

GOAL

Approximate nearest neighbor (ANN) search

I. CONTRIBUTIONS

- A novel quantization approach based on constrained sparse coding for ANN search.
- The proposed approach is demonstrated with an adaptation of two vector quantizers:
 - Product Quantizer (PQ) as "Q α -PQ"
 - Residual Vector Quantizer (RVQ) as "Q α -RVQ"
- Competitive results on four datasets including billion-sized BIGANN dataset.
 - Significant improvement over RVQ.
 - Faster encoding with similar or better accuracy.

II. BACKGROUND

- *ANN search* attempts to find a query's nearest neighbor efficiently by relying on compressed representations.
- A *vector quantizer (VQ)* maps a vector to a codeword,

$$\mathbf{x} \mapsto \mathbf{c}_{k(\mathbf{x})} \in \mathbf{C} \quad (1)$$

where \mathbf{C} is a codebook containing K codewords. The codeword id of vector \mathbf{x} , $k(\mathbf{x}) \in [1, K]$, is often given by:

$$k(\mathbf{x}) = \arg \min_{i \in [1, K]} \|\mathbf{x} - \mathbf{c}_i\|^2 \quad (2)$$

for Euclidean based encoding, or for cosine similarity based encoding as:

$$k(\mathbf{x}) = \arg \max_{i \in [1, K]} \mathbf{x}^\top \mathbf{c}_i \quad (3)$$

- *Product quantizer (PQ)* [2]

There are M sub-codebooks \mathbf{C}^m , each represents a different orthogonal sub-space.

$$Q(\mathbf{x}) = \sum_{m=1}^M \mathbf{c}_{k_m(\mathbf{x})}^m \quad (4)$$

where $\mathbf{c}_i^m \in \mathbf{C}^m$.

An encoded vector takes $M \times \log K$ bits.

- *Residual vector quantizer (RVQ)* [1]

The sub-codebooks form a hierarchy and each represents a different level of residuals.

$$Q(\mathbf{x}) = \sum_{m=1}^M \mathbf{c}_{k_m(\mathbf{r}_m)}^m \quad (5)$$

where $\mathbf{r}_m = \mathbf{x} - \sum_{j=1}^{m-1} \mathbf{c}_{k_j(\mathbf{r}_j)}^j$, with $\mathbf{r}_1 = \mathbf{x}$.

An encoded vector takes $M \times \log K$ bits.

- *Sparse coding view of VQ*. In sparse coding a vector is represented as $\mathbf{C}\alpha$ which is equivalent to

$$Q(\mathbf{x}) = \sum_{i=1}^K \alpha_i \mathbf{c}_i \quad (6)$$

where α is a sparse vector.

PQ and RVQ can be seen as sparse encoding with binary weights and sparsity M (out of MK).

REFERENCES

- [1] Approximate nearest neighbor search by residual vector quantization. Chen, Yongjian and Guan, Tao and Wang, Cheng. *Sensors*, 2010.
- [2] Product quantization for nearest neighbor search. Jégou, Hervé and Douze, Matthijs and Schmid, Cordelia. *IEEE TPAMI*, 2011.
- [3] The inverted multi-index. Babenko, Artem and Lempitsky, Victor. *CVPR*, 2012.
- [4] All about VLAD. Arandjelovic, Relja and Zisserman, Andrew. *CVPR*, 2013.
- [5] Cartesian k-means. Norouzi, Mohammad and Fleet, David J. *CVPR*, 2013.
- [6] Additive quantization for extreme vector compression. Babenko, Artem and Lempitsky, Victor. *CVPR*, 2014.
- [7] Composite Quantization for Approximate Nearest Neighbor Search. Zhang, Ting, *ICML*, 2014.
- [8] Optimized residual vector quantization for efficient approximate nearest neighbor search. Ai, Liefu and Yu, Junqing and Wu, Zebin and He, Yunfeng and Guan, Tao *Multimedia Systems*, 2015.

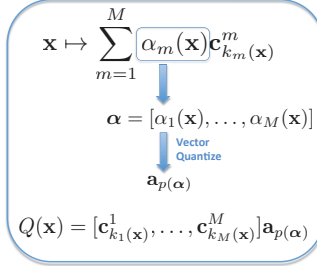
III. APPROACH

- *Proposed representation*

$$Q(\mathbf{x}) = \sum_{m=1}^M \alpha_m(\mathbf{x}) \mathbf{c}_{k_m(\mathbf{x})}^m \quad (7)$$

$\alpha_m(\mathbf{x}) \in \mathbb{R}$ and the codewords are unit norm vectors.

- Coefficient vectors formed as $\alpha = [\alpha_1(\mathbf{x}), \dots, \alpha_M(\mathbf{x})]$ are vector quantized.



where $\mathbf{a}_p(\alpha) \in \mathbf{A}$, the coefficient codebook, with P codewords and $P(\alpha) \in [1, P]$

- *Q α -PQ*

$$k_m(\mathbf{x}) = \arg \max_{i \in [1, K]} \mathbf{x}^\top \mathbf{c}_i^m$$

$$\alpha_m(\mathbf{x}) = \mathbf{x}^\top \mathbf{c}_{k_m(\mathbf{x})}^m$$

$$\alpha \mapsto \mathbf{a}_p(\alpha)$$

\mathbf{x} is encoded as $([k_1(\mathbf{x}), \dots, k_M(\mathbf{x})], \mathbf{a}_p(\alpha))$

- *Q α -RVQ*

$$k_m(\mathbf{r}_m) = \arg \max_{i \in [1, K]} \mathbf{r}_m^\top \mathbf{c}_i^m$$

$$\alpha_m(\mathbf{r}_m) = \mathbf{r}_m^\top \mathbf{c}_{k_m(\mathbf{r}_m)}^m$$

$$\alpha = [\alpha_{k_1}^1(\mathbf{x}), \dots, \alpha_{k_M}^M(\mathbf{x})]^\top \mathbf{x}$$

$$\alpha \mapsto \mathbf{a}_p(\alpha)$$

\mathbf{x} is encoded as $([k_1(\mathbf{r}_1), \dots, k_M(\mathbf{r}_M)], \mathbf{a}_p(\alpha))$ where, $\mathbf{r}_m = \mathbf{x} - \sum_{j=1}^{m-1} \alpha_j \mathbf{r}_j$ and $\mathbf{r}_1 = \mathbf{x}$

- Learning \mathbf{C}^m : Spherical K-means.
- Learning \mathbf{A} : K-means.
- An encoded vector takes $M \times \log K + \log P$ bits.

IV. RESULTS: ANN SEARCH

- *Impact of 1-byte quantization of α* .

- Accuracy with quantized α compared to unquantized α and the associated VQ.
- Q α -RVQ retains the accuracy of α -RVQ.
- Q α -PQ performs well for lower M .

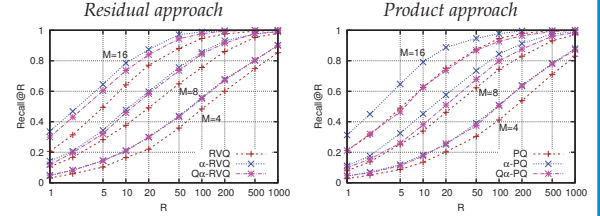


Figure 1: Recall rates on VLAD500K [4] (128 dimensional).

- *Comparison with state-of-the-art*.

- Accuracy of various quantization methods using the same compression rate.
- Q α -RVQ performs the best after additive quantization(AQ) [6], a highly accurate but not scalable approach.

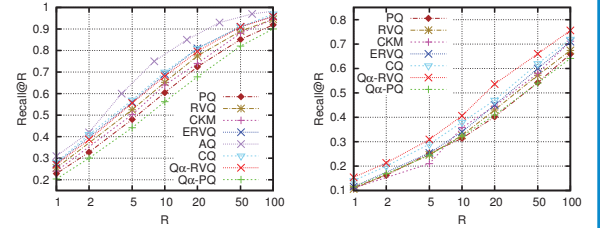


Figure 2: Recall rates on SIFT1M [2] & GIST1M [2] (64-bit encoding) compared with CKM [5], ERVQ [8], AQ [6] and CQ [7].

- *Efficiency*

- Comparing complexity of codebook learning, database encoding and ANN search.

	PQ	Q α -PQ	RVQ	Q α -RVQ
learn	1	0.212	1.250	0.719
encode	1	0.206	1.347	0.613
search	1	1.867	1.220	1.909

In the above experiments, we set $P = 256$ and $K = 256$, except for Q α -PQ in Fig. 2 and Tab. 1 where $K=128$.

Table 1: Complexity w.r.t. PQ on SIFT1M (64-bit encoding).

V. RESULTS: NON-EXHAUSTIVE ANN SEARCH

- *Results on Billion-sized BIGANN [Jégou et al., ICASSP 2011] dataset using two inverted indexing methods.*

- *Inverted index* [2] uses a coarse quantizer and then the residual is quantized using the chosen encoder.

- Q α -RVQ is the most accurate method.
- The hierarchical structure of Q α -RVQ permits to further narrow the exhaustive search leading to much faster search.

Method (bytes)	R@1	R@10	R@100	time
PQ-64 (8)	0.111	0.388	0.756	1.00
PQ-72 (9)	0.144	0.471	0.825	1.03
RVQ (9)	0.124	0.421	0.803	1.02
Q α -PQ (9)	0.139	0.450	0.811	1.69
Q α -RVQ (10)	0.160	0.514	0.868	1.72
Q α -RVQ ₁₂₈	0.160	0.514	0.868	0.89
Q α -RVQ ₈	0.151	0.467	0.730	0.17

Table 2: Performance on BIGANN with relative timings.

- *Inverted multi-index* [3] uses two-fold product quantizers with a large number of inverted lists.

- Q α -RVQ and Q α -PQ outperform others in accuracy.
- Interestingly Q α -RVQ can be more accurate and faster than PQ by scanning fewer candidates.

Method (bytes)	$T = 100K$				$T = 30K$			
	R@1	R@10	R@100	time	R@1	R@10	R@100	time
PQ-64 (8)	0.170	0.535	0.869	29	0.170	0.526	0.823	11
Q α -PQ (9)	0.200	0.587	0.898	30	0.198	0.572	0.848	11
RVQ (9)	0.181	0.553	0.877	37	0.180	0.542	0.831	14
Q α -RVQ (10)	0.227	0.630	0.920	37	0.225	0.613	0.862	14

Table 3: Performance with T candidates exhaustively scanned. Time is in milliseconds.