
Capsule Networks

Rıza Velioğlu
Sinan Kuşçu
Sinan Harms

Content

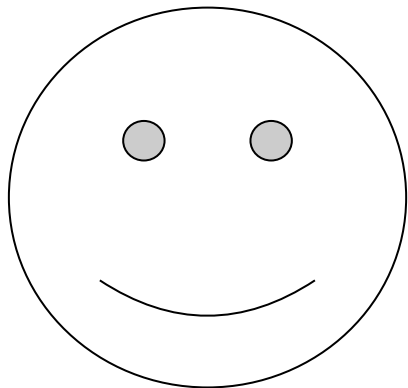
- Main Idea
 - CNN vs CapsNet
 - Concepts
 - Hierarchy of parts
 - Inverse Graphics
 - Capsules
 - Routing by agreement
 - CapsNet Architecture
 - Training & Loss
-

Main Idea

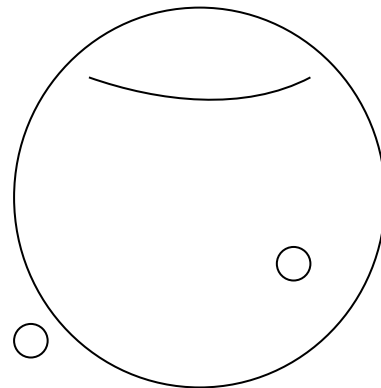
Convolutional Neural Networks

- most commonly applied in analyzing visual imagery
 - main component → convolutional layer
 - detects important features in image pixel
 - deeper layers detect simple features
 - higher layers combine simple features into more complex features
 - max pooling layers → reduce spacial size of data
 - dense layer → combine very high level features and make predictions
-

CNN Drawbacks



This is a smiley face



Is this also a smiley face?

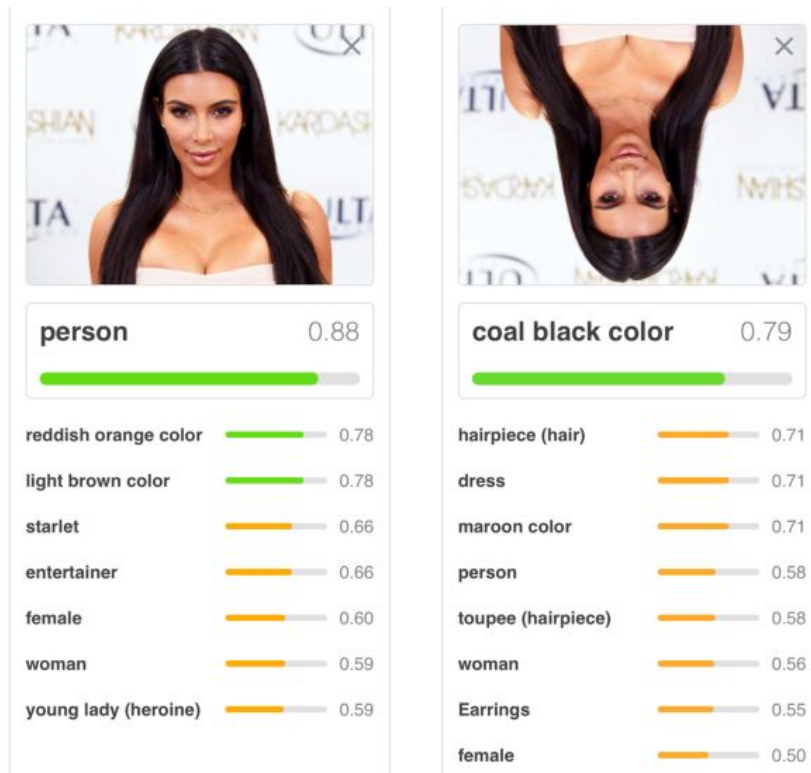
Problem:

if (2 eyes & 1 nose & 1 mouth):
 it's a face

CNN Drawbacks - max pooling

- messenger between two layers
 - tells layers about presence of parts, but not spatial relationship
 - strips information to create translational invariance
 - does not make the model view point invariant
-

CNN Drawbacks - Rotation Invariance



→ no consideration of rotational relationship

Capsule Networks - Intuition

Intuition:

- hierarchical pose relationships make it easy for the model to understand that the object it sees is the same, but from a different point of view

Example:



[source](#)

Capsule Networks

Solution:

if (2 adjacent eyes & nose under eyes & mouth under nose):

 It's a face

- Replace scalar-output feature vectors with vector-output capsules
 - Replace max-pooling layer with “routing-by-agreement”
-

CNN vs. CapsNet

smallINORB benchmark tests:

Test set	Azimuth		Elevation	
	CNN	Capsules	CNN	Capsules
Novel viewpoints	20%	13.5%	17.8%	12.3%
Familiar viewpoints	3.7%	3.7%	4.3%	4.3%

CapsNets vs. CNNs

Pro

- view point invariance
- fraction of data required to learn model
- consideration of pose relationships

Potential Cons:

- slower than modern deep learning models
 - not tested on different data sets
-

Concepts

(Computer) Graphics

$$T_{\lambda, \theta} = \begin{pmatrix} \lambda \cos \theta & \sin \theta \\ -\sin \theta & \lambda \cos \theta \end{pmatrix}$$

T is a transformation matrix
Theta specifies the rotation
Lambda specifies the sizing

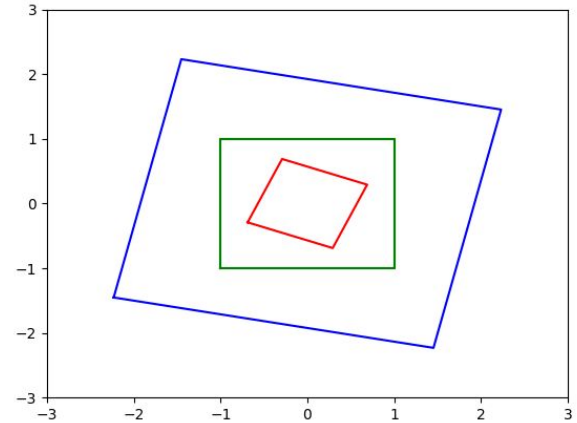
$$\begin{pmatrix} \lambda = 0.2 \\ \theta = 0.5 \end{pmatrix}$$
$$\begin{pmatrix} \lambda = 0.4 \\ \theta = 2 \end{pmatrix}$$



$$T_{\square}(\lambda = 0.2, \theta = 0.5)$$



$$T_{\square}(\lambda = 0.4, \theta = 2)$$



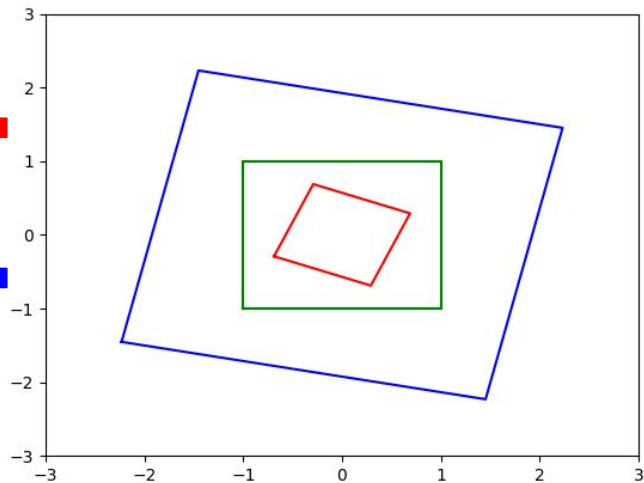
Inverse Graphics

Extracting the pose parameters of an object (rotation,size...) is called inverse graphics

$$\begin{pmatrix} \lambda = 0.2 \\ \theta = 0.5 \end{pmatrix}$$
$$\begin{pmatrix} \lambda = 0.4 \\ \theta = 2 \end{pmatrix}$$

$$\leftarrow T_{\square}^{-1}(\square) \leftarrow$$

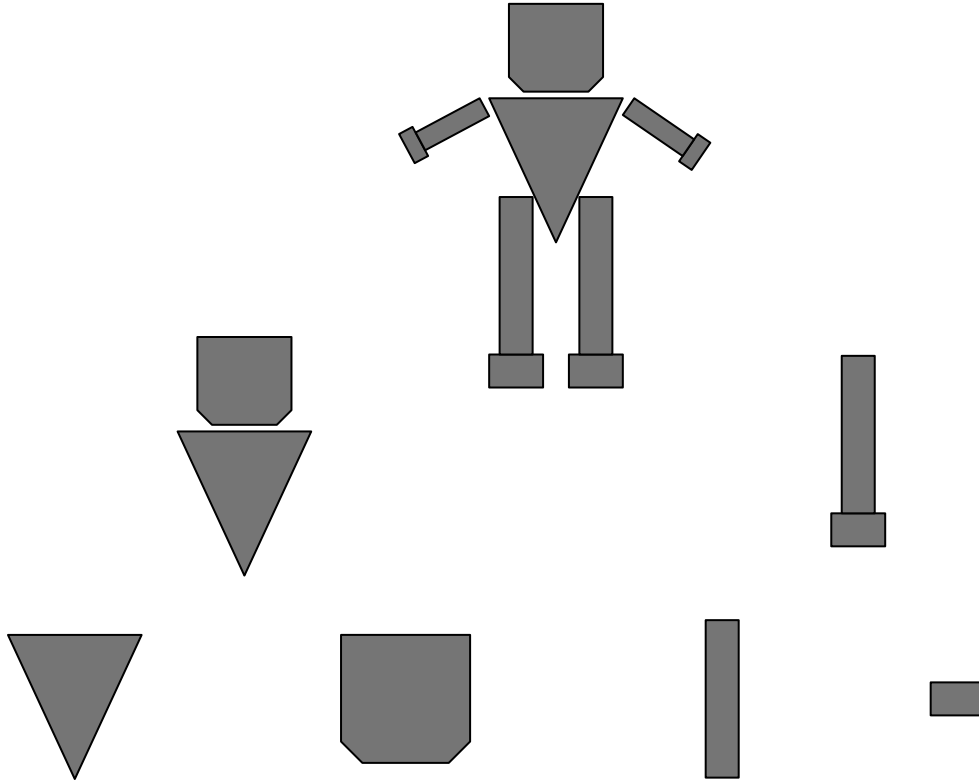
$$\leftarrow T_{\square}^{-1}(\square) \leftarrow$$



Hierarchy of Parts

Higher Level

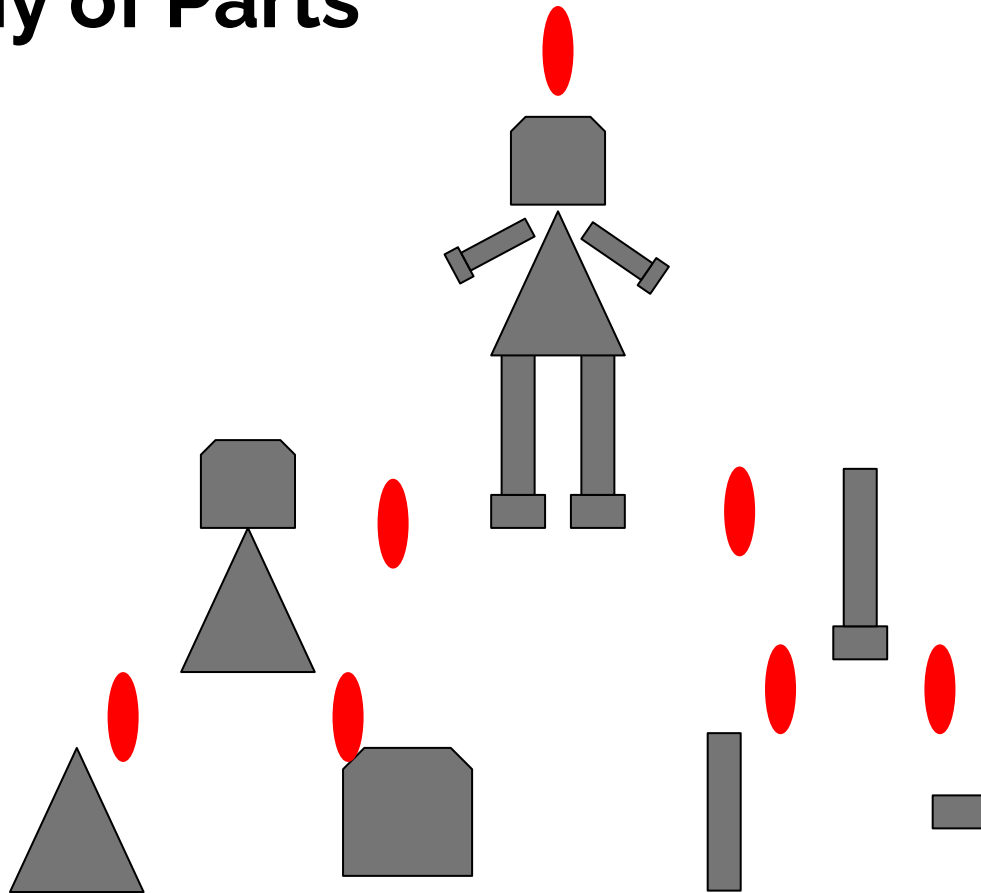
Lower Level



Hierarchy of Parts

Higher Level

Lower Level



The pose of lower level parts also decides about higher level parts

The pose of a part is detected by capsules



Capsules



Feed your Neural Network with capsules and it will see the truth!

Capsules

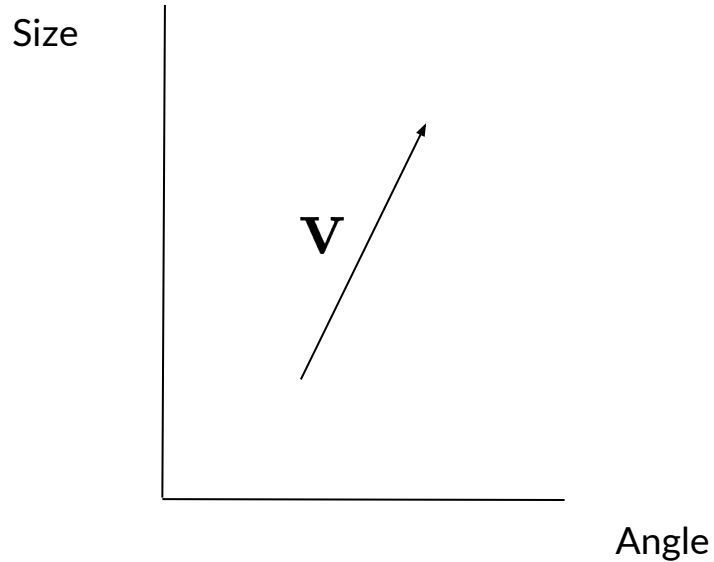
- A capsule is a fully connected neural network
 - It learns to extract the pose and presence for a given object
 - It gets the inputs from a conv layer or lower capsule layers
 - The conv layer converts pixel intensities to the activities of local feature detectors
 - It outputs a activation vector, which encodes
 - The probability that an object is present
 - A pose matrix
-

Capsules

Activation vector

Length = Estimated probability of presence

Direction = Objects estimated pose parameter



Capsules

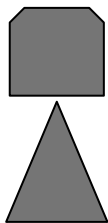
By **Squashing** it can be ensured that the output vector of the capsule j has a length between 0 and 1.

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

Weight factor

Normalized total input of capsule j

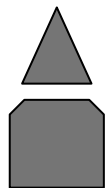
Routing by agreement



Triangle
Size = 0.8 Angle = 0°



Hexagon
Size = 1 Angle = 0°



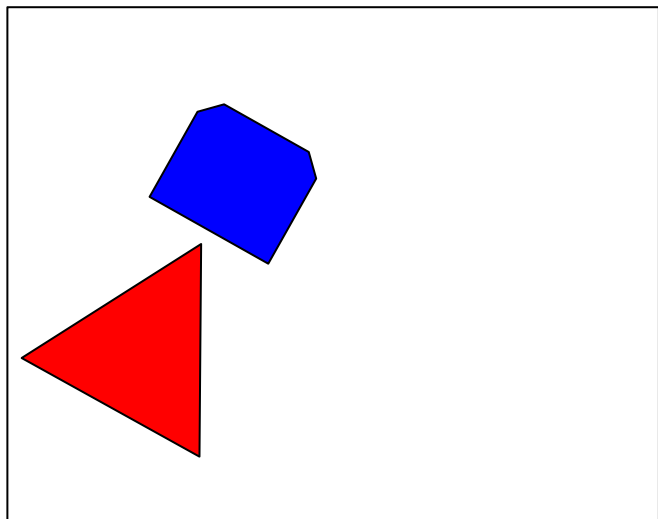
Triangle
Size = 0.5 Angle = 0°



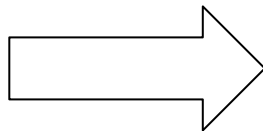
Hexagon
Size = 1 Angle = 180°

Routing by agreement

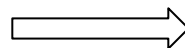
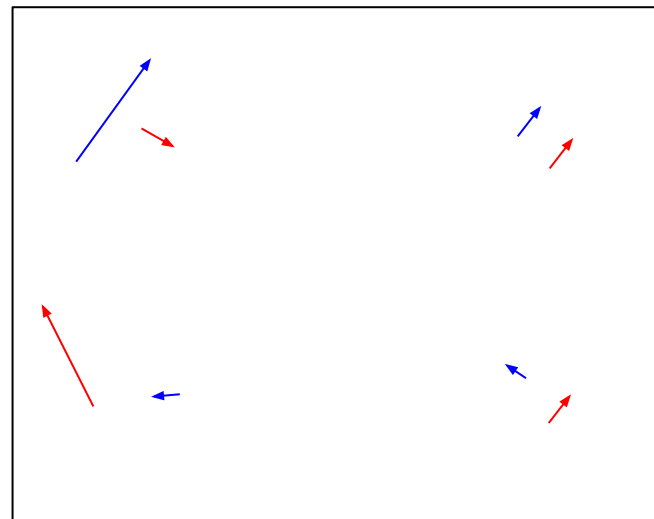
Output of convolutional layer



Hexagon capsule detects a hexagon in the upper-left corner
Triangle capsule detects a triangle in the lower-left corner



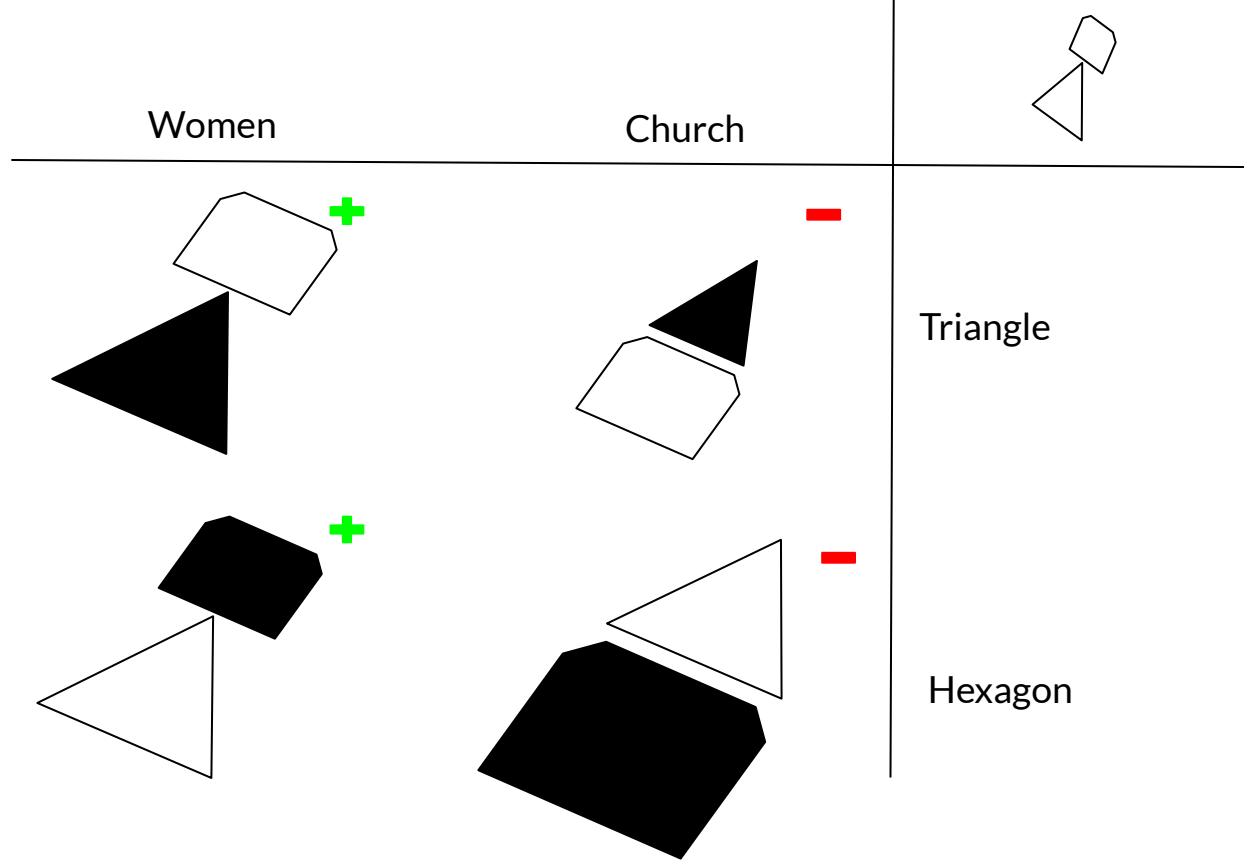
Activity vector of **hexagon** capsule
Activity vector of **triangle** capsule



Does these parts belong to a church or women?

Routing by agreement

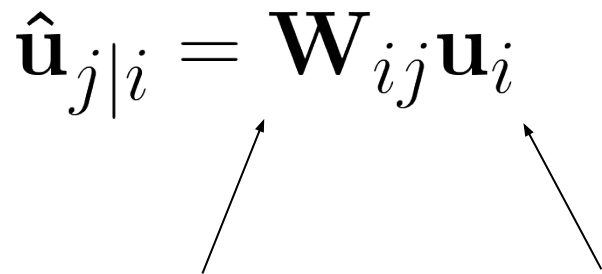
- The output of lower layer capsules is routed to a higher layer capsule only if there is an agreement between both
- An agreement is reached if lower parts predict similar pose and presence about the higher capsule as the higher capsule do



- The body and head capsule have an agreement about the women capsule
- The output of triangle and hexagon capsule is routed to the women capsule

Routing by agreement

The lower capsule i predicts the pose and presence of part j under the assumption, that the lower part i is part of the higher part j

$$\hat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij} \mathbf{u}_i$$


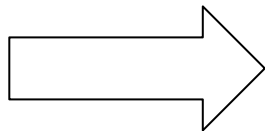
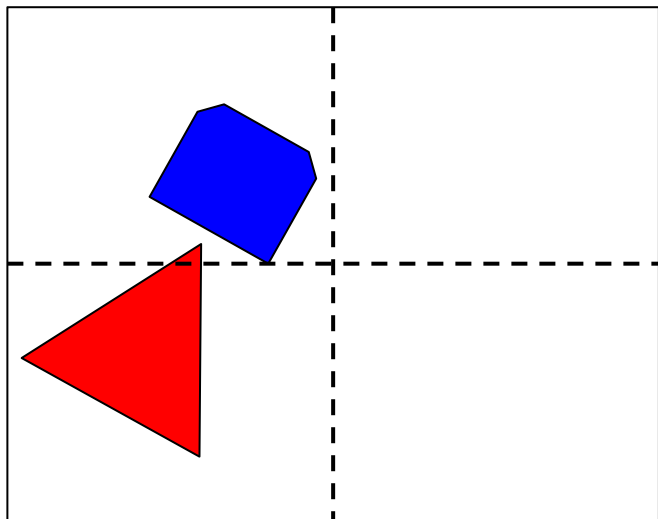
Transformation matrix for
lower part i and higher part j

Input of capsule i

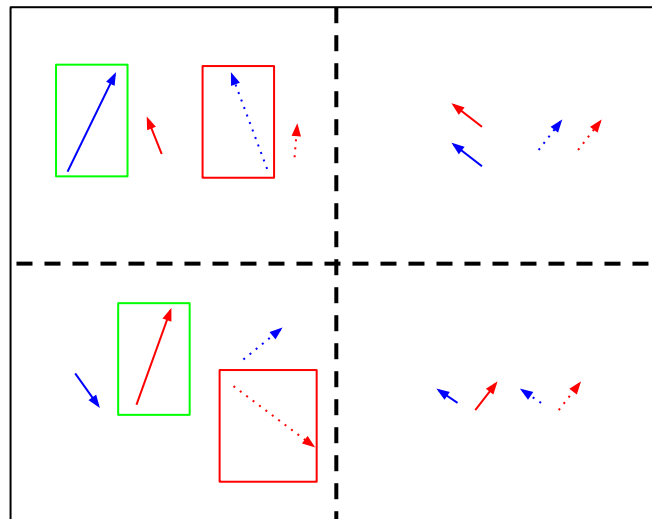
Routing by agreement

---▶ Prediction for church
—▶ Prediction for women

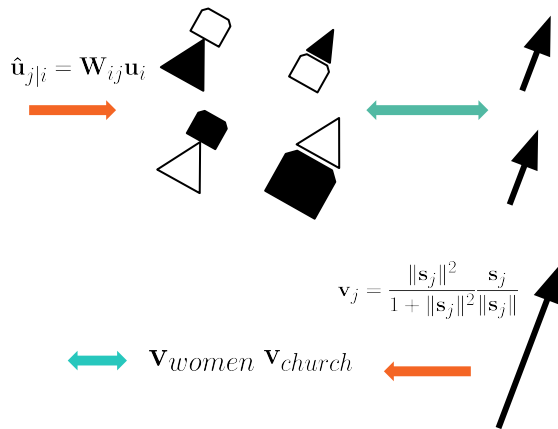
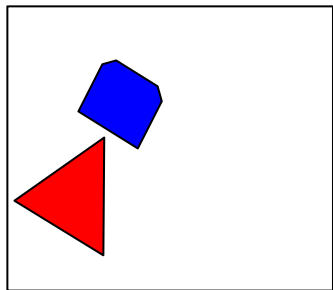
Output of convolutional layer



Activity vector of **hexagon** capsule
Activity vector of **triangle** capsule



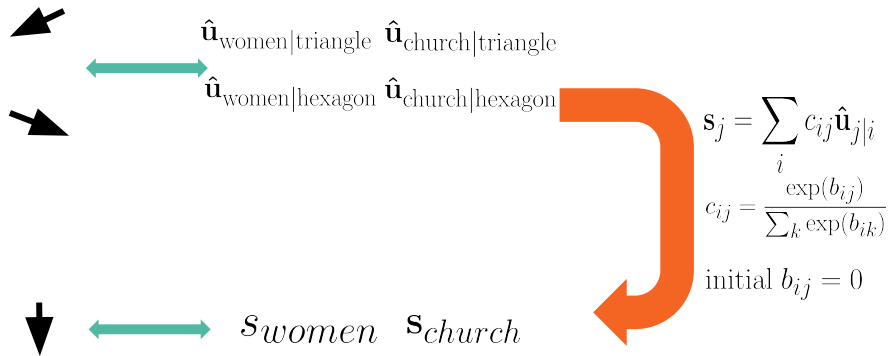
Routing by agreement



$$b_{ij} = b_{ij} + \hat{\mathbf{u}}_{j|i} \mathbf{v}_j$$

$$\mathbf{b}_{\text{triangle women}} = \mathbf{b}_{\text{triangle women}} + \hat{\mathbf{u}}_{\text{women}|\text{triangle}} \mathbf{v}_{\text{women}}$$

$$\mathbf{b}_{\text{hexagon women}} = \mathbf{b}_{\text{hexagon women}} + \hat{\mathbf{u}}_{\text{women}|\text{hexagon}} \mathbf{v}_{\text{women}}$$



$$\mathbf{b}_{\text{triangle church}} = \mathbf{b}_{\text{triangle church}} + \hat{\mathbf{u}}_{\text{church}|\text{triangle}} \mathbf{v}_{\text{church}}$$


$$\mathbf{b}_{\text{hexagon church}} = \mathbf{b}_{\text{hexagon church}} + \hat{\mathbf{u}}_{\text{church}|\text{hexagon}} \mathbf{v}_{\text{church}}$$

Routing by agreement

Total input of capsule j

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

Coupling coefficients can't be calculated in first iteration step



Output of capsule j

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

Routing by agreement

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

$$b_{ij} = 0$$

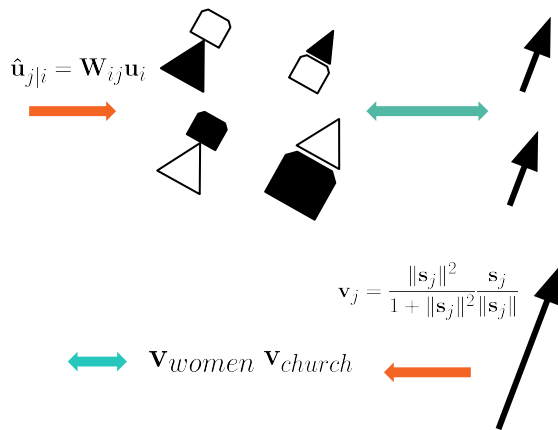
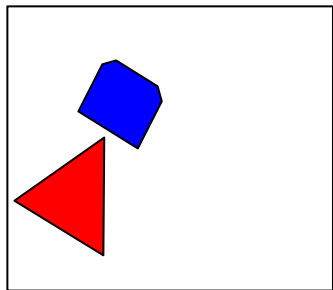


log prior probabilities that capsule i should be coupled to capsule j

Since in the first step no information can be gained from agreements, it is sensible to assume, that each lower capsule has the same probability to be assigned to every higher capsule

$$c_{ij} = \frac{1}{\text{number of capsules in next layer}}$$

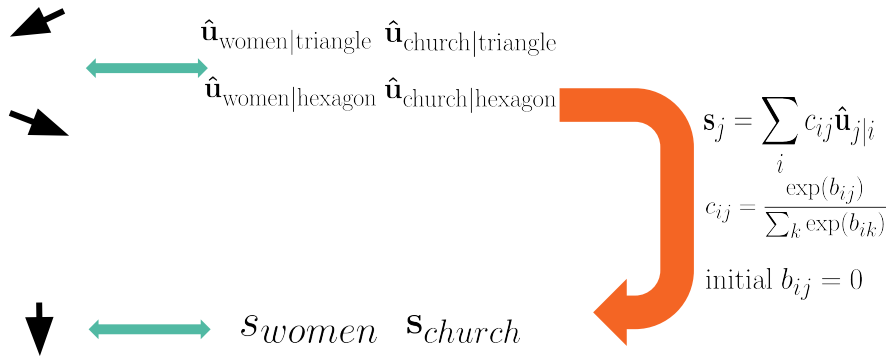
Routing by agreement



$$b_{ij} = b_{ij} + \hat{\mathbf{u}}_{j|i} \mathbf{v}_j$$

$$\mathbf{b}_{\text{triangle women}} = \mathbf{b}_{\text{triangle women}} + \hat{\mathbf{u}}_{\text{women}|\text{triangle}} \mathbf{v}_{\text{women}}$$

$$\mathbf{b}_{\text{hexagon women}} = \mathbf{b}_{\text{hexagon women}} + \hat{\mathbf{u}}_{\text{women}|\text{hexagon}} \mathbf{v}_{\text{women}}$$



$$\mathbf{b}_{\text{triangle church}} = \mathbf{b}_{\text{triangle church}} + \hat{\mathbf{u}}_{\text{church}|\text{triangle}} \mathbf{v}_{\text{church}}$$

$$\mathbf{b}_{\text{hexagon church}} = \mathbf{b}_{\text{hexagon church}} + \hat{\mathbf{u}}_{\text{church}|\text{hexagon}} \mathbf{v}_{\text{church}}$$

Routing by agreement

log prior probabilities that capsule i should be coupled to capsule j

$$b_{ij} = b_{ij} + \hat{\mathbf{u}}_{j|i} \mathbf{v}_j$$

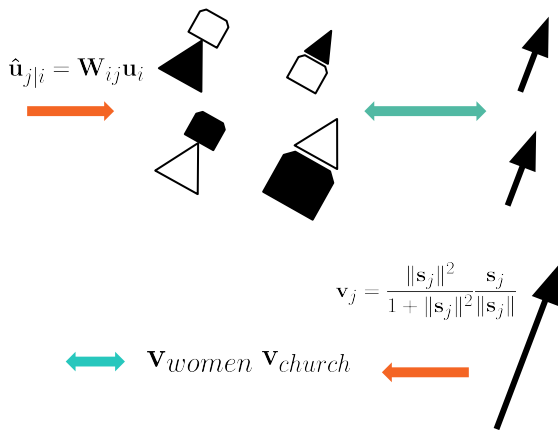
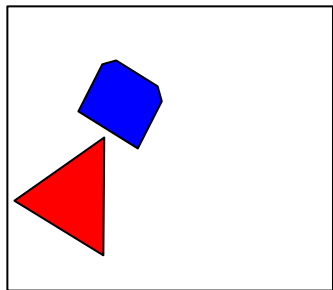
coupling coefficients

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

Agreement between capsule i and j is the scalar product of the prediction from capsule i about the output from capsule j and the output of capsule j

If the agreement between capsule i and capsule j is relatively strong, the coupling coefficient will be greater

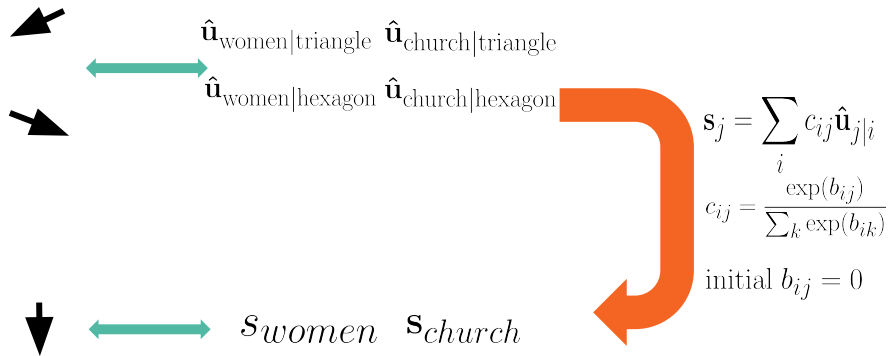
Routing by agreement



$$b_{ij} = b_{ij} + \hat{\mathbf{u}}_{j|i} \mathbf{v}_j$$

$$\mathbf{b}_{\text{triangle women}} = \mathbf{b}_{\text{triangle women}} + \hat{\mathbf{u}}_{\text{women}|\text{triangle}} \mathbf{v}_{\text{women}}$$

$$\mathbf{b}_{\text{hexagon women}} = \mathbf{b}_{\text{hexagon women}} + \hat{\mathbf{u}}_{\text{women}|\text{hexagon}} \mathbf{v}_{\text{women}}$$



$$\mathbf{b}_{\text{triangle church}} = \mathbf{b}_{\text{triangle church}} + \hat{\mathbf{u}}_{\text{church}|\text{triangle}} \mathbf{v}_{\text{church}}$$

$$\mathbf{b}_{\text{hexagon church}} = \mathbf{b}_{\text{hexagon church}} + \hat{\mathbf{u}}_{\text{church}|\text{hexagon}} \mathbf{v}_{\text{church}}$$

Routing by agreement

- In our case the **coupling coefficients** between the lower capsules ...

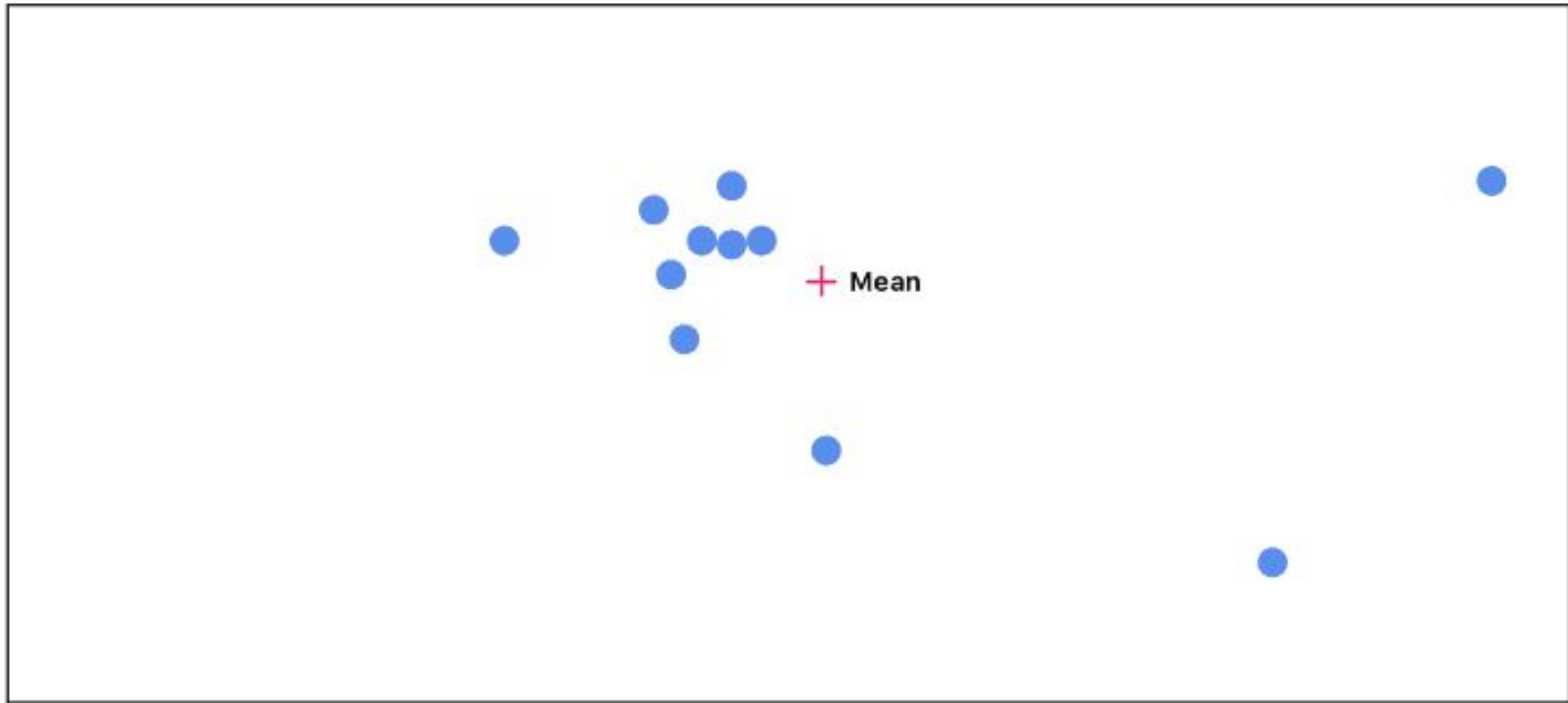
 - ... and the **church** capsule **decrease**

 - ... and the **women** capsule **increase**

with each iteration step

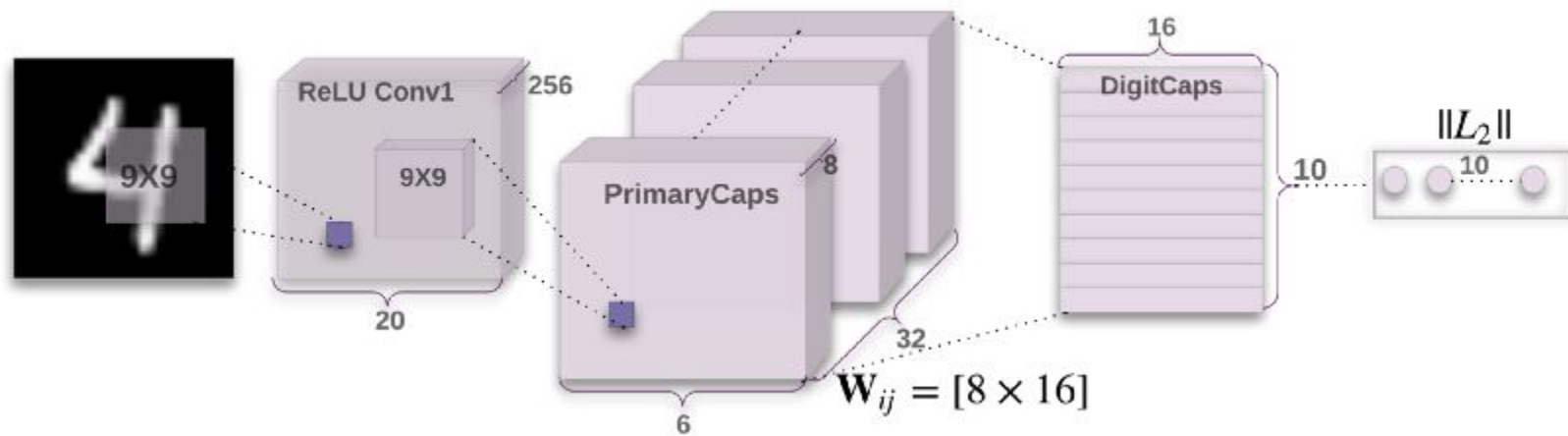
- After a sufficient number of steps, the **coupling coefficients** between the lower capsules and the ...
 - ... **church** capsule **vanishes**
 - ... **women** converges to **1**

-> Winner takes it all

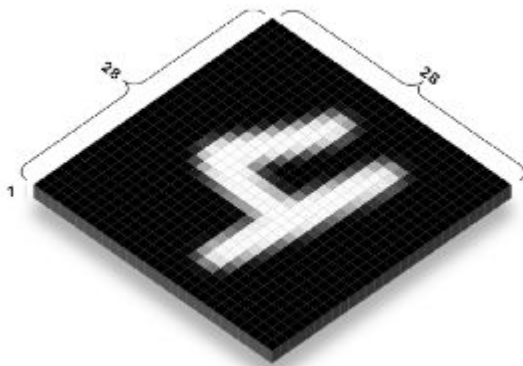


CapsNet Architecture

CapsNet



The Input



[28x28x1]

94	178	124	90	131	0
23	94	135	147	94	138
153	120	140	73	162	6
72	64	10	124	56	64
3	60	75	82	86	129
116	92	165	106	170	89

[0-255]

Convolutions

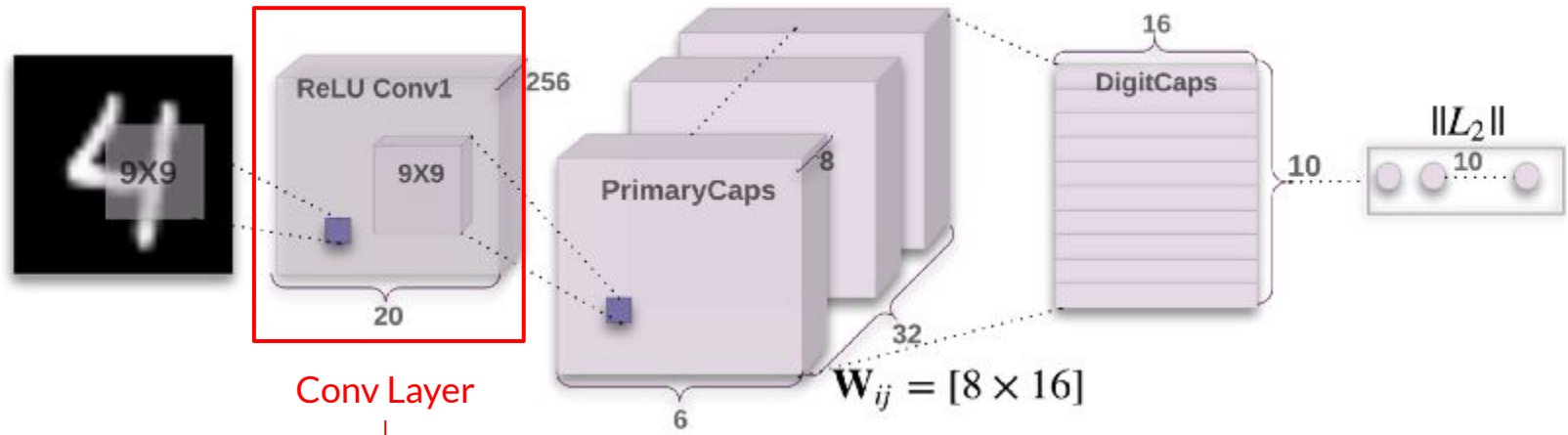
255 x 0	255 x 0	255 x 0	114	114	114	114
255 x 1	255 x 0	255 x -1	114	114	114	114
255 x 0	255 x 0	255 x 0	114	114	114	114
255	255	255	114	114	114	114
255	255	255	255	255	255	255
255	255	255	255	255	255	255
255	255	255	255	255	255	255

0				

$$\lceil (\text{input_size} - \text{kernel_size} + 2 * \text{padding}) / \text{stride} \rceil + 1$$

$$\lceil (28 - 9 + 2 * 0) / 1 \rceil + 1 = 20$$

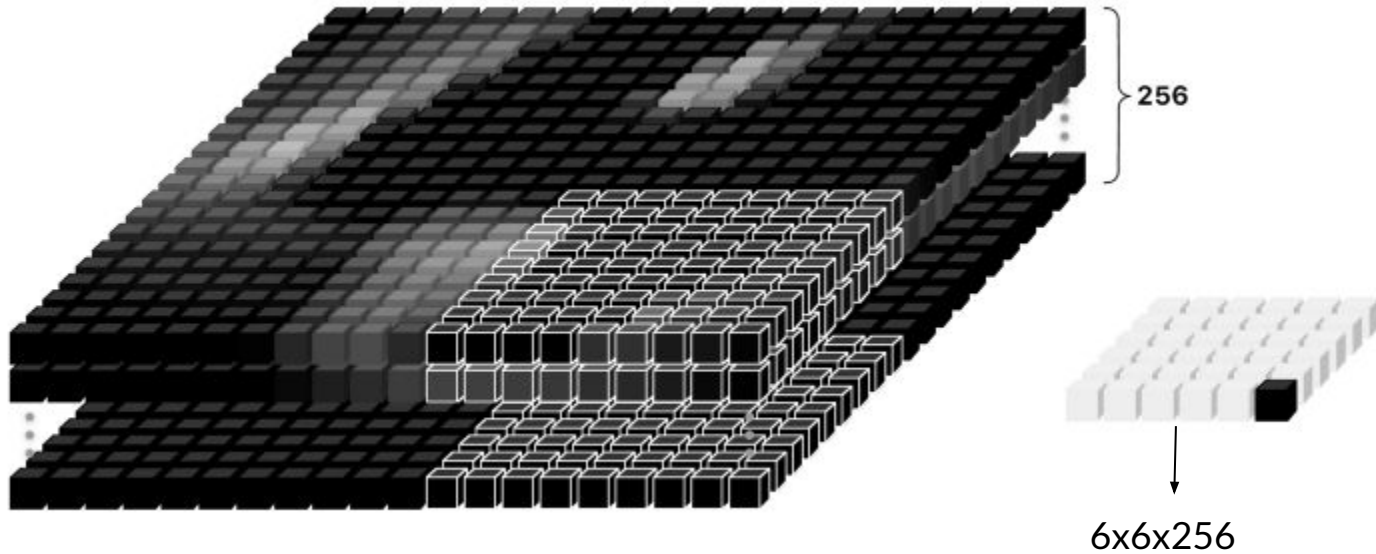
CapsNet-Convolution Layer



Conv Layer

Input: 28x28x1
Output: 20x20x256

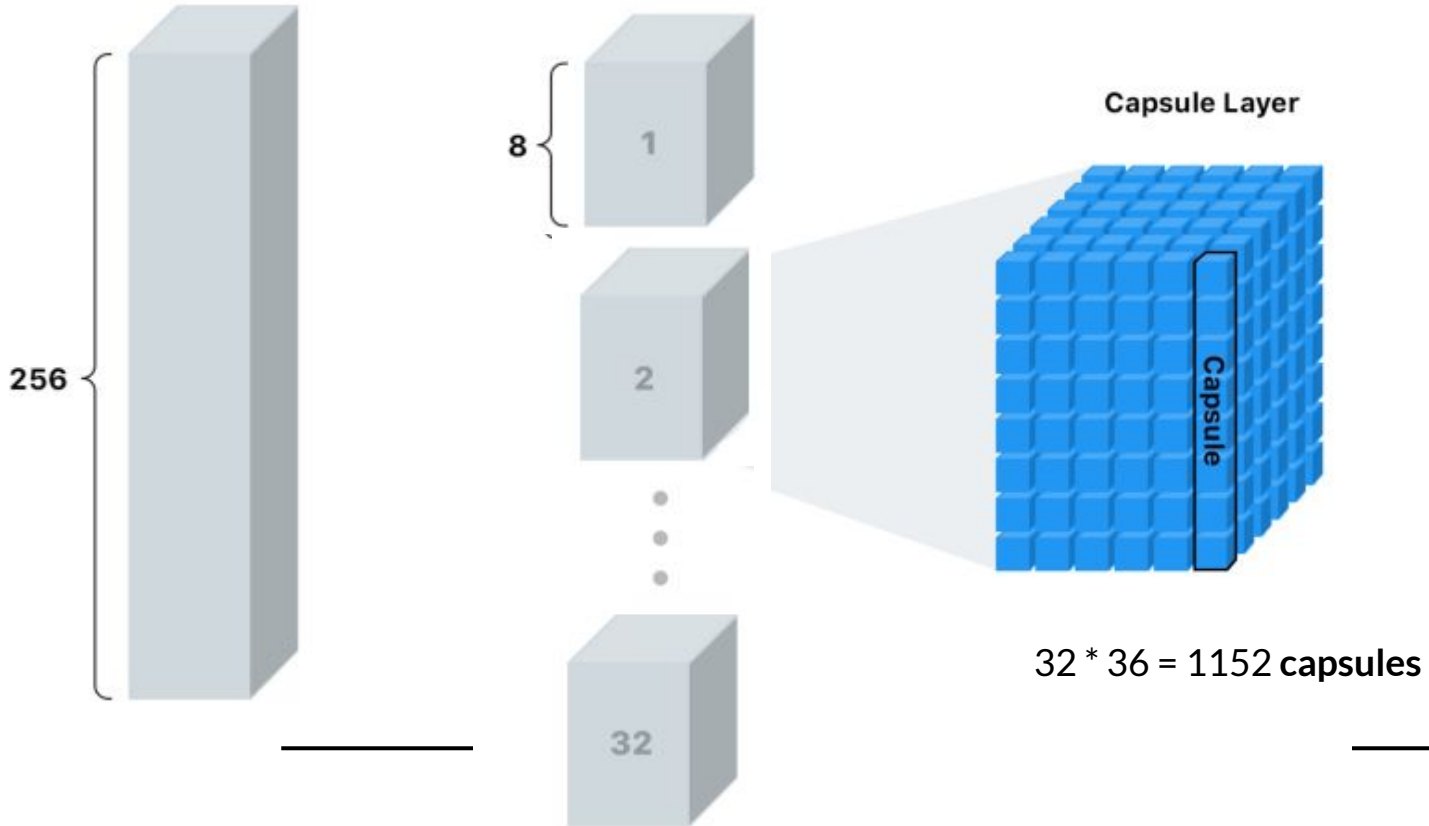
PrimaryCaps



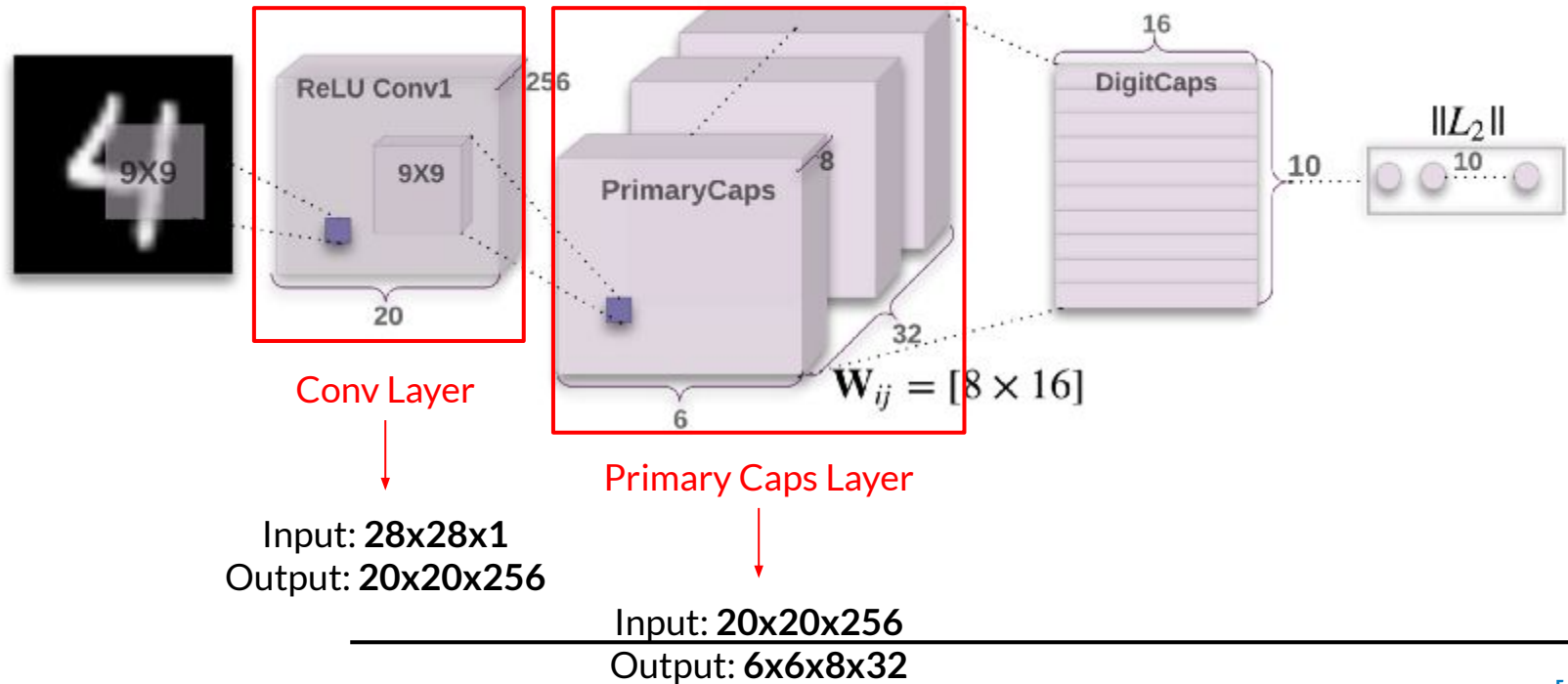
$$\lceil [(20 - 9 + 2 * 0) / 2] + 1 \rceil = 6.5$$

6x6x256

PrimaryCaps cont'd



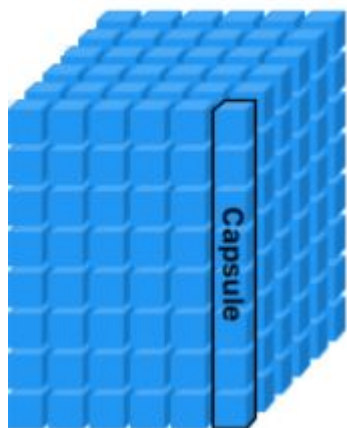
CapsNet-PrimaryCaps Layer



DigitCaps

Capsule Layer

1152 capsules * 10 classes = 11.520 predictions in total

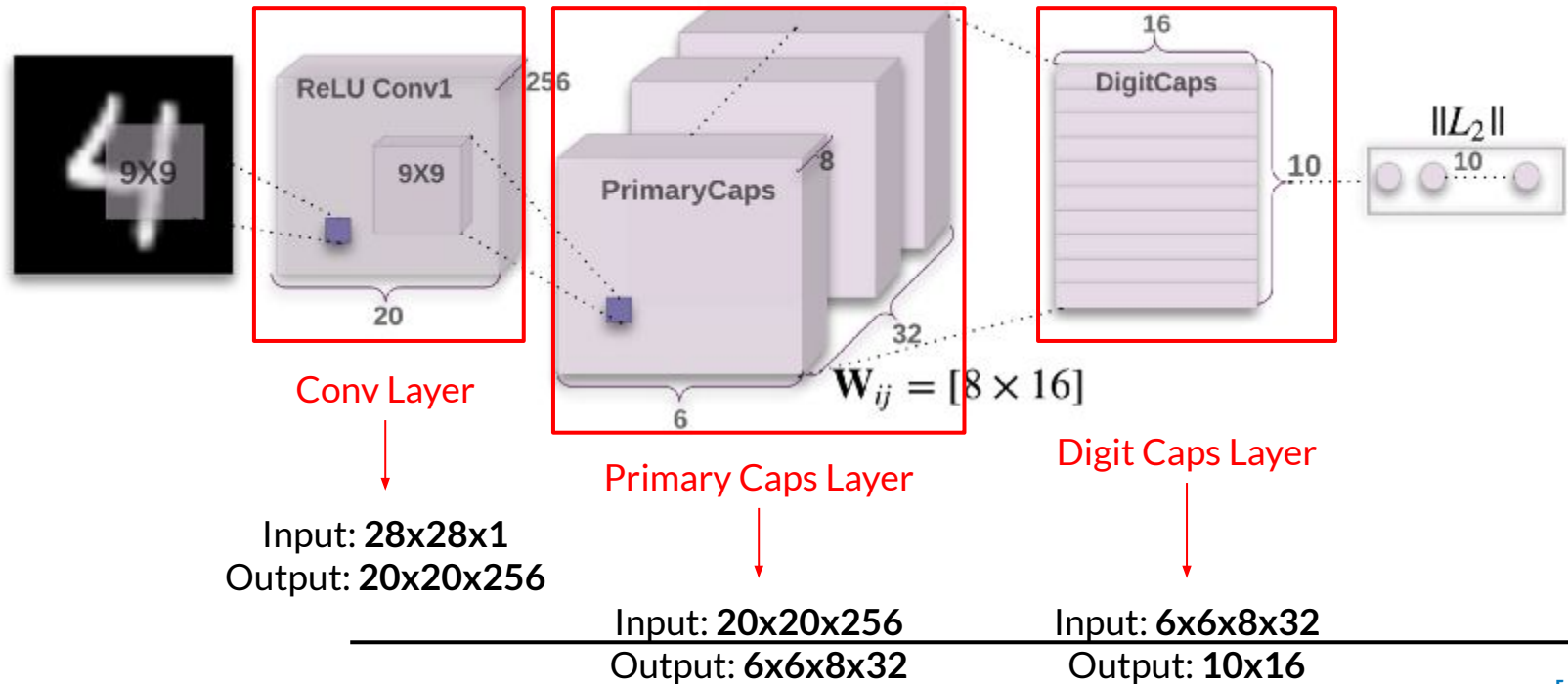


1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0

[1 2 3 4 5 6 7 8] [9 2 3 4 5 6 7 8 0 0 0 0 0 0 0 0 0 0 0 0]

$$(1 \times 1) + (2 \times 0) + (3 \times 0) + (4 \times 2) + (5 \times 0) + (6 \times 0) + (7 \times 0) + (8 \times 0)$$

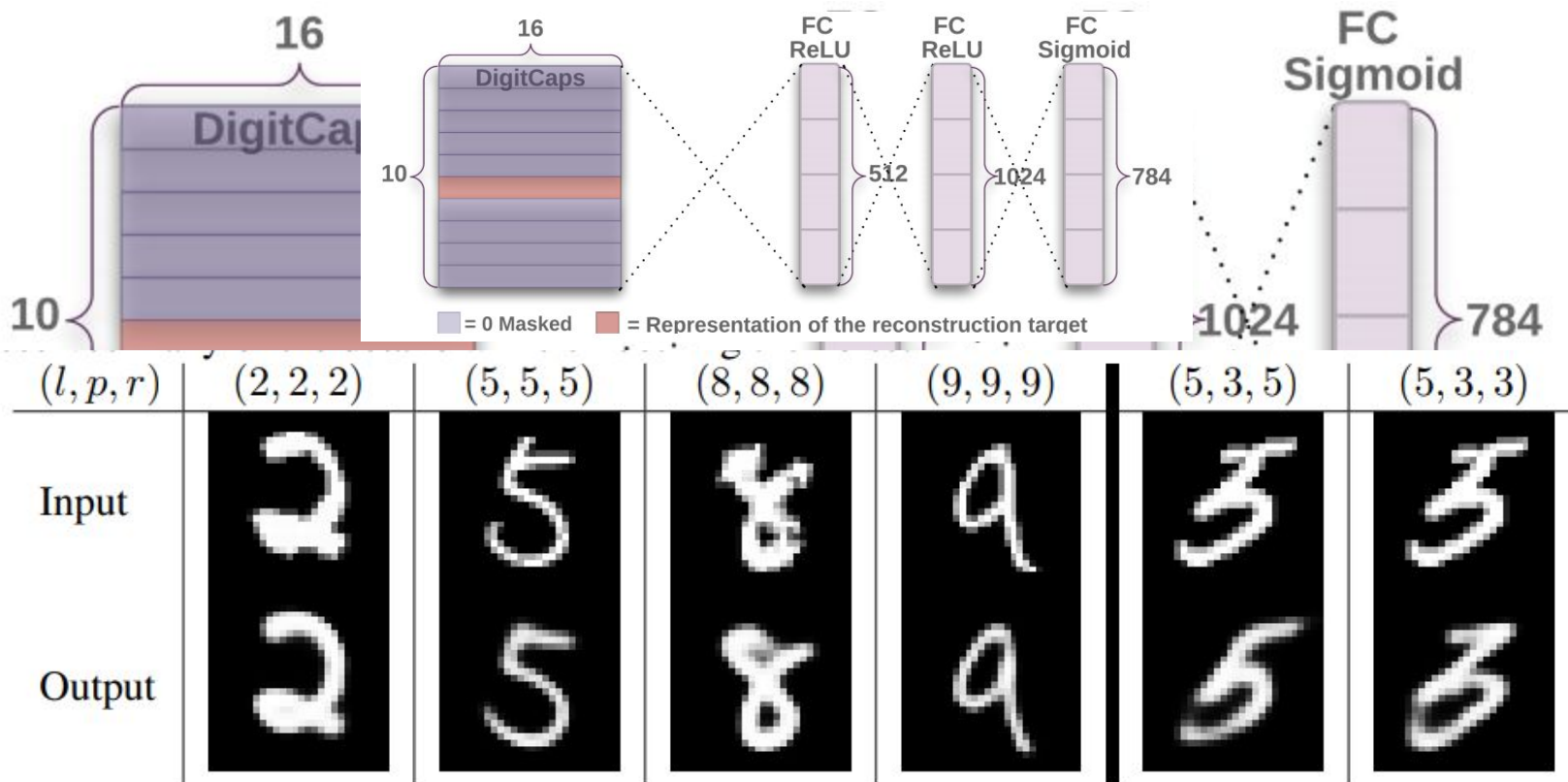
CapsNet-DigitCaps Layer



CapsNet-Predictions



Reconstruction



CapsNet-Loss function

CapsNet Loss Function

loss term for one DigitCap

calculated for correct DigitCap

calculated for incorrect DigitCaps

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

L2 norm

L2 norm

1 when correct DigitCap, 0 when incorrect

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

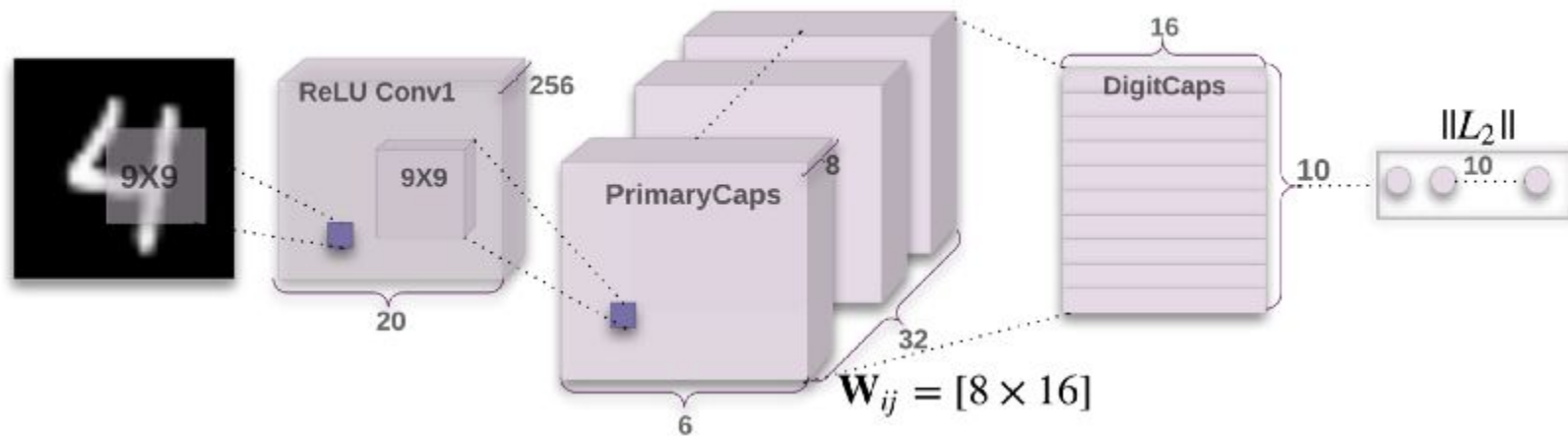
0.5 constant used for numerical stability

1 when incorrect DigitCap, 0 when correct

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

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

CapsNet [reviewed]



References

- Original paper
- BlogPost#1
- BlogPost#2
- BlogPost#3
- Continuation of original paper

Additional Resources

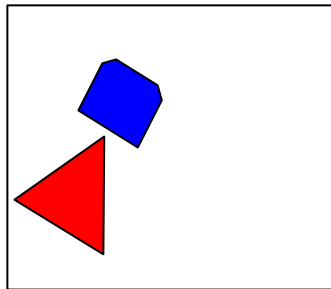
- BlogPost#4
- BlogPost#5
- BlogPost#6
- BlogPost#7
- <https://www.youtube.com/watch?v=pPN8d0E3900&t=763s>
- Zheng et al.

Q&A

—

- What would you say about training time of CapsNets?

Routing by agreement



Calculate the predictions of lower-level parts about the pose and presence of the higher-level parts



Calculate the activity vector of the higher-level parts based on the coupling coefficients and predictions of lower parts

Update coupling coefficients



Squash the output of higher parts

Calculate agreement between lower and higher capsules/parts

Calculate log prior probabilities that lower capsule i should be coupled with higher capsule j based on agreements



Routing by agreement

Procedure 1 Routing algorithm.

```
1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}, r, l$ )
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 computes Eq. 3
5:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$ 
6:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$  ▷ squash computes Eq. 1
7:     for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow b_{ij} + \hat{\mathbf{u}}_{j|i} \cdot \mathbf{v}_j$ 
   return  $\mathbf{v}_j$ 
```

Equations Pool

$$T = \begin{pmatrix} \lambda \cos \theta & \sin \theta \\ -\sin \theta & \lambda \cos \theta \end{pmatrix}$$

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

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

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

$$\hat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij} \mathbf{u}_i$$

$$\lambda \quad \hat{\mathbf{u}}_{i|j} \cdot \mathbf{v}_j$$
