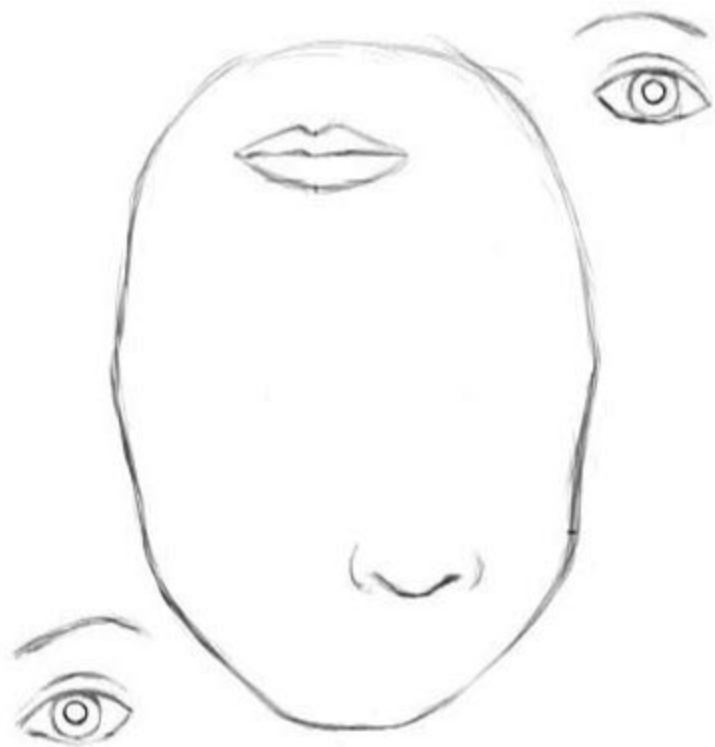
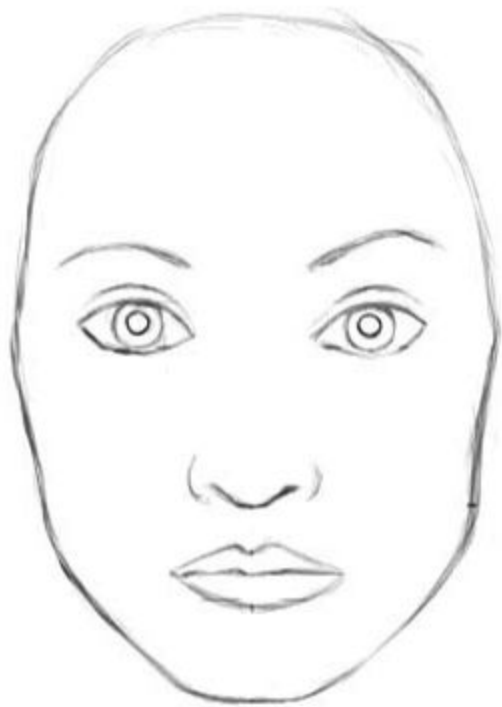


Dynamic Routing Between Capsules

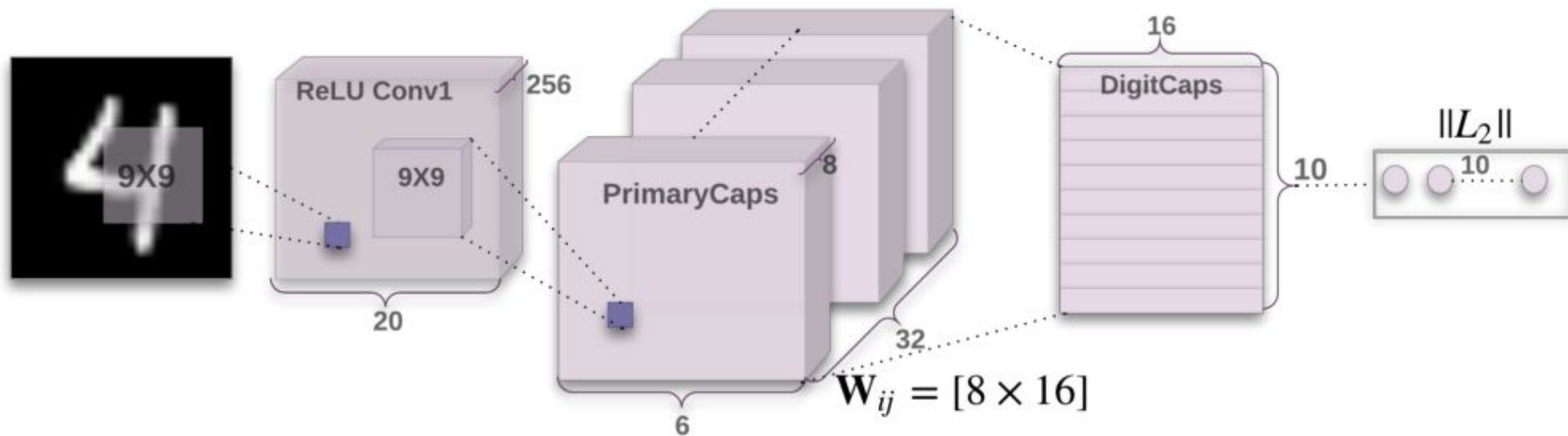
Main Idea

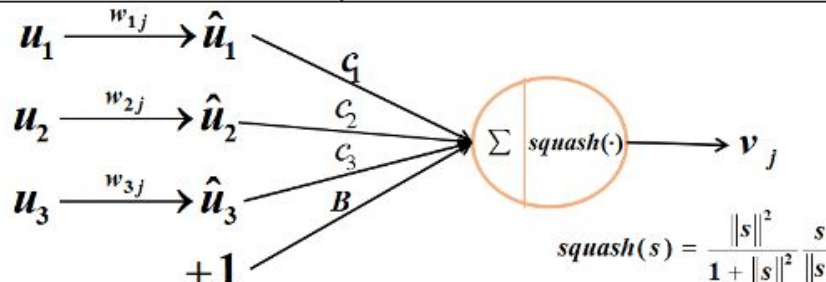
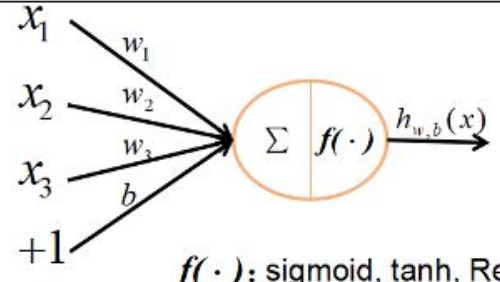
- Conv nets are good at translational invariance, but not great at other invariances (rotation, shadow, etc), room for improvement
- Capsules attempt to capture ‘instantiation parameters’ of features
 - Essentially doing ‘inverse graphics’
- Each feature is a vector rather than a scalar
 - Magnitude is ‘probability of the feature’
 - Direction is ‘parameters of the feature’
- Higher level features can depend on the parameters as well as the probability





Encoder Architecture



		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)
 <p style="text-align: center;">$\text{squash}(s) = \frac{\ s\ ^2}{1 + \ s\ ^2} \frac{s}{\ s\ }$</p>			 <p style="text-align: center;">$f(\cdot)$: sigmoid, tanh, ReLU, etc.</p>	

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

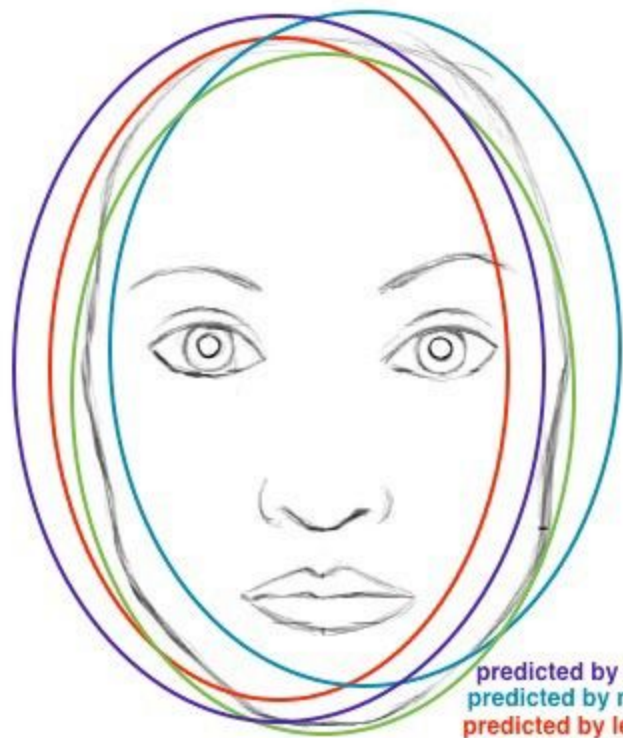
For each capsule i , the sum of c_{ij} is 1

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

additional "squashing" unit scaling

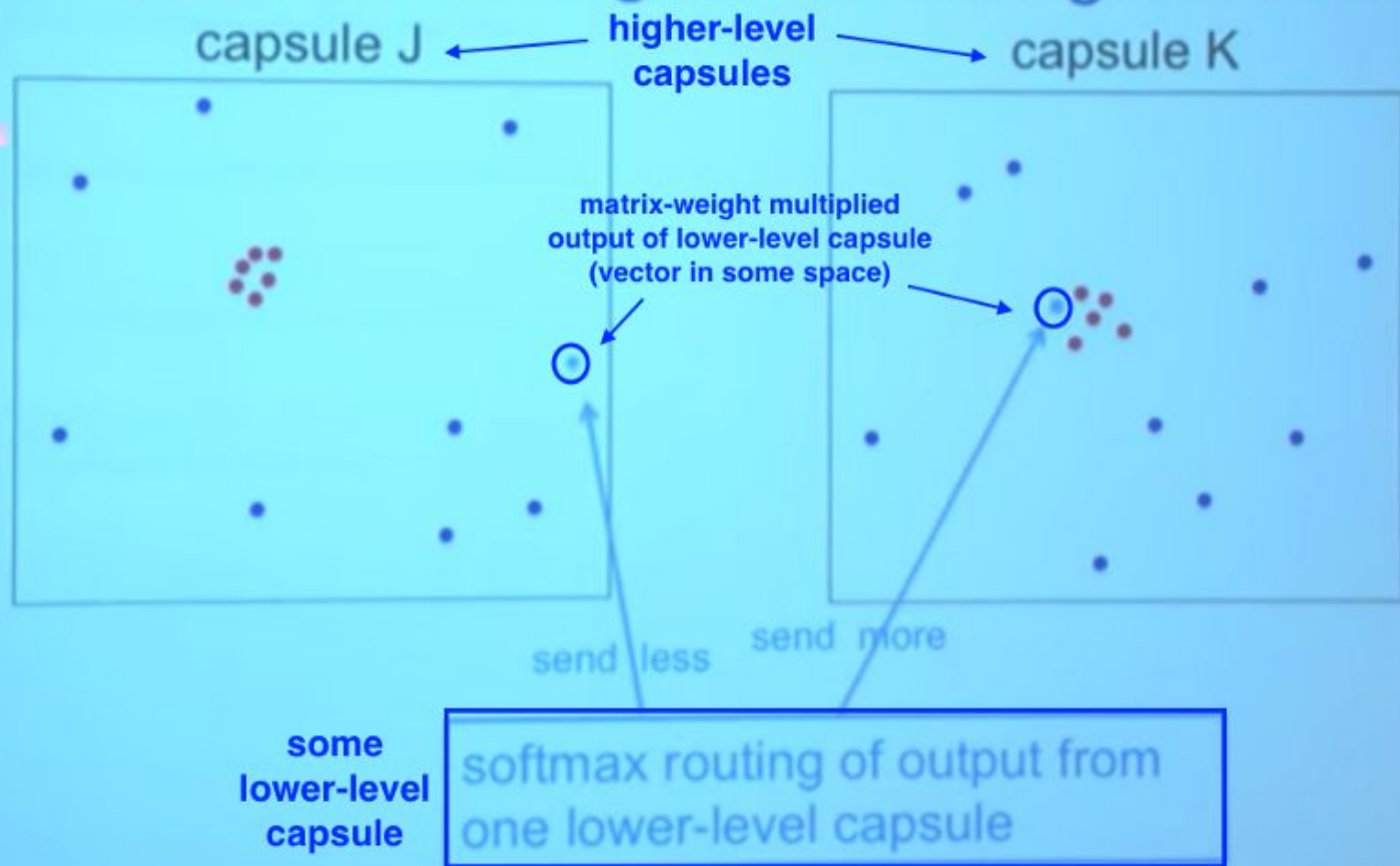
Procedure 1 Routing algorithm.

```
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  $\mathbf{v}_j$ 
```

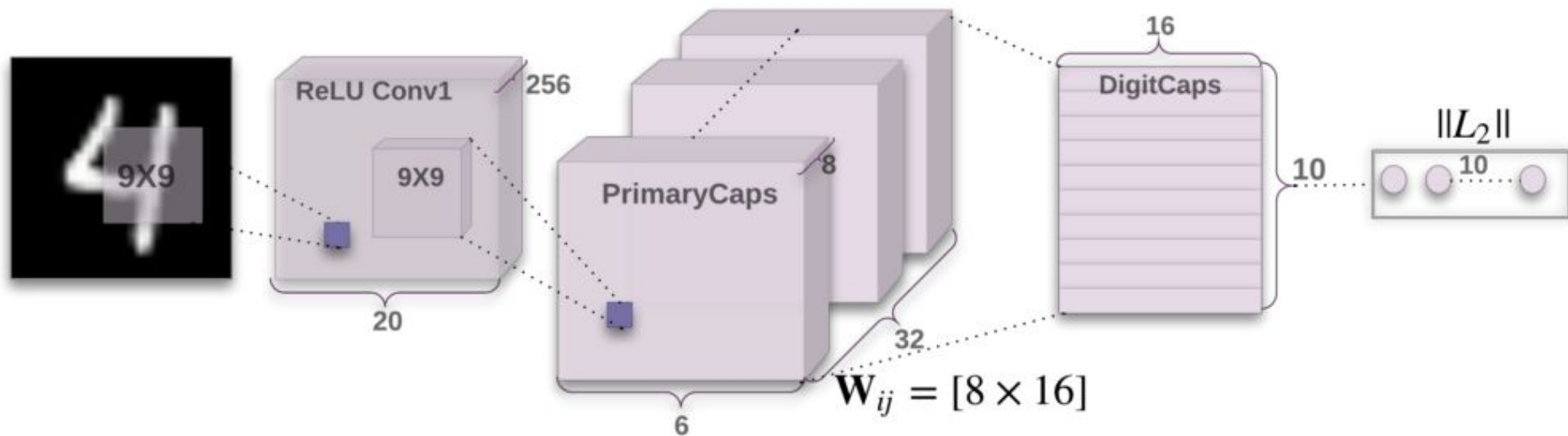


predicted by nose
predicted by mouth
predicted by left eye
predicted by right eye

Dynamic routing based on agreement



Encoder



CapsNet Loss Function

calculated for correct DigitCap

calculated for incorrect DigitCaps

loss term for one DigitCap

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

1 when correct
DigitCap,
0 when incorrect

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

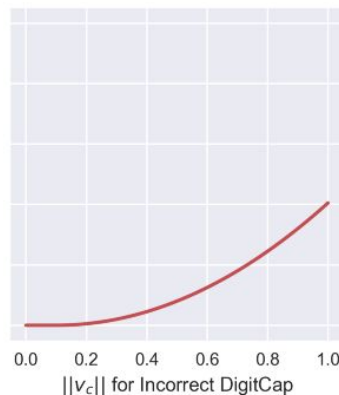
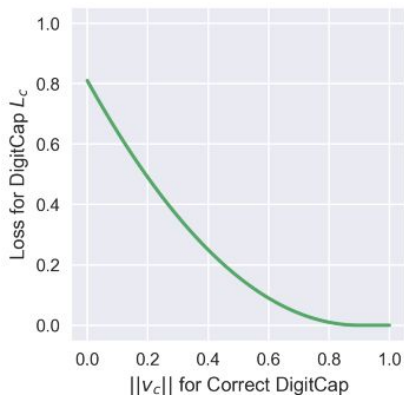
0.5 constant
used for
numerical
stability

1 when incorrect
DigitCap,
0 when correct

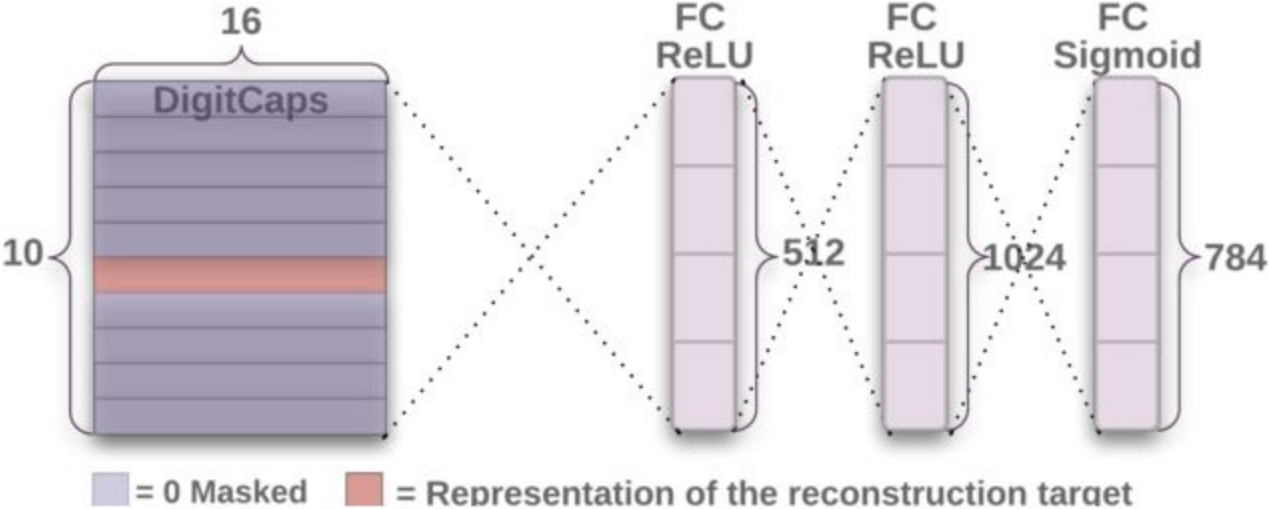
zero loss when incorrect
prediction with probability
less than 0.1, non-zero
otherwise




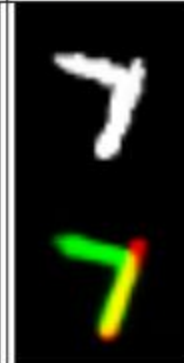









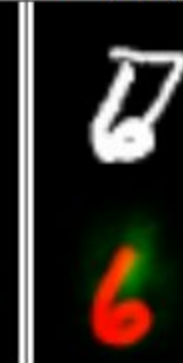
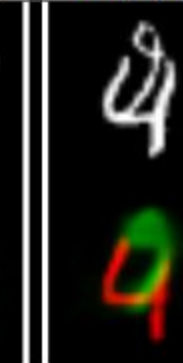
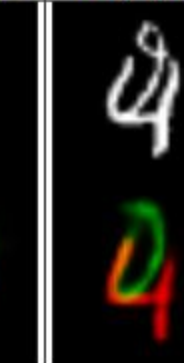
Note: correct DigitCap is one that matches training label, for each training example there will be 1 correct and 9 incorrect DigitCaps

Loss Function Value for Correct and Incorrect DigitCap



Decoder



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)
							

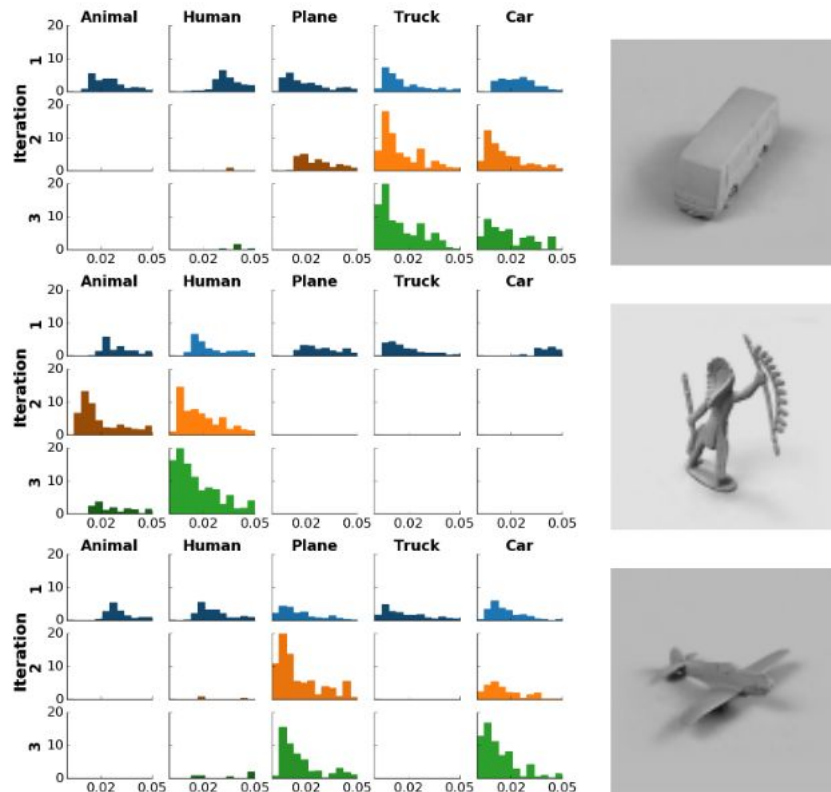


Figure 2: Histogram of distances of votes to the mean of each of the 5 final capsules after each routing iteration. Each distance point is weighted by its assignment probability. All three images are selected from the smallNORB test set. The routing procedure correctly routes the votes in the truck and the human example. The plane example shows a rare failure case of the model where the plane is confused with a car in the third routing iteration. The histograms are zoomed-in to visualize only votes with distances less than 0.05. Fig. B.2 shows the complete histograms for the "human" capsule without clipping the x-axis or fixing the scale of the y-axis.

Procedure 1 Routing algorithm¹ returns **activation** and **pose** of the capsules in layer $L + 1$ given the **activations** and **votes** of capsules in layer L . V_{ich} is an H dimensional vote from capsule i with activation a_i in layer L to capsule c in layer $L + 1$. β_a, β_v are learned discriminatively and the inverse temperature λ increases at each iteration with a fixed schedule.

1: **procedure** EM ROUTING(\mathbf{a}, V)

2: $\forall i, c: R_{ic} \leftarrow 1/\text{size}(L + 1)$

3: **for** t iterations **do**

4: $\forall c: M_{c:}, S_{c:}, \mathbf{a}'_c \leftarrow \text{M-STEP}(R_{:c}, \mathbf{a}, V_{:c:})$

5: $\forall i: R_{ij} \leftarrow \text{E-STEP}(M, S, \mathbf{a}', V_{i::})$

return \mathbf{a}', M

1: **procedure** M-STEP($\mathbf{r}, \mathbf{a}, V'$)

▷ for one higher-level capsule

2: $\forall i: \mathbf{r}'_i \leftarrow \mathbf{r}_i * \mathbf{a}_i$

3: $\forall h: \boldsymbol{\mu}_h \leftarrow \frac{\sum_i \mathbf{r}'_i V'_{ih}}{\sum_i \mathbf{r}'_i}$

4: $\forall h: \boldsymbol{\sigma}_h^2 \leftarrow \frac{\sum_i \mathbf{r}'_i (V'_{ih} - \boldsymbol{\mu}_h)^2}{\sum_i \mathbf{r}'_i}$

5: $\text{cost}_h \leftarrow (\beta_v + \log(\boldsymbol{\sigma}_h)) \sum_i \mathbf{r}'_i$

6: $\mathbf{a}' \leftarrow \text{sigmoid}(\lambda(\beta_a - \sum_h \text{cost}_h))$

7: **return** $\boldsymbol{\mu}, \boldsymbol{\sigma}, \mathbf{a}'$

1: **procedure** E-STEP(\mathbf{a}', S, M, V'')

▷ for one lower-level capsule

2: $\forall c: p_c \leftarrow \frac{1}{\sqrt{\prod_h^H 2\pi S_{ch}^2}} e^{-\sum_h \frac{(V''_{ch} - M_{ch})^2}{2S_{ch}^2}}$

3: $\forall c: \mathbf{r}_c \leftarrow \frac{\mathbf{a}'_c p_c}{\sum_j \mathbf{a}'_j p_j}$

4: **return** \mathbf{r}

Resources

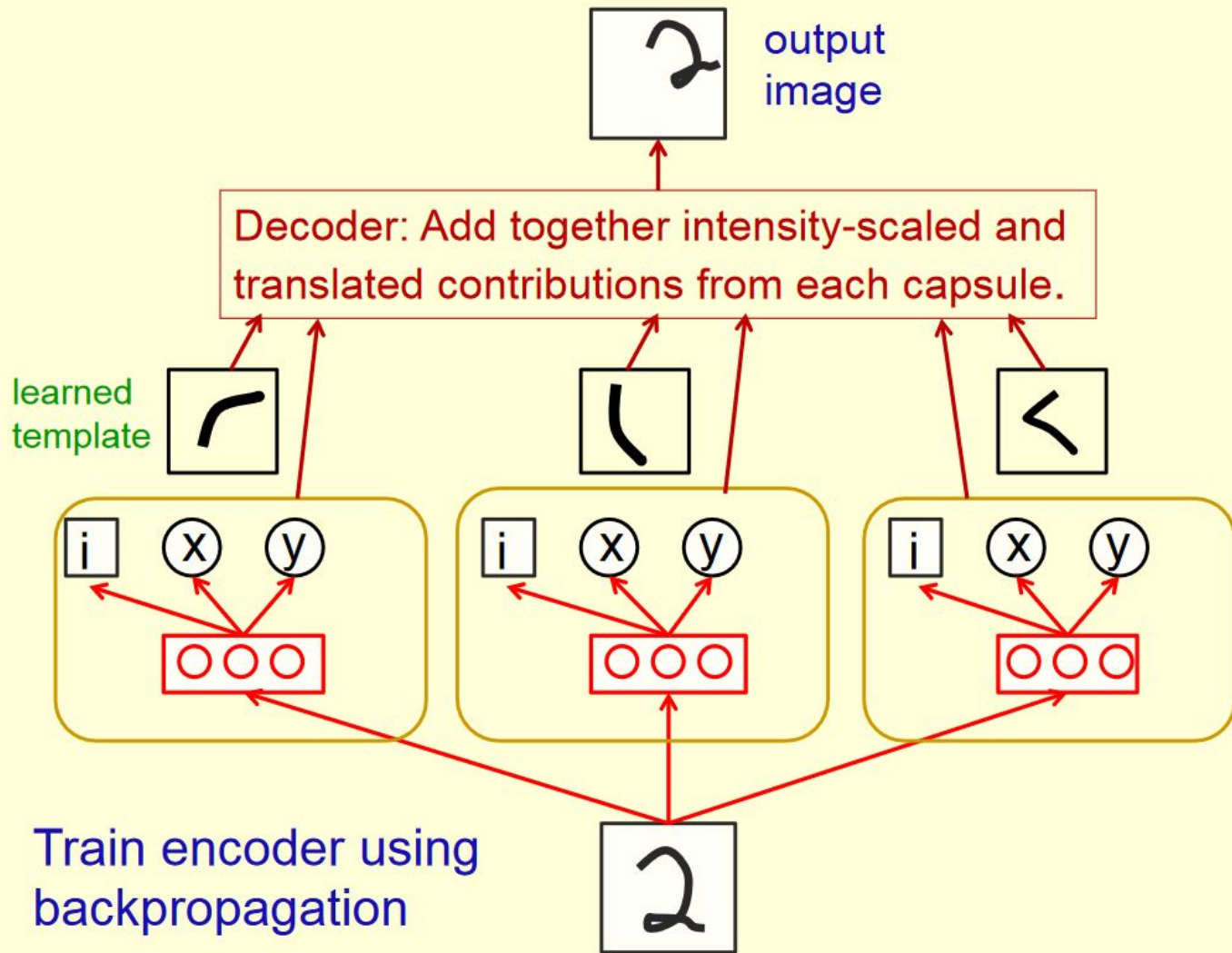
<https://arxiv.org/abs/1710.09829>

<https://pechyonkin.me/capsules-1/>

<https://github.com/naturomics/CapsNet-Tensorflow/>

<https://github.com/sekwiatkowski/awesome-capsule-networks>

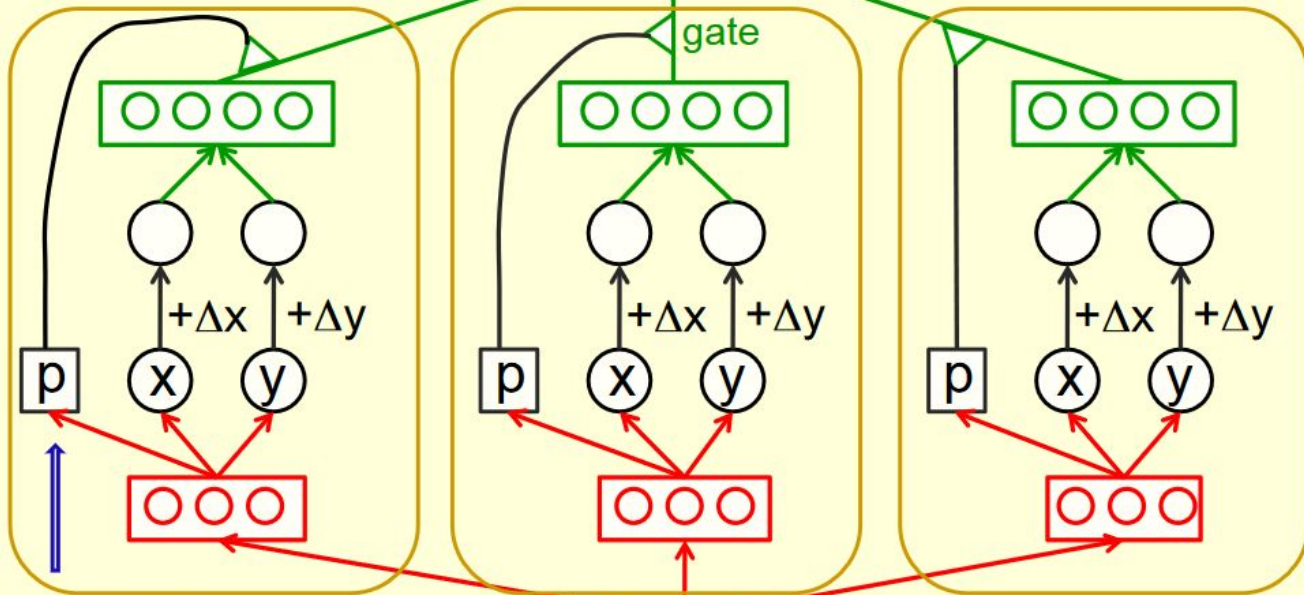
http://helper.ipam.ucla.edu/publications/qss2012/qss2012_10754.pdf



actual
output



target
output



probability that
the capsule's
visual entity is
present



input
image

