

# Attention routing between capsules

Jaewoong Choi  
Seoul National University  
chjw1475@snu.ac.kr

Hyun Seo  
Seoul National University  
hseo0618@snu.ac.kr

Sui Im  
Seoul National University  
a5828167@snu.ac.kr

Myungju Kang  
Seoul National University  
mkang@snu.ac.kr

## Abstract

*In this paper, we propose a new capsule network architecture called Attention Routing CapsuleNet (AR CapsNet). We replace the dynamic routing and squash activation function of the capsule network with dynamic routing (CapsuleNet) with the attention routing and capsule activation. The attention routing is a routing between capsules through an attention module. The attention routing is a fast forward-pass while keeping spatial information. On the other hand, the intuitive interpretation of the dynamic routing is finding a centroid of the prediction capsules. Thus, the squash activation function and its variant focus on preserving a vector orientation. However, the capsule activation focuses on performing a capsule-scale activation function.*

*We evaluate our proposed model on the MNIST, affNIST, and CIFAR-10 classification tasks. The proposed model achieves higher accuracy with fewer parameters ( $\times 0.65$  in the MNIST,  $\times 0.82$  in the CIFAR-10) and less training time than CapsuleNet ( $\times 0.19$  in the MNIST,  $\times 0.35$  in the CIFAR-10). These results validate that designing a capsule-scale operation is a key factor to implement the capsule concept.*

*Also, our experiment shows that our proposed model is transformation equivariant as CapsuleNet. As we perturb each element of the output capsule, the decoder attached to the output capsules shows global variations. Further experiments show that the difference in the capsule features caused by applying affine transformations on an input image is significantly aligned in one direction.*

## 1. Introduction

Convolutional neural networks (CNNs) have had much success in computer vision tasks [8] [13] [3] [6]. The convolutional layer is an effective method to extract local features due to its local connectivity and parameter sharing with spatial location. However, the convolutional layer has a limited

ability to encode a transformation. For example, if the convolutional layer is combined with a max-pooling layer, the extracted feature is local translation invariant. As CNN models become deeper [3] [15], the receptive field of each feature is getting larger. Then, the information loss from the translation invariance also increases.

To overcome the transformation invariance of CNNs, the transforming autoencoder [4] uses the concept of "capsule". A capsule is a vector representation of a feature. Each capsule not only represents a specific type of entity but also describes how the entity is instantiated, such as precise pose and deformation. In other words, the capsules are transformation equivariant.

The CapsuleNet [12] is a novel method that implements the idea of the capsules. By introducing the dynamic routing algorithm and squash activation function 1, CapsuleNet uses vector-output capsules as a basic unit instead of scalar-output features.

$$\text{squash}(\mathbf{s}_j) = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|} \quad (1)$$

where  $\mathbf{s}_j$  is a pre-activation capsule. However, CapsuleNet has a room for development. The number of parameters of CapsuleNet is much larger than that of comparable performance CNN-based models. Also, the dynamic routing is an iterative process. The reported accuracy of CapsuleNet on the benchmark datasets like CIFAR-10 is inferior to state-of-the-art CNN models. [18]

In this paper, we propose a convolutional capsule network architecture comprised of building blocks of CNNs. We substitute the dynamic routing and squash capsule-activation function of CapsuleNet [12] with attention routing and capsule activation. In the attention routing, the log probabilities of agreement coefficients between the  $l$ th layer and the  $(l + 1)$ th layer are learned by a scalar-product between the capsules of the  $l$ th layer and the kernel of convolution. The kernel of convolution serves as an approxi-

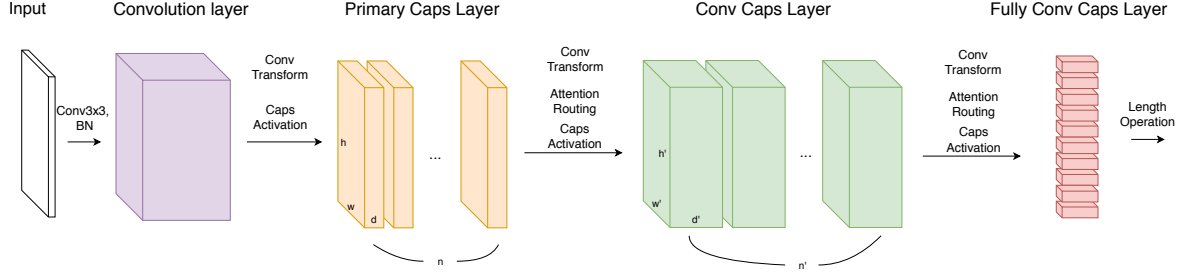


Figure 1. Overview of AR CapsNet. AR CapsNet is composed of primary caps layer, conv caps layer, and fully conv caps layer. BN denotes the batch normalization. Conv Transform and Caps Activation denotes the convolutional transform and capsule activation respectively.

mation of the reference vector to perform routing. By replacing an iterative process of the dynamic routing with forward-pass convolution, the attention routing is fast while maintaining spatial information. Two important properties of squash activation function 1 is that the squash activation function preserves a vector orientation and is a capsule-wise activation function, not an element-wise activation function such as ReLU or tanh. The dynamic routing is an unsupervised algorithm to find a centroid-like output capsule of the prediction capsules. Therefore, the squash activation function and its variant 2 [18] focus on preserving a capsules orientation.

$$\text{squash variant}(s_j) = \left(1 - \frac{1}{\exp(\|s_j\|)}\right) \frac{s_j}{\|s_j\|} \quad (2)$$

However, we focus on the capsule-wise operation rather than preserving orientation. The capsule activation performs an affine transform on the capsules and then applies an element-wise activation function. The capsules on the same capsule channel share parameters used in the affine transformation. Thus, the capsules on the same capsule channel are mapped to the same feature space, and the operation is parameter efficient. Therefore, the capsule activation is a capsule-wise function that does not preserve a vector orientation. Since the capsule activation applies a non-linear transformation to a linear combination of the prediction capsules, parametrizing the routing process through the attention routing is compatible. We refer to our proposed model as *Attention Routing CapsuleNet (AR CapsNet)*.

We evaluate the AR CapsNet on three datasets (MNIST, affNIST, and CIFAR-10). The AR CapsNet significantly outperforms CapsuleNet in the affNIST and CIFAR-10 classification task and shows a comparable performance in the MNIST dataset while being faster and using less than half parameters than CapsuleNet. Moreover, the AR CapsNet preserves the transformation equivariant property of CapsuleNet. As we perturb each element of the output capsule, the decoder attached to the output capsules shows global variations as in [12]. Further experiment showed that the affine transformations on an input image cause the fea-

ture capsules to change in the significantly aligned direction. From these experiments, we prove that the AR CapsNet encodes an affine transformation on the input image in some basis of capsule space. In addition, our proposed architecture is constructed in a convolutional manner so that it can be easily extended to a deeper network structure.

### 1.1. Contribution

- We propose a new architecture called AR CapsNet by introducing two modifications to the CapsuleNet [12]. These modifications are the *attention routing* and *capsule activation*.
- The capsule activation expands the concept of the existing capsule-wise activation functions such as the squash activation. The capsule activation performs an orientation-nonpreserving transform on the pre-activation capsules. The performance of the AR CapsNet demonstrates that the transformation equivariant features can be extracted even if the routing process is not restricted to the clustering approach and the capsule activation is not limited to the normalization.
- The AR CapsNet shows better results on the affNIST, and CIFAR-10 classification tasks and comparable results on the MNIST classification task while using much smaller parameters than CapsuleNet. Also, the AR CapsNet preserves the transformation equivariant property of the CapsuleNet. As we perturb each element of the output capsule, the decoder attached to the output capsule shows global variation as in [12].
- To investigate the transformation equivariance further, we suggest a new experiment. We observe the difference in the output capsule caused by applying transformations on an input image. In the AR CapsNet, these difference vectors are significantly aligned compared to a set of random vectors. These results demonstrate that transformation on an input image is encoded in some basis vector.

## 2. Related Works

The CNN models that consist of convolutional layers and max-pooling layers have a local translation invariance. To overcome transformation invariance, CapsuleNet [12] uses vector-output capsules and the dynamic routing in place of scalar-output features and max-pooling. By demonstrating that the dimension perturbation of digit capsules leads to a global transformation of the reconstruction image, CapsuleNet claims to have transformation equivariance.

A number of methods to improve the performance of CapsuleNet have been proposed in [5] [17] [1] [9] [10]. In [17], they interpreted the routing-by-agreement process as an optimization problem of minimizing clustering loss. They proposed another routing process from the point of view of clustering. Their approach achieved better results on an unsupervised perceptual grouping task compared to [12]. The matrix capsules with EM routing [5] proposed another routing method called EM routing. The EM routing measures compatibility between matrix capsules by clustering matrix capsules through Gaussian distributions. The matrix capsules with EM routing achieved the state-of-the-art performance on a shape recognition task using the small-NORM dataset. The spectral capsule networks [1] is a variation of [5]. Spectral capsule networks use a singular value to compute the activation of each capsule instead of the logistic function in [5]. Spectral capsule networks achieved better performance on a diagnosis dataset compared to [5] and deep GRU networks while showing faster convergence compared to [5].

The SegCaps [9] applied a capsule network to the object segmentation task. The SegCaps introduced two modifications to the CapsuleNet and devised the concept of deconvolutional capsules from these modifications. The two modifications are the locally connected dynamic routing and the sharing of transformation matrices within the same capsules channel. The sharing of transformation matrices is equivalent to the convolutional transform of our conv caps layer except for the addition of biases. The EncapNet [10] performs a one-time pass approximation of the routing process by introducing two branches. The master branch extracts a feature from the locally connected capsules as in [9] and the aide branch combines information from all the remaining capsules. Also, they introduced a Sinkhorn divergence loss which works as a regularizer. The EncapNet achieved competitive results on CIFAR-10/100, SVHN, and a subset of ImageNet.

Our proposed model uses attention architecture as a routing algorithm. The attention architecture learns a compatibility function between low-level features and high-level features. In [2], the output of attention architecture is a weighted sum of input features, and the weights are the compatibilities based on the input features and the RNN hidden state. The compatibility function is a feedforward

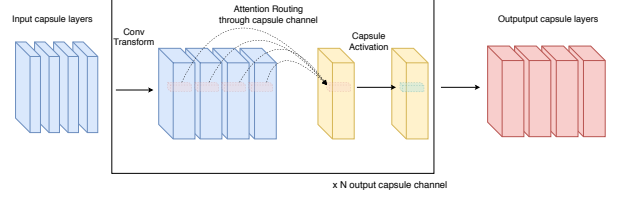


Figure 2. Detailed operation process of conv caps layer. Conv Transform denotes the convolutional transform. The convolutional transform performs a locally connected affine transform on each capsule channel. The attention routing learns the agreement between the convolutional transformed capsules for each spatial location. The capsule activation applies a capsule-wise activation function on each capsule channel.

neural network with a softmax activation function. In [11], they experimented on various kinds of attention architectures from global attention to local attention and three compatibility functions. One of the three compatibility functions was a softmax output of the scalar-products between a target hidden state vector and source hidden state vector. The transformer network [16] uses a similar attention architecture as in [11]. Transformer performs a scaled scalar-product between the keys and values and then applies a softmax activation function. Our proposed attention routing computes the scalar product between capsules and a kernel.

## 3. Proposed Method

Our proposed architecture consists of *primary caps layer*, *conv caps layer*, and *fully conv caps layer*. We denote the  $l$ th capsule layer as  $\mathbf{u}_{w,h,d,n}^l$ , where  $w, h, d$ , and  $n$  index the spatial width axis, spatial height axis, capsule dimension axis, and capsule channel axis, respectively. We refer to the capsules with the same capsule channel index as a *capsule channel*  $\mathbf{u}_{(:, :, :, n)}^l$ <sup>1</sup>.

### 3.1. Primary Caps Layer

We denote the primary capsule layer as the 0th capsule layer. Before entering the primary caps layer, we extract local features  $\tilde{x}$  from the input image  $x$  by performing the convolution blocks composed of convolution layer and batch normalization.[7] We consider the local features  $\tilde{x}$  as a single capsule layer. In our primary caps layer with  $N$  channels of  $D$  dimensional output capsules,  $3 \times 3$  convolution with kernels of filter size  $D$  and stride 2 is performed on the input capsules  $\tilde{x}$   $N$  times independently. Each output of a convolution layer is a capsule channel.

$$\mathbf{u}_{(:, :, :, n_0)}^0 = \text{ReLU}(\text{Conv}_{3 \times 3}(\tilde{x})) \quad (3)$$

Note that this is equivalent to performing a  $3 \times 3$  convolution of  $N \times D$  kernels and then reshaping the features

$$\mathbf{u}_{(:, :, :, n_0)}^l := \{u_{w,h,d,n}^l | n = n_0\}$$

to  $(B, W, H, D, N)$  where  $B$  denotes the batch size and  $(W, H)$  denotes the spatial size of the capsule layer. Then, the capsule activation is applied to each capsule channel instead of the squash activation function in [12].

### 3.2. Capsule Activation

The capsule activation takes an affine transformation on each capsule channel and then applies tanh activation function. The capsules on the same capsule channel share parameters of the affine transformation. Thus, the capsule activation is equivalent to taking 1x1 convolution with a kernel of filter size  $D$  and tanh activation function on each capsule channel.

$$\mathbf{u}_{(:, :, :, n_0)} = \tanh \left( \text{Conv}_{1 \times 1} \left( \mathbf{s}_{(:, :, :, n_0)} \right) \right) \quad (4)$$

Each element of the output capsules of the capsule activation depends on the corresponding input capsule. Therefore, the capsule activation is a capsule-wise activation function. The tanh activation function normalizes each element of capsules, thus stabilizes the lengths of the capsules.

### 3.3. Conv Caps Layer

We denote the input to the  $l$ th conv caps layer as  $\mathbf{u}_{w,h,d,n}^{l-1}$  which is the output of the  $(l-1)$ th conv caps layer. We first perform a *convolutional transform* on each capsule channel. The convolutional transform is a locally-connected affine transformation sharing parameters within the same capsule channel. In particular, the convolutional transform is a 3x3 convolution of  $D^l$  kernels without activation function, where  $D^l$  denotes the capsule dimension of the  $l$ th conv caps layer.

$$\tilde{\mathbf{s}}_{(:, :, :, m)}^{l,n} = \text{Conv}_{3 \times 3} \left( \mathbf{u}_{(:, :, :, m)}^{l-1} \right) \quad (5)$$

Each output of the convolutional transform is fed to the attention routing. The output of attention routing is a linear combination of the convolutional transformed capsules with the same spatial location.

$$\mathbf{s}_{(w,h,:,n)}^l = \sum_{m=1, \dots, N^{l-1}} c_{(w,h,m)}^{l,n} \cdot \tilde{\mathbf{s}}_{(w,h,:,m)}^{l,n} \quad (6)$$

where the capsules  $\mathbf{s}_{(w,h,:,n)}^l, \tilde{\mathbf{s}}_{(w,h,:,m)}^l \in \mathbb{R}^{D^l}$ . The weights  $c_{(w,h,m)}^{l,n} \in \mathbb{R}$  are computed by the attention routing. The log probabilities  $b_{(w,h,m)}^{l,n}$  are the scalar-product between a concatenation of capsules  $[\tilde{\mathbf{u}}_{w,h,:,1}^l, \tilde{\mathbf{u}}_{w,h,:,2}^l, \dots, \tilde{\mathbf{u}}_{w,h,:,N^{l-1}}^l]$  and a parameter vector  $w_n^l \in \mathbb{R}^{D^l \times N^{l-1}}$ . This operation can be implemented efficiently by 3D convolution on the convolutional transformed capsule layers with kernels  $w_n^l \in \mathbb{R}^{1 \times 1 \times D^l \times N^{l-1}}$ , stride=(1,1,1), and valid padding.

$$b_{(:, :, :,)}^{l,n} = \text{Conv}_{3D_{1 \times 1 \times D^l}} \left( \tilde{\mathbf{s}}_{(:, :, :,)}^l \right) \quad (7)$$

where  $[\tilde{\mathbf{u}}_{(:, :, :, 1)}^l, \dots, \tilde{\mathbf{u}}_{(:, :, :, N^{l-1})}^l]$  denotes a concatenation of capsule channels. The weights  $c_{(w,h,m)}^{l,n}$  are softmax outputs of the log probabilities  $b_{(w,h,m)}^{l,n}$  along the capsule channel axis.

$$c_{w,h,m}^{l,n} = \frac{\exp(b_{(w,h,m)}^{l,n})}{\sum_{1 \leq m \leq N^{l-1}} \exp(b_{(w,h,m)}^{l,n})} \quad (8)$$

Note that the attention routing adjusts the weight  $c_{w,h,m}^{l,n}$  for each spatial location  $(w, h)$  corresponding to the convolutional transformed capsules  $\{\tilde{\mathbf{u}}_{(w,h,:,m)}^l\}_m$  with the same spatial location.

Finally, the capsule activation is performed on each capsule channel  $\mathbf{s}_{(:, :, :, n)}^l$ . A set of convolutional transform, attention routing, and capsule activation is performed independently  $N^l$  times. (i.e., each output of the convolutional transform, attention routing, and capsule activation is a capsule channel  $\mathbf{u}_{(:, :, :, n)}^l$ )

$$\mathbf{u}_{(:, :, :, n)}^l = \tanh \left( \text{Conv}_{1 \times 1} \left( \mathbf{s}_{(:, :, :, n)}^l \right) \right) \quad (9)$$

Intuitively, the dynamic routing uses a centroid of the transformed capsules as the reference vector to measure agreement by scalar-product. As the dynamic routing process iterates, the capsule with the higher agreement has a larger weight, and the reference vector evolves in that capsule direction. On the other hand, since the capsule activation in the conv caps layer do not preserve vector orientation, the output capsule  $\mathbf{u}_{(:, :, :, n)}^l$  cannot approximate the centroid of transformed capsules  $\{\tilde{\mathbf{u}}_{(w,h,:,n)}^l\}$ . Instead of measuring agreement between the transformed capsules and the output capsule  $\mathbf{u}_{(:, :, :, n)}^l$ , the attention routing parametrizes the routing process. The parameter vector  $w_n^l$  which is the kernel of convolution serves as an approximation of the reference vector to perform routing.

We propose replacing the dynamic routing of [12] with the convolutional transform and attention routing. Compared to dynamic routing, our proposed operation is faster and more parameter efficient. Since dynamic routing is constructed in a fully connected manner, the transform weight matrices are assigned for each pair of the input capsule and output capsule. We share the weight matrices across the spatial location and keep the translation equivariance by performing 3x3 convolution on the  $l$ th layer in the convolutional transform. (Section 4.4) Besides, the dynamic routing has an iterative routing process to compute the weight  $c_{w,h,m}^{l,n}$ . On the other hand, by introducing a trainable parameter vector, our proposed operation is a fast forward-pass.

---

**Algorithm 1: The process of Attention Routing**


---

```

1 Input:  $\mathbf{u}^{\ell=0} \in \mathbb{R}^{(w,h,D^{\ell=0},N^{\ell=0})}$ 
2 for  $\ell = 1, \dots, L$  do
3   for  $n = 1, \dots, N^\ell$  do
4     /* Convolutional transformation
       for each capsule channel */
5     for  $m = 1, \dots, N^{\ell-1}$  do
6        $\tilde{\mathbf{s}}_{(:, :, :, m)}^{\ell, n} \leftarrow \text{Conv2D}_{3 \times 3}(\mathbf{u}_{(:, :, :, m)}^{\ell-1})$ 
7     end
8     /* Attention through capsule
       channel */
9      $\mathbf{b}^{l,n} \leftarrow \text{Conv3D}_{1 \times 1 \times D^\ell}(\tilde{\mathbf{s}}^\ell)$ 
10    for  $w, h = 1, \dots, W^\ell, H^\ell$  do
11       $\mathbf{c}_{w,h,:N^{\ell-1}}^{n,l} \leftarrow \text{softmax}(\mathbf{b}_{w,h,:N^{\ell-1}}^{\ell,n})$ 
12       $\mathbf{s}_{w,h,:n}^\ell \leftarrow \sum_{m=1, \dots, N^{\ell-1}} \mathbf{c}_{w,h,m}^{n,l} \cdot \tilde{\mathbf{s}}_{(w,h,:m)}^{\ell,n}$ 
13    end
14    /* Capsule activation for each
       capsule channel */
15     $\mathbf{u}_{(:, :, :, n)}^\ell \leftarrow \tanh(\text{Conv2D}_{1 \times 1}(\mathbf{s}_{(:, :, :, n)}^\ell))$ 
16  end
17 end

```

---

### 3.4. Fully Conv Caps Layer

The fully conv caps layer is almost the same as the conv caps layer and serves as the output layer of AR CapsNet. The convolutional transform combines capsule features from the all spatial location by applying a kernel of the same spatial size as the input with valid padding. Therefore, the output of the fully conv caps Layer has a shape of  $(1, 1, D^L, N^L)$ .

### 3.5. Margin Loss and Reconstruction Regularizer

We adopt the margin loss and reconstruction regularizer in [12]. Since the output capsules of capsule activation have a length of up to  $\sqrt{D^L}$  where  $D^L$  denotes the capsule dimension, we use the normalized length to predict the probability of the corresponding class of the dataset.

$$\|\mathbf{u}_n^L\|_{\text{nor}} = \frac{\|\mathbf{u}_n^L\|}{\sqrt{D^L}} \quad (10)$$

where  $\|\mathbf{u}_n^L\|$  denotes the output capsules of the fully conv caps layer and  $n$  indexes the capsule channel axis. We ap-

plied the Margin loss,  $L_n$ , for each class  $n$  on the  $\|\mathbf{u}_n^L\|_{\text{nor}}$ .

$$L_n = T_n \max(0, m^+ - \|\mathbf{u}_n^L\|_{\text{nor}})^2 + \lambda(1 - T_n) \max(0, \|\mathbf{u}_n^L\|_{\text{nor}} - m^-)^2 \quad (11)$$

where  $T_n = 1$  iff the corresponding class of output capsule is present and  $m^+ = 0.9$  and  $m^- = 0.1$ .

The output capsules  $\{\mathbf{u}_n^L\}_{n=1, \dots, N}$  are fed to the reconstruction decoder. We used a decoder consisting of 3 fully connected layers as in [12] except that our decoder has (512, 512, the number of input image pixel) nodes. We refer to the mean of L2 loss between an input image and the decoder output as a reconstruction loss. We add the reconstruction loss that is scaled down by 0.3 to the margin loss as a regularization method.<sup>2</sup>

## 4. Experiments

We evaluate our model on the MNIST, affNIST, and CIFAR-10 datasets. For each dataset, we split the training images into a training set (90%) and a validation set (10%). We choose the model with the lowest validation error and evaluate the model on the test set. Then, we compare the results with CapsuleNet [12]. We use a Keras implementation<sup>3</sup> for CapsuleNet.

Before training the model on the image dataset, we divide each pixel value by 255 so that it is scaled in the range of 0 to 1. Then, we extract the local features  $\tilde{x}$  from an input image through two convolutional layers with batch normalization(BN) [7] and ReLU activation function. These two convolutional layers use 3x3 kernels with a stride 1. Then, the features go through the AR CapsNet to obtain vector outputs. For each conv caps layer and fully conv caps layer, the dropout layer [14] of keep probability 0.5 is applied to the input capsules before the convolutional transform.

We use the RMSprop optimizer with rho of 0.9 and decay of 1e-4 to minimize the loss defined in Section 3.5. We set the learning rate as 0.001 and batch size as 100

### 4.1. Classification Results on MNIST and affNIST

**Dataset** The MNIST dataset is composed of  $28 \times 28$  hand-written digit images. We adopted 0.1 translation as a data augmentation for the MNIST dataset. The affNIST dataset consists of  $40 \times 40$  images, which are obtained by applying various affine transformations such as rotation and expansion to the images from MNIST. For the affNIST classification task, we trained our model with randomly translated MNIST images in horizontal or vertical directions up to shift fraction 0.2 as in [12]. Any other affine transformations like rotations were not used in the training process.

<sup>2</sup>CapsuleNet [12] scaled the reconstruction loss by 0.392. Since we use the mean of L2 loss and CapsuleNet use the sum of L2 loss,  $0.392 = 0.0005 \times 784$ .

<sup>3</sup><https://github.com/XifengGuo/CapsNet-Keras>



Method	MNIST	MNIST+	affNIST	C10	C10+
CapsuleNet[12]	99.45*	99.75 (99.52*)	79.0	63.1*	69.6*
CapsuleNet+ensemble(7)	-	-	-	-	89.4
Ours	99.46	99.46	91.6	87.19	88.61
Ours+ensemble(7)	-	-	-	88.94	90.11

Table 1. Test accuracy (%) on the MNIST, affNIST, and CIFAR-10 classification tasks. C10 represents the CIFAR-10 dataset. + denotes training with data augmentation. We adopted translation for MNIST+ and translation, rotation, and horizontal flip for C10+. \* indicates the results from our experiment.

The affNIST dataset has a separate validation set, thus we chose the model with the lowest validation error based on the affNIST validation set. Then, we tested our model with the affNIST test set.

**Implementation** For the MNIST and affNIST datasets, we used the AR CapsNet which consists of a primary caps layer, one conv caps layer and fully conv caps layer. Before entering the AR CapsNet, an input image goes through two convolutional layers of 64 channels (3x3 Conv - BN - ReLU). The primary caps layer has eight channels of 16-dimensional capsules, the conv caps has eight channels, and the fully conv caps layer has ten channels. Each capsule channels in the conv caps layer and fully conv caps layer has 32 dimensions in the MNIST and 16 dimensions in the affNIST. We decreased the spatial size of the capsule features by applying a 3x3 convolution of stride 2 in the convolutional transform of the conv caps layer. We trained our model for 20 epochs.

**Accuracy** Our model shows a comparable accuracy with the substantial decrease in the number of parameters and training time. Our model with 5.31M parameters achieved 99.45% accuracy on the MNIST dataset without any data augmentation and 99.46% accuracy with data augmentation. (Table 1) The CapsuleNet with 8.21M parameters achieved 99.45% accuracy without any data augmentation and 99.52% with data augmentation. The reported accuracy of CapsuleNet on the MNIST dataset with translation augmentation is 99.75% [12]. Also, the training took 37.2 seconds per epoch for our proposed model and 199.5 seconds per epoch for CapsuleNet when we experimented on GTX 1080 GPUs.

In the affNIST experiments, there are two options to generate training images from the MNIST dataset. The first option is to create a larger dataset by generating a set of all the possible augmented data before training. The second option is to apply translation over the original dataset for each epoch. The reported accuracy of CapsuleNet is 79% and that of the baseline CNN model is 66% in [12]. The experiment is performed on the former option.<sup>4</sup> Our proposed model achieved 91.6% accuracy for the latter option. Under the comparable experiment, our model outper-

formed the CapsuleNet and the baseline CNN model. Since our model is transformation equivariant (Section 4.4), our model is robust to affine transformations.

## 4.2. Classification Results on CIFAR-10

**Dataset** The CIFAR-10 dataset is a  $32 \times 32$  colored natural images in 10 classes. We adopted 0.1 translation, rotation up to 20 degrees, and horizontal flip as a data augmentation for CIFAR-10.

**Implementation** For the CIFAR-10 classification task, we added four conv caps layer between a primary caps layer and fully conv caps layer. We decreased the spatial size of the capsule features in the first conv caps layer as in Section 4.1. Each conv caps layer has eight channels of 32-dimensional capsules and is connected to the next conv caps layer with a residual connection [3]. Note that the residual connection in [3] connects the  $l$ th layer and  $(l + 2)$ th layer, but our residual connection connects the  $l$ th conv caps layer and  $(l + 1)$  conv caps layer. We trained our model for 200 epochs.

**Accuracy** The results in Table 1 show that our model outperforms CapsuleNet with and without data augmentation. CapsuleNet with 11.74M parameters shows 63.1% accuracy in C10 and 69.6% accuracy in C10+. However, our proposed model with 9.6M parameters shows 87.19% accuracy in C10 and 88.61% accuracy in C10+. In [12], an ensemble of 7 models achieves 89.4% accuracy when the models are trained with  $24 \times 24$  patches of images and the introduction of a *none-of-the-above* category. However, an ensemble of 7 AR CapsNet models trained with C10+ achieved 90.11% test accuracy. Note that C10+ only uses rotation, shift, and horizontal flip as data augmentation and not the cropping or the *none-of-the-above* category.

## 4.3. Robustness to hyperparameters

**Implementation** The AR CapsNet requires a set of hyperparameters, such as the number of conv caps layer and the capsule dimension of each capsule layer. To test the robustness to hyperparameters, we evaluate the AR CapsNet in the CIFAR-10 classification tasks according to the various setting of hyperparameters. The evaluated AR CapsNet architecture is the same as the models mentioned in Sec-

<sup>4</sup><https://github.com/Sarasra/models/tree/master/research/capsules>

Conv caps layer	Caps dim	Params	C10	C10+
0	16	7.3M	77.51	81.89
	32	12.6M	77.44	81.97
1	16	3.5M	81.96	82.83
	32	7.7M	82.39	83.92
2	16	3.7M	84.48	84.77
	32	8.4M	85.46	87.01
3	16	3.8M	85.56	86.93
	32	8.9M	86.56	87.91
4	16	4.0M	86.37	87.21
	32	9.6M	87.19	88.61

Table 2. Test accuracy (%) on the MNIST and CIFAR-10 for various hyperparameters. In each experiment, we trained a model for 200 epochs and chose the model with the lowest validation error. For each hyperparameter setting, AR CapsNet shows stable performance without showing severe degradation.

tion 4.1 and 4.2. The primary caps layer has eight channels of 16-dimensional capsules, and the conv caps layer has eight capsule channels. In the setting of hyperparameters, the *Conv caps layer* denotes the number of conv caps layer between the primary caps layer and fully conv caps layer. In every model with at least one conv caps layer, the first conv caps layer decreases the spatial size of the capsule layer by adopting a 3x3 convolution of stride 2 in the convolutional transform. The *Capsule dim* denotes the capsule dimension of the conv caps layer and fully conv caps layer.

**Robustness** All the AR CapsNet models trained with CIFAR-10 dataset show decent performance in Table 2. Increasing capsule dimension and the number of conv caps layer lead to an improvement in the test accuracy. The AR CapsNet model with four conv caps layer shows 86.37% accuracy with 16-dimensional capsules and 87.19% accuracy with 32-dimensional capsules. The AR CapsNet with four conv caps layer and 32-dimensional capsules shows the best results of 87.19% in C10 and 88.61% in C10+. Also, the AR CapsNet model with no conv caps layer has more parameters than the model with four conv cap layer and shows the worst performance. The features in the primary caps layer has a large spatial size. Thus, the fully conv caps layer connected to the primary caps layer assigns excessive parameters, and this causes overfitting.

#### 4.4. Transformation Equivariance

**Dimension perturbation** To prove that our proposed model is transformation equivariant, we executed experiments on the MNIST model as in [12]. We observed the variations of decoder output as we perturbed one scalar element of the output capsules (Figure 3). The experiments in the [12] perturbed one scalar element from -0.25 to 0.25. Since the output capsules of the AR CapsNet have lengths

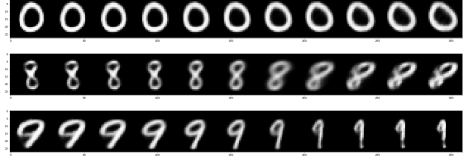


Figure 3. Decoder outputs according to dimension perturbations. We observed the variations of decoder output as we perturbed one dimension of the output capsules by steps of  $0.05\sqrt{D^L}$  from  $-0.25\sqrt{D^L}$  to  $+0.25\sqrt{D^L}$ . The perturbation leads to the combination of variations in the decoder output images. (e.g., rotation, thickness, etc.).

Digit	Rot+	x+	y+	Rot-	x-	y-
8	0.89	0.89	0.91	0.86	0.87	0.86
5	0.89	0.78	0.73	0.84	0.88	0.86
avg	0.88 (0.89)	0.86 (0.88)	0.84 (0.80)	0.83 (0.81)	0.83 (0.84)	0.84 (0.84)

Table 3. The average of relative ratio  $\{r_i\}$  for each combination of digit and transformation. The avg represents the average of all 10,000 test samples for each affine transformation. We report the results of models trained on the MNIST+ dataset in the (). A high relative ratio implies the difference vectors are strongly aligned. For random vectors, the average of relative ratio  $\{r_i\}$  is 0.311 and standard deviation is 0.262.

of up to  $\sqrt{D^L}$  compared to 1 in [12], we perturbed one scalar element from  $-0.25\sqrt{D^L}$  to  $0.25\sqrt{D^L}$  where  $D^L$  denotes the capsule dimension of output capsules. Figure 3 shows that some dimensions of the output capsules represent variations in the way the digit of the corresponding class is instantiated. Some dimensions of the output capsules represent the localized skew in digit 0, the rotation and the size of the higher circle in digit 8, and the rotation, thickness, and skew in digit 9.

**Alignment ratio** Each scalar element of the output capsules represents a combination of variations such as rotation, thickness, and skew. (Digit 9 in Figure 3) Since the length of the output capsules is basis-invariant, the transformation on an input image could be represented in coordinates of any basis. To further test the transformation equivariance of the AR CapsNet, we tested whether the difference in the output capsules caused by applying a transformation on an input image is aligned in one direction.

Let  $\{T_i\}_{i=1,\dots,N}$  be a set of affine transformations on an input image  $x$ . We denote the difference between the output capsules  $\mathbf{u}_n^L(T_i(x))$  and  $\mathbf{u}_n^L(x)$  as  $\mathbf{v}_i(x)$  where  $n$  denotes the corresponding class of  $x$ .

$$\mathbf{v}_i(x) = \mathbf{u}_n^L(T_i(x)) - \mathbf{u}_n^L(x) \quad (12)$$

We denote the concatenation of  $\mathbf{v}_i(x)$  along the row axis as  $\mathbf{V}$ . In order to obtain a representative unit vector  $\tilde{v}$  of

$\{\mathbf{v}_i(x)\}$ , we apply a Singular-Value Decomposition(SVD) on matrix  $\mathbf{V}$ .

$$\mathbf{c}, \tilde{v} = \arg \min_{c_i, \tilde{v}} \sum_i \|\mathbf{v}_i(x) - c_i \cdot \tilde{v}\|_2^2 \quad (13)$$

$$= \arg \min_{c_i, \tilde{v}} \|\mathbf{V} - \mathbf{c} \cdot \tilde{v}^T\|_F^2 \quad (14)$$

where  $F$  denotes the Frobenius norm,  $\mathbf{c} = (c_1, \dots, c_N)^T$ , and  $c_i \in \mathbb{R}$ . The exact solution of this low rank approximation problem is the first right-singular vector  $\tilde{v}$  of  $\mathbf{V}$ . This experiment is similar to the Principal Component Analysis(PCA) except that we do not subtract the mean for each columns of  $\mathbf{V}$ . The align vector  $\tilde{v}$  corresponds to the principal vector of PCA. We observed the relative ratio  $r_i$  of principal component of  $\mathbf{v}_i(x)$  to the vector norms  $\|\mathbf{v}_i(x)\|_2$ .

$$r_i = \frac{|\mathbf{v}_i(x) \cdot \tilde{v}|}{\|\mathbf{v}_i(x)\|_2} \quad (15)$$

We randomly chose 10,000 images from the test set. For each test image, we generated five images by applying an affine transformation and observed the relative ratio  $r_i$ . In Table 3,  $\text{Rot}(\pm)$  represents  $\pm\{5, 10, 15, 20, 25\}$  degrees rotations and  $x(\pm)$  represents a horizontal translation up to  $\pm 5$  pixels.  $y(\pm)$  represents a vertical translation up to  $\pm 5$  pixels as well. We observed the average of relative ratio  $r_i$  for each combination of digit and transformation. Table 3 shows the average of relative ratio  $r_i$  for two digits (highest : digit 8, lowest : digit 5) and the average for 10,000 test samples for each transformation. As a reference, we generated five random vectors from the standard multivariate normal distribution. We conducted the same experiment for random vectors for 1,000 times as well. We obtained an average of 0.311 and a standard deviation of 0.262 for random vectors. Even for the worst-case digit 5, every transformation shows a significantly higher relative ratio  $r_i$  of 0.73 in  $y+$  than the random vectors. This result implies that the difference vectors are strongly aligned in one direction. Therefore, AR CapsNet encodes affine transformations on an input image by some vector components. Also, we report the results of models trained on the MNIST+ dataset in the (). The models trained on the MNIST+ show comparable relative ratio  $r_i$  to those trained on the MNIST. This result shows that AR CapsNet encodes affine transformations even without observing transformations during training.

An interesting observation is that AR CapsNet is transformation equivariant but do not distinguish the positive and negative transformations. Figure 4 shows the histogram of the cosine similarity between the align vectors of positive and negative transformation.<sup>5</sup> We observed two peaks around -1 and 1. The cosine similarity of -1 and 1 imply that

<sup>5</sup>The direction of the first right-singular vector  $\tilde{v}$  is given by  $\mathbf{v}_i(x) \cdot \tilde{v} > 0$  in for each positive and negative transformation.

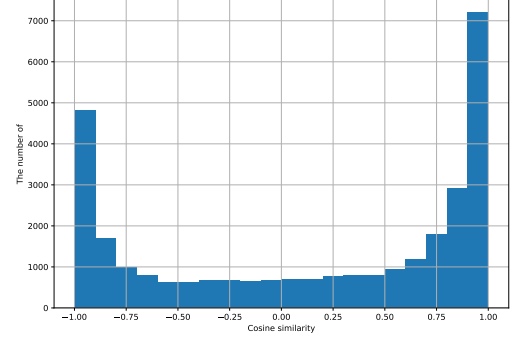


Figure 4. Distribution of cosine similarity between unit align vectors  $\tilde{v}$  of positive and negative affine transformations.

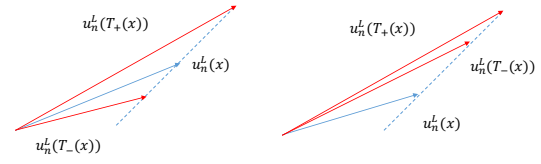


Figure 5. Output capsules  $\mathbf{u}_n^L$  when the cosine similarity between align vectors of positive and negative transformations is -1 or 1.  $T_+$  and  $T_-$  represent the positive and negative transformation. **Left:** cosine similarity -1. **Right:** cosine similarity 1.

positive and negative transformations are encoded in one dimension. However, the cosine similarity of 1 suggests that the difference vectors of positive and negative transformations have the same direction.(Figure 5) We leave the explanation of this observation to future work.

## 5. Conclusion

In this work, we suggested a new capsule network architecture called Attention Routing CapsuleNet (AR CapsNet). By introducing the attention routing and capsule activation, AR CapsNet obtained a higher accuracy compared to CapsuleNet while using fewer parameters and less training time. The attention routing is an effective way to route between capsules because it only compares capsules of the same spatial location. In addition, the attention routing does not require an iterative routing process as the dynamic routing does because it directly learns the weights between capsules. The capsule activation is based on the assumption that the capsule-scale activation can extract transformation equivariant features even if it is not orientation-preserving. This assumption distinguish the capsule activation from the squash activation function and its variant.

While using the building blocks of CNNs, AR CapsNet is transformation equivariant. We showed that capsules have transformation information by manipulating the output capsules and then observing the decoder output images. Also, we observed the difference vectors between the



output capsules of an original image and an affine transformed image. By showing that the difference vectors are strongly aligned in one direction, we proved that AR CapsNet encodes transformation information in some dimensions. There are natural variations of AR CapsNet such as introducing a feature compression by 1x1 convolution to the capsule activation and a transformer network [16] to the attention routing. We plan to study these variations in the future.

## References

- [1] M. T. Bahadori. Spectral capsule networks. 2018. [3](#)
- [2] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014. [3](#)
- [3] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#), [6](#)
- [4] G. E. Hinton, A. Krizhevsky, and S. D. Wang. Transforming auto-encoders. In *International Conference on Artificial Neural Networks*, pages 44–51. Springer, 2011. [1](#)
- [5] G. E. Hinton, S. Sabour, and N. Frosst. Matrix capsules with em routing. 2018. [3](#)
- [6] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017. [1](#)
- [7] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. [3](#), [5](#)
- [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. [1](#)
- [9] R. LaLonde and U. Bagci. Capsules for object segmentation. *arXiv preprint arXiv:1804.04241*, 2018. [3](#)
- [10] H. Li, X. Guo, B. DaiWanli Ouyang, and X. Wang. Neural network encapsulation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 252–267, 2018. [3](#)
- [11] M.-T. Luong, H. Pham, and C. D. Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015. [3](#)
- [12] S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems*, pages 3856–3866, 2017. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [13] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [1](#)
- [14] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014. [5](#)
- [15] R. K. Srivastava, K. Greff, and J. Schmidhuber. Training very deep networks. In *Advances in neural information processing systems*, pages 2377–2385, 2015. [1](#)
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017. [3](#), [9](#)
- [17] D. Wang and Q. Liu. An optimization view on dynamic routing between capsules. 2018. [3](#)
- [18] E. Xi, S. Bing, and Y. Jin. Capsule network performance on complex data. *arXiv preprint arXiv:1712.03480*, 2017. [1](#), [2](#)