

# Unsupervised Neural Generative Semantic Hashing

Casper Hansen  
University of Copenhagen  
c.hansen@di.ku.dk

Christian Hansen  
University of Copenhagen  
chrh@di.ku.dk

Jakob Grue Simonsen  
University of Copenhagen  
simonsen@di.ku.dk

Stephen Alstrup  
University of Copenhagen  
s.alstrup@di.ku.dk

Christina Lioma  
University of Copenhagen  
c.lioma@di.ku.dk

## ABSTRACT

Fast similarity search is a key component in large-scale information retrieval, where semantic hashing has become a popular strategy for representing documents as binary hash codes. Recent advances in this area have been obtained through neural network based models: generative models trained by learning to reconstruct the original documents. We present a novel unsupervised generative semantic hashing approach, *Ranking based Semantic Hashing* (RBSH) that consists of both a variational and a ranking based component. Similarly to variational autoencoders, the variational component is trained to reconstruct the original document conditioned on its generated hash code, and as in prior work, it only considers documents individually. The ranking component solves this limitation by incorporating inter-document similarity into the hash code generation, modelling document ranking through a hinge loss. To circumvent the need for labelled data to compute the hinge loss, we use a weak labeller and thus keep the approach fully unsupervised.

Extensive experimental evaluation on four publicly available datasets against traditional baselines and recent state-of-the-art methods for semantic hashing shows that RBSH significantly outperforms all other methods across all evaluated hash code lengths. In fact, RBSH hash codes are able to perform similarly to state-of-the-art hash codes while using 2-4x fewer bits.

## KEYWORDS

Unsupervised semantic hashing, Deep learning, Generative model, Document ranking

## ACM Reference Format:

Casper Hansen, Christian Hansen, Jakob Grue Simonsen, Stephen Alstrup, and Christina Lioma. 2019. Unsupervised Neural Generative Semantic Hashing. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '19)*, July 21–25, 2019, Paris, France. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3331184.3331255>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

SIGIR '19, July 21–25, 2019, Paris, France

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6172-9/19/07...\$15.00

<https://doi.org/10.1145/3331184.3331255>

## 1 INTRODUCTION

The task of similarity search consists of querying a potentially massive collection to find the content most similar to a query. In Information Retrieval (IR), fast and precise similarity search is a vital part of large-scale retrieval [28], and has applications in content-based retrieval [14], collaborative filtering [13], and plagiarism detection [10, 26]. Processing large-scale data requires solutions that are both computationally efficient and highly effective, and that work in an unsupervised fashion (because manually labelling massive datasets is unfeasible). Semantic hashing [21] is a highly effective class of methods that encode the semantics of a document into a binary vector called a *hash code*, with the property that similar documents have a short Hamming distance between their codes, which is simply the number of differing bits in the codes as efficiently computed by the sum of the XOR operation. For short hash codes of down to a single byte, this provides a very fast way of performing similarity searches [34], while also reducing the storage requirement compared to full text documents.

Originally, work on semantic hashing focused on generating hash codes for a fixed collection [25], but more modern information needs require querying *unseen* documents for retrieving similar documents in the collection. Modern semantic hashing methods are based on machine learning techniques that, once trained, are able to produce the hash code based solely on the document alone. This can be done using techniques similar to Latent Semantic Indexing [33], spectral clustering [29], or two-step approaches of first creating an optimal encoding and then training a classifier to predict this [34]. Recent work has focused on deep learning based methods [4, 5, 23] to create a generative document model. However, none of the methods directly model the end goal of providing an effective similarity search, i.e., being able to accurately rank documents based on their hash codes, but rather just focus solely on generating document representations.

We present a novel unsupervised generative semantic hashing approach, *Ranking based Semantic Hashing* (RBSH) that combines the ideas of a variational autoencoder, via a so-called variational component, together with a ranking component that aims at directly modelling document similarity through the generated hash codes. The objective of the variational component is to maximize the document likelihood of its generated hash code, which is intractable to compute directly, so a variational lower bound is maximized instead. The variational component is modelled via neural networks and learns to sample the hash code from a Bernoulli distribution, thus allowing end-to-end trainability by avoiding a post-processing

step of binarizing the codes. The ranking component aims at learning to rank documents correctly based on their hash codes, and uses weak supervision through an unsupervised document similarity function to obtain pseudo rankings of the original documents, which circumvents the problem of lacking ground truth data in the unsupervised setting. Both components are optimized jointly in a combined neural model, which is designed such that the final model can be used to generate hash codes solely based on a new unseen document, without computing any similarities to documents in the collection. Extensive experimental evaluation on four publicly available datasets against baselines and state-of-the-art methods for semantic hashing, shows that RBSH outperforms all other methods significantly. Similarly to related work [4, 5, 23], the evaluation is performed as a similarity search of the most similar documents via the Hamming distance and measured using precision across hash codes of 8-128 bits. In fact, RBSH outperforms other methods to such a degree, that generally RBSH hash codes perform similarly to state-of-the-art hash codes while using 2-4x less bits, which corresponds to an effective storage reduction of a factor 2-4x.

In summary, we **contribute** a novel generative semantic hashing method, *Ranking based Semantic Hashing* (RBSH), that through weak supervision directly aims to correctly rank generated hash codes, by modelling their relation to weakly labelled similarities between documents in the original space. Experimentally this is shown to significantly outperform all state-of-the-art methods, and most importantly to yield state-of-the-art performance using 2-4x fewer bits than existing methods.

## 2 RELATED WORK

### 2.1 Semantic Hashing

Semantic hashing functions provide a way to transform documents to a low dimensional representation consisting of a sequence of bits. These compact bit vectors are an integral part of fast large-scale similarity search in information retrieval [28], as they allow efficient nearest neighbour look-ups using the Hamming distance. Locality Sensitive Hashing (LSH) [6] is a widely known data-independent hashing function with theoretically founded performance guarantees. However, it is general purpose and as such not designed for semantic hashing, hence it empirically performs worse than a broad range of semantic hashing functions [4, 5]. In comparison to LSH, semantic hashing methods employ machine learning based techniques to learn a *data-dependent* hashing function, which has also been denoted as learning to hash [28].

Spectral Hashing (SpH) [29] can be viewed as an extension of spectral clustering [17], and preserves a global similarity structure between documents by creating balanced bit vectors with uncorrelated bits. Laplacian co-hashing (LCH) [33] can be viewed as a version of binarized Latent Semantic Indexing (LSI) [7, 21] that directly optimizes the Hamming space as opposed to the traditional optimization of Latent Semantic Indexing. Thus, LCH aims at preserving document semantics, just as LSI traditionally does for text representations. Self-Taught Hashing (STH) [34] has the objective of preserving the local similarities between samples found via a  $k$ -nearest neighbour search. This is done through computing the bit vectors by considering document connectivity, however without learning document features. Thus, the objective of preserving

local similarities contrasts the global similarity preservation of SpH. Interestingly, the aim of our RBSH can be considered as the junction of the aims of STH and SpH: the variational component of RBSH enables the learning of local structures, while the ranking component ensures that the hash codes incorporate both local and global structure. Variational Deep Semantic Hashing (VDSH) [5] is a generative model that aims to improve upon STH by incorporating document features by preserving the semantics of each document using a neural autoencoder architecture, but without considering the neighbourhood around each document. The final bit vector is created using the median method [29] for binarization, which means the model is not end-to-end trainable. Chaidaroon et al. [4] propose a generative model with a similar architecture to VDSH, but in contrast incorporate an average document of the neighbouring documents found via BM25 [20] which can be seen as a type of weak supervision. The model learns to also reconstruct the average neighbourhood document in addition to the original document, which has similarities with STH in the sense that they both aim to preserve local semantic similarities. In contrast, RBSH directly models document similarities based on a weakly supervised ranking through a hinge loss, thus enabling the optimization of both local and global structure. Chaidaroon et al. [4] also propose a model that combines the average neighbourhood documents with the original document when generating the hash code. However this model is very computationally expensive in practice as it requires to find the top- $k$  similar documents online at test time, while not outperforming their original model [4]. NASH [23] proposed an end-to-end trainable generative semantic hashing model that learns the final bit vector directly, without using a second step of binarizing the vectors once they have been generated. This binarization is discrete and thus not differentiable, so a straight-through estimator [2] is used when optimizing the model.

The related work described above has focused on unsupervised text hashing. Direct modelling of the hash code similarities as proposed in this paper has not been explored. For the case of *supervised* image hashing, some existing work has aimed at generating hash codes using ranking strategies from labelled data, e.g., based on linear hash functions [27] and convolutional neural networks [30, 36]. In contrast, our work develops a generative model and utilises weak supervision to circumvent the need for labelled data.

### 2.2 Weak Supervision

Weak supervision has showed strong results in the IR community [8, 9, 18, 32], by providing a solution for problems with small amounts of labelled data, but large amounts of unlabelled data. While none of these are applied in a problem domain similar to ours, they all show that increased performance can be achieved by utilizing weak labels. Zamani et al. [8] train a neural network end-to-end for ad-hoc retrieval. They empirically show that a neural model trained on weakly labelled data via BM25 is able to generalize and outperform BM25 itself. A similar approach is proposed by Nie et al. [18], who use a multi-level convolutional network architecture, allowing to better differentiate between the abstraction levels needed for different queries and documents. Zamani et al. [32] present a solution for the related problem of query performance

prediction, where multiple weak signals of clarity, commitment, and utility achieve state-of-the-art results.

### 3 RANKING BASED SEMANTIC HASHING

We first present an overview of our model, Ranking Based Semantic Hashing (RBSH), and then describe in detail the individual parts of the model. RBSH combines the principles of a variational autoencoder with a ranking component using weak supervision and is an unsupervised generative model. For document  $d$ , the variational component of RBSH learns a low dimensional binary vector representation  $z \in \{0, 1\}^m$ , called the hash code, where  $m$  is the number of bits in the code. RBSH learns an encoder and decoder function, modelled by neural networks, that are able to encode  $d$  to  $z$  and  $z$  to  $\hat{d}$ , respectively, where  $\hat{d}$  is an approximation of the original document  $d$ . The goal of the encoder-decoder architecture is to reconstruct the original document as well as possible via the hash code. Additionally, we incorporate a ranking component which aims to model the similarity between documents, such that the resulting hash codes are better suited for finding nearest neighbours. Specifically, during training RBSH takes document triplets,  $(d, d_1, d_2)$ , as inputs with estimated pairwise similarities, and through weak supervision attempts to correctly predict either  $d_1$  or  $d_2$  as being most similar to  $d$ . Training the model on inputs of various similarities (e.g., from the top 200 most similar documents) enables the model to learn both the local and global structure to be used in the hash code generation.

In summary, through the combination of the variational and ranking components the objective of RBSH is to be able to both reconstruct the original document as well as correctly rank the documents based on the produced hash codes. An overview of the model can be seen in Figure 1. In the sections below we describe the generative process of the variational component (Section 3.1), followed by the encoder function (Section 3.2), decoder function (Section 3.3), the ranking component (Section 3.4), and finally the combined model (Section 3.5).

#### 3.1 Variational component

We assume each document  $d$  to be represented as a bag-of-words representation of vocabulary size  $V$  such that  $d \in \mathbb{R}^V$ . We denote the set of unique words in document  $d$  as  $\mathcal{W}_d$ . For each document we sample a binary semantic vector  $z \sim p(z)$  where  $p(z_i) = p_i^{z_i} (1 - p_i)^{1-z_i}$ , which allows the hash codes to be end-to-end trainable, similarly to Shen et al. [23]. For each bit,  $p_i$  corresponds to the probability of sampling a 1 at position  $i$  and  $(1 - p_i)$  is the probability of sampling a 0. Thus,  $z$  is obtained by repeating a Bernoulli trial  $m$  times. Using the sampled semantic vector, we consider each word as  $w_i \sim p(w_i|f(z))$  and define the document likelihood as follows:

$$p(d|z) = \prod_{j \in \mathcal{W}_d} p(w_j|f(z)) \quad (1)$$

that is, a simple product of word probabilities where the product iterates over all unique words in document  $d$  (denoted  $\mathcal{W}_d$ ). In this setting  $f(z)$  is a function that maps the hash code,  $z$ , to a latent vector useful for modelling word probabilities.

**3.1.1 Variational loss.** The first objective of our model is to maximize the document log likelihood:

$$\log p(d) = \log \int_{\{0,1\}^m} p(d|z)p(z)dz \quad (2)$$

However, due to the non-linearity of computing  $p(w_j|f(z))$  from Equation 1 this computation is intractable and the variational lower bound [12] is maximized instead:

$$\log p(d) \geq E_Q[\log p(d|z)] - \text{KL}(Q(z|d)||p(z)) \quad (3)$$

where  $Q(z|d)$  is a learned approximation of the posterior distribution  $p(z|d)$ , the computation of which we describe in Section 3.2, and KL is the Kullback-Leibler divergence. Writing this out using the document likelihood we obtain the model's variational loss:

$$\mathcal{L}_{\text{var}} = E_Q \left[ \sum_{j \in \mathcal{W}_d} \log p(w_j|f(z)) \right] - \text{KL}(Q(z|d)||p(z)) \quad (4)$$

where  $j$  iterates over all unique words in document  $d$ . The purpose of this loss is to maximize the document likelihood under our modelling assumptions, where the  $E_Q$  term can be considered the reconstruction loss. The KL divergence acts as a regularizer by penalizing large differences between the approximate posterior distribution and the Bernoulli distribution with equal probability of sampling 0 and 1 ( $p = 0.5$ ), which can be computed in closed form as:

$$\text{KL}(Q(z|d)||p(z)) = Q(d) \log \frac{Q(d)}{p} + (1 - Q(d)) \log \frac{1 - Q(d)}{1 - p} \quad (5)$$

#### 3.2 Encoder function

The approximate posterior distribution  $Q(z|d)$  can be considered as the encoder function that transforms the original document representation into its hash code of  $m$  bits. We model this using a neural network that outputs the sampling probabilities used for the Bernoulli sampling of the hash code. First, we compute the representation used as input for computing the sampling probabilities:

$$v_1 = \text{ReLU}(W_a(d \odot E_{\text{imp}}) + b_a) \quad (6)$$

$$v_2 = \text{ReLU}(W_b v_1 + b_b) \quad (7)$$

where  $\odot$  corresponds to elementwise multiplication,  $W$  and  $b$  are weight matrices and bias vectors respectively, and  $E_{\text{imp}}$  is an *importance embedding* that learns a scalar for each word that is used to scale the word level values of the original document representation, and the same embedding is also used in the decoder function. The purpose of this embedding is to scale the original input such that unimportant words have less influence on the hash code generation. We transform the intermediate  $v_2$  representation to a vector of the same size as the hash code, such that the  $i^{\text{th}}$  entry corresponds to the sampling probability for the  $i^{\text{th}}$  bit:

$$Q(d) = \sigma(W_m v_2 + b_m) \quad (8)$$

where  $W_m$  and  $b_m$  have the dimensions corresponding to the code length  $m$ , and  $\sigma$  is the sigmoid function used to enforce the values to be within the interval  $[0, 1]$ , i.e., the range of probability values. The final hash code can then be sampled from the Bernoulli distribution. In practice, this is estimated by a vector  $\mu = [\mu_1, \mu_2, \dots, \mu_m]$  of

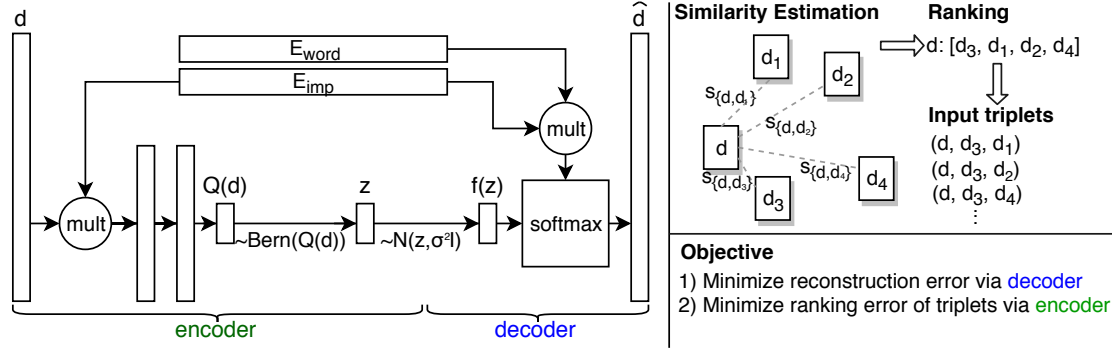


Figure 1: Model overview

values sampled uniformly at random from the interval  $[0, 1]$  and computing each bit value of either 0 or 1 as:

$$z_i = \lceil Q(d)_i - \mu_i \rceil \quad (9)$$

Sampling  $\mu$  uniformly at random corresponds to a stochastic strategy, as the same  $Q(d)$  could result in different hash codes. The opposite deterministic strategy consists of fixing  $\mu_i = 0.5$ , such that the network always generates the same code for a given document. To encourage exploration during training, the stochastic strategy is chosen, while the deterministic is used for testing. To compute the gradient of the sampled  $z$  for back-propagation, we use a straight-through estimator [2].

### 3.3 Decoder function

The purpose of the decoder function is to reconstruct the original document  $d$  given the hash code  $z$ . This is computed using the document log likelihood (Equation 1) as the sum of word log probabilities:

$$\begin{aligned} \log p(d|z) &= \sum_{j \in \mathcal{W}_d} \log p(w_j | f(z)) \\ &= \sum_{j \in \mathcal{W}_d} \log \frac{e^{f(z)^T g(E_{\text{word}}(o_j \odot E_{\text{imp}})) + b_j}}{e^{\sum_{i \in \mathcal{W}_{\text{all}}} f(z)^T g(E_{\text{word}}(o_i \odot E_{\text{imp}})) + b_i}} \end{aligned} \quad (10)$$

where the sums iterate over all unique words in document  $d$ ;  $\odot$  corresponds to elementwise multiplication;  $o_j$  is a one-hot-vector with 1 in the  $j^{\text{th}}$  position and 0 everywhere else;  $E_{\text{imp}}$  is the same importance embedding as in the encoder function;  $E_{\text{word}}$  is a word embedding;  $b$  is a bias vector;  $\mathcal{W}_{\text{all}}$  contains all vocabulary words; and the  $g$  function will be detailed later.  $E_{\text{word}}$  is a mapping from a word to a *word embedding* space, such that  $\log p(d|z)$  is maximized when the hash code is similar to most words in the document. To this end, the importance embedding assists in reducing the need to be similar to all words, as it learns to reduce the value of unimportant words. The word embedding  $E_{\text{word}}$  is made by learning a 300 dimensional embedding matrix, and  $g(E_{\text{word}}(o_j \odot E_{\text{imp}}))$  corresponds to a transformation through a fully connected linear layer to fit the code length. The choice of 300 dimensions was made to be similar in size to standard GloVe and Word2vec word embeddings [16, 19]. This two-step embedding process was chosen to allow the model to learn a code length-independent embedding initially, such that the underlying word representation is not limited by the code length.

**3.3.1 Reduce overfitting through noise injection.** We inject noise into the hash code before decoding, which has been shown to reduce overfitting and to improve generalizability in generative models [3, 12, 24]. For semantic hashing applications, this corresponds to observing significantly more *artificial* documents with small perturbations, which is beneficial for reducing overfitting in the reconstruction step. To this end we choose a Gaussian noise model, which is traditionally done for variational autoencoders [12], such that  $f(z)$  in Equation 10 is sampled as  $f(z) \sim \mathcal{N}(z, \sigma^2 I)$  where  $I$  is the identity matrix and  $\sigma^2$  is the variance. Instead of using a fixed variance, we employ variance annealing, where the variance is reduced over time towards 0. Variance annealing has previously been shown to improve performance for generative models in the image domain [3], as it reduces the uncertainty over time when the model confidence increases. However, the gradient estimate with this noise computation exhibits high variance [12], so we use the reparameterization trick to compute  $f(z)$  as:

$$f(z; \sigma^2) = z + \epsilon \sigma, \quad \epsilon \sim \mathcal{N}(0, I) \quad (11)$$

which is based on a single source of normal distributed noise and results in a gradient estimate with lower variance [12].

### 3.4 Ranking component

The variational loss guides the model towards being able to reconstruct the original document from the hash code, but no hash code similarity is enforced between similar documents. We introduce a ranking component into the model, which ensures that similar documents have a small hash code distance between them. To enable the network to learn the correct document ranking we consider document triplets as inputs,  $(d, d_1, d_2)$  with corresponding pairwise similarities of  $s_{\{d, d_1\}}$  and  $s_{\{d, d_2\}}$ . However, in the unsupervised setting we do not have a ground truth annotated ranking of the documents to extract the similarities. To this end, we generate pseudo pairwise similarities between the documents, such that weak supervision can be used to train the network in a supervised fashion.

**3.4.1 Estimating pairwise similarities.** For estimating pairwise similarities in our setting, one of many traditional ranking functions or document similarity functions could be employed. We assume such a function is chosen such that a similarity between  $d$  and  $d_1$  can be computed.

For concreteness, in this paper we choose to compute document similarities using the hash codes generated by Self-Taught Hashing (STH) [34] as this has been shown to perform well for semantic hashing (see Section 4.5). Using the STH hash codes, document similarity is computed based on the Euclidean distance between two hash codes:

$$s_{\{d, d_1\}} = -\|z^{\text{STH}} - z_1^{\text{STH}}\|_2 \quad (12)$$

where  $z^{\text{STH}}$  corresponds to the STH hash code for document  $d$ , such that  $s_{\{d, d_1\}}$  is highest when two documents are very similar. We use the  $k$ -nearest neighbour algorithm to find the top  $k$  most similar documents for each document.

**3.4.2 Ranking loss.** To train the ranking component we use a modified version of the hinge loss, as the hinge loss has previously been shown to work well for ranking with weak supervision [8]. We first define the following short-hand expressions:

$$\text{sign}_{d, d_1, d_2} = \text{sign}(s_{\{d, d_1\}} - s_{\{d, d_2\}}) \quad (13)$$

$$D_{d, d_1, d_2} = \|z - z_2\|_2^2 - \|z - z_1\|_2^2 \quad (14)$$

such that  $\text{sign}_{d, d_1, d_2}$  corresponds to the sign of the estimated pairwise document similarities, and  $D_{d, d_1, d_2}$  is the difference between the squared Euclidean distance of the hash codes of the document pairs. Using this we can define our modified hinge loss as the following piece-wise function:

$$\mathcal{L}_{\text{rank}} = \begin{cases} \max(0, \epsilon - \text{sign}_{d, d_1, d_2} D_{d, d_1, d_2}) & s_{\{d, d_1\}} \neq s_{\{d, d_2\}} \\ |D_{d, d_1, d_2}| & \text{otherwise.} \end{cases} \quad (15)$$

where  $\epsilon$  determines the margin of the hinge loss, which we fix to 1 to allow a small bitwise difference between hash codes of highly similar documents. Traditionally, the hinge loss consists only of the first part of the piece-wise function, but since the similarity estimates are based on distance computations on hash codes, some document pairs will have the same similarity. In that case the pairwise similarities are equal and the loss is simply the absolute value of  $D_{d, d_1, d_2}$ , as it should be close to 0.

### 3.5 Combining variational and ranking components

We train the variational and ranking components simultaneously by minimizing a combined weighted loss from Equation 4 and 15:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{rank}} - E_Q \left[ \sum_{j \in \mathcal{W}_d} \log p(w_j | f(z)) \right] + \beta \text{KL}(Q(z|d) || p(z)) \quad (16)$$

where  $j$  iterates over all unique words in document  $d$ ,  $\alpha$  is used to scale the ranking loss,  $\beta$  is used to scale the KL divergence of the variational loss, and we keep the unscaled version of the reconstruction part of the variational loss. During training we start with initial weight parameters of 0 and gradually increase the values in order to focus on just being able to reconstruct the input well.

	$n$	multi-class	num. classes	unique words
20news	18,846	No	20	52,447
TMC	28,596	Yes	22	18,196
Reuters	9,848	Yes	90	16,631
AGnews	127,598	No	4	32,154

Table 1: Dataset statistics

## 4 EXPERIMENTAL EVALUATION

### 4.1 Datasets

We use the four publicly available datasets summarized in Table 1. 1) *20 newsgroups*<sup>1</sup> is a dataset of posts from 20 different newsgroups. 2) *TMC*<sup>2</sup> is a dataset of NASA air traffic reports, where each report is labelled with multiple classes. 3) *Reuters21578*<sup>3</sup> is a dataset of news documents from Reuters, where each document is labelled with one or more classes. The Reuters21578 dataset is subsampled such that documents are removed if none of their associated classes are among the 20 most frequent classes. This was done by Chaidaroon and Fang [5] and their subsampled dataset was used by Shen et al. [23]. 4) *AGnews* [35] contains news articles from 4 categories.

The datasets are commonly used in related work [4, 5, 23], but without full details of preprocessing. So, in the following we describe how we preprocess the data. We filter all documents in a dataset by removing *hapax legomena*, as well as words occurring in more than 90% of the documents. In addition, we apply stopwords removal using the NLTK stopwords list<sup>4</sup>, do not apply any stemming, and use TF-IDF [22] as the document representation.

For each dataset we make a training, validation, and testing split of 80%, 10%, and 10% of the data, respectively. In all experiments the training data is used to train an unsupervised model, the validation data is used for early stopping by monitoring when the validation loss starts increasing, and the results are reported on the testing data.

### 4.2 Performance metric

The purpose of generating binary hash codes (of equal length) is to use them to obtain fast similarity searches via the Hamming distance, i.e., computing the number of bits where they differ. If two documents are semantically similar, then the generated semantic hash codes should have small Hamming distance between them. To evaluate the effectiveness of a semantic hashing method we treat each testing document as a query and perform a  $k$ -nearest-neighbour (kNN) search using the Hamming distance on the hash codes. Similarly to previous work [4, 5, 23], we retrieve the 100 most similar documents and measure the performance on a specific test document as the precision among the 100 retrieved documents (Prec@100). The total performance for a semantic hashing method is then simply the average Prec@100 across all test documents. The used datasets are originally created for text classification, but we can define two documents to be similar if they share at least one class in their labelling, meaning that multiclass documents need not to be of exactly the same classes. This definition of similarity is also used by related work [4, 5, 23].

<sup>1</sup>[http://scikit-learn.org/0.19/datasets/twenty\\_newsgroups.html](http://scikit-learn.org/0.19/datasets/twenty_newsgroups.html)

<sup>2</sup><https://catalog.data.gov/dataset/siam-2007-text-mining-competition-dataset>

<sup>3</sup><http://www.nltk.org/book/ch02.html>

<sup>4</sup>[https://www.nltk.org/nltk\\_data/](https://www.nltk.org/nltk_data/)

### 4.3 Baselines

We compare our method against traditional baselines and state-of-the-art semantic hashing methods used in related work as described in Section 2: Spectral Hashing (SpH) [29], Self-Taught Hashing (STH) [34], Laplacian co-hashing (LCH) [33], Variational Deep Semantic Hashing (VDSH) [5], NASH [23], and the neighbourhood recognition model (NbrReg) proposed by Chaidaroon et al. [4]. We tune the hyperparameters of these methods on the validation data as described in their original papers.

### 4.4 Tuning

For the encoder function (Section 3.2) we use two fully connected layers with 1000 nodes in each layer on all datasets. The network is trained using the ADAM optimizer [11]. We tune the learning rate from the set  $\{0.001, 0.0005\}$ , where 0.0005 was chosen consistently for 20news and 0.001 on the other datasets. To improve generalization we add Gaussian distributed noise to the hash code before reconstruction in the decoder function, where the variance of the sampled noise distribution is annealed over time. Initially we start with a variance of 1 and reduce it by  $10^{-6}$  every iteration, which we choose conservatively to not reduce it too fast. For the ranking component we use STH [34] to obtain a ranking of the most similar documents for each document in the training and validation set, where we choose every  $10^{th}$  document from the top 200 most similar documents. This choice was made to limit the number of triplets generated for each document as it scales quadratically in the number of similar documents to consider.

When combining the variational and ranking components of our model (Section 3.5), we added a weight parameter on the ranking loss and the KL divergence of the variational loss. We employ a strategy similar to variance annealing in this setting, however in these cases we start at an initial value and increase the weight parameters with every iteration. For the KL divergence we fix the start value at 0 and increase it by  $10^{-5}$  with every iteration. For the ranking loss we tune the models by considering starting values from the set  $\{0, 0.5, 1, 1.5\}$  and increase from the set  $\{30000^{-1}, 300000^{-1}, 1500000^{-1}, 3000000^{-1}\}$ . The code was implemented using the Tensorflow Python library [1] and the experiments were performed on Titan X GPUs.

### 4.5 Results

The experimental comparison between the methods is summarized in Table 2, where the methods are used to generate hash codes of length  $m \in \{8, 16, 32, 64, 128\}$ . We highlight the best performing method according to the Prec@100 metric on the testing data. We perform a paired two tailed t-test at the 0.05 level to test for statistical significance on the Prec@100 scores from each test document. We apply a Shapiro-Wilk test at the 0.05 level to test for normality, which is passed for all methods across all code lengths.

**4.5.1 Baseline comparison.** On all datasets and across all code lengths (number of bits) our proposed Ranking based Semantic Hashing (RBSH) method outperforms both traditional approaches (SpH, STH, and LCH) and more recent neural models (VDSH, NbrReg, and NASH). Generally, we observe a larger performance variation for the traditional methods depending on the dataset compared

to the neural approaches, which are more consistent in their relative performance across the datasets. For example, STH is among the top performing methods on Agnews, but performs among the worst on 20news. This highlights a possible strength of neural approaches for the task of semantic hashing.

Our RBSH consistently outperforms other methods to such a degree, that it generally allows to use hash codes with a factor of 2-4x fewer bits compared to state-of-the-art methods, while keeping the same performance. This provides a notable benefit on large-scale similarity searches, as computing the Hamming distance between two hash codes scales linearly with the code length. Thus, compared to prior work our RBSH enables both a large speed-up as well as a large storage reduction.

**4.5.2 Performance versus hash code length.** We next consider how performance scales with the hash code length. For all methods 128 bit codes perform better than 8 bit codes, but the performance of scaling from 8 to 128 bits varies. The performance of SpH and STH on Reuters peaks at 32 bit and reduces thereafter, and a similar trend is observed for VDSH on Agnews and TMC. This phenomenon has been observed in prior work [5, 23], and we posit that it is due to longer hash codes being able to more uniquely encode each document, thus resulting in a degree of overfitting. However, generally a longer hash code leads to better performance until the performance flattens after a certain code length, which for most methods happens at 32-64 bits.

**4.5.3 Result differences compared to previous work.** Comparing our experimental results to results reported in previous work [4, 5, 23], we observe some smaller differences most likely due to preprocessing. Previous work have not fully described the preprocessing steps used, thus to do a complete comparison we had to redo the preprocessing as detailed in Section 4.1.

On 20news and TMC the baseline performance scores we report in this paper are slightly larger for most hash code lengths. The vectorized (i.e., bag-of-words format) Reuters dataset released by the VDSH authors<sup>5</sup>, and also used in the NASH [23] paper, only consisted of 20 (unnamed) classes instead of the reported 90 classes, so these results are not directly comparable.

### 4.6 Effect of ranking component

To evaluate the influence of the ranking component in RBSH we perform an experiment where the weighting parameter of the ranking loss was set to 0 (thus removing it from the model), and report the results in Table 3. Generally, we observe that on all datasets across all hash code lengths, RBSH outperforms RBSH without the ranking component. However, it is interesting to consider the ranking component's effect on performance as the hash code length increases. On all datasets we observe the largest improvement on 8 bit hash codes, but then on Reuters, Agnews, and TMC a relatively large performance increase happens that reduces the difference in performance. On 20news the performance difference is even larger at 16 bit than at 8 bit, but as the bit size increases the difference decreases until it is marginal. This highlights that one of the major strengths of RBSH, its performance using short hash codes, can be

<sup>5</sup><https://github.com/unsuthee/VariationalDeepSemanticHashing/blob/master/dataset/reuters.tfidf.mat>

	20news					Agnews				
	8 bits	16 bits	32 bits	64 bits	128 bits	8 bits	16 bits	32 bits	64 bits	128 bits
SpH [29]	0.0820	0.1319	0.1696	0.2140	0.2435	0.3596	0.5127	0.5447	0.5265	0.5566
STH [34]	0.2695	0.4112	0.5001	0.5193	0.5119	0.6573	0.7909	0.8243	0.8377	0.8378
LCH [33]	0.1286	0.2268	0.4462	0.5752	0.6507	0.7353	0.7584	0.7654	0.7800	0.7879
VDSH [5]	0.3066	0.3746	0.4299	0.4403	0.4388	0.6418	0.6754	0.6845	0.6802	0.6714
NbrReg [4]	0.4267	0.5071	0.5517	0.5827	0.5857	0.4274	0.7213	0.7832	0.7988	0.7976
NASH [23]	0.3537	0.4609	0.5441	0.5913	0.6404	0.7207	0.7839	0.8049	0.8089	0.8142
RBSH	<b>0.5190<sup>▲</sup></b>	<b>0.6087<sup>▲</sup></b>	<b>0.6385<sup>▲</sup></b>	<b>0.6655<sup>▲</sup></b>	<b>0.6668<sup>▲</sup></b>	<b>0.8066<sup>▲</sup></b>	<b>0.8288<sup>▲</sup></b>	<b>0.8363<sup>▲</sup></b>	<b>0.8393<sup>▲</sup></b>	<b>0.8381<sup>▲</sup></b>
	Reuters					TMC				
	8 bits	16 bits	32 bits	64 bits	128 bits	8 bits	16 bits	32 bits	64 bits	128 bits
SpH [29]	0.4647	0.5250	0.6311	0.5985	0.5880	0.5976	0.6405	0.6701	0.6791	0.6842
STH [34]	0.6981	0.7555	0.8050	0.7984	0.7748	0.6787	0.7218	0.7695	0.7818	0.7797
LCH [33]	0.5619	0.6235	0.6587	0.6610	0.6586	0.6546	0.7028	0.7498	0.7817	0.7948
VDSH [5]	0.6371	0.6686	0.7063	0.7095	0.7129	0.6989	0.7300	0.7416	0.7310	0.7289
NbrReg [4]	0.5849	0.6794	0.6290	0.7273	0.7326	0.7000	0.7012	0.6747	0.7088	0.7862
NASH [23]	0.6202	0.7068	0.7644	0.7798	0.8041	0.6846	0.7323	0.7652	0.7935	0.8078
RBSH	<b>0.7409<sup>▲</sup></b>	<b>0.7740<sup>▲</sup></b>	<b>0.8149<sup>▲</sup></b>	<b>0.8120<sup>▲</sup></b>	<b>0.8088<sup>▲</sup></b>	<b>0.7620<sup>▲</sup></b>	<b>0.7959<sup>▲</sup></b>	<b>0.8138<sup>▲</sup></b>	<b>0.8224<sup>▲</sup></b>	<b>0.8193<sup>▲</sup></b>

**Table 2: Prec@100 with varying bit size. Bold marks the highest score. ▲ shows statistically significant improvements with respect to the best baseline at the 0.05 level using a paired two tailed t-test. A Shapiro-Wilk test at the 0.05 level is used to test for normality.**

	20news					Agnews				
	8 bits	16 bits	32 bits	64 bits	128 bits	8 bits	16 bits	32 bits	64 bits	128 bits
RBSH	<b>0.5190</b>	<b>0.6087</b>	<b>0.6385</b>	<b>0.6655</b>	<b>0.6668</b>	<b>0.8066</b>	<b>0.8288</b>	<b>0.8363</b>	<b>0.8393</b>	<b>0.8381</b>
RBSH w/o ranking	<u>0.4482</u>	0.5000	<u>0.6263</u>	<u>0.6641</u>	<u>0.6659</u>	<u>0.7986</u>	<u>0.8244</u>	<u>0.8344</u>	0.8332	0.8306
	Reuters					TMC				
	8 bits	16 bits	32 bits	64 bits	128 bits	8 bits	16 bits	32 bits	64 bits	128 bits
RBSH	<b>0.7409</b>	<b>0.7740</b>	<b>0.8149</b>	<b>0.8120</b>	<b>0.8088</b>	<b>0.7620</b>	<b>0.7959</b>	<b>0.8138</b>	<b>0.8224</b>	<b>0.8193</b>
RBSH w/o ranking	<u>0.7061</u>	<u>0.7701</u>	<u>0.8075</u>	<u>0.8099</u>	<u>0.8081</u>	<u>0.7310</u>	<u>0.7804</u>	<u>0.8040</u>	<u>0.8119</u>	<u>0.8172</u>

**Table 3: Effect of including the ranking component. Prec@100 with varying bit size. Bold marks the highest score and underline marks a score better than the best baseline.**

partly attributed to the ranking component. This is beneficial for the application of similarity search, as the Hamming distance scales linearly with the number of bits in the hash codes, and can thus provide a notable speed-up while obtaining a similar performance using fewer bits. Additionally, when comparing the performance of RBSH without the ranking component against the baselines in Table 2, then it obtains a better performance in 17 out of 20 cases, thus highlighting the performance of just the variational component.

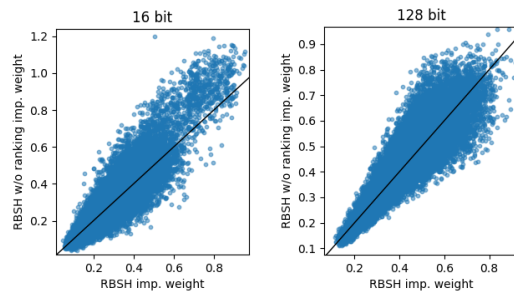
To further investigate the ranking component effect, as well as RBSH in general, in Section 4.7 we consider word level differences in the learned importance embeddings, as well as relations between inverse document frequency (IDF) and the importance embedding weights for each word. In Section 4.8 we investigate what makes a word difficult to reconstruct (i.e., using the decoder function in Section 3.3), which is done by comparing the word level reconstruction log probabilities to both IDF and the learned importance embedding weights. Finally, in Section 4.9 we do a quantitative comparison of RBSH with and without the ranking component. The comparison is based on a t-SNE [15] dimensionality reduction of the hash codes, such that a visual inspection can be performed. In the following sections we consider 16 and 128 bit hash codes generated on 20news, as these provide the largest and one of the

smallest performance difference of RBSH with and without the ranking component, respectively.

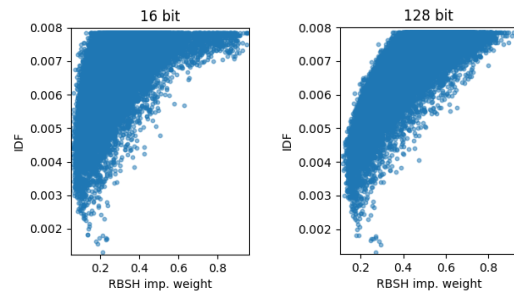
#### 4.7 Investigation of the importance embedding

We posit that the ranking component in RBSH enables the model to better differentiate between the importance of individual words when reconstructing the original document. If we consider the decoder function in Equation 10, then it is maximized when the hash code is similar to most of the importance weighted words, which in the case of equally important words would correspond to a word embedding average. However, if the hash code is short, e.g., 8 bits, then similar documents have a tendency to hash to exactly the same code, as the space of possible codes are considerably smaller than at e.g., 128 bits. This leads to worse generalizability observed on unseen documents when using short hash codes, but the ranking component enables the model to better prioritize which words are the most important. Figure 2 compares the learned importance embedding weights for 16 and 128 bit codes on 20news with and without the ranking component. For 16 bit codes we observe that RBSH without the ranking component tends to estimate a higher importance for most words, and especially for words with an RBSH importance over 0.6. This observation could be explained by the





**Figure 2: Visualization of the learned importance embedding for each word with and without using the ranking component of RBSH. The plot is made on 20news with 16 and 128 bit hash codes, and the black diagonal line corresponds to equal importance weights.**



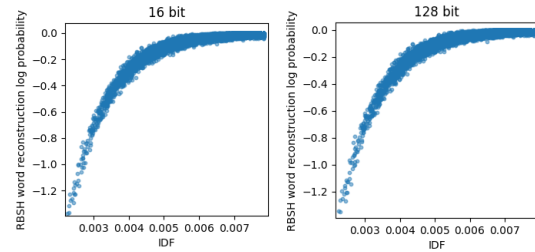
**Figure 3: Visualization of the learned importance embedding for each word compared to the inverse document frequency (IDF). The plot is made on 20news with 16 and 128 bit hash codes.**

ranking component acting as a regularizer, by enabling a direct modelling of which words are important for correctly ranking documents as opposed to just reconstruction. However, as the code length increases this becomes less important as more bits are available to encode more of the occurring words in a document, which is observed from the importance embedding comparison for 128 bits, where the over estimation is only marginal.

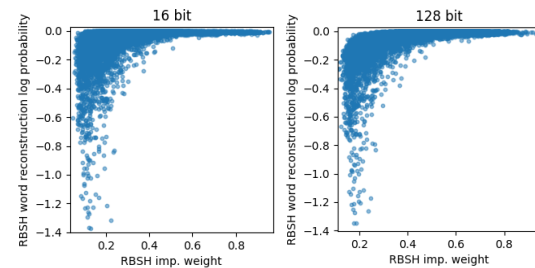
Figure 3 shows the importance embedding weights compared to the inverse document frequency (IDF) of each word. For both 16 and 128 bits we observe a similar trend of words with high importance weight that also have a high IDF; however words with a high IDF do not necessarily have a high importance weight. When we consider low importance weights, then the corresponding IDF is more evenly distributed, especially for 16 bit hash codes. For 128 bit we observe that lower importance weights are more often associated with a low IDF. These observations suggest that the model learns to emphasize rare words, as well as words of various rarity that the model deems important for both reconstruction and ranking.

#### 4.8 Investigation of the difficulty of word reconstruction

To better understand what makes a word difficult to reconstruct we study the word level reconstruction log probabilities, i.e., each summand in Equation 10, where a 0 value represents a word that is always possible to reconstruct while a smaller value corresponds to a word more difficult to reconstruct. Figure 4 compares the word



**Figure 4: Comparison of the word level reconstruction log probability compared to each word's inverse document frequency (IDF). The plot is made on 20news with 16 and 128 bit hash codes.**



**Figure 5: Comparison of the word level reconstruction log probability compared to each word's learned importance weighting. The plot is made on 20news with 16 and 128 bit hash codes.**

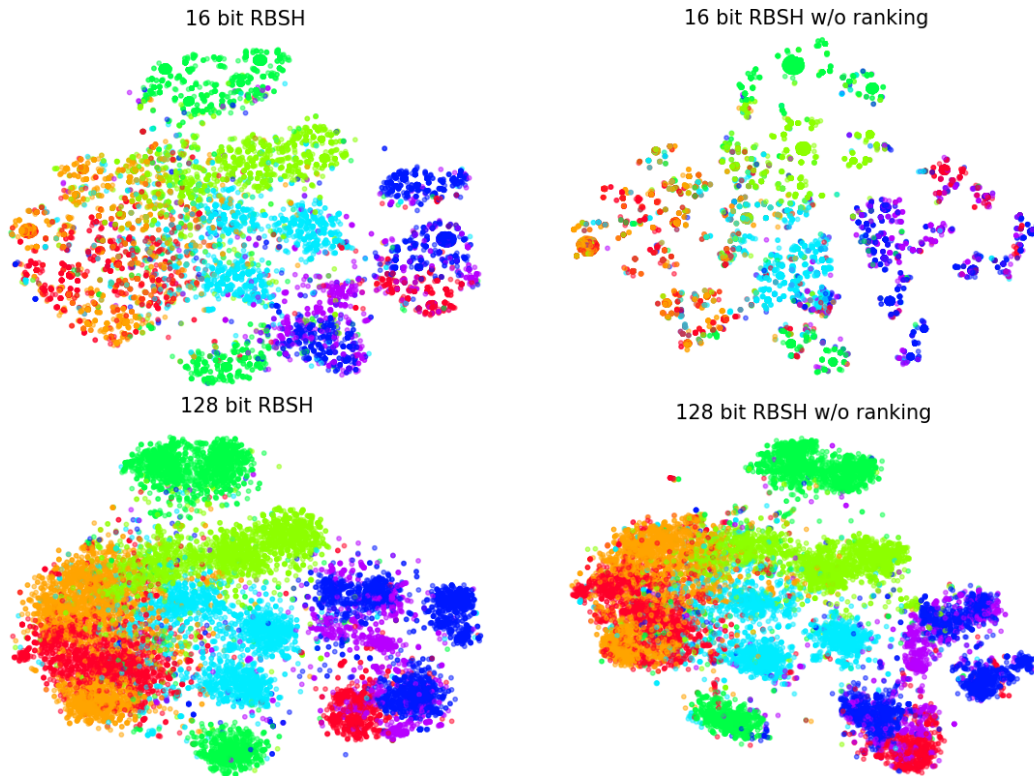
level reconstruction log probabilities to each word's IDF for 16 and 128 bit hash codes. There is no notable difference between the plots, which both show that the model prioritizes being able to reconstruct rare words, while focusing less on words occurring often. This follows our intuition of an ideal semantic representation, as words with a high IDF are usually more informative than those with a low IDF.

Figure 5 shows a comparison similar to above, where the word level reconstruction log probabilities are plotted against the learned importance embedding weights. For both 16 and 128 bit hash codes we observe that words that are difficult to reconstruct (i.e., have a low log probability) are associated with a low importance weight. Words with a low reconstruction log probability are also associated with a low IDF. This shows that the model chooses to ignore often occurring words with low importance weight. When considering words with a reconstruction log probability close to 0, then in the case of 16 bit hash codes the corresponding important weights are very evenly distributed in the entire range. In the case of 128 bit hash codes we observe that words the model reconstructs best have importance weights in the slightly higher end of the spectrum, however for lower log probabilities the two hash code lengths behave similarly. This shows that the model is able to reconstruct many words well irrespective of their learned importance weight, but words with a high importance weight are always able to be reconstructed well.

#### 4.9 Hash code visualization

In Section 4.7 we argued that the ranking component of RBSH enables the model to better prioritize important words for short hash codes, by directly modelling which words were relevant for





**Figure 6: t-SNE [15] visualization of the 16 and 128 bit hash codes from our RBSH with and with the ranking component. 20news was used as the dataset and the same color coding for class labels is used across the plots.**

ranking the documents. To further study this we perform a qualitative visualization using t-SNE [15] of 16 and 128 bit hash codes on 20news (see Figure 6), where we do the visualization for RBSH with and without the ranking component. For 16 bit hash codes we observe that RBSH without the ranking component most often creates very tight clusters of documents, corresponding to the fact that many of the produced hash codes are identical. When the ranking component is included the produced hash codes are more varied. This creates larger, more general clusters of similar documents. This leads to better generalizability as the space is better utilized, such that unseen documents are less likely to hash into unknown regions, which would result in poor retrieval performance. When considering the 128 bit hash codes for RBSH with and without the ranking component, we observe that they are highly similar, which was also expected as the Prec@100 performance was almost identical.

## 5 CONCLUSION

We presented a novel method for unsupervised semantic hashing, *Ranking based Semantic Hashing* (RBSH), which consists of a variational and ranking component. The variational component has similarities with variational autoencoders and learns to encode a input document to a binary hash code, while still being able to reconstruct the original document well. The ranking component is trained on document triplets and learns to correctly rank the documents based on their generated hash codes. To circumvent the need of labelled data, we utilize a weak labeller to estimate the

rankings, and then employ weak supervision to train the model in a supervised fashion. These two components enable the model to encode both local and global structure into the hash code. Experimental results on four publicly available datasets showed that RBSH is able to significantly outperform state-of-the-art semantic hashing methods to such a degree, that RBSH hash codes generally perform similarly to other state-of-the-art hash codes, while using 2-4x fewer bits. This means that RBSH can maintain state-of-the-art performance while allowing a direct storage reduction of a factor 2-4x. Further analysis showed that the ranking component provided performance increases on all code lengths, but especially improved the performance on hash codes of 8-16 bits. Generally, the model analysis also highlighted RBSH’s ability to estimate the importance of rare words for better hash encoding, and that it prioritizes the encoding of rare informative words in its hash code.

Future work includes incorporating multiple weak labellers when generating the hash code ranking, which under certain independence assumptions has been theoretically shown to improve performance of weak supervision [31]. Additionally, it could be interesting to investigate the effect of more expressive encoding functions, such as recurrent or convolutional neural networks, that have been used for image hashing [30, 36].

## ACKNOWLEDGMENTS

Partly funded by Innovationsfonden DK, DABAI (5153-00004A).

## REFERENCES

- [1] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: a system for large-scale machine learning. In *OSDI*, Vol. 16. 265–283.
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432* (2013).
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. 2019. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*.
- [4] Suthae Chaidaroon, Travis Ebesu, and Yi Fang. 2018. Deep Semantic Text Hashing with Weak Supervision. *ACM SIGIR Conference on Research and Development in Information Retrieval*, 1109–1112.
- [5] Suthae Chaidaroon and Yi Fang. 2017. Variational deep semantic hashing for text documents. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 75–84.
- [6] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni. 2004. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*. ACM, 253–262.
- [7] Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American society for information science* 41, 6 (1990), 391–407.
- [8] Mostafa Dehghani, Hamed Zamani, Aliaksei Severyn, Jaap Kamps, and W Bruce Croft. 2017. Neural ranking models with weak supervision. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 65–74.
- [9] Casper Hansen, Christian Hansen, Stephen Alstrup, Jakob Grue Simonsen, and Christina Lioma. 2019. Neural Check-Worthiness Ranking with Weak Supervision: Finding Sentences for Fact-Checking. In *Companion Proceedings of the 2019 World Wide Web Conference*.
- [10] Monika Henzinger. 2006. Finding near-duplicate web pages: a large-scale evaluation of algorithms. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 284–291.
- [11] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *International Conference on Learning Representations*.
- [12] Diederik P Kingma and Max Welling. 2014. Auto-encoding variational bayes. In *International Conference on Learning Representations*.
- [13] Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *ACM SIGKDD international conference on Knowledge discovery and data mining*. 426–434.
- [14] Michael S Lew, Nicu Sebe, Chabane Djeraba, and Ramesh Jain. 2006. Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 2, 1 (2006), 1–19.
- [15] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, Nov (2008), 2579–2605.
- [16] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [17] Andrew Y Ng, Michael I Jordan, and Yair Weiss. 2002. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*. 849–856.
- [18] Yifan Nie, Alessandro Sordani, and Jian-Yun Nie. 2018. Multi-level abstraction convolutional model with weak supervision for information retrieval. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 985–988.
- [19] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. In *Empirical Methods in Natural Language Processing (EMNLP)*. 1532–1543.
- [20] Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at TREC-3. *Nist Special Publication Sp 109* (1995), 109.
- [21] Ruslan Salakhutdinov and Geoffrey Hinton. 2009. Semantic hashing. *International Journal of Approximate Reasoning* 50, 7 (2009), 969–978.
- [22] Gerard Salton and Christopher Buckley. 1988. Term-weighting approaches in automatic text retrieval. *Information processing & management* 24, 5 (1988), 513–523.
- [23] Dinghan Shen, Qinliang Su, Paidamoyo Chapfuwa, Wenlin Wang, Guoyin Wang, Ricardo Henao, and Lawrence Carin. 2018. NASH: Toward End-to-End Neural Architecture for Generative Semantic Hashing. In *Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2041–2050.
- [24] Kihyuk Sohn, Honglak Lee, and Xinchun Yan. 2015. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*. 3483–3491.
- [25] Benno Stein. 2007. Principles of hash-based text retrieval. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 527–534.
- [26] Benno Stein, Sven Meyer zu Eissen, and Martin Potthast. 2007. Strategies for retrieving plagiarized documents. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 825–826.
- [27] Jun Wang, Wei Liu, Andy X Sun, and Yu-Gang Jiang. 2013. Learning hash codes with listwise supervision. In *Proceedings of the IEEE International Conference on Computer Vision*. 3032–3039.
- [28] Jingdong Wang, Ting Zhang, Nicu Sebe, Heng Tao Shen, et al. 2018. A survey on learning to hash. *IEEE transactions on pattern analysis and machine intelligence* 40, 4 (2018), 769–790.
- [29] Yair Weiss, Antonio Torralba, and Rob Fergus. 2009. Spectral hashing. In *Advances in neural information processing systems*. 1753–1760.
- [30] Ting Yao, Fuchen Long, Tao Mei, and Yong Rui. 2016. Deep Semantic-Preserving and Ranking-Based Hashing for Image Retrieval. In *IJCAI*. 3931–3937.
- [31] Hamed Zamani and W Bruce Croft. 2018. On the theory of weak supervision for information retrieval. In *Proceedings of the 2018 ACM SIGIR International Conference on Theory of Information Retrieval*. ACM, 147–154.
- [32] Hamed Zamani, W. Bruce Croft, and J. Shane Culpepper. 2018. Neural Query Performance Prediction Using Weak Supervision from Multiple Signals. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 105–114.
- [33] Dell Zhang, Jun Wang, Deng Cai, and Jinsong Lu. 2010. Laplacian co-hashing of terms and documents. In *European Conference on Information Retrieval*. Springer, 577–580.
- [34] Dell Zhang, Jun Wang, Deng Cai, and Jinsong Lu. 2010. Self-taught hashing for fast similarity search. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 18–25.
- [35] Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in neural information processing systems*. 649–657.
- [36] Fang Zhao, Yongzhen Huang, Liang Wang, and Tieniu Tan. 2015. Deep semantic ranking based hashing for multi-label image retrieval. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1556–1564.