

---

# Appendix: A Lower Bound of Hash Codes' Performance

---

**Xiaosu Zhu**<sup>1</sup>                      **Jingkuan Song**<sup>1\*</sup>                      **Yu Lei**<sup>1</sup>  
xiaosu.zhu@outlook.com      jingkuan.song@gmail.com      leiyu6969@gmail.com

**Lianli Gao**<sup>1</sup>                                      **Heng Tao Shen**<sup>1,2</sup>  
lianli.gao@uestc.edu.cn                      shenhengtao@hotmail.com

<sup>1</sup>Center for Future Media, University of Electronic Science and Technology of China

<sup>2</sup>Peng Cheng Laboratory

In this supplementary material, we discuss the following topics: Firstly, we give explanation in Appendix A to demystify concepts of rank lists. Then, proofs of the propositions in main paper are given in Appendix B. We further discuss why to adopt  $AP$  as a criterion of hash codes' performance in Appendix C. To train hash-models, we treat the posterior of hash codes to be under the multivariate Bernoulli distribution. We explain why and how to perform posterior estimation in Appendices D and E. Additional experiments are finally given in Appendix F.

## A Definitions

**The queries, true positives and false positives.** We demonstrate concepts of queries, true positives and false positives in rank lists in Figs. 1 and 2(a) for easy understanding. As the figure shows, any true positives or false positives are assigned with ranks  $i$ . Meanwhile, any true positives are also tagged by mis-ranks  $m$  introduced in this paper.  $m$  indicates how many false positives have the higher ranks than the current true positive.



Figure 1: An example rank list for demonstration. True positives are green while false positives are orange. All positives have ranks  $i$  placed on upper-right.

We assume distances between query and any positive samples are different with each other. Then, we obtain following inequalities according to the property of a rank list:

$$d(q, tp^1) < d(q, fp^2) < \dots < d(q, fp^8). \tag{1}$$

---

\*Corresponding author.

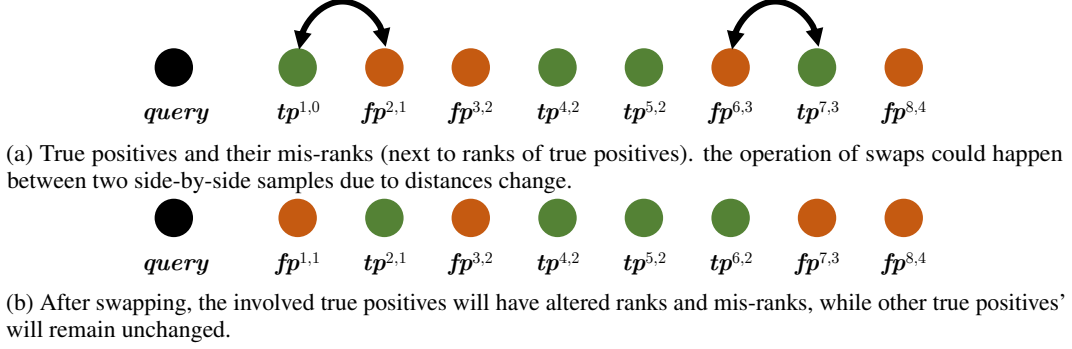


Figure 2: Mis-ranks marked on true positives and swaps that change ranks and mis-ranks.

**The side-by-side swaps.** Naturally, if a true positive and a false positive are placed side-by-side, and a swap happens between them due to the distances change, then mis-rank  $m$  as well as rank  $i$  of the true positive will be altered by 1 (Figs. 2(a) and 2(b)). Meanwhile, mis-ranks of any other true positives will remain unchanged. A special case is a side-by-side swap between two true positives or two false positives, where all ranks and mis-ranks would not change.

**Normal swaps.** More generally, any swaps happen in a rank list would influence ranks and mis-ranks of involved positive samples. To determine results after swaps, we could decompose the normal swaps into a series of side-by-side swaps with involved positive samples.

## B Proof

**Corollary.** Average precision increases iff  $m$  decreases.

*Proof.* From the above demonstration, if any swap happens in a rank list between true and false positives, the mis-rank of that true positive is definitely changed and will only result in increase or decrease of  $m$  and  $i$  by 1. Correspondingly, precision at rank of involved true positive changes:

$$\frac{i \pm 1 - (m \pm 1)}{i \pm 1} = \frac{i - m}{i \pm 1},$$

which is reversely proportional to  $m$ . Since precision of other true positives' are unchanged,  $AP$  is therefore influenced by the above true positive and reversely proportional to  $m$ . According to the derivation of normal swaps in Appendix A, this corollary will cover all cases of  $AP$  calculations.  $\square$

**Proposition.**

$$\bar{m} \propto \frac{\max d(\mathbf{q}, \mathbf{tp})}{\min d(\mathbf{q}, \mathbf{fp})} \forall \mathbf{tp} \in \mathbf{TP}, \mathbf{fp} \in \mathbf{FP}$$

where  $\bar{\cdot}$  denotes upper bounds. Correspondingly,

$$\underline{AP} \propto \frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})} \forall \mathbf{tp} \in \mathbf{TP}, \mathbf{fp} \in \mathbf{FP}$$

where  $\underline{\cdot}$  denotes lower bounds.

*Proof.* Considering any true positive, we are able to derive following inequalities from it according to Eq. (1):

$$\min d(\mathbf{q}, \mathbf{fp})_m < d(\mathbf{q}, \mathbf{tp}^m) < \min d(\mathbf{q}, \mathbf{fp})_{m+1} \quad (2)$$

where  $\min(\cdot)_m$  is the  $m$ -th minimum value among the whole set. For example, from Fig. 1, we could obtain:

$$\min d(\mathbf{q}, \mathbf{fp})_3 < d(\mathbf{q}, \mathbf{tp}^{7,3}) < \min d(\mathbf{q}, \mathbf{fp})_4 \quad (3)$$

where  $\min d(\mathbf{q}, \mathbf{fp})_3 = d(\mathbf{q}, \mathbf{fp}^6)$  and  $\min d(\mathbf{q}, \mathbf{fp})_4 = d(\mathbf{q}, \mathbf{fp}^8)$ .

Then, for the last true positive in rank list, which has the largest  $m$  among all true positives, *i.e.*  $\bar{m}$ , according to above inequalities, we could arrive:

$$\min d(\mathbf{q}, \mathbf{fp})_{\bar{m}} < \max d(\mathbf{q}, \mathbf{tp}) < \min d(\mathbf{q}, \mathbf{fp})_{\bar{m}+1}.$$

Now, if  $\max d(\mathbf{q}, \mathbf{tp})$  becomes smaller and  $\min d(\mathbf{q}, \mathbf{fp})_{\bar{m}}$  becomes larger so that a swap happens between the last true positive and its left most false positive, then we will have:

$$\begin{aligned} \bar{m}' &= \bar{m} - 1, \text{ iff} \\ \max d(\mathbf{q}, \mathbf{tp})' &< \max d(\mathbf{q}, \mathbf{tp}), \\ \min d(\mathbf{q}, \mathbf{fp})'_{\bar{m}} &> \min d(\mathbf{q}, \mathbf{fp})_{\bar{m}} \text{ and} \\ \min d(\mathbf{q}, \mathbf{tp})'_{\bar{m}} &> \max d(\mathbf{q}, \mathbf{tp})' \end{aligned}$$

where  $(\cdot)'$  indicates the new value. Combining our assumption in Appendix A as well as the pigeonhole principle, increase of  $\min d(\mathbf{q}, \mathbf{tp})$  results in increase of  $\min d(\mathbf{q}, \mathbf{fp})_{\bar{m}}$ . Therefore:

$$\left. \begin{array}{l} \max d(\mathbf{q}, \mathbf{tp}) \downarrow \\ \min d(\mathbf{q}, \mathbf{tp}) \uparrow \Rightarrow \min d(\mathbf{q}, \mathbf{fp})_{\bar{m}} \uparrow \end{array} \right\} \bar{m} \downarrow,$$

and vice versa. This proves Eq. (2). Eq. (3) is immediately obtained since  $AP$  is reversely proportional to  $\bar{m}$  according to above corollary.  $\square$

**Proposition.**

$$\underline{AP} \propto \frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})} \geq \text{const} \cdot \frac{\min \mathcal{D}_{inter}}{\max \mathcal{D}_{intra}}. \quad (4)$$

*Proof.* Without loss of generality,  $d(\mathbf{q}, \mathbf{tp})$  is covered by distances between any two codes of the same class  $c$  (denoted by  $\mathcal{D}^c$ ):  $d(\mathbf{q}, \mathbf{tp}) \subseteq \mathcal{D}^c = \{d(\mathbf{b}^i, \mathbf{b}^j), \text{ where } y^i = y^j = c\}$ . In contrast,  $d(\mathbf{q}, \mathbf{fp})$  is covered by:  $d(\mathbf{q}, \mathbf{fp}) \subseteq \mathcal{D}^{\neq c} = \{d(\mathbf{b}^i, \mathbf{b}^j), \text{ where } y^i = c, y^j \neq c\}$ . For the first case, we have:

$$\begin{aligned} d(\mathbf{b}^i, \mathbf{c}^c) &\leq \max \mathcal{D}_{intra}, d(\mathbf{b}^j, \mathbf{c}^c) \leq \max \mathcal{D}_{intra}, \\ &\text{for any } y^i = y^j = c \end{aligned}$$

where  $\mathbf{c}^c$  is the center of class  $c$ . Then, according to the triangle inequality:

$$d(\mathbf{b}^i, \mathbf{b}^j) \leq d(\mathbf{b}^i, \mathbf{c}^c) + d(\mathbf{b}^j, \mathbf{c}^c) \leq 2 \max \mathcal{D}_{intra}.$$

Therefore,

$$\max d(\mathbf{q}, \mathbf{tp}) \subseteq \mathcal{D}^c \leq 2 \max \mathcal{D}_{intra} = \text{const} \cdot \max \mathcal{D}_{intra}.$$

Similarly,  $\min d(\mathbf{q}, \mathbf{fp}) \geq \text{const} \cdot \min \mathcal{D}_{inter}$  could be derived. Combined with two inequalities, the above proposition is proved.  $\square$

## B.1 Analysis on the Proposed Lower Bound

### B.1.1 Is the Introduced Lower Bound Tight?

To determine whether the lower bound is tight is a little bit difficult. We firstly introduce some concepts and assumptions to make it easier. Let us start at the example placed in beginning of Appendix A.

**Asm. 1.** Any positive samples do not have the same distances to query. This ensures Eq. (1).

**Asm. 2.** Noticed that we are working in the Hamming space where the Hamming distances between any two samples are discrete and range from 0 to  $h$ , Eq. (1) becomes:

$$0 \leq d(\mathbf{q}, \mathbf{tp}^1) < d(\mathbf{q}, \mathbf{fp}^2) < \dots < d(\mathbf{q}, \mathbf{fp}^8) \leq h. \quad (5)$$

**Asm. 3.** The above array is strictly no gaps *i.e.* differences between any side-by-side  $d(\mathbf{q}, \cdot)$  are 1.

Then, we would derive the closed form lowest  $AP$  with  $\frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})}$  under the same order of magnitude. This conclusion could be intuitively drawn from the above example where  $tp^{7,3}$  and  $fp^2$  determine  $\max(\mathbf{q}, \mathbf{tp})$ ,  $\min(\mathbf{q}, \mathbf{fp})$ . If we keep two values unchanged (in other words, ranks of  $tp^{7,3}$  and  $fp^2$  unchanged), then the highest  $AP$  will appear as the following figure shows:

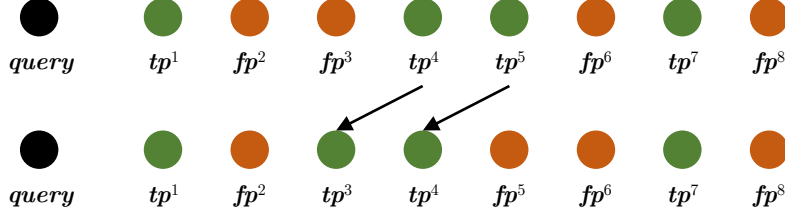


Figure 3: When  $\max(\mathbf{q}, \mathbf{tp})$ ,  $\min(\mathbf{q}, \mathbf{fp})$  are fixed, the highest  $AP$  will appear when all true positives have higher ranks than false positives in-between  $\max(\mathbf{q}, \mathbf{tp})$ ,  $\min(\mathbf{q}, \mathbf{fp})$ .

And the lowest  $AP$  will appear as the following figure shows:

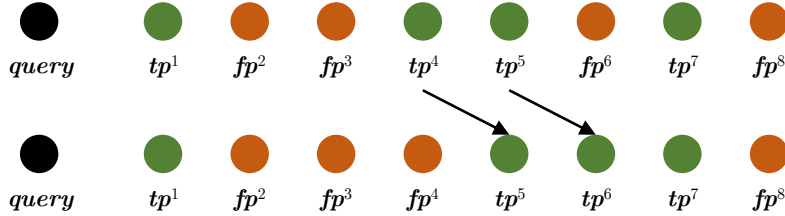


Figure 4: Correspondingly, the lowest  $AP$  will appear when all true positives have lower ranks than false positives in-between  $\max(\mathbf{q}, \mathbf{tp})$ ,  $\min(\mathbf{q}, \mathbf{fp})$ .

Based on this, we could easily derive that the lowest  $AP$  equals to

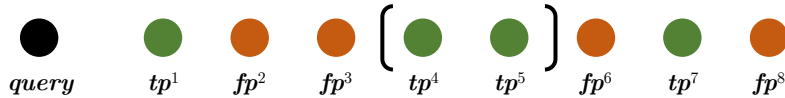
$$\min d(\mathbf{q}, \mathbf{fp}) - 1 + \sum_{i=1}^{|\mathbf{TP}|} \frac{i}{\max d(\mathbf{q}, \mathbf{tp}) - \min d(\mathbf{q}, \mathbf{fp}) + i}$$

which is proportional to  $\frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})}$ . Therefore, the lower bound is tight.

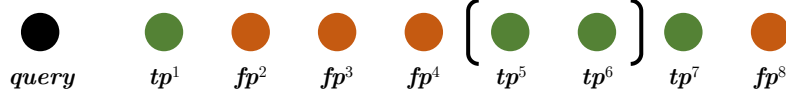
We could further extend  $\frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})}$  to  $\frac{\min \mathcal{D}_{inter}}{\max \mathcal{D}_{intra}}$  for the tight lower bound. From Eq. (4), we find that  $\frac{\min d(\mathbf{q}, \mathbf{fp})}{\max d(\mathbf{q}, \mathbf{tp})} \geq \frac{\min \mathcal{D}_{inter}}{\max \mathcal{D}_{intra}}$ . The equality is achieved when query's code  $\mathbf{q}$  is exactly the same as its class-center  $\mathbf{c}$ 's code. In this condition, the tight lower bound is derived by  $\frac{\min \mathcal{D}_{inter}}{\max \mathcal{D}_{intra}}$ .

Then, could it be applied to general cases? We give our humble discussion for a simple study. There may have untouched complicated cases that will be leaved for future study.

**Case 1: Duplicated positives.** If some samples are hashed to the same binary code (collision), then distances from query to them are all equal. They will appear at the same position of the rank list. If they are all true positives or false positives, then we could directly treat them as a single duplicated sample and follow the above rules. For example:



The above rank list has the duplicated true positives ( $d(\mathbf{q}, \mathbf{tp}^{4,2}) = d(\mathbf{q}, \mathbf{tp}^{5,2})$ ). If a swap happens between them and  $fp^6$ , the rank list will become:



where  $i, m$  of the duplicated true positives are both increased by 1. Obviously, the lower bound is still tight.

**Case 2: Mixed positives.** It is tricky when true positives and false positives have the same distance with query (we call them mixed positives). The sorting algorithm to produce final rank list also has impact to determine ranks of these mixed positives. Whether the lower bound is tight in this case is hard to judge, but our lower bound is still valid since above corollary still makes sense.

**Case 3: Rank list with gaps.** If there are gaps in between two distances, *e.g.*,  $d(q, fp^i) \ll d(q, tp^{i+1})$  (this usually happens on outliers), then  $AP$  will be far from the lower bound. Please refer to Appendix B.2 for details.

In conclusion, under the above assumptions Asm. 1, 2, 3, our lower bound is tight. Meanwhile, our lower bound also covers common cases as in the above discussion and makes a strong connection to  $AP$ .

## B.2 Edge Cases of the Lower Bound

Some edge cases when the lower bound is far way from the  $AP$  are further revealed according to the above analysis.

**Edge Case 1.** A huge amount of samples are hashed to the same binary code. Based on Case 1 and 2 in Appendix B.1.1, we will observe many duplicated or mixed samples. Their ranks will be increased / decreased simultaneously, and  $AP$  significantly changes along with them. The rank list is now be treated as “unstable”. Intuitively, when a hash-model has poor performance, it could not distinguish differences between samples and simply hashes them to the same code, and such a model will produce “unstable” rank lists.

**Edge Case 2.** Gaps in rank list. In the above example, if  $d(q, fp^6) \ll d(q, tp^{7,3}) = \max d(q, TP)$ ,  $\max d(q, TP)$  needs to be significantly decreased until a swap happens between  $fp^6$  and  $tp^{7,3}$  to influence  $AP$ . Therefore,  $AP$  is high but  $\frac{\min(q, fp)}{\max(q, tp)}$  is low. A potential reason leads to this edge case is outliers.

## C Criterion of Hash Codes’ Performance

Hashing techniques are widely applied in common multimedia tasks. To determine hash codes’ performance, the direct way is to adopt evaluation metrics used in these tasks *e.g.* retrieval precision, recall or recognition accuracy. In this section, we give the detailed explanation on how to adopt  $AP$  to cover the above typical metrics.

**Retrieval.** Common evaluation metrics in retrieval can be derived from  $AP$ . Specifically,

- Precision at rank  $i$  equals to  $\frac{i-m}{i}$ . The corollary in Appendix B exactly applies to it.
- Recall at rank  $i$  equals to  $\frac{i-m}{|T|}$  where  $T$  is set of all groundtruth samples. Therefore,  $|T|$  is a constant and recall increases *iff*  $m$  decreases.
- F-score equals to  $\frac{2}{recall^{-1} + precision^{-1}} = \frac{2}{i/(i-m) + |T|/(i-m)} = \frac{2(i-m)}{i+|T|}$ .

All the above metrics are reversely proportional to  $m$ . Therefore, analysis in Appendix B is also valid to them.

**Recognition.** There are many ways to adopt hash codes for recognition, and we discuss two main approaches here. 1)  $k$ NN based recognition is very similar to retrieval. In this scenario, prediction is obtained by voting among top- $k$  retrieved results. If there are more than half true positives in rank list, then the prediction will be correct, otherwise it will be wrong. Therefore, the recognition accuracy can be treated as a special case where we hold a rank list of all samples and expect  $m < k/2$  when  $i = k$ , which is induced to the goal of increasing  $AP$ .

As for logistic regression models for recognition, a learnable weight is adopted to make predictions:  $\mathbf{b} \xrightarrow{\mathbf{w}} \mathbf{cls}$ ,  $\mathbf{w} \in \mathbb{R}^{C \times h}$  where the output  $\mathbf{cls}$  is  $C$ -dim scores for each class. We could treat  $\mathbf{w}$  as  $C$  class-prototypes. Each prototype  $\mathbf{w}^i$  is a  $h$ -dim vector belongs to class  $i$  while  $\mathbf{cls}_i$  measures inner-product similarity between prototype and hash code  $\mathbf{b}$ . Inner-product similarity is an approximation of Hamming distance and  $\mathbf{w}$  is trained by cross-entropy objective, resulting in maximizing  $\mathbf{cls}_c$  while minimizing  $\mathbf{cls}_{other}$ . Such a goal is equivalent to maximizing intra-class compactness and inter-class distinctiveness.

### C.1 Potential Use Cases of the Lower Bound

Exploring use cases of the proposed lower bound would reveal significance and value of this work. Since Figure 1 and 5 in main paper and the above analysis tell us the value of  $\frac{\min \mathcal{D}_{inter}}{\max \mathcal{D}_{intra}}$  partially reflects hash codes' performance, we would quickly evaluate model's performance by such a metric other than calculating  $AP$  or accuracy which is time consuming. Therefore, adopting it as a performance indicator benefits for model tuning or selection, including but not limited to parameter search.

## D Multivariate Bernoulli Distribution

In this paper, we treat the posterior of hash codes  $p(\mathcal{B} | \mathcal{X})$  is under Multivariate Bernoulli distribution MVB. This is based on our observations of hash-model outputs. To demonstrate this claim, a toy 2 bits CSQ model is trained on 2 randomly chosen classes of CIFAR-10. We then plot  $\ell$  extracted by CSQ in Fig. 5.

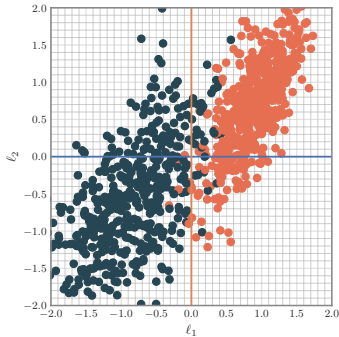


Figure 5:  $\ell$  extracted from a toy 2 bits CSQ model with 2 classes' samples of CIFAR-10. Correlations exist in  $\mathbf{b}_1$ ,  $\mathbf{b}_2$ .

	$p(\mathbf{b}_1)$	$\mathbf{b}_1 < 0$	$\mathbf{b}_1 > 0$
$p(\mathbf{b}_2)$		0.473	0.527
$\mathbf{b}_2 > 0$	0.525	0.079	0.446
$\mathbf{b}_2 < 0$	0.475	0.394	0.081

Figure 6: Estimated joint and marginal probabilities. Obviously, joint probability is not equal to product of marginal probability.

Obviously, features are concentrated on the upper-right and lower-left regions on the figure. This indicates that hash codes of these samples are tend to be (+1+1) or (-1-1) while very few samples are hashed to (+1-1) or (-1+1). Based on this pattern, we could infer that if  $\mathbf{b}_1$  and  $\mathbf{b}_2$  have a positive correlation. We also count frequencies of all 4 possible code combinations and calculate the estimated joint and marginal probability in Fig. 6, which confirms our observation. As the figure and the table show, joint probabilities are not equal to product of marginal probabilities, meaning that correlation exists among codes. Therefore, MVB is suitable to model the posterior of hash codes. Previous works are not sufficient to model it since they process each hash bit independently.

For instance, if  $h = 4$ , joint probability of  $\mathbf{b} = (+1+1-1-1)$  is formulated as:

$$p(\mathbf{b} = (+1+1-1-1)) = p(\mathbf{b}_1 < 0, \mathbf{b}_2 < 0, \mathbf{b}_3 > 0, \mathbf{b}_4 > 0). \quad (6)$$

Unfortunately, estimating this joint probability is difficult [2, 3, 6]. Furthermore, all  $2^h$  joint probabilities are required to be estimated in order to control hash codes precisely. Therefore, we propose our posterior estimation approach to tackle above difficulties.

## E Demonstration of Posterior Estimation

We continue to use the example in Eq. (6) to explain how we build our surrogate model to perform posterior estimation. Firstly, 4 bits hash codes have all 16 choices from (-1-1-1-1) to (+1+1+1+1). Therefore,  $\mathcal{o}$  has 16 entries where each entry in  $\mathcal{o}$  will be used to estimate a specific joint probability. We provide a simple demonstration in Fig. 7 to show how we process it.

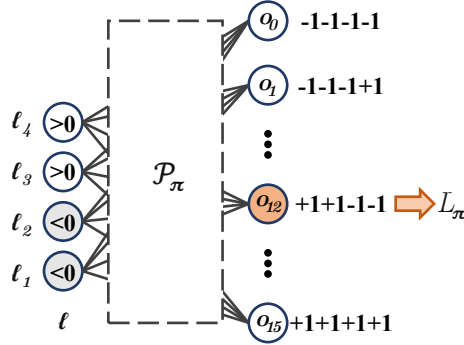


Figure 7: Demonstration of how  $\mathcal{P}_\pi$  performs posterior estimation.

Specifically, if we feed  $\ell$  where  $\ell_1 < 0$ ,  $\ell_2 < 0$ ,  $\ell_3 > 0$ ,  $\ell_4 > 0$  i.e.  $\mathbf{b} = (+1+1-1-1)$ , then we could directly find output  $o_{12} \hat{=} p(\mathbf{b}_1 < 0, \mathbf{b}_2 < 0, \mathbf{b}_3 > 0, \mathbf{b}_4 > 0)$  where  $12 = 1 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0$  and perform MLE on it to train model  $\pi$ .

### E.1 Implementation of Posterior Estimation

Correspondingly, we provide a minimal PyTorch-style implementation of  $\mathcal{P}_\pi$  and its training objective in Fig. 8. We first give one of a implementation of non-linear  $\ell$  to  $\mathcal{o}$  mapping in line № 1~8. Then, the main body of surrogate model (line № 11~55) shows how we process codes in blocks (line № 49) and how we convert every 8 bits inputs to 256 categorical outputs and calculate loss *w.r.t.* targets (line № 53).

## F Additional Experiments

We continue to conduct experiments for more comparisons and ablations. Before this, we firstly provide more implementation details as a supplement to main paper.

As Fig. 8 shows, we choose a normal two layer model with SiLU activation as  $\mathcal{P}_\pi$  to perform posterior estimation. To avoid over-fitting, we insert a dropout layer with 0.5 probability between layers. We train all methods for 100 epoches with batch-size 64, and learning rate is exponential decayed by 0.1 every 10 epoches. Most of the methods could achieve 95% performance after 30 epoches. We do not perform grid-search on hyper-parameters to obtain the highest results since it is not our topic in this paper.

To generate centers of specific classes for training (Alg.1 Line.2), we want a *perfect* generator which could produce finite number of codes with maximized pairwise Hamming distance among each other. Such a problem is formulated as the  $A_q(n, d)$  problem [4, 5] and such perfect codes are recognized to attain the *Hamming bound* [1, 7, 8]. However, this is extremely hard to tackle. Therefore, we conduct ablation study on a series of error-correction codes which aims at approaching the above goal to pick the best one, which will be explained below.

---

```

1  # One of the implementations for non-linear
2  # mapping to be used in surrogate model.
3  _nonLinearNet = lambda: nn.Sequential(
4      nn.Linear(8, 256),
5      nn.SiLU(),
6      nn.Dropout(0.5),
7      nn.Linear(256, 256)
8  )
9
10
11 class Surrogate(nn.Module):
12     """
13     Surrogate model to estimate MVB and further provide gradients.
14
15     Args:
16         bits (int): Length of hash codes.
17     """
18     def __init__(self, bits: int):
19         super().__init__()
20         # Number of blocks.
21         self._u = bits // 8
22         self._bits = bits
23         self._net = nn.ModuleList(
24             _nonLinearNet() for _ in range(self._u)
25         )
26         # binary to decimal multiplier
27         self.register_buffer("_multiplier",
28                               (2 ** torch.arange(8)).long())
29         self._initParameters()
30
31     def forward(self, x: Tensor, t: Tensor = None) -> (Tensor):
32         """
33         Module forward.
34         Args:
35             x (Tensor): [N, h] Model outputs before `.sign()`
36                 (un-hashed values).
37             t (Tensor, Optional): [N, h] Target hash result of x.
38         Return:
39             Tensor: [] If t is None, produces loss to
40                 train surrogate model. Otherwise,
41                 produces loss to calculate surrogate
42                 gradient w.r.t. x.
43         """
44         loss = list()
45         t = t or x
46         # split hash codes into `u` pieces, each 8 bits.
47         for subNet, org, tgt in zip(self._net,
48                                   torch.chunk(x, self._u, -1),
49                                   torch.chunk(t, self._u, -1)):
50             # feed each 8 bits codes into i-th non-linear model
51             # to get [256] predictions.
52             prd = subNet(org)
53             y = ((tgt > 0) * self._multiplier).sum(-1)
54             loss.append(F.cross_entropy(prd, y))
55         return sum(loss)

```

---

Figure 8: Minimal implementation of our proposed surrogate model for posterior estimation.

## F.1 Performance on Large Dataset

To evaluate performance of our proposed method on large-scale data, we conduct a new experiment on the whole ImageNet (ILSVRC 2012). Specifically, this dataset includes 1,000 classes with a relatively large training set (1.2M images,  $\sim 10\times$  larger than datasets adopted in main paper). The



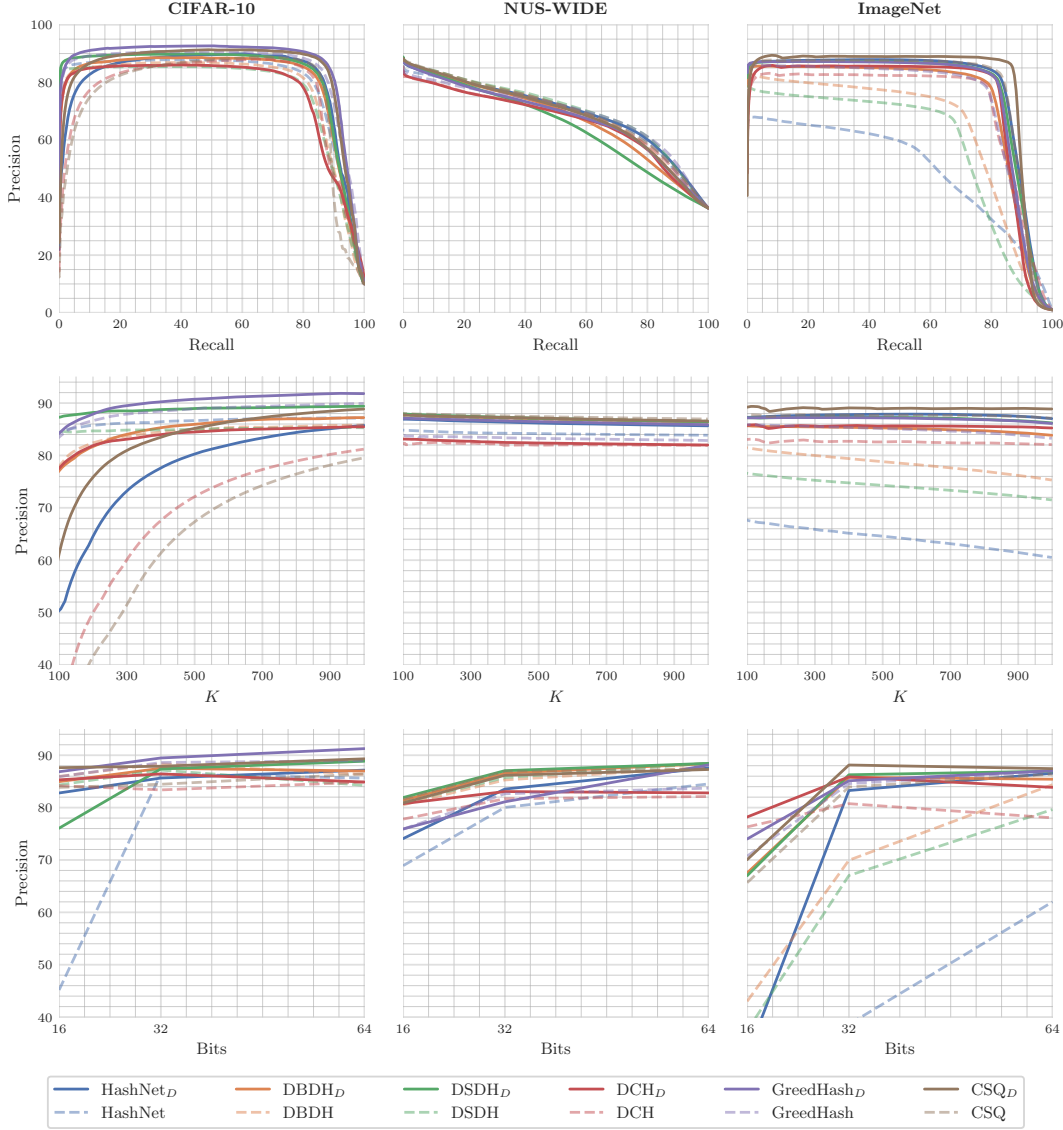


Figure 9: **First row:** Precision-Recall (P-R) curves for three datasets on 64 bits. With our method integrated, most of methods’ P-R curves slightly move to upper right. **Second row:** Precision@ $K$  ( $P@K$ ,  $K$  from 100 to 1,000) curves for three datasets on 64 bits. With our method integrated, most of methods’  $P@K$  curves slightly go up. **Third row:** Precision@ $H = 2$  (retrieval inside the  $H = 2$  Hamming ball) *w.r.t.* code-length curves on three datasets.

validation set has 50,000 images and 50 for each class for retrieval. We randomly pick 5 images of each class as queries (5,000 in total) and the remaining is adopted to formulate the base split (45,000 in total). Networks are trained for 20 epoch and we only update the last hash layers since the backbone is already pre-trained on ImageNet. Other settings are the same with main paper. From Tab. 1, we could see our method is also effective when trained with large dataset.

## F.2 Quantitative Comparisons

**Precision-Recall Curves.** To indicate retrieval performance over the whole rank list, we plot precision-recall curves for all methods in the first row of Fig. 9 for three datasets on 64 bits. We obtain similar results with Tab.2 where all the methods obtain performance gain (curves move to upper-right) after integrating our methods.

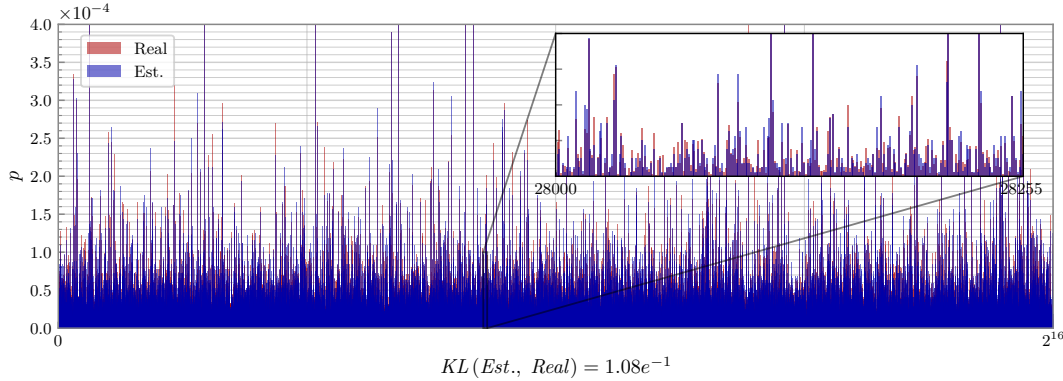


Figure 10: Posterior estimation on 16 bits codes (65, 536 joint probabilities) with a block=2 surrogate model. We randomly show a zoomed-in view with interval of 256 on the upper-right. Ours  $KL$ : 0.1084, naïve’s  $KL$ : 0.4965.

**Precision@K Curves.** To visualize retrieval precision of top retrieved samples, Precision@K ( $K$  from 100 to 1,000) curves for three datasets on 64 bits are placed in the second row of Fig. 9. Similar to P-R curves, most of the methods receive a precision increase and curves slightly move up.

**Precision@H = 2 Curves.** To show retrieval performance inside  $H = 2$  Hamming ball, we draw Precision@H = 2 *w.r.t.* code-length curves on the third row of Fig. 9. With our methods, precision inside Hamming ball is also increased. This indicates that our method pushes more true positives close to queries than the original methods.

### F.3 Ablation Study

#### F.3.1 Posterior Estimation on Longer Bits.

Due to space limitation in main paper, we place explanation on how we estimate MVB by naïve way here (Sec.6.2). Then, an extra experiments on longer bits *i.e.* 16 bits will be conducted to validate the effectiveness of blocked estimation.

For naïve estimation on MVB, dependence among  $\mathbf{b}$  is ignored. Specifically, we directly count frequency of value to be greater than zero on each bit of codes to approximate  $p(\mathbf{b}_i > 0)$  and  $p(\mathbf{b}_i < 0)$ . Then, any joint probability is estimated by the product of marginal probabilities, *e.g.*  $p(\mathbf{b} = +1+1-1-1) \hat{=} p(\mathbf{b}_4 > 0) \cdot p(\mathbf{b}_3 > 0) \cdot p(\mathbf{b}_2 < 0) \cdot p(\mathbf{b}_3 < 0)$ .

To train on 16 bits codes, we adopt a block=2 surrogate model where the first block is used to produce the first 8 bits codes and the second one is for the last 8 bits (as in Fig. 8 implements). Then, for all 65, 536 joint probabilities, we obtain them by the Cartesian product of two models’ predictions, which is shown in Fig. 10. The figure shows similar results as the 8 bits experiment in main paper. Our model achieves a much lower  $KL$  than naïve one, showing that performing block estimation will not introduce significant bias to estimate posterior (Ours: 0.1084, naïve: 0.4965).

Table 1: Retrieval performance on the full ImageNet dataset.

Method	16 bits	32 bits	64 bits
HashNet	36.9	41.2	44.7
HashNet <sub>D</sub>	<b>55.0</b> <sup>↑18.1</sup>	<b>56.1</b> <sup>↑14.9</sup>	<b>61.8</b> <sup>↑17.1</sup>
DBDH	37.4	41.1	49.1
DBDH <sub>D</sub>	<b>57.2</b> <sup>↑19.8</sup>	<b>57.9</b> <sup>↑16.8</sup>	<b>60.4</b> <sup>↑11.3</sup>
DSDH	39.3	48.4	54.2
DSDH <sub>D</sub>	<b>56.0</b> <sup>↑16.7</sup>	<b>59.1</b> <sup>↑10.7</sup>	<b>62.0</b> <sup>↑7.8</sup>
DCH	60.2	62.3	64.1
DCH <sub>D</sub>	<b>61.9</b> <sup>↑1.7</sup>	<b>62.9</b> <sup>↑0.6</sup>	<b>64.4</b> <sup>↑0.3</sup>
GreedHash	62.7	63.1	64.9
GreedHash <sub>D</sub>	<b>63.1</b> <sup>↑0.4</sup>	<b>64.0</b> <sup>↑0.9</sup>	<b>65.9</b> <sup>↑1.0</sup>
CSQ	64.6	65.0	65.7
CSQ <sub>D</sub>	<b>65.7</b> <sup>↑1.1</sup>	<b>66.2</b> <sup>↑1.2</sup>	<b>66.4</b> <sup>↑0.7</sup>

Table 2: mAP comparisons with different pre-defined centers on three benchmark datasets for 16, 32, 48 bits codes. Values on lower right is minimum pairwise Hamming distances among centers.

Center	CIFAR-10			NUS-WIDE			ImageNet		
	16 bits	32 bits	48 bits	16 bits	32 bits	48 bits	16 bits	32 bits	48 bits
Reed-Solomon	87.7 <sub>6</sub>	89.0 <sub>14</sub>	89.8 <sub>21</sub>	83.1 <sub>5</sub>	84.9 <sub>12</sub>	85.0 <sub>19</sub>	88.5 <sub>4</sub>	89.5 <sub>10</sub>	90.2 <sub>16</sub>
BCH	88.3 <sub>7</sub>	89.2 <sub>16</sub>	89.9 <sub>21</sub>	83.0 <sub>4</sub>	84.8 <sub>12</sub>	85.1 <sub>19</sub>	87.3 <sub>2</sub>	89.0 <sub>9</sub>	89.8 <sub>14</sub>
Hadamard	88.7 <sub>8</sub>	89.2 <sub>16</sub>	89.5 <sub>19</sub>	83.3 <sub>8</sub>	85.3 <sub>16</sub>	85.0 <sub>19</sub>	87.9 <sub>3</sub>	89.0 <sub>8</sub>	89.8 <sub>14</sub>

### F.3.2 Impact of Different Kinds of Pre-defined Centers.

To generate pre-defined centers as separate as possible as targets to train model, we compare a series of error-correction codes. The experiment includes Reed-Solomon code, BCH code and Hadamard code. Four kinds of codes have different pairwise Hamming distance. According to our main proposition in paper, if the minimum pairwise Hamming distance between centers are smaller, the lower bound of hash codes' performance is tend to be correspondingly lower. To show this phenomenon, minimum pairwise Hamming distance of codes as well as mAP are measured for  $CSQ_D$ , which is placed in Tab. 2. It is worth noting that Hadamard code is not applicable when number of centers is large or code-length is not order of 2, in this case we report results with randomly generated codes. The table confirms the positive correlation between minimum pairwise Hamming distance and performance. Meanwhile, performance differences among three kinds of codes are small. Therefore in practice, we could choose any of them.

## References

- [1] P. J. Cameron, J. A. Thas, and S. E. Payne. Polarities of generalized hexagons and perfect codes. *Geometriae Dedicata*, 5(4):525–528, 1976. 7
- [2] B. Dai. Multivariate bernoulli distribution models. Technical report, Citeseer, 2012. 7
- [3] B. Dai, S. Ding, and G. Wahba. Multivariate bernoulli distribution. *Bernoulli*, 19(4):1465–1483, 2013. 7
- [4] R. W. Hamming. *Coding and information theory*. Prentice-Hall, Inc., 1986. 7
- [5] F. J. MacWilliams and N. J. A. Sloane. *The theory of error correcting codes*, volume 16. Elsevier, 1977. 7
- [6] J. L. Teugels. Some representations of the multivariate bernoulli and binomial distributions. *Journal of multivariate analysis*, 32(2):256–268, 1990. 7
- [7] A. Tietäväinen. On the nonexistence of perfect codes over finite fields. *SIAM Journal on Applied Mathematics*, 24(1):88–96, 1973. 7
- [8] J. H. Van Lint. A survey of perfect codes. *The Rocky Mountain Journal of Mathematics*, 5(2): 199–224, 1975. 7