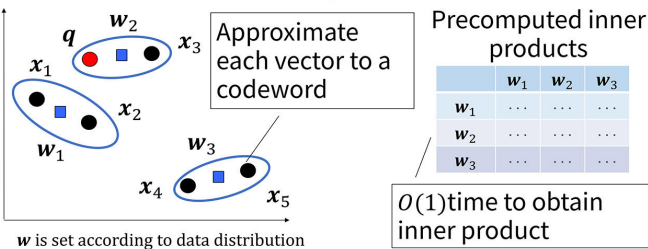
Maximum Inner Product Search

Find vector x of maximum inner product for query q

- Various application (e.g., image search and NLP)
- ✗ High CPU cost for inner products of the query and vectors in the database

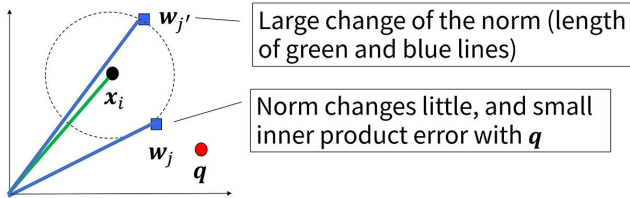
Quantization: approximate vector x with codeword w for search

- Improved efficiency
- ✗ Degraded search accuracy



Existing Approach: ScaNN

Idea: Preserve vector norms and reduce inner product error to improve search accuracy



Approach: Assign to the codeword minimizing error function $l(x_i, w_j)$

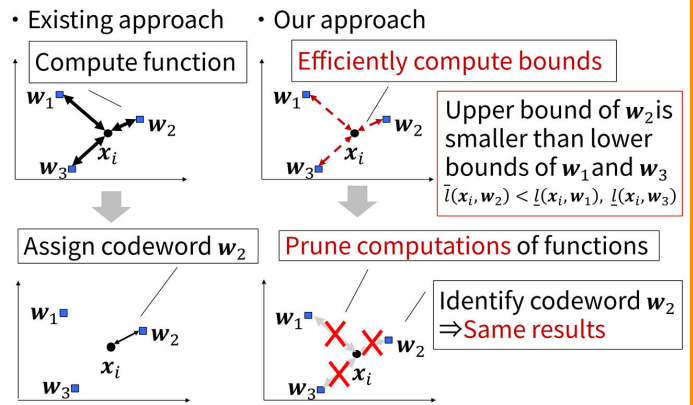
- Error function is low \Leftrightarrow Norm is preserved
- ✗ High CPU cost for error function

$$l(x_i, w_j) = \frac{h_i - 1}{\|x_i\|^2} (p(x_i, w_j))^2 - 2h_i p(x_i, w_j) + h_i \|x_i\|^2 + \|w_j\|^2$$

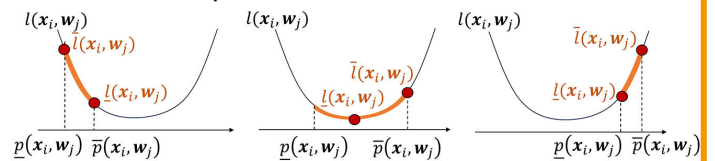
$h_i (> 1)$: hyper parameter, $p(x_i, w_j)$: inner product of x_i and w_j

Our approach: computation pruning

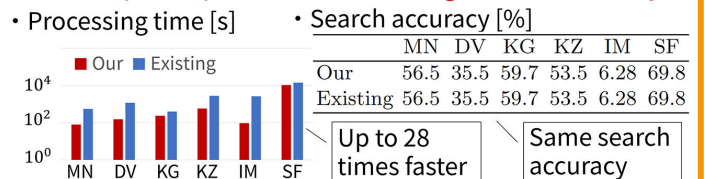
Prune unnecessary computations by efficiently obtaining upper bound $\bar{l}(x_i, w_j)$ and lower bound $\underline{l}(x_i, w_j)$



Function $l(x_i, w_j)$ is quadratic for inner product $p(x_i, w_j) \Rightarrow \bar{l}(x_i, w_j)$ and $\underline{l}(x_i, w_j)$ are obtained from upper/lower bounds of inner product



Achieve speedup without sacrificing search accuracy



Experimental data (Number of vectors \times Length of vectors)

MN : MNIST (70,000 \times 780) KZ : Kuzushiji (270,912 \times 784)
 DV : Devnagari-Script (92,000 \times 1,024) IM : imagenet (2,340,373 \times 150)
 KG : kaggle (200,000 \times 202) SF : sift10m (10,000,000 \times 128)

References

[1] Y. Fujiwara, Á. López, Y. Ida, A. Kumagai, M. Nakano, M. Nakatsuji, A. Kimura, "Fast Vector Quantization Algorithm for ScaNN", in Proc KDD, 2026.

Contact

Yasuhiro Fujiwara, Recognition Research Group, Media Information Laboratory