

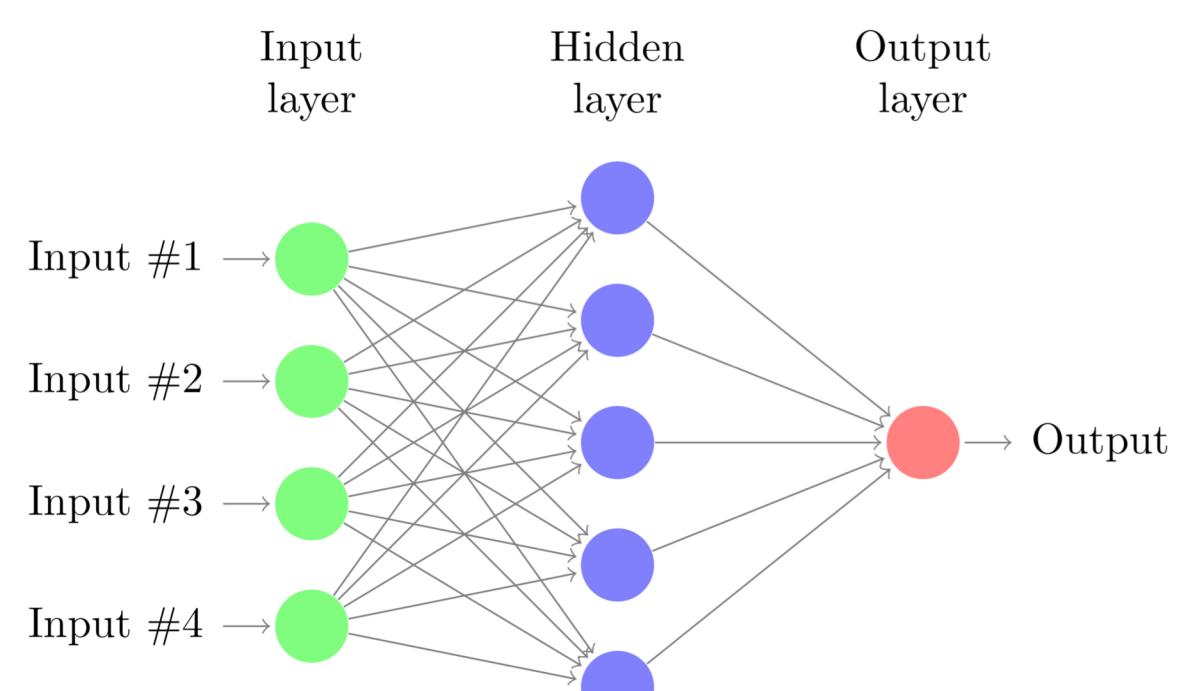
Residual Capsule Network

Can Zhong, Lingfeng Wang
Dept. of Computer Science, CEAS, University of Wisconsin - Milwaukee

Introduction

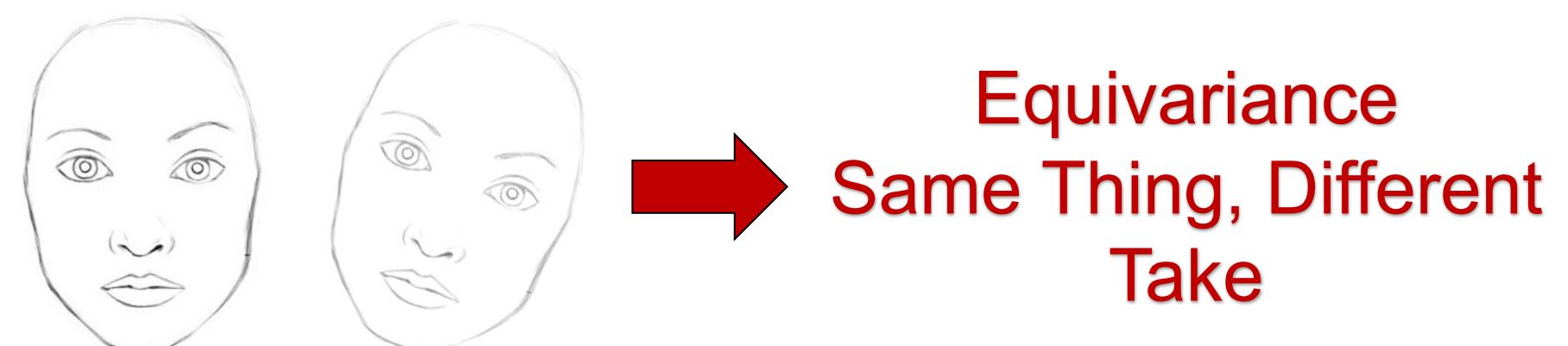
Artificial intelligence is the ability for an Agent to make practical and meaningful decisions in a dynamic or static environment. Machine learning is a means to accomplish artificial intelligence decision making through the implementations of algorithms which can find numerical relationships between data points without being explicitly programmed.

Due to the accumulation of massive amounts of data and the emergence of efficient computational hardware, deep learning has been a popular framework which has gain traction in research and applications in tasks such as classification, decision-making, and prediction.



CapNetv1 – Dynamic Routing

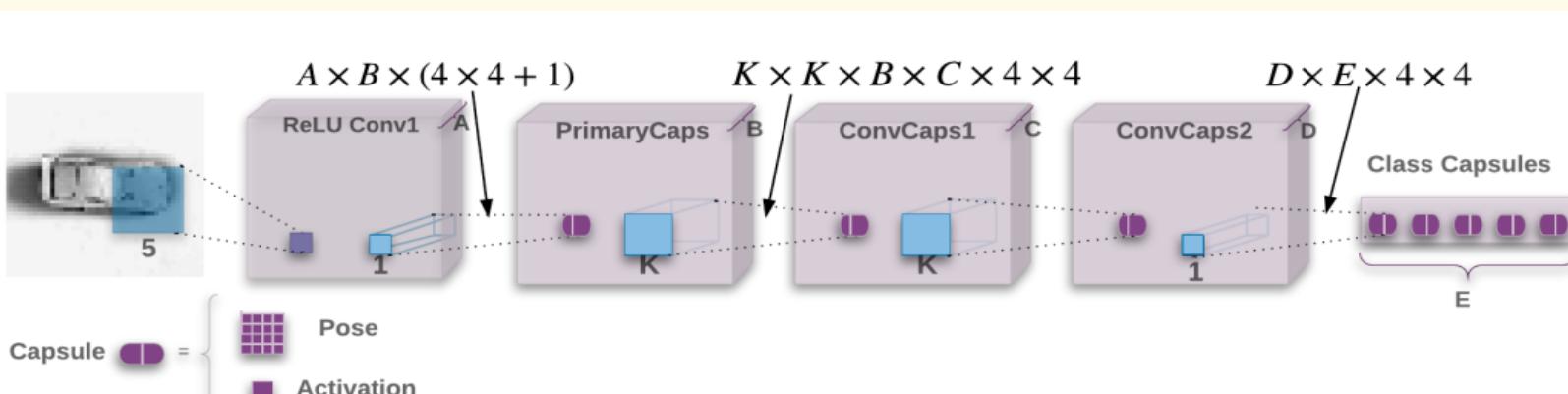
The most prominent architecture in use are convolutional neural networks which learn shared “concepts” throughout the data point by recognizing trends in a smaller space. CNNs excel at understanding the hierarchical properties of information yet one of its challenges is keeping perspective independent from the learned features.



Capsule networks are an expansion of CNN which consider perspective or pose independent of activations. This architecture can be done with less parameters because the network does not need to learn different feature map for varying i.e. angles.

capsule		VS.	traditional neuron
Input from low-level neurons/capsules	vector(u_i)		scalar(x_i)
Operations	Linear/Affine Transformation	$\hat{u}_j = W_j u_i + B_j$ (Eq. 2)	$a_j = w_j \cdot x_i + b_j$
	Weighting	$s_j = \sum_i c_{ij} \hat{u}_j$ (Eq. 2)	$z_j = \sum_i \mathbf{1} \cdot a_{ji}$
	Summation		
	Non-linearity activation	$v_j = \text{squash}(s_j)$ (Eq. 1)	$h_{w,d}(x) = f(z_j)$
output		vector(v_j)	scalar(h)
$\begin{aligned} u_1 &\xrightarrow{w_{1j}} \hat{u}_1 \\ u_2 &\xrightarrow{w_{2j}} \hat{u}_2 \\ u_3 &\xrightarrow{w_{3j}} \hat{u}_3 \\ &+1 \end{aligned}$ $\begin{aligned} \hat{u}_1 &\xrightarrow{c_1} c_1 \\ \hat{u}_2 &\xrightarrow{c_2} c_2 \\ \hat{u}_3 &\xrightarrow{c_3} c_3 \\ &\sum \text{squash}(s_j) = \frac{\ \mathbf{s}_j\ ^2}{1 + \ \mathbf{s}_j\ ^2} \mathbf{s}_j \end{aligned}$ $\begin{aligned} X_1 &\xrightarrow{w_1} \\ X_2 &\xrightarrow{w_2} \\ X_3 &\xrightarrow{w_3} \\ &+1 \end{aligned}$ $\begin{aligned} &\sum f(\cdot) \cdot h_{w,d}(x) \end{aligned}$			
<p>Capsule = New Version Neuron! vector in, vector out VS. scalar in, scalar out</p> $v_j = \frac{\ \mathbf{s}_j\ ^2}{1 + \ \mathbf{s}_j\ ^2} \mathbf{s}_j$			

CapNetv2 – EM Routing



Procedure 1 Routing algorithm returns activation and pose of the capsules in layer L . V_{ij}^h is the h^{th} dimension of the vote from capsule i with activation a_i in layer L to capsule j in layer $L+1$. β_a, β_s are learned discriminatively and the inverse temperature λ increases at each iteration with a fixed schedule.

```

1: procedure EM ROUTING( $a, V$ )
2:    $\forall i \in \Omega_L, j \in \Omega_{L+1}$ :  $R_{ij} \leftarrow 1/|\Omega_{L+1}|$ 
3:   for  $t$  iterations do
4:      $\forall j \in \Omega_{L+1}$ : M-STEP( $a, R, V, j$ )
5:    $\forall i \in \Omega_L$ : E-STEP( $\mu, \sigma, a, V, i$ )
   return  $a, \bar{M}$ 
  
```

```

1: procedure EM ROUTING( $a, V$ )
2:    $\forall i \in \Omega_L, j \in \Omega_{L+1}$ :  $R_{ij} \leftarrow R_{ij} * a_i$ 
3:    $\forall h$ :  $\mu_h^j \leftarrow \frac{\sum_i R_{ij} v_{ij}^h}{\sum_i R_{ij}}$ 
4:    $\forall h$ :  $(\sigma_h^j)^2 \leftarrow \frac{\sum_i R_{ij} (V_{ij}^h - \mu_h^j)^2}{\sum_i R_{ij}}$ 
5:    $cost^h \leftarrow (\beta_a + log(\sigma_h^j)) \sum_i R_{ij}$ 
6:    $a_j \leftarrow logistic(\lambda(\beta_a - \sum_h cost^h))$ 
  
```

```

1: procedure E-STEP( $\mu, \sigma, a, V, i$ )
2:    $\forall j \in \Omega_{L+1}$ :  $R_{ij} \leftarrow \frac{1}{\sqrt{\prod_{k \in \Omega_{L+1}} a_k p_k}} \exp \left( -\sum_h \frac{(V_{ij}^h - \mu_h^j)^2}{2(\sigma_h^j)^2} \right)$ 
3:    $\forall j \in \Omega_{L+1}$ :  $R_{ij} \leftarrow \frac{R_{ij}}{\sum_{k \in \Omega_{L+1}} a_k p_k}$ 
  
```

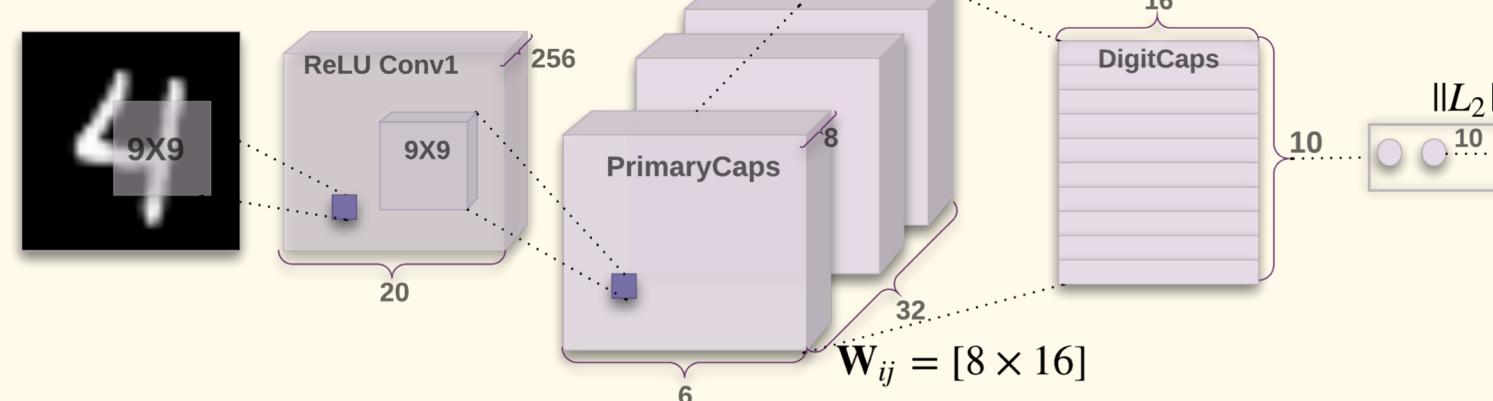
$$P_{ij}^h = \frac{1}{\sqrt{2\pi(\sigma_h^j)^2}} \exp \left(-\frac{(V_{ij}^h - \mu_h^j)^2}{2(\sigma_h^j)^2} \right), \quad ln(P_{ij}^h) = -\frac{(V_{ij}^h - \mu_h^j)^2}{2(\sigma_h^j)^2} - ln(\sigma_h^j) - ln(2\pi)/2$$

$$a_j = logistic \left(\lambda \left(\beta_a - \beta_u \sum_i r_{ij} - \sum_h cost_j^h \right) \right)$$

Methods

- To fully appreciate the intuitions behind neural networks, it was necessary for me to understand the fundamental concepts first. This includes the history, algorithms, architectures, statistics and calculus.
- To explore the recent development in various fields of neural networks, research papers were read and annotated.
- From these insights, 4 models were modified and then trained with Nvidia Tesla K80 and V100 GPUs on the MNIST dataset with good results. MNIST test results were good, a best of 98.82 %.
- Further experiments were trained and tested on a more complex dataset, CIFAR10 with the same models achieving a staggering low of 66.92 %. A conceptually-comparable network (DCNET) achieved 82.68 % with 16 million parameters. Further improved my model and re-ran. Reached a respectful 76.12 % while using only 4.5 million parameters.
- Results from CIFAR10 did not meet benchmarks but were promising. To beat the state of the art, I attempted to design an architecture which uses ODE-based functions for training and testing.

CapNetv1 - Dynamic Routing Architecture + Algorithm

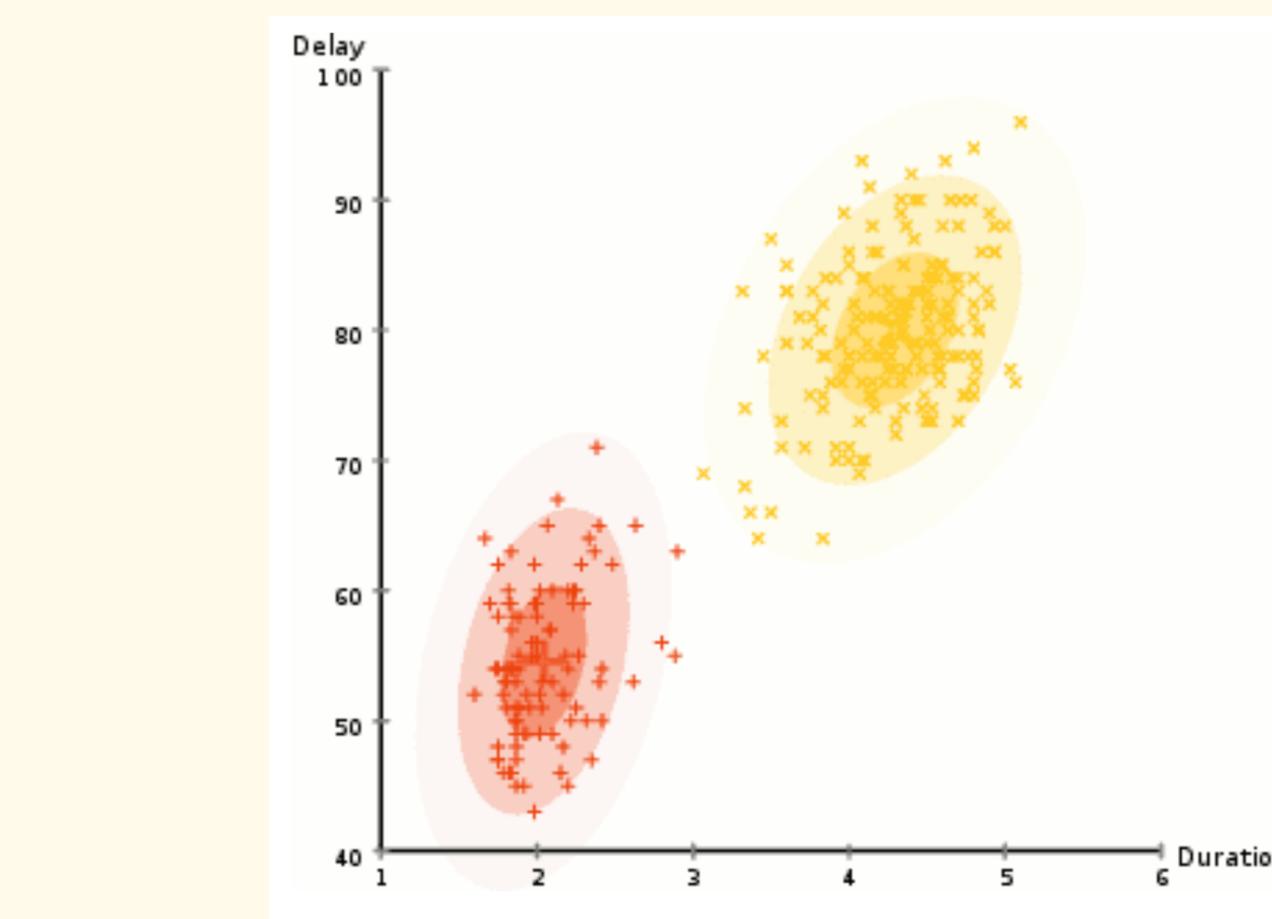


Procedure 1 Routing algorithm.

```

1: procedure ROUTING( $\hat{u}_{ij}, r, l$ )
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  $l$ :  $c_i \leftarrow softmax(b_{ij})$   $\triangleright$  softmax computes Eq. 4.
5:     for all capsule  $j$  in layer  $(l+1)$ :  $s_j \leftarrow \sum_i c_{ij} \hat{u}_{ij}$ 
6:     for all capsule  $i$  in layer  $l$ :  $v_j \leftarrow squash(s_j)$   $\triangleright$  squash computes Eq. 6.
7:     for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l+1)$ :  $b_{ij} \leftarrow b_{ij} + \hat{u}_{ij} \cdot v_j$ 
   return  $v_j$ 
  
```

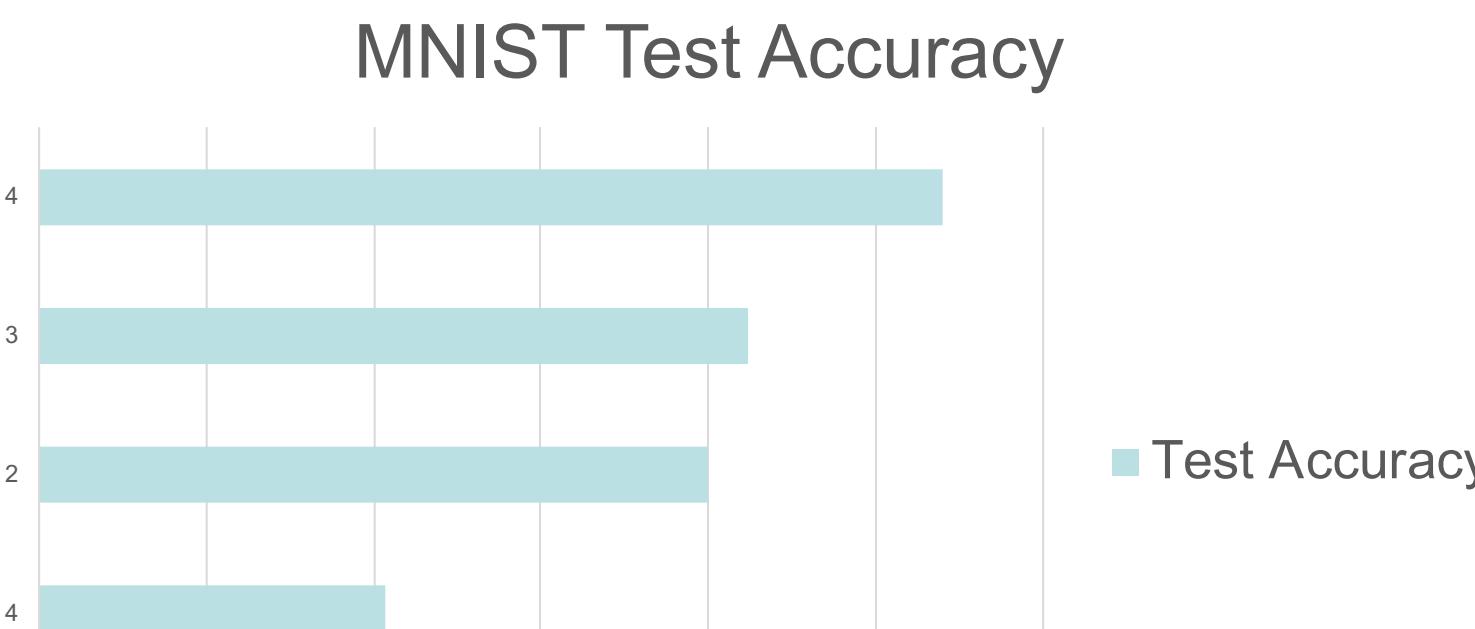
CapNetv2 – Expectation Maximization Clustering



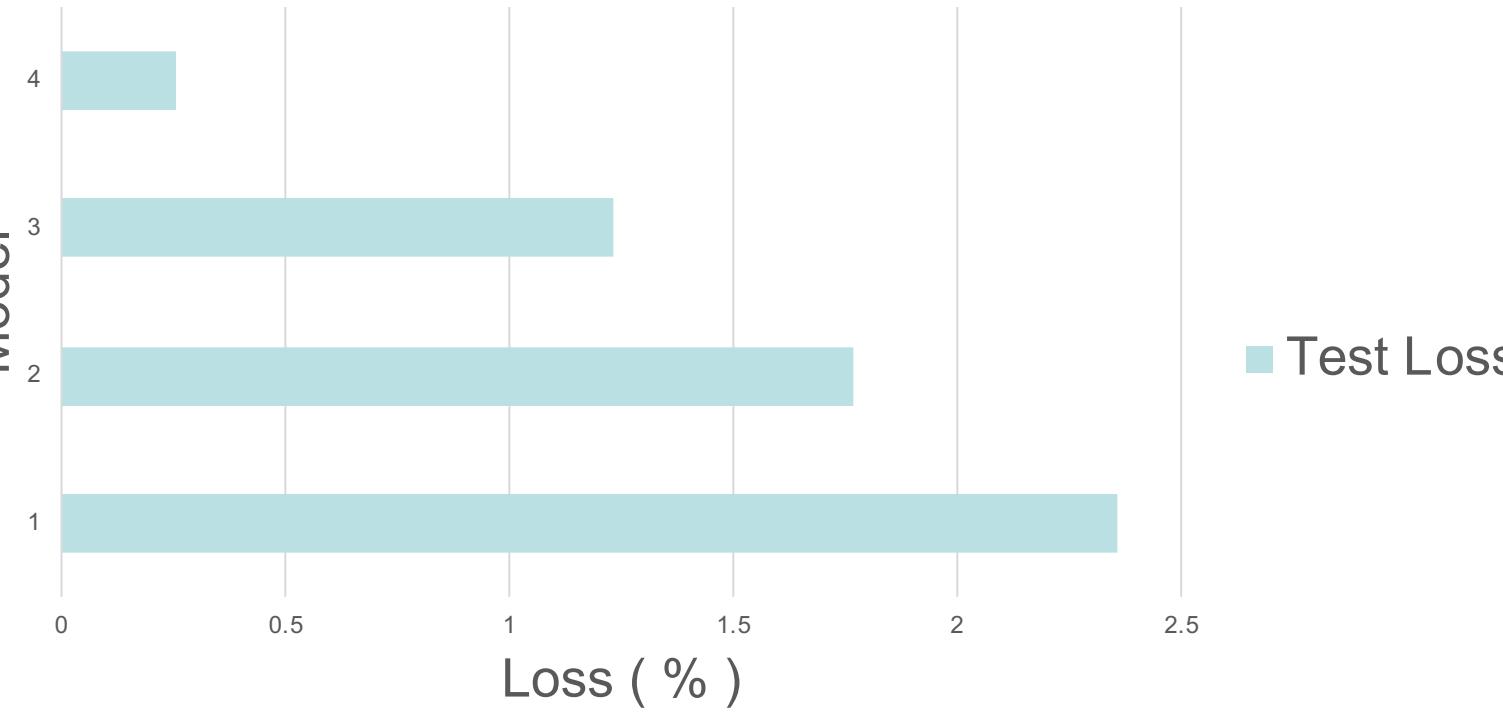
Experiment

- In my experiment, I wanted to explore the hardware and computational ceilings that capsules face. Capsules are a relatively new concept and therefore development is in its infancy.
- I first wanted to test if the capsules inputs had any effects on the accuracy. I started through this route because I find the ability of a capsule to recognize patterns without having access to the whole very enticing.
- I hypothesized that if we changed what the input features were so that the input features described the datapoint more generic, then the classification accuracy would increase.
- I modified the networks in a way where the input to the primary capsules came from various networks. Baseline CNN, Residual, Residual with Cardinality, and an ODE-network.
- I found it quite slow to train my models as compared to other established architectures like CNN, and Dense Networks.
- I was limited to 1 GPU, but trained on the best Tesla V100.
- Residual Networks + Cardinality was an implementation which used more parameters than was necessary.
- ODE-network was the best option, but the problem with capsules is not how real data is encoded into the latent space, but rather how capsules normalize the flow from the latent space to the capsules.
- Feature endeavors include normalizing flow using capsules, and modifying the structure so that the state can use the last residual state. A problem with capsules is features learned earlier in the layer are more often than not to be forgotten in later layers.

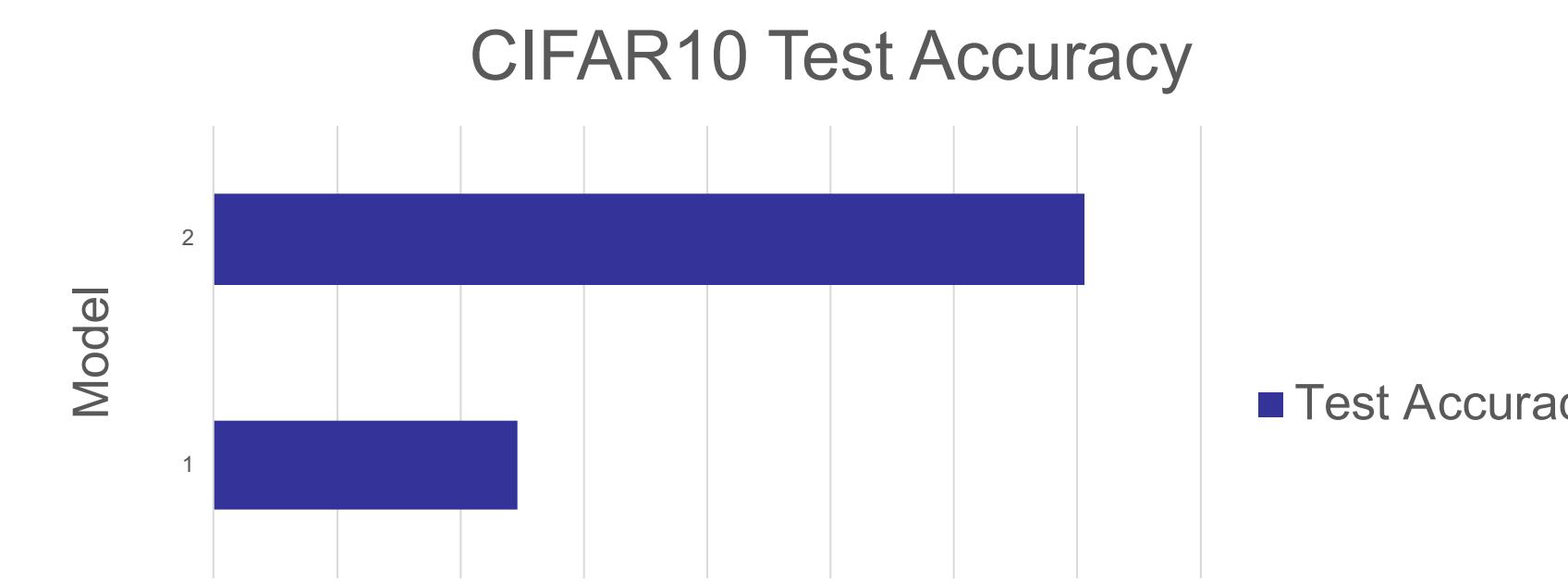
MNIST Results



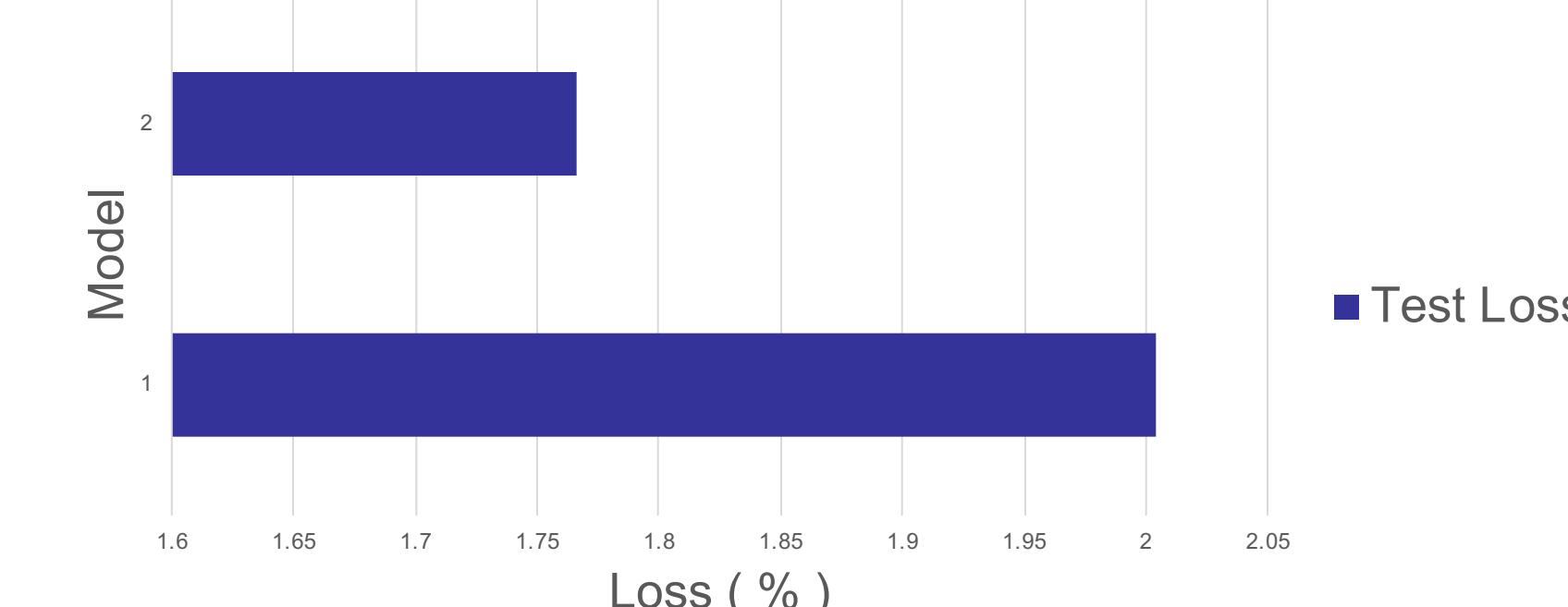
MNIST Test Loss



CIFAR10 Results



CIFAR10 Test Loss



Citations

- Hinton, Geoffrey, Sabour, Sara, Frosst, Nicholas (2018). Matrix Capsules with EM Routing. *ICLR 2018*. 1-15.
- Schmidhuber, Jürgen (2014). Deep Learning in Neural Networks: An Overview. *University of Lugano & SUPSI*. 1-88.
- Xie, Saining, Girshick, Ross, Dollar, Piotr, Tu, Zhuowen, He, Kaiming (2017). Aggregated Residual Transformations for Deep Neural Networks. *UC San Diego, Facebook AI Research*. 1-10.
- Chen, Ricky, Rubanova, Yulia, Bettencourt, Jesse, Duvenaud, David (2018). Neural Ordinary Differential Equations. *University of Toronto, Vector Institute*. 1-19.
- Sabour, Sara, Frost, Nicholas, Hinton, Geoffrey (2017). Dynamic Routing Between Capsules. *Cornell University, Computer Vision and Pattern Recognition*. 1-11.