

Fast Dynamic Routing Based on Weighted Kernel Density Estimation

Suofei Zhang¹, Wei Zhao², Xiaofu Wu¹, Quan Zhou¹

¹Nanjing University of Post and Telecommunication

²SIAT, Chinese Academy of Sciences

Abstract. Capsules as well as dynamic routing between them are most recently proposed structures for deep neural networks. A capsule groups data into vectors or matrices as poses rather than conventional scalars to represent specific properties of target instance. Besides of pose, a capsule should be attached with a probability (often denoted as activation) for its presence. The dynamic routing helps capsules achieve more generalization capacity with many fewer model parameters. However, the bottleneck that prevents widespread applications of capsule is the expense of computation during routing. To address this problem, we generalize existing routing methods within the framework of weighted kernel density estimation, and propose two fast routing methods with different optimization strategies. Our methods prompt the time efficiency of routing by nearly 40% with negligible performance degradation. By stacking a hybrid of convolutional layers and capsule layers, we construct a network architecture to handle inputs at a resolution of 64×64 pixels. The proposed models achieve a parallel performance with other leading methods in multiple benchmarks.

Keywords: capsule, dynamic-routing, clustering, kernel-density-estimation, deep-learning

1 Introduction

During the last decade, deep learning algorithms, especially Convolutional Neural Networks (CNNs) have achieved remarkable progress on numerous practical vision tasks [1,2,3]. However, the stack of convolutional filters and non-linearity units still implies difficulty of understanding the internal organization of neural networks. It brings a reputation of “black box” to current neural networks [4]. Conversely, human vision systems always show much higher interpretability during the procedure of recognition. We can explicitly tell the specific cues such as shape, color, texture or intermediate semantic concepts, and build the part-whole relationship as proofs to our final decisions. The capsule structure introduces an analogous way to construct neural networks with higher interpretability. It uses grouped scalars as pose to represent specific properties of current part. Based on the multi-dimensional representation, it finds clusters within a “routing-by-agreement” framework as an interpretable representation during forward inference of network.

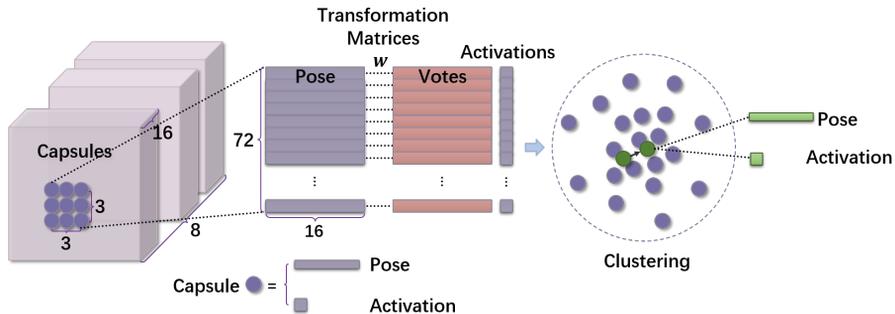


Fig. 1. Capsule structures and dynamic routing between them. A capsule consists of a group of 16 scalars as pose, as well as an 1D activation as magnitude. The pose of capsule can have the format as either vector [5] or matrix [6]. Here routing between capsules occurs in a [3,3] grid including a sum of 72 capsules. It finally results in new capsules with their activations from clusters of candidate samples as output of whole procedure. Here only one output capsule is illustrated.

Dynamic routing has been proven as an effective approach with higher generalization capacity and fewer parameters [6,7]. However, it relies on intensive computation of clustering during inference. The expense of computation prevents capsules from widespread applications in practical vision tasks. To address this problem, we model dynamic routing from the perspective of nonparametric clustering and Kernel Density Estimation (KDE), proposing a target function to explain the procedure of routing. The main contributions of this paper are twofold: (1) We equip two simplified routing methods with different optimization strategies. Comparisons in different benchmarks empirically prove that our fast methods significantly prompt the efficiency of routing with negligible performance degradation. (2) We propose a hybrid network structure consisting of both convolutional layers and capsules. The network can efficiently handle images with a resolution of 64×64 pixels. Experiments in multiple benchmarks show that our models achieve a parallel performance with other leading methods. We believe our research of simplified routing methods prompts the potentiality of capsule networks to be extensively applied in practical vision tasks.

1.1 Dynamic routing between capsules

Capsule structure is originally proposed in [8]. Differing from a convolutional neural unit, each layer here is divided into small groups of neurons called *capsules*. As an instance in Figure 1, an 128D vector of data is grouped into 8 capsules with 16D data as poses. Normally, an 1D activation is combined with the pose to represent the magnitude of current capsule. Based on such multi-dimensional representation, a dynamic routing mechanism over capsules [5] is designed between layers.

Dynamic routing is a reminiscent of the unsupervised clustering, which is purely based on the divergence between samples. Hence, different clustering methods were considered here as nonparametric routing framework [6]. Finally, the activations of capsules can be exploited for solving typical machine learning problems directly. Meanwhile, the pose of capsule can also be mapped to some specific properties as regularization. The capsule structure comes with attractive features such as higher interpretability and generalization capacity. However, the high complexity of clustering method also leads to low efficiency for both training and inference. This drawback prevents the structure from applications on large scale vision tasks such as ImageNet [9].

2 Kernel density estimation

KDE, also known as Parzen Window method, is a nonparametric technique for describing underlying empirical distribution over given samples. Since the calculation of KDE is only related to specific form of kernel function and distance metric, it normally leads to less computation as well as higher efficiency. Given n samples $\{\mathbf{u}_i | i = 1, \dots, n\}$ in the D -dimensional space \mathbb{R}^D , the multivariate KDE of random variable \mathbf{v} can be defined as

$$\hat{f}(\mathbf{v}) \triangleq \frac{1}{nz_k} \sum_{i=1}^n k(d(\mathbf{v} - \mathbf{u}_i)), \quad (1)$$

where $k(x)$ is a bounded function called *profile* function with support in univariate space \mathbb{R} . z_k is the normalization constant only related to specific $k(\cdot)$. Some instances of $k(x)$ can be illustrated as [10]

$$\text{Gaussian} : k(x) \triangleq \exp\left(-\frac{x}{2}\right), \quad \text{Epanechnikov} : k(x) \triangleq \begin{cases} 1-x & x \in [0, 1) \\ 0 & x \geq 1. \end{cases} \quad (2)$$

The distance metric between samples $d(\mathbf{v} - \mathbf{u}_i)$ can be replaced by different definitions of distance, e.g., ℓ_2 norm: $d(\mathbf{v} - \mathbf{u}_i) = \|\mathbf{v} - \mathbf{u}_i\|^2$, and Mahalanobis distance: $(\mathbf{v} - \mathbf{u}_i)^T \Sigma^{-1} (\mathbf{v} - \mathbf{u}_i)$. Intuitively, $\hat{f}(\mathbf{v})$ reflects the average distance between current point \mathbf{v} and surrounding samples.

3 Dynamic routing based on weighted KDE

KDE is the basic technique of various clustering methods [11]. Despite the relationship between dynamic routing and clustering, there still exists two major problems preventing using KDE for dynamic routing directly. First, eq. 1 only considers the case of density estimation for one cluster, rather than routing between samples and multiple clusters. Second, there is no mechanism for including

activation in the framework of KDE. To address these problems, we extend the density estimation from one cluster to mixture of clusters as

$$\hat{f}(\mathbf{v}, \mathbf{r}) \triangleq \frac{1}{n_l z_k} \sum_{j=1}^{n_{l+1}} \sum_{i=1}^{n_l} r_{ij} a_i^u k(d(\mathbf{v}_j - \mathbf{u}_i)), \quad (3)$$

where n_l and n_{l+1} are number of capsules at layer l and $l + 1$, respectively. We rewrite $\hat{f}(\mathbf{v})$ as $\hat{f}(\mathbf{v}, \mathbf{r})$ given weights \mathbf{r} as parameters of model. $\{\mathbf{v}_j | j = 1, \dots, n_{l+1}\}$ are poses of capsule at layer $l + 1$, i.e. the resulting clusters, while $\{\mathbf{u}_i | i = 1, \dots, n_l\}$ are candidate samples at layer l . Note that as shown in Figure 1, we actually use the transformed votes \mathbf{u}_{ij} instead of \mathbf{u}_i for clustering. However, since the clustering only takes place in the scope of \mathbf{v}_j and its corresponding votes $\{\mathbf{u}_{ij} | i | j = 1, \dots, n_l\}$, such change of notation will not break the following derivation. We will still use \mathbf{u}_i for simplicity. Activation a_i^u of input capsule is introduced here as prior knowledge from below layer. It is a straightforward way to let samples with higher activations give more impact to final position of cluster. We introduce r_{ij} here to measures how much \mathbf{u}_i contributes to the position of \mathbf{v}_j , namely routing weight between 2 capsules.

Since in the case of dynamic routing, clustering jointly takes place between samples and multiple clusters now. For the diversity of clusters, we want one sample can contribute to clusters in different proportions. Furthermore, the total contribution from each sample to final mixture of capsules should be equivalent. Therefore, we propose to model the procedure of dynamic routing as solving the following optimization question.

$$\mathbf{v}, \mathbf{r} = \arg \max_{\mathbf{v}, \mathbf{r}} \hat{f}(\mathbf{v}, \mathbf{r}) \quad s.t. \quad \forall i, j : r_{ij} > 0, \sum_{j=1}^{n_{l+1}} r_{ij} = 1. \quad (4)$$

We will propose two different strategies to solve eq. 4 in the following parts, and discuss the relationship between the proposed methods and other existing routing methods.

3.1 Routing based on mean shift

Mean shift [12] is a typical clustering method based on the framework of KDE. It is extensively applied in the fields of feature analysis and related vision tasks such as segmentation and object tracking [13]. The original mean shift method iteratively maximizes $\hat{f}(\mathbf{v})$ in eq. 1 by solving the following equation

$$\nabla \hat{f}(\mathbf{v}) = \frac{1}{n z_k} \sum_{i=1}^n k'(d(\mathbf{v} - \mathbf{u}_i)) \frac{\partial d(\mathbf{v} - \mathbf{u}_i)}{\partial \mathbf{v}} = 0. \quad (5)$$

By replacing $\hat{f}(\mathbf{v})$ with $\hat{f}(\mathbf{v}, \mathbf{r})$, we propose to optimize variables \mathbf{v} and \mathbf{r} alternately. First, \mathbf{v}_j^r can be updated by analogously solving eq. 5 with fixed r_{ij}^r .

as

$$\mathbf{v}_j^{\tau+1} = \frac{\sum_{i=1}^{n_l} r_{ij}^{\tau} a_i^u k'(d(\mathbf{v}_j^{\tau} - \mathbf{u}_i)) \mathbf{u}_i}{\sum_{i=1}^{n_l} r_{ij}^{\tau} a_i^u k'(d(\mathbf{v}_j^{\tau} - \mathbf{u}_i))}. \quad (6)$$

This form of \mathbf{v}_j intuitively shows that the cluster can be explained as normalized weighted summation of candidate samples. The weight consists of the coefficient r_{ij} , the prior knowledge of candidate sample and the derivate of kernel function.

Then, with updated \mathbf{v}_j^{τ} , we optimize r_{ij} with the standard gradient descent method as

$$r_{ij}^{\tau+1} = r_{ij}^{\tau} + \alpha a_i^u k(d(\mathbf{v}_j^{\tau} - \mathbf{u}_i)), \quad (7)$$

where α is the hyper parameter to control step size. In experiments we directly use $\alpha = 1$ as a common configuration. To satisfy the constraints in eq. 4, we simply normalize r_{ij} as $r_{ij} = r_{ij} / \sum_{j=1}^{n_l+1} r_{ij}$ for further calculation. Combining eq. 6 and eq. 7, a dynamic routing method can be obtained as algorithm 1.

Algorithm 1 Dynamic routing based on mean shift.

Require: poses \mathbf{u}_i , activations a_i^u

Initialize $\forall i, j: r_{ij} = 1/n_{l+1}$

for r iterations **do**

1. $\forall i, j: r'_{ij} \leftarrow \frac{r_{ij}}{\sum_j r_{ij}}$
2. $\forall j: \mathbf{v}_j \leftarrow \frac{\sum_i r'_{ij} a_i^u k'(d(\mathbf{v}_j - \mathbf{u}_i)) \mathbf{u}_i}{\sum_i r'_{ij} a_i^u k'(d(\mathbf{v}_j - \mathbf{u}_i))}$
3. $\forall i, j: r_{ij} \leftarrow r_{ij} + a_i^u k(d(\mathbf{v}_j - \mathbf{u}_i))$

end for

return capsules with poses \mathbf{v}_j

Here we omit α in algorithm 1. Note that the proposed routing method maximizes $\hat{f}(\mathbf{v}, \mathbf{r})$ by following a framework of coordinate descent [14]. Variables \mathbf{v} and \mathbf{r} are optimized alternately towards the direction of partial gradient. The value of $\hat{f}(\mathbf{v}, \mathbf{r})$ is ensured to get increased or unchanged after each iteration. Exploiting different kernel functions and distance metrics in algorithm 1 can lead to specific instances of routing methods.

3.2 Routing based on expectation maximization

Eq. 4 can be optimized within an Expectation Maximization (EM) framework [15] as well. Due to the symmetry of kernel function, $k(d(\mathbf{v}_j - \mathbf{u}_i))$ can also be explained as the likelihood of sample \mathbf{u}_i given the assumption of variable \mathbf{v}_j . From this point of view, Eq. 3 can be treated as an approximation of the log-likelihood

function given samples from the mixture model which consists of \mathbf{v}_j as components and r_{ij} as hidden weights. The EM algorithm can be exploited here to maximize $\hat{f}(\mathbf{v}, \mathbf{r})$ by alternately optimizing \mathbf{v}_j and updating r_{ij} as its expectation. In analog to the standard EM algorithm, we explicitly introduce the mixture coefficient π_j to calculate the expectation of r_{ij} , getting algorithm 2 as another strategy to solve eq. 4.

Algorithm 2 Dynamic routing based on EM algorithm.

Require: poses \mathbf{u}_i , activations a_i^u

Initialize $\forall i, j: r_{ij} = 1/n_{l+1}$

for r iterations **do**

1. $\forall i, j: r'_{ij} \leftarrow \frac{r_{ij}}{\sum_j r_{ij}}$
2. $\forall j: \mathbf{v}_j \leftarrow \frac{\sum_i r'_{ij} a_i^u \mathbf{u}_i}{\sum_i r'_{ij} a_i^u}$
3. $\forall j: \pi_j \leftarrow \frac{\sum_i r_{ij}}{\sum_j \sum_i r_{ij}}$
4. $\forall i, j: r_{ij} \leftarrow \pi_j k(d(\mathbf{v}_j - \mathbf{u}_i))$

end for

return capsules with poses \mathbf{v}_j

Algorithm 2 basically follows the standard EM algorithm to maximize $\hat{f}(\mathbf{v}, \mathbf{r})$. Comparing with another well-known application scenario of EM algorithm, the Gaussian Mixture Model (GMM), our proposed mixture model based on KDE can be treated as a simplified version of standard GMM. The simplification mainly comes from that KDE is based on nonparametric kernel function, rather than the normal distribution configured by expectation $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$. Hence the calculation of $k(d(\mathbf{v}_j - \mathbf{u}_i))$ requires much less computation than Gaussian function $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Also, the update of model in step 2 of algorithm 2 only requires calculation about \mathbf{v} without the variance $\boldsymbol{\Sigma}$ as in GMM. We will empirically prove that such simplification can still provide acceptable performance for dynamic routing.

Activation of capsule. For generic machine learning tasks, the activation of capsule is required as final result to reflect the magnitude of current capsule. It also appears in the routing procedure at above layer as prior knowledge. For the resulting capsule at layer $l + 1$, we propose a unified form of activation a_j^v for both routing methods as

$$a_j^v \triangleq \text{softmax}\left(\sum_{i=1}^{n_l} r'_{ij} a_i^u \left(k\left(\sum_{d=1}^D d(u_{id} - \beta_{jd} v_{jd}) + \beta_{j0}\right)\right)\right), \quad (8)$$

where r'_{ij} is the normalized version of r_{ij} as the result of step 1 in algorithm 1 and 2. Here the distance metric is calculated at each dimension d separately. If we ignore the linear coefficients $\beta_j \in \mathbb{R}^{D+1}$, one can see from eq. 3 that activation is the absolute value of resulting density at \mathbf{v}_j after routing. It is consistent with the original purpose of combining activation to capsule. Parameters $\boldsymbol{\beta} \in$

\mathbb{R}^{D+1} at each dimension are learned by standard back propagation to provide a linear combination rather than the rigid connection between pose and activation. Finally, a softmax function is exploited here as the guarantee of a probability distribution over activations.

Relationship between two routing methods. We propose two dynamic routing methods by maximizing the weighted KDE function $\hat{f}(\mathbf{v}, \mathbf{r})$ with different strategies. The proposed methods can be instantiated with different kernel functions and distance metrics. For generic cases without any specific background knowledge, we propose to directly adopt Epanechnikov kernel for simplification. Note that the Epanechnikov kernel can result in an identical format of step 2 in both algorithm 1 and 2. From this point of view, the only difference between two routing methods is the strategy for update of weights \mathbf{r} . One is based on gradient descent while the other is expectation.

4 Network architecture

Dynamic routing has already been proven as an effective mechanism for classification tasks on datasets with small images such as MNIST [16] and small-NORB [17]. Our main concern here is about the efficiency and generality of dynamic routing. According to the observation of our implementation, they are the main bottleneck preventing dynamic routing from widespread utilization in practical vision tasks. Hence, we designed a relatively big network architecture for 64×64 image as in Figure 2.

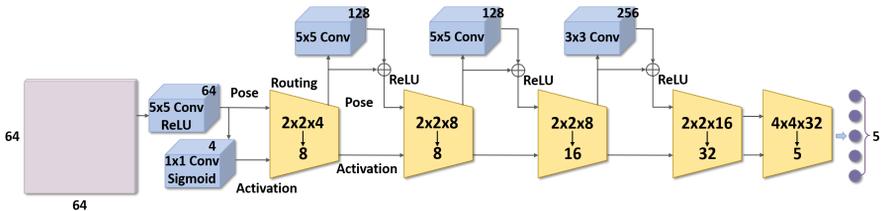


Fig. 2. The proposed network architecture for a 64×64 image. It contains 5 capsule layers where dynamic routing occurs. Except the last one, the dynamic routing at each capsule layer takes place within a 2×2 field with a stride of 2. At each block for dynamic routing, the number of input capsules is listed as width \times height \times number of capsules at each position. The number of output capsules is listed under the arrow. The side lengths of feature maps after each capsule layer are 32, 16, 8 and 4, respectively. The final dynamic routing takes all capsules within a 4×4 feature maps into account, resulting in 5 capsules as output. Here we take the smallNORB dataset for instance.

The proposed network is a hybrid architecture of convolutional layers and capsules with dynamic routing. For all convolutional layers, we adopt a stride of 1 for sliding kernels, and a padding of features to keep the feature map size.

The network starts from convolutional layers for initialization of poses and activations of capsules. Then the primary dynamic routing takes place to represent 16 capsules with 8 capsules at every 2×2 field with a stride of 2. The feature maps are downsampled here with doubled number of capsules at every position.

For the rest part of the network, we stacked a homogeneous structure sequentially as shown in Figure 2. The proposed structure consists of a capsule layer and a residual block. The residual block is analog to the structure in ResNet [18]. It takes poses from below layer as input. The poses are summed with their residual maps to ReLU functions as output. We exploit such structure to integrate convolutional layer and capsules as a generic building block for deeper and larger networks. The feature map will only be downsampled at capsule layer, along with the increase of capsules at each position. The residual structure ensures an identical or better version of pose to adapt the propagation of corresponding activation in the following dynamic routing. We have tested different structures in our experiments, founding that the proposed residual block ensures most stable training for stack of capsules, especially in deeper networks.

We compose the pose of all capsules with 4×4 matrix as [6]. It requires many fewer parameters for transformation matrices during routing. This structure of pose leads to a network as in Figure 2 with 1.2M parameters, in which nearly 90K parameters come from dynamic routing procedure. Differing from the stack of ‘‘ConvCaps’’ in [6], our network shares parameters for routing at different positions without overlap between neighboring receptive fields. So although our proposed network has more parameters than network in [6], it consumes less time for both training and inference.

For comparison, we also implement an analogous CNN to the architecture in Figure 2 as baseline. The network consists of 5 convolutional layers without residual structure and a global average pooling to 5-way softmax output. Except for the 1×1 filters for initialization of activations, the structures of other 4 convolutional layers in Figure 2 are all shared into the baseline CNN. The 5th convolutional layer contains a 3×3 kernel with 5 channels and a stride of 1. Capsule layers are replaced by direct max pooling between convolutional layers. All hidden layers are connected to the ReLU non-linearity. We designed the baseline CNN with an approximately equivalent amount of parameters as well as similar structures to our proposed network, ensuring a fair comparison between them.

5 Experiments

5.1 Implementation details

We compared three routing methods and a CNN baseline in different benchmarks. Here the Fast Routing based on Mean Shift (FRMS), Fast Routing based on EM (FREM), and EM routing from [6] are implemented within the proposed network architecture for consideration. During the training of these routing methods, we also integrate the reconstruction error from pose of resulting capsule as regularization in loss function. The structure of decoder from pose

to original image is analogous to [5]. We resize the input to 32×32 image as ground truth in reconstruction to control the scale of fully connected layers.

In practical implementation, we use softmax function instead of standard normalization for the calculation of r'_{ij} in step 1 of above algorithms. This modification can relax the kernel bandwidth restriction in KDE with only trivial increase of computation. With these modifications, we adopt the Epanechnikov kernel and ℓ_1 norm: $d(\mathbf{v} - \mathbf{u}_i) = |\mathbf{v} - \mathbf{u}_i|$ as the default configuration for all experiments.

5.2 Evaluation on smallNORB

We started our experiments from comparison of different routing methods on smallNORB dataset, which contains 5 categories of targets with an amount of 24,300 96×96 images. During training we randomly cropped 64×64 patches and augment data with random brightness and contrast. During test we directly cropped a 64×64 patch from the center of image. The image is centralized with unit deviation as input to the network. We implemented all the algorithms with TensorFlow 1.2.0 [19] on a hardware with Intel Xeon E5-2680 CPU at 2.4GHz and NVIDIA Tesla P40 GPU. We train the network for 50 epochs with mini-batches of size 50.

We report the results of different methods in Table 1. All the models are trained with 2 routing iterations for clustering. The best result is achieved by exploiting FREM within our proposed network architecture. The test error 2.2% is in-par with state-of-the-art performance gained by another routing method [6]. Note that our method achieves the results at a higher resolution with much higher efficiency.

From the perspective of efficiency, one can see that although routing methods consume more time for both training and inference than baseline CNN, they significantly reduce the test error by 80% with similar network architecture. Moreover, the FREM and FRMS methods reduce the time consumption by nearly 40% from EM routing. FREM slightly outperforms EM routing with our architecture.

Table 1. Performance evaluation on smallNORB.

Method	Inference time	Training time	Test error rate
FRMS	$0.158 \pm 0.001s$	$0.470 \pm 0.002s$	2.6%
FREM	$0.158 \pm 0.003s$	$0.471 \pm 0.002s$	2.2%
EM routing	$0.252 \pm 0.003s$	$0.744 \pm 0.003s$	2.3%
Baseline CNN	$0.043 \pm 0.003s$	$0.064 \pm 0.001s$	11.3%

Ablation study. We evaluated the influences of varying different configurations to routing methods on smallNORB in Table 2. One can see that different number of iterations for dynamic routing severely impacts the final performance

of models. Setting iteration number to 1 leads to failures of training for all methods. For the number of iterations as 3, we observed a side effect brought by long paths between capsules. We assume that this is because the training of transformation matrices relies on the back propagation of gradients from every \mathbf{u}_i at different step of iterations. Routing with more iterations tends to be impacted by the vanishing gradient problem [18]. Hence, we recommend to use 2 iterations as the default configuration for routing.

Meanwhile, we also evaluated the influence of different initialization configurations. We adopt the Truncated Normal Distribution (TND) with 0 mean and 0.01 standard deviation as the common initialization for all trainable weights. However, for transformation matrices, we tried higher standard deviation as initialization. One can see from Table 2 that, TND with standard deviation as 0.1 or 1.0 provides much better performances than 0.01. We assume that higher deviation can lead to more differences between transformation matrices, which can serve as better initialization of the clustering.

Table 2. Ablation study on smallNORB. Here we omit variances of time consumption for simplicity. “Stddev” represents the standard deviation of TND for initialization.

Method	Routing iterations	Stddev	Test error rate
FRMS	1	0.1	76.9%
FRMS	2	0.01	5.9%
FRMS	2	0.1	2.6%
FRMS	2	1.0	2.7%
FRMS	3	0.1	7.1%
FREM	1	0.1	77.8%
FREM	2	0.01	5.6%
FREM	2	0.1	2.2%
FREM	2	1.0	2.3%
FREM	3	0.1	6.0%
EM routing	1	0.1	70.2%
EM routing	2	0.01	6.2%
EM routing	2	0.1	2.3%
EM routing	2	1.0	2.3%
EM routing	3	0.1	5.8%

Training loss. We illustrate the curves of training loss in Figure 3. Our training loss consists of output loss defined on activations and reconstruction loss defined on poses. For the output loss function, we tried margin loss [5] and spread loss [6] for all candidate methods. We implement the spread loss by increasing the margin from 0.2 to 0.9 within 5 epochs. It turns out FREM can be well trained with both kinds of loss functions, while the training of FRMS and EM routing seriously relies on the relaxation from spread loss to ensure the diversity of capsules.

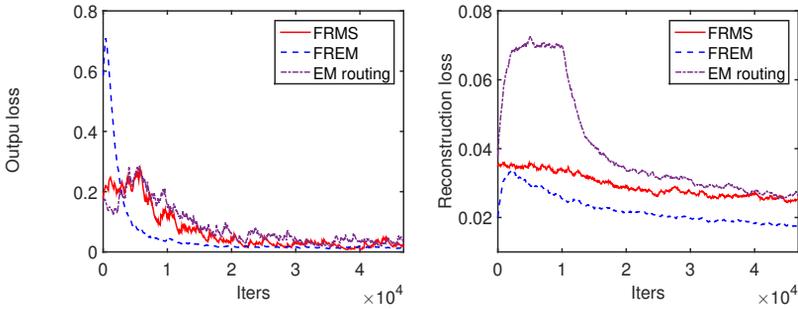


Fig. 3. The tendencies of output loss and reconstruction loss.

5.3 Other datasets

We also tested the proposed methods on MNIST, Fashion-MNIST [20] and CIFAR10 [21]. To adapt the size of images in these datasets, we removed the second capsule layer as well as the residual block from the network in Figure 2. With four capsule layers the network can process an input image with 32×32 pixels. We resize all images to this resolution in the experiments.

We listed all results on different datasets for comparison in Table 3. For MNIST, all methods approximately achieve the same accuracy, since the results are nearly saturated. For Fashion-MNIST, our methods outperform another implementation of routing-by-agreement method with 8M parameters [20]. The FRMS method slightly outperforms the FREM on this dataset. For CIFAR10, we modified the first 5×5 convolutional layer of network as 256 output channels with 3 color channels as input. The reported results are achieved by an ensemble of 5 models with the same hyper parameters as in the experiments on small-NORB. We omit the comparison of time consumptions of methods here since it is basically consistent to Table 1. In the case of lower resolution, there is a trivial gap between our methods and the EM routing. However, our efficient methods for routing still show high potentiality to prompt the performance of baseline CNN.

Table 3. Results of proposed methods on datasets.

Method	MNIST	Fashion-MNIST	CIFAR10
FRMS	0.42%	6.0%	15.6%
FREM	0.38%	6.2%	14.3%
EM routing	0.32%	5.8%	11.6%
Baseline CNN	0.65%	7.6%	19.2%

We also illustrate some reconstruction samples on different datasets in Figure 4. One can see that routing methods result in reasonable reconstruction on

MNIST and Fashion-MNIST. In contrast, they can only reconstruct an implicit prior of input on smallNORB, despite that the reconstruction errors do reduce in Figure 3. The training of model on smallNORB seems not impacted by this failure. We assume that smallNORB mainly consists of target models from different azimuths without any noise from background. The mechanism of routing can handle such affine transformation fairly well. We met the same failures of reconstruction on CIFAR10 as well. It seriously impacts the final performance of capsule networks. Structures with higher complexity rather than only fully connected layers in decoder could be a potential way to mitigate this problem.

6 Discussion and Related Work

We treat dynamic routing as an advanced compression of knowledge from capsules at layer l to $l + 1$. Features from candidate capsules are represented by a mixture of less capsules with corresponding magnitudes. The compression mainly relies on the unsupervised clustering as well as transformation matrices. Conventional clustering methods for mixture models normally start from randomized samples as initialization, while dynamic routing relies on transformation matrices as guarantee of diversity over different clusters. The learning of transformation matrices by back propagation thus can be modeled as the procedure of compressing knowledge into new clusters for a more efficient representation. Based on these principles, we proposed two simplified routing methods for better cooperation with back propagation learning. Our methods can also be treated as a generalized version of existing routing methods.

FRMS and routing-by-agreement [5]. We use the term ‘‘Routing-by-agreement’’ to represent the original routing method in [5]. By a non-linear ‘‘squashing’’ function, all lengths of capsules here are compressed into $[0, 1]$ in analog to probabilities. Since routing-by-agreement is only deployed once at the ‘‘DigitCaps’’ layer, no prior from below layer is considered, i.e. a_i^u can be omitted from algorithm 1. With the same modification, it can be shown that by using the Epanechnikov kernel function and a non-strict distance metric from the cosine similarity as

$$\langle \mathbf{u}, \mathbf{v} \rangle \triangleq 1 - \mathbf{u}^T \mathbf{v}, \quad (9)$$

a quite similar routing method to routing-by-agreement can be derived from FRMS as algorithm 3. Here we replace $d(\mathbf{u} - \mathbf{v})$ by $\langle \mathbf{u}, \mathbf{v} \rangle$ due to its specific form.

Routing-by-agreement directly assigns length of pose as the activation of capsule without any linear coefficient β . We can analogously modify the definition of activation in eq. 8 as

$$a_j^v \triangleq \text{softmax} \left(\sum_{i=1}^{n_i} r'_{ij} \mathbf{u}_i^T \mathbf{v}_j \right). \quad (10)$$

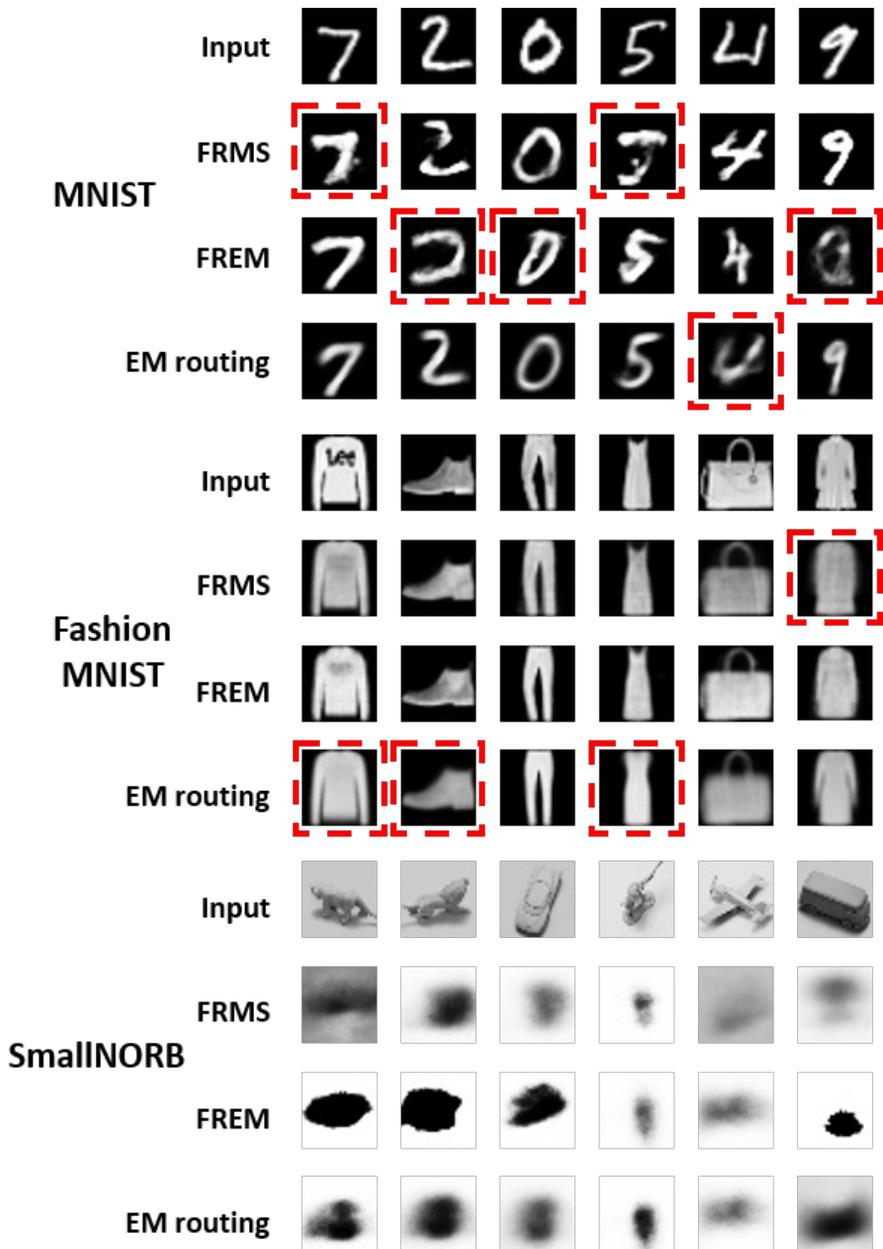


Fig. 4. Samples of reconstructions with different routing methods on datasets. For MNIST and Fashion-MNIST, we subjectively signed some reconstruction results with red dashed rectangle as reference. Some of these samples lack of details while others show confusing features for classification.

Algorithm 3 Variant of FRMS.

Require: poses \mathbf{u}_i , activations a_i^u
 Initialize $\forall i, j: r_{ij} = 1/n_{l+1}$
for r iterations **do**
 1. $\forall i, j: r'_{ij} \leftarrow \text{softmax}(r_{ij})$
 2. $\forall j: \mathbf{v}_j \leftarrow \sum_i r'_{ij} \mathbf{u}_i$
 3. $\forall i, j: r_{ij} \leftarrow r_{ij} + \mathbf{u}_i^T \mathbf{v}_j$
end for
return capsules with poses \mathbf{v}_j

Here we also remove a_i^u from the equation for an easier comparison. Note that since we update \mathbf{v}_j here as $\mathbf{v}_j \leftarrow \sum_{i=1}^{n_l} r'_{ij} \mathbf{u}_i$, eq. 8 can be rewritten as the squared length of \mathbf{v}_j directly.

However, from eq. 5 we can see that the derivation of step 2 in algorithm 3 is suspicious due to the specific form of $\langle \mathbf{u}, \mathbf{v} \rangle$. The convergence of whole algorithm is not well ensured. During our experiments, we also observed that stack of algorithm 3 can easily lead to failures of training. So we remove the method from comparison.

FREM and EM routing [6]. The proposed FREM can be treated as a simplified version of EM routing. The detailed analysis of the relationship between KDE based mean shift and EM could be referred to [22]. EM routing addresses the clustering question completely within the framework of GMM, maximizing the log-likelihood function with Gaussian function as the divergence metric. The differences between FREM and EM routing can be mainly summarized as follows: (1) EM routing exploits activation a_j^v in the update of weight r_{ij} as prior knowledge of cluster. We use the mixture coefficient π_j as standard EM algorithm for simplification. According to our implementation, this modification only brings trivial impact to routing performance. (2) Gaussian function provides asymmetrical kernels with the variance Σ at each dimension of cluster. Our KDE based kernel is symmetrical at all dimensions. This simplification leads to non-trivial influence to the performance of our routing method. However, with the utilization of convolutional layers, the impact is significantly mitigated to an acceptable level.

7 Conclusion

In this paper, we propose two efficient routing methods by generalizing existing methods within the framework of weighted KDE. Rather than constructing network with capsule structures independently for higher performance, we propose to exploit capsules and dynamic routing as effective complements with convolutional units. Experimental results show that such hybrid structure is promising to provide efficient solutions with much higher capacity for large scale vision tasks. In our future work, we plan to further prompt the performance of capsule networks on practical image datasets with higher resolution, e.g., STL-10 [23].

References

1. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems* **25**(2) (2012) 2012
2. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. (2014) 1725–1732
3. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. *nature* **521**(7553) (2015) 436
4. Alain, G., Bengio, Y.: Understanding intermediate layers using linear classifier probes. arXiv preprint arXiv:1610.01644 (2016)
5. Sabour, S., Frosst, N., Hinton, G.E.: Dynamic routing between capsules. In: *Advances in Neural Information Processing Systems*. (2017) 3859–3869
6. Hinton, G.E., Sabour, S., Frosst, N.: Matrix capsules with em routing. In: *ICLR 2018 Conference*. (2018) accepted
7. Wei, Z., Jianbo, Y., Min, Y., Zeyang, L., Suofei, Z., Zhou, Z.: Investigating capsule networks with dynamic routing for text classification. In: *Conference on Empirical Methods on Natural Language Processing*. (2018)
8. Hinton, G.E., Krizhevsky, A., Wang, S.D.: Transforming auto-encoders. In: *International Conference on Artificial Neural Networks*, Springer (2011) 44–51
9. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* **115**(3) (2015) 211–252
10. Wand, M.P., Jones, M.C.: *Kernel smoothing*. Crc Press (1994)
11. Schwander, O., Nielsen, F.: Learning mixtures by simplifying kernel density estimators. In: *Matrix Information Geometry*. Springer (2013) 403–426
12. Comaniciu, D., Meer, P.: Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on pattern analysis and machine intelligence* **24**(5) (2002) 603–619
13. Comaniciu, D., Ramesh, V., Meer, P.: Real-time tracking of non-rigid objects using mean shift. In: *IEEE Conference on Computer Vision and Pattern Recognition*. Volume 2. (2000) 142–149
14. Wright, S.J.: Coordinate descent algorithms. *Mathematical Programming* **151**(1) (2015) 3–34
15. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)* (1977) 1–38
16. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11) (1998) 2278–2324
17. LeCun, Y., Huang, F.J., Bottou, L.: Learning methods for generic object recognition with invariance to pose and lighting. In: *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*. Volume 2., IEEE (2004) II–104
18. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. (2016) 770–778
19. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., et al.: Tensorflow: Large-scale machine learning on heterogeneous distributed systems. arXiv preprint arXiv:1603.04467 (2016)

20. Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms (2017)
21. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. (2009)
22. Carreira-Perpinan, M.A.: Gaussian mean-shift is an em algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **29**(5) (May 2007) 767–776
23. Coates, A., Ng, A., Lee, H.: An analysis of single-layer networks in unsupervised feature learning. In: *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. (2011) 215–223