

Vector and Line Quantization for Billion-scale Similarity Search on GPUs

Wei Chen^a, Jincan Chen^{a,b,*}, Fuhao Zou^{c,*}, Yuan-Fang Li^d, Ping Lu^{a,b},
Qiang Wang^a, Wei Zhao^b

^a*Wuhan National Laboratory for Optoelectronics, Huazhong University of Science and Technology, Wuhan 430074, China*

^b*Key Laboratory of Information Storage System of Ministry of Education, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan 430074, China*

^c*School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan 430074, China*

^d*Faculty of Information Technology, Monash University, Clayton 3800, Australia*

Abstract

Billion-scale high-dimensional approximate nearest neighbour (ANN) search has become an important problem for searching similar objects among the vast amount of images and videos available online. The existing ANN methods are usually characterized by their specific indexing structures, including the inverted index and the inverted multi-index structure. The inverted index structure is amenable to GPU-based implementations, and the state-of-the-art systems such as Faiss are able to exploit the massive parallelism offered by GPUs. However, the inverted index requires high memory overhead to index the dataset effectively. The inverted multi-index structure is difficult to implement for GPUs, and also ineffective in dealing with database with

*Corresponding author

Email addresses: jcchen@hust.edu.cn (Jincan Chen), fuhao_zou@hust.edu.cn (Fuhao Zou)

different data distributions. In this paper we propose a novel hierarchical inverted index structure generated by vector and line quantization methods. Our quantization method improves both search efficiency and accuracy, while maintaining comparable memory consumption. This is achieved by reducing search space and increasing the number of indexed regions.

We introduce a new ANN search system, VLQ-ADC, that is based on the proposed inverted index, and perform extensive evaluation on two public billion-scale benchmark datasets SIFT1B and DEEP1B. Our evaluation shows that VLQ-ADC significantly outperforms the state-of-the-art GPU- and CPU-based systems in terms of both accuracy and search speed. The source code of VLQ-ADC is available at <https://github.com/zjuchenwei/vector-line-quantization>.

Keywords: Quantization; Billion-scale similarity search; high dimensional data; Inverted index; GPU

1. Introduction

In the age of the Internet, the amount of images and videos available online increases incredibly fast and has grown to an unprecedented scale. Google processes over 40,000 various queries per second, and handles more than 400 hours of YouTube video uploads every minute [1]. Every day, more than 100 million photos/videos are uploaded to Instagram, more than 300 million uploaded to Facebook, and a total of 50 billion photos have been shared to Instagram¹. As a result, scalable and efficient search for similar

¹<https://www.omnicoreagency.com/instagram-statistics/>

9 images and videos on the billion scale has become an important problem and
10 it has been under intense investigation.

11 As online images and videos are unstructured and usually unlabeled, it
12 is hard to compare them directly. A feasible solution is to use real-valued,
13 high-dimensional vectors to represent images and videos, and compare the
14 distances between the vectors to find the nearest ones. Due to the curse of
15 dimensionality [2], it is impractical for multimedia applications to perform
16 exhaustive search in billion-scale datasets. Thus, as an alternative, many
17 *approximate nearest neighbor* (ANN) search algorithms are now employed
18 to tackle the billion-scale search problem for high-dimensional data. Recent
19 best-performing billion-scale retrieval systems [3–8] typically utilize two main
20 processes: indexing and encoding.

21 To avoid expensive exhaustive search, these systems use *index structures*
22 that can partition the dataset space into a large number of disjoint regions,
23 and the search process only collects points from the regions that are closest to
24 the query point. The collected points then form a short list of candidates for
25 each query point. The retrieval system then calculates the distance between
26 each candidate and the query point, and sort them accordingly.

27 To guarantee query speed, the indexed points need to be loaded into
28 RAM. For large datasets that do not fit in RAM, dataset points are *encoded*
29 into a compressed representation. Encoding has also proven to be critical for
30 memory-limited devices such as GPUs that excel at handling data-parallel
31 tasks. A high-performance CPU like Intel Xeon Platinum 8180 (2.5 GHz,
32 28 cores) performs 1.12 TFLOP/s single precision peak performance². In

²<https://ark.intel.com/content/www/us/en/ark/products/120496/>

33 contrast, GPUs like NVdia Tesla P100 can provide up to 10T FLOP/s sin-
34 gle precision peak performance³, and are good choices for high performance
35 similarity search systems. Many encoding methods have been proposed, in-
36 cluding hashing methods and quantization methods. Hashing methods en-
37 code data points to compact binary codes through a hash function [9, 10],
38 and quantization methods, typically product quantization (PQ), map data
39 points to a set of centroids and use the indices of the centroids to encode the
40 data points [11, 12]. By hashing methods, the distance between two data
41 points can be approximated by the Hamming distance between their binary
42 code. By quantization methods, the Euclidean distance between the query
43 and compressed points can be computed efficiently. It has been shown in
44 the literature that quantization encoding can be more accurate than various
45 hashing methods [11, 13, 14].

46 Jégou et al. [11] first introduced an index structure that is able to handle
47 billion-scale datasets effieciently. It is based on the inverted index structure
48 that partitions the high dimensional vector space into Voronoi regions for a
49 set of centroids obtained by a quantization method called vector quantization
50 (VQ) [15]. This system, called IVFADC, achieves reasonable recall rates in
51 several tens of milliseconds. However, the VQ-based index structure needs
52 to store a large set of full dimensional centroids to produce a huge number
53 of regions, which would require a large amount of memory.

54 An improved inverted index structure called the inverted multi-index

intel-xeon-platinum-8180-processor-38-5m-cache-2-50-ghz.html

³<https://images.nvidia.com/content/tesla/pdf/nvidia-tesla-p100-PCIe-datasheet.pdf>

55 (IMI) was later proposed by Babenko and Lempitsky [16]. The IMI is based
56 on product quantization (PQ), which divides the point space into several
57 orthogonal subspaces and clusters the subspaces into Voronoi regions inde-
58 pendently. The Cartesian product of regions in each subspace forms regions
59 in the global point space. The strength of the IMI is that it can produce a
60 huge number of regions with much smaller codebooks than that of the in-
61 verted index. Due to the huge number of indexed regions, the point space
62 is finely partitioned and each regions contains fewer points. Hence the IMI
63 can provide accurate and concise candidate lists with memory and runtime
64 efficiency.

65 However, it has been observed that for some billion-scale datasets, the
66 majority of the IMI regions contain no points [5], which is a waste of index
67 space and has a negative impact on the final retrieval performance. The rea-
68 sons for this deficiency is that the IMI learns the centroids independently on
69 the subspaces which are not statistically independent [7]. In fact, some con-
70 volutional neural networks (CNN) produce feature vectors with considerable
71 correlations between the subspaces [10, 17, 18].

72 The high level of parallelism provided by GPUs has recently been lever-
73 aged to accelerate similarity search of high-dimensional data, and it has been
74 demonstrated that GPU-based systems are more efficient than CPU-based
75 systems by a large margin [4, 6]. Comparing to IMI structure, the inverted in-
76 dexing structure proposed by Jégou et al. [11] is more straightforward to par-
77 allelize, because the IMI structure depends on a complicated multi-sequence
78 algorithm, which is sequential in nature [4] and hard to parallelize.

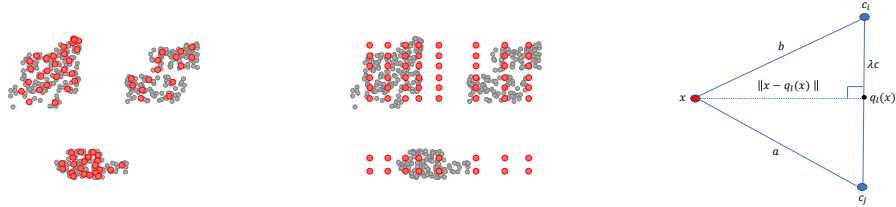
79 To the best of our knowledge, there are two high performance systems

80 that are able to handle ANN search for billion-scale datasets on the GPU:
81 PQT [4] and Faiss [6]. PQT proposes a novel quantization method call line
82 quantization (LQ) and is the first billion-scale similarity retrieval system on
83 the GPU. Subsequently Faiss implements the idea of IVFADC on GPUs and
84 currently has the state-of-the-art performance on GPUs. We compare our
85 method against Faiss and two other systems in Section 5.

86 In this paper, we present VLQ-ADC, a novel billion-scale ANN similarity
87 search framework. VLQ-ADC includes a two-level hierarchical inverted in-
88 dexing structure based on Vector and Line Quantization (VLQ), which can
89 be implemented on GPU efficiently. The main contributions of our solution
90 are threefold.

- 91 1. We demonstrate how to increase the number of regions with memory
92 efficiency by the novel inverted index structure. The efficient indexing
93 contributes to high accuracy for approximate search.
- 94 2. We describe how to encode data points via a novel algorithm for high
95 runtime/memory efficiency.
- 96 3. Our evaluation shows that our system consistently and significantly
97 outperforms state-of-the-art GPU- and CPU-based retrieval systems
98 on both recall and efficiency on two public billion-scale benchmark
99 datasets with single- and multi-GPU configurations.

100 The rest of the paper is organized as follows. Section 2 introduces related
101 works on indexing with quantization methods. Section 3 presents VLQ, our
102 approach for approximate nearest neighbor (ANN)-based similarity search
103 method. Section 4 introduces the details of GPU implementation. Section 5
104 provides a series of experiments, and compares the results to the state of the



(a) Vector Quantization. (b) Product Quantization. (c) Line Quantization.

Figure 1: Three different quantization methods. Vector and Product quantization methods are both with $k = 64$ clusters. The red dots in plot (a) and (b) denote the centroids and the grey dots denote the dataset points in both plots. Vector quantization (a) maps the dataset points to the closest centroids. Product quantization (b) performs clustering in each subspace independently (here axes). In plot (c), a 2-dimensional point x (red dot) is projected on line $l(c_i, c_j)$ with the anchor point $q_l(x)$ (black dot). The a, b, c denote the values of $\|x - c_i\|^2, \|x - c_j\|^2$ and $\|c_i - c_j\|^2$ respectively. We use the parameter λ to represent the value of $\|c_i - q_l(x)\| / c$. The anchor point $q_l(x)$ can be represented by c_i, c_j and λ . The distance from x to $l(c_i, c_j)$ can be calculated by a, b, c and λ .

106 2. Related work

107 In this section, we briefly introduce some quantization methods and sev-
 108 eral retrieval systems related to our approach. Table 1 summarizes the com-
 109 mon notations used throughout this paper. For example, we assume that
 110 $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^D$ is a finite set of N data points of dimension D .

111 2.1. Vector quantization (VQ)

112 In vector quantization [15] (Figure 1 a), a quantizer is a function q_v that
 113 maps a D -dimensional vector x to a vector $q_v(x) \in C$, where C is a finite
 114 subset of \mathbb{R}^D , of k vectors. Each vector $c \in C$ is called a centroid, and C is

Table 1: Commonly used notations.

Notation	Description
x_i, D	data points, their dimension and the number of data points
\mathcal{X}, N	a set of data points and its size, $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^D$
$c, s, l(c, s)$	centroids, nodes and edges
m	encoding length
k	the number of first-level centroids
n	the number of edges of each first-level centroid
w_1	the number of first-layer nearest regions for a query
α	the portion of the nearest of the $w \cdot n$ second-level regions
w_2	the number of second-level nearest regions for a query, $w_2 = \alpha \cdot w_1 \cdot n$
λ	a scalar parameter for line quantization
r	displacement from data points to the approximate points

115 a codebook of size k . We can use Lloyd iterations [19] to efficiently obtain a
 116 codebook C on a subset of the dataset. For a finite dataset, \mathcal{X} , $q_v(x)$ induces
 117 quantization error E :

$$E = \sum_{x \in \mathcal{X}} \|x - q_v(x)\|^2. \quad (1)$$

118 According to Lloyd's first condition, to minimize quantization error a
 119 quantizer should map vector x to its nearest codebook centroid.

$$q_v(x) = \arg \min_{c \in C} \|x - c\|. \quad (2)$$

120 Hence, the set of points $\mathcal{X}_i = \{x \in \mathbb{R}^D \mid q_v(x) = c_i\}$ is called a cluster or
 121 a region for centroid c_i .

122 The **inverted index structure** based on VQ [11] can split the dataset
 123 space into k regions that correspond to the k centroids of the codebook.
 124 Since the ratio of regions to centroids is 1:1, it requires a large amount of
 125 space to store the D -dimensional centroids when k is large. This would give
 126 a negative effect on the performance of the retrieval system. Our hierarchical
 127 index structure based on VLQ increase the ratio by n times, i.e., n times more
 128 regions can be generated by our indexing structure with the same number of
 129 centroids as the VQ based indexing structure.

130 2.2. Product quantization (PQ)

131 Product quantization (Figure 1 (b)) is an extension of vector quantization.
 132 Assuming that the dimension D is a multiple of m , any vector $x \in \mathbb{R}^D$ can be
 133 regarded as a concatenation (x^1, \dots, x^m) of m sub-vectors, each of dimension
 134 D/m . Suppose that C^1, \dots, C^m are m codebooks of subspace $\mathbb{R}^{D/m}$, each
 135 owns k D/m -dimensional sub-centroids. A codebook of a product quantizer
 136 q_p is thus a Cartesian product of sub-codebooks.

$$C = C^1 \times \dots \times C^m. \quad (3)$$

Hence the codebook C contains a total of k^m centroids, each is a form of
 $c = (c^1, \dots, c^m)$, where each sub-centroid $c^i \in C^i$ for $i \in \mathcal{M} = \{1, \dots, m\}$.
 A product quantizer q_p should minimize the quantization error E defined in
 Formula 1. Hence, for $x \in \mathbb{R}^D$, the nearest centroid in codebook C is

$$q_p(x) = (q_p^1(x^1), \dots, q_p^m(x^m)), \quad (4)$$

137 where q^i is a sub-quantizer of q and $q_p^i(x)$ is the nearest sub-centroid for
 138 sub-vector x^i , i.e., the nearest centroid $q_p(x)$ for x is the concatenation of the
 139 nearest sub-centroids for sub-vector x^i .

140 The **inverted multi-index structure** (IMI) applies the idea of PQ for
 141 indexing and can generate k^m regions with m codebooks of k sub-centroids
 142 each. The benefit of inverted multi-index is thus it can easily generate a much
 143 larger number of regions than that of VQ-based inverted index structure with
 144 moderate values of m and k . The drawback of IMI is that it produces a lot of
 145 empty regions when the distributions of subspaces are not independent [5].
 146 This will affect the system’s performance when handling datasets which have
 147 significant correlations between different subspaces, such as CNN-produced
 148 feature point dataset [5].

149 The PQ-based indexing structure later has been improved by OPQ [20]
 150 and LOPQ [12]. OPQ make a rotation on dataset points by a global $D \times D$
 151 rotation matrix and LOPQ rotates the points which belong to the same cell
 152 by a same local $D \times D$ rotation matrix to minimize correlations between two
 153 subspaces [20]. OPQ and LOPQ can both improve the indexing efficiency of
 154 PQ but slow down the query speed by a large margin as well.

Additionally, PQ can also be used to compress datasets. Typically each
 sub-codebook of PQ contains 256 sub-centroids and each vector x is mapped
 to a concatenation of m sub-centroids $(c_{j_1}^1, \dots, c_{j_m}^m)$, for j_i is a value between
 1 and 256. Hence the vector x can be encoded into an m -byte code of
 sub-centroid index (j_1, \dots, j_m) . With the approximate representation by
 PQ, the Euclidean distances between the query vector and the large number
 of compressed vectors can be computed efficiently. According to the ADC

procedure [11], the computation is performed based on lookup tables.

$$\|y - x\|^2 \approx \|y - q_p(x)\|^2 = \sum_{i=1}^m \|y^i - c_{j_i}^i\|^2 \quad (5)$$

155 where y^i is the i th subvector of a query y . The Euclidean distances between
 156 the query sub-vector y^i and each sub-centroids $c_{j_i}^i$ can be precomputed and
 157 stored in lookup tables that reduce the complexity of distance computation
 158 from $\mathcal{O}(D)$ to $\mathcal{O}(m)$. Due to the high compression quality and efficient dis-
 159 tance computation approach, PQ is considered the top choice for compact
 160 representation of large-scale datasets[3, 7, 12, 14, 20].

161 2.3. Line quantization (LQ)

162 Line quantization (LQ) [4] owns a codebook C of k centroids like VQ. As
 163 shown in Figure 1 (c), with any two different centroids $c_i, c_j \in C$, a line is
 164 formed and denoted by $l(c_i, c_j)$. A line quantizer q_l quantizes a point x to
 165 the nearest line as follows:

$$q_l(x) = \arg \min_{l(c_i, c_j)} d(x, l(c_i, c_j)), \quad (6)$$

where $d(x, l(c_i, c_j))$ is the Euclidean distance from x to the line $l(c_i, c_j)$, and
 the set $\mathcal{X}_{i,j} = \{x \in \mathbb{R}^D | q_l(x) = l(c_i, c_j)\}$ is called a cluster or a region for line
 $l(c_i, c_j)$. The squared distance $d(x, l(c_i, c_j))$ can be calculated as following :

$$\begin{aligned} d(x, l(c_i, c_j))^2 &= (1 - \lambda) \|x - c_i\|^2 + (\lambda^2 - \lambda) \|c_j - c_i\|^2 \\ &\quad + \lambda \|x - c_j\|^2 \end{aligned} \quad (7)$$

Because the values of $\|x - c_j\|^2$, $\|x - c_i\|^2$, $\|c_j - c_i\|^2$ can be pre-computed
 between x and all centroids, Equation 7 can be calculated efficiently. The

anchor point of x is represented by $(1 - \lambda) \cdot c_i + \lambda \cdot c_j$, where λ is a scalar parameter that can be computed as following:

$$\lambda = 0.5 \cdot \frac{(\|x - c_i\|^2 + \|c_j - c_i\|^2 - \|x - c_j\|^2)}{\|c_j - c_i\|^2}. \quad (8)$$

When x is quantized to a region of $l(c_i, c_j)$, then the displacement of x from $l(c_i, c_j)$ can be computed as following:

$$r_{ql}(x) = x - ((1 - \lambda) \cdot c_i + \lambda \cdot c_j). \quad (9)$$

166 Here we regard $l(c_i, c_j)$ and $l(c_j, c_i)$ as two different lines. So LQ-based
 167 indexing structure can produce $k \cdot (k - 1)$ regions with a codebook of k
 168 centroids, The benefit of LQ-based indexing structure is that it can produce
 169 many more regions than that of VQ-based regions. However it is considerably
 170 more complicated to find the nearest line for a point x when k is large. So
 171 we use LQ as an indexing approach with a codebook of a few lines.

Table 2: A summary of current state-of-the-art retrieval systems based on quantization method. N is the size of the dataset \mathcal{X} , m is the number of sub-vectors in product quantization (PQ), k is the size of the codebook, and n is the number of second-level regions. In the last column of each row, the first term is the complexity for encoding, and the second term is the complexity for indexing.

System	Index structure	Encoding	CPU/GPU	Space complexity
Faiss [6]	VQ	PQ	GPU	$\mathcal{O}(N \cdot m) + \mathcal{O}(k \cdot D)$
Ivf-hnsw [7]	2-level VQ	PQ	CPU	$\mathcal{O}(N \cdot m) + \mathcal{O}(k \cdot (D + n))$
Multi-D-ADC [16]	IMI (PQ)	PQ	CPU	$\mathcal{O}(N \cdot m) + \mathcal{O}(k \cdot (D + k))$
VLQ-ADC (our system)	VLQ	PQ	GPU	$\mathcal{O}(N \cdot m) + \mathcal{O}(k \cdot (D + n))$

172 *2.4. The applications of VQ-based and PQ-based indexing structures for billion-*
173 *scale dataset*

174 In this subsection we introduce several billion-scale similarity retrieval
175 systems that apply VQ- or PQ-based indexing structure and encoded by
176 PQ, and discuss their strengths and weaknesses.

177 All the systems discussed below are best-performing, state-of-the-art sys-
178 tems for billion-scale high-dimensional ANN search. Their indexing structure
179 and encoding method are summarized in Table 2. Since all these systems em-
180 ploy the same encoding method based on PQ, we will mainly focus on their
181 indexing structures in the discussions below.

Faiss [6] is a very efficient GPU-based retrieval approach, by realizing the idea of IVFADC [11] on GPUs. Faiss uses the inverted index based on VQ [21] for non-exhaustive search and compresses the dataset by PQ. The inverted index of IVFADC owns a vector quantizer q with a codebook of k centroids. Thus there are k regions for the data space. Each point $x \in \mathcal{X}$ is quantized to a region corresponding to a centroid by a VQ quantizer q_v . The displacement of each point from the centroid of a region it belongs to is defined as

$$r_q(x) = x - q(x), \quad (10)$$

182 where the displacement $r_q(x)$ is encoded by PQ with m codebooks shared
183 by all regions. For each region, an inverted list of data points is maintained,
184 along with PQ-encoded displacements.

185 The search process of Faiss/IVFADC proceeds as follows:

- 186 1. A query point y is quantized to its w nearest regions, extracting a list
187 of candidates $\mathcal{L}_c \subset \mathcal{X}$ which have a high probability of containing the

- 188 nearest neighbor.
- 189 2. The displacement of the query point y from the centroid of each sub-
190 region is computed as $r_q(y)$.
 - 191 3. The distances between $r_q(y)$ and PQ-encoded displacements in \mathcal{L}_c are
192 then computed according to Formula 5.
 - 193 4. Sort the list \mathcal{L}_c to be \mathcal{L}_s based on the distances computed above. The
194 first points in \mathcal{L}_s are returned as the search result for query point y .

195 **Ivf-hnsw** [7] is a retrieval system based on a two-level inverted index
196 structure. Ivf-hnsw first splits the data space into k regions like IVFADC.
197 Then each region is further split into several sub-regions that correspond to
198 n sub-centroids. Each sub-centroid of a region can be represented by the
199 centroid of the region and another centroid of a neighbor region. Assume
200 that each region has n neighbor regions, thus each region can be split into
201 n regions. Each data point is first quantized to a region and then further
202 quantized to a sub-region of the region. The displacement of each point from
203 the sub-centroid of a sub-region it belongs to is encoded by PQ. An inverted
204 list of data point is maintained for each sub-regions similar to IVFADC.

205 The search process of Ivf-hnsw proceeds as follows:

- 206 1. A query point y is quantized to its w first-level nearest regions, giving
207 $w \cdot n$ sub-regions.
- 208 2. Among the $w \cdot n$ sub-regions, y is secondly quantized to $0.5 \cdot w \cdot n$ nearest
209 sub-regions, generating a list of candidates $\mathcal{L}_c \subset \mathcal{X}$.
- 210 3. The displacement of the query point y from the sub-centroid of each
211 sub-region is computed as $r_q(y)$.

- 212 4. The distances between $r_q(y)$ and PQ-encoded displacements in \mathcal{L}_c are
 213 then computed according to Formula 5.
- 214 5. The re-ordering process of Ivf-hnsw is similar to IVFADC/Faiss.

215 **Multi-D-ADC** [16] is based on the inverted multi-index which is cur-
 216 rently the state-of-the-art indexing method for high-dimensional large-scale
 217 datasets. An inverted multi-index of Multi-D-ADC usually owns a product
 218 quantizer with two sub-quantizers q^1, q^2 for subspace $\mathbb{R}^{D/2}$, each of k sub-
 219 centroids. A region in the D-dimensional space is now a Cartesian product
 220 of two corresponding subspace regions. So the IMI can produce k^2 regions.
 221 For each point $x = (x^1, x^2) \in \mathcal{X}$, sub-vectors $x^1, x^2 \in \mathbb{R}^{D/2}$ are separately
 222 quantized to subspace regions of $q^1(x^1), q^2(x^2)$ respectively, and x is then
 223 quantized to the region of $(q^1(x^1), q^2(x^2))$. The displacement of each point
 224 x from the centroid $(q^1(x^1), q^2(x^2))$ is also encoded by PQ, and an inverted
 225 list of points is again maintained for each region.

226 The search process of Multi-D-ADC proceeds as follows:

- 227 1. For a query point $y = (y^1, y^2)$, The Euclidean distances of each of sub-
 228 vectors y^1, y^2 to all sub-centroids of q^1, q^2 are computed respectively.
 229 The distance of y to a region can be computed according to Formula 5
 230 for $m = 2$.
- 231 2. Regions are traversed in ascending order of distance to y by the multi-
 232 sequence algorithm [16] to generate a list of candidates $\mathcal{L}_c \subset \mathcal{X}$.
- 233 3. The displacement of the query point y from the centroid (c^1, c^2) of each
 234 region is computed as $r_q(y)$ as well.
- 235 4. The re-ordering process of Multi-D-ADC is similar to IVFADC/Faiss.

236 The VQ-based indexing structure requires a large full-dimensional code-
237 book to produce regions when k is large. The PQ-based indexing structure
238 are not suitable for all datasets, especially for those produced by convolu-
239 tional neural networks (CNN) [7]. The novel VQ-based indexing structure
240 proposed by Ivf-hnsw can produce more regions than the prior VQ-based
241 indexing structure. However its performance on the codebook of small size
242 is not good enough. We will discuss that in Sec.5. In comparison, our index-
243 ing structure is efficient with a small size of codebook which can accelerate
244 query speed and at the same time is suitable for any dataset irrespective of
245 the presence/absence of correlations between subspaces.

246 **3. The VLQ-ADC System**

247 In this section we introduce our GPU-based similarity retrieval system,
248 VLQ-ADC, that contains a two-layer hierarchical indexing structure based
249 on vector and line quantization and an asymmetric distance computation
250 method. VLQ-ADC incorporates a novel index structure that can index the
251 dataset points efficiently (Sec. 3.1). The indexing and encoding process will
252 be presented in Sec. 3.2, and the querying process is discussed in Sec. 3.3.

253 Comparing with the existing systems above, One major advantage of our
254 system is that our indexing structure can generate shorter and more accu-
255 rate candidate list for the query point, which will accelerate query speed by
256 a large margin. Another advantage of our system is that the improved asym-
257 metric distance computation method base on PQ encoding method provide a
258 higher search accuracy. In the remainder of this section we will use Figure 2
259 to illustrate our framework. We recall that commonly used notations are

260 summarized in Table 1.

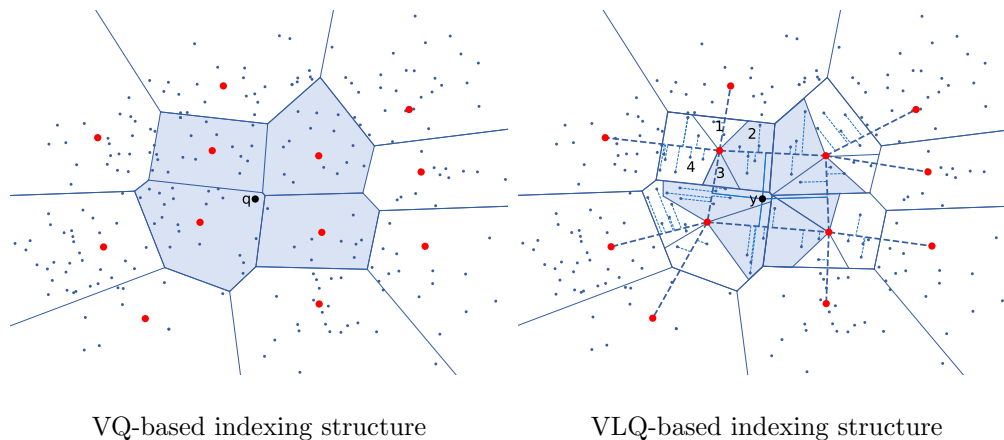


Figure 2: A comparison of the indexing structure and search process of the VQ-based indexing structure (**left**) and our VLQ-based indexing structure (**right**) on data points (small blue dots) of dimension 2 ($D = 2$). The large red dots denote the (first-level) same cell centroids in both figures. **Left:** The 4 shaded areas in the left figure represent the first-level regions, one for each centroid, and they make up the areas that need to be traversed for the query point q . **Right:** For each centroid in the right figure, $n = 4$ nearest neighboring centroids are found. Thus the n -NN graph consists of all the centroids and the edges (thick dashed lines) between them. Each first-level region in the right figure consists of 4 second-level regions, each of which represent the data points closet to the corresponding edge in the n -NN graph as denoted by the line quantizer q_i . Given the query point q and parameter $\alpha = 0.5$, only half of the second-level subregions (shaded in blue) need to be traversed. As can be seen, VLQ allows search to process substantially smaller regions in the dataset than a VQ-based approach.

261 *3.1. The VLQ-based index structure*

262 For billion-scale datasets with a moderate number of regions (e.g., 2^{16})
 263 produced by vector quantization (VQ), the number of data points in most
 264 regions is too large, which negatively affects search accuracy. To alleviate

265 this problem, we propose a hierarchical indexing structure. In our structure,
266 each list is split into several shorter lists, i.e., each region is divided into
267 several subregions, using line quantization (LQ).

268 Our indexing structure is a two-layer hierarchical structure which consists
269 of two levels of quantizers. The first level contains a vector quantizer q_v
270 with a codebook of k centroids. The vector quantizer q_v partitions the data
271 point space \mathcal{X} into k regions. The second level contains a line quantizer q_l
272 with an n -nearest neighbor (n -NN) graph. The n -NN graph is a directed
273 graph in which nodes are first-level centroids and edges connect a centroid
274 to its n nearest neighbors. In each first-level region, the line quantizer q_l
275 then quantizes each data point to the closest edge in the n -NN graph, thus
276 splitting the region into n second-level regions.

277 As an example, in the right side of Figure 2, given $n = 4$, the top left
278 first-level region is further divided into 4 subregions by q_l , enclosed by solid
279 lines and denoted 1, 2, 3, and 4. Each subregion contains all the data points
280 that are closest to a given edge of the n -NN graph, as calculated by the line
281 quantization q_l .

282 **Training the codebook.** We use Lloyd iteration in the fashion of the
283 Linde-Buzo-Gray algorithm [15] to obtain the codebook of the VQ quantizer
284 q_v . The n -NN graph is then built on the centroids of the codebook.

285 **Memory overhead of indexing structure.** One advantage of our
286 indexing structure is its ability to produce substantially more subregions with
287 little additional memory consumption. Same as VQ, our first layer codebook
288 needs $k \cdot D \cdot \text{sizeof}(\text{float})$ bytes. In addition, for second-level indexing, for
289 each of the k first-layer centroids, the n -NN graph only needs to store (1)

290 the indices of its n nearest neighbors and (2) the distances to its n nearest
 291 neighbors, which amounts to $k \cdot n \cdot (\text{sizeof}(\mathbf{int}) + \text{sizeof}(\mathbf{float}))$ bytes. Note
 292 we do not need to store the full-dimensional points. For a typical values of
 293 $k = 2^{16}$ centroids and $n = 32$ subcentroids, the additional memory overhead
 294 for storing the graph is $2^{16} \cdot 32 \cdot (32 + 32)$ bits (16 MB), which is acceptable
 295 for billion-scale datasets.

296 One way to produce the subregions is by utilizing vector quantization
 297 (VQ) again in each region. However, that would require storing full-dimensional
 298 subcentroids and thus consume too much additional memory. For the same
 299 configuration ($k = 2^{16}$ centroids and $n = 32$ subcentroids) and a dimension
 300 of $D = 128$, the additional memory overhead for a VQ-based hierarchical in-
 301 dexing structure would be $2^{16} \cdot 32 \cdot 128 \cdot \text{sizeof}(\mathbf{float})$ additional bits (1,024
 302 MB). As can be seen, our VLQ-based hierarchical indexing structure is sub-
 303 stantially more compact, only consuming 1/64 of the memory required by a
 304 VQ-based approach for the second-level codebook.

305 We note that the PQ-based indexing structure requires $\mathcal{O}(k \cdot (D + k))$
 306 memory to maintain the indexing structure (Table 2), which is memory in-
 307 efficient as it is quadratic in k . This is a limitation of PQ-based indexing
 308 structure. In contrast, the space complexity of our hierarchical indexing
 309 structure is $\mathcal{O}(k \cdot (D + n))$, where typically $n \ll k$ (n is much smaller than
 310 k), hence making our index much more memory efficient.

311 3.2. Indexing and encoding

312 In this subsection, we will describe the indexing and encoding process
 313 and summarize both processes in Algorithm 1 and 2 respectively.

314 For our two-level index structure, the indexing process comprises two

Algorithm 1 VLQ-ADC batch indexing process

```
1: function INDEX( $[x_1, \dots, x_N]$ )
2:   for  $t \leftarrow 1 : N$  do
3:      $x_t \mapsto q_v(x) = \arg \min_{c \in C} \|x_t - c\|^2$  // VQ
4:      $S_i = n\text{-arg} \min_{c \in C} \|c - c_i\|^2$  // Construct the  $n$ -NN graph
5:      $x_t \mapsto q_l(x) = \arg \min_{l(c_i, s_{ij}), s_{ij} \in S_i} d(x, l(c_i, s_{ij}))$  // LQ
6:   end for
7: end function
```

315 different quantization procedures, one for each layer. Similar to the IVFADC
316 scheme, each dataset point is quantized by the vector quantizer q_v to the first-
317 level regions surrounded by the dotted lines in Figure 2. These regions form
318 a set of inverted lists as search candidates.

We describe the second-level indexing process as follows. Let \mathcal{X}^i be a region of $\{x_1, \dots, x_l\}$ that corresponds to a centroid c_i , for $i \in \{1, \dots, k\}$. In constructing the n -NN graph, let $S_i = \{s_{i1}, \dots, s_{in}\}$ denote the set of the n centroids closest to c_i and $l(c_i, s_{ij})$ denote an edge between c_i and s_{ij} , for $j \in \{1, \dots, n\}$. The points in \mathcal{X}^i are quantized to the subregions by a line quantizer q_l with a codebook \mathcal{E}_i of n edges $\{l(c_i, s_{i1}), \dots, l(c_i, s_{in})\}$. Thus the region \mathcal{X}^i is split into n subregions $\{\mathcal{X}_1^i, \dots, \mathcal{X}_n^i\}$ and each point $x \in \mathcal{X}^i$ is quantized to a second-level subregion \mathcal{X}_j^i . So the entire space \mathcal{X} are divided into $k \times n$ second-level subregions.

$$\mathcal{X}_j^i = \{x \in \mathcal{X}^i \mid q_l(x) = l(c_i, s_{ij})\}, \text{ for all } i \in \{1 \dots k\} \quad (11)$$

Each data point in the dataset \mathcal{X} is assigned to one of the $k \cdot n$ cells. When the data point x is quantized to the sub-region of edge $l(c_i, s_{ij})$, according to

Algorithm 2 VLQ-ADC batch encoding process

```
1: function ENCODE( $[x_1, \dots, x_N]$ )
2:   for  $t \leftarrow 1 : N$  do
3:      $r_{q_l}(x_t) = x_t - ((1 - \lambda_{ij}) \cdot c_i + \lambda_{ij} \cdot s_{ij})$            // Equation 12
4:     let  $r_t = r_{q_l}(x_t)$                                            // displacement
5:      $r_t = [r_t^1, \dots, r_t^m]$                                        // divide  $r_t$  into  $m$  subvectors
6:     for  $p \leftarrow 1 : m$  do
7:        $r_t^p \mapsto c_{j_p} = \arg \min_{c_{j_p} \in C^p} \| r_t^p - c_p \|^2$ 
8:     end for
9:      $Code_t = (j_1, \dots, j_m)$ 
10:  end for
11: end function
```

the Equation 9 and 8 the displacement of x from the corresponding anchor point can be computed as following:

$$r_{q_l}(x) = x - ((1 - \lambda_{ij}) \cdot c_i + \lambda_{ij} \cdot s_{ij}), \text{ where} \quad (12)$$

$$\lambda_{ij} = -0.5 \cdot \frac{(\|x - s_{ij}\|^2 - \|x - c_i\|^2 - \|s_{ij} - c_i\|^2)}{\|s_{ij} - c_i\|^2}. \quad (13)$$

319 .

320 As shown in Algorithm 2, the value of $r_{q_l}(x)$ is first computed by Equa-
321 tion 12 and encoded into m bytes using PQ [11]. The PQ codebooks are
322 denoted by C^1, \dots, C^m , each containing 256 sub-centroids. The vector $r_{q_l}(x)$
323 is mapped to a concatenation of m sub-centroids $(c_{j_1}^1, \dots, c_{j_m}^m)$, for j_i is a
324 value between 1 and 256. Hence the vector $r_{q_l}(x)$ is encoded into an m -byte
325 code of sub-centroid index (j_1, \dots, j_m) . In Figure 1(c), we assume that c_i
326 is the closest centroid to x and can observe that the anchor point of each

327 point x is closer to x than c_i . So the dataset points can be encoded more
328 accurately with the same code length. This will improve the recall rate of
329 search, as can be seen in our evaluation in Section 5.

330 From Equation 13, the value of λ_{ij} for each point can be computed. it
331 is a `float` type value and requires 4 bytes for each data point. To further
332 improve memory efficiency, we quantize it into 256 values and encode it by a
333 byte. Empirically we find that the encoded λ_{ij} still exhibits high recall rates.

334 3.3. Query

335 One important advantage of our indexing structure is that at query time,
336 a specific query point only needs to traverse a small number of cells whose
337 edges are closest to the query point, as shown in the right side of Figure 2.
338 There are three steps for query processing: (1) region traversal, (2) distance
339 computation and (3) re-ranking.

340 3.3.1. Region traversal

341 The region traversal process consists of two steps: first-level regions
342 traversal and second-level regions traversal. During first-level regions traver-
343 sal, a query point y is quantized to its w_1 nearest first-level regions, which
344 correspond to $w_1 \cdot n$ second-level regions produced by quantizer q_v . The
345 subregions traversal is performed within only the $w_1 \cdot n$ second-level regions.
346 Moreover, y is quantized again to w_2 nearest second-level regions by quan-
347 tizer q_l . Then the candidate list of y is formed by the data points only within
348 the w_2 nearest second-level regions. Because the w_2 second-level regions is
349 obviously smaller than the w_1 first-level regions, the candidate list produced
350 by our VLQ-based indexing structure is shorter than that produced by the

351 VQ-based indexing structure. This will result in a faster query speed.

352 We use parameter α to determine the percentage of $w_1 \cdot n$ second-level
 353 regions to be traversed give a query, such that $w_2 = \alpha \cdot w_1 \cdot n$. We conduct a
 354 series of experiments in Section 5 to discuss the performance of our system
 355 with different values of α .

356 3.3.2. Distance computation

Distance computation is a prerequisite condition for re-ranking. In this section, we describe how to compute the approximate distance between a query point y to a candidate point x . According to [11], the distance from y to x can be evaluated by asymmetric distance computation (ADC) as follows:

$$\| y - q_1(x) - q_2(x - q_1(x)) \|^2 \quad (14)$$

357 where $q_1(x) = (1 - \lambda_{ij}) \cdot c_i + \lambda_{ij} \cdot s_{ij}$ and $q_2(\dots)$ is the PQ approximation of
 358 the x_i displacement.

Expression 14 can be further decomposed as follows [3]:

$$\begin{aligned} & \| y - q_1(x) \|^2 + \| q_2(\dots) \|^2 + 2\langle q_1(x), q_2(\dots) \rangle - \\ & 2\langle y, q_2(\dots) \rangle. \end{aligned} \quad (15)$$

359 where $\langle \cdot, \cdot \rangle$ denotes the inner product between two points.

If $l(c_i, s_{ij})$ is the closest edge to x , i.e., $q_1(x) = (1 - \lambda_{ij})c_i + \lambda_{ij}s_{ij}$, Expression 15 can be transformed in the following way:

$$\begin{aligned} & \underbrace{\| y - ((1 - \lambda_{ij})c_i + \lambda_{ij}s_{ij}) \|^2}_{\text{term1}} + \underbrace{\| q_2(\dots) \|^2}_{\text{term2}} + \\ & 2(1 - \lambda_{ij}) \underbrace{\langle c_i, q_2(\dots) \rangle}_{\text{term3}} + 2\lambda_{ij} \underbrace{\langle s_{ij}, q_2(\dots) \rangle}_{\text{term4}} - \underbrace{2\langle y, q_2(\dots) \rangle}_{\text{term5}}. \end{aligned} \quad (16)$$

According to Equation 7, term1 in Expression 16 can be computed in the following way:

$$\begin{aligned} \|y - ((1 - \lambda_{ij})c_i + \lambda_{ij}s_{ij})\|^2 = & (1 - \lambda_{ij}) \underbrace{\|y - c_i\|^2}_{\text{term6}} + \\ & (\lambda_{ij}^2 - \lambda_{ij}) \underbrace{\|c_i - s_{ij}\|^2}_{\text{term7}} + \lambda_{ij} \underbrace{\|y - s_{ij}\|^2}_{\text{term8}}. \end{aligned} \quad (17)$$

360 In Expression 16 and Equation 17, some computations can be done in
361 advance and stored in lookup table as follows:

- 362 • All of term2, term3, term4 and term7 are independent of the query.
363 They can be precomputed from the codebooks. Term2 is the squared
364 norm of the displacement approximation and can be stored in a table
365 of size $256 \times m$. Term7 is the square of the length of the edge that the
366 point x belongs to and is already computed in the codebook learning
367 process. Term3 and term4 are scalar products of the PQ sub-centroids
368 and the corresponding first-level centroid subvectors and can be stored
369 in a table of size $k \times 256 \times m$.
- 370 • Term6 and term8 are the distances from the query point to the first-
371 layer centroids. They are the by-product of first-layer traversal.
- 372 • Term5 is the scalar product of the PQ sub-centroids and the corre-
373 sponding query subvectors and can be computed independently before
374 the search. Its computation costs $256 \times D$ multiply-adds [6].

375 The proposed decomposition is used to simplify the distance computation.
376 With the lookup tables, the distance computation only requires $256 \times D$
377 multiply-adds and $2 \times m$ lookup-adds. In comparison, the classic IVFADC

378 distance computation requires $256 \times D$ multiply-adds and m lookup-adds
379 [6]. The additional m lookup-adds in our framework improves the distance
380 computation accuracy with a moderate increase of time overhead. We will
381 discuss this trade-off in detail in Section 5.

382 3.3.3. *Re-ranking*

383 Re-ranking is a step of re-sorting the candidate list of data points accord-
384 ing to the distances from candidate points to the query point. It is the last
385 step of the query process. The purpose of re-ranking is to find out the near-
386 est neighbours to the query point among the candidate points by distance
387 comparing. We apply the fast sorting algorithm of [6] to our re-ranking step.
388 Due to the shorter candidate list and more accurate approximate distances,
389 the re-ranking step of our system is both faster and more accurate than that
390 of Faiss.

391 4. GPU Implementation

392 One advantage of our VLQ-ADC framework is that it is amenable to
393 implementations on GPUs. It is mainly because our searching and distance
394 computing algorithm that applied during query can be efficiently parallelized
395 on GPUs. In this work we have implemented our framework in CUDA.

396 There are three different levels of granularity of parallelism on GPU:
397 threads, blocks and grids. A block is composed of multiple threads, and a
398 grid is composed of multiple blocks. Furthermore, there are three memory
399 types on GPU. Global memory is typically 4–32 GB in size with 5–10×
400 higher bandwidth than CPU main memory [6], and can be shared by different
401 blocks. Shared memory is similar to CPU L1 cache in terms of speed and is

Algorithm 3 VLQ-ADC batch search process

```
1: function SEARCH( $[y_1, \dots, y_{n_q}], \mathcal{L}_1, \dots, \mathcal{L}_{k \times n}$ )
2:   for  $t \leftarrow 1 : n_q$  do
3:      $C_t \leftarrow w_1\text{-arg min}_{c \in C} \| y_t - c \|^2$ 
4:      $L_{LQ}^t \leftarrow w_2\text{-arg min}_{c_i \in C_t, s_{ij} \in S_i} \| y - (1 - \lambda_{ij}) \cdot c_i - \lambda_{ij} \cdot s_{ij} \|^2$  //
    described in Sec. 3.3.1
5:     Store values of  $\| y_t - c \|^2$ 
6:   end for
7:   for  $t \leftarrow 1 : n_q$  do
8:      $L_t \leftarrow []$ 
9:     Compute  $\langle y_t, q_2(\dots) \rangle$  // See term5 in Equation 16
10:    for  $L$  in  $L_{LQ}^t$  do
11:      for  $t'$  in  $\mathcal{L}_L$  do
12:        // distance evaluation described in Sec. 3.3.2
13:         $d \leftarrow \| y_t - q_1(x_{t'}) - q_2(x_{t'} - q_1(x_{t'})) \|^2$ 
14:        Append  $(d; L; j)$  to  $L_t$ 
15:      end for
16:    end for
17:  end for
18:   $R_t \leftarrow$  K-smallest distance-index pairs  $(d, t')$  from  $L_t$  // Re-ranking
19:  return  $R_t$ 
20: end function
```

402 only shared by threads within the same block. GPU register file memory has
403 the highest bandwidth and the size of register file memory on GPU is much
404 larger than that on CPU [6].

405 VLQ-ADC is able to utilize GPU efficiently for indexing and search.
406 For example, we use blocks to process D -dimensional query points and the
407 threads of a block to traverse the inverted lists. We use global memory to
408 store the indexing structures and compressed dataset that is shared by all
409 blocks and grids, and load part of lookup tables in the shared memory to
410 accelerate distance computation. As the GPU register file memory is very
411 large, we store structured data in the register file memory to increase the
412 performance of the sorting algorithm.

413 Algorithm 3 summarizes our search process that is implemented on GPU.
414 We use four arrays to store the information of the inverted index lists. The
415 first array stores the length of each index list, the second one stores the
416 sorted vector IDs of each list, and the third the fourth store the correspond-
417 ing codes and λ values of each list respectively. For an NVIDIA GTX Titan
418 X GPU with a 12GB of RAM, we load part of the dataset indexing struc-
419 ture in the global memory for different kernels, i.e., region size, data points
420 compressed codes and λ values of each list. A kernel is the unit of work (in-
421 struction stream with arguments) scheduled by the host CPU and executed
422 by GPUs [6]. We load the vector IDs on the CPU side, because vector IDs
423 are resolved only if re-ranking step determines K -nearest membership. This
424 lookup produces a few sparse memory reads in a large array, thus the IDs
425 stored on CPU can only cause a tiny performance cost.

426 Our implementation makes use of some basic functions from the Faiss
427 library, including matrix multiplication and the K -selection algorithm to im-

428 prove the performance of our approach⁴.

429 *K-selection.* The K-selection algorithm is a high-performance GPU-based
430 sorting method proposed by Faiss [6] and GSKNN [8]. The K-selection keep
431 intermediate data in the register file memory. It exchanges register data using
432 the warp shuffle instruction, enabling warp-wide parallelism and storage. The
433 warp is a 32-wide vector of GPU threads, each thread in the warp has up
434 to 255 32-bit registers in a shared register file. All the threads in the same
435 warp can exchange register data using the warp shuffle instruction.

436 *List search.* We use two kernels for inverted list search. The first kernel is
437 responsible for quantizing each query point to w_1 nearest first-level regions
438 (line 3 in Algorithm 3). The second kernel is responsible for finding out the w_2
439 nearest second-level regions for the query point (line 4 in Algorithm 3). The
440 distances between each query point and its w nearest centroids are stored
441 for further calculation. In the two kernels, we use a block of threads to
442 process one query point, thus a batch of n_q query points can be processed
443 concurrently.

444 *Distance computation and re-ranking.* After the inverted lists \mathcal{L}_i of each
445 query point are collected, there are up to $n_q \times w_2 \times \max |\mathcal{L}_i|$ candidate points
446 to process. During the distance computation and re-ranking process, pro-
447 cessing all the query points in a batch yields high parallelism, but can exceed
448 available GPU global memory. Hence, we choose a tile size $t_q < n_q$ based
449 on amount of available memory to reduce memory overhead, bounding its

⁴The source code will be released upon publication.

450 complexity by $\mathcal{O}(t_q \times w_2 \times \max|\mathcal{L}_i|)$.

451 We use one kernel to compute the distances from each query point to the
452 candidate points according to Expression 16, and sort the distances via the
453 K-selection algorithm in a separate kernel. The lookup tables are stored in
454 the global memory. In the distance computation kernel, we use a block to
455 scan all w_1 inverted lists for a single query point, and the significant portion
456 of the runtime is the $2 \times w_2 \times m$ lookups in the lookup tables and the linear
457 scanning of the \mathcal{L}_i from global memory.

458 In the re-ranking kernel, we refer to Faiss by using a two-pass K-selection.
459 First reduce $t_q \times w_2 \times \max|\mathcal{L}_i|$ to $t_q \times \tau \times K$ partial results, where τ is some
460 subdivision factor, then the partial results are reduced again via k-selection
461 to the final $t_q \times K$ results.

462 Due to the limited amount of GPU’s memory, if an index instance with
463 long encoding length cannot fit in the memory of a single GPU, it cannot
464 be processed on the GPU efficiently. Our framework supports multi-GPU
465 parallelism to process indexing instance of a long encoding length. For b
466 GPUs, we split the index instance into b parts, each of which can fit in the
467 memory of a single GPU. We then process the local search of n_q queries on
468 each GPU, and finally join the partial results on one GPU. Our multi-GPU
469 system is based on MPI, which can be easily extended to multiple GPUs on
470 multiple servers.

471 5. Experiments and Evaluation

472 In this section, we evaluate the performance of our system VLQ-ADC and
473 compare it to three state-of-the-art billion-scale retrieval systems that are

474 based on different indexing structures and implemented on CPUs or GPUs:
475 Faiss [6], Ivf-hnsw [7] and Multi-D-ADC [16]. All the systems are evalu-
476 ated on the standard metrics: accuracy and query time, with different code
477 lengths. All the experiments are conducted on a machine with two 2.1GHz
478 Intel Xeon E5-2620 v4 CPUs and two NVIDIA GTX Titan X GPUs with 12
479 GB memory each.

480 The evaluation is performed on two public benchmark datasets that are
481 commonly used to evaluate billion-scale ANN search: SIFT1B [22] of 10^9
482 128-D vectors and DEEP1B [5] of 10^9 96-D vectors. Each dataset has a
483 10,000 query set with the precomputed ground-truth nearest neighbors. For
484 our system, we sample 2×10^6 vectors from each dataset for learning all the
485 trainable parameters. We evaluate the search accuracy by the test result
486 Recall@ K , which is the rate of queries for which the nearest neighbors is in
487 the top K results.

488 Here we choose nprobe =64 for all the inverted indexing systems (Faiss,
489 Ivf-hnsw and VLQ-ADC), as 64 is a typical value for nprobe in the Faiss
490 system. The parameter max_codes that means the max number of candidate
491 data points for a query is only useful to CPU-based system (max_codes is
492 set to 100,000), hence for GPU-based systems like Faiss and VLQ-ADC,
493 max_codes parameter is not configured. In fact, we compute the distances of
494 query point to all the data points that are contained in the neighbor regions.

495 *5.1. Evaluation without re-ranking*

496 In experiment 1, we evaluate the index quality of each retrieval system.
497 We compare three different inverted index structures and two inverted multi-
498 index schemes with different codebooks sizes without the re-ranking step.

- 499 1. **Faiss**. We build a codebook of $k = 2^{18}$ centroids by k-means, and find
500 proposed inverted lists of each query by Faiss.
- 501 2. **Ivf-hnsw**. We use a codebook of $k = 2^{16}$ centroids by k-means, and
502 set 64 sub-centroids for each first-level centroid.⁵
- 503 3. **Multi-D-ADC**. We use two IMI schemes with two codebook sizes
504 $k = 2^{10}$ and $k = 2^{12}$ and choose the implementation from the Faiss
505 library for all the experiments.
- 506 4. **VLQ-ADC**. For our approach, we use the same codebook as Ivf-hnsw,
507 and a 64-edge k-NN graph with indexing and querying as described in
508 Section 3.2 and 3.3.⁶

509 The recall curves of each indexing approach are presented in Figure 3. On
510 both datasets, our proposed system VLQ-ADC (blue curve) outperforms the
511 other two inverted index systems and the Multi-D-ADC scheme with small
512 codebooks ($k = 2^{10}$) for all the reasonable range of X . Compared with the
513 Multi-D-ADC scheme with a larger codebook ($k = 2^{12}$), our system performs
514 better on DEEP1B, and almost equally well on SIFT1B.

515 On the DEEP1B dataset, the recall rate of our system is consistently
516 higher than that of all the other indexing structures. With a codebook
517 that is only 1/4 the size of Faiss’ codebook, the recall rate of our inverted
518 index is higher than Faiss. This demonstrates that the line quantization
519 procedure can further improve the index quality than the previous inverted

⁵We use the implementation of Ivf-hnsw that is available online
(<https://github.com/dbaranchuk/ivf-hnsw>) for all the experiments.

⁶The VLQ-ADC source code is available at <https://github.com/zjuchenwei/vector-line-quantization>.

520 index methods.

521 Even on the SIFT1B dataset, the performance of our indexing structure
522 is almost the same as that of IMI with much larger codebook $k = 10^{12}$ and
523 much better than other inverted index structures.

524 As shown in Figure 3, for the SIFT1B dataset, the IMI with $k = 2^{12}$ can
525 generate better candidate list than the inverted indexing structures. While
526 for the DEEP1B dataset, the performance of the IMI falls behind that of
527 the inverted indexing structures. The reasons are that SIFT vectors are
528 histogram-based and the subvectors are corresponding to the different sub-
529 spaces, which describe disjoint image parts that have weak correlations in the
530 subspace distributions. On the other hand, the DEEP vectors are produced
531 by CNN that have a lot of correlations between the subspaces. It can be
532 observed that the performances of our indexing structure is consistent across
533 the two datasets. This demonstrates that our indexing structure’s suitability
534 for different data distributions.

535 5.2. Evaluation with re-ranking

536 In experiment 2, we evaluate the recall rates with the re-ranking step. In
537 all systems the dataset points are encoded in the same way: indexing and
538 encoding. (1) **Indexing**: displacements from data points to the nearest cell
539 centroids are calculated. For VLQ-ADC the displacements are calculated
540 from data points to the nearest anchor points on the line. (2) **Encoding**:
541 the residual values are encoded into 8 or 16 bytes by PQ with the same
542 codebooks shared by all the cells. Here we compare the same four retrieval
543 systems as in experiments 1. All the configurations for the retrieval systems
544 are the same as in experiment 1. For the GPU-based systems, we evaluate

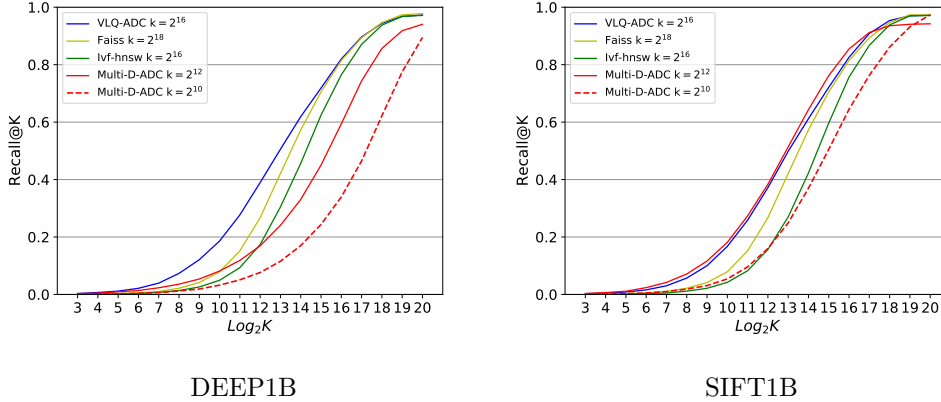


Figure 3: Recall rate comparison of our system, VLQ-ADC, without the re-ranking step, against two inverted index systems, Faiss, Ivf-hnsw, and one inverted multi-index scheme, Multi-D-ADC (with two different codebook sizes: $k = 2^{10}$ and $k = 2^{12}$).

545 performance with 8-byte codes on 1 GPU and 16-byte codes on 2 GPUs.

546 The Recall@ K values for different values $K = 1/10/100$ and the average
 547 query times on both datasets in milliseconds (ms) are presented in Table 3.
 548 From Table 3 we can make the following important observations.

549 **Overall best recall performance.** Our system VLQ-ADC achieves best
 550 recall performance for both datasets and the two codebooks (8-byte and
 551 16-byte) in most cases. For the twelve recall values (Recall@1/10/100 \times
 552 two codebooks \times two datasets), VLQ-ADC achieves best values in nine
 553 cases and second best in two cases. The second-best system is Faiss,
 554 obtaining best results in two cases. Multi-D-ADC (with $k = 2^{12} \times 2^{12}$
 555 regions) obtains best results in one case.

556 **Substantial speedup.** VLQ-ADC is consistently and significantly faster
 557 than all the other systems in all experiments. For all configurations,

558 VLQ-ADC’s query time is within 0.054–0.068 milliseconds, while the
559 other systems’ query times vary greatly. In the most extreme case,
560 VLQ-ADC is $125\times$ faster than Multi-D-ADC (0.068 vs 8.54). At the
561 same time, VLQ-ADC is also consistently faster than the second fastest
562 system, the GPU-based Faiss, by an average $5\times$ speedup.

563 **Comparison with Faiss.** VLQ-ADC outperforms the current state-of-the-
564 art GPU-based system Faiss in terms of both accuracy and query time
565 by a large margin, except for only three out of sixteen cases (R@10 with
566 16-byte codes for SIFT1B, and R@100 with 16-byte codes for SIFT1B
567 and DEEP1B). E.g., as a GPU-based system, VLQ-ADC outperforms
568 Faiss in terms of accuracy by 17%, 14%, 4% of R@1, R@10 and R@100
569 respectively on the SIFT1B dataset and 8-byte codes. At the same
570 time, the query time is consistently and significantly faster than Faiss,
571 with a speedup of up to $5.7\times$. Faiss outperforms VLQ-ADC in recall
572 values in three cases, all with 16-byte codes. However, the difference is
573 negligible ($\sim 0.02\%$). Similarly, though less pronounced, characteristics
574 can be observed on DEEP1B.

575 The main reason for this improvement is that the index quality and
576 encoding precision in VLQ-ADC is better than those of Faiss. Due
577 to the better indexing quality, the inverted list of our system is much
578 shorter than that of Faiss, which results in a much shorter query time.

579 Additionally, although the codebook size of our system ($k = 2^{16}$) is
580 only $1/4$ of that of Faiss ($k = 2^{18}$), our system produces more regions
581 (2^{22}) than Faiss (2^{18}). Therefore, our system achieves better accuracy

582 as well as memory and runtime efficiency than Faiss.

583 **Comparison with Multi-D-ADC.** The proposed system also outperforms
584 the IMI based system Multi-D-ADC both in terms of accuracy and
585 query time on both datasets. For example, VLQ-ADC leads Multi-D-
586 ADC with codebooks $k = 2^{12}$ by 14.2%, 7.4%, 1.3% of R@1, R@10
587 and R@100 respectively on the SIFT1B dataset and 8-byte codes with
588 up to $6.8\times$ speedup. On the DEEP1B dataset, the advantage of our
589 system is even more pronounced. Similarly, VLQ-ADC outperforms
590 Multi-D-ADC scheme with smaller codebooks $k = 2^{10}$ even more signif-
591 icantly, especially in terms of query time, where VLQ-ADC consistently
592 achieves speedups of at least one order of magnitude while obtaining
593 better recall values.

594 **Comparison with Ivf-hnsw.** Similarly, VLQ-ADC outperforms Ivf-hnsw,
595 another CPU-based retrieval system in both recall and query time. Al-
596 though Ivf-hnsw can also produce more regions with a small codebook,
597 it still cannot outperform the VQ-based indexing structure with larger
598 size of codebook.

599 **Effects on recall of indexing and encoding.** The improvement of R@10
600 and R@100 shows that the second-level line quantization provides more
601 accurate short-list of candidates than the previous inverted index struc-
602 ture, and the improvement of R@1 shows that it can also improves
603 encoding accuracy.

604 **Multi-D-ADC.** From Table 3, we can also observe that Multi-D-ADC
605 scheme with $k = 2^{12}$ outperforms the scheme with $k = 2^{10}$ in query

606 time by a large margin. It is mainly because Multi-D-ADC with larger
 607 codebooks can produce more regions, which can extract more concise
 608 and accurate short-lists of candidates.

Table 3: Performance comparison between VLQ-ADC (with the re-ranking step) against three other state-of-the-art retrieval systems of recall@1/10/100 and retrieval time on two public datasets. For each system the number of total regions is specified beneath each system’s name. VLQ-ADC consistently achieves higher recall values and significantly lower query time than all other systems. Best result in each column is **bolded**, and second best is underlined. For the two GPU-based systems, Faiss and VLQ-ADC, we experiments are performed on 1 GPU for 8-byte encoding length, and on 2 GPUs for 16-byte encoding length.

System	SIFT1B								DEEP1B							
	8 bytes				16 bytes				8 bytes				16 bytes			
	R@1	R@10	R@100	t (ms)	R@1	R@10	R@100	t (ms)	R@1	R@10	R@100	t (ms)	R@1	R@10	R@100	t (ms)
Faiss 2 ¹⁸	0.1383	0.4432	0.7978	<u>0.31</u>	0.3180	0.7825	0.9618	<u>0.280</u>	<u>0.2101</u>	<u>0.4675</u>	<u>0.7438</u>	<u>0.32</u>	<u>0.3793</u>	0.7650	0.9509	<u>0.33</u>
Ivf-hnsw 2 ¹⁶	<u>0.1599</u>	<u>0.496</u>	0.778	2.35	0.331	0.737	0.8367	2.77	0.217	0.467	0.7195	2.30	0.3646	0.7096	0.828	3.07
Multi-D-ADC 2 ¹⁰ × 2 ¹⁰	0.1255	0.4191	0.7843	1.65	0.3064	0.7716	0.9782	8.54	0.1716	0.3834	0.6527	3.28	0.324	0.6918	0.9258	6.152
Multi-D-ADC 2 ¹² × 2 ¹²	0.1420	0.4720	<u>0.8183</u>	0.367	<u>0.3324</u>	<u>0.8029</u>	<u>0.9752</u>	1.603	0.1874	0.4240	0.6979	0.839	0.3557	0.7087	0.9059	1.52
VLQ-ADC 2 ¹⁶	0.1620	0.507	0.829	0.054	0.345	0.8033	0.9400	0.068	0.2227	0.4855	0.7559	0.059	0.394	<u>0.7644</u>	<u>0.9272</u>	0.067

609 *5.3. Data point distributions of different indexing structures*

610 The space and time efficiency of an indexing structure is impacted by
 611 the distribution of data points produced by the structure. To analyse the
 612 distribution produced by the structures studied in this paper, we plot in
 613 Figure 4 the percentages of regions by the discretized number of data points
 614 in each region.

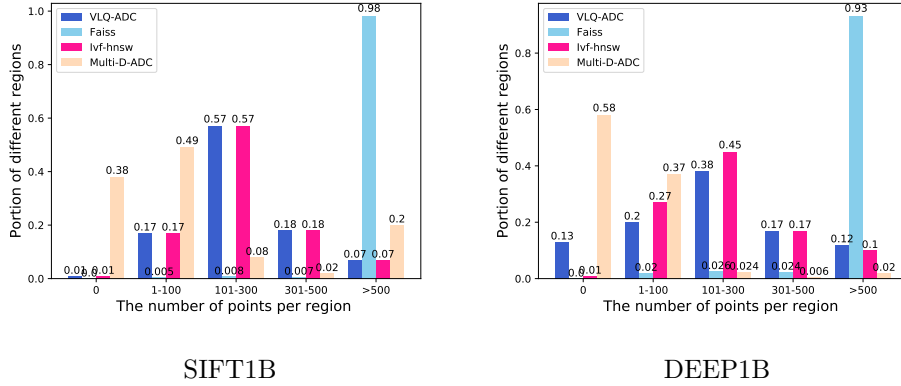


Figure 4: The distributions of data points in regions produced by different indexing structures. The x axis is five categories representing the discretized numbers of data points in each region (0, 1–100, 101–300, 301–500 and > 500). The y axis is the percentage of regions in each different categories.

615 As shown in Figure 4, the portion of empty regions produced by the
 616 inverted indexing structures (Faiss, Ivf-hnsw and VLQ-ADC) is much less
 617 than that produced by the inverted multi-index structure (Multi-D-ADC).
 618 For Multi-D-ADC, there are 38% empty regions for SIFT1B and 58% empty
 619 regions for DEEP1B (left most group in each plot). This result empirically
 620 validates the space inefficiency of inverted multi-index structure [16].

621 For Faiss, which is based on the inverted indexing structure using VQ,
 622 over 98% and 93% of regions contain more than 500 data points for SIFT1B
 623 and DEEP1B respectively. This will possibly produce long candidate lists
 624 for queries, thus negatively impacting query speed. For VLQ-ADC (and
 625 Ivf-hnsw), the regions are much more evenly distributed. The majority of
 626 the regions on both datasets contain less than 500 data points, and more
 627 regions contain 101–300 data points than others. This is a main reason why

628 VLQ-ADC can provide shorter candidate lists and thus a faster query speed.

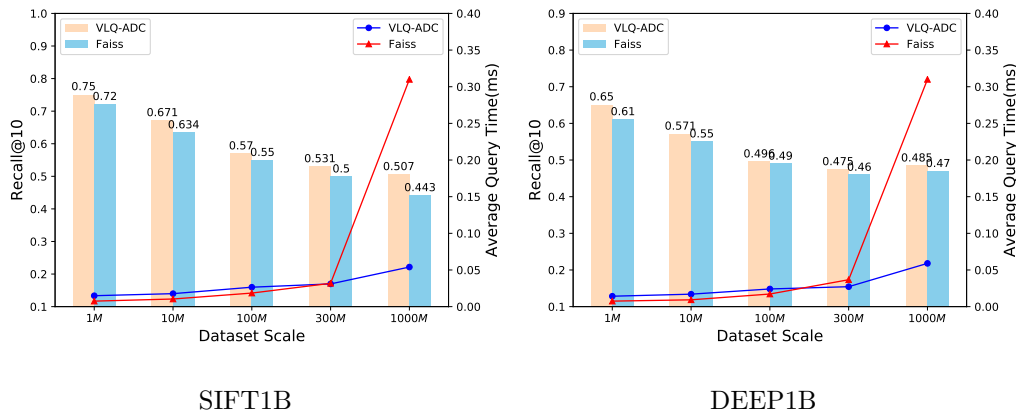


Figure 5: Comparison of recall@10 and average query time between VLQ-ADC and Faiss under different dataset scales. The two systems are compared with an 8-byte encoding length. The x axis indicates the five data scales (1M/10M/100M/300M/1000M). The left y axis is the recall@10 value (represented by the bars) and the right y axis is the average query time (in ms, represented by the lines).

629 5.4. Performance comparison under different dataset scales

630 In this section we evaluate the performance of our system under differ-
 631 ent dataset scales. Figure 5 shows, for SIFT1B and DEEP1B, the recall
 632 and query time values for Faiss and VLQ-ADC for subsets of SIFT and
 633 DEEP1B of different sizes: 1M, 10M, 100M, 300M and 1000M (full dataset)
 634 respectively. As can be seen in the figure, the recall values of VLQ-ADC is
 635 always higher than that of Faiss under all dataset scales. When the scale of
 636 dataset is under 300M, the query speed of Faiss is slightly faster than that
 637 of VLQ-ADC. When the scale of the datasets is over 300M, the query speed
 638 of VLQ-ADC matches that of Faiss.

639 It can also be observed from the figure that for the full datasets of
640 SIFT1B and DEEP1B (1000M), Faiss takes 0.31ms and 0.32ms respectively
641 (see Table 3 too). Compared to the 100M subsets of these two datasets, Faiss
642 suffers an approx. $15\times$ slowdown when data scale grows $10\times$. On the other
643 hand, for these two datasets, VLQ-ADC takes 0.054ms and 0.059ms respec-
644 tively, representing only an approx. $2\times$ slowdown when data scale grows $10\times$.
645 The superior scalability and robustness of VLQ-ADC over Faiss is evident
646 from this experiment.

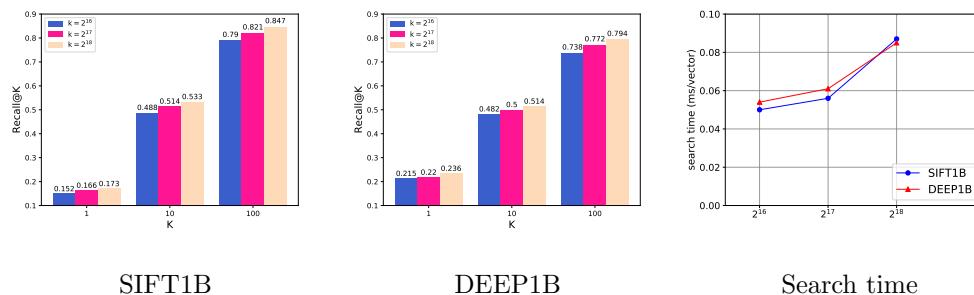


Figure 6: The performance of VLQ-ADC on different numbers of centroids $k = 2^{16}/2^{17}/2^{18}$. The results are collected on the same two datasets with an 8-byte encoding length and $n = 32$ edges of each centroids. The right plot shows the average search time with different values of k .

647 *5.5. Evaluation on impact of parameter values*

648 **Number of centroids k and edges n .** We evaluate the performance of
649 VLQ-ADC on different k and n values with 8-byte codes. We first fix the
650 value of n to 64 and compare the performance of our system for different
651 k centroids. In Figure 6, we present the evaluation of VLQ-ADC for $k =$
652 $2^{16}/2^{17}/2^{18}$. Then we fix $k = 2^{16}$ and increase the number of edge n from 32

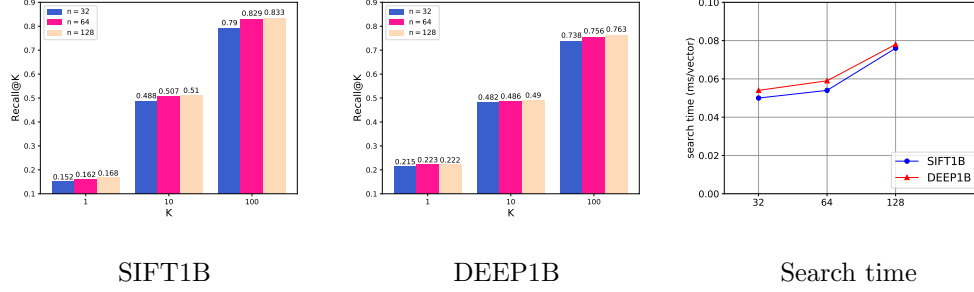


Figure 7: The performance of VLQ-ADC on different numbers of graph edges $n = 32/64/128$. The results are collected on the same two datasets with an 8-byte encoding length and $k = 2^{16}$ number of centroids. The right plot shows the average search time with different values of n .

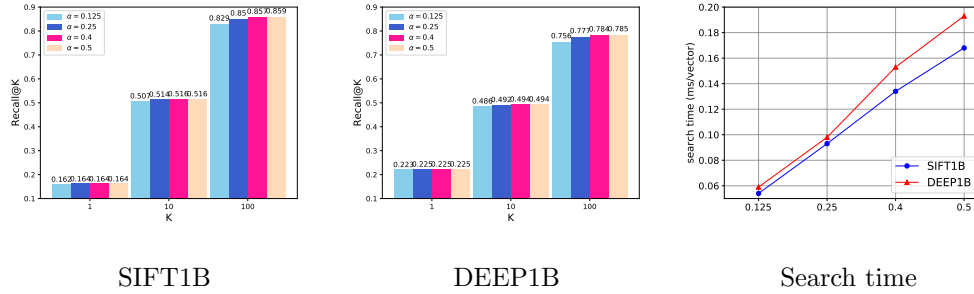


Figure 8: The performance of VLQ-ADC on different values of parameter $\alpha = 0.25/0.4/0.5$, with values of k, n and w fixed at $k = 2^{16}, n = 64, w = 64$. The result are collected on the same two datasets with an 8-byte encoding length and 64 edges of each centroids. The right plot shows the average search time with different values of α .

653 to 64 and 128. In Figure 7, we present the evaluation of the VLQ-ADC for
 654 different edge numbers.

655 From Figure 6 and 7 we can observe that the increase in the number
 656 of centroids and edges can improve search accuracy, while slightly increas-
 657 ing query time. This is because the indexing scheme with more centroids

658 and more edges can represent the dataset points more accurately and hence
659 provide more accurate short inverted lists.

660 **Value of portion α .** Now, we discuss how to determine the value of pa-
661 rameter α for subregions pruning, as described in Section 3.3.1. As shown in
662 Figure 8, we test several values of α on both datasets. A lower α value means
663 fewer subregions will be traversed, hence lower query time. At the same time,
664 we can observe that higher α values only moderately increase recall values,
665 while significantly increases query time (up to $3.7\times$ times). Hence we choose
666 $\alpha = 0.25$.

667 **Time and memory consumption.** Because the billion-scale dataset do
668 not fit on the GPU, the database is built in batches of 2M vectors, then
669 aggregating the information on the CPU. With file I/O, it takes about 150
670 minutes to build the whole database on a single GPU.

671 Here we analyze the memory consumption of each system. As shown in
672 Table 2, for a database of $N = 10^9$ points, the basic memory consumption for
673 all systems is $4 \cdot N$ bytes for point IDs that are Integer type and $m \cdot N$ bytes
674 for point codes. In addition to that, Multi-D-ADC consumes $4 \cdot k^2$ bytes to
675 store the region boundaries. Faiss consumes $4 \cdot k \cdot D$ bytes for the codebooks
676 and $4 \cdot k \cdot m \cdot 256$ bytes for the lookup tables. Ivf-hnsw requires N bytes for
677 quantized norm items $4 \cdot k \cdot (D + n)$ bytes for its indexing structure[7]. For
678 our system, we require N bytes for quantized λ values and $4 \cdot k \cdot (D + 2n +$
679 $m \cdot 256)$ bytes for the codebook, the n -NN graph and the lookup tables. We
680 summarize the total memory consumption for all systems in Table 4 with
681 8-byte encoding length on both datasets.

682 As presented in Table 4, the memory consumption of our system is less

683 than that of Faiss, and about 10% more than that of Multi-D-ADC with 2^{12}
 684 codebook, which is acceptable for most realistic setups.

Table 4: The memory consumption of all systems for SIFT1B of 10^9 128-dimensional data points.

System (codebook size)	Memory consumption (GB)
Faiss (2^{18})	14
Ivf-hnsw (2^{16})	13.04
Multi-D-ADC ($2^{12} \times 2^{12}$)	12.25
VLQ-ADC (2^{16})	13.55

685 6. Conclusion

686 Billion-scale approximate nearest neighbor (ANN) search has become an
 687 important task as massive amounts of visual data becomes available online.
 688 In this work, we proposed VLQ-ADC, a simple yet scalable indexing struc-
 689 ture and a retrieval system that is capable of handling billion-scale datasets.
 690 VLQ-ADC combines line quantization with vector quantization to create a
 691 hierarchical indexing structure. Search space is further pruned to a por-
 692 tion of the closest regions, further improving ANN search performance. The
 693 proposed indexing structure can partition the billion-scale database in large
 694 number of regions with a moderate size of codebook, which solved the draw-
 695 back of prior VQ-based indexing structures.

696 We performed comprehensive evaluation on two billion-scale benchmark
 697 datasets: SIFT1B and DEEP1B and three state-of-the-art ANN search sys-

698 tems: Multi-D-ADC, Ivf-hnsw, and Faiss. Our evaluation shows that VLQ-
699 ADC consistently outperforms all three systems on both recall and query
700 time. VLQ-ADC achieves a recall improvement over Faiss, the state-of-the-
701 art GPU-based system, of up to 17% and a query time speedup of up to $5\times$
702 times.

703 Moreover, VLQ-ADC takes the data distribution into account in the in-
704 dexing structure. As a result, it performs well on datasets with different
705 distributions. Our evaluation shows that VLQ-ADC is the best performer
706 on both SIFT1B and DEEP1B, demonstrating its robustness with respect to
707 data with different distributions.

708 We conclude by pointing out a number of future work directions. We plan
709 to investigate further improvements to the indexing structure. Moreover,
710 a more systematic and principled method for hyperparameter selection is
711 worthy investigation.

712 **Acknowledgment**

713 This work is supported in part by the National Natural Science Foun-
714 dation of China under Grant No.61672246, No.61272068, No.61672254 and
715 the Fundamental Research Funds for the Central Universities under Grant
716 HUST:2016YXMS018. In addition, we gratefully acknowledge the support of
717 NVIDIA Corporation with the donation of the Titan Xp GPUs used for this
718 research. The authors appreciate the valuable suggestions from the anony-
719 mous reviewers and the Editors.

720 **References**

- 721 [1] T. H. H. Chan, A. Guerqin, M. Sozio, Fully dynamic k-center clustering,
722 in: World Wide Web Conference, 2018, pp. 579–587.
- 723 [2] R. Weber, H.-J. Schek, S. Blott, A quantitative analysis and perfor-
724 mance study for similarity-search methods in high-dimensional spaces,
725 in: VLDB, volume 98, 1998, pp. 194–205.
- 726 [3] A. Babenko, V. Lempitsky, Improving bilayer product quantization for
727 billion-scale approximate nearest neighbors in high dimensions, arXiv
728 preprint arXiv:1404.1831 (2014).
- 729 [4] P. Wieschollek, O. Wang, A. Sorkine-Hornung, H. Lensch, Efficient
730 large-scale approximate nearest neighbor search on the GPU, in: Pro-
731 ceedings of the IEEE Conference on Computer Vision and Pattern
732 Recognition, 2016, pp. 2027–2035.
- 733 [5] A. B. Yandex, V. Lempitsky, Efficient indexing of billion-scale datasets
734 of deep descriptors, in: Computer Vision and Pattern Recognition, 2016,
735 pp. 2055–2063.
- 736 [6] J. Johnson, M. Douze, H. Jégou, Billion-scale similarity search with
737 GPUs, arXiv preprint arXiv:1702.08734 (2017).
- 738 [7] D. Baranchuk, A. Babenko, Y. Malkov, Revisiting the inverted indices
739 for billion-scale approximate nearest neighbors, in: Proceedings of the
740 European Conference on Computer Vision (ECCV), 2018, pp. 202–216.

- 741 [8] C. D. Yu, J. Huang, W. Austin, B. Xiao, G. Biros, Performance op-
742 timization for the k-nearest neighbors kernel on x86 architectures, in:
743 Proceedings of the International Conference for High Performance Com-
744 puting, Networking, Storage and Analysis, ACM, 2015, p. 7.
- 745 [9] Y. Gong, S. Lazebnik, A. Gordo, F. Perronnin, Iterative quantization:
746 A procrustean approach to learning binary codes for large-scale image
747 retrieval, *IEEE Transactions on Pattern Analysis and Machine Intelli-*
748 *gence* 35 (2013) 2916–2929.
- 749 [10] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image
750 recognition, in: Proceedings of the IEEE conference on computer vision
751 and pattern recognition, 2016, pp. 770–778.
- 752 [11] H. Jégou, M. Douze, C. Schmid, Product quantization for nearest neigh-
753 bor search, *IEEE Transactions on Pattern Analysis and Machine Intel-*
754 *ligence* 33 (2011) 117.
- 755 [12] Y. Kalantidis, Y. Avrithis, Locally optimized product quantization for
756 approximate nearest neighbor search, in: *Computer Vision and Pattern*
757 *Recognition*, 2014, pp. 2329–2336.
- 758 [13] M. Muja, D. G. Lowe, Scalable nearest neighbor algorithms for high
759 dimensional data, *IEEE Transactions on Pattern Analysis and Machine*
760 *Intelligence* 36 (2014) 2227–2240.
- 761 [14] M. Norouzi, D. J. Fleet, Cartesian k-means, in: *IEEE Conference on*
762 *Computer Vision and Pattern Recognition*, 2013, pp. 3017–3024.

- 763 [15] Y. Linde, A. Buzo, R. Gray, An algorithm for vector quantizer design,
764 IEEE Transactions on communications 28 (1980) 84–95.
- 765 [16] A. Babenko, V. Lempitsky, The inverted multi-index, in: Computer
766 Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on,
767 IEEE, 2012, pp. 3069–3076.
- 768 [17] A. S. Razavian, H. Azizpour, J. Sullivan, S. Carlsson, CNN features
769 off-the-shelf: An astounding baseline for recognition, in: 2014 IEEE
770 Conference on Computer Vision and Pattern Recognition Workshops,
771 2014, pp. 512–519.
- 772 [18] Y. Gong, L. Wang, R. Guo, S. Lazebnik, Multi-scale orderless pooling
773 of deep convolutional activation features, in: European conference on
774 computer vision, Springer, 2014, pp. 392–407.
- 775 [19] S. Lloyd, Least squares quantization in pcm, IEEE transactions on
776 information theory 28 (1982) 129–137.
- 777 [20] T. Ge, K. He, Q. Ke, J. Sun, Optimized product quantization for approx-
778 imate nearest neighbor search, in: Proceedings of the IEEE Conference
779 on Computer Vision and Pattern Recognition, 2013, pp. 2946–2953.
- 780 [21] J. Sivic, A. Zisserman, Video Google: a text retrieval approach to
781 object matching in videos, in: Proceedings Ninth IEEE International
782 Conference on Computer Vision, 2003, pp. 1470–1477.
- 783 [22] H. Jégou, R. Tavenard, M. Douze, L. Amsaleg, Searching in one billion
784 vectors: re-rank with source coding, in: Acoustics, Speech and Signal

785 Processing (ICASSP), 2011 IEEE International Conference on, IEEE,
786 2011, pp. 861–864.