

seoulai.com

## Dynamic Routing Between Capsules.

Geoffrey Hinton et al. (2017)



Dmitriy Khvan

# Intro.

- I am...
- Why Capsule Nets?
- Prerequisites
- My goal:
  - ★ encourage to read the paper
  - ★ give intuitive explanation behind the idea
  - ★ practice public speaking :) and mac Keynote... :(
  - ★ make it fun for everyone!

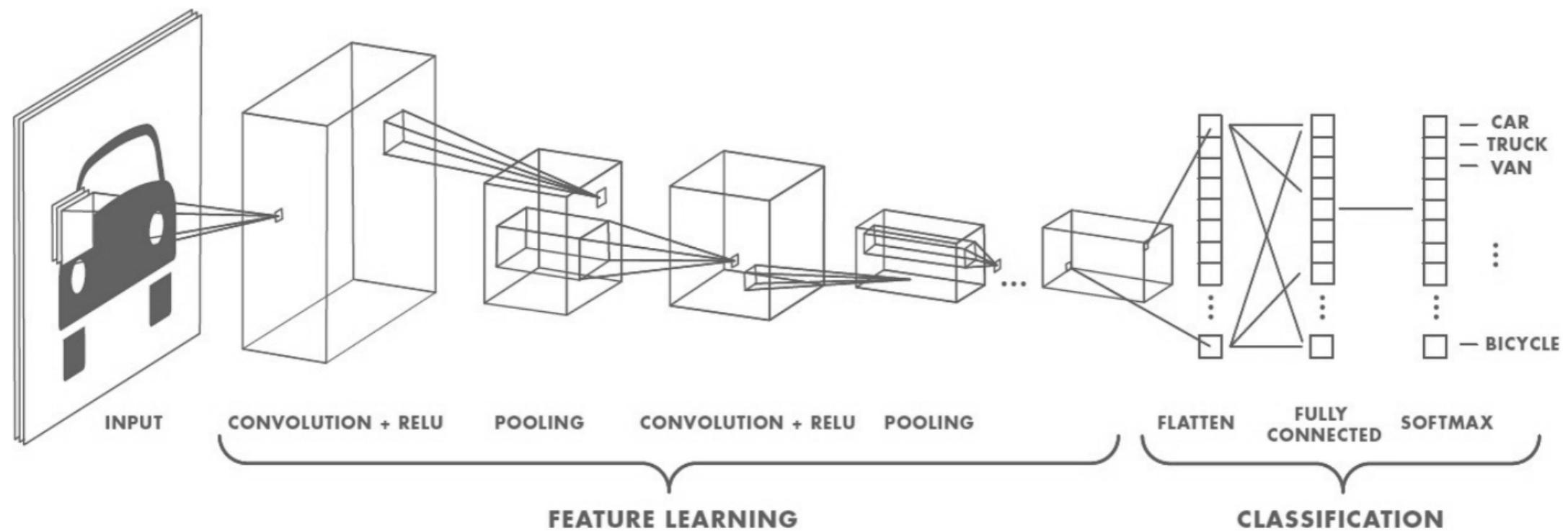
# References.

- ✓ Dynamic Routing Between Capsules G. Hinton et al.
- ✓ Understanding Hinton's Capsule Networks Max Pechyonkin
- ✓ Beginner's Guide to Capsule Networks Zafar (Kaggle)
- ✓ What is wrong with convolutional neural nets? G.Hinton
- ✓ Capsule Networks Tutorial Aurellien Geron

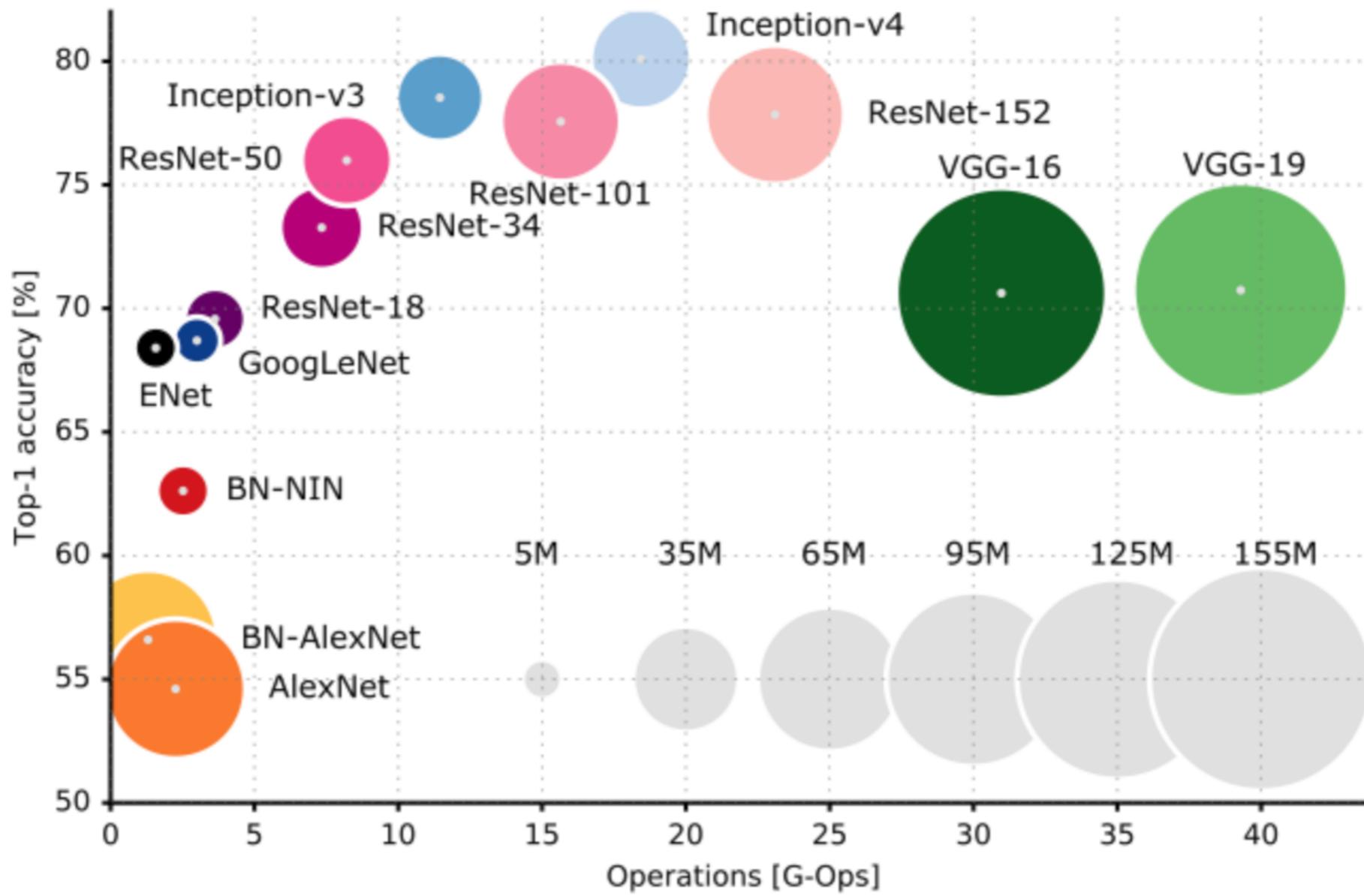
# Motivation.

- Need for intuitive approach closely reflecting HVRS
- Reduce data needed for training
- Robustness to affine transformations
- Need for recognition of overlapping objects

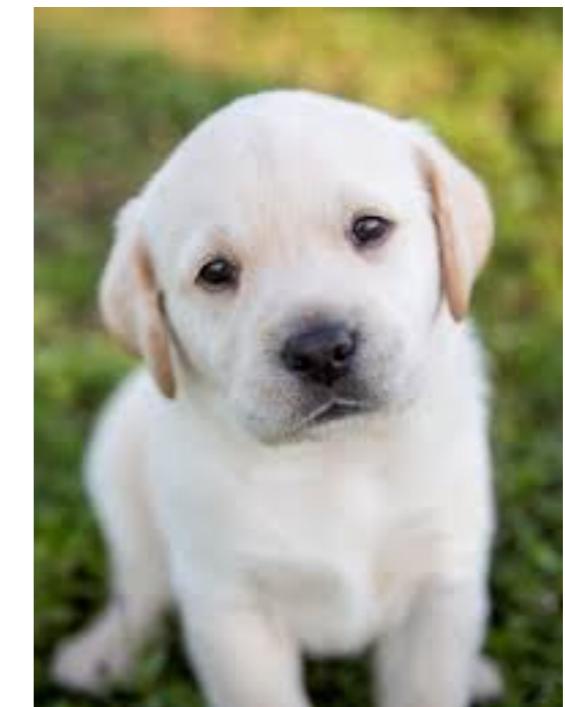
# CNN.



# CNN.

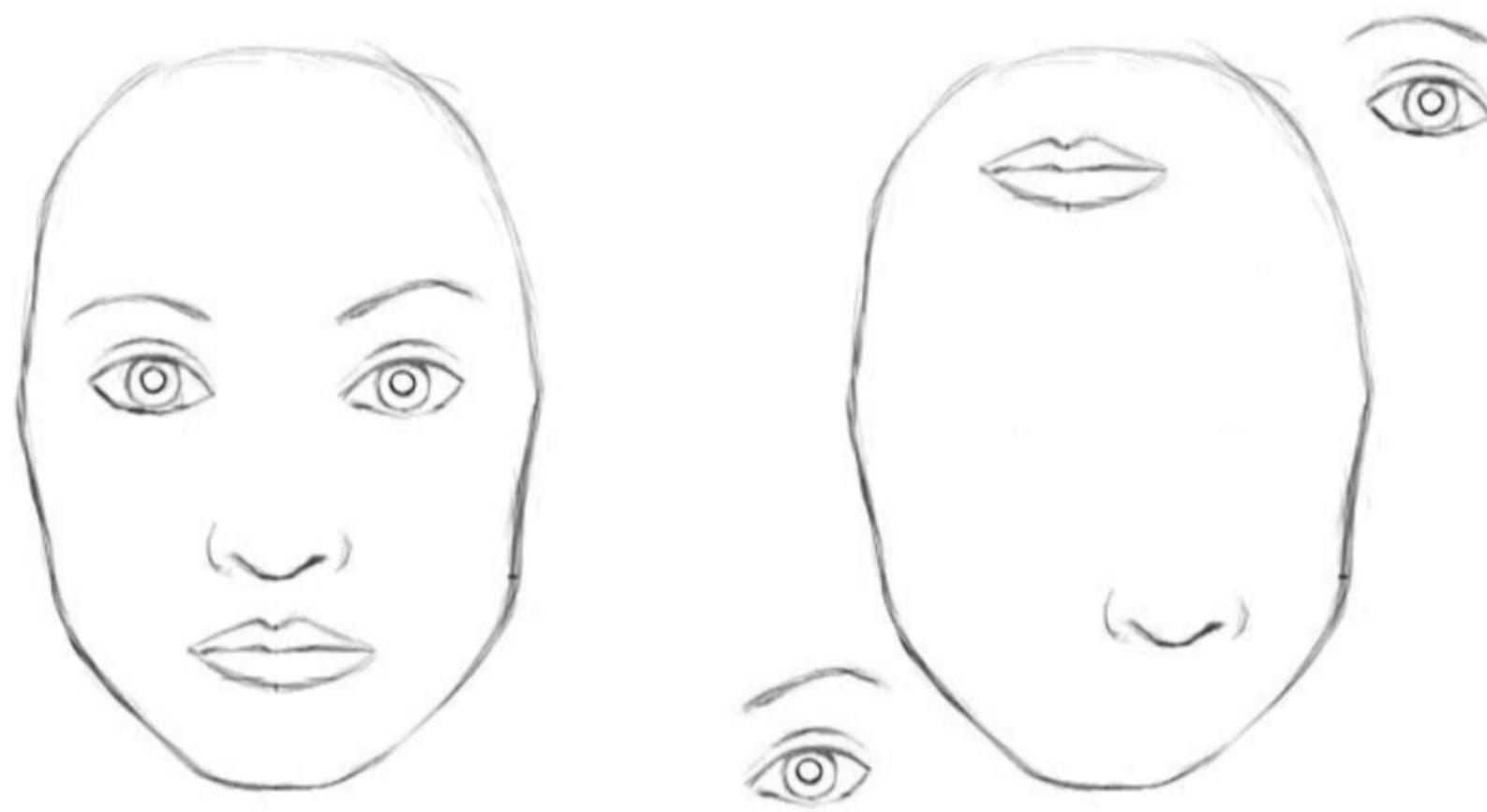


# What is wrong with convolutional neural networks?



CNN fails to recognize variations of an image.

# What is wrong with convolutional neural networks?



Max pooling takes highest activations and **forgets where they come from!** By doing this we lose information about spatial relations between different features.

# What is wrong with convolutional neural networks?

CNNs can be black boxes

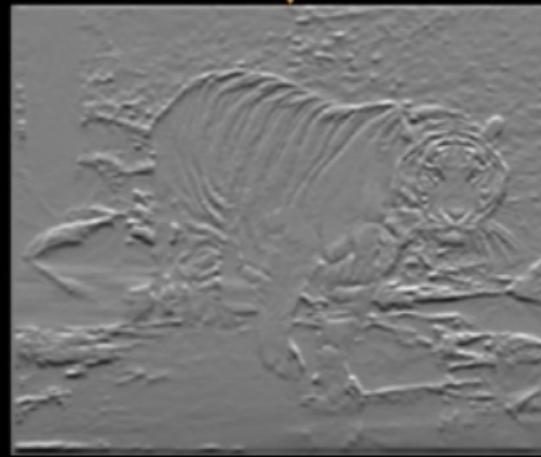
Input



Rotate  
→



Features



Rotate  
→

?

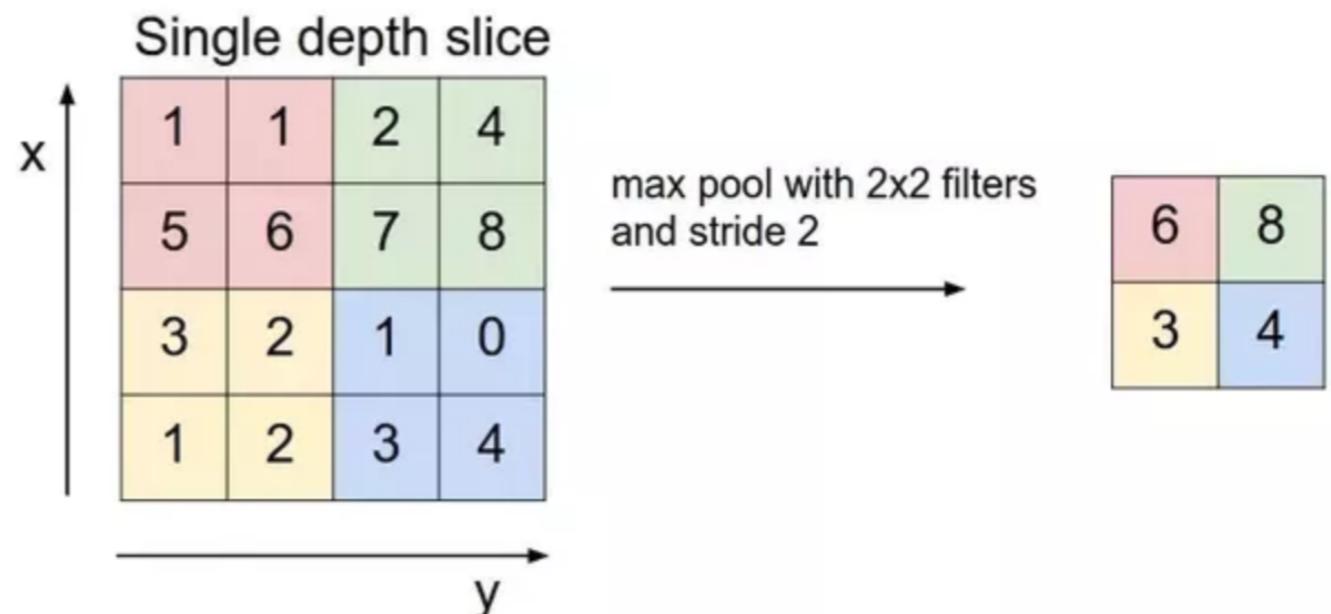
# What is wrong with convolutional neural networks?

## ◆ Too few levels of structure (neuron, layer, nn).

- **Capsule** with 2 instantiation parameters:  
whether an entity present & entity properties.

## ◆ 4 arguments against pooling:

- Does not reflect psychology of **shape perception**;



- We want **equivariance** not invariance;

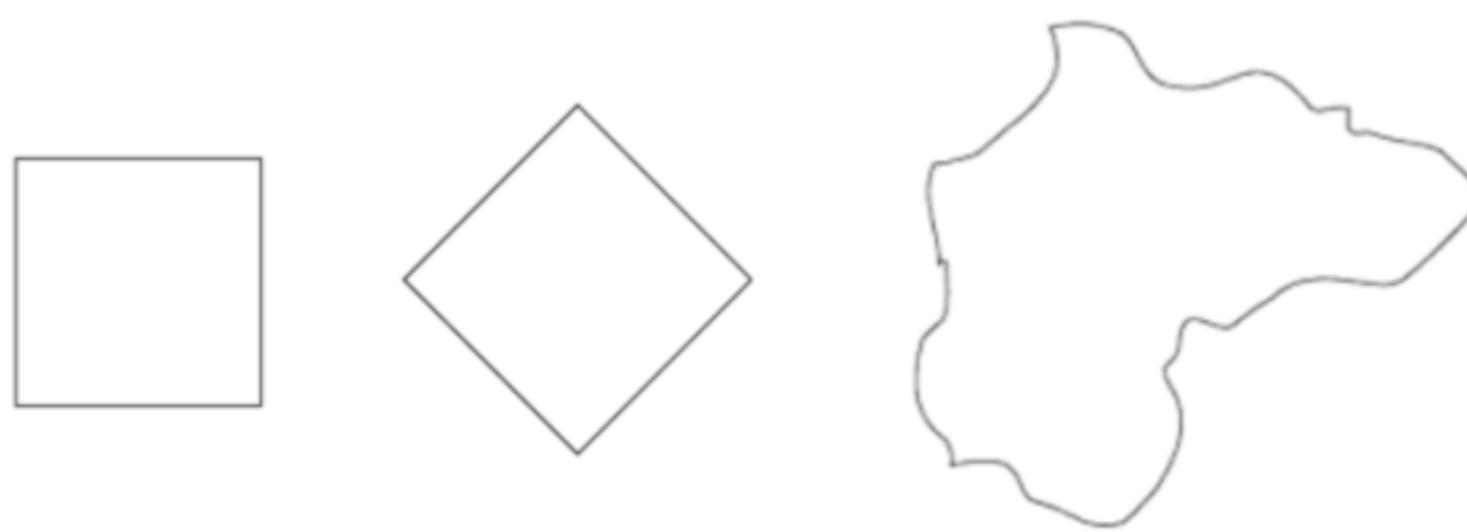
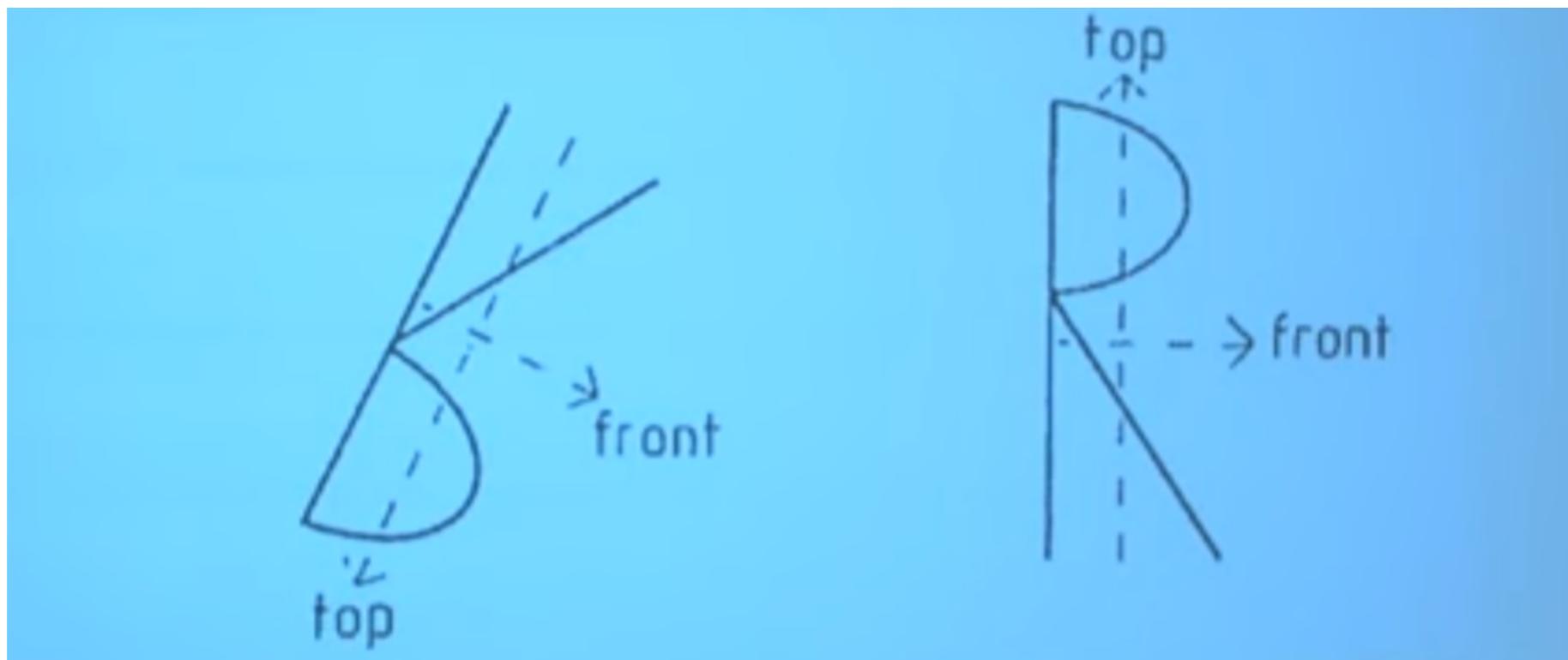
- Ignores underlying **linear structure**;

- Poor way of **routing**;

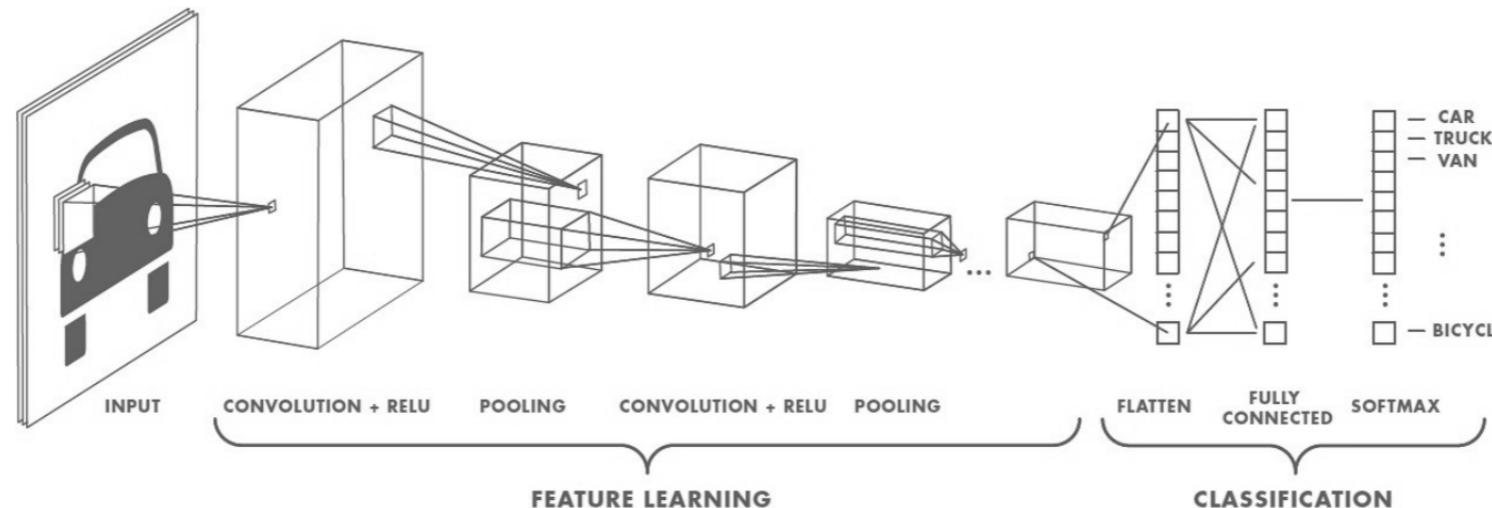
“I believe in convolution but I don't believe in pooling. The fact pooling works so well is a disaster.” G.Hinton.

## Argument 1: CNN does not reflect psychology of shape perception.

It is a bad fit to the psychology of shape perception: It does not explain why we **assign intrinsic coordinate frames to objects** and why they have such huge effects.



## Argument 2: We want equivariance, not invariance.

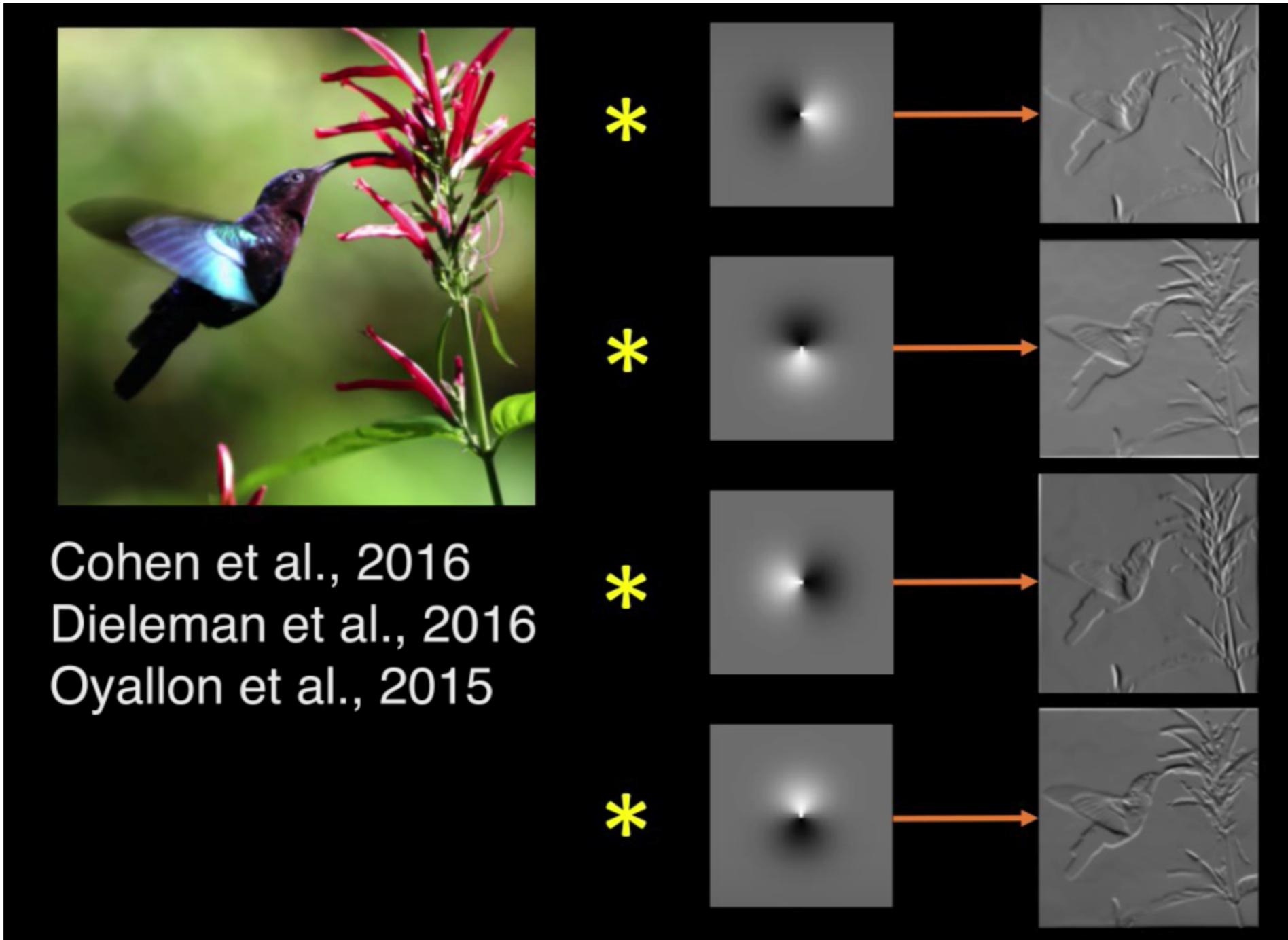


- Translation invariant
- Rotation invariant
- Scale invariant

Several works aim to introduce rotation, scale invariance to CNNs.

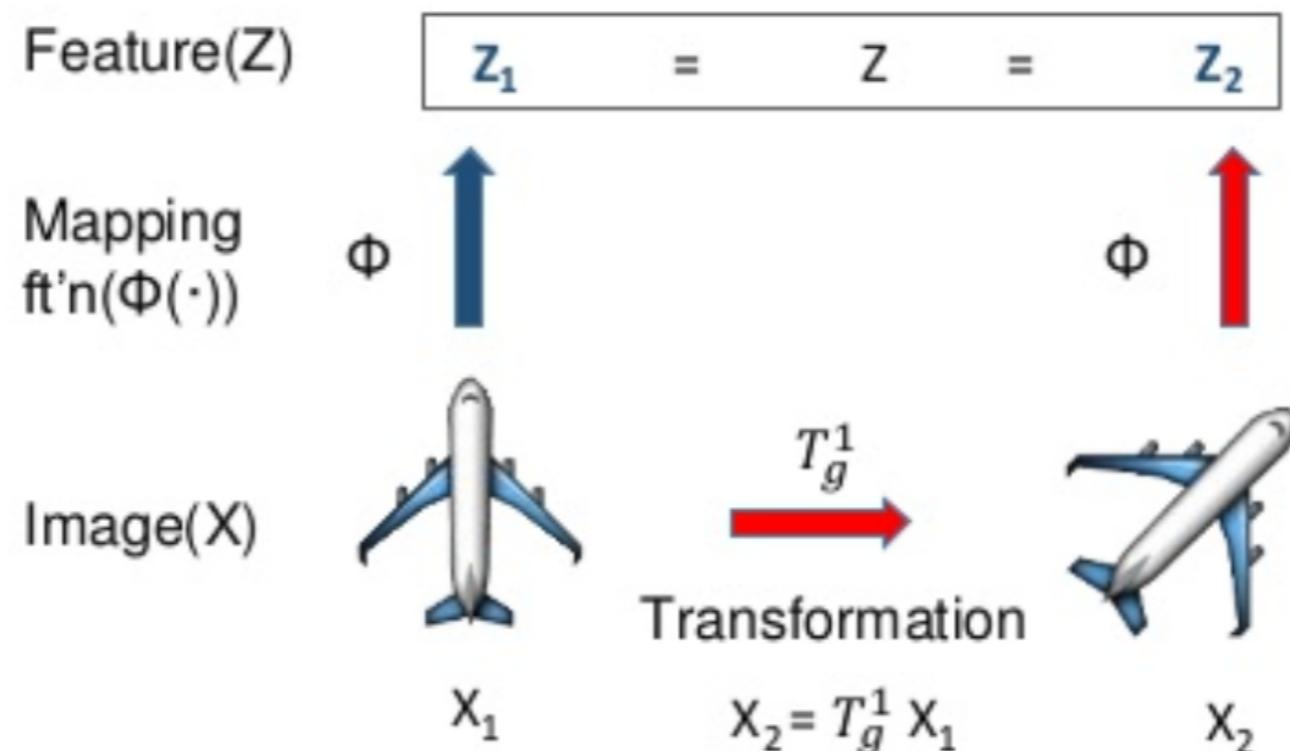
But they **trying to solve wrong problem** (invariance), we want equivariance instead!

# Rotational Invariance.



# Invariance.

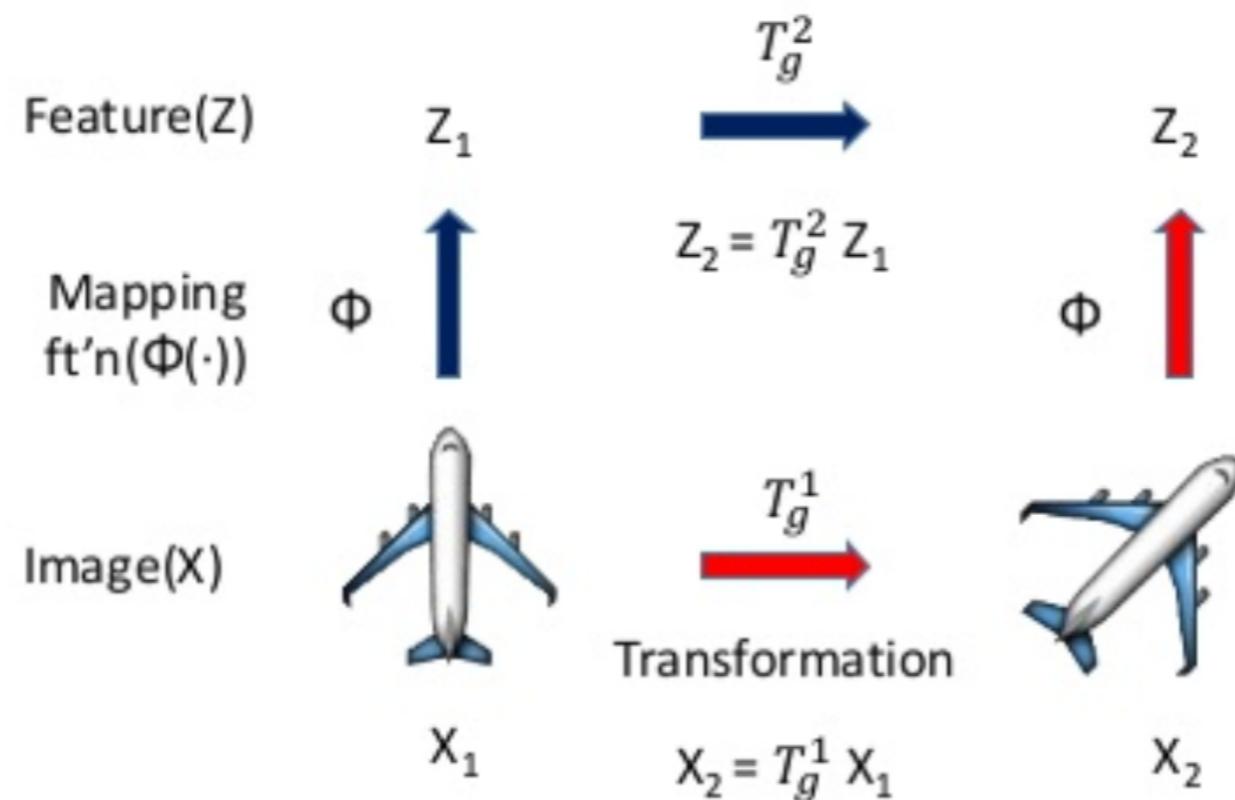
: Mapping independent of transformation,  $T_g$ , for all  $T_g$



$$Z = z_1 = \Phi(x_1) = z_2 = \Phi(x_2) = \Phi(T_g^{-1} x_1)$$

# Equivariance.

: Mapping preserves algebraic structure of transformation



$$z_1 \neq z_2 \text{ but keeps the relationship } z_2 = T_g^2 z_1 = T_g^2 \Phi(x_1) = \Phi(T_g^1 x_1)$$

: Invariance is special case of equivariance where  $T_g^2$  is the identity.

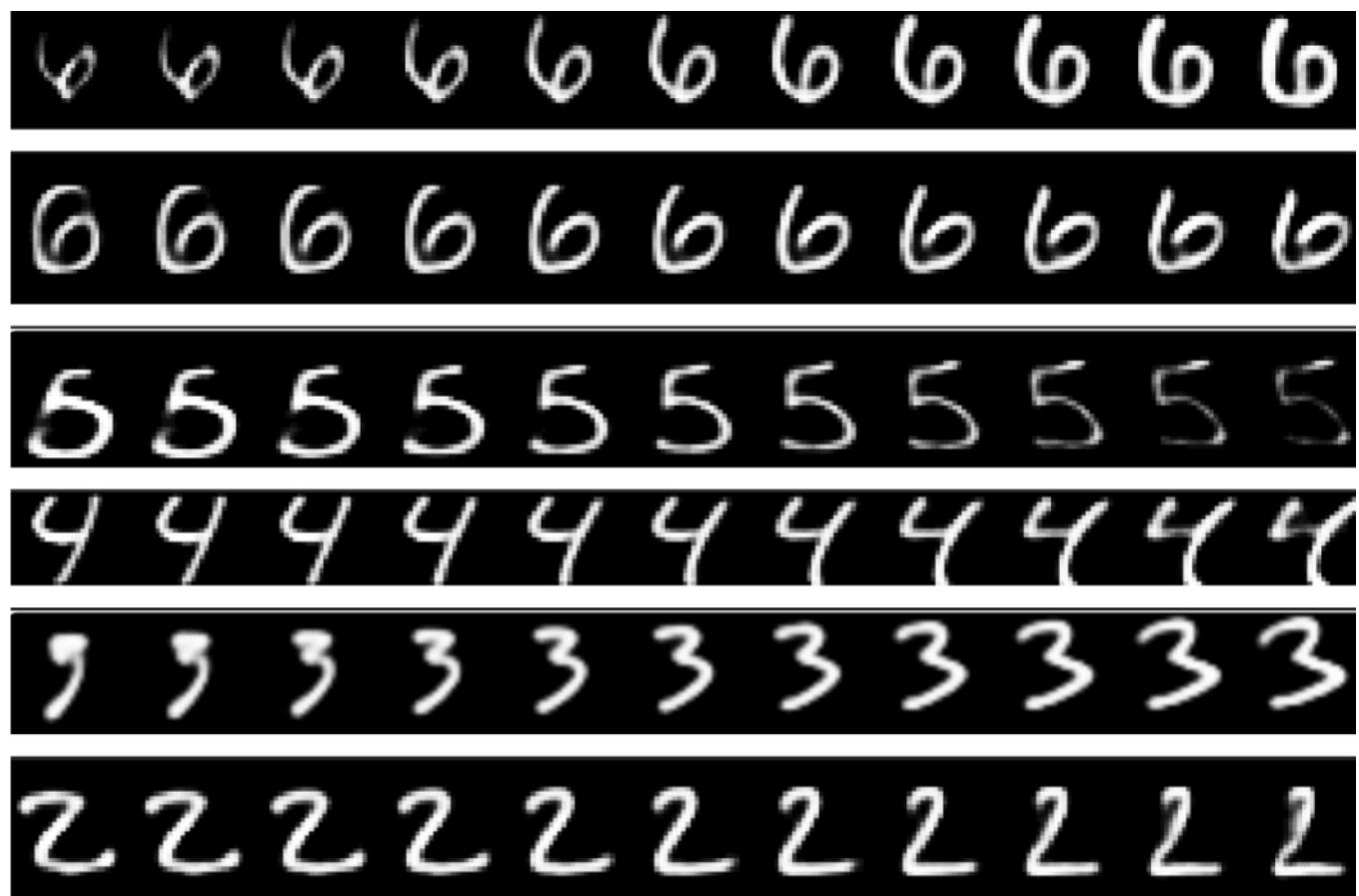
# Viewpoint Equivariance.



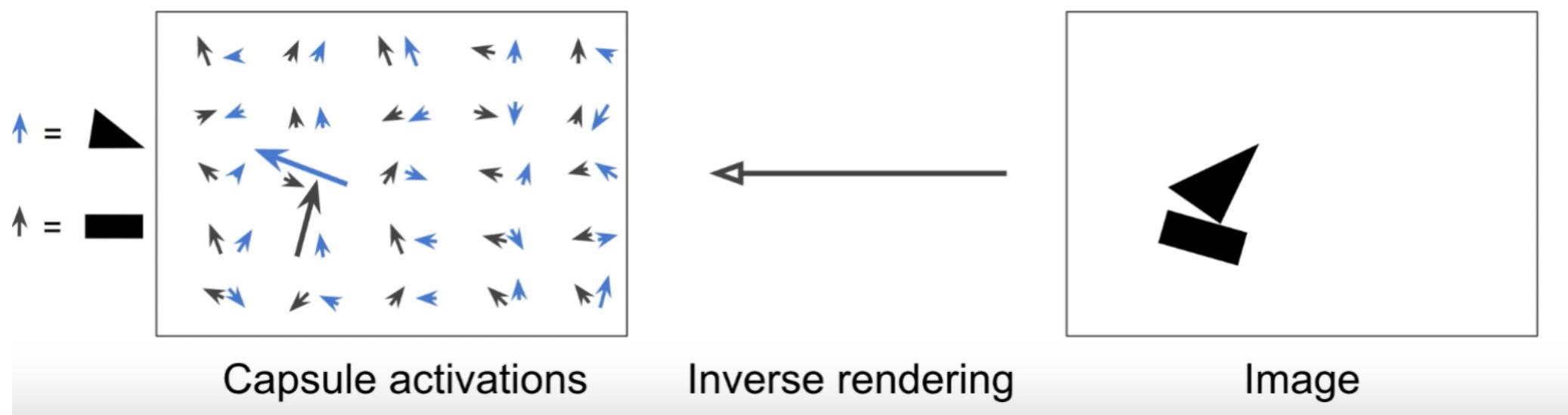
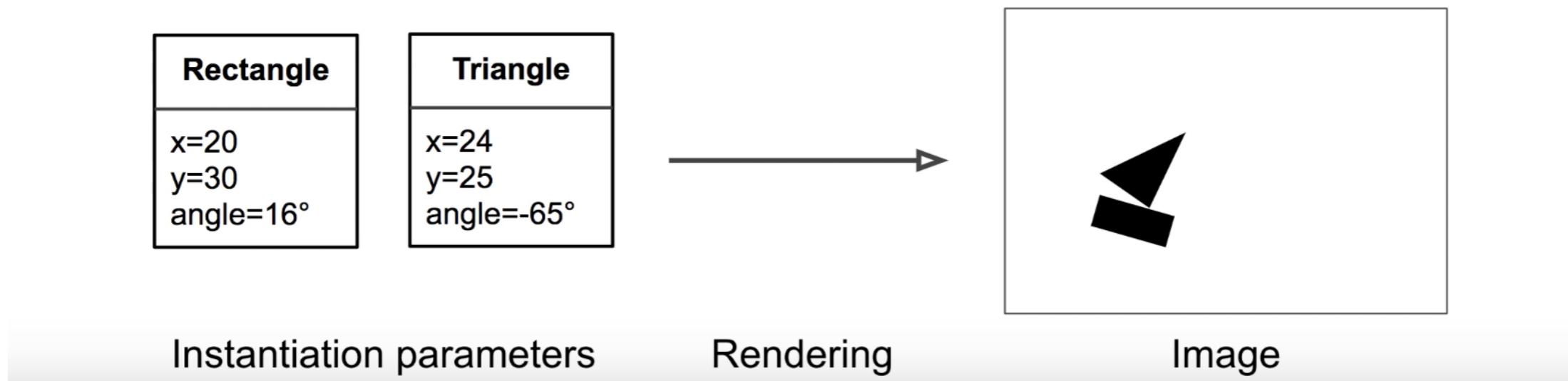
The capsule network is much better than other models at telling that images in top and bottom rows belong to the same classes, only the view angle is different.

## Argument 3: Ignores underlying linear structure.

- CNNs try to conquer the variance of the viewpoint by feeding a lot of various images
- A better way to do that is to transform the image into a space in which the manifold is globally linear.
- Hinton proposed a study called **“inverse-graphics”** in order to reverse the 2D image into the desired space so that we can learn from a small amount of data and manipulate it linearly in that space.



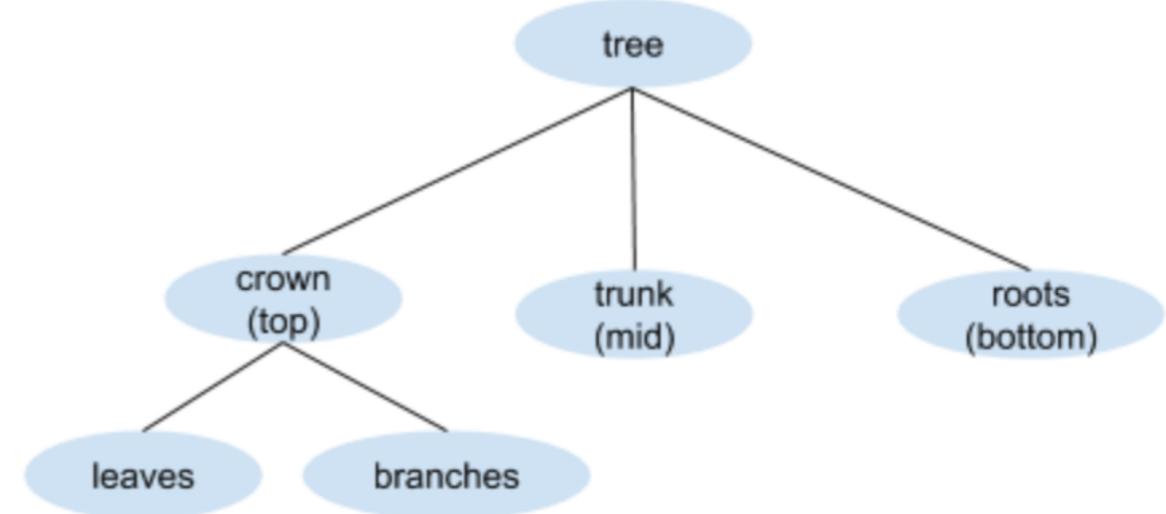
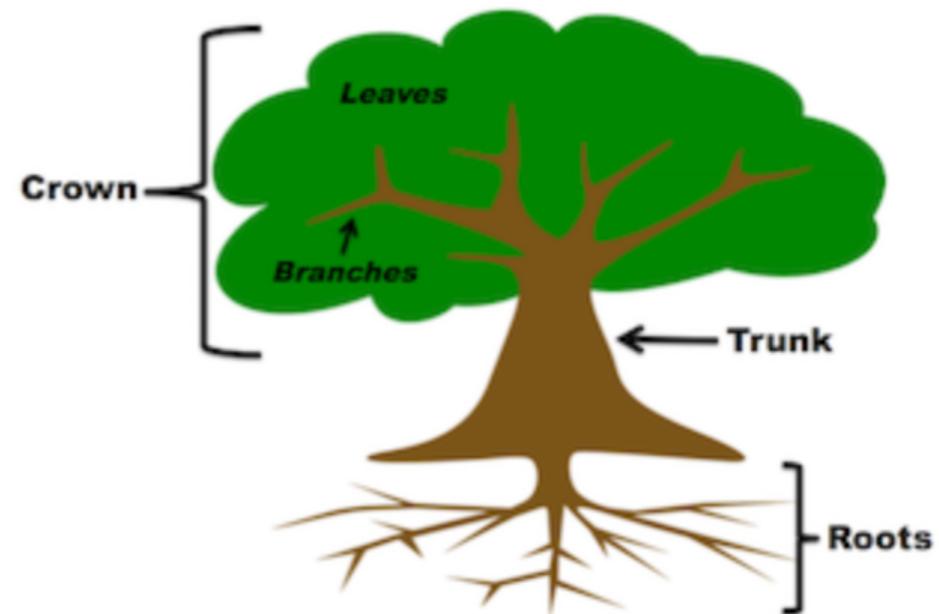
# Inverse Rendering.



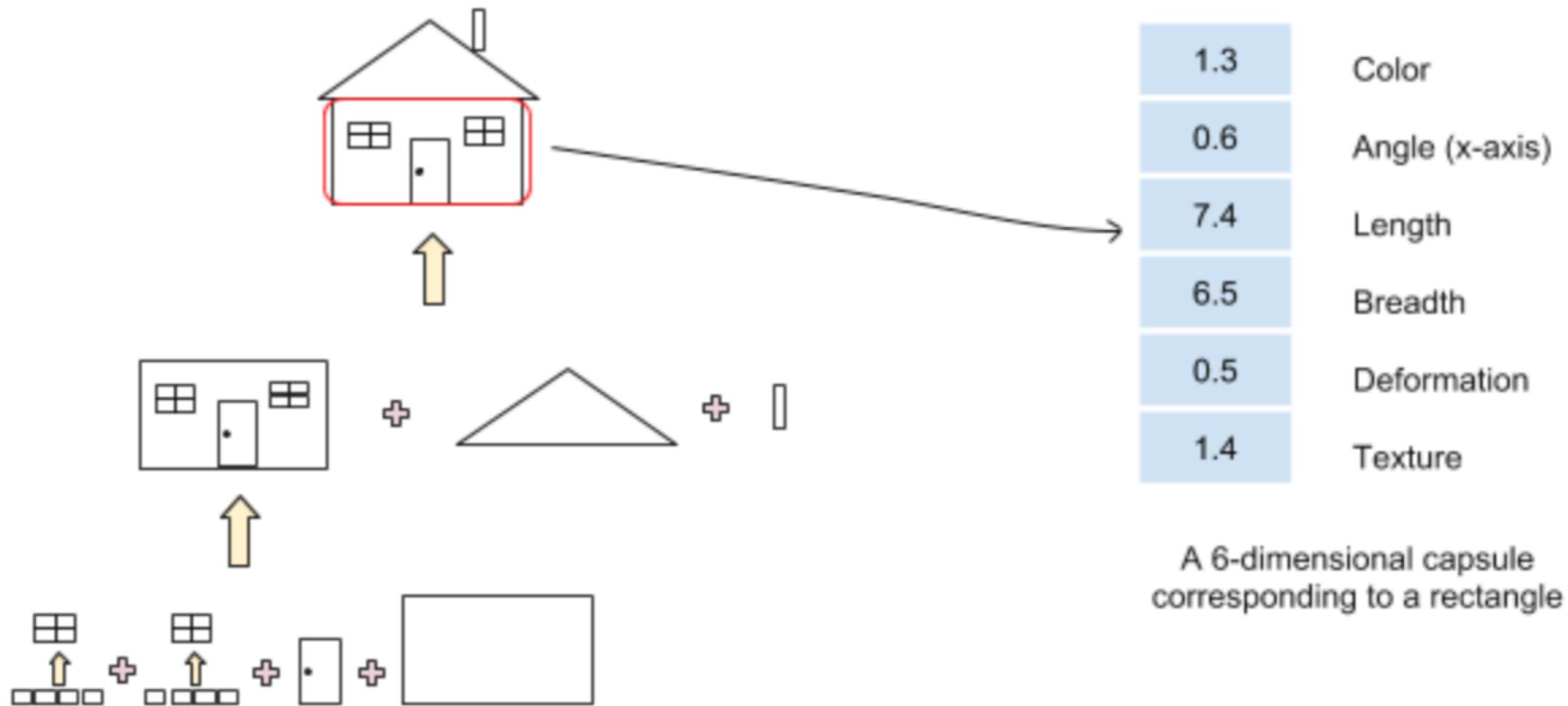
# Human Visual Recognition.

- Tree is composite of crown, trunk and roots;
- crown is composite of leaves and branches;
- HVS recognizes an object (entity) as composition of simpler (more primitive) entities, creating hierarchy of entities.

Q: Does CNN reflect HVS recognition?

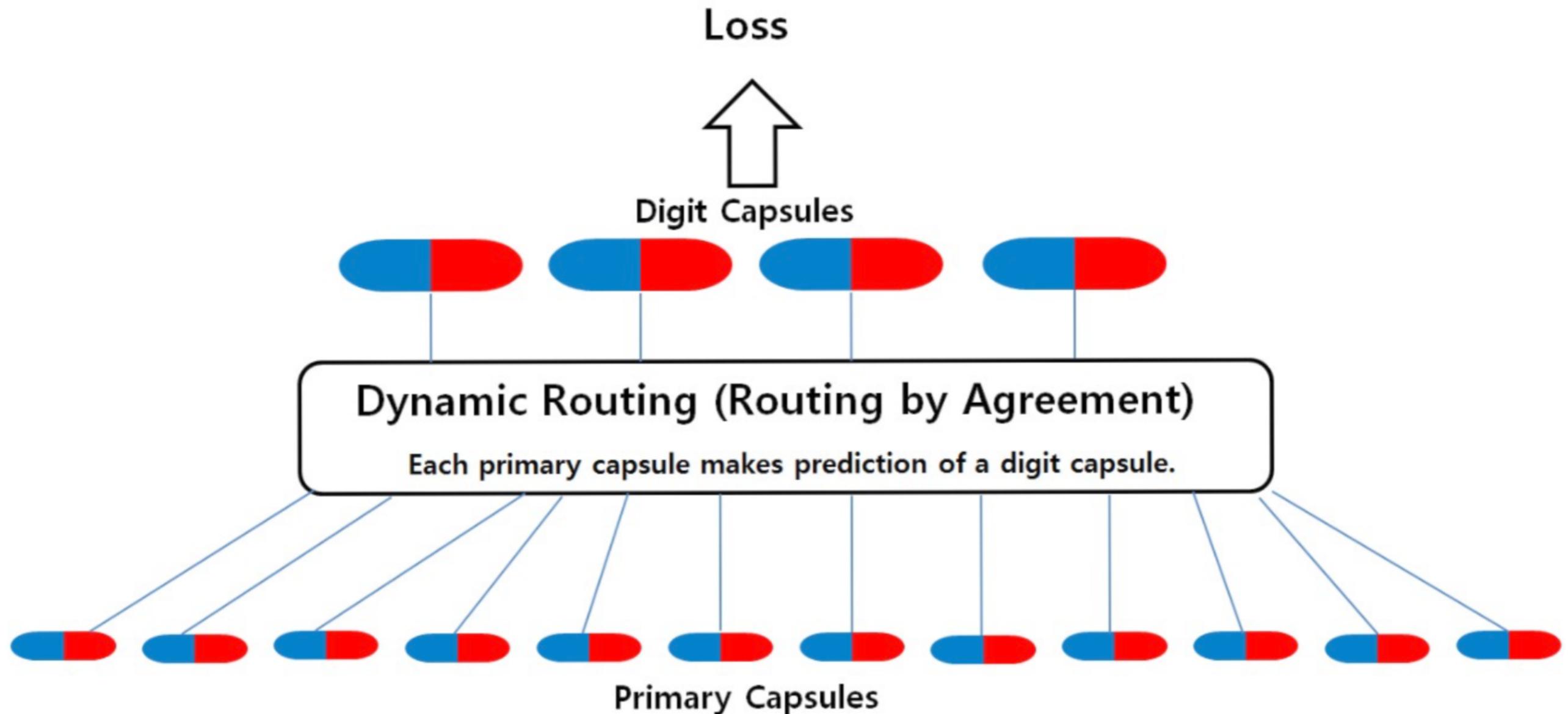


# Human Visual Recognition. Capsule.

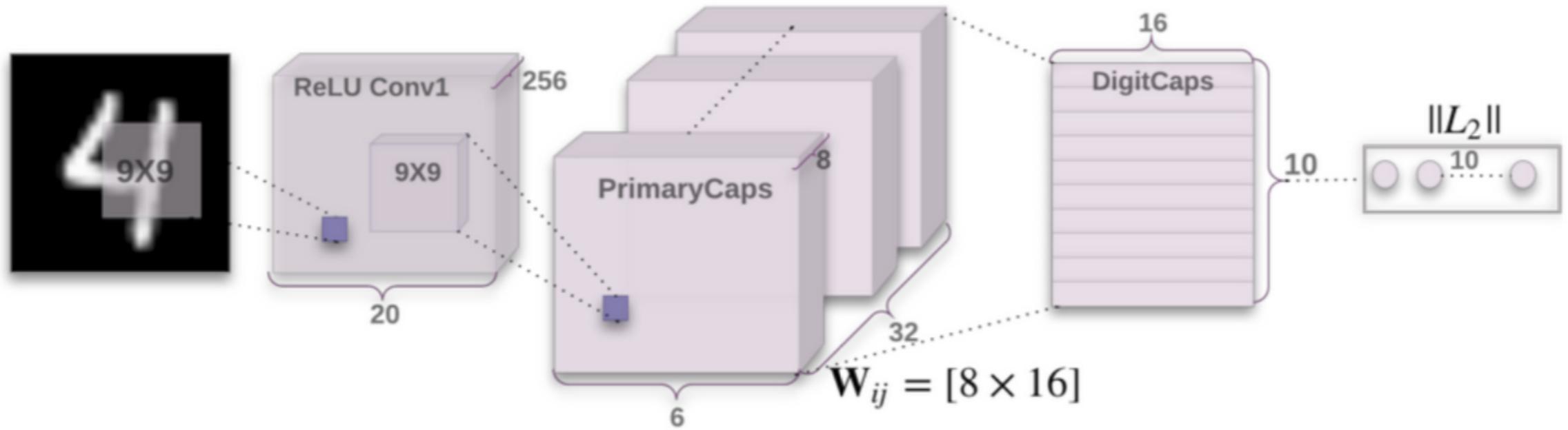


$$\sqrt{1.3^2 + 0.6^2 + 7.4^2 + 6.5^2 + 0.5^2 + 1.4^2} = 10.06$$

# Capsule Hierarchy.



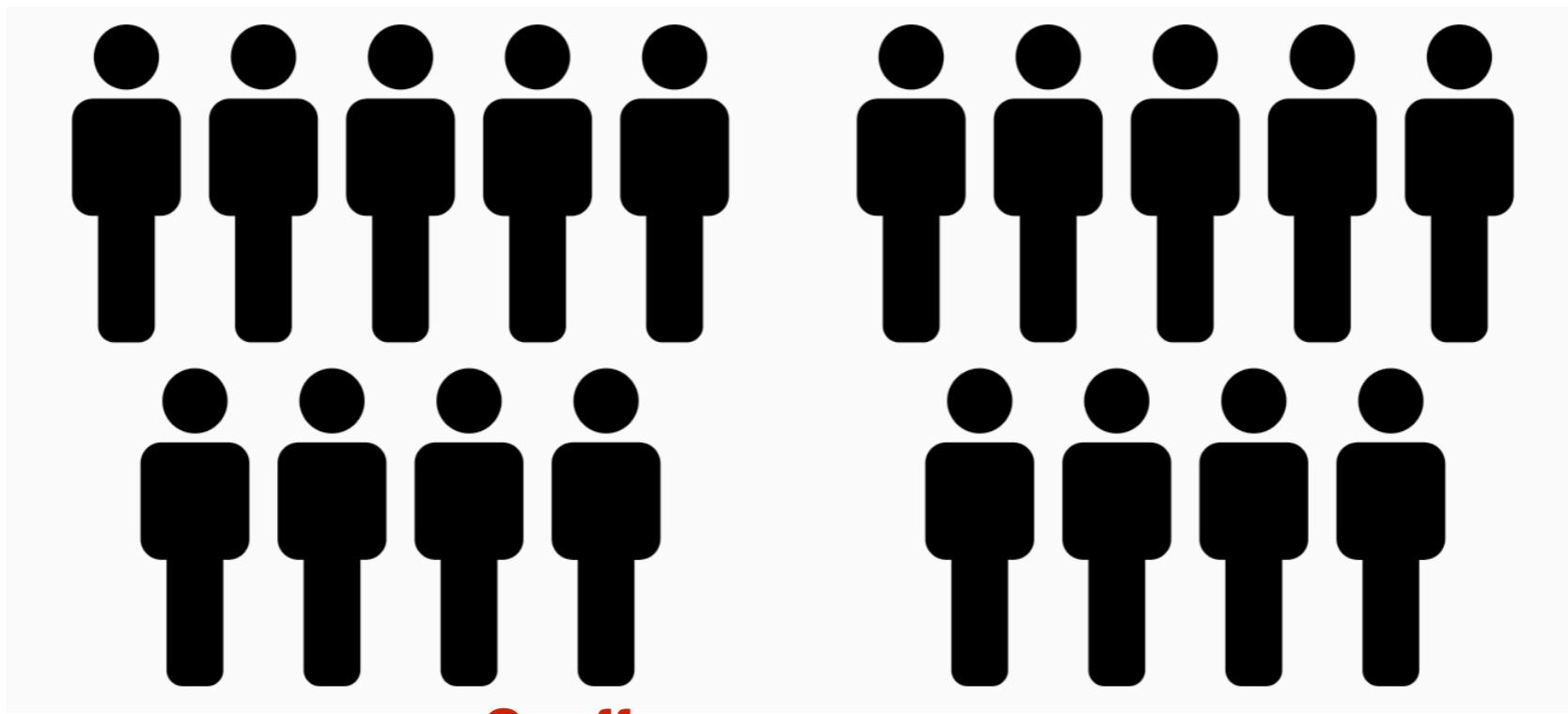
# Capsule Network.



- ▶ Convolution (x 2)
- ▶ Reshape feature maps to 32 groups of 8 feature maps each of size 6 by 6  
**(6x6x32=1152 primary capsules)**
- ▶ Dynamic Routing
- ▶ Digit Caps (higher level capsules of size 16x1)
- ▶ Compute Loss
- ▶ Backpropagate

# Real Estate Agents Analogy.

1152 real estate agents



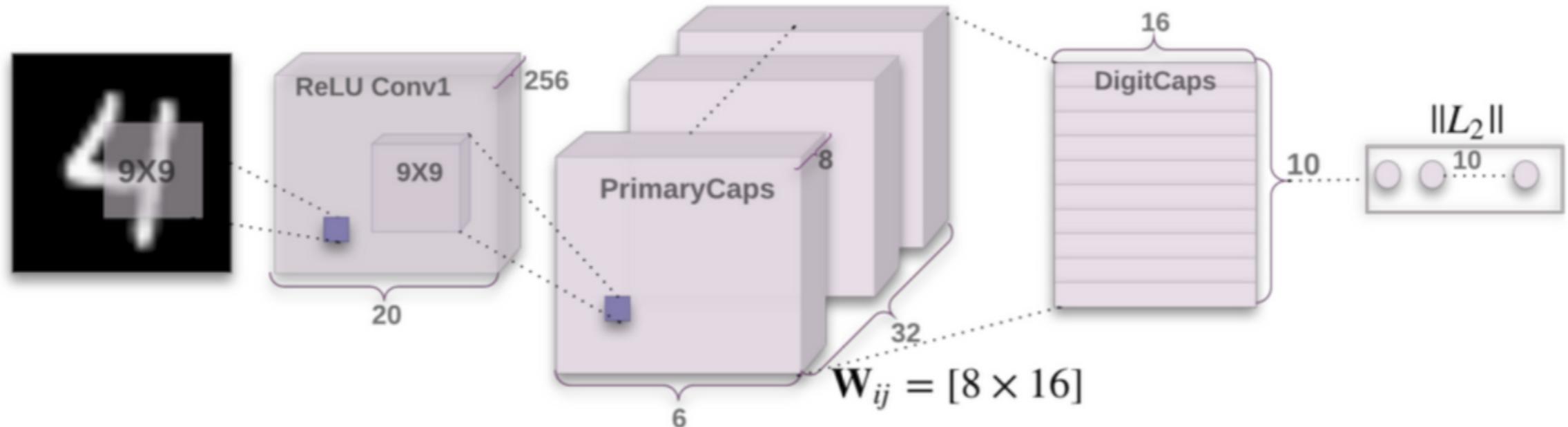
# Dynamic Routing.

$$\mathbf{s}_j = \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}, \quad \hat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij} \mathbf{u}_i \quad c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})}$$

**Procedure 1** Routing algorithm.

```

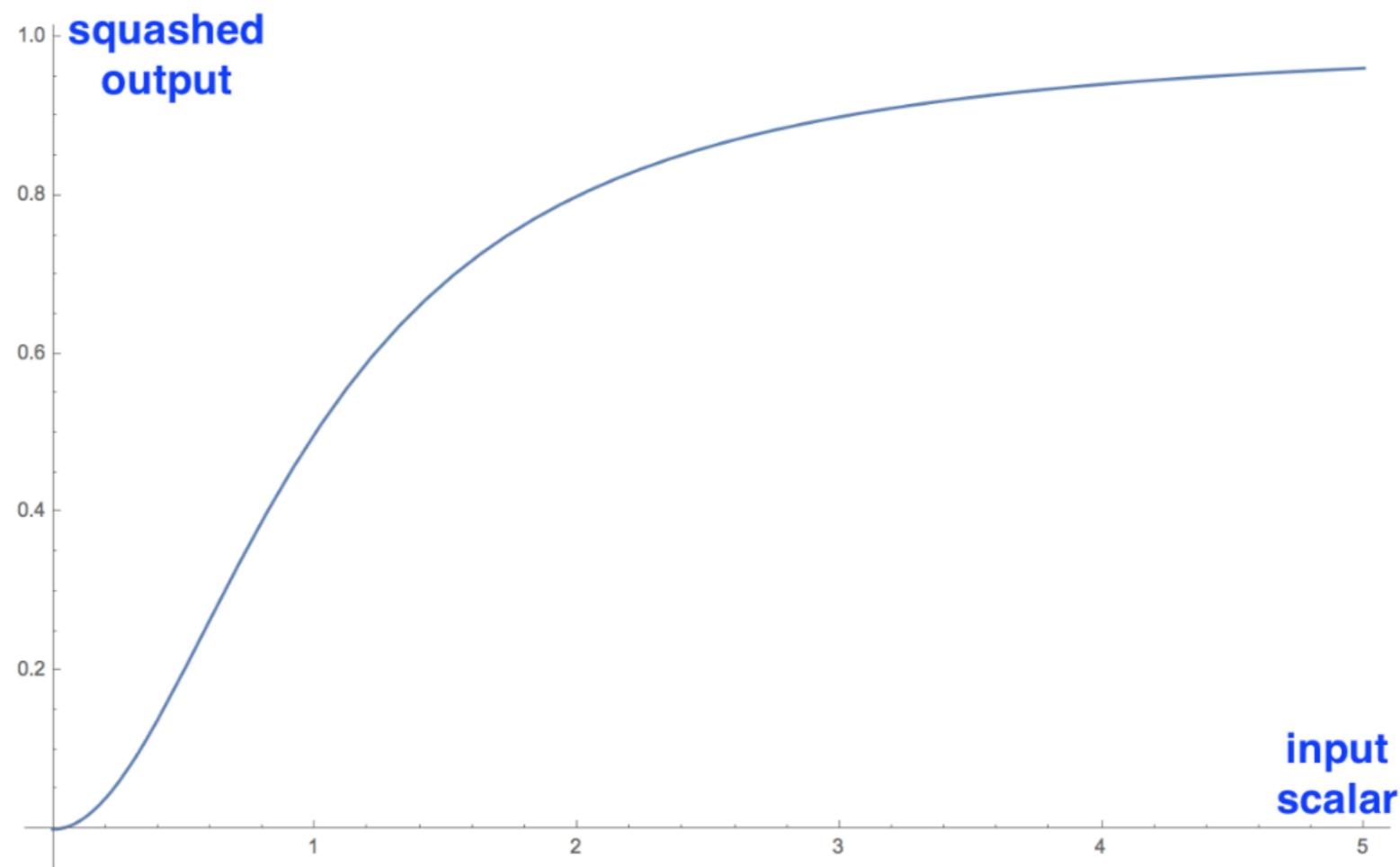
1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}$ ,  $r$ ,  $l$ )
2:   for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow 0$ .
3:   for  $r$  iterations do
4:     for all capsule  $i$  in layer  $l$ :  $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ 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



# Squashing Function.

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

additional "squashing"      unit scaling



# Dynamic Routing.

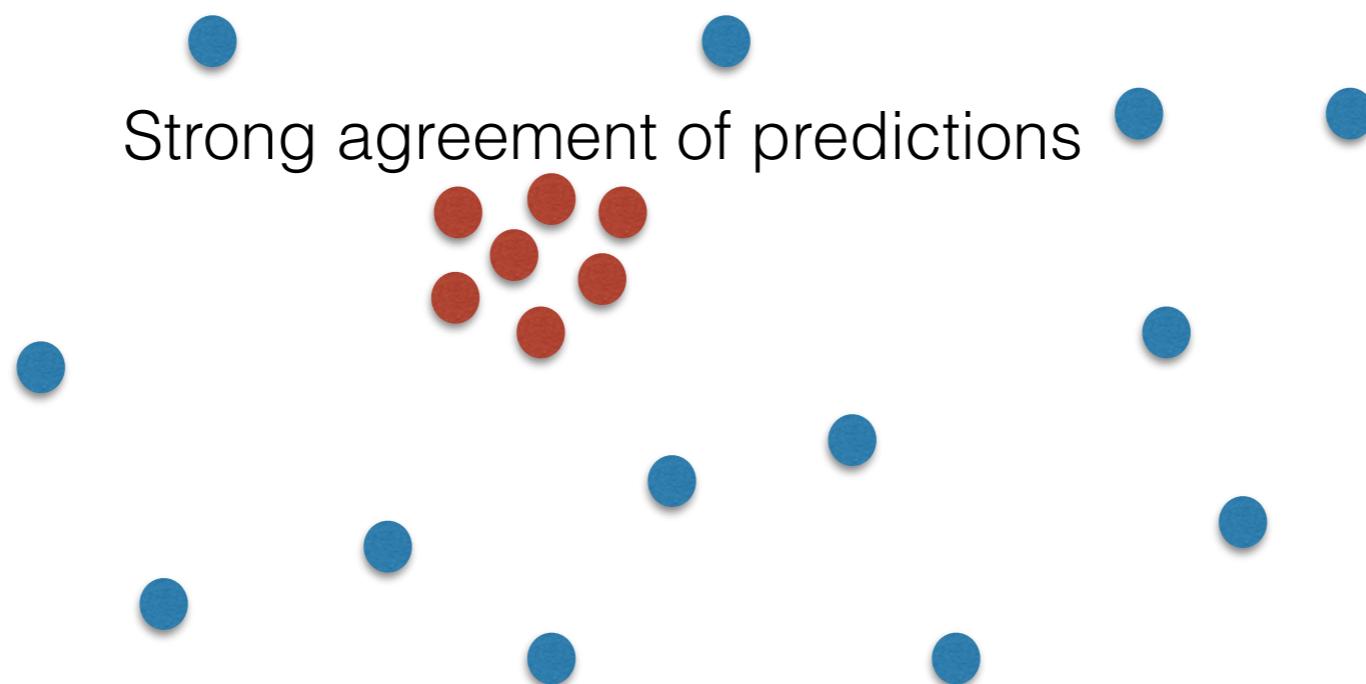
---

## Procedure 1 Routing algorithm.

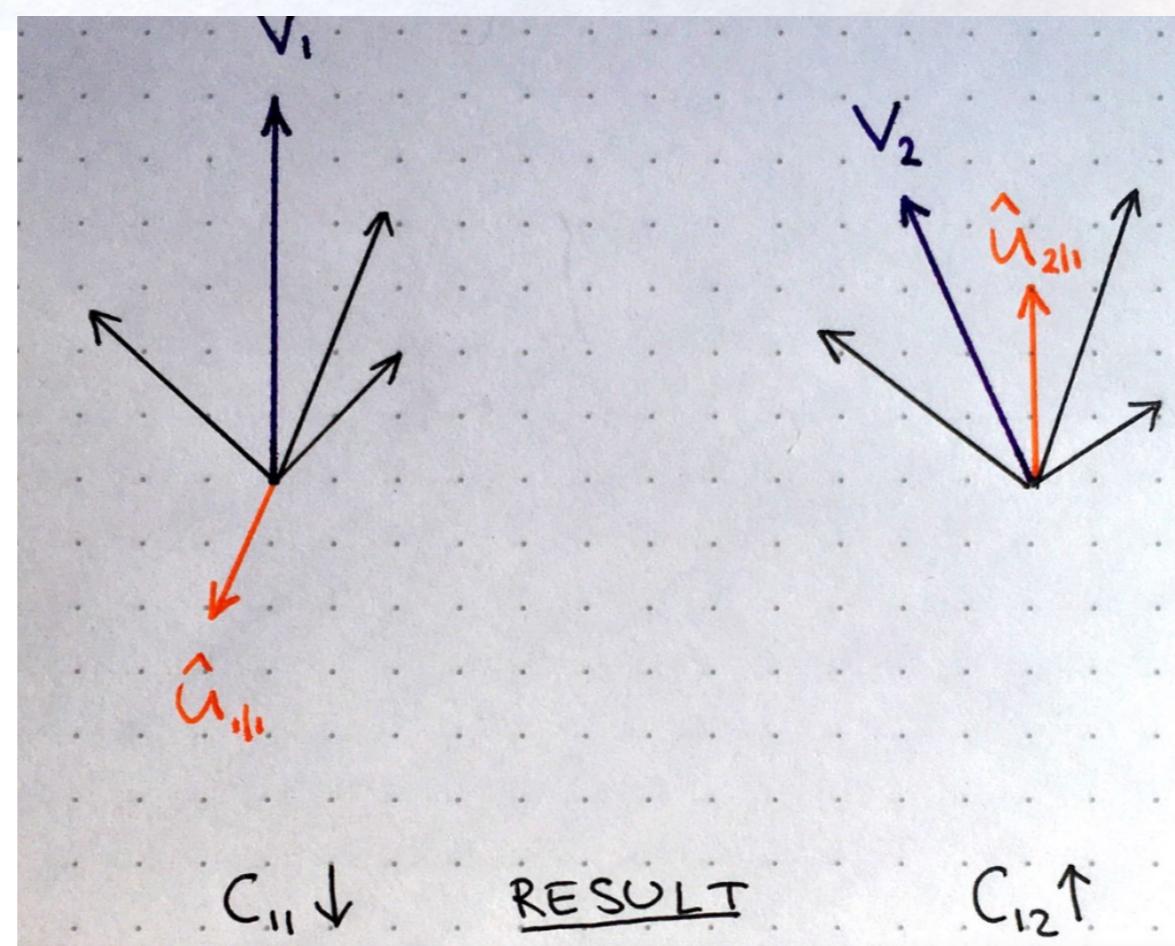
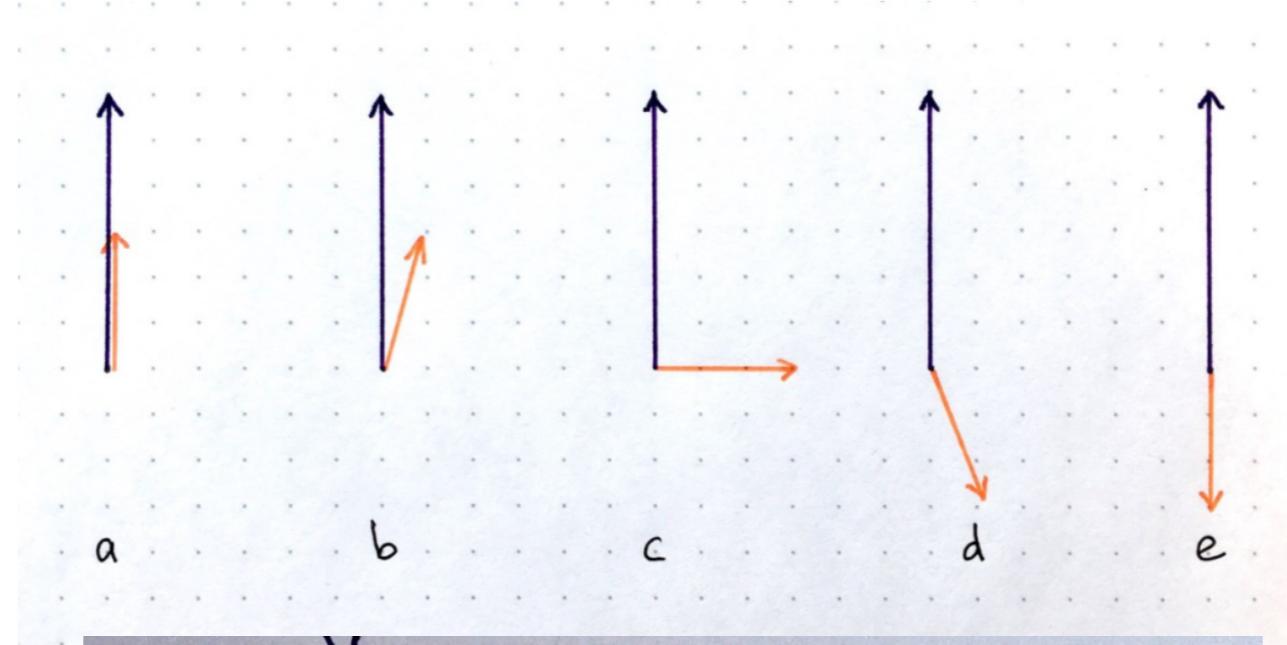
---

```
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  $l$ :  $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$                                  $\triangleright$  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  $(l + 1)$ :  $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$                                  $\triangleright$  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  $\mathbf{v}_j$ 
```

---

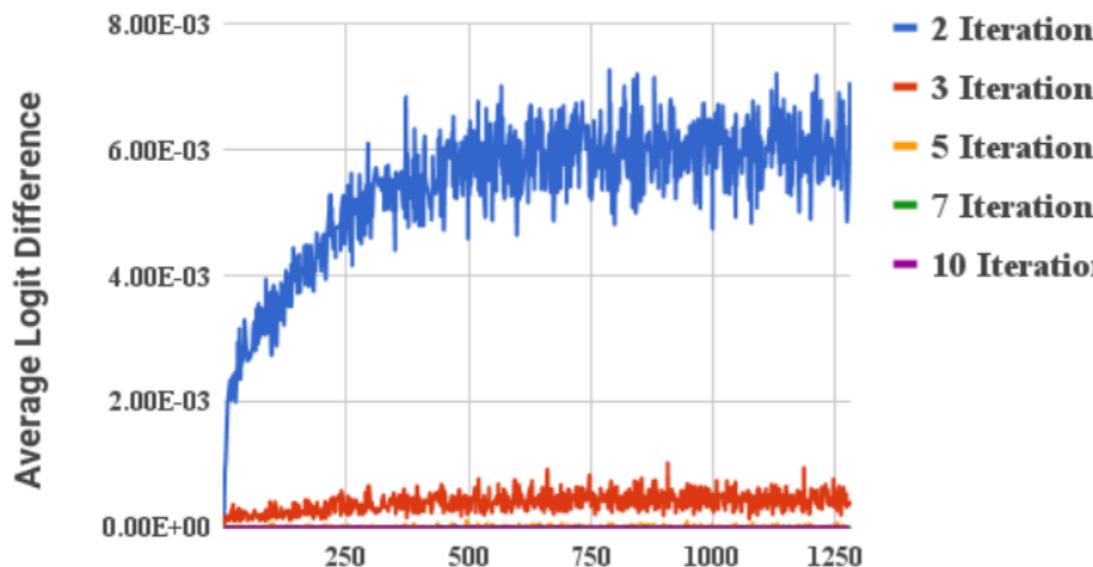


# Dynamic Routing.

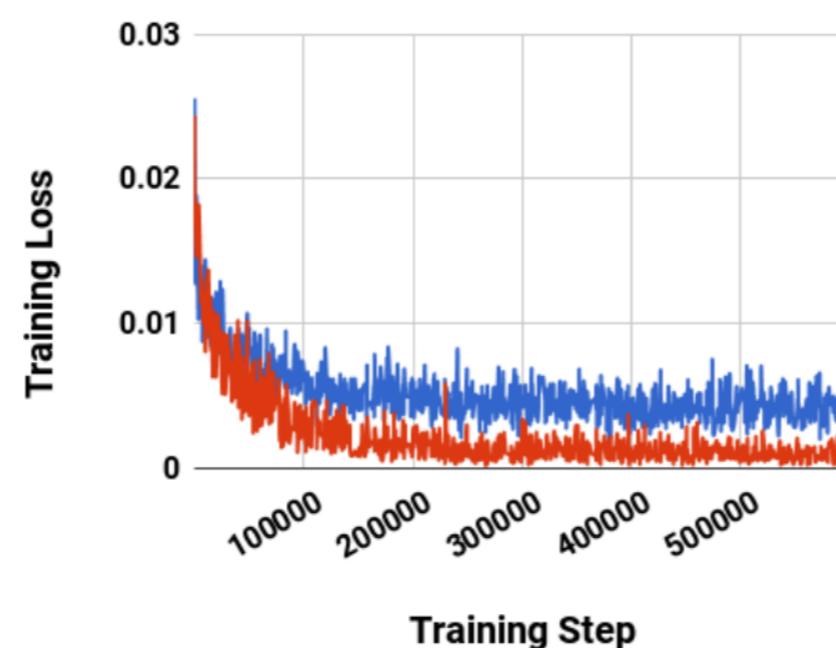
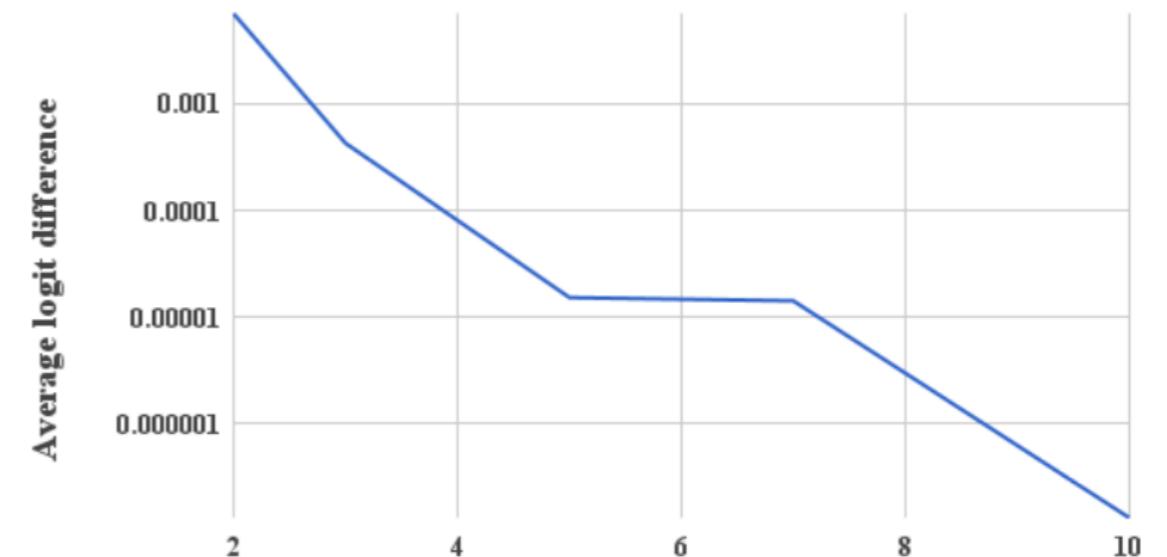


# Training: How many routing iterations to use?

(a) During training.

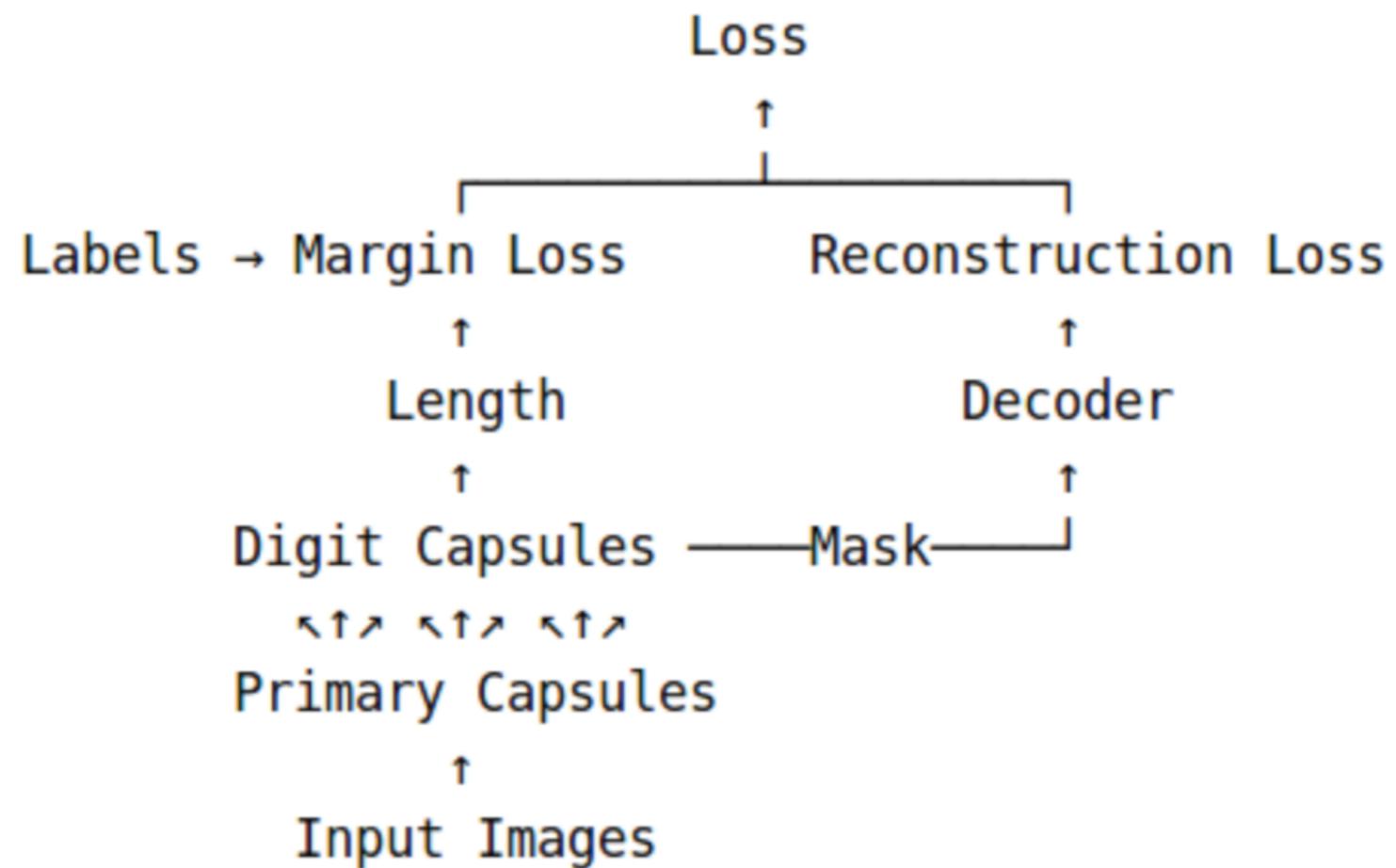


(b) Log scale of final differences.



we suggest 3 iteration of routing for all experiments.

# Training.



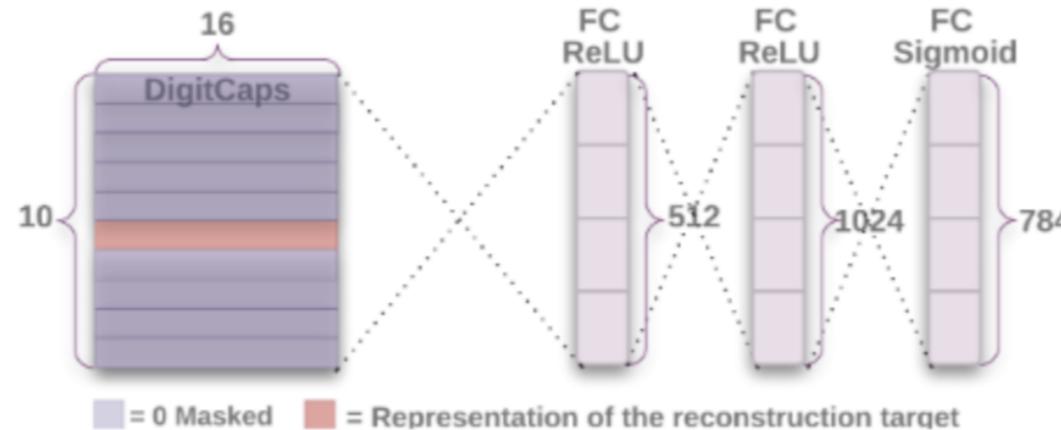
# Training.

Loss: Marginal + Reconstruction

Marginal loss:

$$L_k = T_k \max(0, m^+ - \|\mathbf{v}_k\|)^2 + \lambda (1 - T_k) \max(0, \|\mathbf{v}_k\| - m^-)^2$$

Reconstruction (MSE) loss:



total\_loss = marginal + 0.0005\*recon\_loss.

# What do individual dimensions of capsule represent?

Scale and thickness	0 0 0 0 0 0 0 0 0 0
Localized part	6 6 6 6 6 6 6 6 6 6
Stroke thickness	5 5 5 5 5 5 5 5 5 5
Localized skew	4 4 4 4 4 4 4 4 4 4
Width and translation	3 3 3 3 3 3 3 3 3 3
Localized part	2 2 2 2 2 2 2 2 2 2

Each row shows the reconstruction when one of the 16 dimensions in the DigitCaps representation is tweaked by intervals of 0.05 in the range [-0.25,0.25].

# Segmenting highly overlapping digits.

R:(2, 7) L:(2, 7)	R:(6, 0) L:(6, 0)	R:(6, 8) L:(6, 8)	R:(7, 1) L:(7, 1)	*R:(5, 7) L:(5, 0)	*R:(2, 3) L:(4, 3)	R:(2, 8) L:(2, 8)	R:P:(2, 7) L:(2, 8)
R:(8, 7) L:(8, 7)	R:(9, 4) L:(9, 4)	R:(9, 5) L:(9, 5)	R:(8, 4) L:(8, 4)	*R:(0, 8) L:(1, 8)	*R:(1, 6) L:(7, 6)	R:(4, 9) L:(4, 9)	R:P:(4, 0) L:(4, 9)

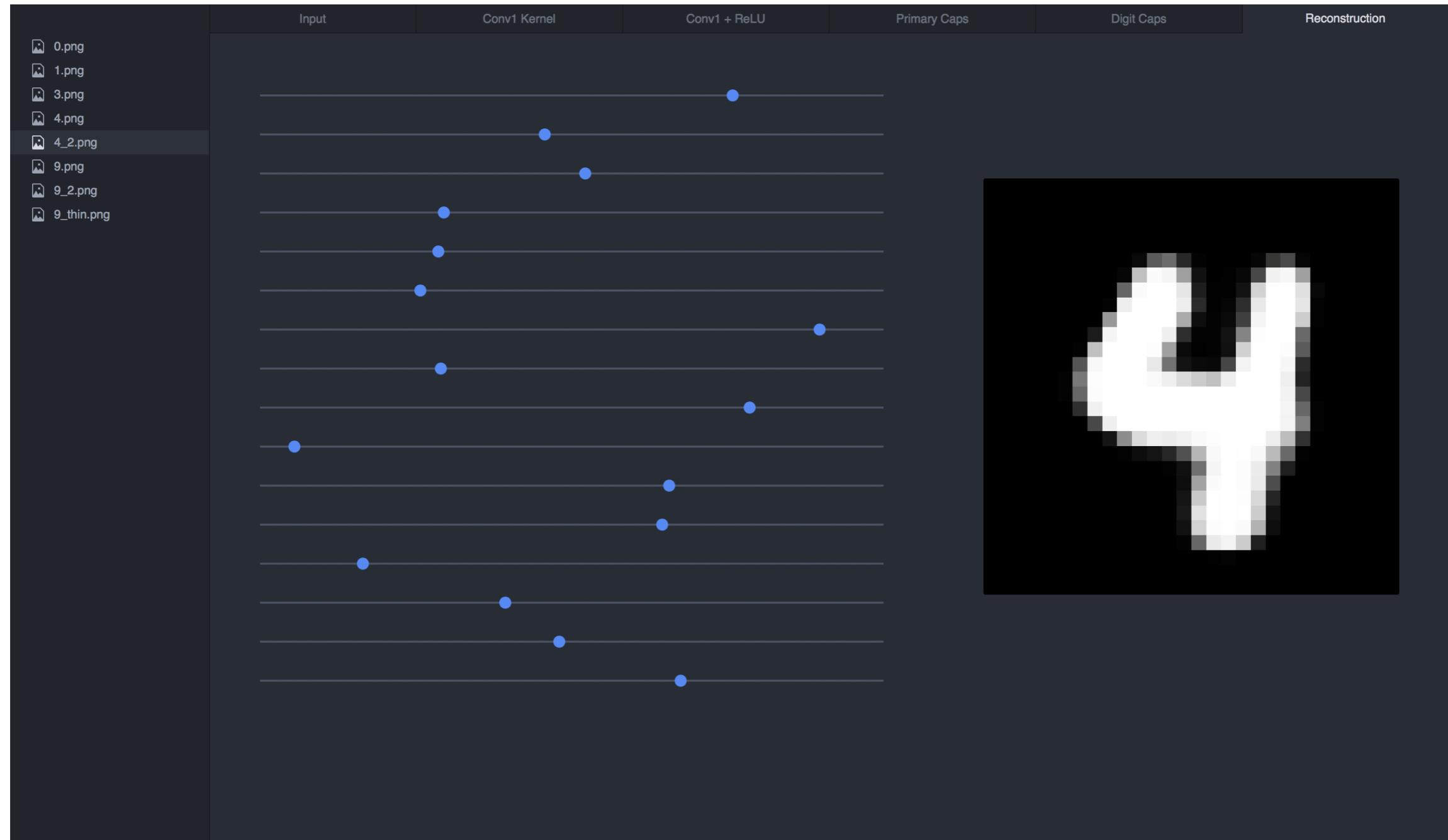
## PROS.

- Reaches high accuracy in MNIST, promising results for CIFAR10;
- Requires less training data
- Position and pose information is preserved
- Promising for image segmentation and detection
- Dynamic Routing works well for overlapping objects
- Capsule activations nicely map the hierarchy of parts

# Arguments against Capsule Networks.

- How about more complex data?
- Complexity
- No comparison with other architectures.

# Visualization.



<https://github.com/bourdakos1/CapsNet-Visualization>