

Dynamic Routing Between Capsules

by S. Sabour, N. Frosst and G. Hinton (NIPS 2017)

presented by Karel Ha

27th March 2018

Pattern Recognition and Computer Vision Reading Group

Outline

Motivation

Capsule


Routing by an Agreement

Capsule Network

Experiments


Conclusion

Motivation



The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster.

G. Hinton

A portrait of Geoffrey Hinton, an older man with grey hair and blue eyes, wearing a dark V-neck sweater over a light-colored collared shirt. He is standing outdoors with green trees and a blue sky in the background.

The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster.

G. Hinton

[...] it makes much more sense to represent a pose as a small matrix that converts a vector of positional coordinates relative to the viewer into positional coordinates relative to the shape itself.

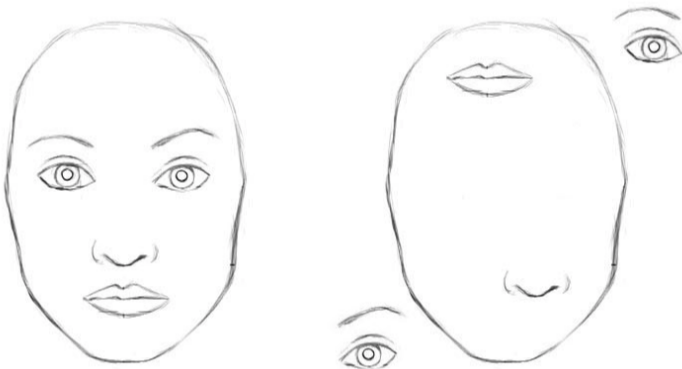
G. Hinton

Part-Whole Geometric Relationships

"What is wrong with convolutional neural nets?"

Part-Whole Geometric Relationships

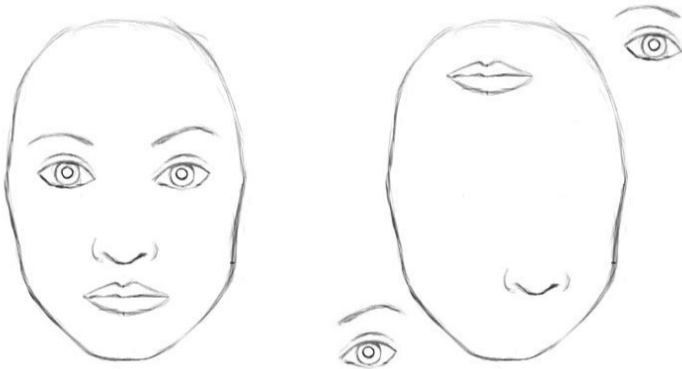
"What is wrong with convolutional neural nets?"



To a CNN (with MaxPool)...

Part-Whole Geometric Relationships

"What is wrong with convolutional neural nets?"

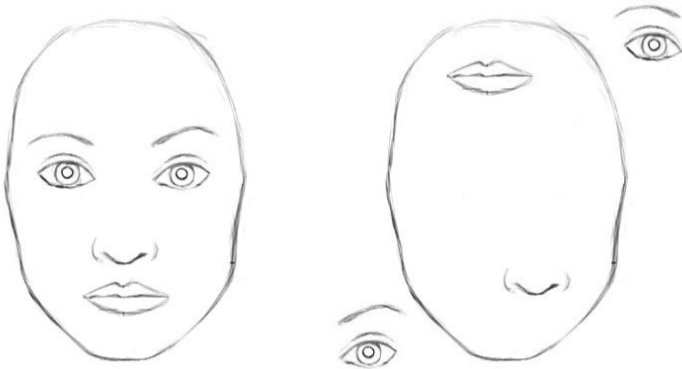


To a CNN (with MaxPool)...

■ ...both pictures are similar, since they both contain similar elements.

Part-Whole Geometric Relationships

"What is wrong with convolutional neural nets?"

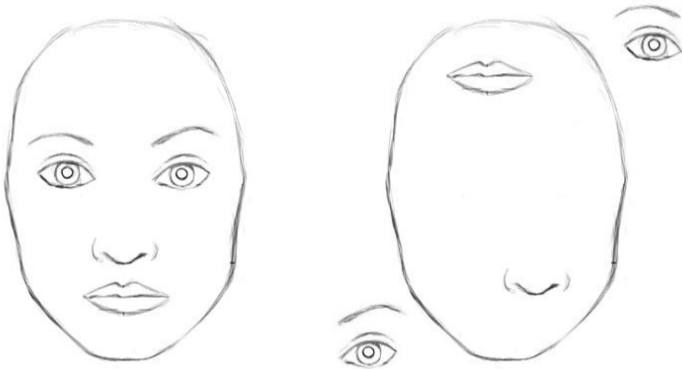


To a CNN (with MaxPool)...

- ...both pictures are similar, since they both contain similar elements.
- ...a mere presence of objects can be a very **strong indicator** to consider a face in the image.

Part-Whole Geometric Relationships

"What is wrong with convolutional neural nets?"



To a CNN (with MaxPool)...

- ...both pictures are similar, since they both contain similar elements.
- ...a mere presence of objects can be a very **strong indicator** to consider a face in the image.
- ...orientational and relative spatial relationships are not very important.

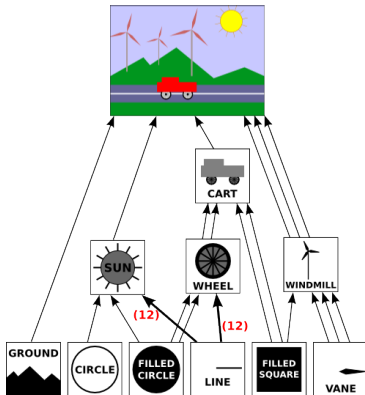
Part-Whole Geometric Relationships

Scene Graphs from Computer Graphics

Part-Whole Geometric Relationships

Scene Graphs from Computer Graphics

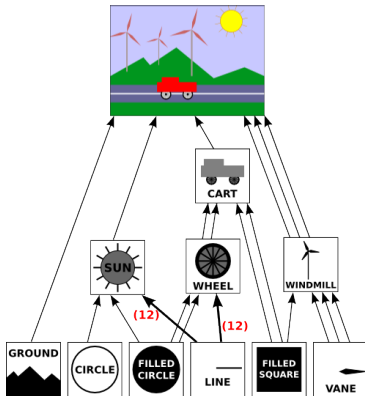
...takes into account relative positions of objects.



Part-Whole Geometric Relationships

Scene Graphs from Computer Graphics

...takes into account relative positions of objects.

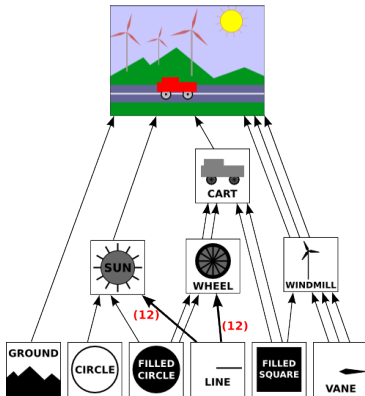


The internal representation in computer memory:

Part-Whole Geometric Relationships

Scene Graphs from Computer Graphics

...takes into account relative positions of objects.



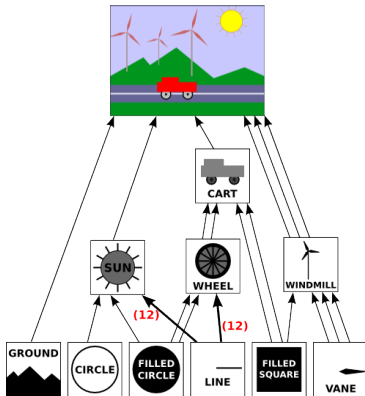
The internal representation in computer memory:

a) arrays of geometrical objects

Part-Whole Geometric Relationships

Scene Graphs from Computer Graphics

...takes into account relative positions of objects.

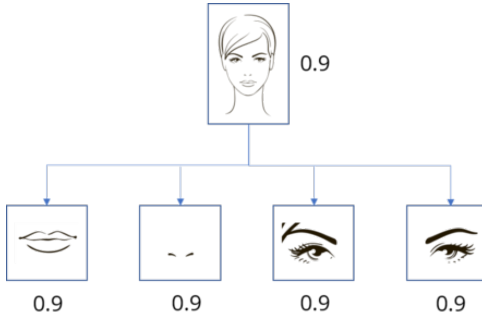


The internal representation in computer memory:

- a) arrays of geometrical objects
- b) matrices representing their **relative positions** and **orientations**

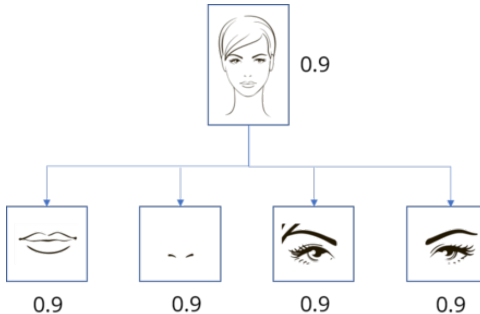
Part-Whole Geometric Relationships

Inverse (Computer) Graphics



Part-Whole Geometric Relationships

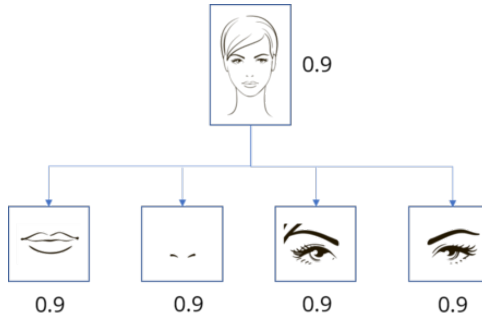
Inverse (Computer) Graphics



Inverse graphics:

Part-Whole Geometric Relationships

Inverse (Computer) Graphics

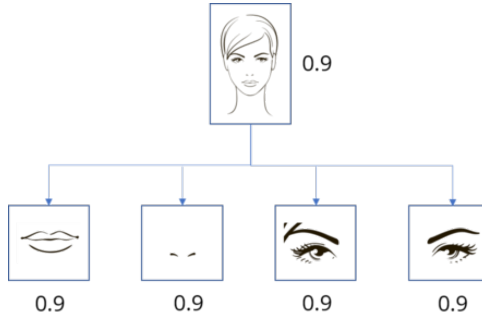


Inverse graphics:

- from visual information received by eyes

Part-Whole Geometric Relationships

Inverse (Computer) Graphics

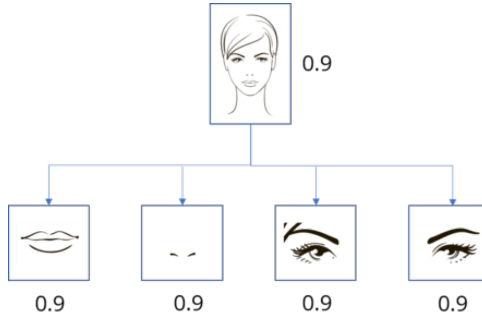


Inverse graphics:

- from visual information received by eyes
- deconstruct a hierarchical representation of the world around us

Part-Whole Geometric Relationships

Inverse (Computer) Graphics

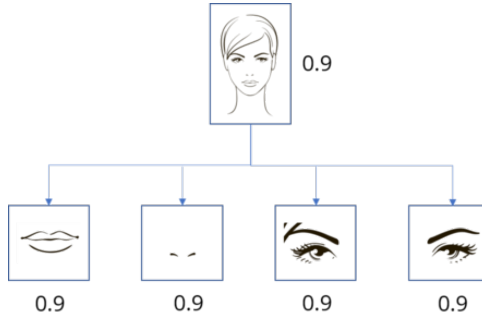


Inverse graphics:

- from visual information received by eyes
- deconstruct a hierarchical representation of the world around us
- and try to match it with already learned patterns and relationships stored in the brain

Part-Whole Geometric Relationships

Inverse (Computer) Graphics



Inverse graphics:

- from visual information received by eyes
- deconstruct a hierarchical representation of the world around us
- and try to match it with already learned patterns and relationships stored in the brain
- relationships between 3D objects using a “pose” (= translation plus rotation)

Pose Equivariance and the Viewing Angle

Pose Equivariance and the Viewing Angle

We have probably never seen these exact pictures, but we can still immediately recognize the object in it...



Pose Equivariance and the Viewing Angle

We have probably never seen these exact pictures, but we can still immediately recognize the object in it...



- internal representation in the brain: independent of the viewing angle

Pose Equivariance and the Viewing Angle

We have probably never seen these exact pictures, but we can still immediately recognize the object in it...



- internal representation in the brain: independent of the viewing angle
- quite hard for a CNN: no built-in understanding of 3D space

Pose Equivariance and the Viewing Angle

We have probably never seen these exact pictures, but we can still immediately recognize the object in it...

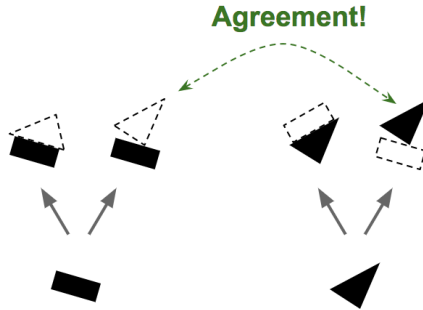


- internal representation in the brain: independent of the viewing angle
- quite hard for a CNN: no built-in understanding of 3D space
- much easier for a CapsNet: these relationships are explicitly modeled

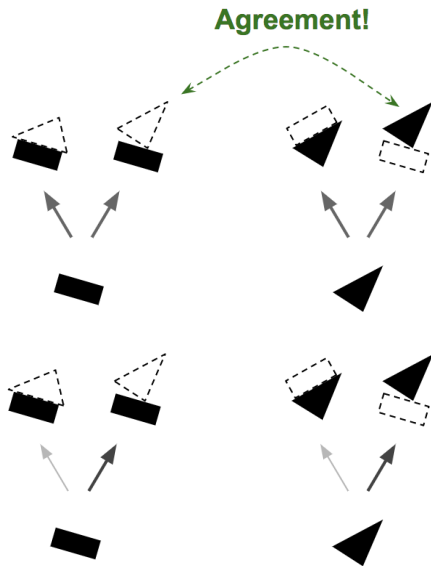
<https://medium.com/ai-theory-practice-business/understanding-hintons-capsule-networks-part-i-intuition-b4b559d1159b>

Routing by an Agreement: High-Dimensional Coincidence

Routing by an Agreement: High-Dimensional Coincidence

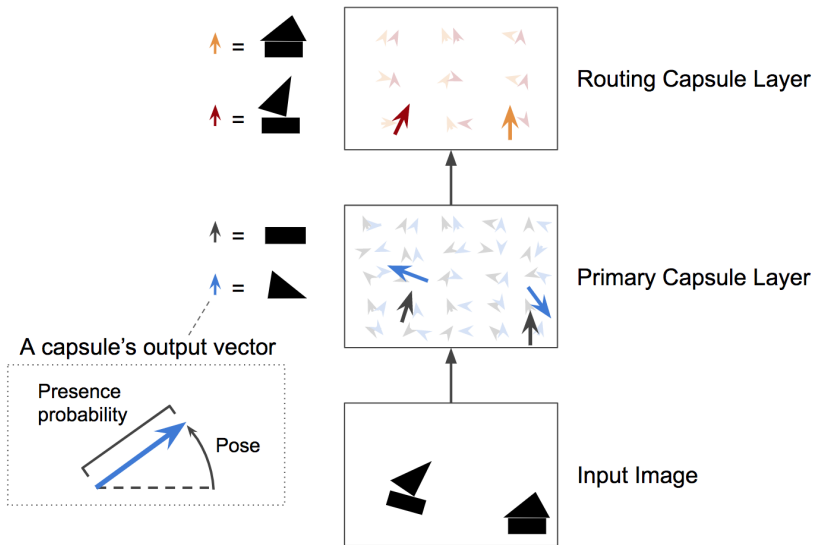


Routing by an Agreement: High-Dimensional Coincidence

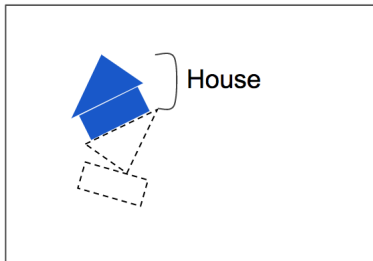
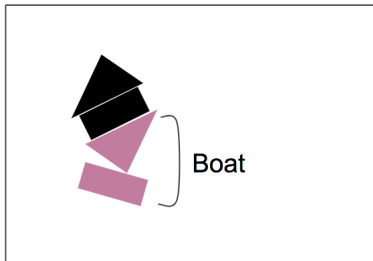
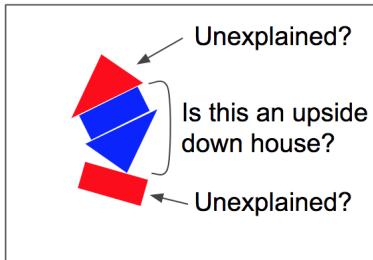


Routing by an Agreement: Illustrative Overview

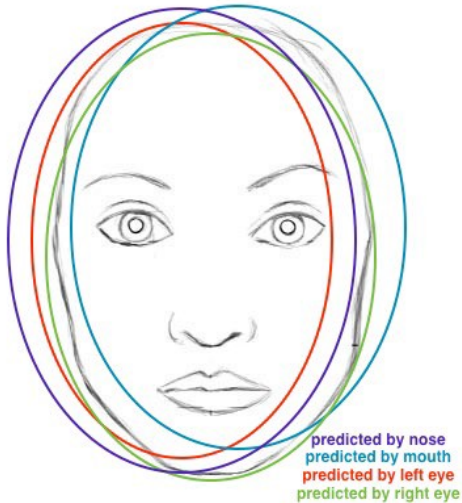
Routing by an Agreement: Illustrative Overview



Routing by an Agreement: Recognizing Ambiguity in Images



Routing: Lower Levels Voting for Higher-Level Feature



How to do it?
(mathematically)

Capsule

What Is a Capsule?

a group of neurons that:

- perform some complicated **internal computations** on their inputs

What Is a Capsule?

a group of neurons that:

- perform some complicated **internal computations** on their inputs
- encapsulate their results into a **small vector** of highly informative outputs

What Is a Capsule?

a group of neurons that:

- perform some complicated **internal computations** on their inputs
- encapsulate their results into a **small vector** of highly informative outputs
- recognize an implicitly defined visual entity (over a limited domain of viewing conditions and deformations)

What Is a Capsule?

a group of neurons that:

- perform some complicated **internal computations** on their inputs
- encapsulate their results into a **small vector** of highly informative outputs
- recognize an implicitly defined visual entity (over a limited domain of viewing conditions and deformations)
- encode the **probability** of the entity being present

What Is a Capsule?

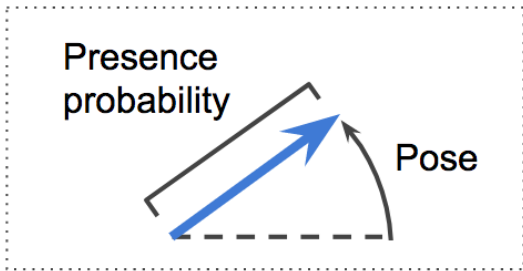
a group of neurons that:

- perform some complicated **internal computations** on their inputs
- encapsulate their results into a **small vector** of highly informative outputs
- recognize an implicitly defined visual entity (over a limited domain of viewing conditions and deformations)
- encode the **probability** of the entity being present
- **encode instantiation parameters**

pose, lighting, deformation relative to entity's (implicitly defined) canonical version

Output As A Vector

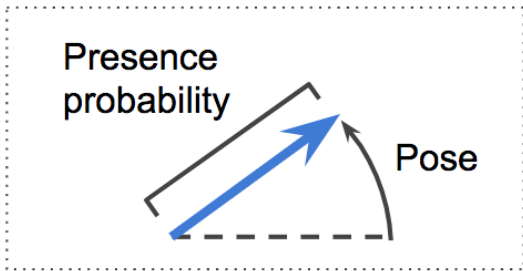
A capsule's output vector



- probability of presence: locally **invariant**

E.g. if 0, 3, 2, 0, 0 leads to 0, 1, 0, 0, then 0, 0, 3, 2, 0 should also lead to 0, 1, 0, 0.

A capsule's output vector



- probability of presence: locally **invariant**

E.g. if 0, 3, 2, 0, 0 leads to 0, 1, 0, 0, then 0, 0, 3, 2, 0 should also lead to 0, 1, 0, 0.

- instantiation parameters: **equivariant**

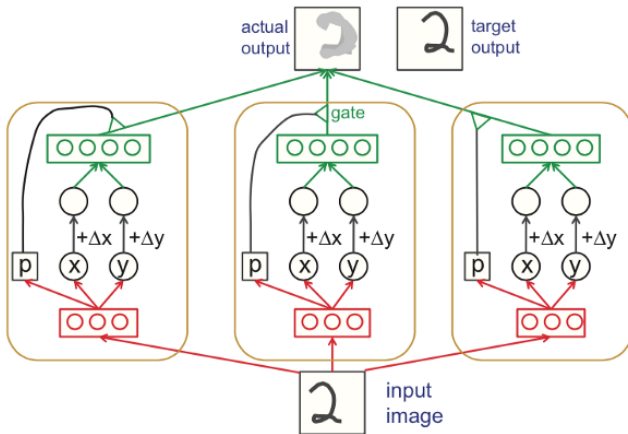
E.g. if 0, 3, 2, 0, 0 leads to 0, 1, 0, 0, then 0, 0, 3, 2, 0 might lead to 0, 0, 1, 0.

Previous Version of Capsules

for illustration taken from “Transforming Auto-Encoders” (Hinton, Krizhevsky and Wang [2011])

Previous Version of Capsules

for illustration taken from "Transforming Auto-Encoders" (Hinton, Krizhevsky and Wang [2011])



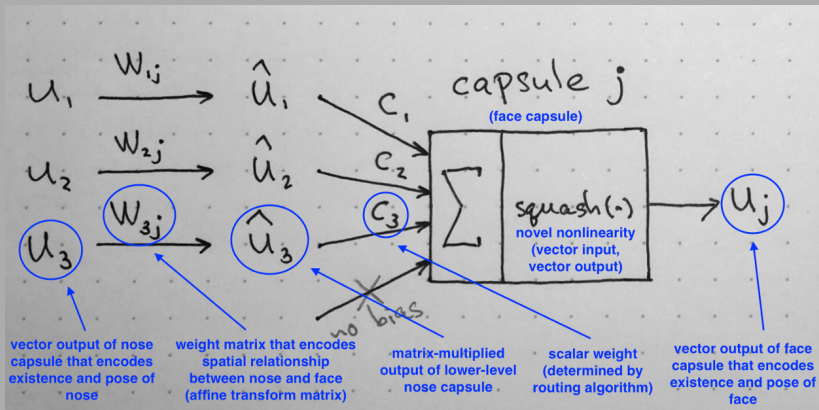
three capsules of a transforming auto-encoder (that models translation)

(Hinton, Krizhevsky and Wang [2011])

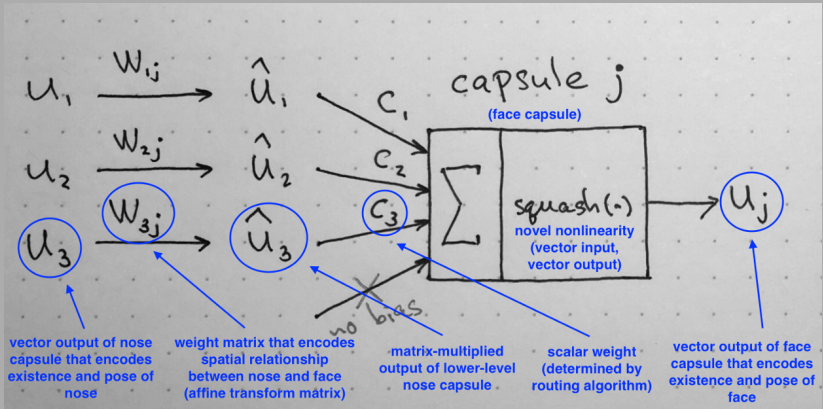
Capsule's Vector Flow

https://cdn-images-1.medium.com/max/1250/1*GbmQ2X9NQoGuJ1M-EOD67g.png

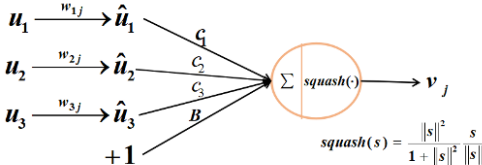
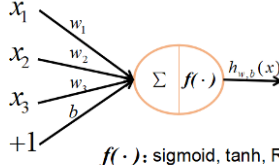
Capsule's Vector Flow



Capsule's Vector Flow



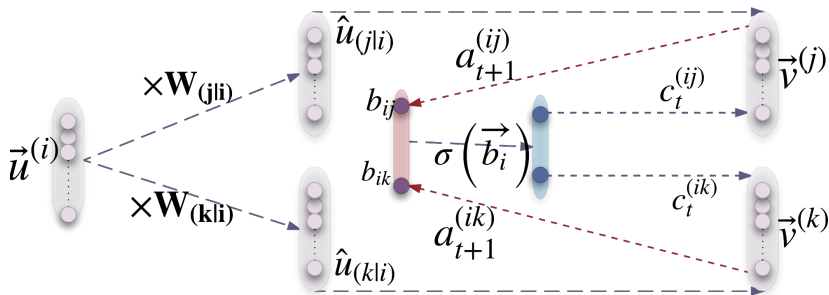
Note: no bias (included in affine transformation matrices W_{ij} 's)

capsule		VS.	traditional neuron
Input from low-level neurons/capsules		vector(u_i)	scalar(x_i)
Operations	Linear/Affine Transformation	$\hat{u}_{ji} = W_{ij} u_i + B_j$ (Eq. 2)	$a_{ji} = w_{ij} x_i + b_j$
	Weighting	$s_j = \sum_i c_{ij} \hat{u}_{ji}$ (Eq. 2)	$z_j = \sum_{i=1}^3 1 \cdot a_{ji}$
	Summation		
	Non-linearity activation	$v_j = \text{squash}(s_j)$ (Eq. 1)	$h_{w,b}(x) = f(z_j)$
output		vector(v_j)	scalar(h)
			

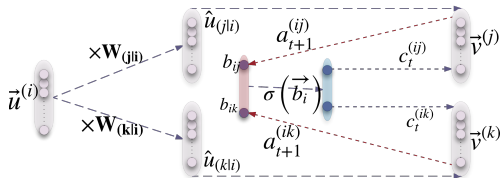
Capsule = New Version Neuron!
vector in, vector out VS. scalar in, scalar out

Routing by an Agreement

Capsule Schema with Routing

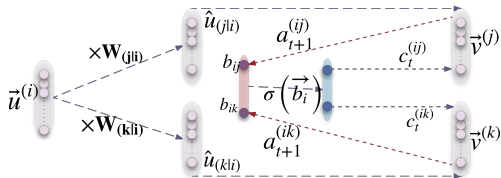


Routing Softmax



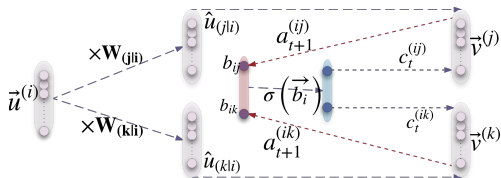
$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})} \quad (1)$$

Prediction Vectors



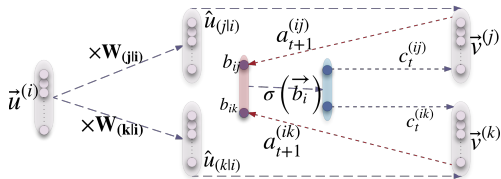
$$\hat{u}_{j|i} = \mathbf{W}_{ij} \mathbf{u}_i \quad (2)$$

Total Input



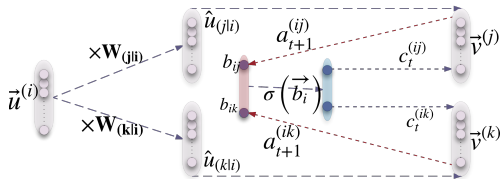
$$\mathbf{s}_j = \sum_i c_{ij} \hat{\mathbf{u}}_{j|i} \quad (3)$$

Squashing: (vector) non-linearity



$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|} \quad (4)$$

Squashing: (vector) non-linearity



$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|} \quad (4)$$

Squashing: Plot for 1-D input

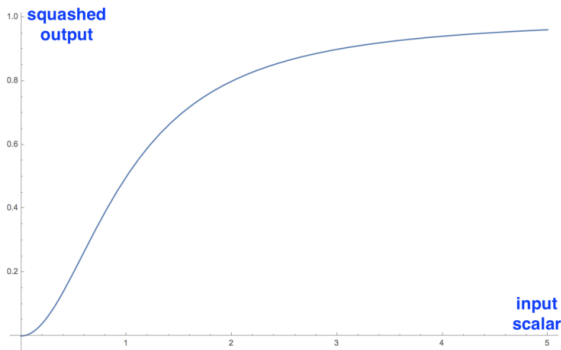
$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}$$

additional "squashing" unit scaling

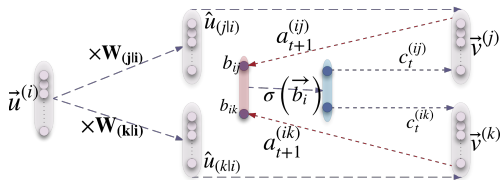
Squashing: Plot for 1-D input

$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}$$

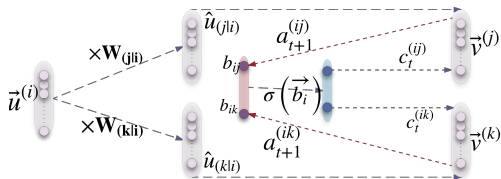
additional "squashing" unit scaling



Routing Algorithm



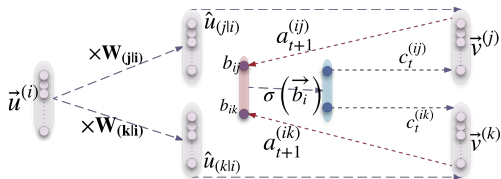
Routing Algorithm



Algorithm Dynamic Routing between Capsules

1: **procedure** ROUTING($\hat{u}_{j|i}, r, l$)

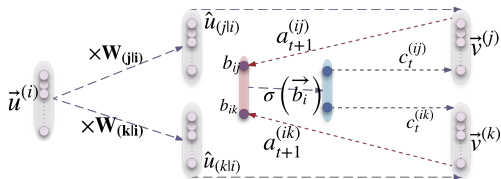
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}, r, l$)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.

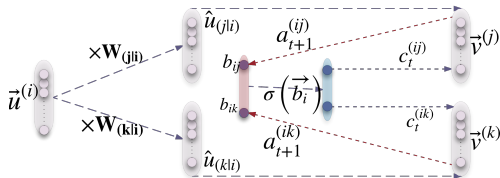
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}$, r , l)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
- 3: **for** r iterations **do**

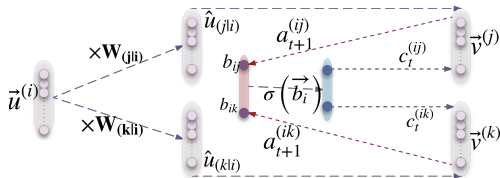
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}, r, l$)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
- 3: **for** r iterations **do**
- 4: for all capsule i in layer l : $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ ▷ softmax from Eq. 1

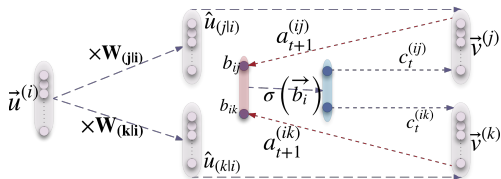
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}, r, l$)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
- 3: **for** r iterations **do**
- 4: for all capsule i in layer l : $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ \triangleright softmax from Eq. 1
- 5: for all capsule j in layer $(l + 1)$: $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{u}_{j|i}$ \triangleright total input from Eq. 3

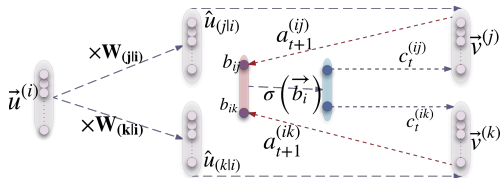
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}$, r , l)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
- 3: **for** r iterations **do**
- 4: for all capsule i in layer l : $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ ▷ softmax from Eq. 1
- 5: for all capsule j in layer $(l + 1)$: $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{u}_{j|i}$ ▷ total input from Eq. 3
- 6: for all capsule j in layer $(l + 1)$: $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$ ▷ squash from Eq. 4

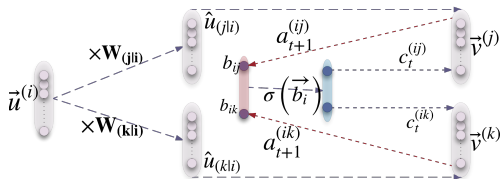
Routing Algorithm



Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}$, r , l)
- 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
- 3: **for** r iterations **do**
- 4: for all capsule i in layer l : $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ ▷ softmax from Eq. 1
- 5: for all capsule j in layer $(l + 1)$: $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{u}_{j|i}$ ▷ total input from Eq. 3
- 6: for all capsule j in layer $(l + 1)$: $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$ ▷ squash from Eq. 4
- 7: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow b_{ij} + \hat{u}_{j|i} \cdot \mathbf{v}_j$

Routing Algorithm

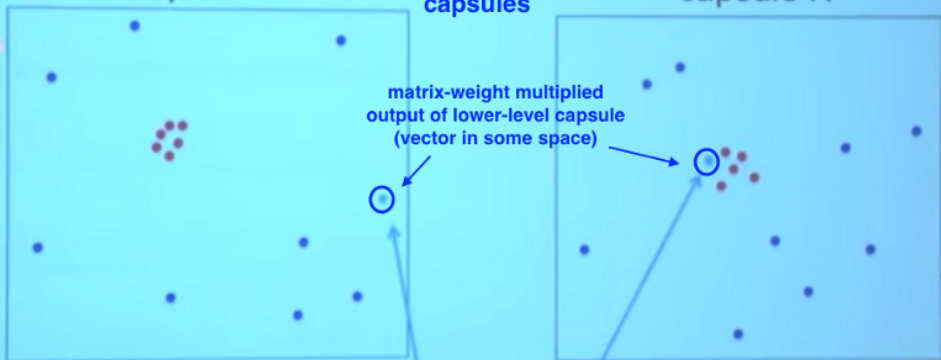


Algorithm Dynamic Routing between Capsules

- 1: **procedure** ROUTING($\hat{u}_{j|i}$, r , l)
 - 2: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow 0$.
 - 3: **for** r iterations **do**
 - 4: for all capsule i in layer l : $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ ▷ softmax from Eq. 1
 - 5: for all capsule j in layer $(l + 1)$: $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{u}_{j|i}$ ▷ total input from Eq. 3
 - 6: for all capsule j in layer $(l + 1)$: $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$ ▷ squash from Eq. 4
 - 7: for all capsule i in layer l and capsule j in layer $(l + 1)$: $b_{ij} \leftarrow b_{ij} + \hat{u}_{j|i} \cdot \mathbf{v}_j$
 - return** \mathbf{v}_j
-

Dynamic routing based on agreement

capsule J ← higher-level capsules → capsule K



send less send more

some
lower-level
capsule

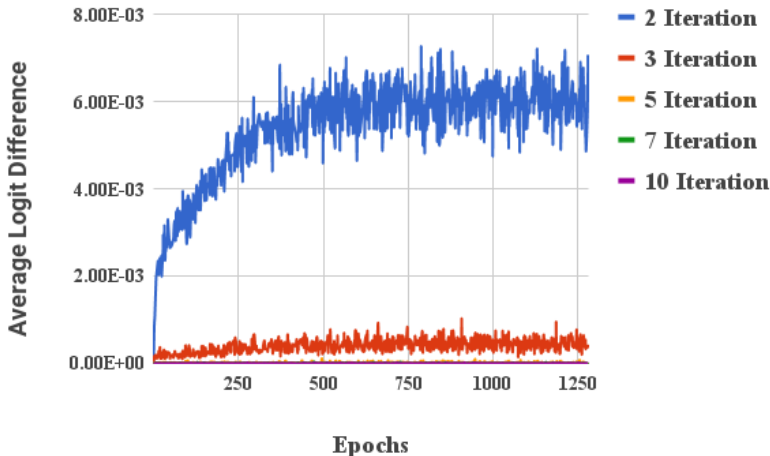
softmax routing of output from
one lower-level capsule

Average Change of Each Routing Logit b_{ij}

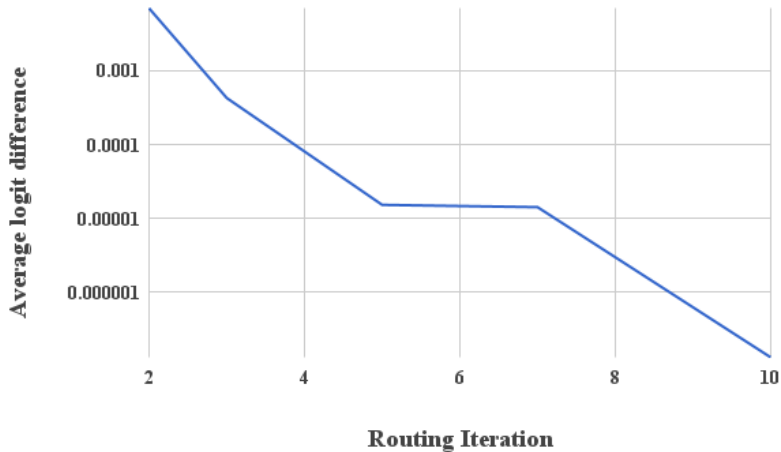
(by each routing iteration during training)

Average Change of Each Routing Logit b_{ij}

(by each routing iteration during training)



Log Scale of Final Differences



Training Loss of CapsNet on CIFAR10

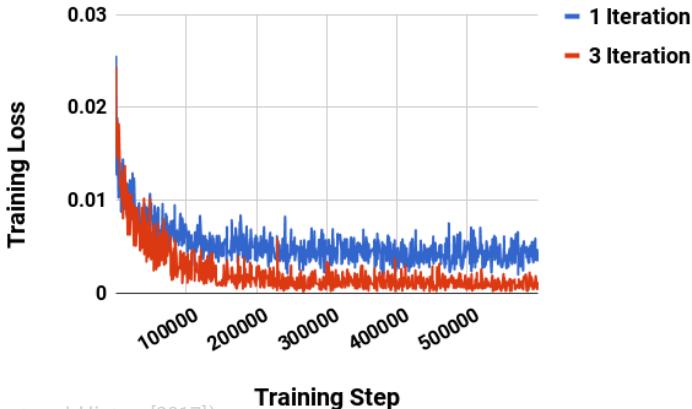
(batch size of 128)

The CapsNet with 3 routing iterations optimizes the loss faster and converges to a lower loss at the end.

Training Loss of CapsNet on CIFAR10

(batch size of 128)

The CapsNet with **3 routing iterations** optimizes the loss faster and converges to a lower loss at the end.

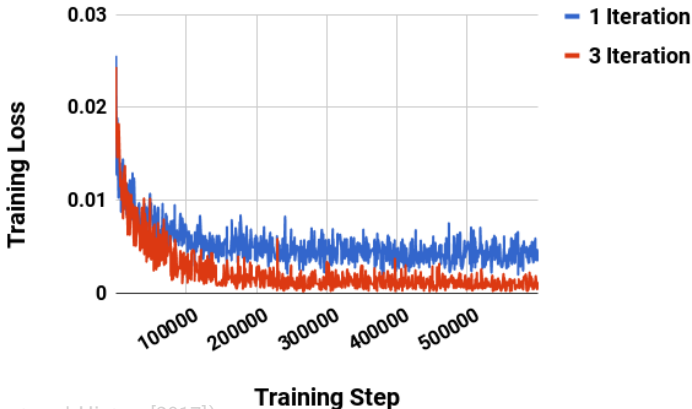


(Sabour, Frosst and Hinton [2017])

Training Loss of CapsNet on CIFAR10

(batch size of 128)

The CapsNet with **3 routing iterations** optimizes the loss faster and converges to a lower loss at the end.

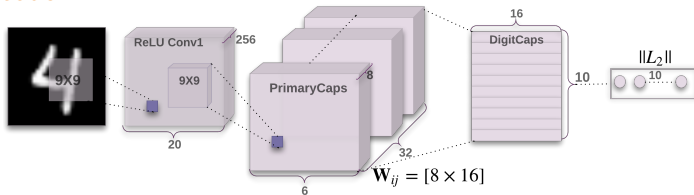


(Sabour, Frosst and Hinton [2017])

Capsule Network

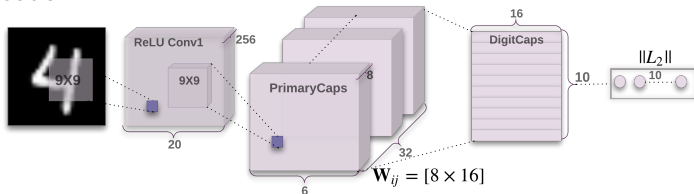
Architecture: Encoder-Decoder

■ encoder:

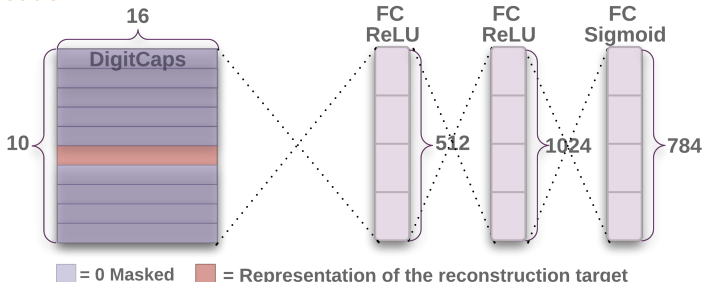


Architecture: Encoder-Decoder

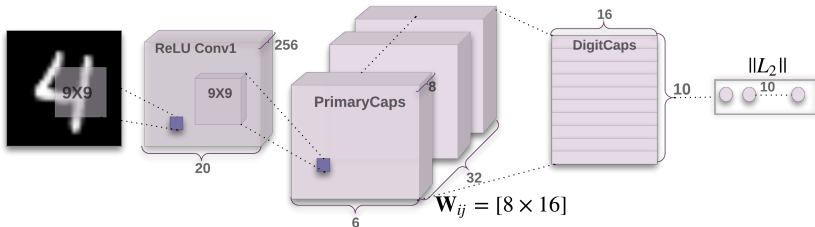
■ encoder:



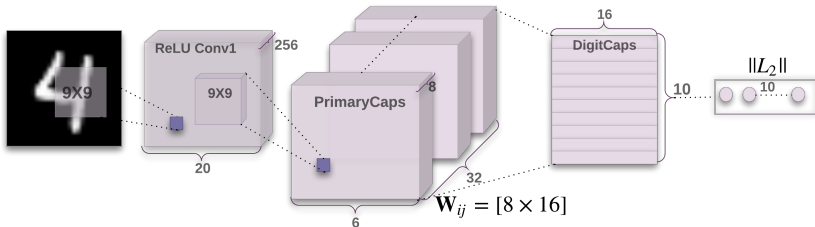
■ decoder:



Encoder: CapsNet with 3 Layers

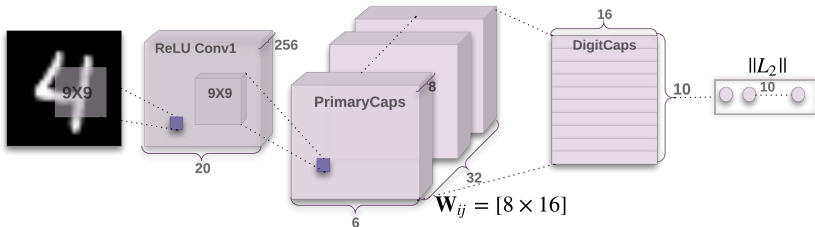


Encoder: CapsNet with 3 Layers



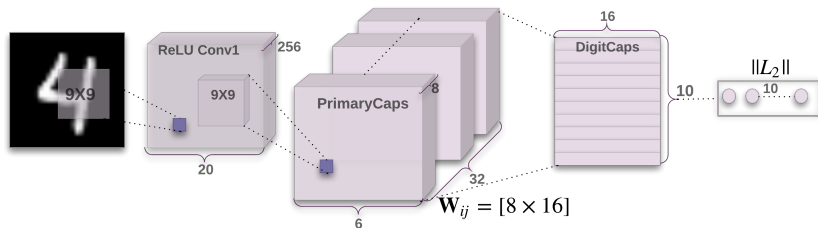
- input: 28 by 28 MNIST digit image

Encoder: CapsNet with 3 Layers

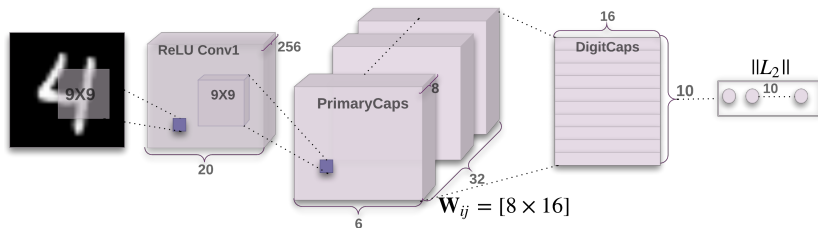


- input: 28 by 28 MNIST digit image
- output: 16-dimensional vector of instantiation parameters

Encoder Layer 1: (Standard) Convolutional Layer

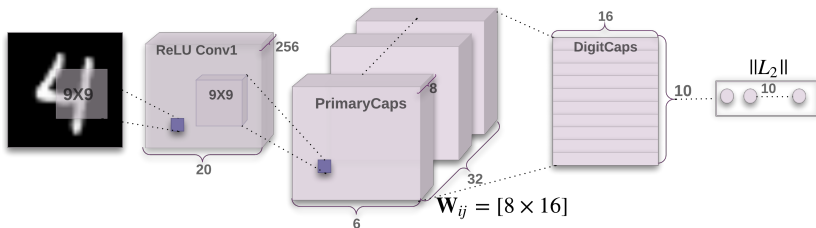


Encoder Layer 1: (Standard) Convolutional Layer



- input: 28×28 image (one color channel)

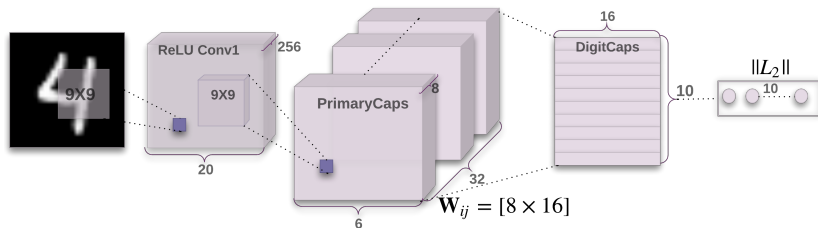
Encoder Layer 1: (Standard) Convolutional Layer



■ input: 28×28 image (one color channel)

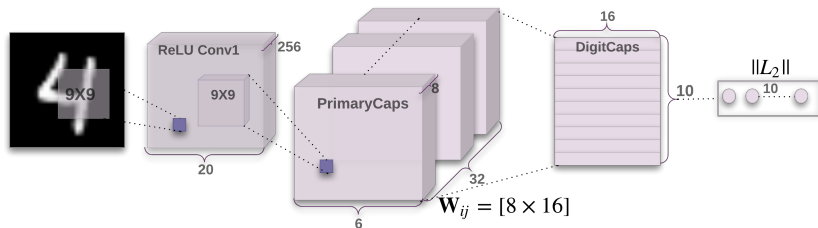
■ output: $20 \times 20 \times 256$

Encoder Layer 1: (Standard) Convolutional Layer



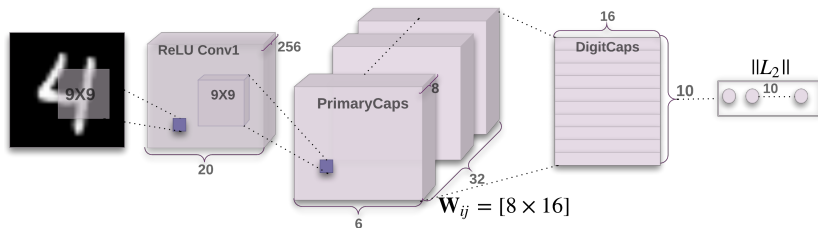
- input: 28×28 image (one color channel)
- output: $20 \times 20 \times 256$
- 256 kernels with size of $9 \times 9 \times 1$

Encoder Layer 1: (Standard) Convolutional Layer



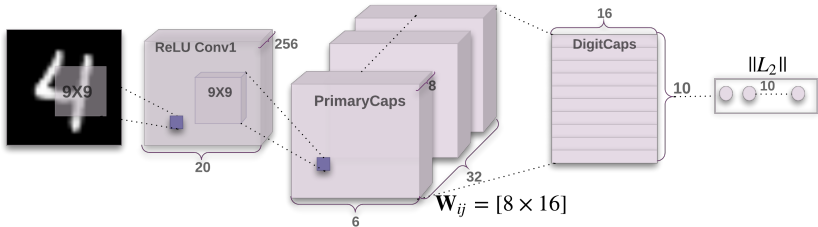
- input: 28×28 image (one color channel)
- output: $20 \times 20 \times 256$
- 256 kernels with size of $9 \times 9 \times 1$
- stride 1

Encoder Layer 1: (Standard) Convolutional Layer

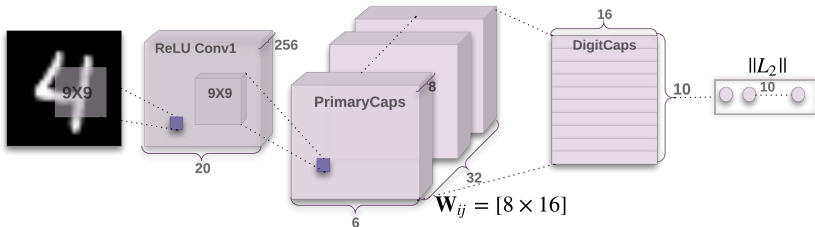


- input: 28×28 image (one color channel)
- output: $20 \times 20 \times 256$
- 256 kernels with size of $9 \times 9 \times 1$
- stride 1
- ReLU activation

Encoder Layer 2: PrimaryCaps



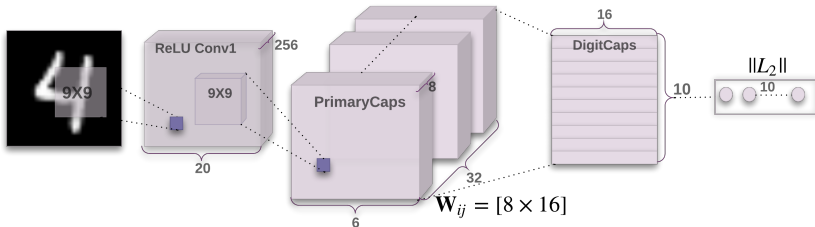
Encoder Layer 2: PrimaryCaps



■ input: $20 \times 20 \times 256$

basic features detected by the convolutional layer

Encoder Layer 2: PrimaryCaps



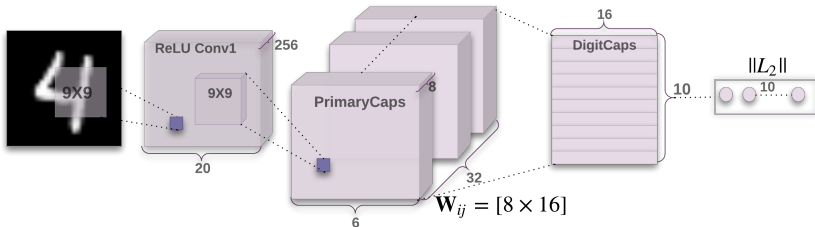
■ input: $20 \times 20 \times 256$

basic features detected by the convolutional layer

■ output: $6 \times 6 \times 8 \times 32$

vector (activation) outputs of primary capsules

Encoder Layer 2: PrimaryCaps



■ input: $20 \times 20 \times 256$

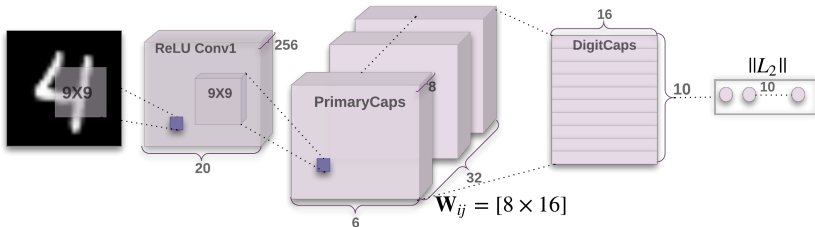
basic features detected by the convolutional layer

■ output: $6 \times 6 \times 8 \times 32$

vector (activation) outputs of primary capsules

■ 32 primary capsules

Encoder Layer 2: PrimaryCaps



■ input: $20 \times 20 \times 256$

basic features detected by the convolutional layer

■ output: $6 \times 6 \times 8 \times 32$

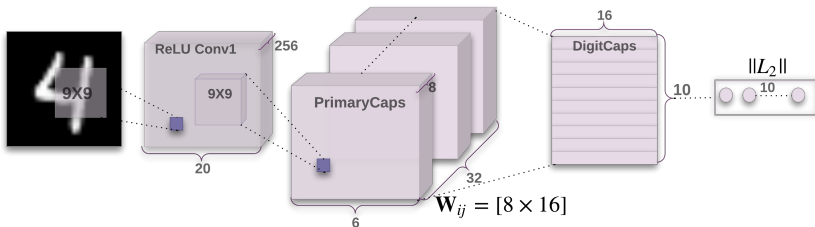
vector (activation) outputs of primary capsules

■ 32 primary capsules

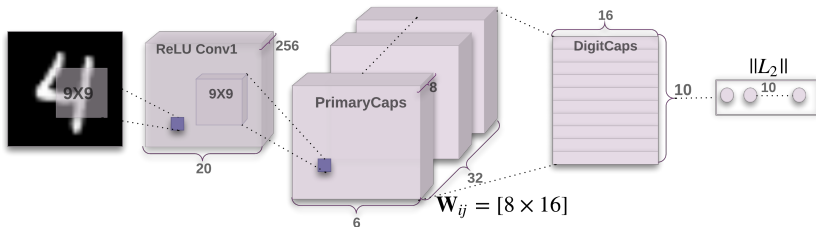
■ each applies eight $9 \times 9 \times 256$ convolutional kernels

to the $20 \times 20 \times 256$ input to produce $6 \times 6 \times 8$ output

Encoder Layer 3: DigitCaps



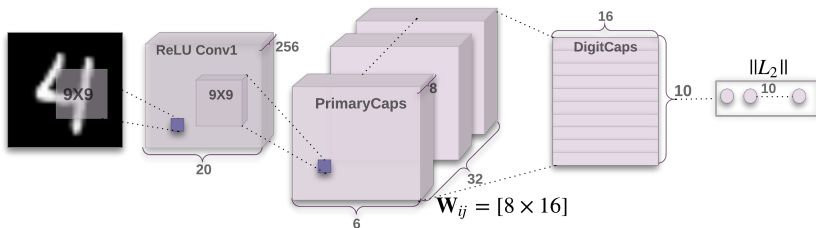
Encoder Layer 3: DigitCaps



■ input: $6 \times 6 \times 8 \times 32$

($6 \times 6 \times 32$)-many 8-dimensional vector activations

Encoder Layer 3: DigitCaps

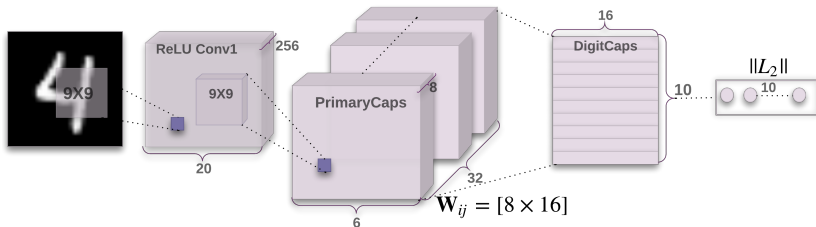


■ input: $6 \times 6 \times 8 \times 32$

($6 \times 6 \times 32$)-many 8-dimensional vector activations

■ output: 16×10

Encoder Layer 3: DigitCaps



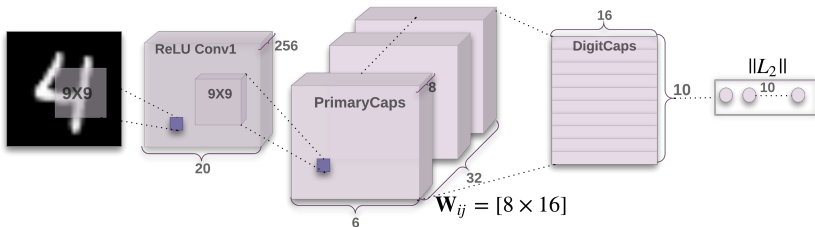
■ input: $6 \times 6 \times 8 \times 32$

($6 \times 6 \times 32$)-many 8-dimensional vector activations

■ output: 16×10

■ 10 digit capsules

Encoder Layer 3: DigitCaps



■ input: $6 \times 6 \times 8 \times 32$

($6 \times 6 \times 32$)-many 8-dimensional vector activations

■ output: 16×10

■ 10 digit capsules

■ input vectors gets their own 8×16 weight matrix W_{ij}

that maps 8-dimensional input space to the 16-dimensional capsule output space

Margin Loss

for a Digit Existence

CapsNet Loss Function

loss term for one DigitCap

calculated for correct DigitCap

calculated for incorrect DigitCaps

$$L_c = T_c \max(0, m^+ - ||\mathbf{v}_c||)^2 + \lambda (1 - T_c) \max(0, ||\mathbf{v}_c|| - m^-)^2$$

1 when correct DigitCap, 0 when incorrect

zero loss when correct prediction with probability greater than 0.9, non-zero otherwise

0.5 constant used for numerical stability

1 when incorrect DigitCap, 0 when correct

zero loss when incorrect prediction with probability less than 0.1, non-zero otherwise

L2 norm

L2 norm

Note: correct DigitCap is one that matches training label, for each training example there will be 1 correct and 9 incorrect DigitCaps

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

■ $m^+ = 0.9$:

The loss is 0 iff the **correct** DigitCap predicts the correct label with probability ≥ 0.9 .

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

■ $m^+ = 0.9$:

The loss is 0 iff the **correct** DigitCap predicts the correct label with probability ≥ 0.9 .

■ $m^- = 0.1$:

The loss is 0 iff the **mismatching** DigitCap predicts an incorrect label with probability ≤ 0.1 .

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

■ $m^+ = 0.9$:

The loss is 0 iff the **correct** DigitCap predicts the correct label with probability ≥ 0.9 .

■ $m^- = 0.1$:

The loss is 0 iff the **mismatching** DigitCap predicts an incorrect label with probability ≤ 0.1 .

■ $\lambda = 0.5$ is down-weighting of the loss for absent digit classes.

It stops the initial learning from shrinking the lengths of the activity vectors.

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

■ $m^+ = 0.9$:

The loss is 0 iff the **correct** DigitCap predicts the correct label with probability ≥ 0.9 .

■ $m^- = 0.1$:

The loss is 0 iff the **mismatching** DigitCap predicts an incorrect label with probability ≤ 0.1 .

■ $\lambda = 0.5$ is down-weighting of the loss for absent digit classes.

It stops the initial learning from shrinking the lengths of the activity vectors.

■ **Squares?** Because there are L_2 norms in the loss function?

Margin Loss

to Train the Whole Encoder

In other words, each DigitCap c has loss:

$$L_c = \begin{cases} \max(0, m^+ - \|\mathbf{v}_c\|)^2 & \text{iff a digit of class } c \text{ is present,} \\ \lambda \max(0, \|\mathbf{v}_c\| - m^-)^2 & \text{otherwise.} \end{cases}$$

■ $m^+ = 0.9$:

The loss is 0 iff the **correct** DigitCap predicts the correct label with probability ≥ 0.9 .

■ $m^- = 0.1$:

The loss is 0 iff the **mismatching** DigitCap predicts an incorrect label with probability ≤ 0.1 .

■ $\lambda = 0.5$ is down-weighting of the loss for absent digit classes.

It stops the initial learning from shrinking the lengths of the activity vectors.

■ Squares? Because there are L_2 norms in the loss function?

■ The total loss is the sum of the losses of all digit capsules.

Margin Loss

Function Value for Positive and for Negative Class

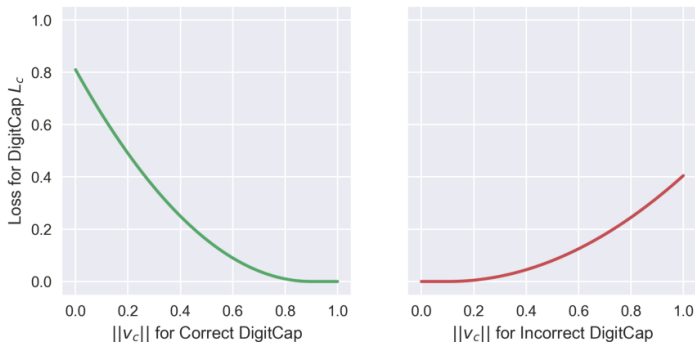
- For the **correct** DigitCap, the loss is 0 iff it predicts the correct label with probability ≥ 0.9 .
- For the **mismatching** DigitCap, the loss is 0 iff it predicts an incorrect label with probability ≤ 0.1 .

Margin Loss

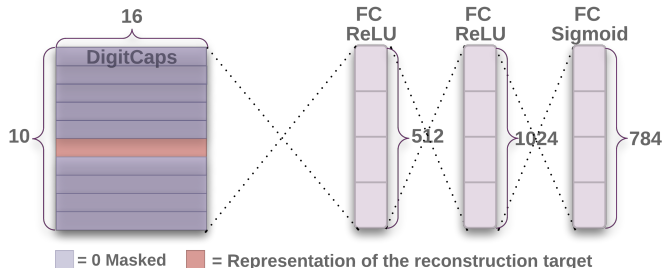
Function Value for Positive and for Negative Class

- For the **correct** DigitCap, the loss is 0 iff it predicts the correct label with probability ≥ 0.9 .
- For the **mismatching** DigitCap, the loss is 0 iff it predicts an incorrect label with probability ≤ 0.1 .

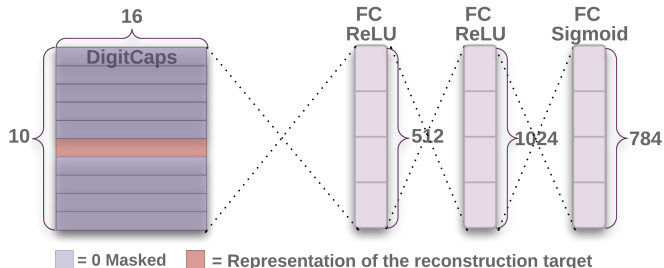
Loss Function Value for Correct and Incorrect DigitCap



Decoder: Regularization of CapsNets



Decoder: Regularization of CapsNets

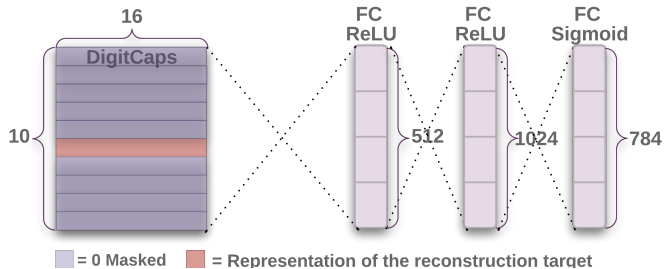


Decoder is used for regularization:

- decodes input from DigitCaps

to recreate an image of a (28×28) -pixels digit

Decoder: Regularization of CapsNets



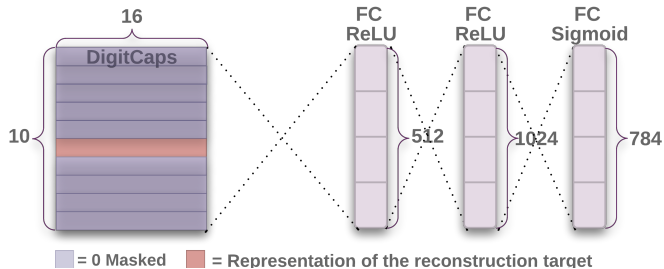
Decoder is used for regularization:

- decodes input from DigitCaps

to recreate an image of a (28×28) -pixels digit

- with the loss function being the Euclidean distance

Decoder: Regularization of CapsNets



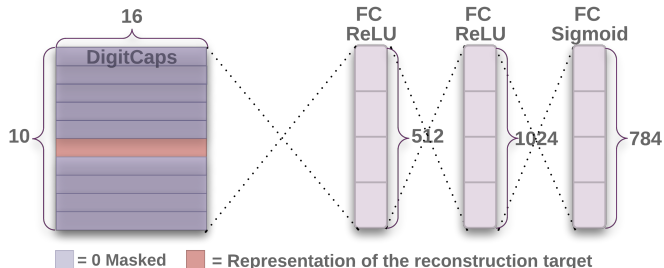
Decoder is used for regularization:

- decodes input from DigitCaps

to recreate an image of a (28×28) -pixels digit

- with the loss function being the Euclidean distance
- ignores the negative classes

Decoder: Regularization of CapsNets



Decoder is used for regularization:

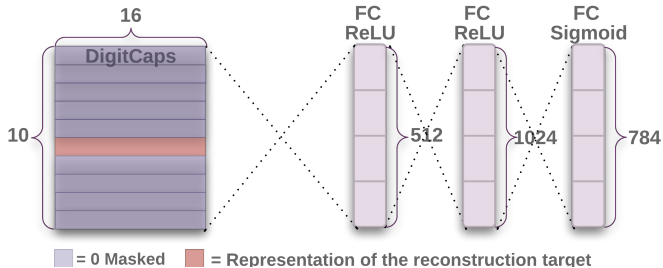
- decodes input from DigitCaps

to recreate an image of a (28×28) -pixels digit

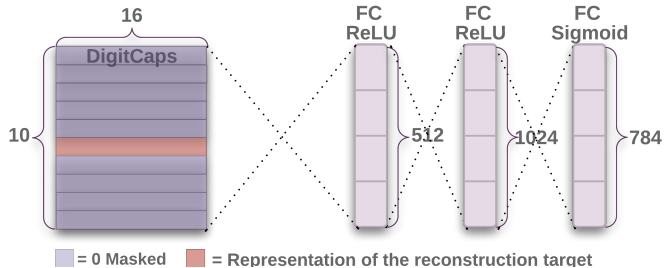
- with the loss function being the Euclidean distance
- ignores the negative classes
- forces capsules to learn features useful for reconstruction

(Sabour, Frosst and Hinton [2017])

Decoder: 3 Fully Connected Layers

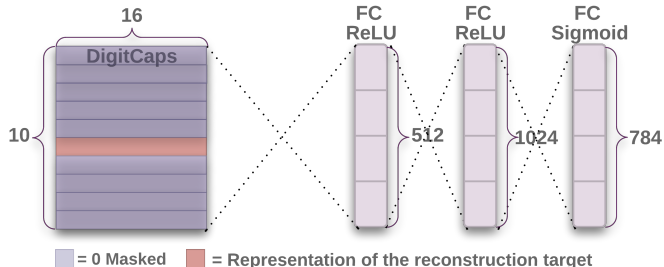


Decoder: 3 Fully Connected Layers



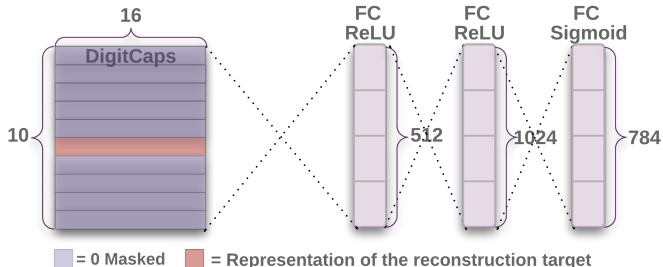
- Layer 4: from 16×10 input to 512 output, ReLU activations

Decoder: 3 Fully Connected Layers



- Layer 4: from 16×10 input to 512 output, ReLU activations
- Layer 5: from 512 input to 1024 output, ReLU activations

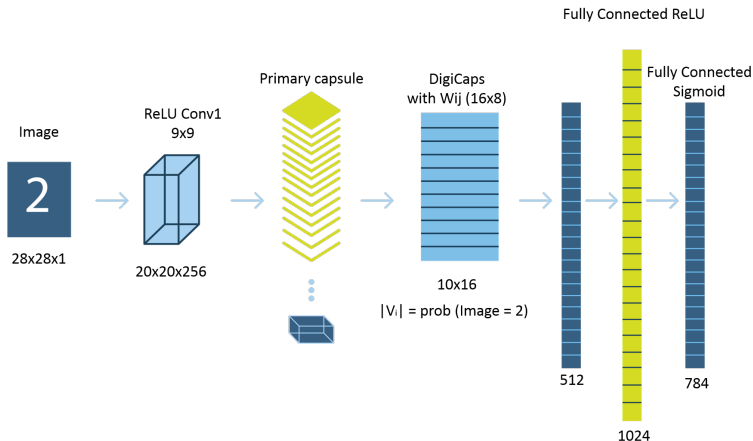
Decoder: 3 Fully Connected Layers



- Layer 4: from 16×10 input to 512 output, ReLU activations
- Layer 5: from 512 input to 1024 output, ReLU activations
- Layer 6: from 1024 input to 784 output, sigmoid activations

(after reshaping it produces a (28×28) -pixels decoded image)

Architecture: Summary



Experiments

MNIST Reconstructions (CapsNet, 3 routing iterations)

Label:	8	5	5	5
Prediction:	8	5	3	3
Reconstruction:	8	5	5	3
Input:				
Output:				

Dimension Perturbations

one of the 16 dimensions, by intervals of 0.05 in the range $[-0.25, 0.25]$:





Dimension Perturbations

one of the 16 dimensions, by intervals of 0.05 in the range $[-0.25, 0.25]$:

Interpretation	Reconstructions after perturbing
“scale and thickness”	
“localized part”	






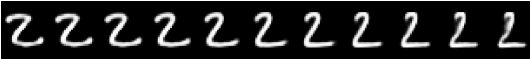
Dimension Perturbations

one of the 16 dimensions, by intervals of 0.05 in the range $[-0.25, 0.25]$:

Interpretation	Reconstructions after perturbing
“scale and thickness”	
“localized part”	
“stroke thickness”	
“localized skew”	

Dimension Perturbations

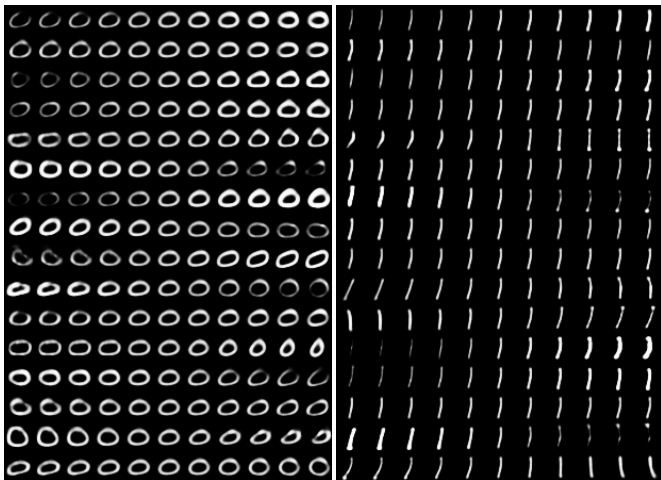
one of the 16 dimensions, by intervals of 0.05 in the range $[-0.25, 0.25]$:

Interpretation	Reconstructions after perturbing
“scale and thickness”	
“localized part”	
“stroke thickness”	
“localized skew”	
“width and translation”	
“localized part”	

Dimension Perturbations: Latent Codes of 0 and 1

rows: DigitCaps dimensions

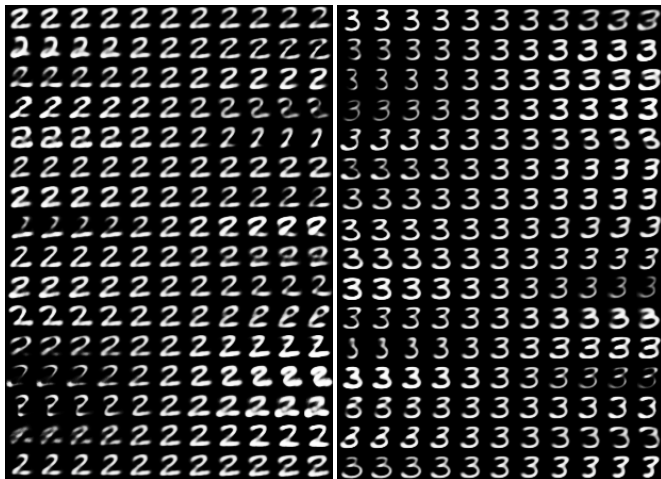
columns (from left to right): $+ \{-0.25, -0.2, -0.15, -0.1, -0.05, 0, 0.05, 0.1, 0.15, 0.2, 0.25\}$



Dimension Perturbations: Latent Codes of 2 and 3

rows: DigitCaps dimensions

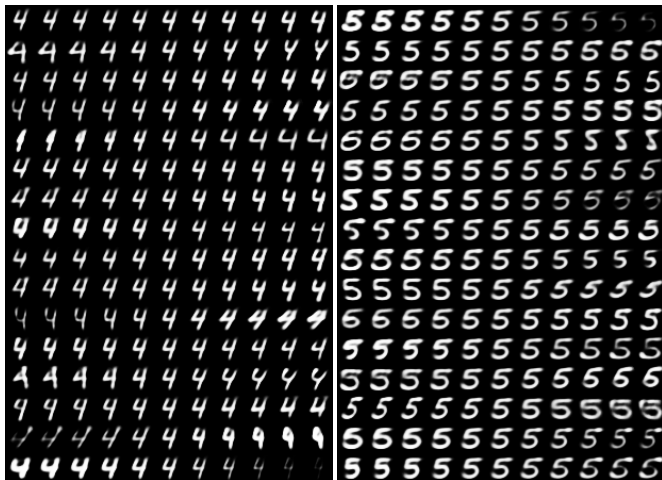
columns (from left to right): $+ \{-0.25, -0.2, -0.15, -0.1, -0.05, 0, 0.05, 0.1, 0.15, 0.2, 0.25\}$



Dimension Perturbations: Latent Codes of 4 and 5

rows: DigitCaps dimensions

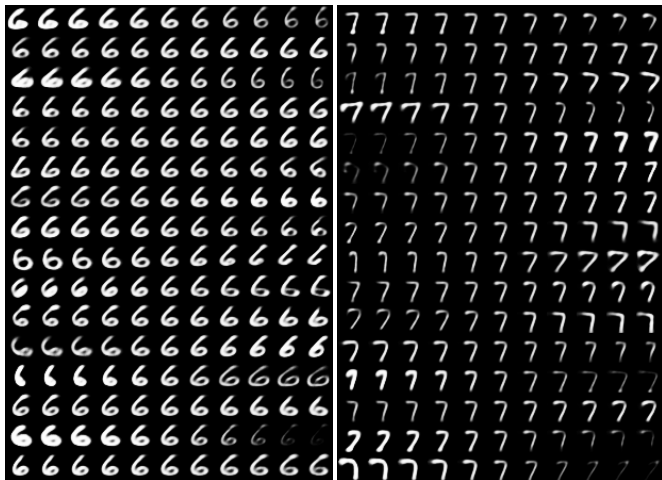
columns (from left to right): $+ \{-0.25, -0.2, -0.15, -0.1, -0.05, 0, 0.05, 0.1, 0.15, 0.2, 0.25\}$



Dimension Perturbations: Latent Codes of 6 and 7

rows: DigitCaps dimensions

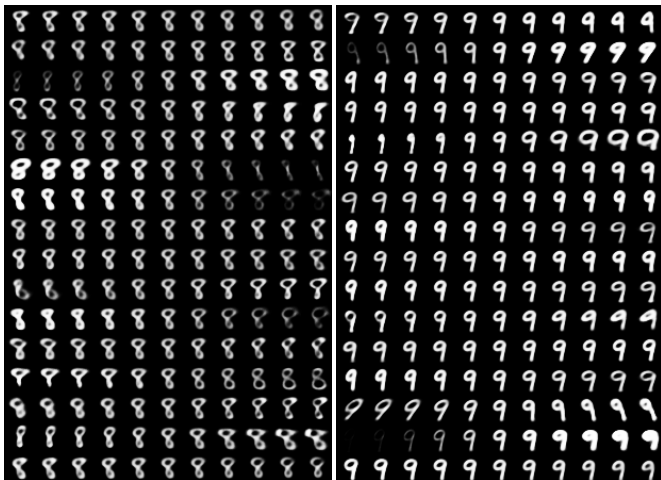
columns (from left to right): $+ \{-0.25, -0.2, -0.15, -0.1, -0.05, 0, 0.05, 0.1, 0.15, 0.2, 0.25\}$



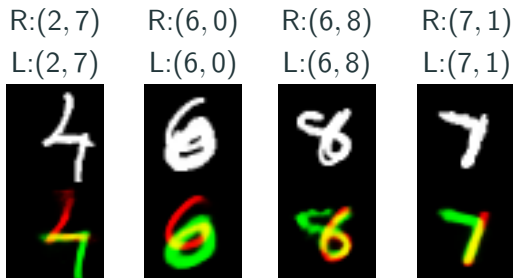
Dimension Perturbations: Latent Codes of 8 and 9

rows: DigitCaps dimensions

columns (from left to right): $+ \{-0.25, -0.2, -0.15, -0.1, -0.05, 0, 0.05, 0.1, 0.15, 0.2, 0.25\}$



MultiMNIST Reconstructions (CapsNet, 3 routing iterations)

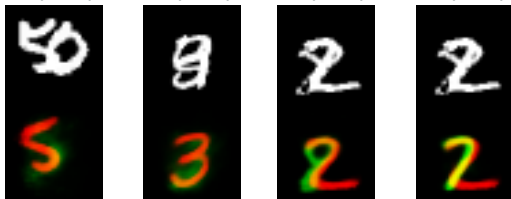


MultiMNIST Reconstructions (CapsNet, 3 routing iterations)

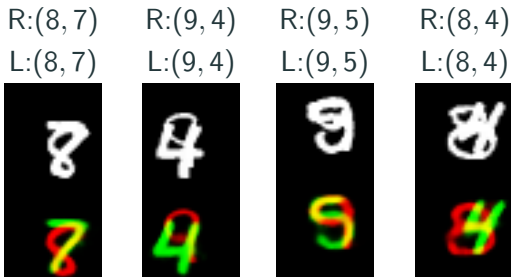
R:(2, 7)	R:(6, 0)	R:(6, 8)	R:(7, 1)
L:(2, 7)	L:(6, 0)	L:(6, 8)	L:(7, 1)



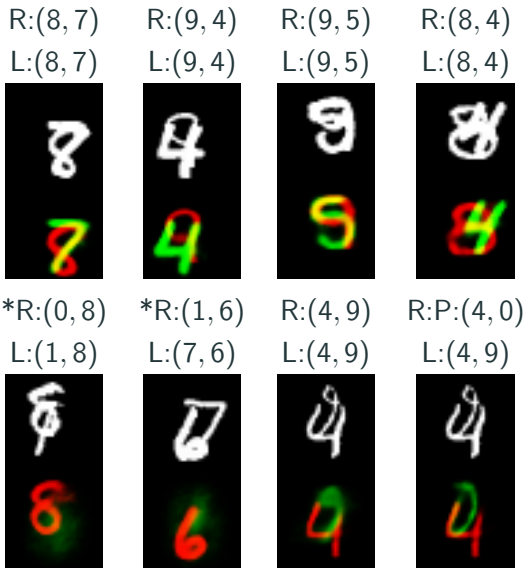
*R:(5, 7)	*R:(2, 3)	R:(2, 8)	R:P:(2, 7)
L:(5, 0)	L:(4, 3)	L:(2, 8)	L:(2, 8)



MultiMNIST Reconstructions (CapsNet, 3 routing iterations)



MultiMNIST Reconstructions (CapsNet, 3 routing iterations)



(Sabour, Frosst and Hinton [2017])

Results on MNIST and MultiMNIST

CapsNet classification test accuracy:

Results on MNIST and MultiMNIST

CapsNet classification test accuracy:

Method	Routing	Reconstruction	MNIST (%)	MultiMNIST (%)
Baseline	-	-	0.39	8.1
CapsNet	1	no	$0.34_{\pm 0.032}$	-
CapsNet	1	yes	$0.29_{\pm 0.011}$	7.5
CapsNet	3	no	$0.35_{\pm 0.036}$	-
CapsNet	3	yes	$0.25_{\pm 0.005}$	5.2

(The MNIST average and standard deviation results are reported from 3 trials.)

Results on Other Datasets

CIFAR10

- 10.6% test error

Results on Other Datasets

CIFAR10

- 10.6% test error
- ensemble of 7 models

Results on Other Datasets

CIFAR10

- 10.6% test error
- ensemble of 7 models
- 3 routing iterations

Results on Other Datasets

CIFAR10

- 10.6% test error
- ensemble of 7 models
- 3 routing iterations
- 24×24 patches of the image

Results on Other Datasets

CIFAR10

- 10.6% test error
- ensemble of 7 models
- 3 routing iterations
- 24×24 patches of the image
- about what standard convolutional nets achieved when they were first applied to CIFAR10 (Zeiler and Fergus [2013])

Conclusion

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era
- capable of learning to achieve **state-of-the art** performance by only using **a fraction of the data** compared to CNN

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era
- capable of learning to achieve **state-of-the art** performance by only using **a fraction of the data** compared to CNN
 - task of telling digits apart: the human brain needs a couple of dozens of examples (hundreds at most).

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era
- capable of learning to achieve **state-of-the art** performance by only using **a fraction of the data** compared to CNN
 - task of telling digits apart: the human brain needs a couple of dozens of examples (hundreds at most).
 - On the other hand, CNNs typically need tens of thousands of them.

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era
- capable of learning to achieve **state-of-the art** performance by only using **a fraction of the data** compared to CNN
 - task of telling digits apart: the human brain needs a couple of dozens of examples (hundreds at most).
 - On the other hand, CNNs typically need tens of thousands of them.

Benefits:

- a new building block usable in deep learning to better model hierarchical relationships
- representations similar to scene graphs in computer graphics
- no algorithm to implement and train a capsule network:
 - now **dynamic routing algorithm**
 - one of the reasons: computers not powerful enough in the pre-GPU-based era
- capable of learning to achieve **state-of-the art** performance by only using **a fraction of the data** compared to CNN
 - task of telling digits apart: the human brain needs a couple of dozens of examples (hundreds at most).
 - On the other hand, CNNs typically need tens of thousands of them.

Downsides:

- **current implementations: much slower than other modern deep learning models**

Thank you!

Questions?

Backup Slides

Results on Other Datasets

smallINORB (LeCun et al. [2004])

■ 2.7% test error

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

smallINORB (LeCun et al. [2004])

- 2.7% test error

- 96×96 stereo grey-scale images

resized to 48×48 ; during training processed random 32×32 crops; during test the central 32×32 patch

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

smallNORB (LeCun et al. [2004])

- 2.7% test error
- 96×96 stereo grey-scale images

resized to 48×48 ; during training processed random 32×32 crops; during test the central 32×32 patch

- same CapsNet architecture as for MNIST

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

smallNORB (LeCun et al. [2004])

- 2.7% test error
- 96×96 stereo grey-scale images
 - resized to 48×48 ; during training processed random 32×32 crops; during test the central 32×32 patch
- same CapsNet architecture as for MNIST
- on-par with the state-of-the-art (Ciresan et al. [2011])

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

SVHN (Netzer et al. [2011])

■ 4.3% test error

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

SVHN (Netzer et al. [2011])

- 4.3% test error
- only 73257 images!

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

SVHN (Netzer et al. [2011])

- 4.3% test error
- only 73257 images!
- the number of first convolutional layer channels reduced to 64

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

SVHN (Netzer et al. [2011])

- 4.3% test error
- only 73257 images!
- the number of first convolutional layer channels reduced to 64
- the primary capsule layer: to 16 $6D$ -capsules

(Sabour, Frosst and Hinton [2017])

Results on Other Datasets

SVHN (Netzer et al. [2011])

- 4.3% test error
- only 73257 images!
- the number of first convolutional layer channels reduced to 64
- the primary capsule layer: to 16 $6D$ -capsules
- final capsule layer: $8D$ -capsules

(Sabour, Frosst and Hinton [2017])

Further Reading i

