
Revisiting Multi-Codebook Quantization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-Codebook Quantization (MCQ) is a generalized version of existing codebook-
2 based quantizations for Approximate Nearest Neighbor (ANN) search. Therefore,
3 MCQ theoretically has the potential to achieve the best performance because so-
4 lutions of other codebook-based quantization methods are all covered by MCQ’s
5 solution space under the same codebook size setting. However, finding the opti-
6 mal solution to MCQ is proved to be NP-hard due to its encoding process, *i.e.*,
7 converting an input vector to a binary code. To tackle this, researchers apply
8 constraints to it to find near-optimal solutions, or employ heuristic algorithms
9 which are still time-consuming for encoding. Different from previous approaches,
10 this paper takes the first attempt to find a deep solution to MCQ. The encoding
11 network is designed to be as simple as possible, so the very complex encoding
12 problem becomes simply a feed-forward. Compared with other methods on three
13 datasets, our method shows state-of-the-art performance. Notably, our method
14 is 11×-38× faster than heuristic algorithms for encoding, which makes it more
15 practical for real scenery of large-scale retrieval. Our code is publicly available:
16 <https://github.com/DeepMCQ/DeepQ>.

17 1 Introduction

18 Rapidly increasing multimedia contents in recent years raise an urgent request for retrieval in a
19 short time. Unlike the exhaustive routine [31, 20], Approximate Nearest Neighbor (ANN) search
20 significantly reduces retrieval time while preserving high recall. It has been widely applied to various
21 scenarios, such as database indexing, fast image retrieval, and recommender systems.

22 As a typical approach, vector quantization (VQ) [7] is at first developed as a compression technique,
23 which uses a codebook to approximate vectors. People further find the power of VQ to preserve
24 similarities between quantized features and enable VQ to perform ANN search. In order to achieve
25 low quantization errors with limited codebook size, a multi-codebook structure is introduced. The
26 proposal of the Multi-Codebook Quantization (MCQ) [2] describes the approach as a combination
27 of one codeword for each sub-codebook, and previous methods [9, 6, 19, 30, 10, 3] are summarized
28 as exceptional cases of MCQ or *constrained MCQs*. The quantization codes are designed to be
29 compacted, which results in negligible storage cost and high-quality results.

30 However, the optimization of MCQ without any constraints is formally NP-hard. [14] models
31 it as the minimization on several fully-connected Markov Random Fields (MRFs). As a result,
32 current researches aim at solving MCQ under acceptable computational costs. Other than applying
33 constraints on it [34, 4, 15], another approach designs algorithms in a heuristic way [2, 14, 16]. The
34 latter achieves better performance but suffers from slow encoding.

35 There are chances to employ neural networks’ power to solve MCQ, where people expect to obtain
36 higher performance and encoding efficiency than previous methods. [11, 5, 28, 33, 27] already give
37 the way to treat codebook as network parameter and update it by gradient-descent, but they are

38 all still under constraints that hinder performance. Morozov and Babenko [18] and Sablayrolles *et*
 39 *al.* [22] map datapoints to learned space, which are not flexible, especially when performing the
 40 reconstruction. Therefore in this paper, we give our first attempt to solve MCQ in a deep learning
 41 approach, without constraints and work-arounds. Our contributions can be summarized as three-folds:

- 42 • Our novel approach, Deep Multi-Codebook Quantization (DeepQ), fully considers encoding
 43 difficulty and time complexity in MCQ. With the high efficient and parallelized encoding networks,
 44 our method significantly reduces encoding time.
- 45 • To tackle the NP-hard encoding problem and non-differentiable gradient estimation, we employ and
 46 further revise a policy gradient method. Value-Corrected Proximal Policy Optimization (VC-PPO)
 47 is proposed to speed up convergence in the training phase.
- 48 • Experiments conducted on a benchmark dataset validate our proposed method. Furthermore, to
 49 evaluate the scalability of the method, it is tested on million-scale datasets to show the effectiveness
 50 of our proposed algorithm.

51 2 Related Works

52 Vector quantization is a routine to approximate vectors by a codebook. Typical applications include
 53 clustering, compression, and Approximate Nearest Neighbor (ANN) search. The famous proposal
 54 k -means [7], also known as Lloyd’s algorithm [13], clusters the dataset into uniformly sized convex
 55 cells. When it is applied to ANN search, datapoints from the base set are quantized into their
 56 nearest centroids and represented by indices. The distance from a given query to any datapoint
 57 is approximated by the distance from the query to the datapoint’s centroid, which is effectively
 58 pre-computed and stored in a lookup table. To perform fine-grained clustering as well as reducing the
 59 space and time complexity, they [9, 6, 19, 10, 30] divide the feature space orthogonally by performing
 60 k -means in each subspace concurrently. Meanwhile, the introduced sub-codebook structure reveals
 61 the prototype of MCQ. Formally, [2] gives a well definition of MCQ, and previous works are all
 62 summarized into constrained MCQs. Specifically, subspace k -means must keep orthogonality among
 63 sub-codebooks. Zhang *et al.* [34] loosens the orthogonality constraint, but sub-codebooks are still
 64 weakly-orthogonal. Chen *et al.* [4] and Martinez *et al.* [15] propose hierarchical k -means, where
 65 vectors are quantized coarse-to-fine. If constraints are moved, MCQ is not easy to solve. Current
 66 state-of-the-art methods develop heuristic algorithms to help to encode. Specifically, Babenko and
 67 Lempitsky [2] employs beam search, Martinez *et al.* [14, 16] give algorithm based on Iterated
 68 Conditional Modes (ICM). However, the above methods do not achieve satisfied time complexity in
 69 encoding yet.

70 When neural networks and gradient descent become a fashion, a few attempts to integrate quantization
 71 into deep retrieval networks are proposed. Klein and Wolf [11] and Song *et al.* [5] propose Deep
 72 Product Quantization (DPQ) and Deep Progressive Quantization (DPgQ) which update codebook by
 73 soft relaxation, but they are still under the same constraints as [9, 15]. Sablayrolles *et al.* [22] and
 74 Morozov and Babenko [18] give pipelines to encode compact representations for compressed-domain
 75 search, but they do not strictly follow the paradigm of MCQ.

76 3 Preliminaries

77 Given a vector $\mathbf{x} \in \mathbb{R}^D$, its quantized vector $\tilde{\mathbf{x}}$ are composed by several codewords in a codebook
 78 \mathbf{C} . More Specifically, $\mathbf{C} = (\mathbf{C}_m)$, $\mathbf{C}_m \in \mathbb{R}^{K \times D}$, $1 \leq m \leq M$ contains M sub-codebooks and K
 79 codewords for each. Quantization codes are formed by $\mathbf{b} = (\mathbf{b}_m)$, $\mathbf{b}_m \in \{1, 2, \dots, K\}$, $1 \leq m \leq$
 80 M , which indicates the picked codeword in each sub-codebook. For the whole training set $\mathbf{X} = \{\mathbf{x}\}$
 81 with N datapoints, MCQ aims at finding the optimal quantization codes $\mathbf{B} = \{\mathbf{b}\}$ and codebook \mathbf{C}
 82 to minimize following objective:

$$\min_{\mathbf{C}, \mathbf{B}} \mathbb{E}_{\mathbf{x} \in \mathbf{X}} \mathbb{Q}(\mathbf{x}, \mathbf{b}, \mathbf{C}) = \min_{\mathbf{C}, \mathbf{B}} \mathbb{E}_{\mathbf{x} \in \mathbf{X}} \left\| \mathbf{x} - \sum_{m=1}^M \mathbf{C}_m \mathbf{b}_m \right\|_2 \quad (1)$$

83 where $\mathbf{C}_m \mathbf{b}_m \in \mathbb{R}^D$ is the \mathbf{b}_m -th codeword of the m -th sub-codebook. The sum of picked codewords
 84 $\sum \mathbf{C}_m \mathbf{b}_m$ tries to approximate \mathbf{x} . \mathbf{C} and \mathbf{b} are stored for further retrieval. Some of the previously
 85 mentioned methods [9, 6, 4, 15, 34] are treated as *constrained MCQs*, as they are all represented

86 as special cases of (1). Specifically, when $M = 1$, (1) becomes VQ. Or if any two sub-codebooks
 87 C_i, C_j are orthogonal, it will be PQ or OPQ.

88 The optimization of (1) without any constraints is proved to be NP-hard [14]. To tackle this, we
 89 propose a Expectation-Maximization style solution. Following sections will explain the deep neural
 90 network for encoding \mathbf{b} (Section 4.1), the way to solve C (Section 4.2), and how to conduct retrieval
 91 (Section 4.3), respectively.

92 4 Methodology

93 4.1 Expectation: Encoding B with neural networks

94 Our first step, is to find a potential code \mathbf{b} by given \mathbf{x} and
 95 a fixed C . A policy π parameterized by θ is employed to
 96 take possible solution of \mathbf{b} by feeding \mathbf{x} :

$$\pi = (\pi_m) = \pi(\mathbf{x} | \theta_m), 1 \leq m \leq M. \quad (2)$$

97 More specifically, π produces M Categorical distributions
 98 $\text{Categorical}(K, \mathbf{p}_{m1}, \dots, \mathbf{p}_{mK})$, where \mathbf{p}_{mj} is
 99 the probability to pick the j -th codeword in the m -th
 100 sub-codebook. A potential encoding \mathbf{b}_m is generated by
 101 drawing samples from π_m , which then helps us to pick
 102 codeword $C_m \mathbf{b}_m$. Therefore:

$$\mathbf{b}_m \sim \pi_m(\mathbf{x} | \theta_m) = \text{Categorical}(K, \mathbf{p}_m). \quad (3)$$

103 Since the independence among different sub-codebooks
 104 is a prerequisite of MCQ, \mathbf{b}_m should be drawn from π_m
 105 *independently*. Intuitively, the probability of \mathbf{b} to be a
 106 specific \mathbf{b}^* is derived by conditional independence:

$$\Pr(\mathbf{b} = \mathbf{b}^*) = \prod_{m=1}^M \Pr(\mathbf{b}_m = \mathbf{b}_m^*) = \prod_{m=1}^M \mathbf{p}_{m\mathbf{b}_m^*}. \quad (4)$$

107 We adopt the power of neural networks to model
 108 π_m . Specifically, θ_m produces K unnormalized log-
 109 probabilities ℓ_m and \mathbf{p}_{mj} is obtained by Softmax. To
 110 keep the independence, θ_m will not share parameters with
 111 each other.

112 Therefore, θ , or our proposed *IndepNet* is illustrated in
 113 Figure 1. We first build a basic structure called *IndepBlock* and duplicate this block for M times as
 114 $\theta_1, \theta_2, \dots, \theta_M$. We try to keep the basic structure really simple to achieve high efficiency during
 115 training and encoding. As the figure shows, *IndepBlock* is an hourglass network contains 6 layer-
 116 groups (consists of a linear layer with ReLU activation and layer-normalization) with skip-connections.
 117 The last three outputs are concatenated and further fed into a final linear layer with K outputs as
 118 $\ell_m = (\ell_{m1}, \dots, \ell_{mK})$, and therefore:

$$\mathbf{p}_{mj} = \text{Softmax}(\ell_m)_j, \text{ where } \ell_m = \theta_m(\mathbf{x}). \quad (5)$$

119 4.1.1 Gradient estimation

120 The objective of training θ is formed as:

$$\min_{\pi} \mathbb{E}_{\mathbf{x} \in \mathbf{X}} \mathbb{Q}(\mathbf{x}, \mathbf{b}, C). \quad (6)$$

121 However, the optimization faces two problems: 1) The encoding of \mathbf{b} involves sampling from discrete
 122 distributions, which is non-differentiable, 2) All possible encoding of \mathbf{b} is $\mathcal{O}(K^M)$. Exhaustive
 123 search becomes impracticable.

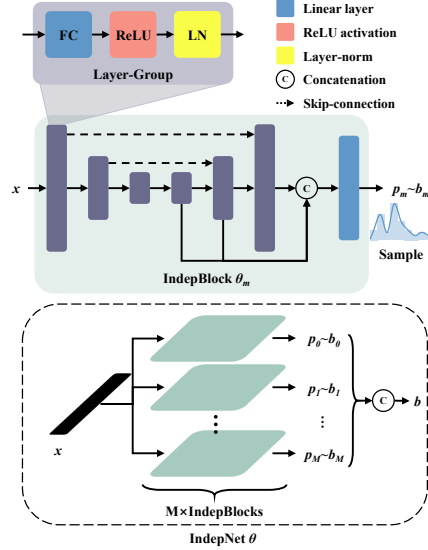


Figure 1: Our proposed *IndepNet* for producing probabilities of choosing each codeword. *IndepBlock* is duplicated for M times without shared parameters, in order to keep independence between different *IndepBlocks*. Categorical distribution is built upon output from the *IndepBlock*. Then, quantization code \mathbf{b}_m associated with sub-codebook C_m is sampled from distribution.

124 Therefore, gradient estimation over discrete, stochastic computation graph is required to train θ .
 125 Mainstream methods [23, 32, 17] include score function gradient estimator, pathwise gradient
 126 estimator, *etc.* Meanwhile, minimizing (6) is also faced with the high-variance problem during
 127 gradient estimation. To tackle this, the advantage function is introduced [12, 25]. Specifically in
 128 our work, a value network called *QENet* parameterized by τ is proposed to model a value function
 129 $v = V(\cdot | \tau)$. It performs a regression task to minimize the following objectives:

$$\min_{\tau} \mathbb{E}_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{b} \sim \pi(\mathbf{x} | \theta)}} \|\mathbf{Q}(\mathbf{x}, \mathbf{b}, \mathbf{C}) - V(\mathbf{x}, \mathbf{b}, \mathbf{C} | \tau)\|_2. \quad (7)$$

130 Advantages \hat{A} is then estimated by

$$\hat{A} = \mathbf{Q}(\mathbf{x}, \mathbf{b}, \mathbf{C}) - V(\mathbf{x}, \mathbf{b}, \mathbf{C} | \tau). \quad (8)$$

131 The detailed architecture of *QENet* is shown in Figure 2.
 132 We reuse the *IndepBlock* to generate v by $M + 1$ blocks:
 133 $\tau = (\tau_1, \dots, \tau_M, \tau_x)$. Specifically, latent representation
 134 for each selected-codeword \mathbf{C}_{mb_m} is obtained by:

$$\mathbf{l}_m = \tau_m(\mathbf{C}_{mb_m}). \quad (9)$$

135 The last *IndepBlock* τ_x is introduced to transform \mathbf{x} . Then,
 136 all the outputs from *IndepBlocks* are summed up to get
 137 scalar value v (denoted as “reduce-sum”):

$$v = \text{sum}(\mathbf{l}_1, \dots, \mathbf{l}_M, \mathbf{l}_x). \quad (10)$$

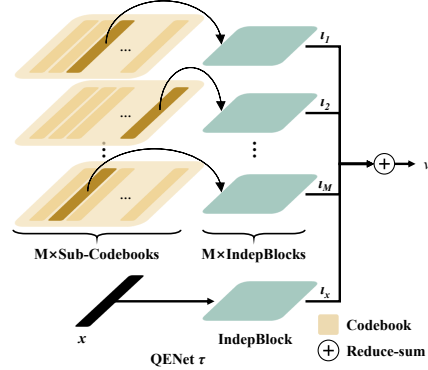


Figure 2: Our proposed *QENet* for advantage estimation. First M *IndepBlocks* are fed by M selected codewords and the last one is fed by \mathbf{x} . Outputs are summed up to get scalar value v .

138 **Value-corrected proximal policy optimization** We
 139 propose a variant of score function gradient estimator
 140 called Value Corrected Proximal Policy Optimization (VC-
 141 PPO) based on PPO to get simple but efficient Trust Re-
 142 gion updates [26, 24]. In the real scenario of large-scale
 143 ANN search, the training size N is usually larger than $10k$. Conventional PPO still does not satisfy
 144 us due to the speed of convergence. Therefore, we revise and propose the Value-Corrected PPO
 145 (VC-PPO) to achieve fast training. Firstly in the sampling stage, \mathbf{b}_o and v_o is produced from datapoint
 146 \mathbf{x} over whole training set \mathcal{X} by freezing current policy network and value network as θ_o, τ_o :

$$\begin{aligned} \mathbf{b}_o &\sim \pi(\mathbf{x} | \theta_o), \\ v_o &= V(\mathbf{x}, \mathbf{b}_o, \mathbf{C} | \tau_o). \end{aligned} \quad (11)$$

147 The probability of producing the sampled \mathbf{b}_o is denoted as $p_o = \Pr(\mathbf{b}_o | \mathbf{x}, \theta_o)$, calculated by
 148 equation (4). Finally, our surrogate objectives of VC-PPO is defined as [8]:

$$\begin{aligned} \mathcal{L}_\theta &= \min \left(\frac{\Pr(\mathbf{b}_o | \mathbf{x}, \theta)}{\Pr(\mathbf{b}_o | \mathbf{x}, \theta_o)} \hat{A}, \right. \\ &\quad \left. \text{clip}_{1-\epsilon}^{1+\epsilon} \left(\frac{\Pr(\mathbf{b}_o | \mathbf{x}, \theta)}{\Pr(\mathbf{b}_o | \mathbf{x}, \theta_o)} \right) \hat{A} \right), \end{aligned} \quad (12)$$

149

$$\begin{aligned} \mathcal{L}_\tau &= \max \left((\mathbf{Q}(\mathbf{x}, \mathbf{b}_o, \mathbf{C}) - V(\mathbf{x}, \mathbf{b}_o, \mathbf{C} | \tau))^2, \right. \\ &\quad \left. (\mathbf{Q}(\mathbf{x}, \mathbf{b}_o, \mathbf{C}) - v_o - \text{clip}_{-\epsilon}^{+\epsilon}(V(\mathbf{x}, \mathbf{b}_o, \mathbf{C} | \tau) - v_o))^2 \right). \end{aligned} \quad (13)$$

150 Here, The $\text{clip}(\cdot)$ forces the policy and value to be not too far from old ones and ϵ is the clip-range.
 151 In both equations, it prevents a large update ratio leading to an unstable policy. The key difference
 152 between the original PPO and our VC-PPO is, we use $V(\mathbf{x}, \mathbf{b}_o, \mathbf{C} | \tau)$ other than the recorded old
 153 value v_o from sampling stage to estimate advantage. This modification is treated as a value-correction
 154 process. Correcting value leads to a precise estimation on advantage, which is based on two reasons:
 155 a) Biases are introduced into advantage estimation if we use v_o , since the policy is getting better and
 156 better during training but τ_o is froze, and 2) The calculation of $V(\mathbf{x}, \mathbf{b}_o, \mathbf{C} | \tau)$ can be done instantly

157 without introducing significant computational overhead. To further encourage the network choose
 158 codewords uniformly, a regularization is applied to θ to maximize the entropy of π :

$$e_\theta = - \sum_{m=1}^M \sum_{j=1}^K p_{mj} \log p_{mj} \quad (14)$$

159 which forces network to try more codeword combinations.

160 4.2 Maximization: Solve C by least-squares

161 To give the closed-form derivation of solving C by given X and B , We will firstly rewrite Equation
 162 (1) to a matrix formulation. Since $\mathbf{b} = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_M)$ and $\mathbf{b}_m \in \{1, 2, \dots, K\}$ is the index of
 163 selected codeword in the i -th sub-codebook, a one-hot encoding and a concatenation on each \mathbf{b}_m :
 164 $\mathbf{b}'_m = \text{one-hot}(\mathbf{b}_m)$, $\mathbf{b}' = (\mathbf{b}'_1, \dots, \mathbf{b}'_M)$ will convert the quantization code to a M -hot vector *i.e.* a
 165 vector that contains M segments, and each segment contains exactly one 1 and remaining 0, where

166 1 is the entry of picked codeword. Correspondingly, a reshape is applied to C : $C' = \begin{pmatrix} C_1 \\ C_2 \\ \vdots \\ C_M \end{pmatrix} \in$

167 $\mathbb{R}^{(M \times K) \times D}$. (1) will become:

$$\min_{C'} \|X - B'C'\|_2^2. \quad (15)$$

168 This equation is formally a linear least-squares regression, where $B' \in \{0, 1\}^{N \times (M \times K)}$ is known
 169 and X is target. Although there is a bunch of algorithms to solve it, we finally choose *gelsy* [1],
 170 which in our experiments shows the best results. The solution is to first apply a QR factorization with
 171 column permutation on B' :

$$B' = Q \begin{pmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{pmatrix} P^\top \quad (16)$$

172 where Q and $R = \begin{pmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{pmatrix}$ is the factorization matrix and P is an orthogonal matrix that
 173 permutes columns of B' until R_{11} is well-conditioned (its estimated condition number approaches
 174 0). With the permutation, R_{22} becomes negligible. Moreover, R_{12} is erased by another orthogonal
 175 transformation:

$$\begin{pmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{pmatrix} \rightarrow \begin{pmatrix} R_{11} & R_{12} \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} T_{11} & 0 \\ 0 & 0 \end{pmatrix} Z \quad (17)$$

176 where T and Z are from the orthogonal transformation of R . Then, C' is derived by:

$$B' = Q \begin{pmatrix} T_{11} & 0 \\ 0 & 0 \end{pmatrix} Z P^\top, \quad (18)$$

$$C \leftarrow C' \leftarrow P Z^\top \begin{pmatrix} T_{11}^{-1} Q_1^\top X \\ 0 \end{pmatrix}$$

177 where Q_1 is the top $\text{rank}(B')$ columns of Q .

178 In brief, our overall training approach is summarized into algorithm 1.

179 4.3 Fast retrieval

180 After training, we are able to encode the base set for retrieval. Other than sampling from π , codewords
 181 are simply rolled out by greedy assignments:

$$\mathbf{b}_m^g = \arg \max \theta_m(\mathbf{x}). \quad (19)$$

182 We firstly use the greedy roll-out strategy to obtain B in the training set in order to solve the final
 183 codebook. Then, we employ the same strategy to encode the base set.

184 To further refine assignments, we add an extra step that randomly selects and alters \mathbf{b}_i while fixing
 185 others:

$$\mathbf{b}_i^g \leftarrow \arg \min_{\mathbf{b}_i^g} Q(\mathbf{x}, \mathbf{b}^g, \mathbf{C}), \quad (20)$$

$$i \sim \mathcal{U}[1, M].$$

186 Since this refinement only causes negligible overhead referred to the implementation by [14], in
 187 practice, we benefit from it not only to get lower quantization error but also to obtain acceptable
 188 performance from a fast training, *i.e.*, training within a very few steps before the network is converged.

189 The encoded and refined base set, combined with the codebook, is finally employed for retrieval. The
 190 LSQ-style lookup table [14] is utilized to speed up similarity search.

191 4.4 Discussion

192 Our work aims at solving Multi-Codebook Quantization via neural networks. Similar works include
 193 Unsupervised Neural Quantization (**UNQ**) [18] and
 194 **Spreading Vectors** [22]. But ours has several key
 195 advantages compared to previous works: **1)** Unlike
 196 UNQ, which reconstructs features by an encoder-
 197 decoder structure, we follow the paradigm of MCQ
 198 to directly give binary codes and codebooks for the
 199 benefit of speed and storage, for UNQ needs an
 200 extra decoding stage during retrieval. **2)** UNQ and
 201 Spreading Vectors both project original features
 202 into a learned space. Although similarities between
 203 features are preserved, they still have biases in quan-
 204 tized results. This causes several issues, especially
 205 when we want to perform a reconstruction to ap-
 206 proximate original features, *e.g.* data compression.
 207

208 Compared to LSQ [14], the state-of-the-art heur-
 209 istic algorithm, our work is the first to tackle MCQ in
 210 a deep learning fashion. The policy network is de-
 211 signed to be very simple to get fast encoding speed
 212 and comparable retrieval performance.

213 5 Experiments

214 Our proposed Deep Multi-Codebook Quantization (**DeepQ**) is compared against the state-of-the-arts
 215 on a visual-feature dataset (**LabelMe22K**) to evaluate retrieval performance and encoding speed.
 216 Then, we scale up to make comparisons on commonly used large-scale datasets (**SIFT1M** and
 217 **DEEP1M**), whose base sets include 1 million vectors for retrieval. Furthermore, ablation study on
 218 **SIFT1M** investigates the effectiveness of each component in our proposed pipeline.

219 5.1 Datasets and evaluation metrics

220 **LabelMe22K** [29]: This dataset collects images by the LabelMe annotation tool¹ and uses Convolu-
 221 tional Neural Network (CNN) to extract them into 512-d features. It has 22, 019 vectors for training
 222 and 2, 000 vectors for test.

223 **SIFT1M**² and **DEEP1M**³: Both datasets contain 10^4 , 10^5 , 10^6 vectors in query, training and base
 224 set, respectively. Vectors from SIFT1M is extracted by Scale-Invariant Feature Transform (128-d)
 225 while DEEP1M contains 96-d vectors from outputs of a CNN.

¹<https://github.com/CSAILVision/LabelMeAnnotationTool>

²<http://corpus-texmex.irisa.fr/>

³<http://sites.skoltech.ru/compvision/noimi/>

Algorithm 1: VC-PPO for Training

Inputs: Training set \mathbf{X} , max step T ,
 hyper-
 parameters α, ϵ , learning rates η_1, η_2 .
Outputs: Policy π .
 Initialize codebook \mathbf{C} , parameters θ and τ ;
 $i \leftarrow 0$;
while $i < T$ **do**
/ Training loop */*
for \mathbf{x} in \mathbf{X} **do**
/ Sampling stage */*
 Sample $\mathbf{b}_o \sim \pi(\mathbf{x} | \theta_o)$ into \mathbf{B} ;
 Compute v_o, p_o into \mathbf{V}, \mathbf{Pr} ;
end
for $\mathbf{x}, \mathbf{b}_o, v_o, p_o$ in $\mathbf{X}, \mathbf{B}, \mathbf{V}, \mathbf{Pr}$ **do**
/ Updating stage */*
 $\tau \leftarrow \tau - \eta_1 \nabla_{\tau} \mathcal{L}_{\tau}$;
 Compute \hat{A} by (8);
 $\theta \leftarrow \theta + \eta_2 \nabla_{\theta} (\mathcal{L}_{\theta} + \alpha \cdot e_{\theta})$;
end
 $\mathbf{C} \leftarrow$ Solved by (15) ~ (18);
 $i \leftarrow i + 1$;
end
return $\pi(\cdot | \theta)$

226 Recall@{1, 10, 100} and quantization error are adopted as evaluation metrics. These two metrics
 227 indicate not only the retrieval performance but also the reconstruction accuracy. Because LabelMe22K
 228 does not have a base set, its training set is adopted as a base set. We train on the training set, and then
 229 encode the base set for evaluations with queries. When calculating recall, groundtruth is defined as
 230 the nearest neighbor of each query in the base set (sorted by l_2 distance). As for quantization error,
 231 the average value of $Q(x, b, C)$ is reported over all x in the base set.

232 We compare our proposal with both shallow and deep methods, including three classic quantiza-
 233 tion: **OPQ** [6], **SQ** [15] and **LSQ++** [14, 16] (denoted as **LSQ** for simplicity. Also, these two in
 234 our experiments have similar performance), as well as three gradient-based methods: **DPQ** [11],
 235 **DPgQ** [5] and **DRQ** [28]. **DPQ** and **PQNet** [33] have basically the same architecture that extend
 236 **PQ** with gradient-descent, so we only report the performance of **DPQ**. Additionally, **UNQ** [18] is
 237 also included, although they introduce an extra decoder and re-ranking trick for retrieval.

238 5.2 Implementation details

239 Our method is implemented with Py-
 240 Torch,⁴ the popular deep learning
 241 package in Python. Codebook C is
 242 solved by Intel MKL that has been
 243 fully optimized for speed. As for net-
 244 work training, we adopt Adam opti-
 245 mizer with AMSGrad [21] and hyper-
 246 parameters are tuned by grid search.
 247 Specifically, learning rates $\eta_1 = \eta_2 =$
 248 2×10^{-4} , with an exponential learning
 249 rate decay $\gamma = 0.9999$. Batch-size
 250 in updating stage is 2000, while other
 251 hyper-parameters $\epsilon = 0.2$, $\alpha = 0.05$.
 252 Additionally, during training, we in-
 253 sert dropout layers after every layer-
 254 normalization in all layer-groups to
 255 tackle overfitting. More detailed set-
 256 tings as well as specifications of *In-*
 257 *depNet* θ and *QENet* τ on each dataset
 258 (LabelMe22K, SIFT1M, DEEP1M) can be found in supplementary material.

259 As for quantization code-lengths, $K = 256$ codewords for each sub-codebook and $M = \{4, 8\}$
 260 sub-codebooks are employed in total. We follow [2] to report “effective” code-lengths (additional
 261 code-length for storing $\|x\|$ for lookup table is ignored). Therefore code-lengths become $\{32, 64\}$
 262 bits, respectively.

263 For a fair comparison, experiments are conducted on a single machine, equipped with Intel Xeon
 264 E5-2678v3 CPU, 256 GiB RAM, and NVIDIA RTX 3090 GPU. For other methods, we re-run on all
 265 datasets under unified settings with implementations provided by the authors.

266 5.3 Comparisons with state-of-the-arts

267 Under the small training set and base set settings on LabelMe22K, we get the results placed in
 268 Table 1. Our method takes the highest recall on this dataset, outperforming the state-of-the-art
 269 by 2.20%, 7.85%, 3.30% on 32 bits for R@1, R@10 and R@100. It also outperforms the best
 270 competitor by 2.70%, 1.45%, 0.65% on 64 bits. In brief, All methods except for UNQ are generally
 271 split into three styles: **1) PQ-like**: OPQ and DPQ. **2) SQ-like**: SQ, DPgQ and DRQ. **3) MCQ**: LSQ
 272 and ours. Generally, DPQ, DPgQ, and DRQ achieve similar results compared to their shallow
 273 versions. However, since they are still constrained MCQs, they show worse performances than 3).
 274 The performance of LSQ is worse than ours, shows the effectiveness of neural networks for modeling
 275 the MCQ encoding problem. As for UNQ, it takes several extra tricks *i.e.*, another network for
 276 decoding and re-ranking in retrieval. Although it beats LSQ, our network still shows the power of
 277 MCQ to win the competition.

Method	LabelMe22K					
	32 bits			64 bits		
	R@1	R@10	R@100	R@1	R@10	R@100
OPQ	18.70	57.25	90.10	32.30	80.40	98.00
SQ	18.45	57.60	90.85	32.65	82.05	99.05
LSQ	21.20	60.85	<u>94.35</u>	36.45	<u>86.25</u>	<u>99.15</u>
DPQ	8.60	32.80	77.50	15.35	48.75	90.75
DPgQ	19.85	57.80	90.70	35.05	84.10	98.90
DRQ	9.65	34.15	80.15	30.75	77.35	97.10
UNQ	<u>22.25</u>	<u>61.20</u>	89.30	<u>37.10</u>	85.55	98.80
Ours	24.45	69.05	97.65	39.60	87.60	99.80

Table 1: Recall(R)@{1, 10, 100} on **LabelMe22K** dataset (%). Ours outperforms state-of-the-arts by at least 2.20%, 7.85%, 3.30% (32 bits), and 2.70%, 1.45%, 0.65% (64 bits), respectively.

⁴<https://pytorch.org/>

Method	SIFT1M						DEEP1M					
	32 bits			64 bits			32 bits			64 bits		
	R@1	R@10	R@100	R@1	R@10	R@100	R@1	R@10	R@100	R@1	R@10	R@100
OPQ	5.34	22.03	56.72	22.84	60.27	92.19	3.07	15.39	48.40	15.34	50.06	87.96
SQ	9.45	34.88	70.07	24.41	65.48	93.17	6.41	26.79	70.25	19.95	56.31	91.27
LSQ	11.43	40.48	80.52	33.23	78.37	98.72	<u>7.29</u>	<u>28.96</u>	72.93	21.12	<u>61.47</u>	<u>93.98</u>
DPQ	5.41	22.97	58.57	21.87	59.39	91.66	1.59	8.96	33.09	9.53	33.45	72.80
DPgQ	9.71	35.03	74.19	27.96	69.98	96.04	6.36	26.16	70.02	18.98	55.80	90.95
DRQ	1.40	8.87	35.27	18.56	53.06	88.45	4.48	22.46	62.57	16.10	52.76	89.31
UNQ	10.01	33.92	73.39	<u>28.37</u>	69.15	95.99	5.19	23.55	65.09	16.12	52.06	90.10
Ours	<u>11.02</u>	<u>37.73</u>	<u>76.79</u>	28.02	<u>70.22</u>	<u>96.43</u>	7.43	30.03	<u>72.48</u>	<u>20.87</u>	62.06	94.07

Table 2: Quantitative comparisons with state-of-the-arts on **SIFT1M** and **DEEP1M** datasets. Recall(R)@{1, 10, 100} are reported (%). Ours shows comparable performance with staet-of-the-arts on SIFT1M, while achieving the highest recall in most cases on DEEP1M.

278 5.3.1 Large-scale retrieval performance

279 Our evaluations on SIFT1M and DEEP1M datasets is presented in Table 2. The training set and
 280 base set are scaled up, and retrievals on these datasets become more difficult. We observe expected
 281 results on two datasets. Compared to our main competitor, LSQ, our method achieves comparable
 282 performance on SIFT1M, and outperforms LSQ on DEEP1M in most cases. Our method achieves
 283 higher recall on DEEP1M than SIFT1M. A potential reason is that DEEP1M is under a nearly normal
 284 distribution that, in practice, is easier to converge than SIFT1M, which has a larger variance between
 285 datapoints. The performance of UNQ in our experiments is lower than expected, possibly due to
 286 different dataset settings.

287 Another key advantage of our method is that, different from shallow methods, which are hand-crafted
 288 algorithms that find possible solutions manually or with constraints, our DeepQ encodes vectors by
 289 only a feed-forward.

290 5.3.2 Encoding efficiency

291 In order to verify the encoding efficiency of
 292 our method, evaluations of encoding time on
 293 SIFT1M with the 10^6 base set are conducted
 294 by checking the total time spent. All of them
 295 are run under GPU-acceleration. Additionally,
 296 we evaluate the time with and without the extra
 297 codewords refinement that introduced in section
 298 4.3 (128 bits results are simulated). As Figure
 299 3 shows, our network is significantly faster than
 300 LSQ since it needs to perform local search it-
 301 eratively for 25 or even 100 rounds. Specifi-
 302 cally, to encode SIFT1M base set, LSQ takes
 303 52.84s, 96.99s, 256.86s and 639.18s for 16, 32,
 304 64 and 128 bits respectively. By contrast, our
 305 method takes 4.46s, 5.46s, 8.26s and 16.64s,
 306 which is $11.8\times$, $17.8\times$, $31.1\times$ and $38.4\times$
 307 faster than LSQ. Moreover, our method is even faster
 308 than most of the constrained MCQs. We also
 309 notice that the refinement takes negligible over-
 310 head. Although UNQ takes the fastest encoding
 311 speed, it still needs to decode and re-rank during
 312 retrieval, which slows down its retrieval speed.

313 5.3.3 Reconstruction accuracy

314 The comparisons of quantization error on three
 315 datasets are stated in Table 3. Basically, when the quantization error gets lower, recall will be higher.

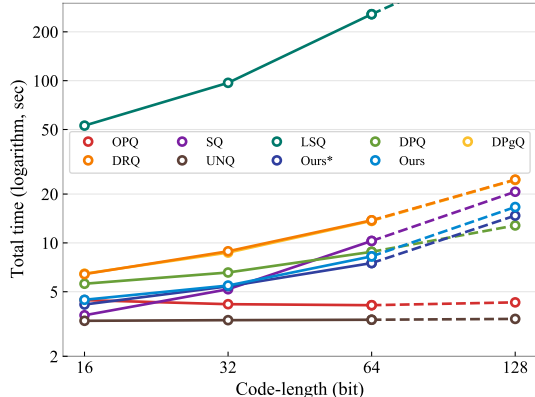


Figure 3: Total encoding time w.r.t. code-length on SIFT1M dataset. For 128 bits, we illustrate the simulated results. The variant *Ours** removes extra refinement step to show its overhead. Our two variants are significantly faster than LSQ while achieving similar performance. Furthermore, our method is slightly faster than most of the constrained MCQs. Our method achieves high performance as well as superior encoding efficiency. UNQ has the shortest time to encode the whole set, however during retrieval, they still need to decode and re-rank that slow down the speed.

Method	SIFT1M		DEEP1M		LabelMe22K	
	32 bits	64 bits	32 bits	64 bits	32 bits	64 bits
OPQ	4.03×10^4	2.51×10^4	4.25×10^{-1}	2.70×10^{-1}	1.25×10^{-1}	9.25×10^{-2}
SQ	3.42×10^4	2.13×10^4	3.24×10^{-1}	2.10×10^{-1}	1.25×10^{-1}	9.10×10^{-2}
LSQ	2.90×10^4	1.12×10^4	3.04×10^{-1}	1.99×10^{-1}	1.21×10^{-1}	8.57×10^{-2}
DPQ	4.01×10^4	2.48×10^4	4.58×10^{-1}	3.54×10^{-1}	1.77×10^{-1}	1.60×10^{-1}
DPgQ	3.30×10^4	2.10×10^4	3.29×10^{-1}	2.12×10^{-1}	1.31×10^{-1}	8.74×10^{-2}
DRQ	4.75×10^4	2.88×10^4	3.52×10^{-1}	2.54×10^{-1}	1.61×10^{-1}	1.01×10^{-1}
UNQ	4.14×10^4	2.33×10^4	3.52×10^{-1}	2.39×10^{-1}	1.48×10^{-1}	1.08×10^{-1}
Ours	<u>2.92×10^4</u>	<u>1.91×10^4</u>	2.92×10^{-1}	1.93×10^{-1}	1.02×10^{-1}	6.72×10^{-2}

Table 3: Comparisons of quantization error with state-of-the-arts on three datasets (*lower is better*). Ours achieves the lowest quantization error in most cases. This gives us benefits of feature reconstruction. Observe that UNQ performs poorly, we believe it focuses more on ranking and similarity preservation, other than reconstruction.

316 Ours get the 2nd place on SIFT1M, and the lowest on remaining datasets in most cases. Quantization
317 error indicates reconstruction accuracy and further shows the quality of codebook generation and
318 quantization codes selection. Notably, ours significantly outperforms UNQ, which has a strong bias
319 on the reconstruction task. This is because they focus more on ranking, not the quantization error.
320 The result shows that our method can be applied to other areas, *e.g.* vector compression.

321 5.4 Ablation study

322 Our ablation study is conducted on the
323 SIFT1M dataset, with the code-length of
324 32 bits, which in our experiments is suf-
325 ficient to show how does each component
326 affects our model. We choose the following
327 variants to perform ablation:

328 **w/o regularization:** which removes e_θ in
329 the losses, and the output distributions will
330 not be forced to be uniform.

331 **w/o return-norm:** which does not normal-
332 ize R , and therefor advantage is computed
333 by R other than \bar{R} .

334 **w/o correction:** which removes value cor-
335 rection. So our VC-PPO falls back to the original PPO.

336 **w/o refinement:** which directly encode the base set without extra refinement.

337 Quantization error and recall are evaluated and placed in Table 4. We report the best value they
338 ever met during the training procedure. Specifically, when regularization is removed, it seems that
339 the network is trapped in local-optima and the performance drops. Meanwhile, although return
340 normalization and value correction give us only subtle improvements, we find they help the network
341 to converge quickly. The extra refinement gives us lower quantization error and higher recall, specially
342 when we want to perform fast training before the network is converged.

343 6 Conclusion and Future Work

344 In this paper, we first review previous works of constrained MCQs, and investigate solutions to
345 unconstrained ones. Since finding the global-optima of MCQ is NP-hard, researchers apply constraints
346 to find near-optimal solutions or employ heuristic algorithms that are still time-consuming. This
347 paper takes the first attempt to find a *deep* solution to MCQ. The proposed IndepNet is designed to be
348 simple enough to encode vectors extremely fast. Furthermore, our network shows state-of-the-art
349 performance in retrieval and reconstruction tasks. Our method is slow to converge in a large dataset,
350 which hinders our performance. So, our future work will focus on training speedup.

Method	SIFT1M@32 bits			
	QE	R@1	R@10	R@100
w/o regularization	3.38×10^4	7.60	29.96	68.73
w/o return-norm	<u>3.06×10^4</u>	<u>10.57</u>	<u>36.44</u>	<u>76.04</u>
w/o correction	3.10×10^4	10.09	35.30	75.16
w/o refinement	3.17×10^4	9.91	30.39	68.28
DeepQ	2.92×10^4	11.02	37.73	76.79

Table 4: Ablation study conducted on SIFT1M with 32 bits code-length. Entropy regularization forces network to try more codeword combinations, which help to jump out of local-optima. Return normalization and value correction help for fast convergence. The extra refinement leads to low quantization error and high recall with negligible costs.

351 **References**

- 352 [1] Anderson, E., Bai, Z., Dongarra, J. J., Greenbaum, A., McKenney, A., Croz, J. D., Hammarling, S.,
353 Demmel, J., Bischof, C. H., and Sorensen, D. C. (1990). LAPACK: a portable linear algebra library for
354 high-performance computers. In *SC*, pages 2–11.
- 355 [2] Babenko, A. and Lempitsky, V. (2014). Additive quantization for extreme vector compression. In *CVPR*,
356 pages 931–938.
- 357 [3] Babenko, A. and Lempitsky, V. (2015). Tree quantization for large-scale similarity search and classification.
358 In *CVPR*, pages 4240–4248.
- 359 [4] Chen, Y., Guan, T., and Wang, C. (2010). Approximate nearest neighbor search by residual vector
360 quantization. *Sensors*, 10(12):11259–11273.
- 361 [5] Gao, L., Zhu, X., Song, J., Zhao, Z., and Shen, H. T. (2019). Beyond product quantization: Deep progressive
362 quantization for image retrieval. In *IJCAI*, pages 723–729.
- 363 [6] Ge, T., He, K., Ke, Q., and Sun, J. (2013). Optimized product quantization for approximate nearest neighbor
364 search. In *CVPR*, pages 2946–2953.
- 365 [7] Gray, R. (1984). Vector quantization. *IEEE Assp Magazine*, 1(2):4–29.
- 366 [8] Ilyas, A., Engstrom, L., Santurkar, S., Tsipras, D., Janoos, F., Rudolph, L., and Madry, A. (2020). A closer
367 look at deep policy gradients. In *ICLR*.
- 368 [9] Jégou, H., Douze, M., and Schmid, C. (2010). Product quantization for nearest neighbor search. *IEEE Trans.*
369 *Pattern Anal. Mach. Intell.*, 33(1):117–128.
- 370 [10] Kalantidis, Y. and Avrithis, Y. (2014). Locally optimized product quantization for approximate nearest
371 neighbor search. In *CVPR*, pages 2329–2336.
- 372 [11] Klein, B. and Wolf, L. (2019). End-to-end supervised product quantization for image search and retrieval.
373 In *CVPR*, pages 5041–5050.
- 374 [12] Konda, V. R. and Tsitsiklis, J. N. (2000). Actor-critic algorithms. In *NeurIPS*, pages 1008–1014.
- 375 [13] Lloyd, S. (1982). Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–
376 137.
- 377 [14] Martinez, J., Clement, J., Hoos, H. H., and Little, J. J. (2016). Revisiting additive quantization. In *ECCV*,
378 pages 137–153. Springer.
- 379 [15] Martinez, J., Hoos, H. H., and Little, J. J. (2014). Stacked quantizers for compositional vector compression.
380 *arXiv preprint arXiv:1411.2173*.
- 381 [16] Martinez, J., Zakhmi, S., Hoos, H. H., and Little, J. J. (2018). Lsq++: Lower running time and higher
382 recall in multi-codebook quantization. In *ECCV*, pages 491–506.
- 383 [17] Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A. (2020). Monte carlo gradient estimation in machine
384 learning. *J. Mach. Learn. Res.*, 21:132:1–132:62.
- 385 [18] Morozov, S. and Babenko, A. (2019). Unsupervised neural quantization for compressed-domain similarity
386 search. In *ICCV*, pages 3036–3045.
- 387 [19] Norouzi, M. and Fleet, D. J. (2013). Cartesian k-means. In *CVPR*, pages 3017–3024.
- 388 [20] Radenovic, F., Iscen, A., Tolias, G., Avrithis, Y., and Chum, O. (2018). Revisiting oxford and paris:
389 Large-scale image retrieval benchmarking. In *CVPR*, pages 5706–5715.
- 390 [21] Reddi, S. J., Kale, S., and Kumar, S. (2018). On the convergence of adam and beyond. In *ICLR*.
- 391 [22] Sablayrolles, A., Douze, M., Schmid, C., and Jégou, H. (2019). Spreading vectors for similarity search. In
392 *ICLR*.
- 393 [23] Schulman, J., Heess, N., Weber, T., and Abbeel, P. (2015a). Gradient estimation using stochastic computa-
394 tion graphs. In *NeurIPS*, pages 3528–3536.
- 395 [24] Schulman, J., Levine, S., Abbeel, P., Jordan, M. I., and Moritz, P. (2015b). Trust region policy optimization.
396 In *ICML*, pages 1889–1897.
- 397 [25] Schulman, J., Moritz, P., Levine, S., Jordan, M. I., and Abbeel, P. (2016). High-dimensional continuous
398 control using generalized advantage estimation. In *ICLR*.
- 399 [26] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization
400 algorithms. *arXiv preprint arXiv:1707.06347*.
- 401 [27] Song, J., Lang, R., Zhu, X., Xu, X., Gao, L., and Shen, H. T. (2020). 3d self-attention for unsupervised
402 video quantization. In *ACM SIGIR*, pages 1061–1070.
- 403 [28] Song, J., Zhu, X., Gao, L., Xu, X.-S., Liu, W., and Shen, H. T. (2019). Deep recurrent quantization for
404 generating sequential binary codes. In *IJCAI*, pages 912–918.

- 405 [29] Torralba, A., Fergus, R., and Weiss, Y. (2008). Small codes and large image databases for recognition. In
406 *CVPR*, pages 1–8. IEEE.
- 407 [30] Wang, J., Wang, J., Song, J., Xu, X., Shen, H. T., and Li, S. (2015). Optimized cartesian k-means. *IEEE*
408 *Transactions on Knowledge and Data Engineering*, 27(1):180–192.
- 409 [31] Weyand, T., Araujo, A., Cao, B., and Sim, J. (2020). Google landmarks dataset v2 - A large-scale
410 benchmark for instance-level recognition and retrieval. In *CVPR*, pages 2572–2581.
- 411 [32] Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement
412 learning. *Machine learning*, 8(3-4):229–256.
- 413 [33] Yu, T., Yuan, J., Fang, C., and Jin, H. (2018). Product quantization network for fast image retrieval. In
414 *ECCV*, pages 186–201.
- 415 [34] Zhang, T., Du, C., and Wang, J. (2014). Composite quantization for approximate nearest neighbor search.
416 In *ICML*, volume 2, page 3.

417 Checklist

418 The checklist follows the references. Please read the checklist guidelines carefully for information on
419 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
420 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
421 the appropriate section of your paper or providing a brief inline description. For example:

- 422 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 423 • Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- 424 • Did you include the license to the code and datasets? **[N/A]**

425 Please do not modify the questions and only use the provided macros for your answers. Note that the
426 Checklist section does not count towards the page limit. In your paper, please delete this instructions
427 block and only keep the Checklist section heading above along with the questions/answers below.

428 1. For all authors...

- 429 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
430 contributions and scope? **[Yes]**
- 431 (b) Did you describe the limitations of your work? **[Yes]** See Section 6.
- 432 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 433 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?
434 **[Yes]**

435 2. If you are including theoretical results...

- 436 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 437 (b) Did you include complete proofs of all theoretical results? **[N/A]**

438 3. If you ran experiments...

- 439 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
440 results (either in the supplemental material or as a URL)? **[Yes]** See [https://github.com/](https://github.com/DeepMCQ/DeepQ)
441 [DeepMCQ/DeepQ](https://github.com/DeepMCQ/DeepQ).
- 442 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
443 chosen)? **[Yes]** See Section 5 and supplementary materials.
- 444 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
445 multiple times)? **[No]**
- 446 (d) Did you include the total amount of compute and the type of resources used (e.g., type of
447 GPUs, internal cluster, or cloud provider)? **[Yes]** See Section 5.

448 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 449 (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Section 5.
- 450 (b) Did you mention the license of the assets? **[N/A]**
- 451 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** See
452 <https://github.com/DeepMCQ/DeepQ>.
- 453 (d) Did you discuss whether and how consent was obtained from people whose data you’re
454 using/curating? **[N/A]**

- 455 (e) Did you discuss whether the data you are using/curating contains personally identifiable
456 information or offensive content? [N/A]
- 457 5. If you used crowdsourcing or conducted research with human subjects...
- 458 (a) Did you include the full text of instructions given to participants and screenshots, if applica-
459 ble? [N/A]
- 460 (b) Did you describe any potential participant risks, with links to Institutional Review Board
461 (IRB) approvals, if applicable? [N/A]
- 462 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
463 participant compensation? [N/A]