# DEEP HASHING USING ENTROPY REGULARISED PRODUCT QUANTISATION NETWORK

**Jo Schlemper, Jose Caballero, Andy Aitken, Joost van Amersfoort**
Cortex Applied Machine Learning, Twitter
`jo.schlemper11@imperial.ac.uk`, {`jcaballero,aaitken`}`@twitter.com`,
`joost.van.amersfoort@cs.ox.ac.uk`

## ABSTRACT

In large scale systems, approximate nearest neighbour search is a crucial algorithm to enable efficient data retrievals. Recently, deep learning-based hashing algorithms have been proposed as a promising paradigm to enable data dependent schemes. Often their efficacy is only demonstrated on data sets with fixed, limited numbers of classes. In practical scenarios, those labels are not always available or one requires a method that can handle a higher input variability, as well as a higher granularity. To fulfil those requirements, we look at more flexible similarity measures. In this work, we present a novel, flexible, end-to-end trainable network for large-scale data hashing. Our method works by transforming the data distribution to behave as a uniform distribution on a product of spheres. The transformed data is subsequently hashed to a binary form in a way that maximises entropy of the output, (i.e. to fully utilise the available bit-rate capacity) while maintaining the correctness (i.e. close items hash to the same key in the map). We show that the method outperforms baseline approaches such as locality-sensitive hashing and product quantisation in the limited capacity regime.

## 1 INTRODUCTION

In the modern era where we have an increasingly large amount of high-dimensional data to handle, it can be useful to have a system that can efficiently retrieve information that we care about. Examples of such systems are content-based image retieval (CBIR) (Datta et al., 2008; Babenko et al., 2014) and document/information retrieval (Mitra & Craswell, 2018). In large scale systems, linear search through the dataset is prohibitive. Therefore, one often resorts to approximate methods, which allow trading off accuracy for speed. These methods are commonly called approximate nearest neighbour (ANN) methods. Another important aspect of modern systems is data locality – ideally the data is stored a local fast disk, which restricts our representation to be quantised due to memory restrictions. Notable examples commonly used are locality-sensitive hashing (LSH) (Datar et al., 2004) and product quantisation (PQ) (Jegou et al., 2011; Ge et al., 2013).

Recently, deep learning has become an increasingly powerful tool for learning embeddings, due to the success of deep embedding learning (Oh Song et al., 2016; Hermans et al., 2017; Wu et al., 2017). These advances have motivated an approach called *deep hashing* (Wang et al., 2016; Erin Liong et al., 2015; Zhu et al., 2016), where one attempts to directly obtain a hash code from an image that can be used for content-based image retrieval tasks. These methods have been shown to greatly outperform traditional approaches. However, most methods rely on explicitly incorporating the class label prediction (as opposed to constructing an affinity matrix) to improve performance, which leads to the following issues. Firstly, while exploiting the class labels can improve the discriminative capability, it makes incorporating new labels a non-trivial task. Secondly, the methods do not directly account for semantic similarities at a granular level, making it unsuitable for certain tasks such as a duplication detection. Lastly, it is common to only demonstrate the efficacy of the methods for dataset with a small number of classes ($n \leq 100$), and the generalisation for the large scale dataset seems yet to be proven.

In this work, we propose a novel network architecture for end-to-end semantic hashing, which can be used for both deep hashing and learning an index structure (Kraska et al., 2018). Our network is

inspired by a *catalyser* network (Sablayrolles et al., 2018) and a supervised structured binary code (SUBIC) (Jain et al., 2017): it explores the idea of transforming an input distribution to a uniform distribution, but directly learns to generate the hash code. Our method is also flexible such that it relies on a similarity distance, which can be neighbour ranking or class labels. We show the applicability of our model for retrieval task using publicly available data set and we experimentally show our approach outperforms baseline methods such as LSH and PQ, in particular when the available bit-rate is limited.
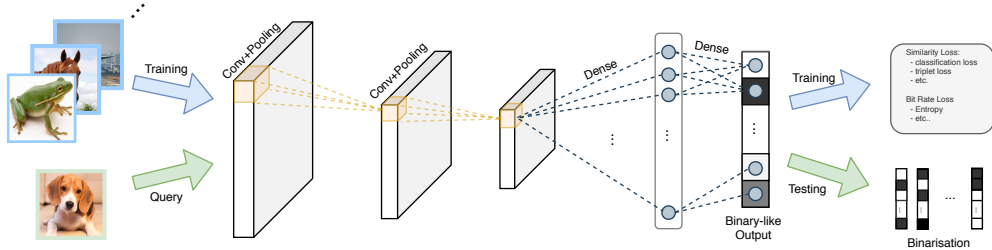
## 2 RELATED WORK



Figure 1: Schematics of Supervised Deep Hashing.

In the literature of hashing, common methods include: LSH, Iterative Quantisation (ITQ) (Gong et al., 2013) and PQ. While the first two aim to generate code in hamming space, PQ aims to represent data using a code book for a set of sub vectors. Recently deep learning-based hashing has become an active area of research (Lin et al., 2015; Liu et al., 2016). In particular, *supervised-hashing* is a research area that is concerned with hashing and retrieving objects (e.g. text or images) belonging to specific categories. In deep-learning based hashing, often three aspects are considered for the objective function, which are:

1. How to preserve semantic similarities of the inputs in their generated hash codes?

2. How to devise a continuous representation that can be trained using a neural network, which simultaneously minimises the discrepancy from test-time discretisation/binarisation?

3. How can we optimise the available bit-rate of a database (the output space), which can help minimise the collision probability?

The first challenge is often addressed by using a metric learning approach, such as contrastive loss (Hadsell et al., 2006), triplet loss (Schroff et al., 2015; Hermans et al., 2017) and their $n$-way extensions (Chen et al., 2017). The main drawback of these losses is that it is difficult to optimise them in a high dimensional space due to the curse of dimensionality (Friedman et al., 2001). Therefore, whenclass labels are available, it can be more effective to utilise those (Jain et al., 2017). Doing so induces a dependence on quality and availability of the class labels. In fact, as pointed out in Sablayrolles et al. (2017), retrieving objects based on classification puts an upper bound on the performance for the recall value. It can be more desirable to optimise the model on a more flexible objective which allows granular hashing based on semantic similarity, rather than solely the label information. The second challenge is a product of the non-differentiability of a naive discretisation step. Therefore at train time, one can resort to non-linearities such as $sigmoid$ and $tanh$, or continuous functions with better properties (Cao et al., 2017; 2018). The last point is usually handled by variants of entropy-based regularisation (Jain et al., 2017).

We also provide a more in-depth account of SUBIC (Jain et al., 2017) and catalyser network (Sablayrolles et al., 2018) as these networks form a foundation of our network architecture.

**SUBIC** Let $x \in \mathbb{R}^N$ be an input data (e.g. an image) and $y \in \mathcal{Y} = \{1, \ldots, C\}$ is the the class label. Given an image $z$, SUBIC outputs a structured binary code $b$, which is expressed as $M$ $K$-blocks: $b = [b^{(1)}; \ldots; b^{(M)}]$ with $b^{(i)} \in \mathcal{K}_K = \{d \in \{0,1\}^K \text{ s.t. } \|d\|_1 = 1\}$. This is achieved by the following network:

$$b = f_{\text{SUBIC}}(x) = \sigma \circ f_{\text{enc}} \circ f_{\text{feat}}(x) \tag{1}$$

where $f_{\text{feat}} : \mathbb{R}^N \to \mathbb{R}^{N'}$ is a feature extractor, $f_{\text{enc}} : \mathbb{R}^{N'} \to \mathbb{R}^{MK}$ is a hash encoder, $\sigma$ is a nonlinearity which applies $K$-softmax function to each of the $M$ block and "$\circ$" is a composition operator. During training, the $K$-blocks are relaxed into $K$-simplicies: $\widetilde{b} = [\widetilde{b}^{(1)}; \ldots; \widetilde{b}^{(M)}]$ with $\widetilde{b}^{(i)} \in \Delta_K = \{d \in [0,1]^K \text{ s.t. } \|d\|_1 = 1\}$. The novelty of SUBIC is to fit a classification layer $f_{\text{clf}} : \mathbb{R}^{MK} \to \mathbb{R}^C$ to learn a discriminative binary code. Given minibatch $\mathcal{B} = \{(x_i, y_i)\}_{i=1}^B$, the network is trained by minimising the following loss function:

$$\mathcal{L}(\mathcal{B}) = \underbrace{\frac{1}{B} \sum_{i=1}^B \log s[y_i]}_{\text{Cross Entropy}} + \gamma \underbrace{\frac{1}{B} \sum_{i=1}^B E(\widetilde{b}_i)}_{\text{Entropy}} - \mu \underbrace{E\left(\frac{1}{B} \sum_{i=1}^B \widetilde{b}_i\right)}_{\text{Batch Entropy}} \tag{2}$$

where $s = f_{\text{clf}}(\widetilde{b})$, $\widetilde{b} = f_{\text{SUBIC}}(x)$ and $E$ is mean entropy of $K$ blocks:

$$E(b) = -\frac{1}{M} \sum_{m=1}^M \sum_{k=1}^K b^{(m)}[k] \log_2 b^{(m)}[k] \tag{3}$$

The idea of the *Entropy* term is to encourage the network output to become one-hot like. On the other hand, the *Negative Batch Entropy* term encourages the uniform block support so that the available bit rate is fully exploited. At test time, $\sigma$ is replaced by a block-wise $\arg\max$ operation, where the binary code is obtained by setting the maximum activated entry in each block to 1, the rest to 0.

**Catalyser network** Let $x$ be the input, $z = f_{\text{cat}}(x)$ be the network embedding before quantisation is applied. Let $b = f_{\text{quant}}(z)$ be the quantised representation. The idea of the catalyser network is to embed data uniformly on an $\ell_2$-sphere, i.e. $z \in S^N$, which is subsequently encoded by an efficient lattice quantiser. The network is trained by minimising teh *triplet rank* loss (Hermans et al., 2017) and maximising the entropy loss. Given a triplet of input $(x_a, x_p, x_n)$ (anchor point, positive sample and negative sample respectively), the loss is defined as:

$$\mathcal{L}_{\text{tri}} = [\|f_{cat}(x_a) - f_{cat}(x_p)\|_2 - \|f_{cat}(x_a) - f_{cat}(x_n)\|_2 + \alpha]_+ \tag{4}$$

for margin $\alpha$. For entropy regularisation, Kozachenko and Leonenko (KoLeo) entropy estimator is used as a surrogate function:

$$H_n = \frac{\alpha}{n} \sum_i^n \log(\rho_{n,i}) + \beta, \quad \rho_{n,i} = \min_{i \neq j} \|f_{cat}(x_i) - f_{cat}(x_j)\| \tag{5}$$

The equation is simplified to:

$$\mathcal{L}_{KoLeo} = \sum_i^n \log(\rho_{n,i}) \tag{6}$$

The geometric idea is to ensure any two points are sufficiently far from each other, where the penalty decays logarithmically. The network then quantises the output using a Gosset Code, we refer the interested reader to Sablayrolles et al. (2018) for more details.
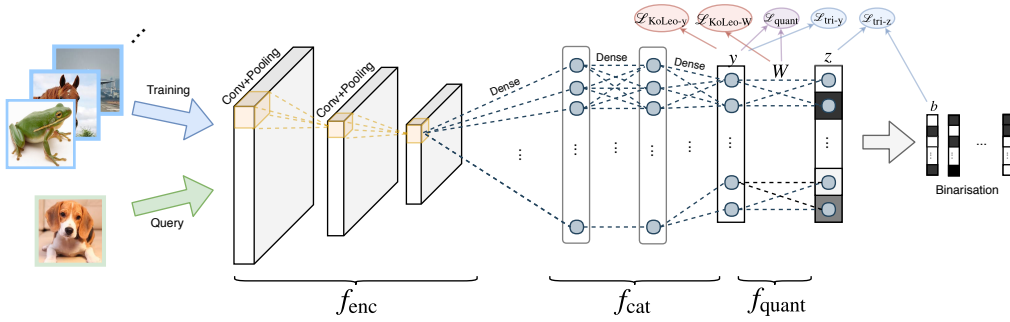
Figure 2: Schematic of the proposed architecture.

## 3 PROPOSED APPROACH

While SUBIC is effective, its application is limited to cases where classification labels are available. On the other hand, catalyser networks are more flexible, but not end-to-end trainable. In this work, we propose a flexible approach which incorporates the benefit of both and mitigates their limitations.

The proposed network is composed of three components: a feature extractor $f_{\text{enc}}$, a catalyser $f_{\text{cat}}$, and a quantiser $f_{\text{quant}}$, see Fig. 2. Given an input image $x$, the network directly generates a hash, which is a structured binary code as in SUBIC: $b \in [0,1]_K^M$. Let $y = f_{\text{cat}} \circ f_{\text{enc}}(x)$, $b = f_{\text{quant}}(y)$. The difference between catalyser networks and our work is that we learn the quantisation network $f_{\text{quant}}$, making the architecture end-to-end trainable. Our quantiser $f_{\text{quant}}(y)$ is given by $z^{(m)} = \sigma_k(Wy)$, where $\sigma_k$ is a block-wise K-softmax.

For simplicity, consider each $K$-block separately. Our key insight is the following: a fully connected layer is simply a dot product between $y$ and the row vectors of $W$. We have $z^{(m)} = \sigma_k(\langle w_1, y \rangle, \ldots, \langle w_K, y \rangle)$. Since at test-time, binarisation is done by selecting the maximally activated entry (within each $K$-block), this is equivalent to selecting the row vector with the smallest angular difference. This can be visualised using row vectors, $w_1, \ldots, w_k$ which linearly partition the output space, and the decision boundary is extending from the origin (Fig. 3).
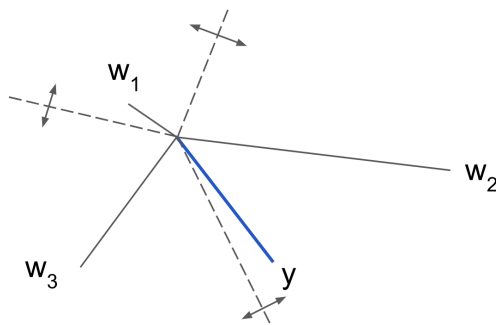


Figure 3: The geometry of the linear layer. The dashed line defines the assignment boundary for each $w_i$;s

Let $p(y)$ denote the probability distribution of the catalyser output taking a specific value in $y \in S^{K-1}$. To achieve the maximal entropy of $f_{\text{quant}}(y)$ (i.e. to maximise the used bit-rate), one can assign $y$ to each row vector $w_i$'s with an equal probability. Geometrically, this can be seen as equally partitioning the support of $p(y)$ by $w_i$'s, where $\text{supp}(p(y)) = \{y \in S^{K-1} | p(y = f_{cat} \circ f_{enc}(x)) > 0\}$. If $y$ is *uniformly* distributed on $S^{K-1}$, then it is sufficient to uniformly distribute $w_i$'s to partition the $S^{K-1}$ equally. This gives us the following strategy: we (1) encourage the distribution of $p(y)$

4

to be uniform on a sphere and (2) uniformly distribute $w_i$'s on the sphere. We can achieve both by using KoLeo entropy estimators in Eq. 6:

$$\mathcal{L}_{\text{KoLeo-y}} = \sum_i^n \log(\rho_{n,i}), \quad \rho_{n,i} = \min_{j \neq i} \|y_i - y_j\| \tag{7}$$

$$\mathcal{L}_{\text{KoLeo-W}} = \sum_i^n \log(\tau_{n,i}), \quad \tau_{n,i} = \min_{j \neq i} \|w_i - w_j\| \tag{8}$$

However, it is likely that a perfect uniform distribution cannot be achieved by the training, especially once combined with other deep embedding losses. To circumvent this, we add the following:

$$\mathcal{L}_{\text{quant}} = \sum_i \|y_i - w_{\iota(y_i)}\| \tag{9}$$

where $j = \iota(y_i) = \arg\min_{j'} \|y_i - w_{j'}\|$ (i.e. the index of the closest $w_j$ to $y$). By minimising Eq. 9, the row vectors $w_i$'s will be gravitated towards the probability mass of $p(y)$.

The remaining aspect is similar to the previous approaches: we minimise triplet loss to ensure similar points are embedded closely. For this, we can minimise triplet loss either in the output space of catalyser $\underbrace{S^{K-1} \times \cdots \times S^{K-1}}_{M}$ or the relaxed output space $\Delta_K^M$. Indeed, it is beneficial to directly optimise in the final target space. However, interestingly, it turns out that minimising triplet rank loss in simplex is difficult due to the fact that most points are very close to each other in high dimensional simplicies (see Appendix), yielding training instability. We mitigate this issue by using *asymmetrical* triplet loss:

$$\mathcal{L}_{\text{tri-z}} = [\|z_a - b_p\| - \|z_a - b_n\| + \alpha]_+ \tag{10}$$

where $z_a, z_p, z_n$ are the anchor, positive and negative points in the embedded space respectively, $\alpha$ is a margin, $b_*$'s are the discretised point (i.e. by replacing softmax by argmax). Note that sampling $b_*$'s is non-differentiable due to argmax, but we can nevertheless backpropagate the information using straight-through estimator (STE) proposed by Bengio et al. (2013), which has resemblance to stochastic graph computation approaches (Maddison et al., 2016). Secondly, the loss becomes zero if $b_p$ and $b_n$ share the same binary representation. We empirically found incorporating the triplet loss in $\underbrace{S^{K-1} \times \cdots \times S^{K-1}}_{M}$ was a useful additional loss to overcome this issue. The final objective is thus:

$$\mathcal{L} = \mathcal{L}_{\text{tri-z}} + \lambda_1 \mathcal{L}_{\text{tri-y}} + \lambda_2 \mathcal{L}_{\text{KoLeo-y}} + \lambda_3 \mathcal{L}_{\text{KoLeo-W}} + \lambda_4 \mathcal{L}_{\text{quant}} \tag{11}$$

where $\lambda_i$'s are hyper-parameter to be optimised. Note that for triplet loss, $\ell_2$-normalising $y$ and the rows of $W$ is important as otherwise arbitrary scaling can make the training unstable. Secondly, to reduce the parameters, we partition $W$ to only take each $K$-blocks and learn $W_1, \ldots, W_k$, where $W_i \in \mathbb{R}^{K \times K}$.[1] Finally, note that feature extractor and catalyser are only optimised with respect to $\mathcal{L}_{\text{tri-z}}, \mathcal{L}_{\text{tri-y}}$, and $\mathcal{L}_{\text{KoLeo-y}}$, whereas the quantiser weight $W$ is optimised only with respect to $\mathcal{L}_{\text{tri-z}}$, $\mathcal{L}_{\text{KoLeo-w}}$ and $\mathcal{L}_{\text{quant}}$. In particular, Eq. 9 is only minimised by $W$.

## 3.1 ENCODING AND DISTANCE COMPUTATION

Given a set of data points, we encode via $b = f_{\text{quant}} \circ f_{\text{cat}} \circ f_{\text{feat}}(x)$. The resulting vector can be compressed by storing indices of one-of-$K$ vectors, which only requires $M \log(K)$ bits. The

---

[1]Ideally, $M$ $K$-blocks are decorrelated to remove the redundancy. This is left as a future work.

distance between two compressed points can be given by Euclidean distance: $\|b_i - b_j\|_2$, which can be efficiently computed by $M$ *look up*: $\sum_i^M b_{\text{query}}^{(i)}[i(b^{(i)})]$, where $i(.)$ is the index of $b^{(i)}$ having one. One can also perform *asymmetric distance comparison* (ADC), which in case $b_{\text{query}}$ is replaced by $y_{\text{query}} = f_{\text{cat}} \circ f_{\text{feat}}(x)$, the data representation prior to quantisation.

## 3.2  NETWORK IMPLEMENTATION

For the feature encoder, a pre-trained network can be used, such as VGG or Resnet architectures. The catalyser was implemented using a fully connected network with 2 hidden layers, each having 256 features, and a final layer which maps the dimension to $d = MK$. We used batch-normalisation and Rectified Linear Unit (ReLU) for non-linearity, on all layers except the final one. The quantiser are $M$ separate fully connected layers with $K$ features. The overall network was trained using Adam with $\alpha = 10^{-3}, \beta_1 = 0.9, \beta_2 = 0.999$. The convergence speed of the network depends on the size of $K$, but usually sufficient performance can be obtained within 3 hours of training.
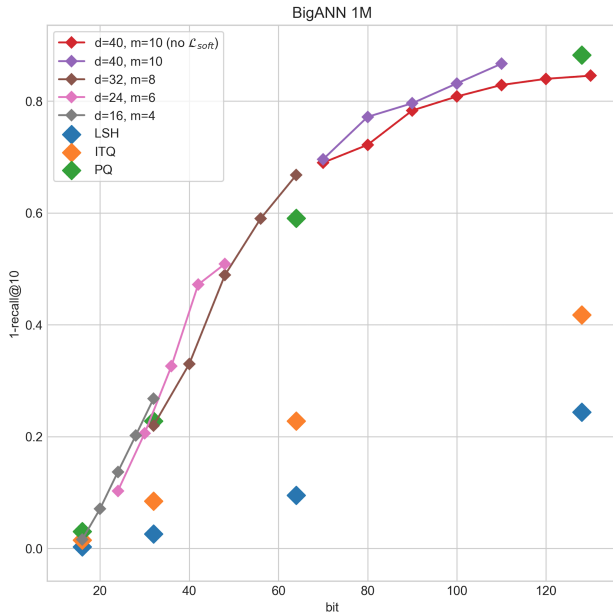


Figure 4: 1-Recall@10 for BigANN1M dataset.

## 4  EXPERIMENT

### 4.1  BIGANN1M

We evaluate our proposed approach using BigANN1M dataset[2]: the dataset contains a collection of 128 dimensional SIFT feature vectors. As the input is already feature vectors, we set $f_{\text{enc}} = \text{id}$. Training data contains 30,000 points, test data contains 10,000 query points and 1 million database points. For each point, we labelled the top $k = 10$ nearest points in terms of Euclidean distance to be the neighbours for triplet loss.

For evaluation, we used the metric 1-Recall@K=10, which measures the probability of retrieving the true first neighbour within the first 10 candidates. We compare to LSH, ITQ and PQ for the baseline methods. For PQ, we chose $K = 256$ for each sub-vector and varied the values of $M$ to achieve the desired bit-length $B \in \{16, 32, 64, 128\}$.

The result is summarised in Fig. 4. For the proposed method, we varied the number of $d$, $M$ and $K$ to get different number of bits. One can see that the performance of the proposed approach is

---

[2]Publicly available at http://corpus-texmex.irisa.fr/

comparable to PQ, but better for lower number of bits. Note that ITQ and LSH uses symmetric distance comparison so it is an unfair comparison. We also compared the proposed model with and without $\mathcal{L}_{quant}$ and we see a noticeable improvement. We speculate that this is because, while even without the loss, since the points are uniformly distributed it can achieve sufficient level of reconstruction, by minimising the quantisation loss, we remove the "gaps" in $\text{supp}(p(y))$.

### 4.1.1 VISUALISATION

We visualise the learnt weight vectors of the quantiser $f_{\text{quant}}$. For each $M$ sub-block, we randomly select 500 row vectors. Then we visualise 2 axes of these vectors (i.e. a projection onto 2 dimensional plane rather than using dimensionality reduction techniques). Without using $\mathcal{L}_{\text{quant}}$, the weights are uniformly distributed on $k$-sphere (Fig 5). However, when the loss is introduced, we see the mass of the rows concentrates on a more local area.
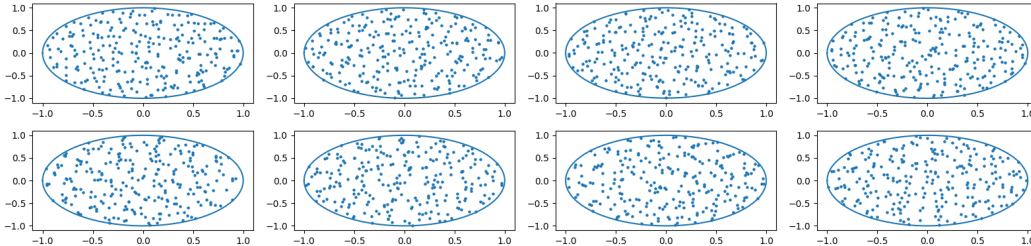


Figure 5: The distribution of the weights $w_i$ of the quantiser without $\mathcal{L}_{\text{quant}}$. Each sphere corresponds to one sub-block. For each sub-block, we randomly select 2 axes for visualization.
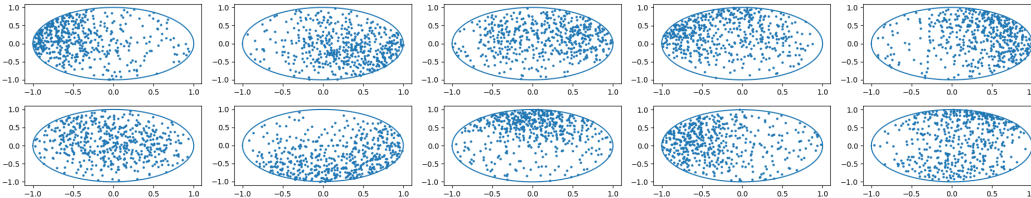


Figure 6: The distribution of weights $w_i$ of the quantiser with $\mathcal{L}_{\text{quant}}$. Each sphere corresponds to one sub-block. For each sub-block, we randomly select 2 axes for visualization.

## 5 CONCLUSION

In this work, we proposed a deep neural network which can perform end-to-end hashing of input, which only requires the knowledge of similarity graph, which is a slightly more relaxed constraint than class labels. The network operates by transforming the input space into a uniform distribution by penalising the cost given by KoLeo differential estimator, which was quantised by weight vectors uniformly distributed in its support. The network performs comparatively to the baseline methods, however, there is plenty of room for improvement. In the future, it will be interesting to impose a different prior on the distribution of the simplex, e.g. via Dirichlet distribution, to help control the output distribution, rather than relying on a uniform distribution.

REFERENCES

Artem Babenko, Anton Slesarev, Alexandr Chigorin, and Victor Lempitsky. Neural codes for image retrieval. In *European conference on computer vision*, pp. 584–599. Springer, 2014.

Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013. URL http://arxiv.org/abs/1308.3432.

Yue Cao, Mingsheng Long, Bin Liu, Jianmin Wang, and MOE KLiss. Deep cauchy hashing for hamming space retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1229–1237, 2018.

Zhangjie Cao, Mingsheng Long, Jianmin Wang, and S Yu Philip. Hashnet: Deep learning to hash by continuation. In *ICCV*, pp. 5609–5618, 2017.

Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, 2017.

Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*, pp. 253–262. ACM, 2004.

Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys (Csur)*, 40(2):5, 2008.

Venice Erin Liong, Jiwen Lu, Gang Wang, Pierre Moulin, and Jie Zhou. Deep hashing for compact binary codes learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2475–2483, 2015.

Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, NY, USA:, 2001.

Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. Optimized product quantization for approximate nearest neighbor search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2946–2953, 2013.

Yunchao Gong, Svetlana Lazebnik, Albert Gordo, and Florent Perronnin. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, 2013.

Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *null*, pp. 1735–1742. IEEE, 2006.

Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.

Himalaya Jain, Joaquin Zepeda, Patrick Pérez, and Rémi Gribonval. Subic: A supervised, structured binary code for image search. In *Proc. Int. Conf. Computer Vision*, volume 1, pp. 3, 2017.

Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2011.

Tim Kraska, Alex Beutel, Ed H Chi, Jeffrey Dean, and Neoklis Polyzotis. The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data*, pp. 489–504. ACM, 2018.

Kevin Lin, Huei-Fang Yang, Jen-Hao Hsiao, and Chu-Song Chen. Deep learning of binary hash codes for fast image retrieval. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 27–35, 2015.

Haomiao Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Deep supervised hashing for fast image retrieval. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2064–2072, 2016.

Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.

Bhaskar Mitra and Nick Craswell. An introduction to neural information retrieval. *Foundations and Trends® in Information Retrieval (to appear)*, 2018.

Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4004–4012, 2016.

Alexandre Sablayrolles, Matthijs Douze, Nicolas Usunier, and Hervé Jégou. How should we evaluate supervised hashing? In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pp. 1732–1736. IEEE, 2017.

Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, and Hervé Jégou. A neural network catalyzer for multi-dimensional similarity search. *arXiv preprint arXiv:1806.03198*, 2018.

Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.

Jun Wang, Wei Liu, Sanjiv Kumar, and Shih-Fu Chang. Learning to hash for indexing big dataa survey. *Proceedings of the IEEE*, 104(1):34–57, 2016.

Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krähenbühl. Sampling matters in deep embedding learning. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017.

Han Zhu, Mingsheng Long, Jianmin Wang, and Yue Cao. Deep hashing network for efficient similarity retrieval. In *AAAI*, pp. 2415–2421, 2016.

## 6 APPENDIX

### 6.1 DISTRIBUTION OF PAIRWISE DISTANCES ON SURFACES IN N DIMENSION

In the main manuscript, we argued that it is difficult to directly train triplet rank loss on high-dimensional simplex. Here we show how points on $n$-dimensional objects are distributed in high dimension as a part of the argument.
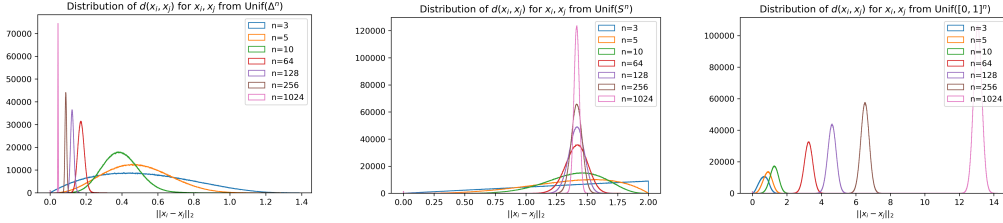


Figure 7: The distribution of $\ell_2$ distance between two random points on $n$-dimensional objects

#### 6.1.1 INTERIOR OF SIMPLEX

We use Dirichlet distribution with concentration parameter $a_1 = \cdots = a_n = 1$ to sample points uniformly in the interior of $n$-dimensional simplex. As one can see from Fig 8, as the dimension increases, the the points become more concentrated around the center of simplex $\mu = (1/n, \ldots, 1/n)$. The distribution of $\ell_2$ distance between two uniformly sampled points on $n$-simplex also sharply concentrates around small value, as it can be seen in 7. However, in the case of triplet rank loss, we would like to guarantee sufficiently high margin to ensure the separation between different classes. For example, the distance from any of the vertices to the centre of simplex is $\frac{\sqrt{n-1}}{\sqrt{n}}$ and the distance between two vertices is $\sqrt{2}$). We empirically saw that often the network collapses to predicting just $\mu$ and it is difficult to satisfy meaningful margin $\alpha$ as well as pushing the points approach towards one of the vertices.
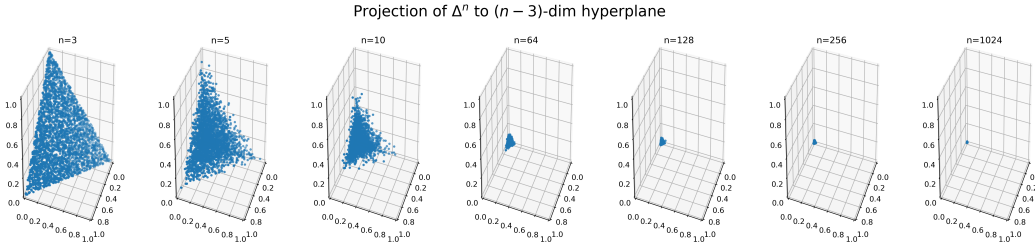


Figure 8: 3D projection of high-dimensional simplex

#### 6.1.2 SURFACE OF N-SPHERE

We sample points uniformly on $n$-sphere by first sampling $n$ from $\mathcal{N}(0,1)$, followed by $\ell_2$-normalisation. In this case, the distribution of distances between two uniformly sampled points on $n$-sphere is given by: $p(d) \propto d^{n-2}[1 - \frac{1}{4}d^2]^{\frac{n-3}{2}}$ (Wu et al., 2017). As $n \to \infty$, the probability distribution converges to $\mathcal{N}(\sqrt{2}, \frac{1}{2n})$. In this case, there is sufficient space left between majority of points, which is why we speculate that it is easier to train with triplet rank loss. Note that this however also means that since all points are already $\sqrt{2}$ far, careful *negative example mining* becomes very important to yield useful gradient.

#### 6.1.3 INTERIOR OF N-CUBE

We also study the distribution of distances between two random points in the interior of a hypercube. Here, the distances gets increasingly large as $n \to \infty$. Therefore, hypercube would have been an
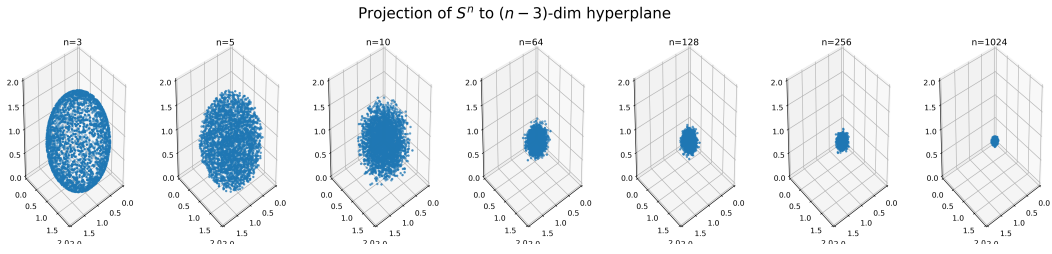
Projection of $S^n$ to $(n-3)$-dim hyperplane



Figure 9: 3D projection of high-dimensional $n$-sphere

alternative shape we could use as a domain for hashing, which could be interesting for future work. However, In this case, we could use sigmoid function to set the range, but this could result in gradient saturation.
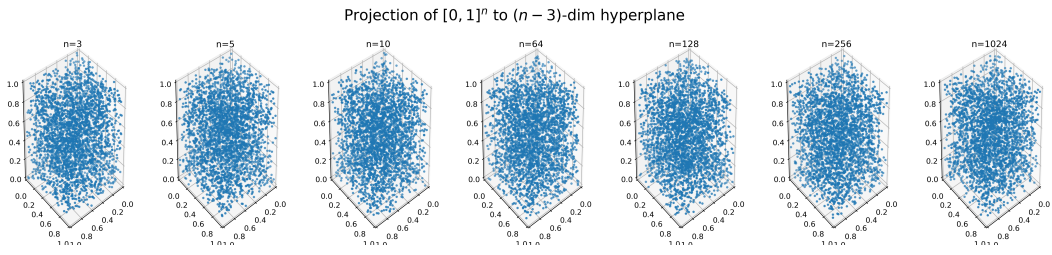
Projection of $[0,1]^n$ to $(n-3)$-dim hyperplane



Figure 10: 3D projection of high-dimensional $n$-cube