Functional Hashing for Compressing Neural Networks

Lei Shi, Shikun Feng, Zhifan Zhu Baidu, Inc. {shilei06, fengshikun, zhuzhifan}@baidu.com

Abstract

As the complexity of deep neural networks (DNNs) trend to grow to absorb the increasing sizes of data, memory and energy consumption has been receiving more and more attentions for industrial applications, especially on mobile devices. This paper presents a novel structure based on functional hashing to compress DNNs, namely FunHashNN. For each entry in a deep net, FunHashNN uses multiple low-cost hash functions to fetch values in the compression space, and then employs a small reconstruction network to recover that entry. The reconstruction network is plugged into the whole network and trained jointly. FunHashNN includes the recently proposed HashedNets [7] as a degenerated case, and benefits from larger value capacity and less reconstruction loss. We further discuss extensions with dual space hashing and multi-hops. On several benchmark datasets, FunHashNN demonstrates high compression ratios with little loss on prediction accuracy.

1 Introduction

Deep Neural networks (DNNs) have been receiving ubiquitous success in wide applications, ranging from computer vision [22], to speech recognition [17], natural language processing [8], and domain adaptation [13]. As the sizes of data mount up, people usually have to increase the number of parameters in DNNs so as to absorb the vast volume of supervision. High performance computing techniques are investigated to speed up DNN training, concerning optimization algorithms, parallel synchronisations on clusters w/o GPUs, and stochastic binarization/ternarization, etc [10, 12, 27].

On the other hand the memory and energy consumption is usually, if not always, constrained in industrial applications [21, 35]. For instance, for commercial search engines (e.g., Google and Baidu) and recommendation systems (e.g., NetFlix and YouTube), the ratio between the increased model size and the improved performance should be considered given limited online resources. Compressing the model size becomes more important for applications on mobile and embedded devices [15, 21]. Having DNNs running on mobile apps owns many great features such as better privacy, less network bandwidth and real time processing. However, the energy consumption of battery-constrained mobile devices is usually dominated by memory access, which would be greatly saved if a DNN model can fit in on-chip storage rather than DRAM storage (c.f. [15, 16] for details).

A recent trend of studies are thus motivated to focus on compressing the size of DNNs while mostly keeping their predictive performance [15, 21, 35]. With different intuitions, there are mainly two types of DNN compression methods, which could be used in conjunction for better parameter savings. The first type tries to revise the training target into more informative supervision using dark knowledge. In specific, Hinton et al. [18] suggested to train a large network ahead, and distill a much smaller model on a combination of the original labels and the soft-output by the large net. The second type observes the redundancy existence

in network weights [7, 11], and exploits techniques to constrain or reduce the number of free-parameters in DNNs during learning. This paper focuses on the latter type.

To constrain the network redundancy, efforts [11, 25, 35] formulated an original weight matrix into either low-rank or fast-food decompositions. Moreover [15, 16] proposed a simple-yet-effective pruning-retraining iteration during training, followed by quantization and fine-tuning. Chen et al. [7] proposed HashedNets to efficiently implement parameter sharing prior to learning, and showed notable compression with much less loss of accuracy than low-rank decomposition. More precisely, prior to training, a hash function is used to randomly group (virtual) weights into a small number of buckets, so that all weights mapped into one hash bucket directly share a same value. HashedNets was further deliberated in frequency domain for compressing convolutional neural networks in [6].

In applications, we observe HashedNets compresses model sizes greatly at marginal loss of accuracy for some situations, whereas also significantly loses accuracy for others. After revisiting its mechanism, we conjecture this instability comes from at least three factors. First, hashing and training are disjoint in a two-phase manner, i.e., once inappropriate collisions exist, there may be no much optimization room left for training. Second, one single hash function is used to fetch a single value in the compression space, whose collision risk is larger than multiple hashes [4]. Third, parameter sharing within a buckets implicitly uses identity mapping from the hashed value to the virtual entry.

This paper proposes an approach to relieve this instability, still in a two-phase style for preserving efficiency. Specifically, we use multiple hash functions [4] to map per virtual entry into multiple values in compression space. Then an additional network plays in a mapping function role from these hashed values to the virtual entry before hashing, which can be also regarded as "reconstructing" the virtual entry from its multiple hashed values. Plugged into and jointly trained within the original network, the reconstruction network is of a comparably ignorable size, i.e., at low memory cost.

This functional hashing structure includes HashedNets as a degenerated special case, and facilitates less value collisions and better value reconstruction. Shortly denoted as Fun-HashNN, our approach could be further extended with dual space hashing and multi-hops. Since it imposes no restriction on other network design choices (e.g. dropout and weight sparsification), FunHashNN can be considered as a standard tool for DNN compression. Experiments on several datasets demonstrate promisingly larger reduction of model sizes and/or less loss on prediction accuracy, compared with HashedNets.

2 Background

Notations. Throughout this paper we express scalars in regular (A or b), vectors in bold (\mathbf{x}) , and matrices in capital bold (\mathbf{X}) . Furthermore, we use x_i to represent the i-th dimension of vector \mathbf{x} , and use X_{ij} to represent the (i, j)-th entry of matrix \mathbf{X} . Occasionally, $[\mathbf{x}]_i$ is also used to represent the i-th dimension of vector \mathbf{x} for specification clarity. Notation $\mathbb{E}[\cdot]$ stands for the expectation operator.

Feed Forward Neural Networks. We define the forward propagation of the ℓ -th layer

$$a_i^{\ell+1} = f(z_i^{\ell+1}), \quad \text{with} \quad z_i^{\ell+1} = b_i^{\ell+1} + \sum_{j=1}^{d^{\ell}} V_{ij}^{\ell} a_j^{\ell}, \quad \text{for } \forall i \in [1, d^{\ell+1}].$$
 (1)

For each ℓ -th layer, d^{ℓ} is the output dimensionality, \mathbf{b}^{ℓ} is the bias vector, and \mathbf{V}^{ℓ} is the (*virtual*) weight matrix in the ℓ -th layer. Vectors \mathbf{z}^{ℓ} , $\mathbf{a}^{\ell} \in \mathbb{R}^{d^{\ell}}$ denote the units before and after the activation function $f(\cdot)$. Typical choices of $f(\cdot)$ include rectified linear unit (ReLU) [29], sigmoid and tanh [2].

Feature Hashing has been studied as a dimension reduction method for reducing model storage size without maintaining the mapping matrices like random projection [32, 34]. Briefly, it maps an input vector $\mathbf{x} \in \mathbb{R}^n$ to a much smaller feature space via $\phi : \mathbb{R}^n \to \mathbb{R}^K$ with

 $K \ll n$. Following the definition in [34], the mapping ϕ is a composite of two approximate uniform hash functions $h: \mathbb{N} \to \{1, \dots, K\}$ and $\xi: \mathbb{N} \to \{-1, +1\}$. The *j*-th element of $\phi(\mathbf{x})$ is defined as:

$$[\phi(\mathbf{x})]_j = \sum_{i:h(i)=j} \xi(i)x_i. \tag{2}$$

As shown in [34], a key property is its inner product preservation, which we quote and restate below.

Lemma [Inner Product Preservation of Original Feature Hashing] With the hash defined by Eq. (2), the hash kernel is unbiased, i.e., $\mathbb{E}_{\phi}[\phi(\mathbf{x})^{\top}\phi(\mathbf{y})] = \mathbf{x}^{\top}\mathbf{y}$. Moreover, the variance is $\operatorname{var}_{\mathbf{x},\mathbf{y}} = \frac{1}{K} \sum_{i \neq j} \left(x_i^2 y_j^2 + x_i y_i x_j y_j \right)$, and thus $\operatorname{var}_{\mathbf{x},\mathbf{y}} = \mathcal{O}(\frac{1}{K})$ if $||\mathbf{x}||_2 = ||\mathbf{y}||_2 = \operatorname{const.}$

HashedNets in [7]. As illustrated in Figure 1(a), HashedNets randomly maps network weights into a smaller number of groups prior to learning, and the weights in a same group share a same value thereafter. A naive implementation could be trivially achieved by maintaining a secondary matrix that records the group assignment, at the expense of additional memory cost however. HashedNets instead adopts a hash function that requires no storage cost with the model. Assume there is a finite memory budge K^{ℓ} per layer to represent \mathbf{V}^{ℓ} , with $K^{\ell} \ll (d^{\ell}+1)d^{\ell+1}$. We only need to store a weight vector $\mathbf{w}^{\ell} \in \mathbb{R}^{K^{\ell}}$, and assign V^{ℓ}_{ij} an element in \mathbf{w}^{ℓ} indexed by a hash function $h^{\ell}(i,j)$, namely

$$V_{ij}^{\ell} = \xi^{\ell}(i,j) \cdot w_{h^{\ell}(i,j)}^{\ell}, \tag{3}$$

where hash function $h^{\ell}(i,j)$ outputs an integer within $[1,K^{\ell}]$. Another independent hash function $\xi^{\ell}(i,j):(d^{\ell+1}\times d^{\ell})\to \pm 1$ outputs a sign factor, aiming to reduce the bias due to hash collisions [34]. The resulting matrix \mathbf{V}^{ℓ} is *virtual*, since d^{ℓ} could be increased without increasing the *actual* number of parameters in \mathbf{w}^{ℓ} once the compression space size K^{ℓ} is determined and fixed.

Substituting Eq. (3) into Eq. (1), we have $z_i^{\ell+1} = b_i^{\ell+1} + \sum_{j=1}^{d^\ell} \xi^\ell(i,j) w_{h^\ell(i,j)}^\ell a_j^\ell$. During training, \mathbf{w}^ℓ is updated by back propagating the gradient via $\mathbf{z}^{\ell+1}$ (and the virtual \mathbf{V}^ℓ). Besides, the activation function $f(\cdot)$ in Eq. (1) was kept as ReLU in [7] to further relieve the hash collision effect through a sparse feature space. In both [7] and this paper, the open source $xxHash^1$ is adopted as an approximately uniform hash implementation with low cost.

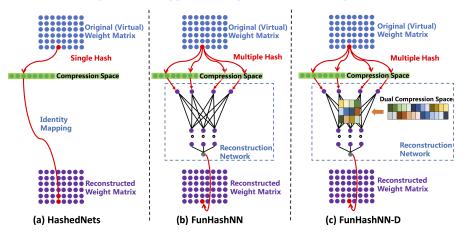


Figure 1: Illustrations of hashing approaches for neural networks compression. (a) HashedNets [7]. (b) our FunHashNN. (c) our FunHashNN with dual space hashing. (Best viewed in color)

¹http://cyan4973.github.io/xxHash/

3 Functional Hashing for Neural Network Compression

3.1 Structure Formulation

For clarity, we will focus on a single layer throughout and drop the super-script ℓ . Still, vector $\mathbf{w} \in \mathbb{R}^K$ denotes parameters in the compression space. The key difference between FunHashNN and HashedNets [7] lies in (i) how to employ hash functions, and (ii) how to map from \mathbf{w} to \mathbf{V} :

- Instead of adopting one pair of hash function (h,ξ) in Eq. (3), we use a set of multiple pairs of independent random hash functions. Let's say there are U pairs of mappings $\{h_u, \xi_u\}_{u=1}^U$, each $h_u(i,j)$ outputs an integer within [1,K], and each $\xi_u(i,j)$ selects a sign factor.
- Eq. (3) of HashedNets employs an identity mapping between one element in **V** and one hashed value, i.e., $V_{ij} = \xi(i,j)w_{h(i,j)}$. In contrast, we use a multivariate function $g(\cdot)$ to describe the mapping from multiple hashed values $\{\xi_u(i,j)w_{h_u(i,j)}\}_{u=1}^U$ to V_{ij} . Specifically,

$$V_{ij} = g\left(\left[\xi_1(i,j) w_{h_1(i,j)}, \dots, \xi_U(i,j) w_{h_U(i,j)} \right]; \alpha \right). \tag{4}$$

Therein, α is referred to as the parameters in $g(\cdot)$. Note that the input $\xi_u(i,j)w_{h_u(i,j)}$ is order sensitive from u=1 to U. We choose $g(\cdot)$ to be a multi-layer feed forward neural network, and other multivariate functions may be considered as alternatives.

As a whole, Figure 1(b) illustrates our FunHashNN structure, which can be easily plugged in any matrices of DNNs. Note that α in the reconstruction network $g(\cdot)$ is of a much smaller size compared to \mathbf{w} . For instance, a setting with U=4 and a 1-layer $g(\cdot; \alpha)$ of $\alpha \in \mathbb{R}^4$ performs already well enough in experiments. In other words, Eq. (4) just uses an ignorable amount of additional memory to describe a functional \mathbf{w} -to- \mathbf{V} mapping, whose properties will be further explained in the sequel.

3.2 Training Procedure

The parameters in need of updating include **w** in the compression and α in $g(\cdot)$. Training FunHashNN is equivalent to training a standard neural network, except that we need to forward/backward-propagate values related to **w** through $g(\cdot)$ and the virtual matrix **V**.

Forward Propagation. Substituting Eq. (4) into Eq. (1), we still omit the super-script ℓ and get

$$z_{i} = b_{i} + \sum_{j=1}^{d} a_{j} V_{ij} = b_{i} + \sum_{j=1}^{d} a_{j} \cdot g\left(\left[\xi_{1}(i, j) w_{h_{1}(i, j)}, \dots, \xi_{U}(i, j) w_{h_{U}(i, j)}\right]; \alpha\right).$$
 (5)

Backward Propagation. Denote \mathcal{L} as the final loss function, e.g., cross entropy or squared loss, and suppose $\delta_i = \frac{\partial \mathcal{L}}{\partial z_i}$ is already available by back-propagation from layers above. The derivatives of \mathcal{L} with respect to \mathbf{w} and $\boldsymbol{\alpha}$ are computed by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{i} \sum_{j} a_{j} \delta_{i} \frac{\partial V_{ij}}{\partial \mathbf{w}}, \qquad \frac{\partial \mathcal{L}}{\partial \alpha} = \sum_{i} \sum_{j} a_{j} \delta_{i} \frac{\partial V_{ij}}{\partial \alpha}, \tag{6}$$

where, since we choose $g(\cdot)$ as a multilayer neural network, derivatives $\frac{\partial V_{ij}}{\partial \mathbf{w}}$ and $\frac{\partial V_{ij}}{\partial \boldsymbol{\alpha}}$ can be calculated through the small network $g(\cdot)$ in a standard back-propagation manner.

Complexity. Concerning time and memory cost, FunHashNN roughly has the same complexity as HashedNets, since the small network $g(\cdot)$ is quite light-weighted. One key variable factor is the way to implement multiple hash functions. On one hand, if they are calculated online, then FunHashNN requires little additional time if tackling them in parallel. On the other, if they are pre-computed and stored in dicts to avoid hashing time cost, the multiple hash functions of FunHashNN demand more storage space. In application, we suggest to pre-compute hashes during offline training for speedup, and to compute hashes in parallel during online prediction for saving memory under limited budget.

3.3 Property Analysis

In this part, we try to depict the properties of our FunHashNN from several aspects to help understanding it, especially in comparison with HashedNets [7].

Value Collision. It should be noted, both HashedNets and FunHashNN conduct hashing prior to training, i.e., in a two-phase manner. Consequently, it would be unsatisfactory if hashing collisions happen among important values. For instance in natural language processing tasks, one may observe wired results if there are many hashing collisions among embeddings (which form a matrix) of frequent words, especially when they are not related at all. In the literature, multiple hash functions are known to perform better than one single function [1, 4, 5]. In intuition, when we have multiple hash functions, the items colliding in one function are hashed differently by other hash functions.

Value Reconstruction. In both HashedNets and FunHashNN, the hashing trick can be viewed as a reconstruction of the original parameter \mathbf{V} from $\mathbf{w} \in \mathcal{R}^K$. In this sense, the approach with a lower reconstruction error is preferred². Then we have at least the following two observations:

- The maximum number of possible distinct values output by hashing intuitively explains the modelling capability [32]. For HashedNets, considering the sign hashing function $\xi(\cdot)$, we have at most 2K possible distinct values of Eq. (3) to represent elements in \mathbf{V} . In contrast, since there are multiple ordered hashed inputs, FunHashNN has at most $(2K)^U$ possible distinct values of Eq. (4). Note that the memory size K is the same for both.
- The reconstruction error may be difficult to analyzed directly, since the hashing mechanism is trained jointly within the whole network. However, we observe $g\left(\left[\xi_1(i,j)w_{h_1(i,j)},\ldots,\xi_U(i,j)w_{h_U(i,j)}\right];\alpha\right)$ degenerates to $g(\xi_1(i,j)w_{h_1(i,j)})$ if we assign zeros to all entries in α unrelated to the 1st input dimension. Since $g(\xi_1(i,j)w_{h_1(i,j)})$ depends only on one single pair of hash functions, it is conceptually equivalent to HashedNets. Consequently, including HashedNets as a special case, FunHashNN with freely adjustable α is able to reach a lower reconstruction error to fit the final accuracy better.

Feature Hashing. In line with previous work [32, 34], we compare HashedNets and FunHashNN in terms of feature hashing. For specification clarity, we drop the sign hashing functions $\xi(\cdot)$ below for both methods, the analysis with which is straightforward by replacing K hereafter with 2K.

• For HashedNets, one first defines a hash mapping function $\phi_i^{(1)}(\mathbf{a})$, whose k-th element is

$$\left[\boldsymbol{\phi}_{i}^{(1)}(\mathbf{a})\right]_{k} \triangleq \sum_{j:h(i,j)=k} a_{j}, \quad \text{for} \quad k=1,\dots,K.$$
 (7)

Thus z_i by HashedNets can be computed as the inner product (details c.f. Section 4.3 in [7])

$$z_i = \mathbf{w}^{\top} \boldsymbol{\phi}_i^{(1)}(\mathbf{a}). \tag{8}$$

• For FunHashNN, we first define a hash mapping function $\phi_i^{(2)}(\mathbf{a})$. Different from a K-dim output in Eq. (7), it is of a much larger size K^U , with $\left(\sum_{u=1}^U k_u K^{(u-1)}\right)$ -th element as

$$\left[\phi_{i}^{(2)}(\mathbf{a})\right]_{\sum_{u=1}^{U} k_{u}K^{(u-1)}} \triangleq \sum_{\substack{j: h_{1}(i,j)=k_{1} \\ h_{2}(i,j)=k_{2} \\ h_{1}(i,j)=k_{2}}} a_{j}, \quad \text{for} \quad \forall u, \ k_{u}=1,\dots,K.$$
(9)

²One might argue that there exists redundancy in \mathbf{V} , whereas here we could imagine \mathbf{V} is already structured and filled by values with least redundancy.

Second, we define vector $\mathbf{g}_{\alpha}(\mathbf{w})$ still of length K^{U} , whose $\left(\sum_{u=1}^{U} k_{u} K^{(u-1)}\right)$ -th entry is

$$\left[\mathbf{g}_{\boldsymbol{\alpha}}(\mathbf{w})\right]_{\sum_{u=1}^{U} k_{u} K^{(u-1)}} \stackrel{\triangle}{=} g\left(w_{k_{1}}, w_{k_{2}}, \dots, w_{k_{U}}; \; \boldsymbol{\alpha}\right), \quad \text{for} \quad \forall u, \; k_{u} = 1, \dots, K. \tag{10}$$

Thus z_i by FunHashNN can be computed as the following inner product

$$z_i = \mathbf{g}_{\alpha}(\mathbf{w})^{\top} \boldsymbol{\phi}_i^{(2)}(\mathbf{a}). \tag{11}$$

The difference between Eq. (8) and Eq. (11) further explains the above discussion about "the maximum number of possible distinct values".

3.4 Extensions

Hashing on Dual Space. If considering a linear model $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{x}$, one can not only deliver analysis like Bayesian or hashing on input feature space of \mathbf{x} , but also do similarly on the dual space of $\boldsymbol{\theta}$ [3]. We now revisit the "reconstruction" network $g(\mathbf{x}_{ij}; \boldsymbol{\alpha})$ in Eq. (4), where vector \mathbf{x}_{ij} concatenates the hashed values $\xi_u(i,j)w_{h_u(i,j)}$ for $u=1,\ldots,U$. What we did in Eq. (4) is in fact hashing (i,j) through \mathbf{w} to get the input feature of $g(\cdot)$. In analogy, we can also hash (i,j) to fetch parameters of $g(\cdot)$, namely we have a new "reconstruction" network in the following form:

$$V_{ij} = g(\mathbf{x}_{ij}; \boldsymbol{\alpha}_{ij}), \quad \text{with} \quad [\mathbf{x}_{ij}]_u = \xi_u(i, j) w_{h_u(i, j)} \quad \text{and} \quad [\boldsymbol{\alpha}_{ij}]_r = \xi_r'(i, j) w'_{h'(i, j)}, \quad (12)$$

where $\{\xi'_r(\cdot), h'_r(\cdot)\}$ are additional multiple pairs of hash functions applied on α , and \mathbf{w}' is an additional vector in the compression space of α . The size of α_{ij} remains the same as previous. Using this trick, the maximum number of possible distinct values of \mathbf{V} further increases exponentially, so that FunHashNN has more potential ability to fit the prediction well. We denote FunHashNN with dual space hashing shortly as FunHashNN-D, and illustrate its structure in Figure 1(c).

Multi-hops. We conjecture that FunHashNN could be used in a multi-hops structure, by imagining \mathbf{w} in the compression space plays a *virtual* role similar to \mathbf{V} . Specifically, we can build another level of hash functions $\left\{\xi_u^{(1)}(\cdot), h_u^{(1)}(\cdot)\right\}$ and compression space $\mathbf{w}^{(1)}$.

Thereafter, each entry in **w** is hashed into multiple values in $\mathbf{w}^{(1)}$ via $\left\{\xi_u^{(1)}(\cdot), h_u^{(1)}(\cdot)\right\}$. Then another reconstruction network $g^{(1)}(\cdot)$ is used to learn the mapping from the hashed values in $\mathbf{w}^{(1)}$ to the corresponding entry in \mathbf{w} .

This procedure can be implemented recursively. If there are in total M-hops, what we need to save in fact just includes a (possibly much more smaller) vector $\mathbf{w}^{(M)}$ at the final hop, a series of M small reconstruction networks $\{g^{(m)}(\cdot)\}_{m=1}^{M}$, and a series of hashing functions. In contrast, the multi-hops version of HashedNets is equivalent to just adjusting the compression ratio, or say the size K.

4 Related Work

Recent studies have confirmed the redundancy existence in the parameters of deep neural networks. Denil et al. [11] decomposed a matrix in a fully-connected layers as the product of two low-rank matrices, so that the number of parameters decreases linearly as the latent dimensionality decreases. More structured decompositions Fastfood [25] and Deep Fried [35] were proposed not only to reduce the number of parameters, but also to speed up matrix multiplications. More recently, Han et al. [15, 16] proposed to iterate pruning-retraining during training DNNs, and used quantization and fine-tuning as a post-processing step. Huffman coding and hardware implementation were also considered. In order to mostly keep accuracy, the authors suggested multiple rounds of pruning-retraining. That is, for little accuracy loss, we have to prune slowly enough and thus suffer from increased training time. Again, the most related work to ours is HashedNets [7], which was then extended in [6] to random hashing in frequency domain for compressing convolutional neural networks. Either HashedNets or FunHashNN could be combined in conjunction with other techniques for better compression.

Extensive studies have been made on constructing and analyzing multiple hash functions, which have shown better performances over one single hash function [4]. One multi-hashing algorithm, d-random scheme [1], uses only one hash table but d hash functions, pretty similar to our settings. One choice alternative to d-random is the d-left algorithm proposed in [5], used for improving IP lookups. Hashing algorithms for natural language processing are also studied in [14]. Papers [32, 34] investigated feature hashing (a.k.a. the hashing trick), providing useful bounds and feasible results.

5 Experiments

We conduct extensive experiments to evaluate FunHashNN on DNN compression. Codes for fully reproducibility will be open source soon after necessary polishment.

5.1 Environment Descriptions

Datasets. Three benchmark datasets [24] are considered here, including (1) the original MNIST hand-written digit dataset, (2) dataset BG-IMG as a variant to MNIST, and (3) binary image classification dataset CONVEX. For all datasets, we use prespecified training and testing splits. In particular, the original MNIST dataset has #train=60,000 and #test=10,000, while the remaining both have #train=12,000 and #test=50,000. Moreover, collected from a commercial search engine, a large scale dataset with billions of samples is used to learn DNNs for pairwise semantic ranking. We randomly split out 20% samples from the training data to form the validation set.

Methods and Settings. In [7], the authors compared HashedNets against several DNN compression approaches, and showed HashedNets performs consistently the best, including the low-rank decomposition [11]. Under the same settings, we compare FunHashNN with HashedNets³ and a standard neural network without compression. All activation functions are chosen as ReLU.

The settings of FunHashNN are tested in two scenarios. First, we will fix to use FunHashNN in Figure 1(b) without extensions, and then compare the effects of compression by FunHashNN and HashedNets. Second, we compare different configurations of FunHashNN itself, including the number U of seeds, the layer of reconstruction network $g(\cdot)$, and extension with the dual space hashing. Hidden layers within $g(\cdot)$ keep using tanh as activation functions. Results by the multi-hops extension of FunHashNN will be included in another ongoing paper for systematic comparisons.

5.2 Varying Compression Ratio

To test robustness, we vary the compression ratio with (1) a fixed virtual network size (i.e., the size of \mathbf{V}^{ℓ} in each layer), and then with (2) a fixed memory size (i.e., the size of \mathbf{w}^{ℓ} in each layer). Three-layer (1 hidden layer) and five-layer (3 hidden layers) networks are investigated. In experiments, we vary the compression ratio geometrically within $\{1, \frac{1}{2}, \frac{1}{4}, \dots, \frac{1}{64}\}$. For FunHashNN, this comparison sticked to use 4 hash functions, 3-layer $g(\cdot)$, and without dual space hashing.

With Virtual Network Size Fixed. The hidden layer for 3-layer nets initializes at 1000 units, and for 5-layer nets starts at 100 units per layer. As the compression ratio ranges from 1 to 1/64 with a fixed virtual network size, the memory decreases and it becomes increasingly difficult to preserve the classification accuracy. The testing errors are shown in Figure 2, where standard neural networks with equivalent parameter sizes are included in comparison. FunHashNN shows robustly effective compression against the compression ratios, and persistently produces better prediction accuracy than HashedNets. It should be noted, even when the compression ratio equals to one, FunHashNN with the reconstruction network structure is still not equivalent to HashedNets and performs better.

With Memory Storage Fixed. We change to vary the compression ratio from 1 to 1/64 with a fixed memory storage size, i.e., the size of the virtual network increases while the

³HashedNets code downloaded from http://www.cse.wustl.edu/~wenlinchen/project/HashedNets/index.html

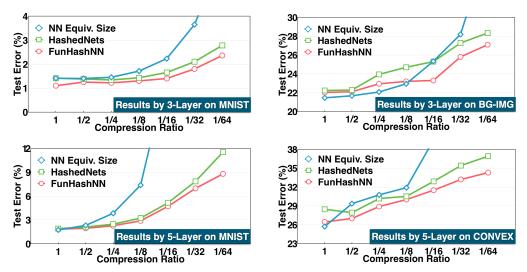


Figure 2: Testing errors by varying compression ratio with a fixed virtual network size.

number of free parameters remains unchanged. In this sense, we'd better call it expansion instead of compression. Both 3-layer and 5-layer nets initialize at 50 units per hidden layer. The testing errors in this scenario are shown in Figure 3. At all compression (expansion) ratios on each dataset, FunHashNN performs better than or at least comparably well compared to HashedNets.

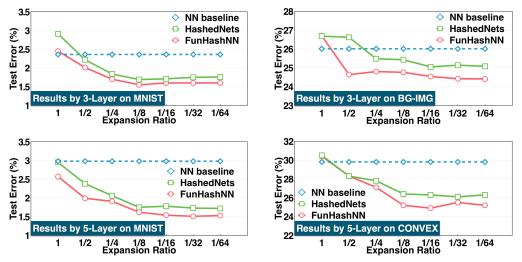


Figure 3: Testing errors by varying compression (expansion) ratio with a fixed memory storage.

5.3 Varying Configurations of FunHashNN

On 3-layer nets with compression ratio 1/8, we vary the configuration dimensions of Fun-HashNN, including the number of hash functions (U), the structure of layers of the reconstruction network $g(\cdot)$, and whether dual space hashing is turned on. Since it is impossible to enumerate all probable choices, U is restricted to vary in $\{2,4,8,16\}$. The structure of $g(\cdot)$ is chosen from $2 \sim 4$ layers, with $U \times 1$, $U \times 0.5U \times 1$, $U \times U \times 0.5U \times 1$ layerwise widths, respectively. We denote Ux-Gy as x hash functions and y layers of $g(\cdot)$, and a suffix -D indicates the dual space hashing.

Table 1 shows the performances of FunHashNN with different configurations on MNIST. The observations are summarized below. First, the series from index (0) to (1.x) fixes a 3-layer $g(\cdot)$ and varies the number of hash functions. As listed, more hash functions do not ensure

a better accuracy, and instead U4-G3 performs the best, perhaps because too many hash functions potentially brings too many partial collisions. Second, the series from (0) to (2.x) fixes the number of hash functions and varies the layer number in $g(\cdot)$, where three layers performs the best mainly due to its strongest representability. Third, indices (3.x) show further improved accuracies using dual space hashing.

Table 1: Performances on MNIST by various configurations of Fun-HashNN.

Index	Config	$\mathrm{Test}\ \mathrm{Error}(\%)$
(0)	U4-G3	1.32
(1.1)	U2-G3	1.42
(1.2)	U8-G3	1.39
(1.3)	U16-G3	1.40
(2.1)	U4-G2	1.34
(2.2)	U4-G3	1.28
(3.1)	U2-G3-D	1.36
(3.2)	U4-G3-D	1.24
(3.3)	U8-G3-D	1.27

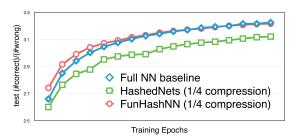


Figure 4: Performances for pairwise semantic ranking. Testing correct-to-wrong pairwise ranking ratios (the larger the better) are plotted versus the number of training epochs.

5.4 Pairwise Semantic Ranking

Finally, we evaluate the performance of FunHashNN on semantic learning-to-rank DNNs. The data is collected from logs of a commercial search engine, with per clicked query-url being a positive sample and per non-clicked being a negative sample. There are totally around 45B samples. We adopt a deep convolutional structured semantic model similar to [19, 31], which is of a siamese structure to describe the semantic similarity between a query and a url title. The network is trained to optimize the cross entropy for each pair of positive and negative samples per query.

The performance is evaluated by correct-to-wrong pairwise ranking ratio on testing set. In Figure 4, we plot the performance by a baseline network as training proceeds, compared to FunHashNN and HashNet both with 1/4 compression ratio. With U=4 hash functions, FunHashNN performs better than HashedNets throughout the training epochs, and even comparable to the full network baseline which requires 4 times of memory storage. The deterioration of HashedNets probably comes from many inappropriate collisions on word embeddings, especially for words of high frequencies.

6 Conclusion and Future Work

This paper presents a novel approach FunHashNN for neural network compression. Briefly, after adopting multiple low-cost hash functions to fetch values in compression space, FunHashNN employs a small reconstruction network to recover each entry in an matrix of the original network. The reconstruction network is plugged into the whole network and learned jointly. The recently proposed HashedNets [7] is shown as a degenerated special case of FunHashNN. Extensions of FunHashNN with dual space hashing and multi-hops are also discussed. On several datasets, FunHashNN demonstrates promisingly high compression ratios with little loss on prediction accuracy.

As future work, we plan to further systematically analyze the properties and bounds of FunHashNN and its extensions. More industrial applications are also expected, especially on mobile devices. This paper focuses on the fully-connected layer in DNNs, and the compression performance on other structures (such as convolutional layers) is also planned to be studied. As a simple and effective approach, FunHashNN is expected to be a standard tool for DNN compression.

References

- Y. Azar, A. Broder, A. Karlin, and E. Upfal. Balanced allocations. In *STOC*, 1994. C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, Inc., 1995.
- C. M. Bishop. Pattern Recognition and Machine Learning. Springer, 2006. A. Broder and A. Karlin. Multilevel adaptive hashing. In SODA, 1990.
- [5] A. Broder and M. Mitzenmacher. Using multiple hash functions to improve IP lookups. In
- INFOCOM, 2001.W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen. Compressing convolutional
- neural networks. In NIPS, 2015. W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen. Compressing neural networks
- with the hashing trick. In *ICML*, 2015.
 [8] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural
- language processing (almost) from scratch. *JMLR*, 12:2493–2537, 2011.
 [9] G. Cormode and S. Muthukrishnan. An improved data stream summary: The Count-Min sketch and its application. J. Algorithms, 55:29-38, 2005.
- [10] M. Courbariaux, Y. Bengio, and J.-P. David. BinaryConnect: Training deep neural networks
- with binary weights during propagations. In NIPS, 2015.
 [11] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. de Freitas. Predicting parameters in deep
- learning. In NIPS, 2013. T. Dettmers. 8-bit approximations for parallelism in deep learning. In ICLR, 2016. X. Glorot, A. Bordes, and Y. Bengio. Domain adaptation for large scale sentiment classification: a deep learning approach. In *ICML*, 2011.
 [14] A. Goyal, H. I. Daume, and G. Cormode. Sketch algorithms for estimating point queries in
- NLP. In EMNLP/CoNLL, 2012.
- [15] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In *ICLR*, 2016. [16] S. Han, J. Pool, J. Tran, and W. J. Dally. Learning both weights and connections for efficient
- neural networks. In NIPS, 2015.
- G. Hinton, L. Deng, D. Yu, G. E. Dahl, A. rahman Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine, 29(6):82-97, 2012.
- [18] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In NIPS
- Deep Learning Workshop, 2014.
 [19] P.-S. Huang, X. He, J. Gao, L. Deng, A. Acero, and L. Heck. Learning deep structured semantic models for web search using clickthrough data. In CIKM, 2013.
- Y. Kang, S. Kim, and S. Choi. Deep learning to hash with multiple representations. In IEEE ICDM, 2012
- [21] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin. Compression of deep convolutional
- neural networks for fast and low power mobile applications. In *ICLR*, 2016.
 [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.
 [23] B. Kulis and T. Darrell. Learning to hash with binary reconstructive embeddings. In NIPS,
- H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio. An empirical evaluations of deep architectures on problems with many factors of variation. In ICML, 2007.
- [25] Q. V. Le, T. Sarlos, and A. J. Smola. Fastfood approximating kernel expansions in loglinear time. In *ICML*, 2013.
- M. Lin, Q. Chen, and S. Yan. Network in network. In arXiv:1312.4400, 2013. Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio. Neural networks with few multiplications.
- Z. Mariet and S. Sra. Diversity networks. In ICLR, 2016.
- V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In ICML, 2010.
- [30] M. Ranzato, Y.-L. Boureau, and Y. LeCun. Sparse feature learning for deep belief networks. In *NIPS*, 2007
- [31] Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil. A latent semantic model with convolutionalpooling structure for information retrieval. In CIKM, 2014.
- Q. Shi, J. Petterson, G. Dror, J. Langford, A. Smola, and S. Vishwanathan. Hash kernels for structured data. JMLR, 10:2615–2637, 2009.
- [33] J. Wang, W. Liu, S. Kumar, and S.-F. Chang. Learning to hash for indexing big data a survey. *Proceedings of IEEE*, 104(1):34–57, 2016.
- [34] K. Weinberger, A. Dasgupta, J. Langford, A. Smola, and J. Attenberg. Feature hashing for large scale multitask learning. In ICML, 2009.
- [35] Z. Yang, M. Moczulski, M. Denil, N. de Freitas, A. Smola, L. Song, and Z. Wang. Deep fried convnets. In ICCV, 2015.