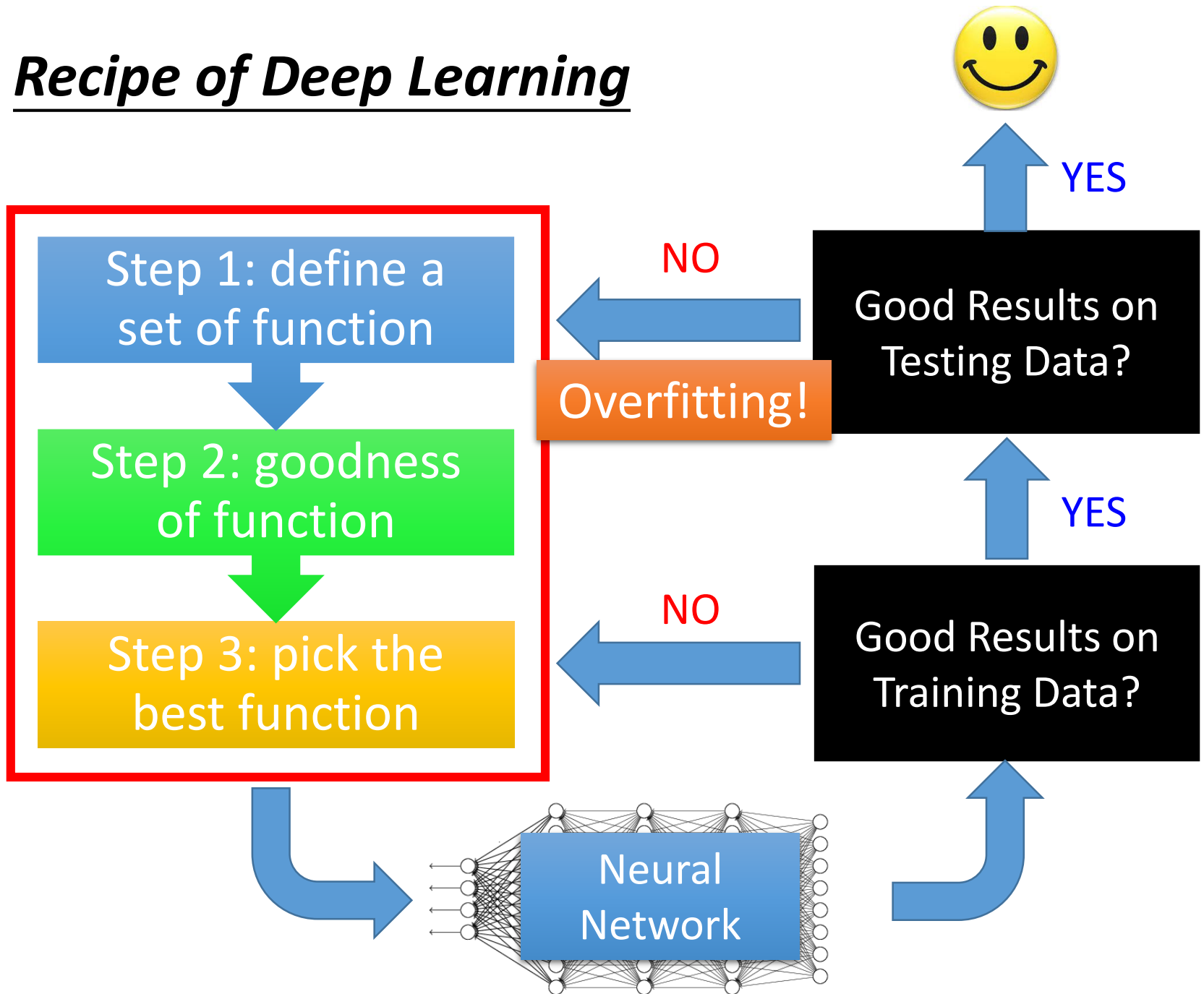


Advanced Tips for Deep Learning

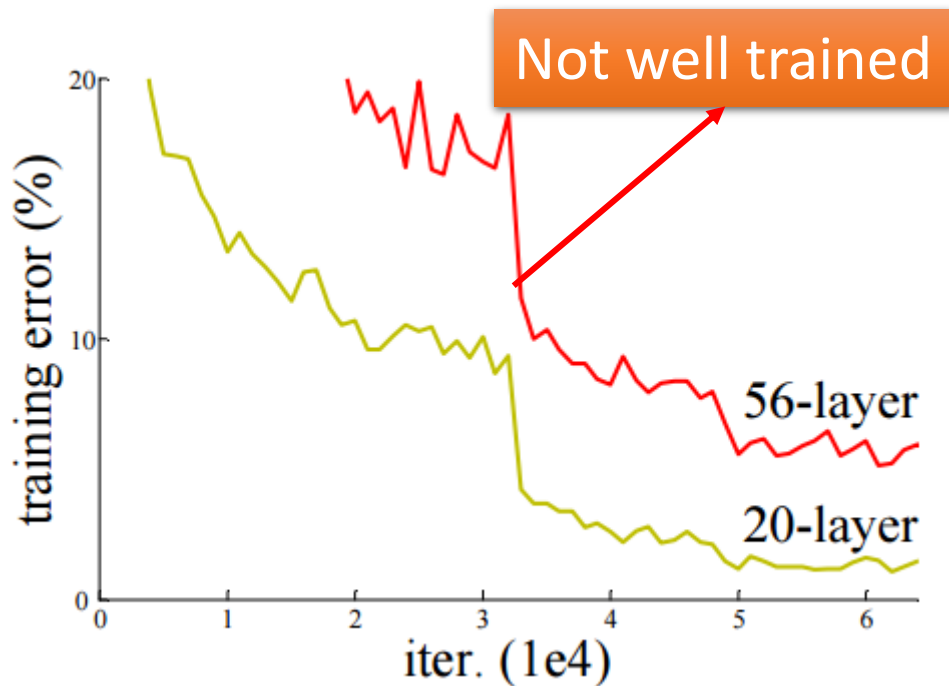
Hung-yi Lee

Prerequisite: <https://www.youtube.com/watch?v=xki61j7z-30>

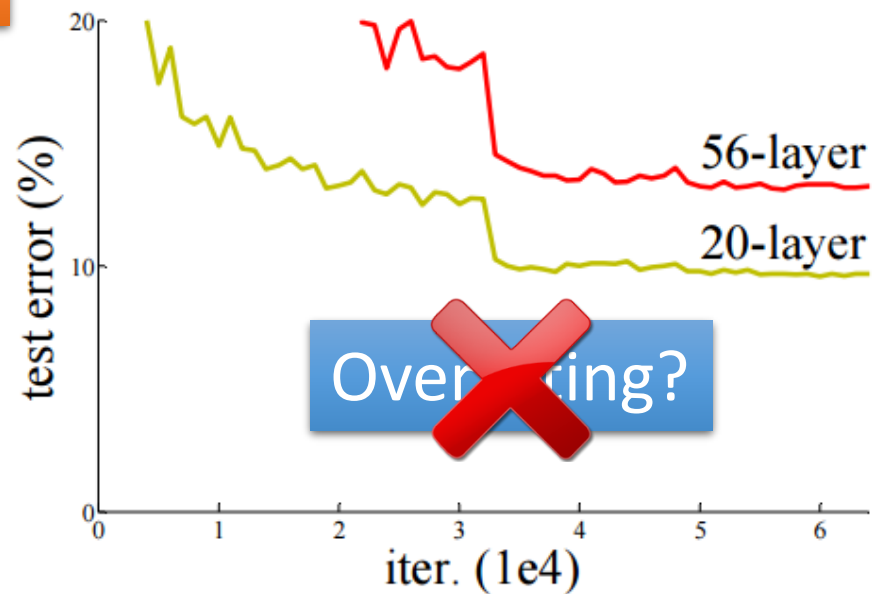
Recipe of Deep Learning



Do not always blame Overfitting



Training Data

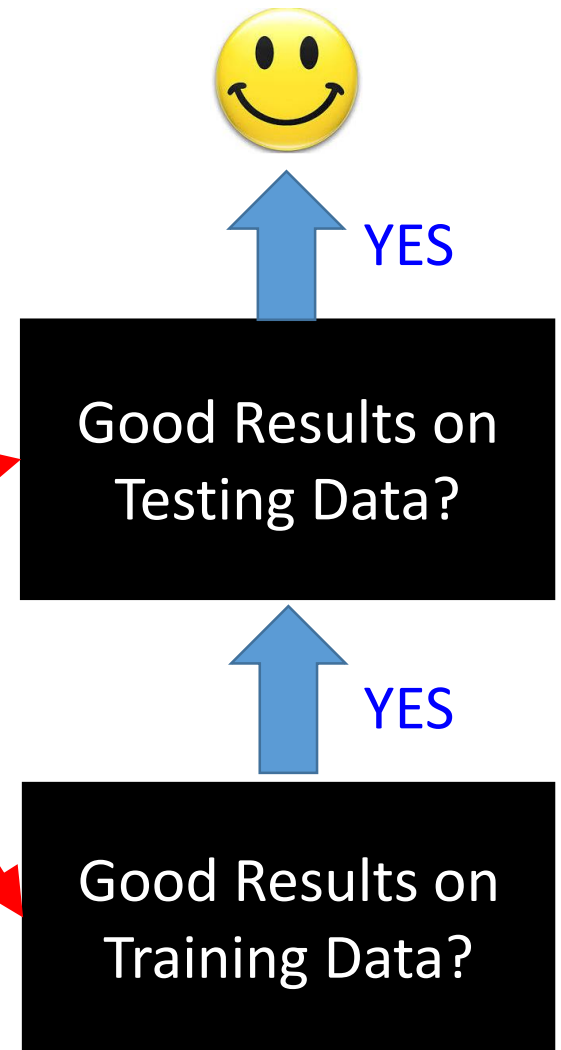
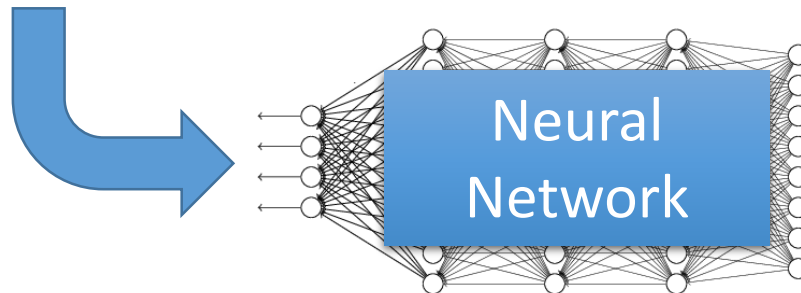


Testing Data

Recipe of Deep Learning

Different approaches for different problems.

e.g. dropout for good results on testing data



Outline

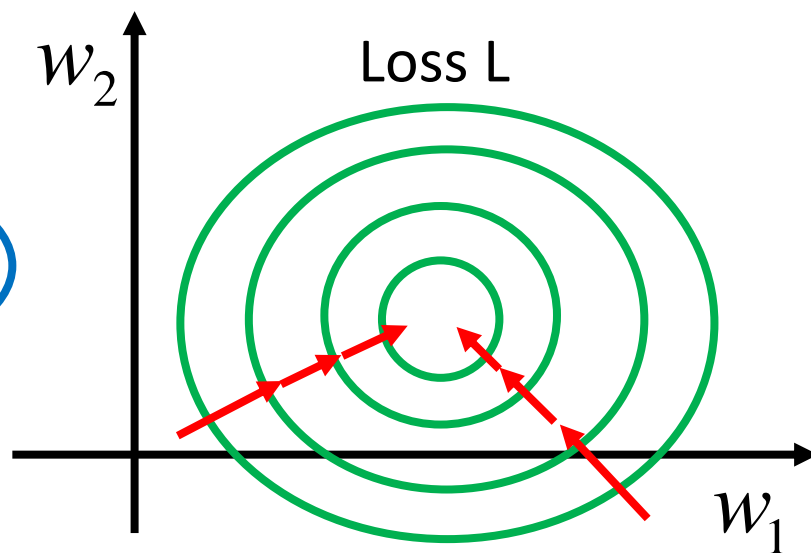
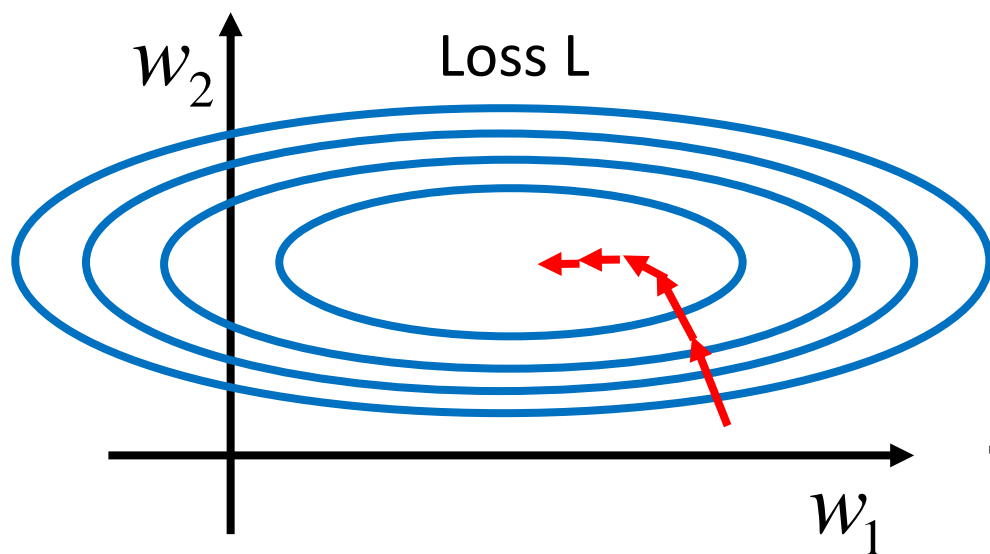
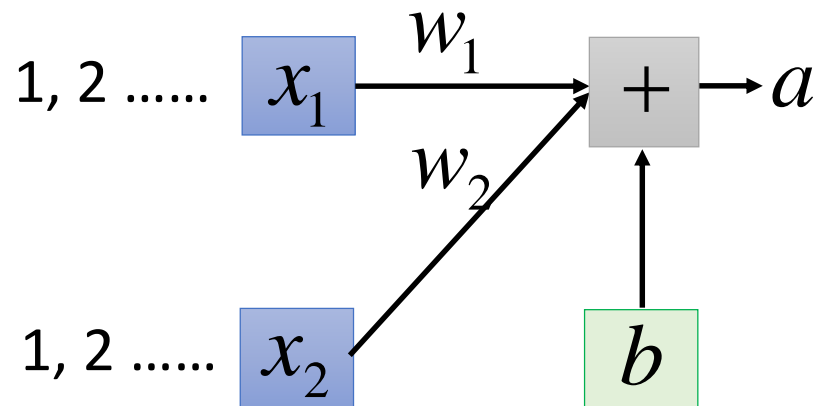
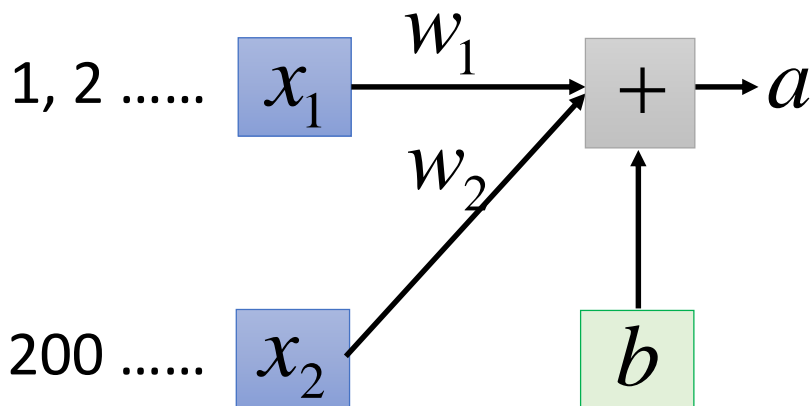
- Batch Normalization
- New Activation Function
- Tuning Hyperparameters
- Interesting facts (?) about deep learning
- Capsule
- New models for QA

Batch Normalization

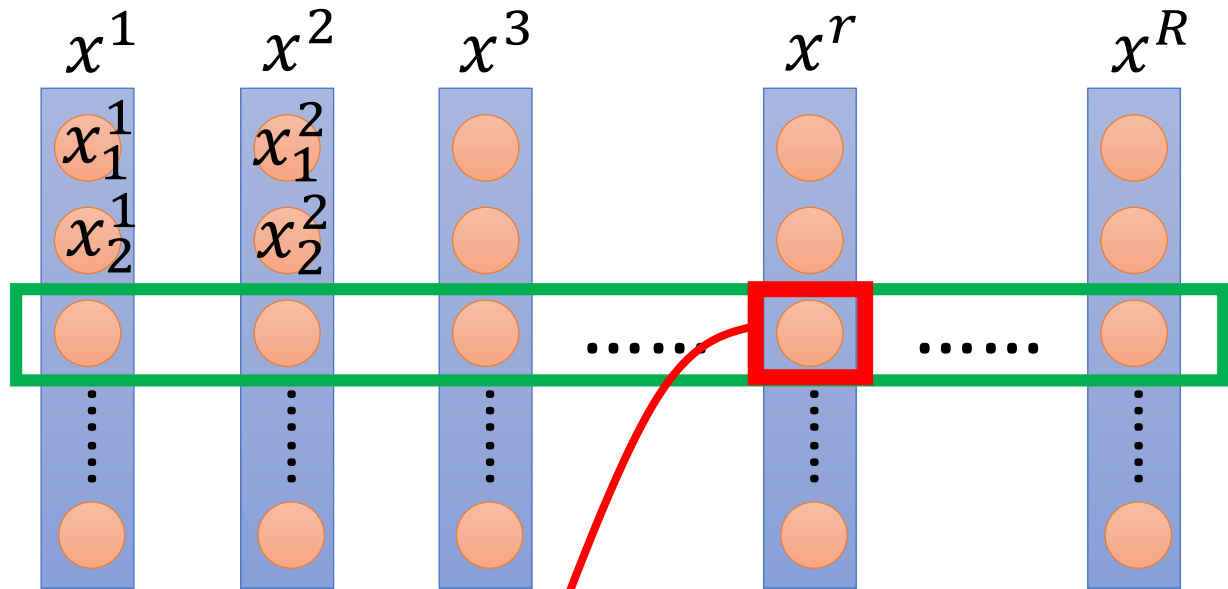
Sergey Ioffe, Christian Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, 2015

Feature Scaling

Make different features have the same scaling



Feature Scaling



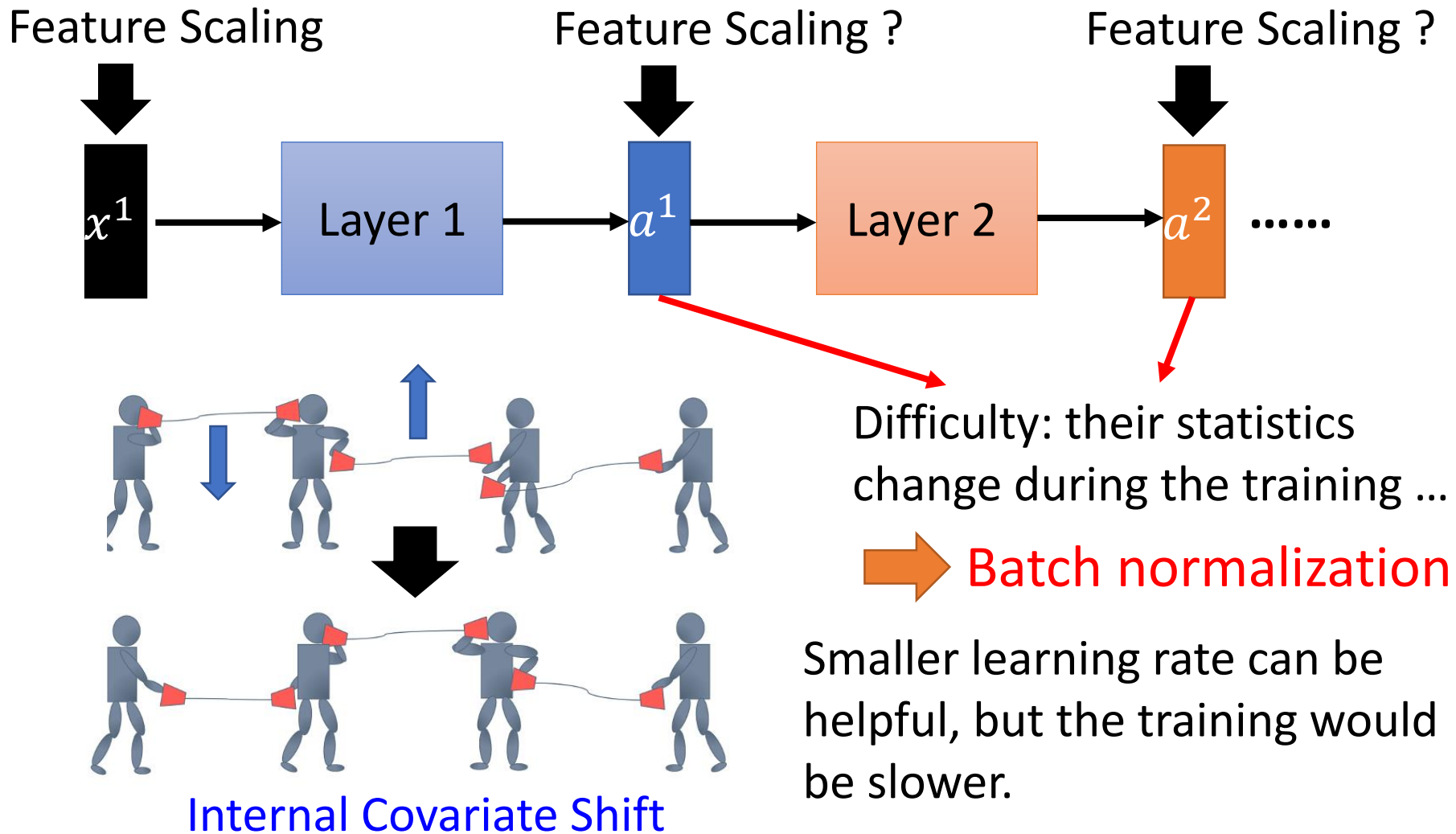
For each dimension i :
mean: m_i
standard deviation: σ_i

$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

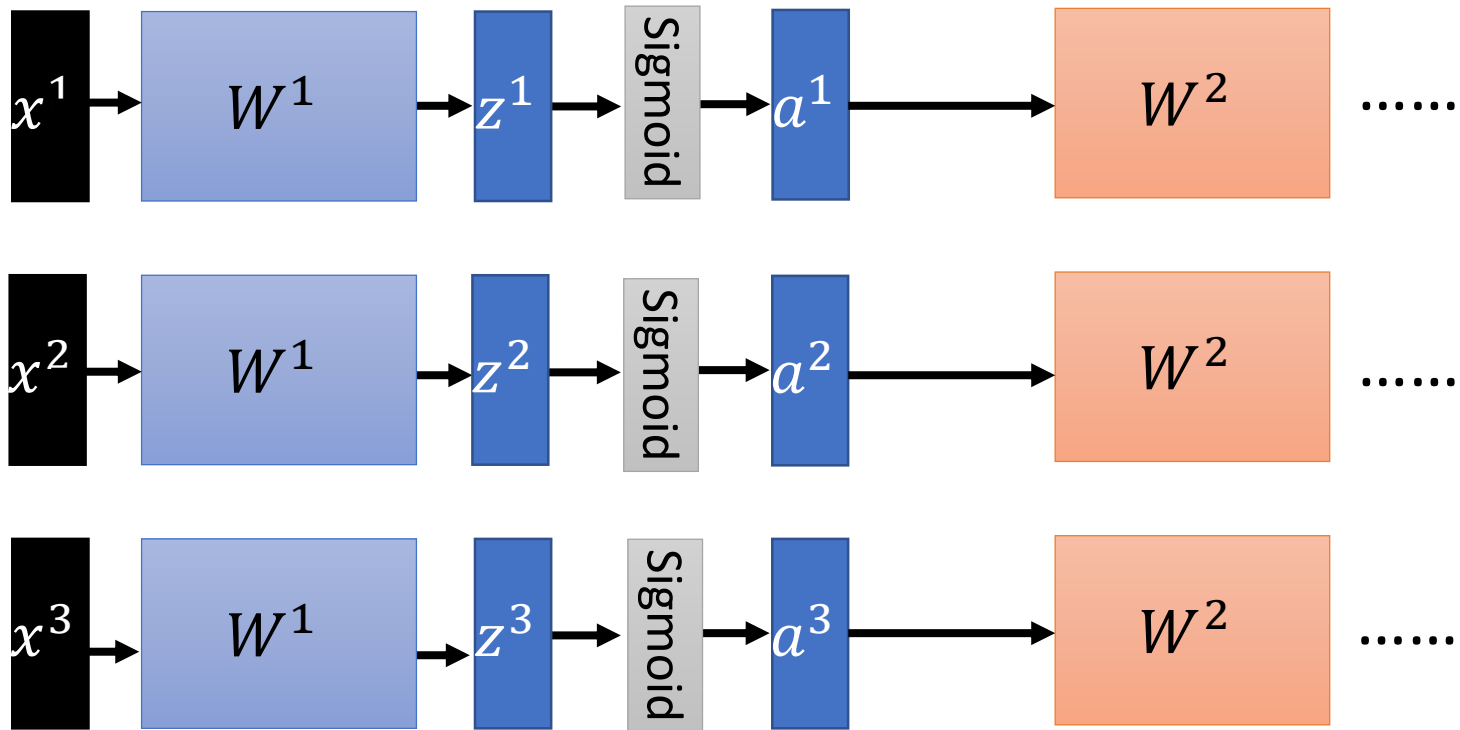
The means of all dimensions are 0, and the variances are all 1

In general, gradient descent converges much faster with feature scaling than without it.

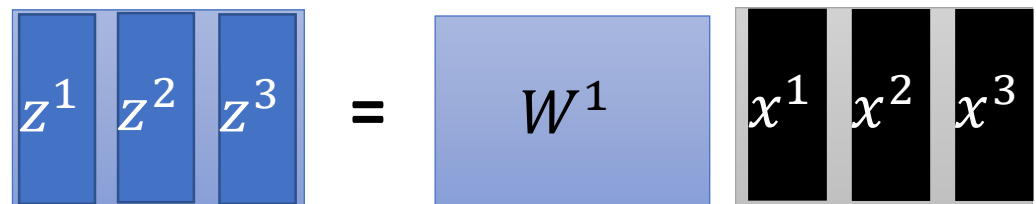
How about Hidden Layer?



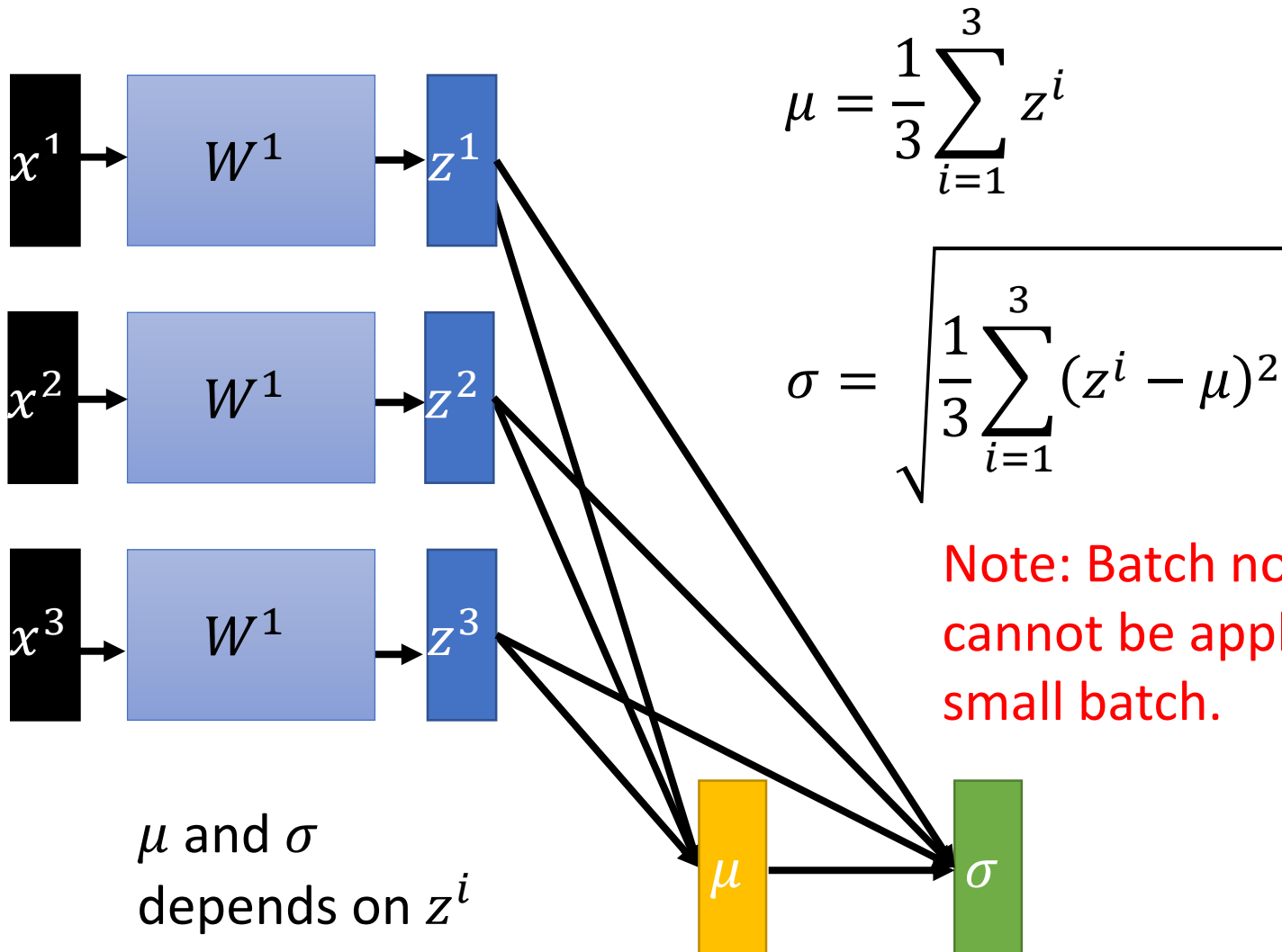
Batch



Batch

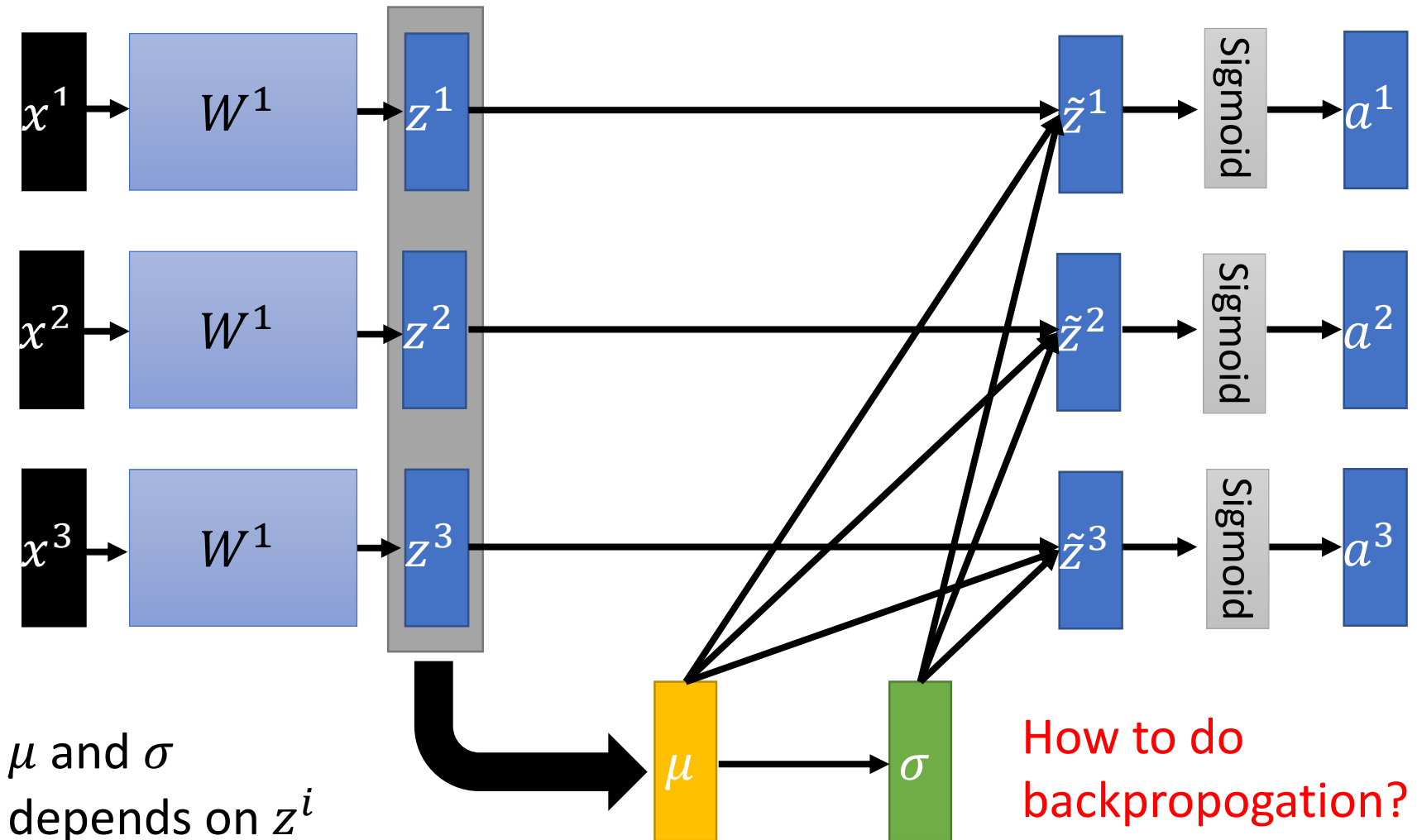


Batch normalization



Batch normalization

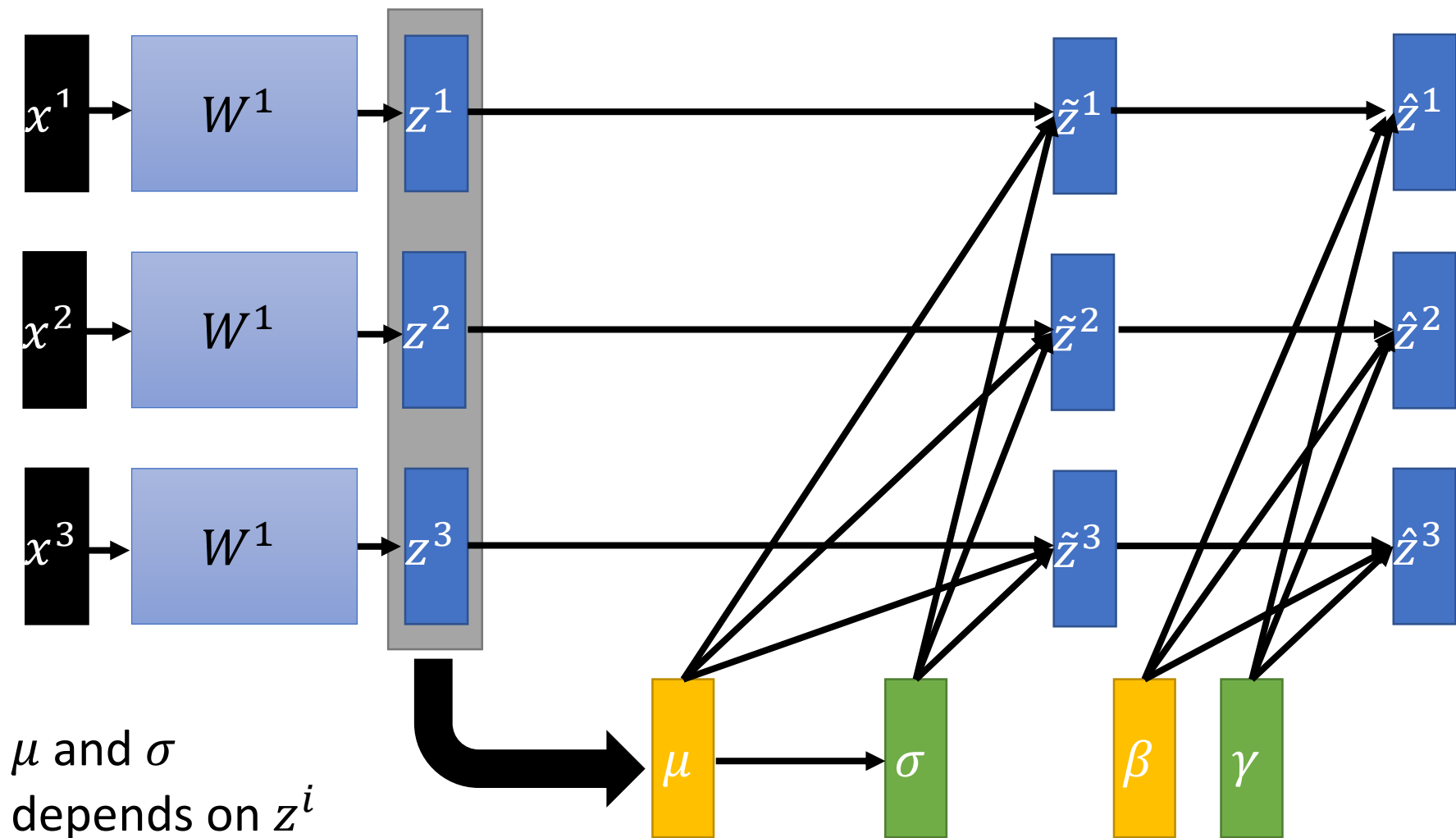
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$



Batch normalization

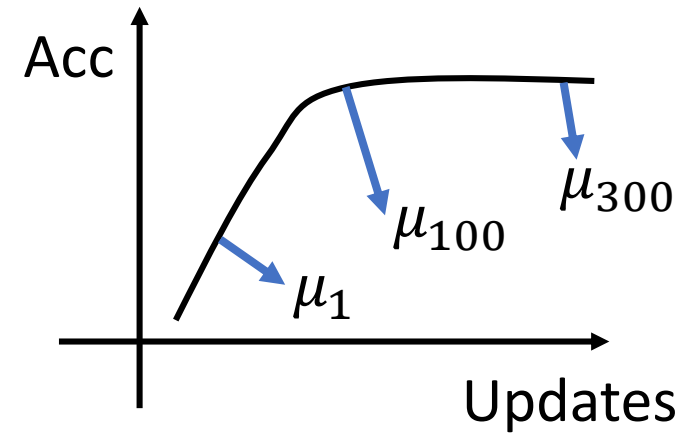
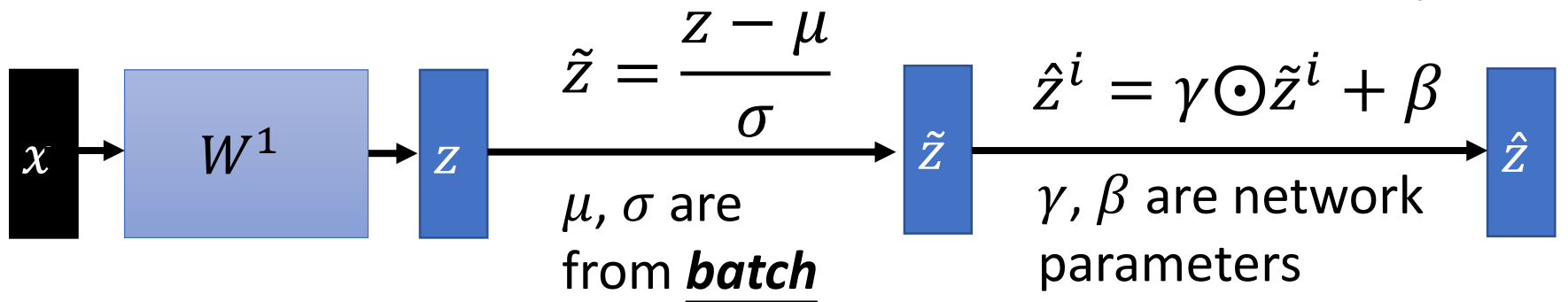
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$



Batch normalization

- At testing stage:



We do not have batch at testing stage.

Ideal solution:

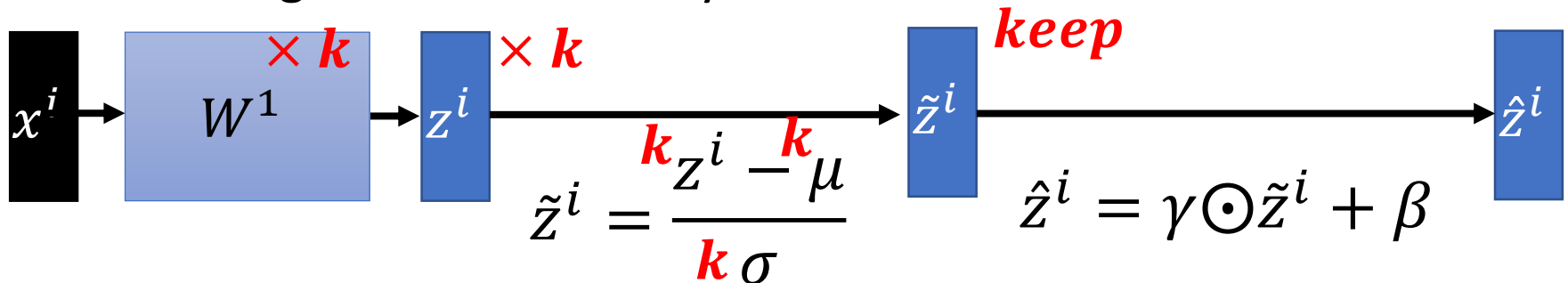
Computing μ and σ using the whole training dataset.

Practical solution:

Computing the moving average of μ and σ of the batches during training.

Batch normalization - Benefit

- BN reduces training times, and make very deep net trainable.
 - Because of less Covariate Shift, we can use larger learning rates.
 - Less exploding/vanishing gradients
 - Especially effective for sigmoid, tanh, etc.
- Learning is less affected by initialization.



- BN reduces the demand for regularization.

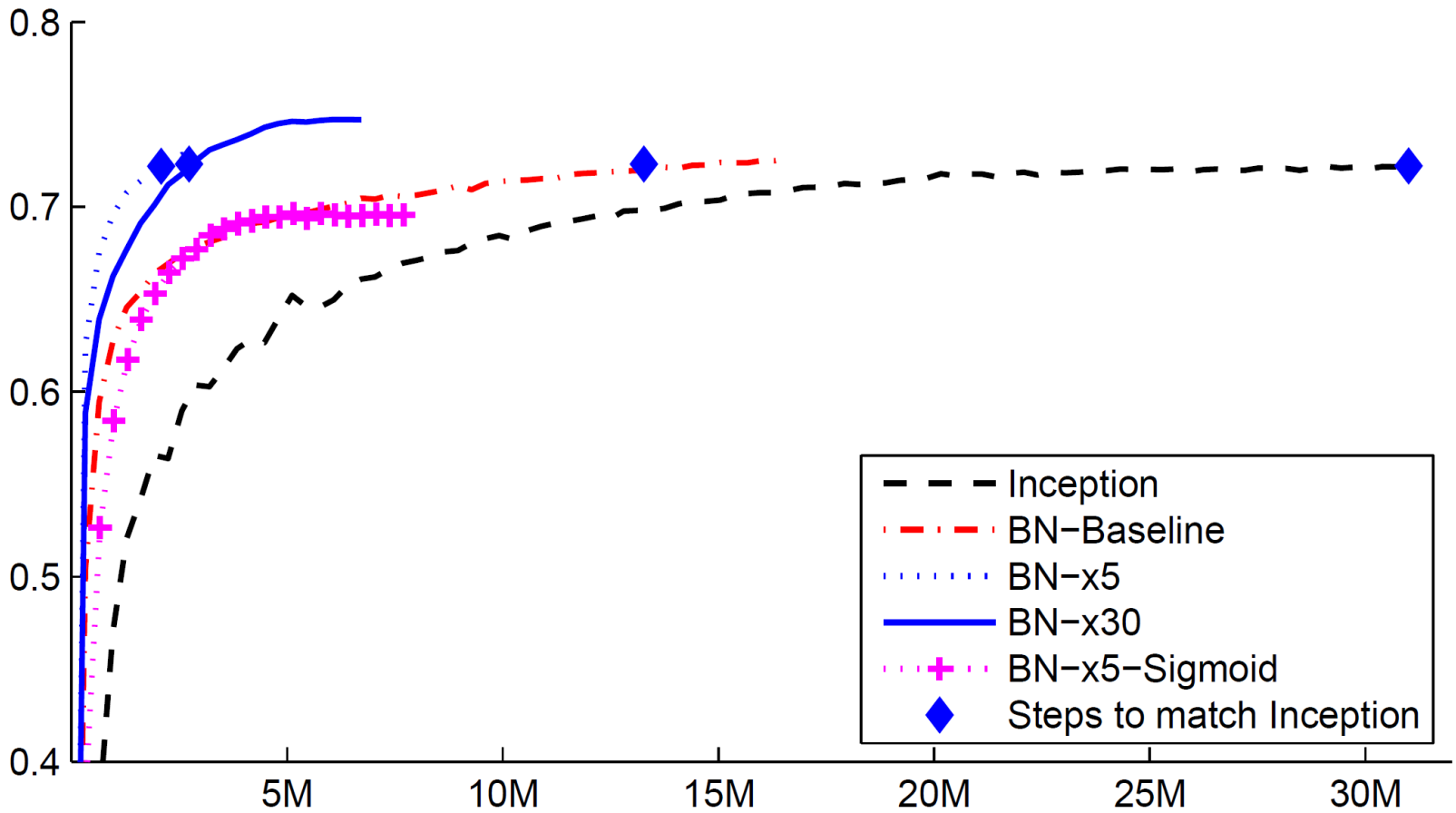


Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*

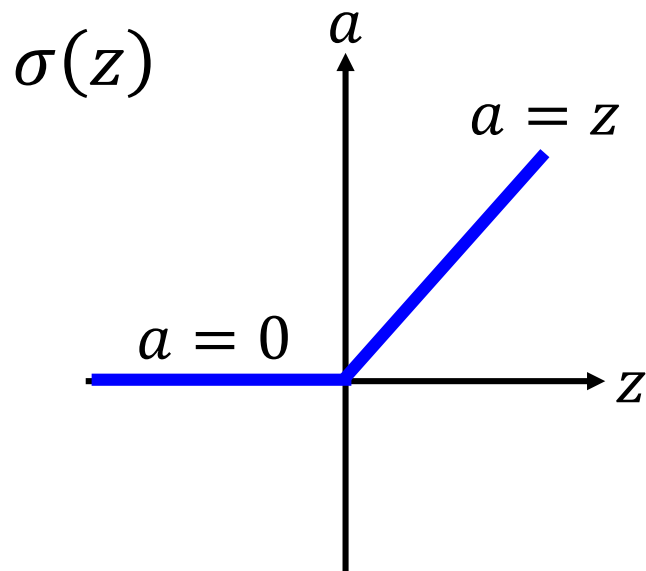
Demo

Activation Function

Günter Klambauer, Thomas Unterthiner, Andreas Mayr, Andreas Mayr,
“Self-Normalizing Neural Networks”, NIPS, 2017

ReLU

- Rectified Linear Unit (ReLU)



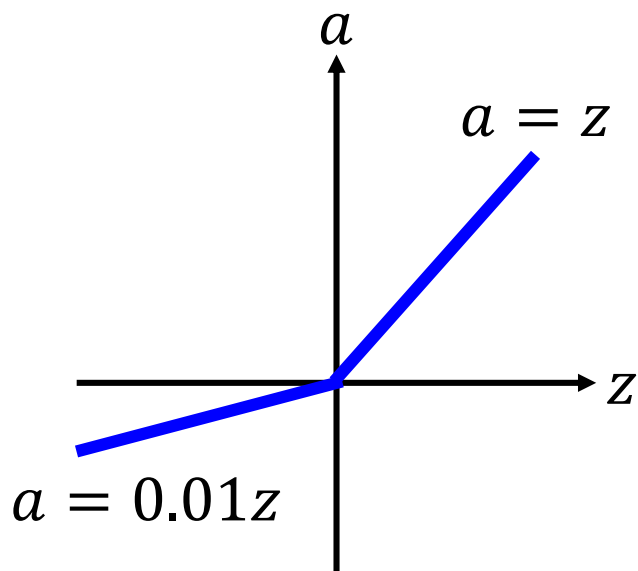
[Xavier Glorot, AISTATS'11]
[Andrew L. Maas, ICML'13]
[Kaiming He, arXiv'15]

Reason:

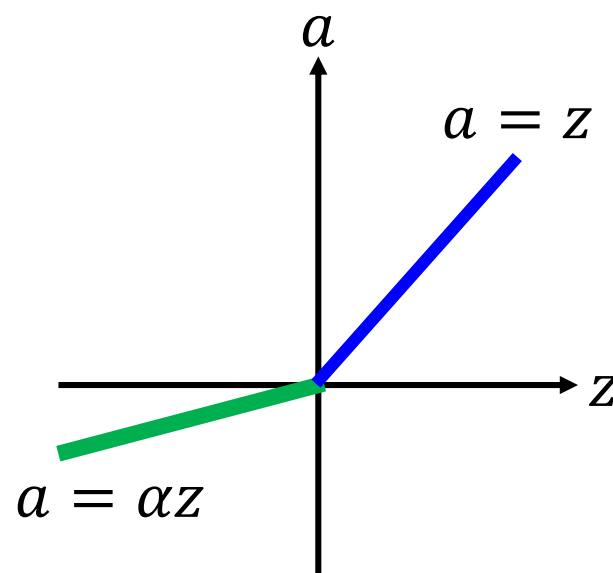
1. Fast to compute
2. Biological reason
3. Infinite sigmoid with different biases
4. Vanishing gradient problem

ReLU - variant

Leaky ReLU



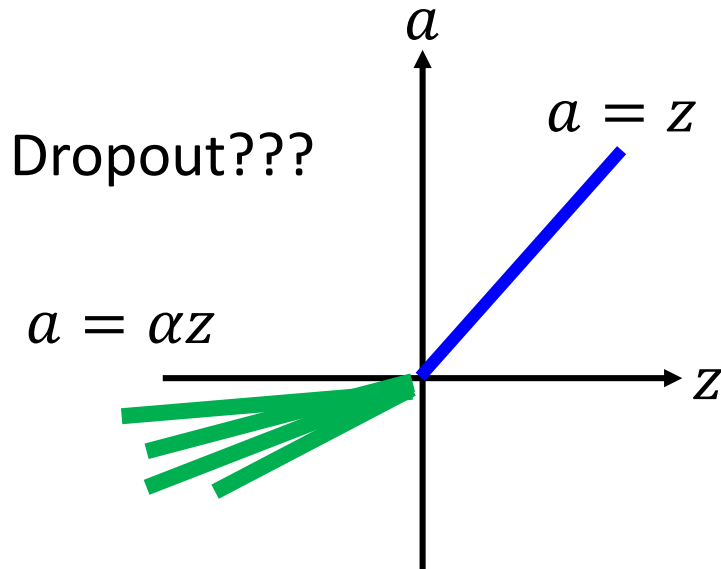
Parametric ReLU



α also learned by
gradient descent

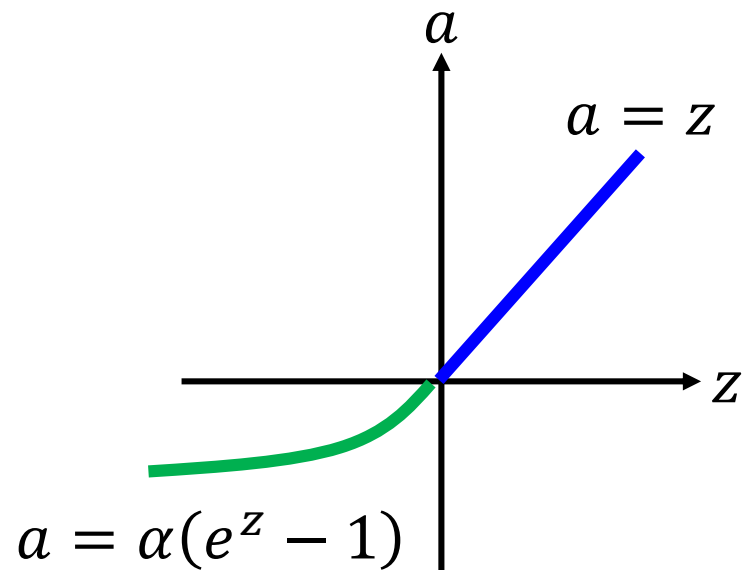
ReLU - variant

Randomized ReLU



α is sampled from a distribution during training.
Fixed during testing.

Exponential Linear Unit (ELU)



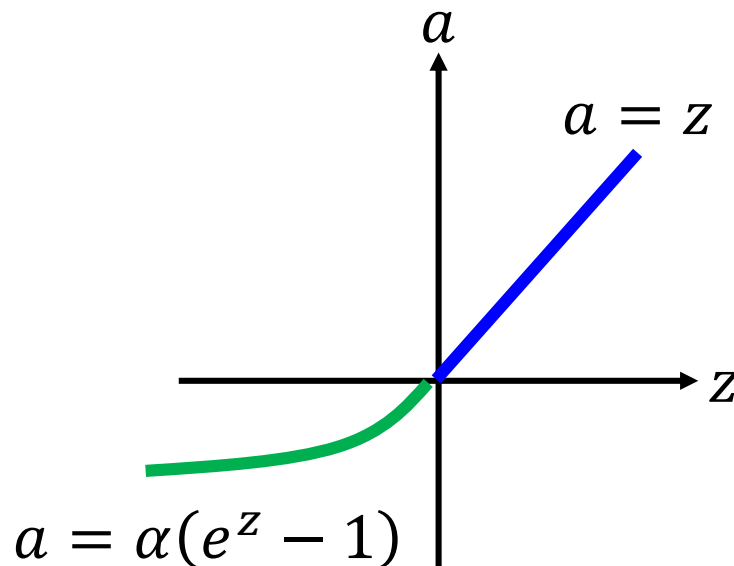
(1) Definition of scaled exponential linear units (SELUs)

In [3]:

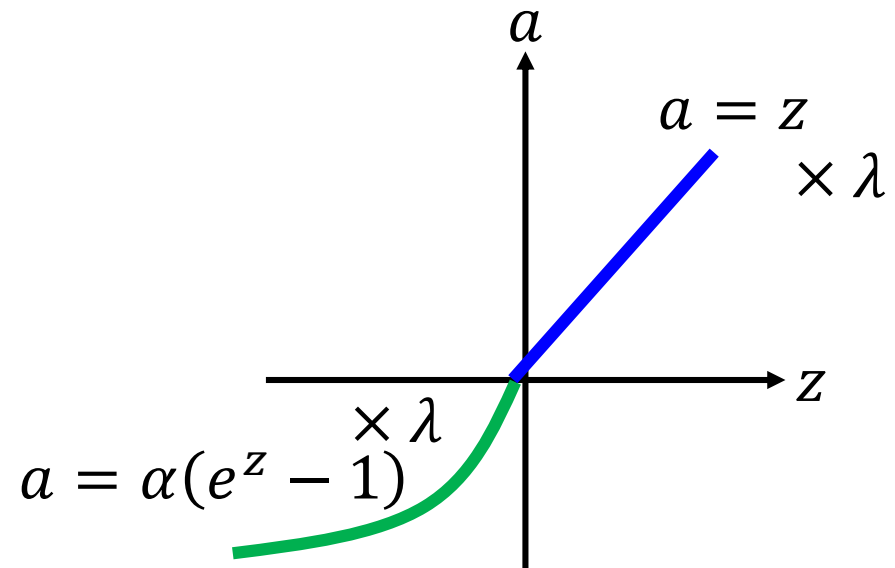
```
def selu(x):  
    with ops.name_scope('elu') as scope:  
        alpha = 1.6732632423543772848170429916717  
        scale = 1.0507009873554804934193349852946  
        return scale*tf.where(x>=0.0, x, alpha*tf.nn.elu(x))
```

<https://github.com/bioinf-jku/SNNs>

Exponential Linear Unit (ELU)



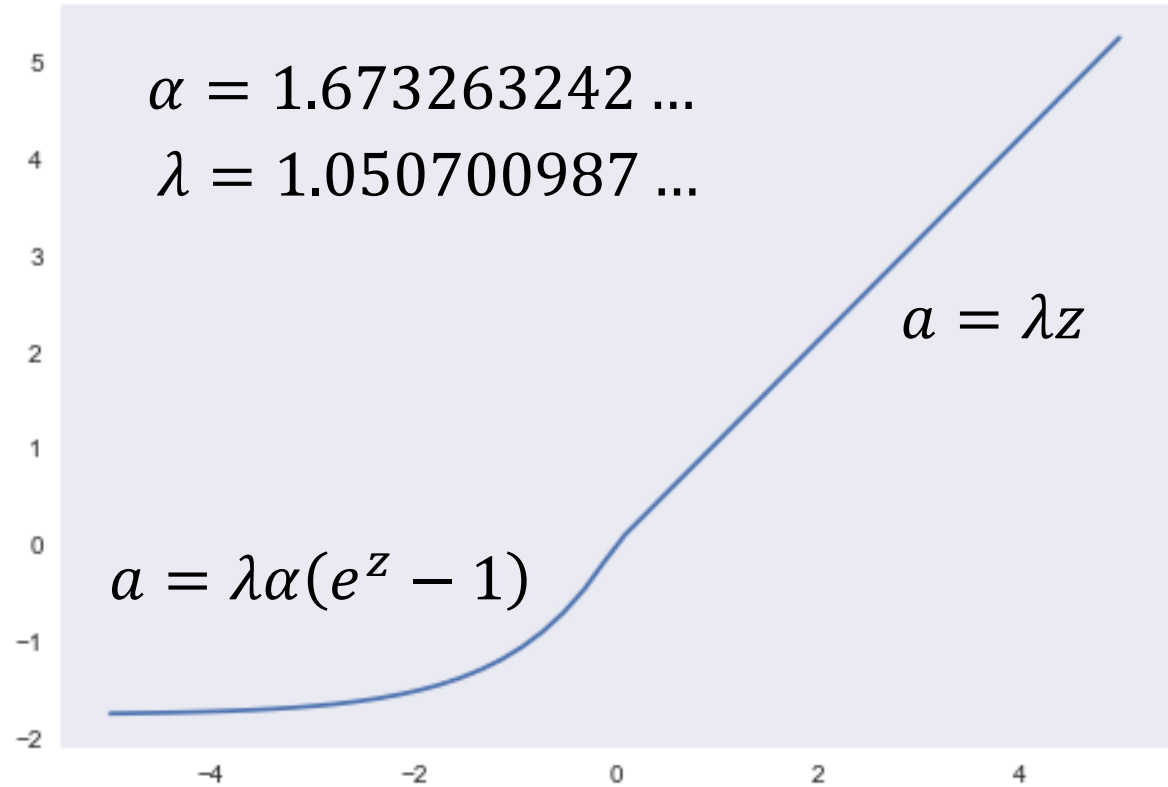
Scaled ELU (SELU)



$$\alpha = 1.6732632423543772848170429916717$$

$$\lambda = 1.0507009873554804934193349852946$$

SELU



Positive and negative values

➡ The whole ReLU family has this property except the original ReLU.

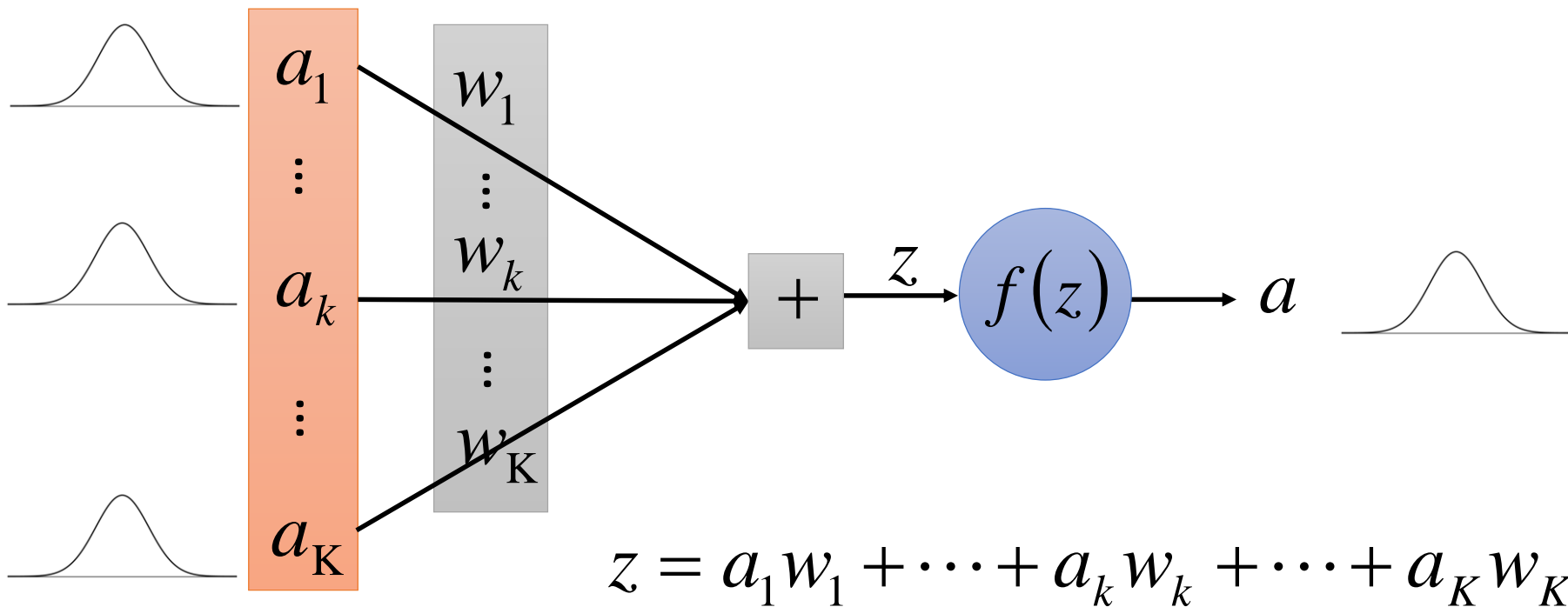
Saturation region ➡ ELU also has this property

Slope larger than 1 ➡ Only SELU also has this property

SELU

The inputs are i.i.d random variables with mean μ and variance σ^2 . $\mu = 0$ $\sigma^2 = 1$

$$\begin{aligned}\mu_z &= E[z] \\ &= \sum_{k=1}^K \frac{E[a_k]}{\mu} w_k = \mu \sum_{k=1}^K w_k = \mu \cdot K \mu_w \\ & \qquad \qquad \qquad = 0 \qquad \qquad = 0\end{aligned}$$



Do not have to be Gaussian

SELU

The inputs are i.i.d random variables with mean μ and variance σ^2 . $\mu = 0$, $\sigma^2 = 1$

$$\mu_z = 0 \quad \mu_w = 0$$

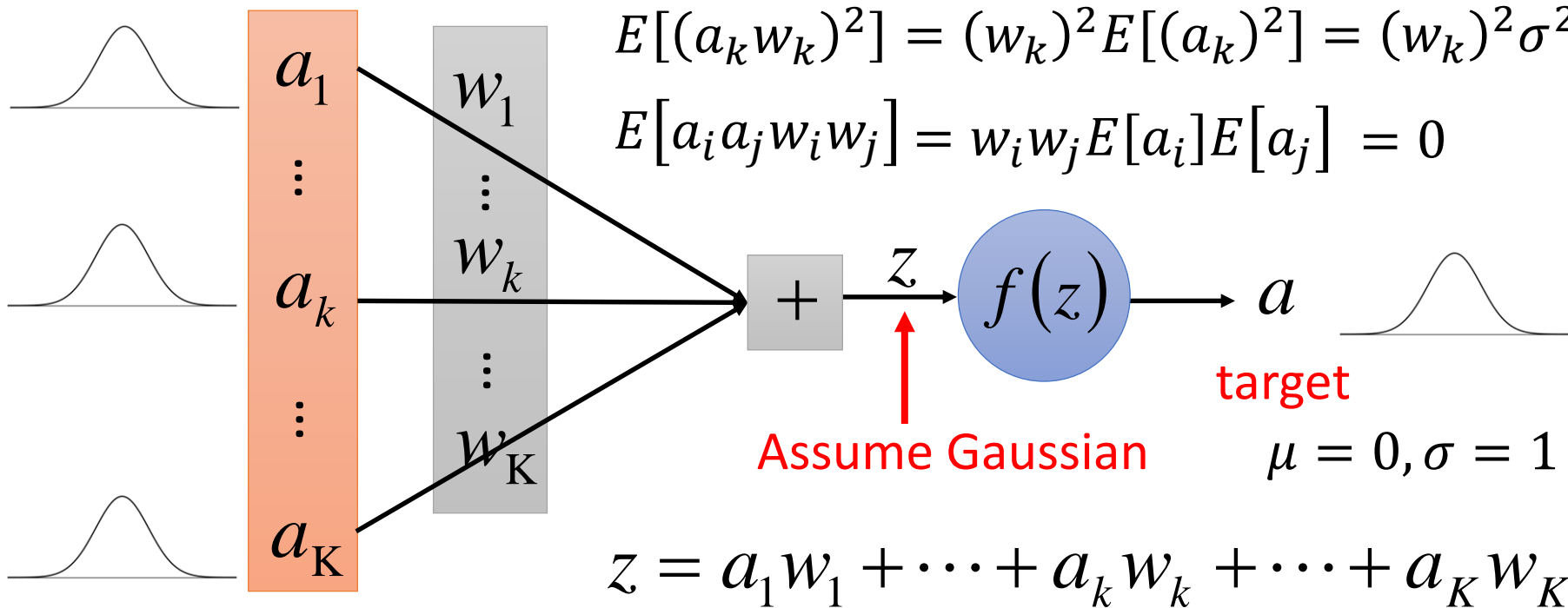
$$\sigma_z^2 = E[(z - \mu_z)^2] = E[z^2]$$

$$= E[(a_1 w_1 + a_2 w_2 + \dots)^2]$$

$$= \sum_{k=1}^K (w_k)^2 \sigma^2 = \sigma^2 \cdot K \sigma_w^2 = 1$$

$$E[(a_k w_k)^2] = (w_k)^2 E[(a_k)^2] = (w_k)^2 \sigma^2$$

$$E[a_i a_j w_i w_j] = w_i w_j E[a_i] E[a_j] = 0$$



$$z = a_1 w_1 + \dots + a_k w_k + \dots + a_K w_K$$

Demo

93 頁的證明

Source of joke:
<https://zhuanlan.zhihu.com/p/27336839>

SELU is actually more general.

$$\left. \frac{2(2x-y)(2x+y)2.911}{(\sqrt{2}\sqrt{x}) \left(\sqrt{\pi \left(\frac{2x+y}{\sqrt{x}} \right)^2 + 2.911^2 + \frac{(2.911-1)\sqrt{\pi}(2x+y)}{\sqrt{2}\sqrt{x}}} \right)} \right) \sqrt{x} - 0.0003 -$$
$$(3x-y) + \left(\frac{(\sqrt{2}\sqrt{2.911})(x-y)(x+y)}{(\sqrt{\pi(x+y)^2 + 2 \cdot 2.911^2 x + (2.911-1)(x+y)\sqrt{\pi}}) (\sqrt{2}\sqrt{x})} - \right.$$
$$\left. \frac{2(2x-y)(2x+y)(\sqrt{2}\sqrt{2.911})}{(\sqrt{2}\sqrt{x}) \left(\sqrt{\pi(2x+y)^2 + 2 \cdot 2.911^2 x + (2.911-1)(2x+y)\sqrt{\pi}} \right)} \right) \sqrt{x} - 0.0003 -$$
$$(3x-y) + 2.911 \left(\frac{(x-y)(x+y)}{(2.911-1)(x+y) + \sqrt{(x+y)^2 + \frac{2 \cdot 2.911^2 x}{\pi}}} - \right.$$
$$\left. \frac{2(2x-y)(2x+y)}{(2.911-1)(2x+y) + \sqrt{(2x+y)^2 + \frac{2 \cdot 2.911^2 x}{\pi}}} \right) - 0.0003 \geq$$
$$(3x-y) + 2.911 \left(\frac{(x-y)(x+y)}{(2.911-1)(x+y) + \sqrt{\left(\frac{2.911^2}{\pi} \right)^2 + (x+y)^2 + \frac{2 \cdot 2.911^2 x}{\pi} + \frac{2 \cdot 2.911^2 x}{\pi}}} - \right.$$
$$\left. \frac{2(2x-y)(2x+y)}{(2.911-1)(2x+y) + \sqrt{(2x+y)^2 + \frac{2 \cdot 2.911^2 x}{\pi}}} \right) - 0.0003 -$$
$$(3x-y) + 2.911 \left(\frac{(x-y)(x+y)}{(2.911-1)(x+y) + \sqrt{(x+y + \frac{2.911^2}{\pi})^2}} - \right.$$
$$\left. \frac{2(2x-y)(2x+y)}{(2.911-1)(2x+y) + \sqrt{(2x+y)^2 + \frac{2 \cdot 2.911^2 x}{\pi}}} \right) - 0.0003 -$$



Andrej Karpathy ✓

@karpathy

Following

maybe it's all generated by a char-rnn. I suspect we will never know.

RETWEETS LIKES

4

41



2:54 AM - 10 Jun 2017



5

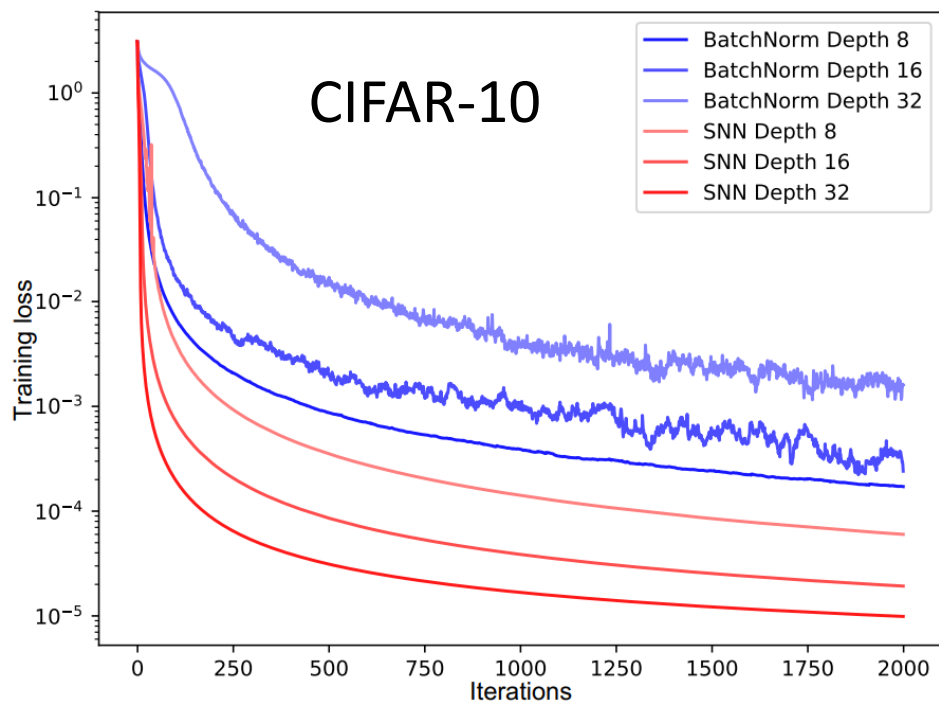
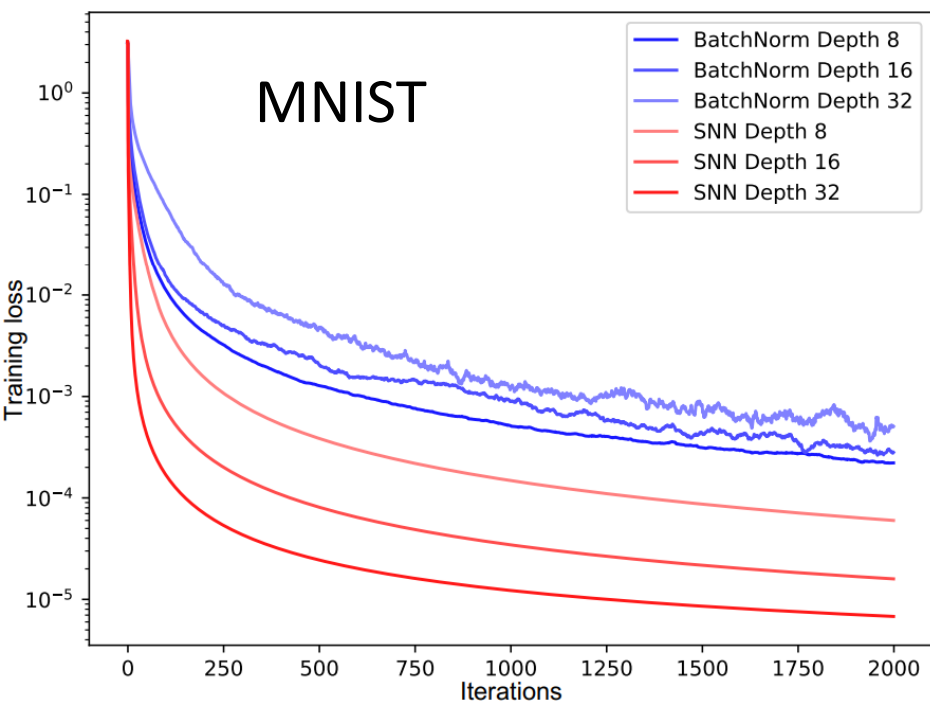


4



41





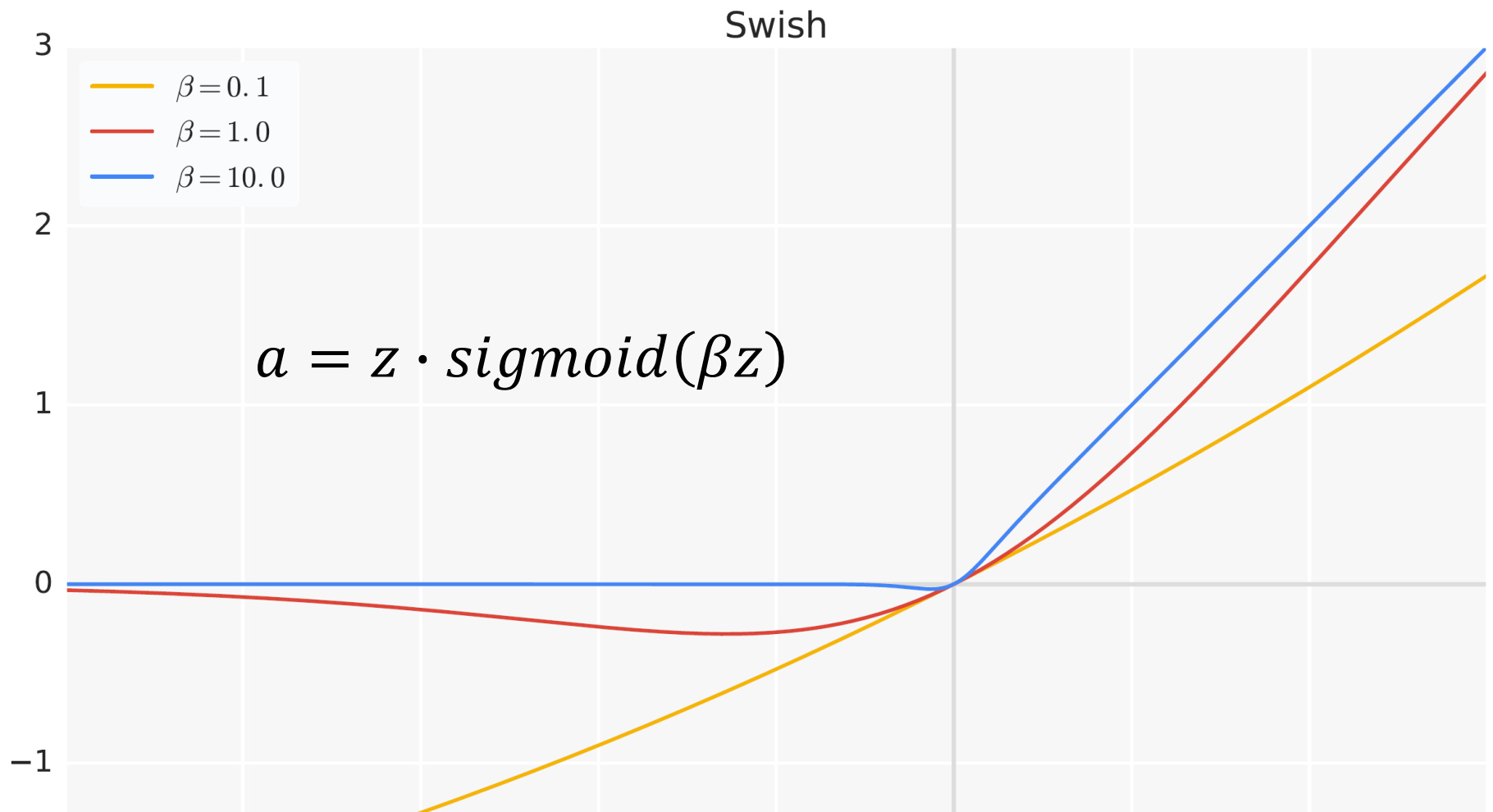
FNN method comparison

Method	avg. rank diff.	<i>p</i> -value
SNN	-0.756	
MSRAinit	-0.240*	2.7e-02
LayerNorm	-0.198*	1.5e-02
Highway	0.021*	1.9e-03
ResNet	0.273*	5.4e-04
WeightNorm	0.397*	7.8e-07
BatchNorm	0.504*	3.5e-06

ML method comparison

Method	avg. rank diff.	<i>p</i> -value
SNN	-6.7	
SVM	-6.4	5.8e-01
RandomForest	-5.9	2.1e-01
MSRAinit	-5.4*	4.5e-03
LayerNorm	-5.3	7.1e-02
Highway	-4.6*	1.7e-03
...

Demo



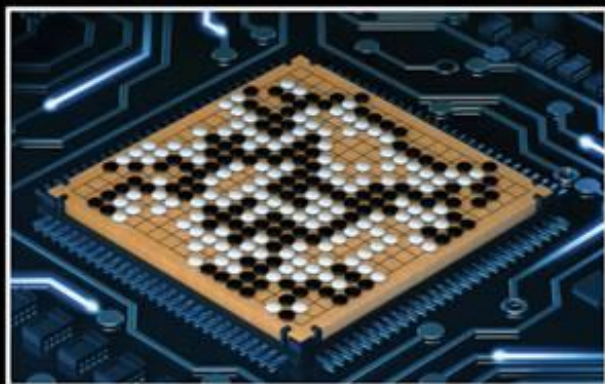
Baselines	ReLU	LReLU	PReLU	Softplus	ELU	SELU	GELU
Swish > Baseline	9	7	6	6	8	8	8
Swish = Baseline	0	1	3	2	0	1	1
Swish < Baseline	0	1	0	1	1	0	0

Figure 4: The Swish activation function.

Hyperparameters

Source of image: <https://medium.com/intuitionmachine/the-brute-force-method-of-deep-learning-innovation-58b497323ae5> (Denny Britz's graphic)

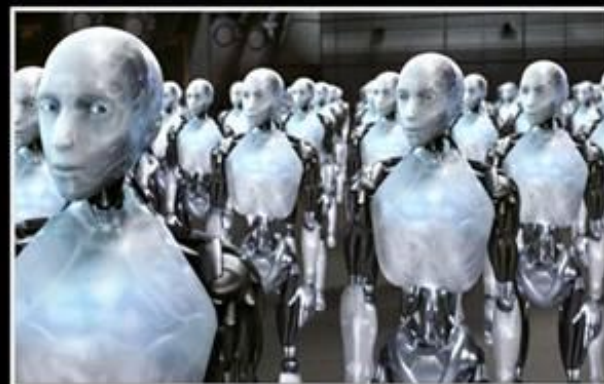
Deep Learning 研究生



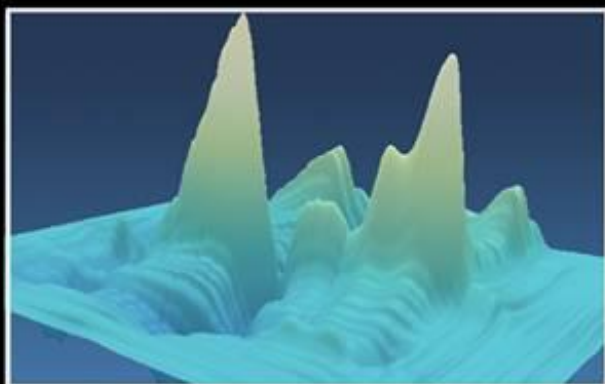
朋友覺得我在



我媽覺得我在



大眾覺得我在



指導教授覺得我在



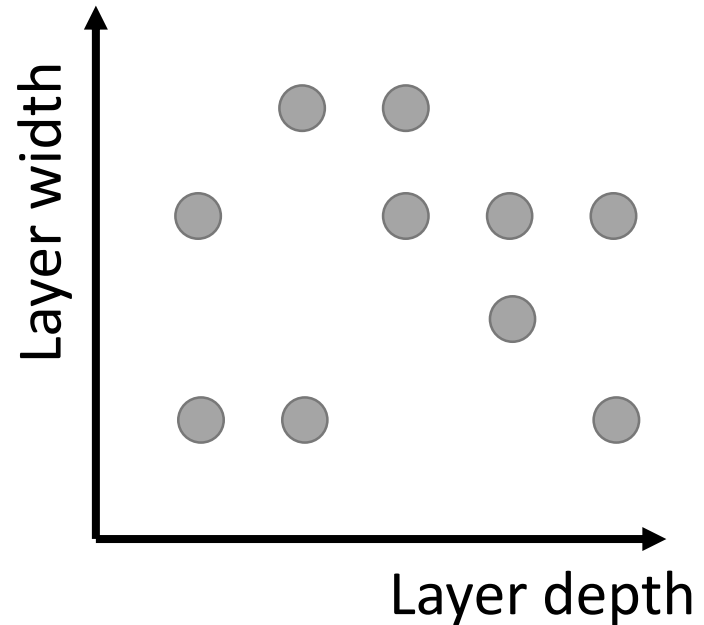
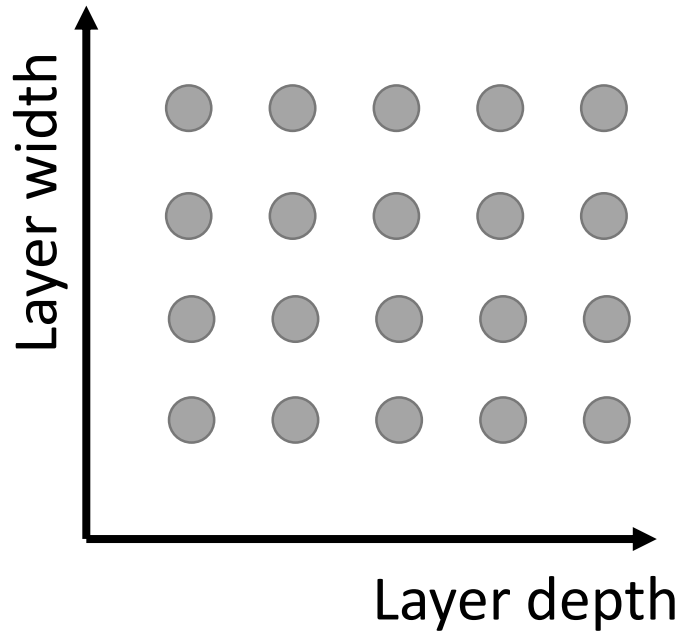
我以為我在

HYPERPARAMETER Spreadsheet			
Hyper-Param	Value	Config	Max Eval
num_epochs	100	100	0.000000000000000000
num_epochs	200	200	0.000000000000000000
num_epochs	300	300	0.000000000000000000
num_epochs	400	400	0.000000000000000000
num_epochs	500	500	0.000000000000000000
num_epochs	600	600	0.000000000000000000
num_epochs	700	700	0.000000000000000000
num_epochs	800	800	0.000000000000000000
num_epochs	900	900	0.000000000000000000
num_epochs	1000	1000	0.000000000000000000
num_epochs	1100	1100	0.000000000000000000
num_epochs	1200	1200	0.000000000000000000
num_epochs	1300	1300	0.000000000000000000
num_epochs	1400	1400	0.000000000000000000
num_epochs	1500	1500	0.000000000000000000
num_epochs	1600	1600	0.000000000000000000
num_epochs	1700	1700	0.000000000000000000
num_epochs	1800	1800	0.000000000000000000
num_epochs	1900	1900	0.000000000000000000
num_epochs	2000	2000	0.000000000000000000
num_epochs	2100	2100	0.000000000000000000
num_epochs	2200	2200	0.000000000000000000
num_epochs	2300	2300	0.000000000000000000
num_epochs	2400	2400	0.000000000000000000
num_epochs	2500	2500	0.000000000000000000
num_epochs	2600	2600	0.000000000000000000
num_epochs	2700	2700	0.000000000000000000
num_epochs	2800	2800	0.000000000000000000
num_epochs	2900	2900	0.000000000000000000
num_epochs	3000	3000	0.000000000000000000
num_epochs	3100	3100	0.000000000000000000
num_epochs	3200	3200	0.000000000000000000
num_epochs	3300	3300	0.000000000000000000
num_epochs	3400	3400	0.000000000000000000
num_epochs	3500	3500	0.000000000000000000
num_epochs	3600	3600	0.000000000000000000
num_epochs	3700	3700	0.000000000000000000
num_epochs	3800	3800	0.000000000000000000
num_epochs	3900	3900	0.000000000000000000
num_epochs	4000	4000	0.000000000000000000
num_epochs	4100	4100	0.000000000000000000
num_epochs	4200	4200	0.000000000000000000
num_epochs	4300	4300	0.000000000000000000
num_epochs	4400	4400	0.000000000000000000
num_epochs	4500	4500	0.000000000000000000
num_epochs	4600	4600	0.000000000000000000
num_epochs	4700	4700	0.000000000000000000
num_epochs	4800	4800	0.000000000000000000
num_epochs	4900	4900	0.000000000000000000
num_epochs	5000	5000	0.000000000000000000
num_epochs	5100	5100	0.000000000000000000
num_epochs	5200	5200	0.000000000000000000
num_epochs	5300	5300	0.000000000000000000
num_epochs	5400	5400	0.000000000000000000
num_epochs	5500	5500	0.000000000000000000
num_epochs	5600	5600	0.000000000000000000
num_epochs	5700	5700	0.000000000000000000
num_epochs	5800	5800	0.000000000000000000
num_epochs	5900	5900	0.000000000000000000
num_epochs	6000	6000	0.000000000000000000
num_epochs	6100	6100	0.000000000000000000
num_epochs	6200	6200	0.000000000000000000
num_epochs	6300	6300	0.000000000000000000
num_epochs	6400	6400	0.000000000000000000
num_epochs	6500	6500	0.000000000000000000
num_epochs	6600	6600	0.000000000000000000
num_epochs	6700	6700	0.000000000000000000
num_epochs	6800	6800	0.000000000000000000
num_epochs	6900	6900	0.000000000000000000
num_epochs	7000	7000	0.000000000000000000
num_epochs	7100	7100	0.000000000000000000
num_epochs	7200	7200	0.000000000000000000
num_epochs	7300	7300	0.000000000000000000
num_epochs	7400	7400	0.000000000000000000
num_epochs	7500	7500	0.000000000000000000
num_epochs	7600	7600	0.000000000000000000
num_epochs	7700	7700	0.000000000000000000
num_epochs	7800	7800	0.000000000000000000
num_epochs	7900	7900	0.000000000000000000
num_epochs	8000	8000	0.000000000000000000
num_epochs	8100	8100	0.000000000000000000
num_epochs	8200	8200	0.000000000000000000
num_epochs	8300	8300	0.000000000000000000
num_epochs	8400	8400	0.000000000000000000
num_epochs	8500	8500	0.000000000000000000
num_epochs	8600	8600	0.000000000000000000
num_epochs	8700	8700	0.000000000000000000
num_epochs	8800	8800	0.000000000000000000
num_epochs	8900	8900	0.000000000000000000
num_epochs	9000	9000	0.000000000000000000
num_epochs	9100	9100	0.000000000000000000
num_epochs	9200	9200	0.000000000000000000
num_epochs	9300	9300	0.000000000000000000
num_epochs	9400	9400	0.000000000000000000
num_epochs	9500	9500	0.000000000000000000
num_epochs	9600	9600	0.000000000000000000
num_epochs	9700	9700	0.000000000000000000
num_epochs	9800	9800	0.000000000000000000
num_epochs	9900	9900	0.000000000000000000
num_epochs	10000	10000	0.000000000000000000

事實上我在

感謝 沈昇勳 同學提供圖檔

Grid Search v.s. Random Search



Assumption: top K results are good enough

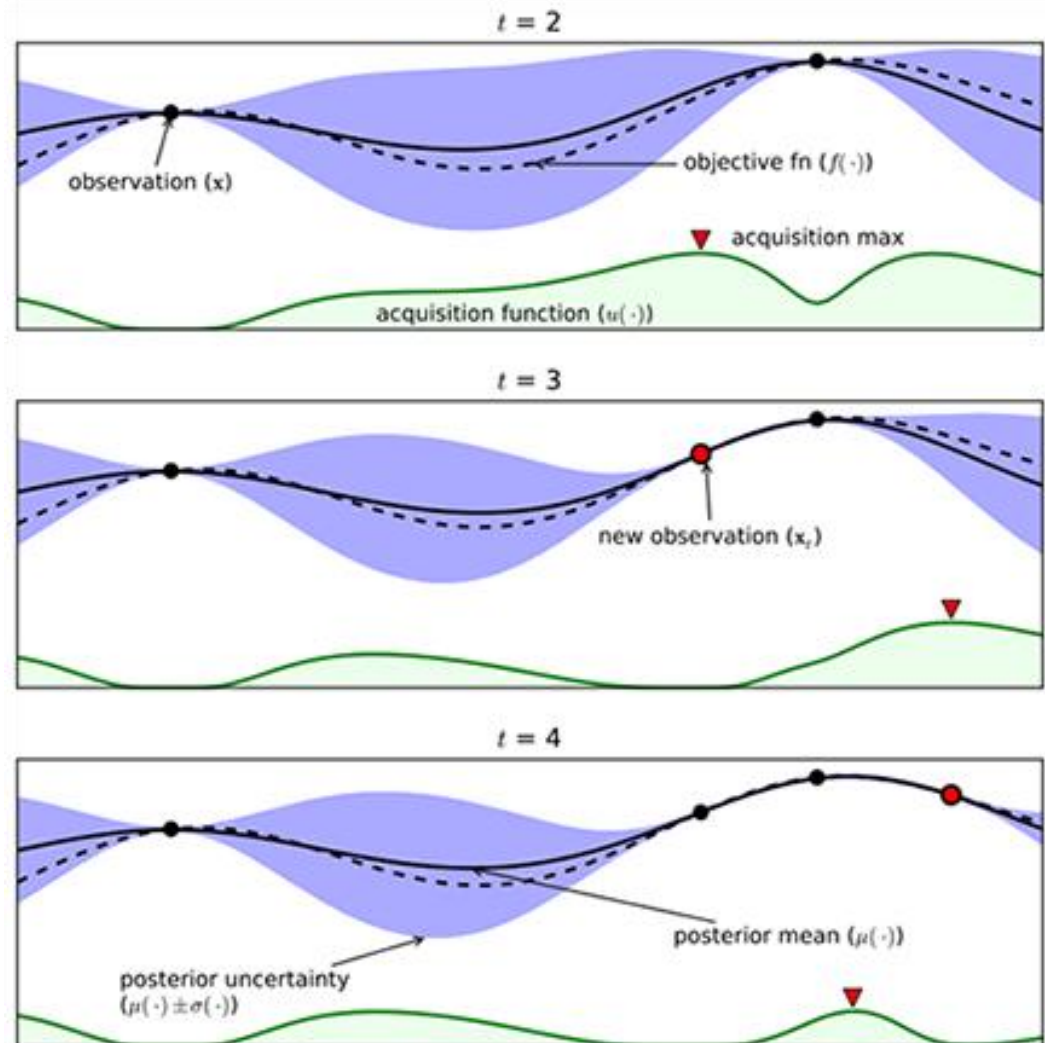
If there are N points, probability K/N that your sample is in top K

Sample x times: $1 - (1 - K/N)^x > 90\%$

If $N = 1000$, $K = 10$ \longrightarrow $x = 230$

$K = 100$ \longrightarrow $x = 22$

Model-based Hyperparameter Optimization

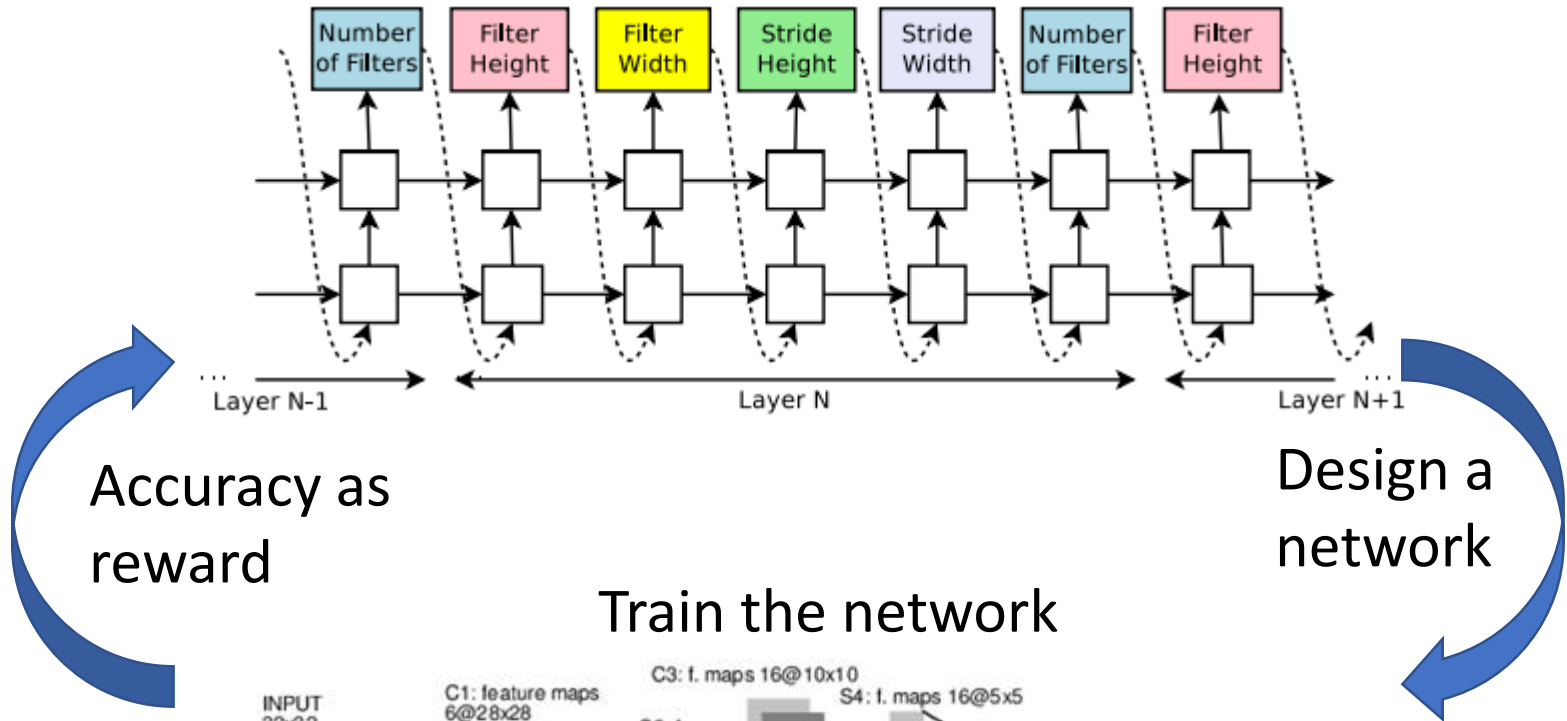


<https://cloud.google.com/blog/big-data/2017/08/hyperparameter-tuning-in-cloud-machine-learning-engine-using-bayesian-optimization>

Reinforcement Learning

It can design LSTM as shown in the previous lecture.

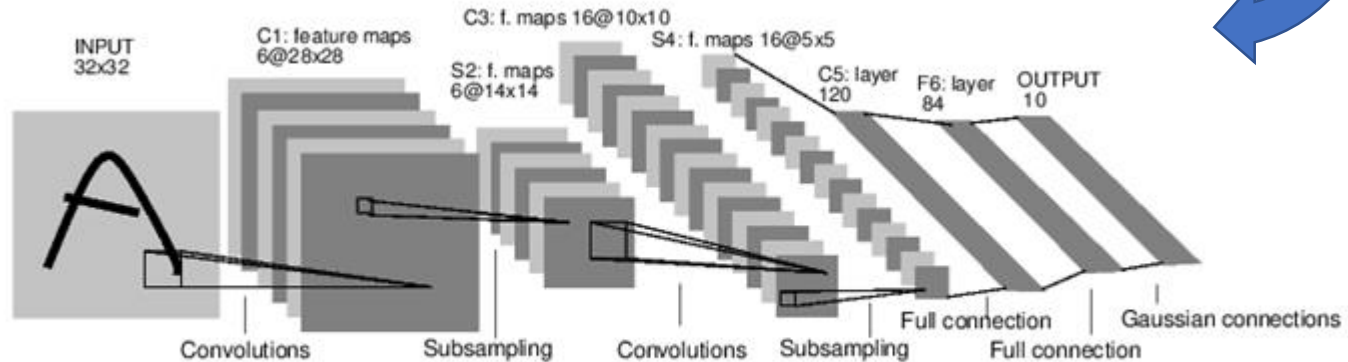
One kind of meta learning (or learn to learn)



Accuracy as reward

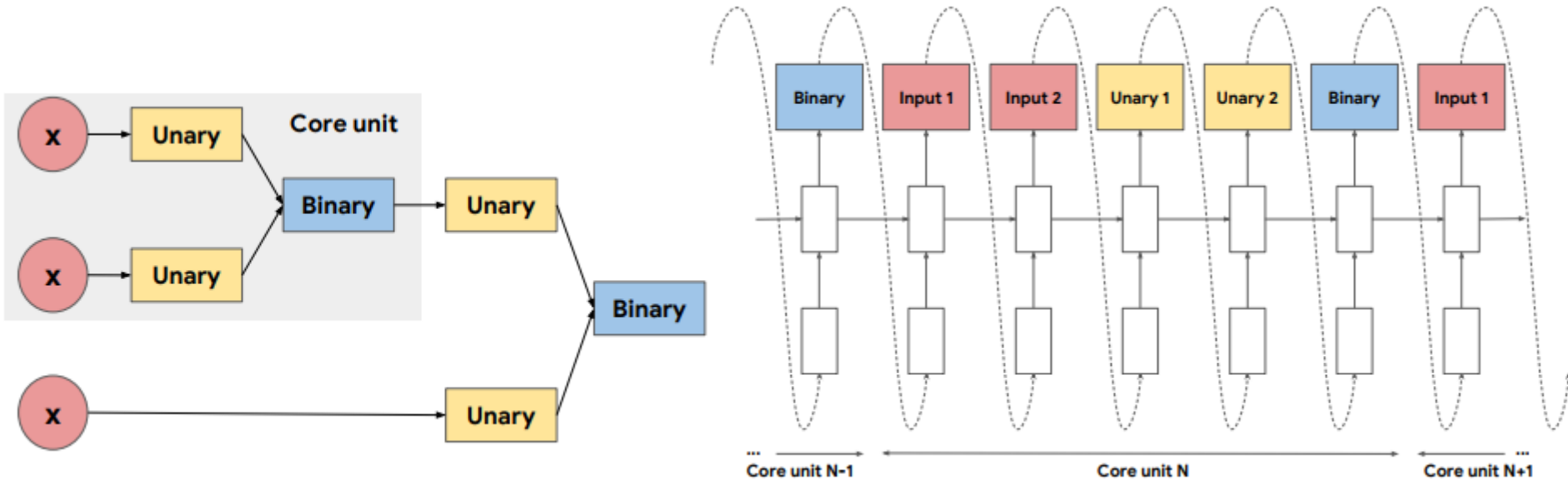
Design a network

Train the network



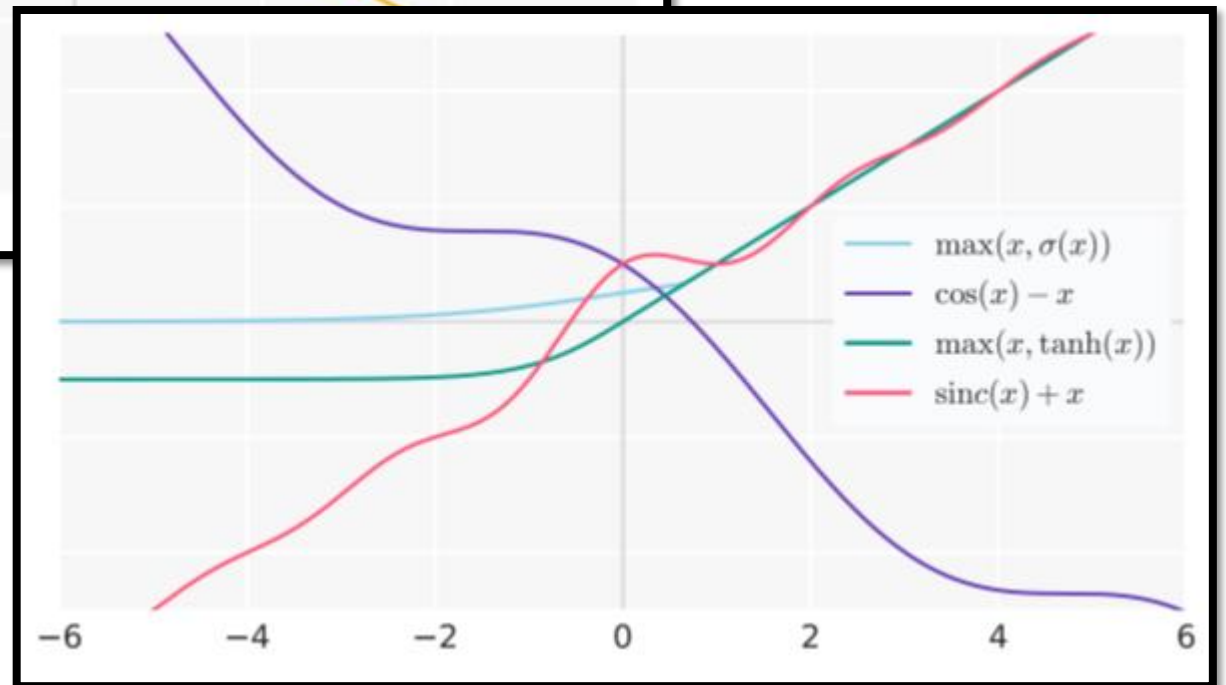
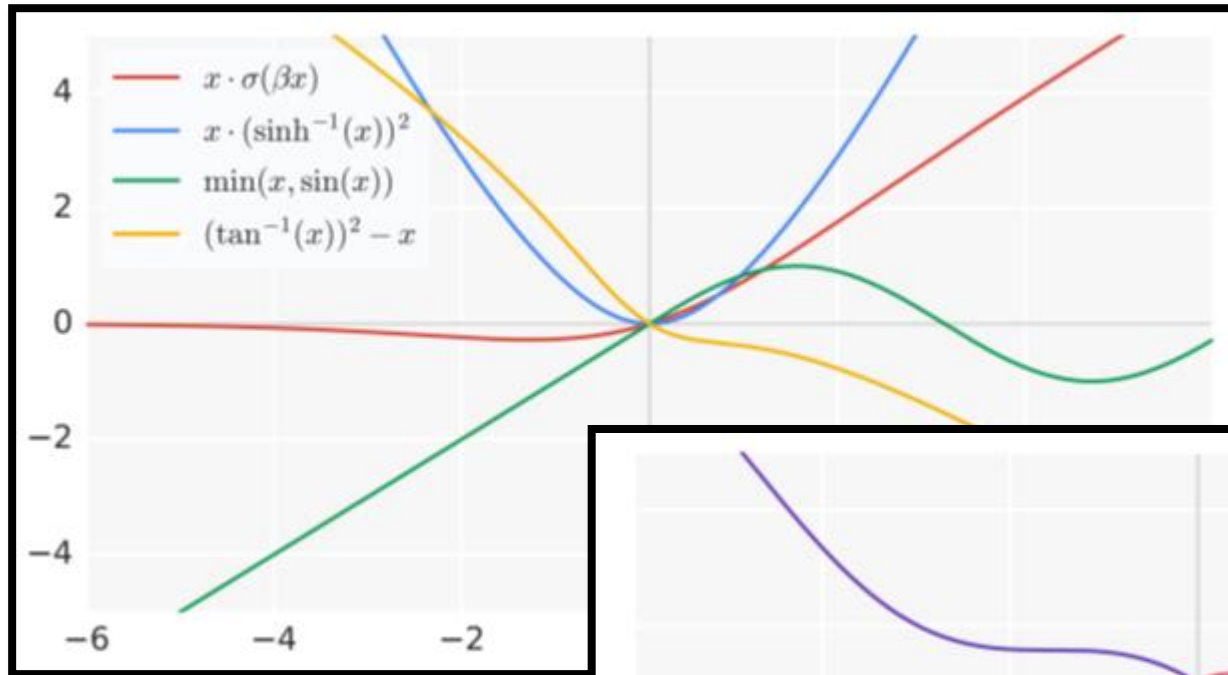
A Full Convolutional Neural Network (LeNet)

SWISH



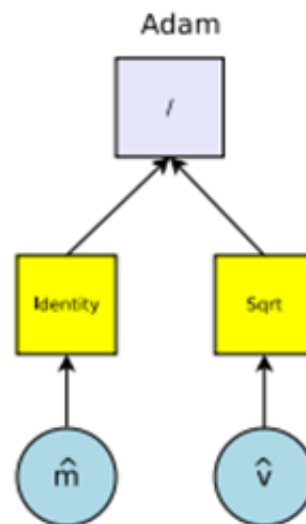
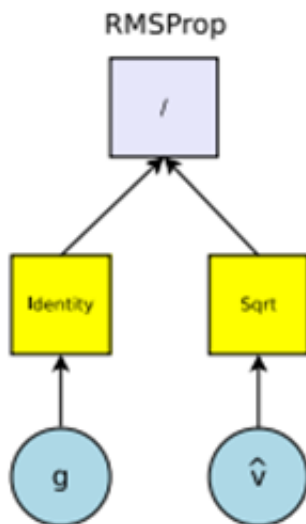
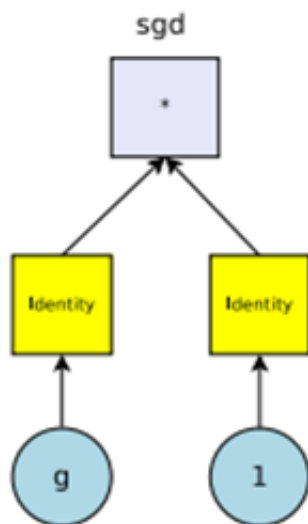
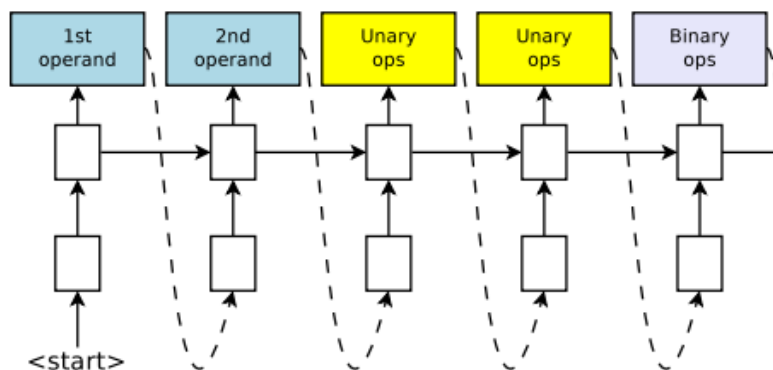
- Unary functions:** x , $-x$, $|x|$, x^2 , x^3 , \sqrt{x} , βx , $x + \beta$, $\log(|x| + \epsilon)$, $\exp(x) \sin(x)$, $\cos(x)$, $\sinh(x)$, $\cosh(x)$, $\tanh(x)$, $\sinh^{-1}(x)$, $\tan^{-1}(x)$, $\text{sinc}(x)$, $\max(x, 0)$, $\min(x, 0)$, $\sigma(x)$, $\log(1 + \exp(x))$, $\exp(-x^2)$, $\text{erf}(x)$, β
- Binary functions:** $x_1 + x_2$, $x_1 \cdot x_2$, $x_1 - x_2$, $\frac{x_1}{x_2 + \epsilon}$, $\max(x_1, x_2)$, $\min(x_1, x_2)$, $\sigma(x_1) \cdot x_2$, $\exp(-\beta(x_1 - x_2)^2)$, $\exp(-\beta|x_1 - x_2|)$, $\beta x_1 + (1 - \beta)x_2$

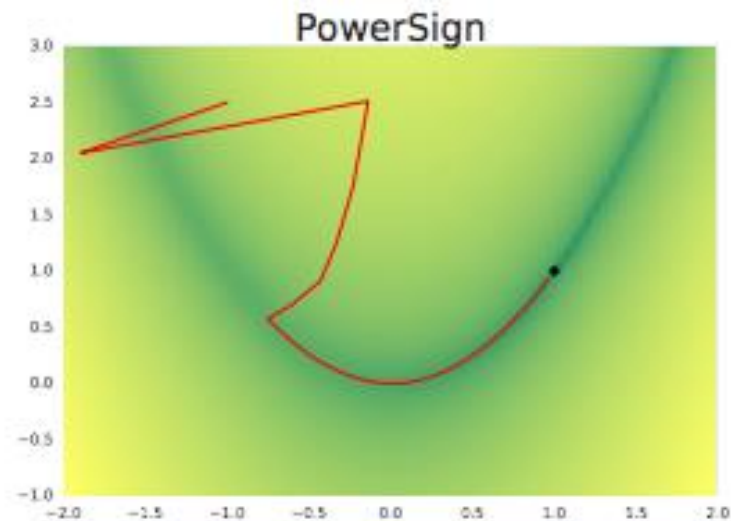
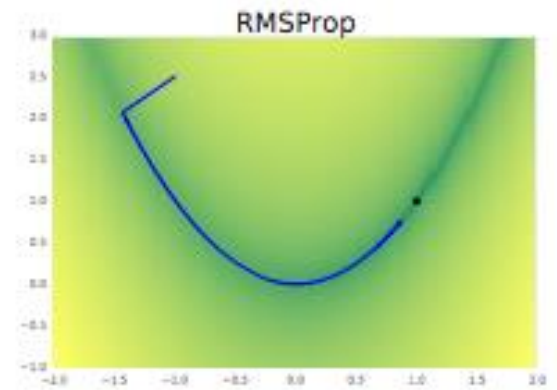
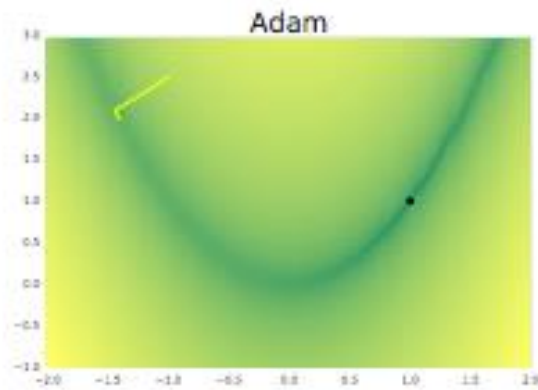
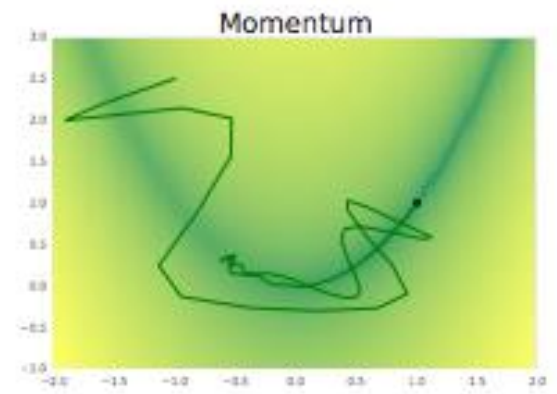
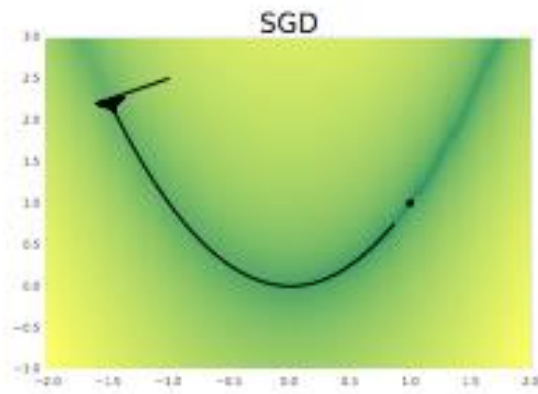
SWISH



Learning Rate

- **Operands:** g , g^2 , g^3 , \hat{m} , \hat{v} , $\hat{\gamma}$, $\text{sign}(g)$, $\text{sign}(\hat{m})$, 1, 2, $\epsilon \sim N(0, 0.01)$, $10^{-4}w$, $10^{-3}w$, $10^{-2}w$, $10^{-1}w$, Adam and RMSProp.
- **Unary functions** which map input x to: x , $-x$, e^x , $\log|x|$, $\sqrt{|x|}$, $\text{clip}(x, 10^{-5})$, $\text{clip}(x, 10^{-4})$, $\text{clip}(x, 10^{-3})$, $\text{drop}(x, 0.1)$, $\text{drop}(x, 0.3)$, $\text{drop}(x, 0.5)$ and $\text{sign}(x)$.
- **Binary functions** which map (x, y) to $x + y$ (addition), $x - y$ (subtraction), $x * y$ (multiplication), $\frac{x}{y + \delta}$ (division), x^y (exponentiation) or x (keep left).





$$e^{\text{sign}(g) * \text{sign}(m)} * g$$

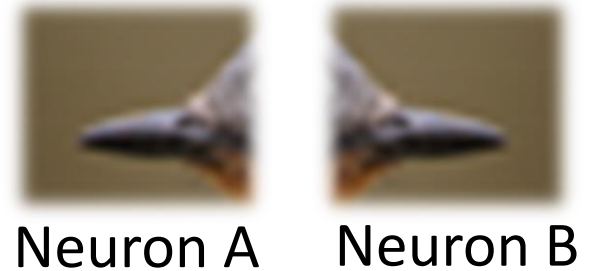
Can transfer to
new tasks

Capsule

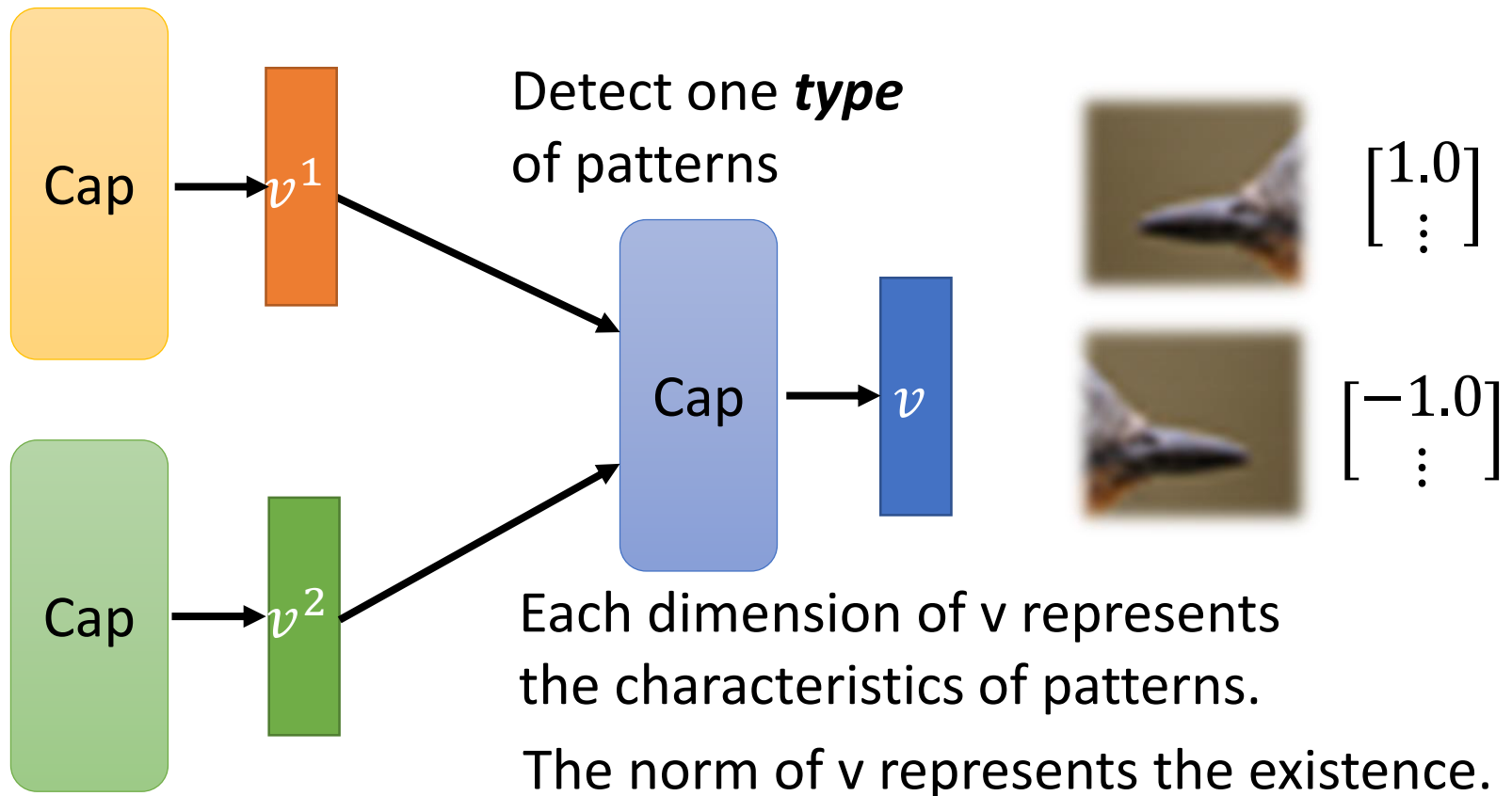
Sara Sabour, Nicholas Frosst, Geoffrey E. Hinton, “Dynamic Routing Between Capsules”, NIPS, 2017

Capsule

A neuron detects a specific pattern.



- Neuron: output a value, Capsule: output a vector



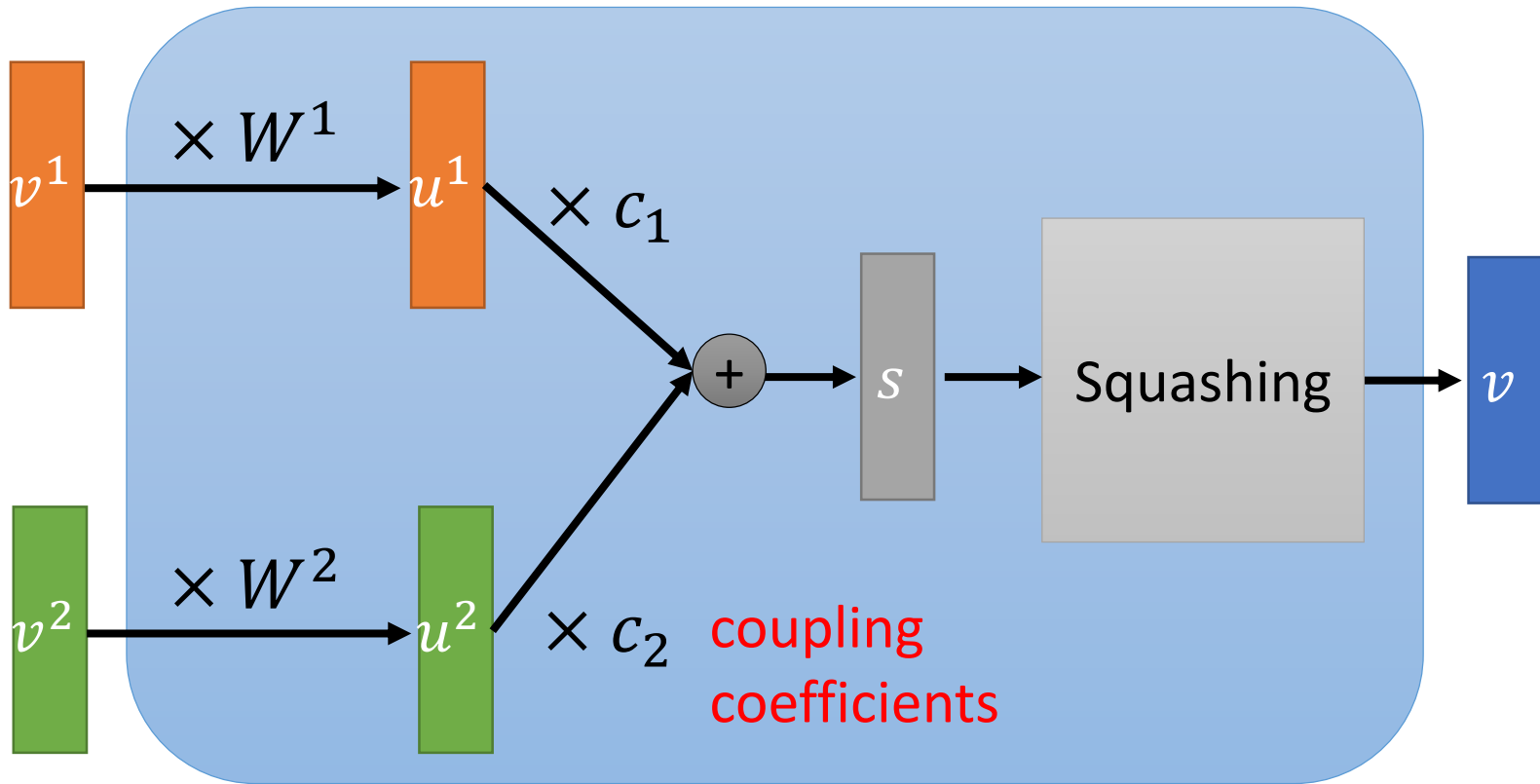
Capsule

$$u^1 = W^1 v^1 \quad u^2 = W^2 v^2$$

$$s = c_1 u^1 + c_2 u^2$$

$$v = \text{Squash}(s)$$

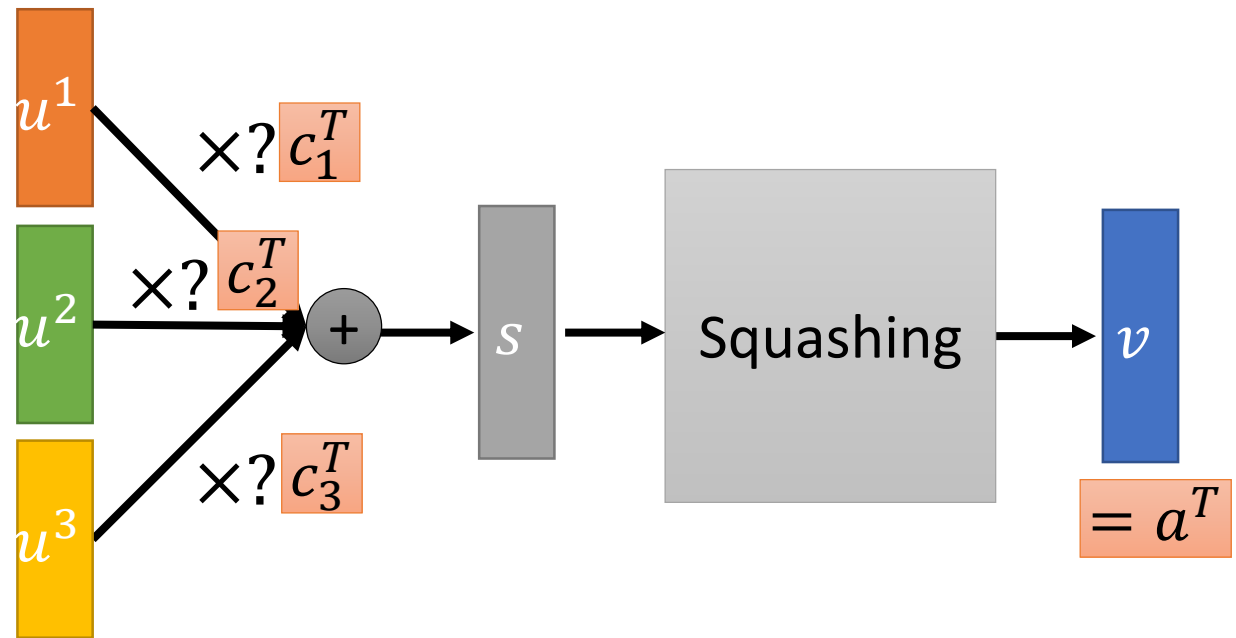
$$v = \frac{\|s\|^2}{1 + \|s\|^2} \frac{s}{\|s\|}$$



c are determined by dynamic routing during the testing stage.

c.f. pooling

Dynamic Routing



$$b_1^0 = 0, b_2^0 = 0, b_3^0 = 0$$

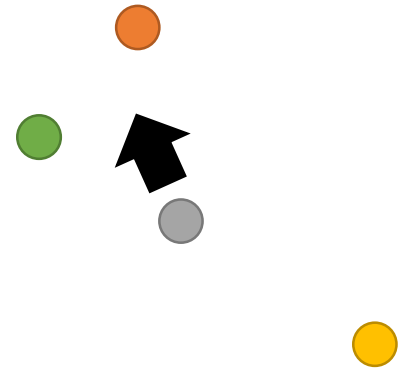
For $r = 1$ to T do

$$c_1^r, c_2^r, c_3^r = \text{softmax}(b_1^{r-1}, b_2^{r-1}, b_3^{r-1})$$

$$s^r = c_1 u^1 + c_2 u^2 + c_3 u^3$$

$$a^r = \text{Squash}(s^r)$$

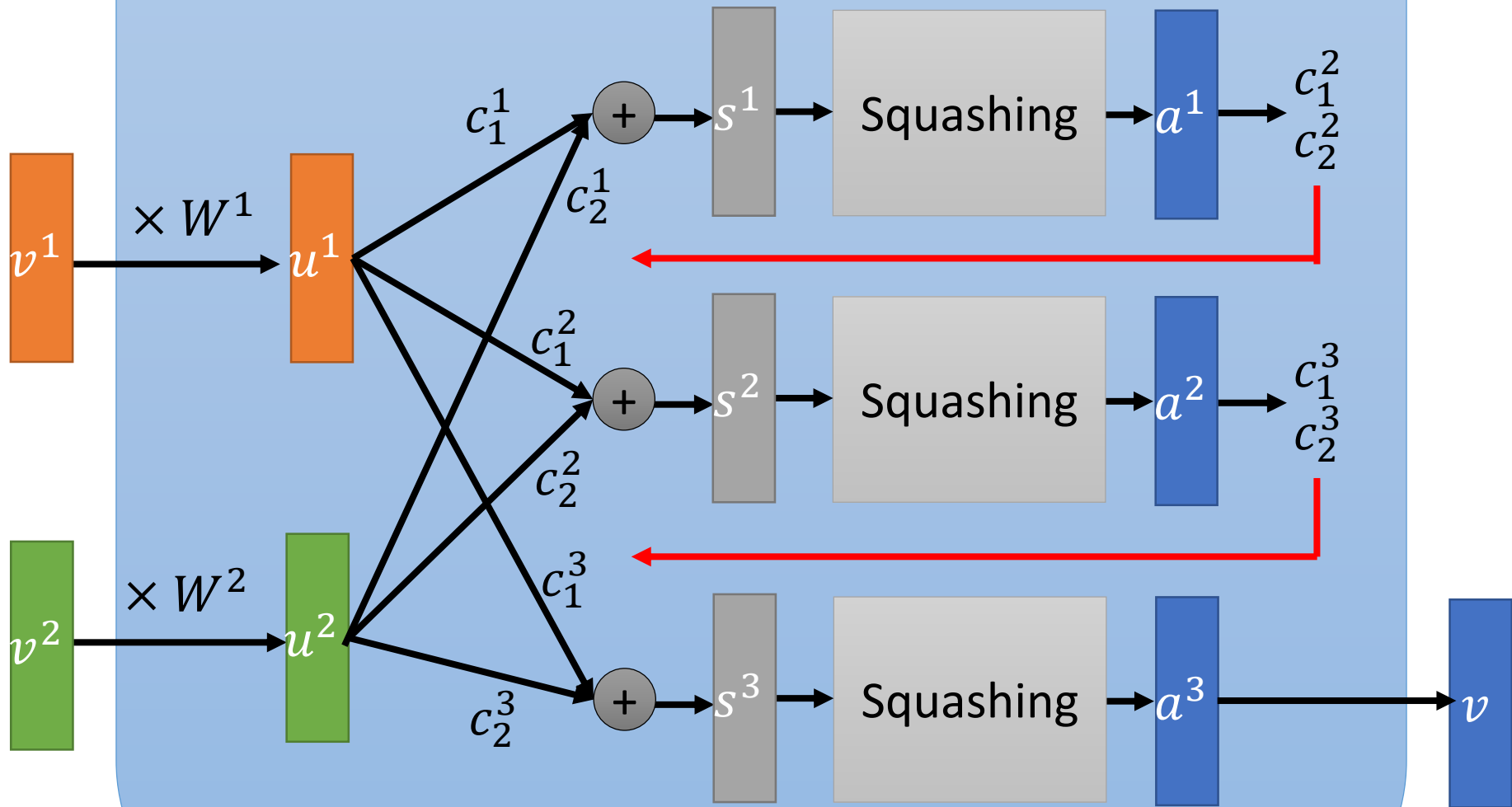
$$b_i^r = b_i^{r-1} + a^r \cdot u^i$$



$$b_1^0 = 0, b_2^0 = 0$$

$$c_1^1, c_2^1 = \text{softmax}(b_1^0, b_2^0)$$

$$b_i^r = b_i^{r-1} + a^r \cdot u^i$$



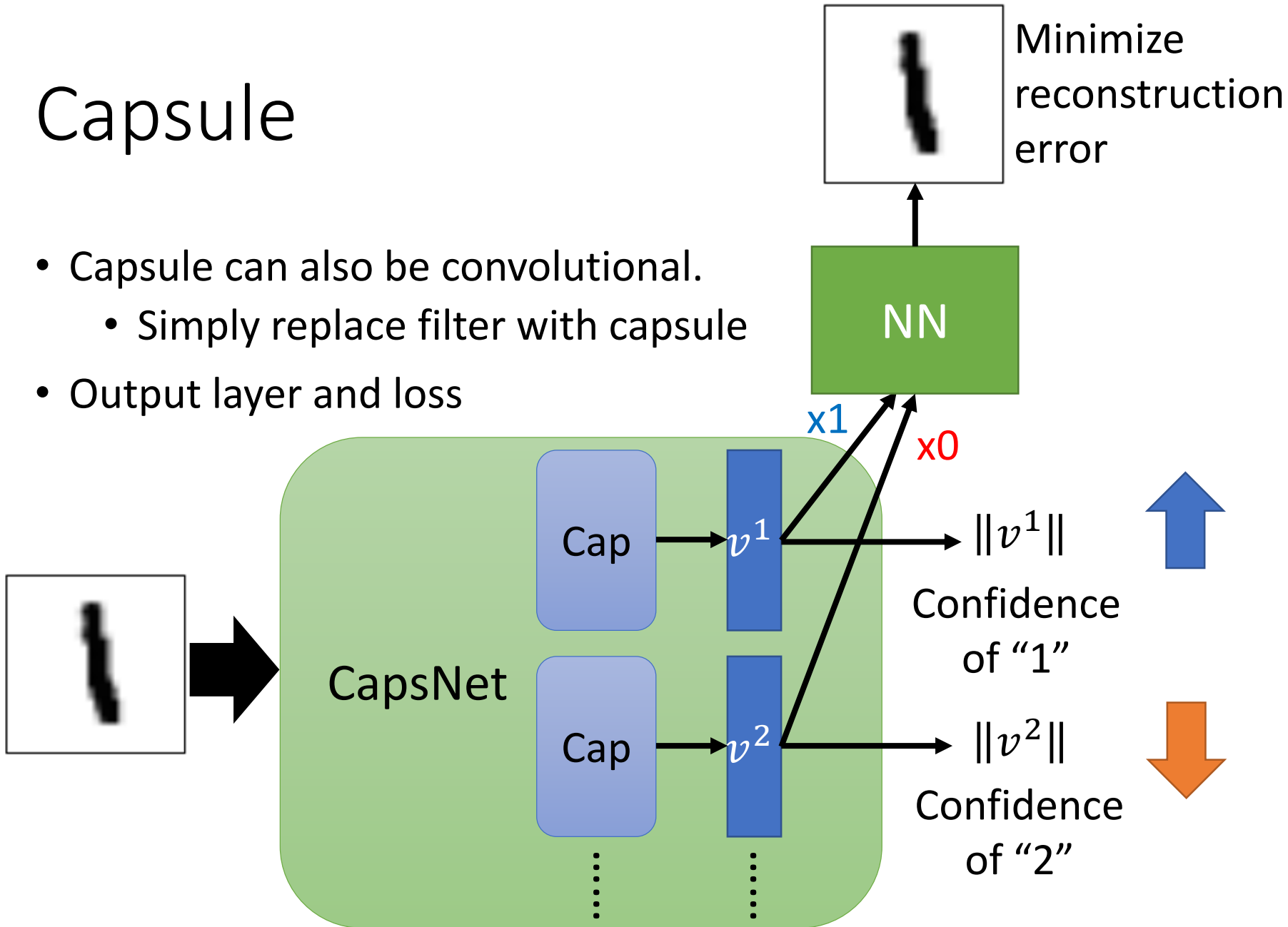
T=3

Like RNN

Also learned by backprop

Capsule

- Capsule can also be convolutional.
 - Simply replace filter with capsule
- Output layer and loss



Experimental Results

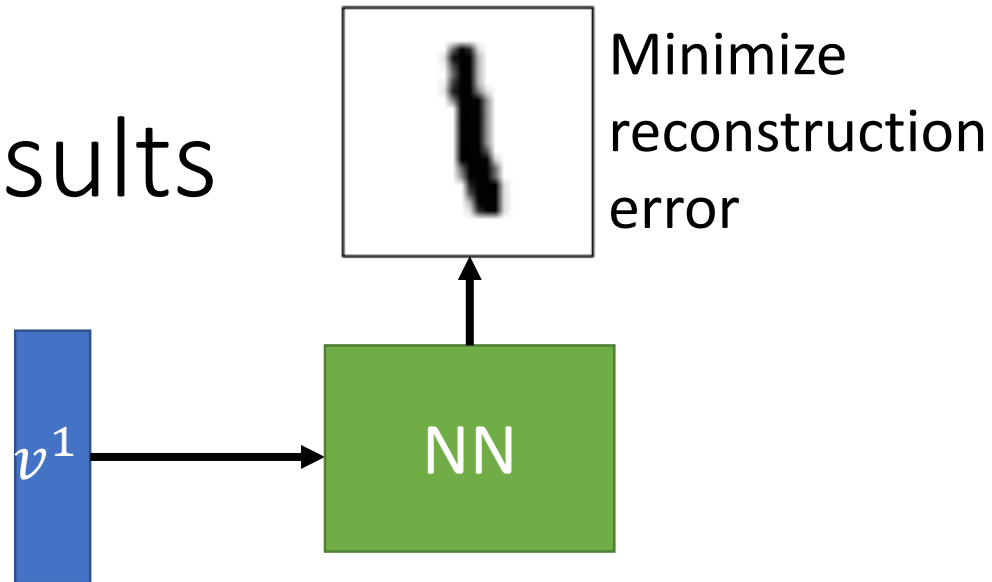
- MNIST

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

- Each example is an MNIST digit with a random small **affine transformation**.
- However, models were **never trained with affine transformations**
- **CapsNet** achieved **79%** accuracy on the affnist test set.
- A **traditional convolutional model** with a similar number of parameters which achieved **66%**.

Experimental Results

- Each dimension contains specific information.



Scale and thickness	
Localized part	
Stroke thickness	
Localized skew	
Width and translation	
Localized part	

Experimental Results

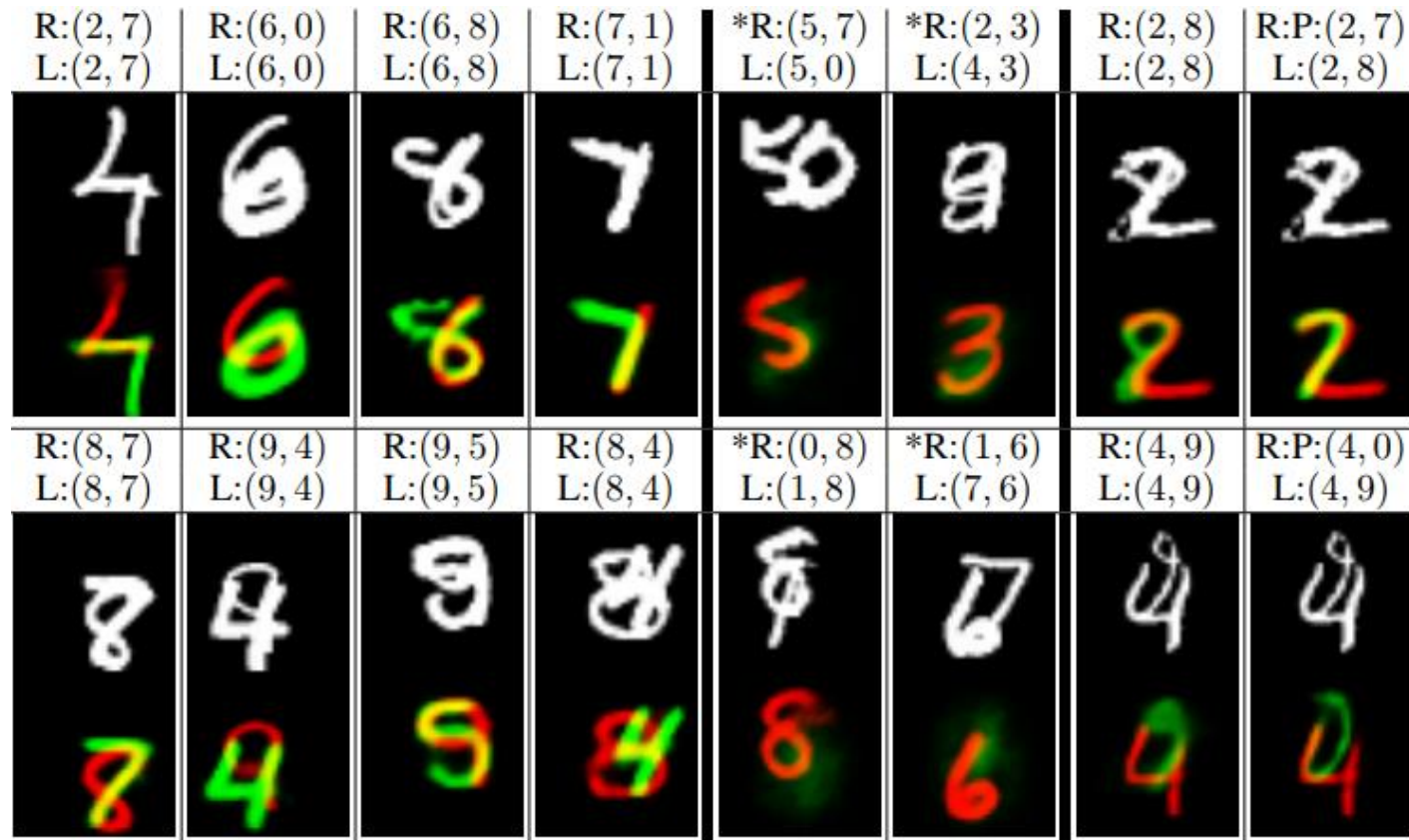
Top: input

Bottom: reconstructed

R: reconstructed digits

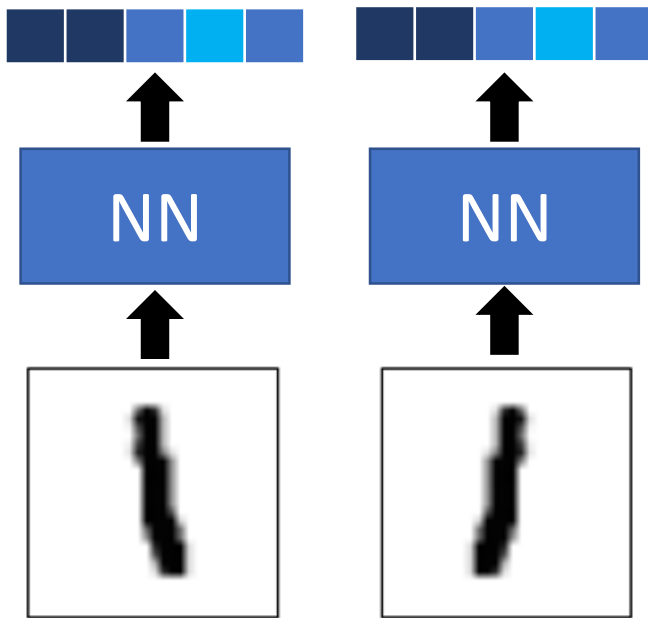
L: true labels

- MultiMNIST

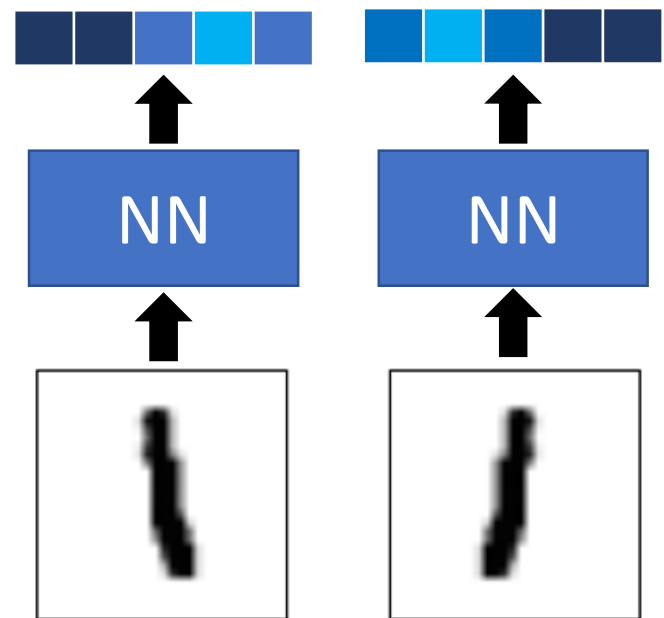


Discussion

- Invariance v.s. Equivariance



Invariance



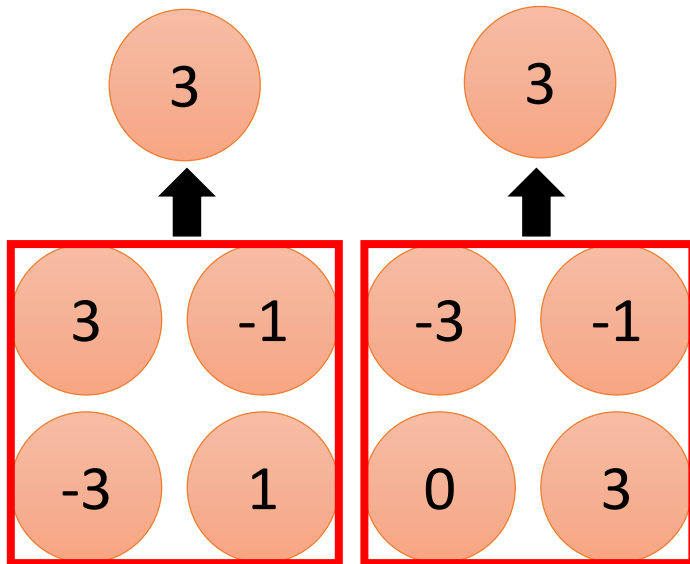
Equivariance

I know the difference, but I do not react to it.

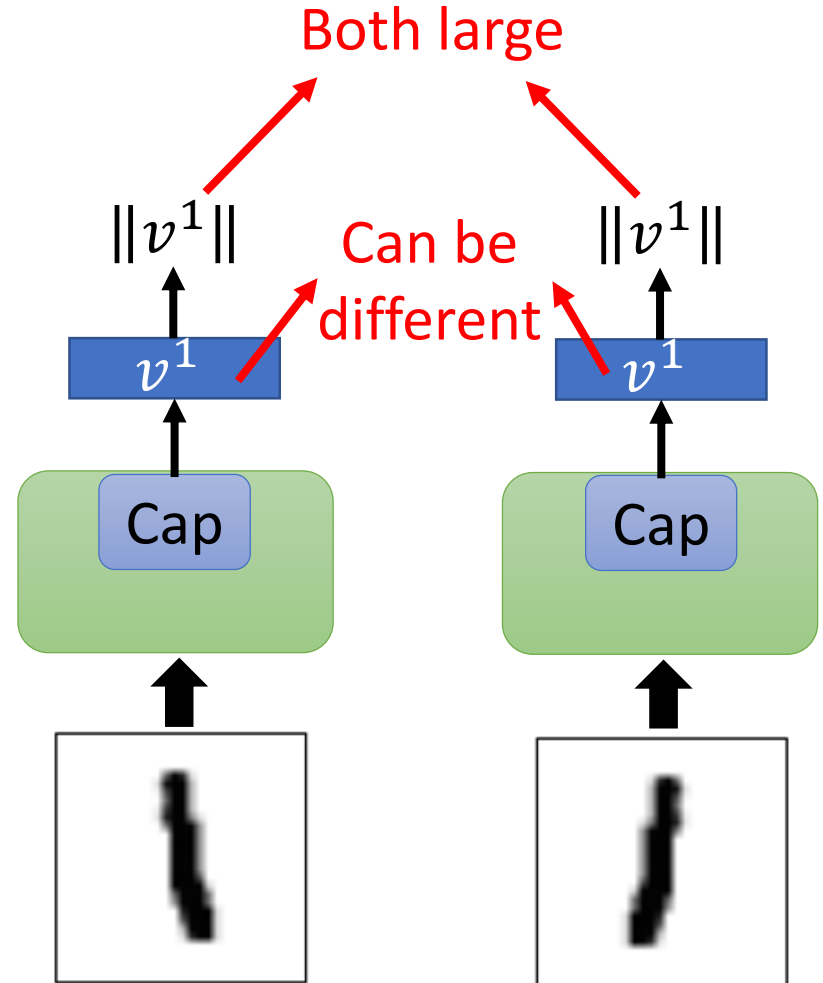
Discussion

- Invariance v.s. Equivariance

I don't know the difference.

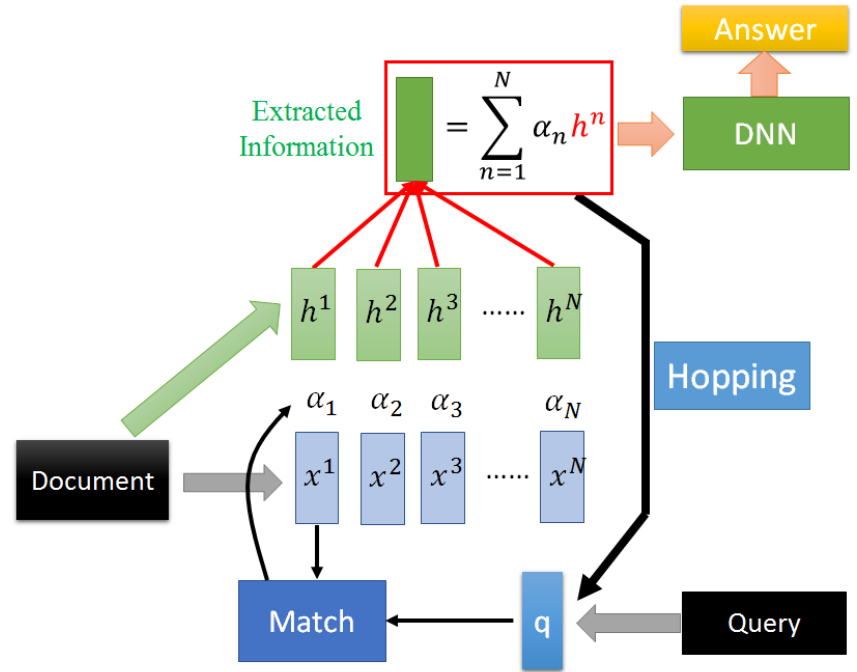
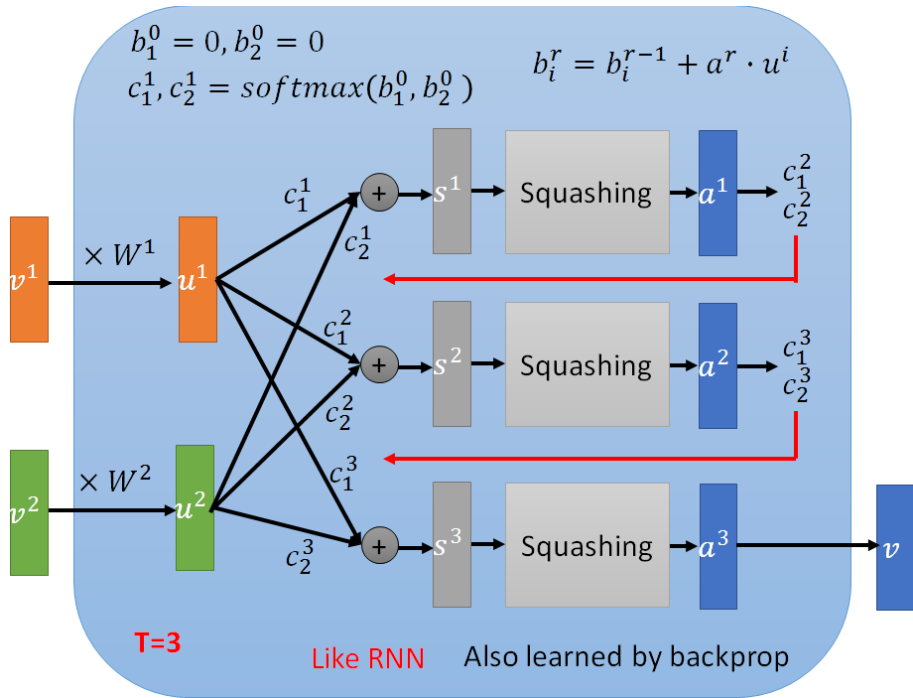


Max pooling has invariance, but not equivariance.



Capsule has both invariance and equivariance.

Dynamic Routing



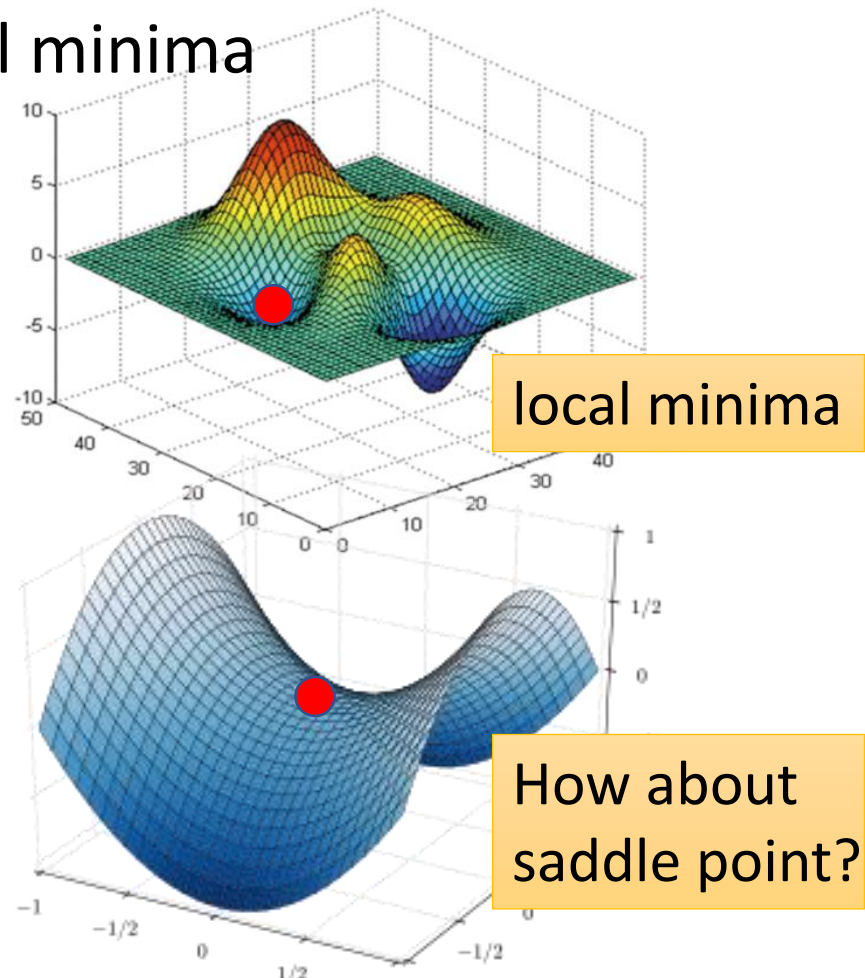
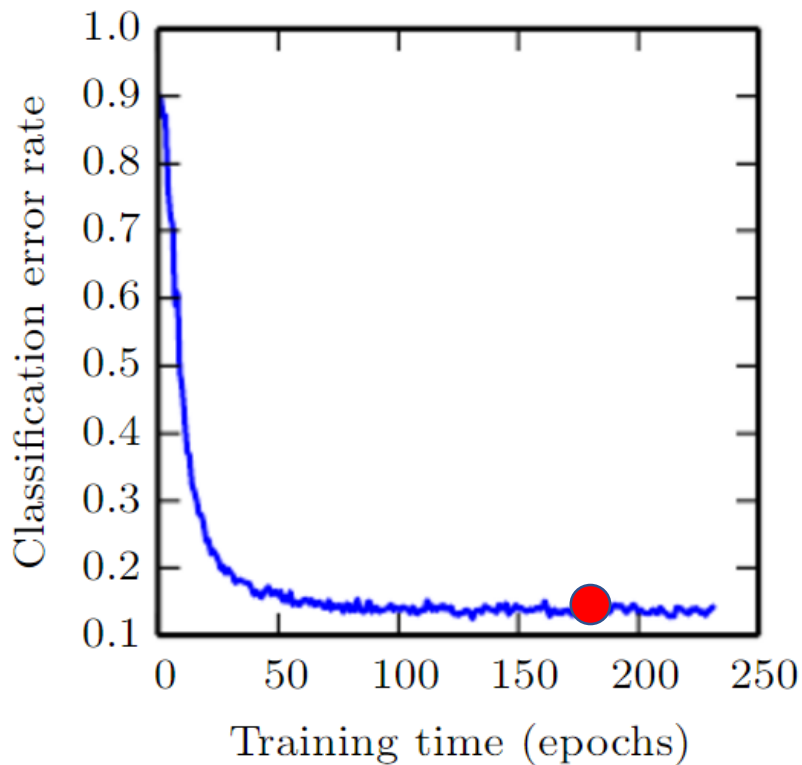
To Learn More

- Hinton's talk:
<https://www.youtube.com/watch?v=rTawFwUvnLE>
- Keras:
 - <https://github.com/XifengGuo/CapsNet-Keras>
- Tensorflow:
 - <https://github.com/naturomics/CapsNet-Tensorflow>
- PyTorch
 - <https://github.com/gram-ai/capsule-networks>
 - <https://github.com/timomernick/pytorch-capsule>
 - <https://github.com/nishnik/CapsNet-PyTorch>

Interesting Facts (?) about Deep Learning

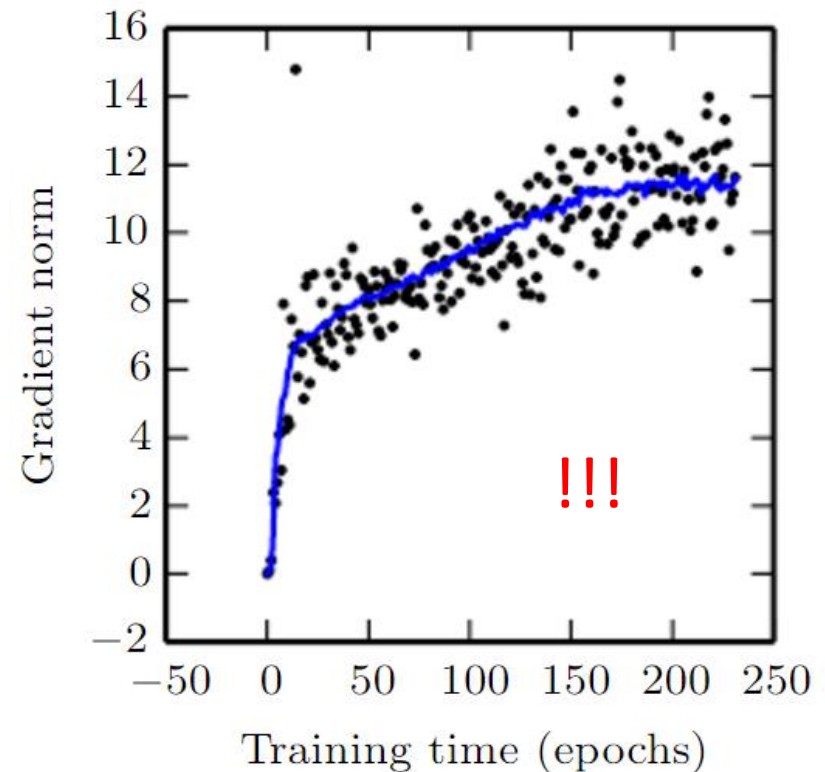
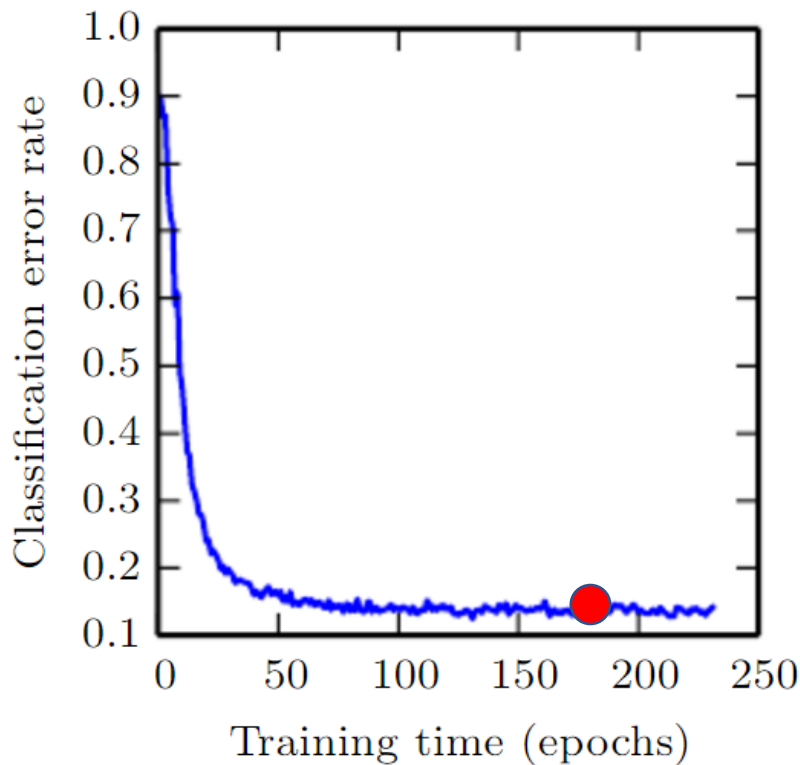
Training stuck because ?

- People believe training stuck because the parameters are near a local minima

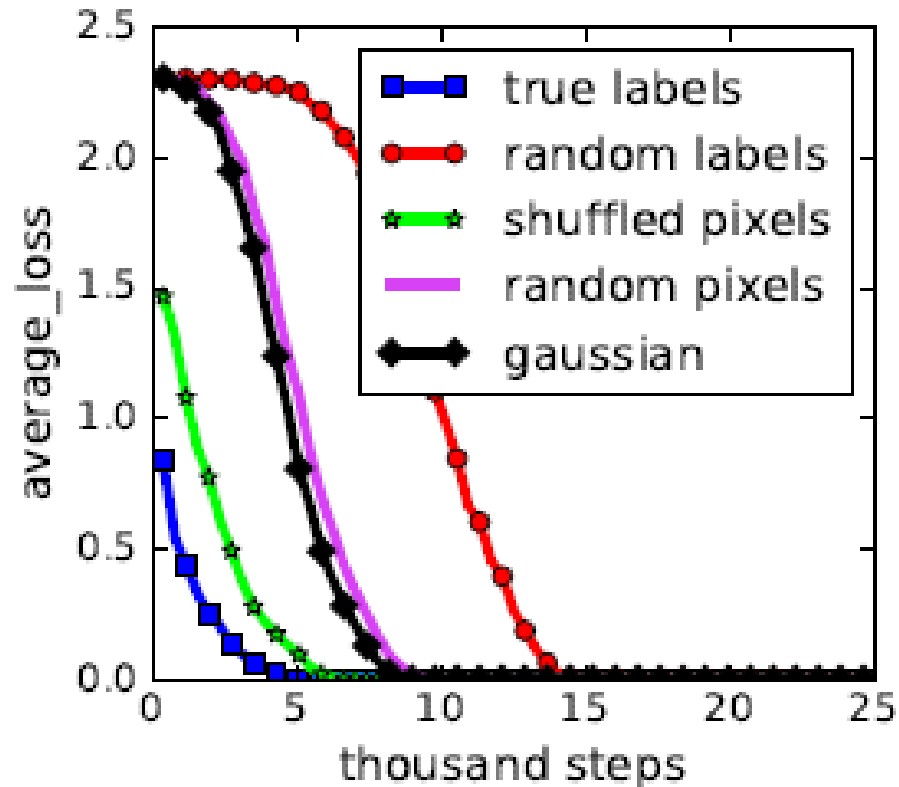


Training stuck because ?

- People believe training stuck because the parameters are around a critical point



Brute-force Memorization ?



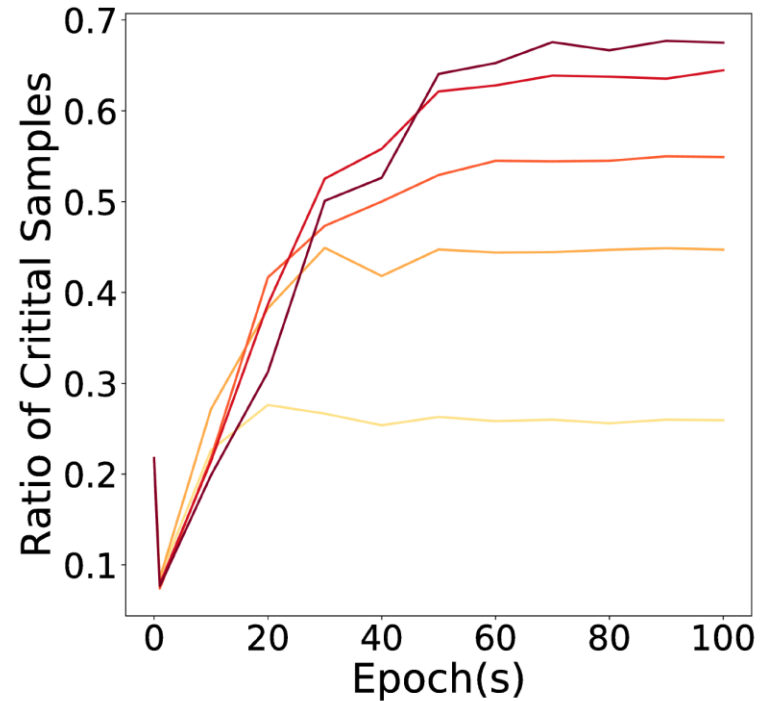
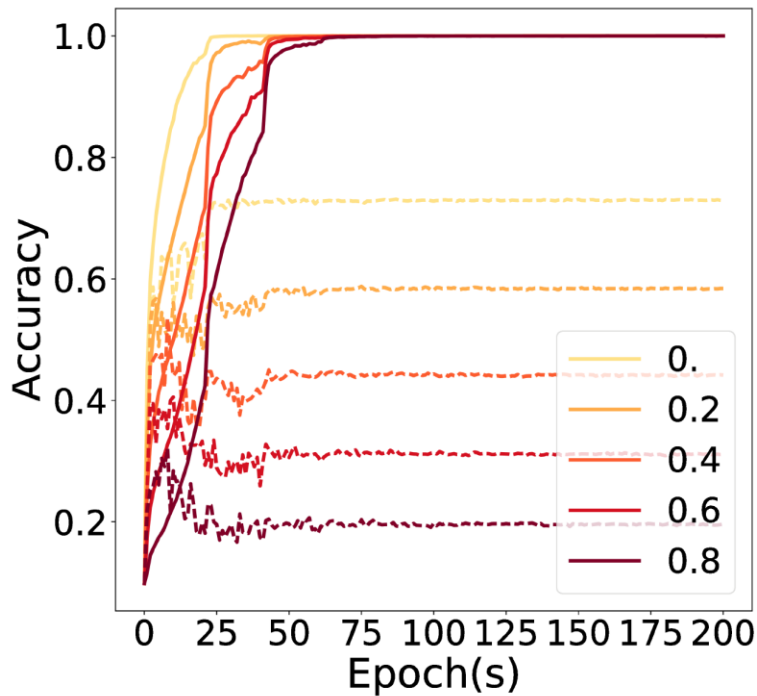
Final of 2017 Spring:

https://ntumlds.wordpress.com/2017/03/27/r05922018_drliao/

Demo

Brute-force Memorization ?

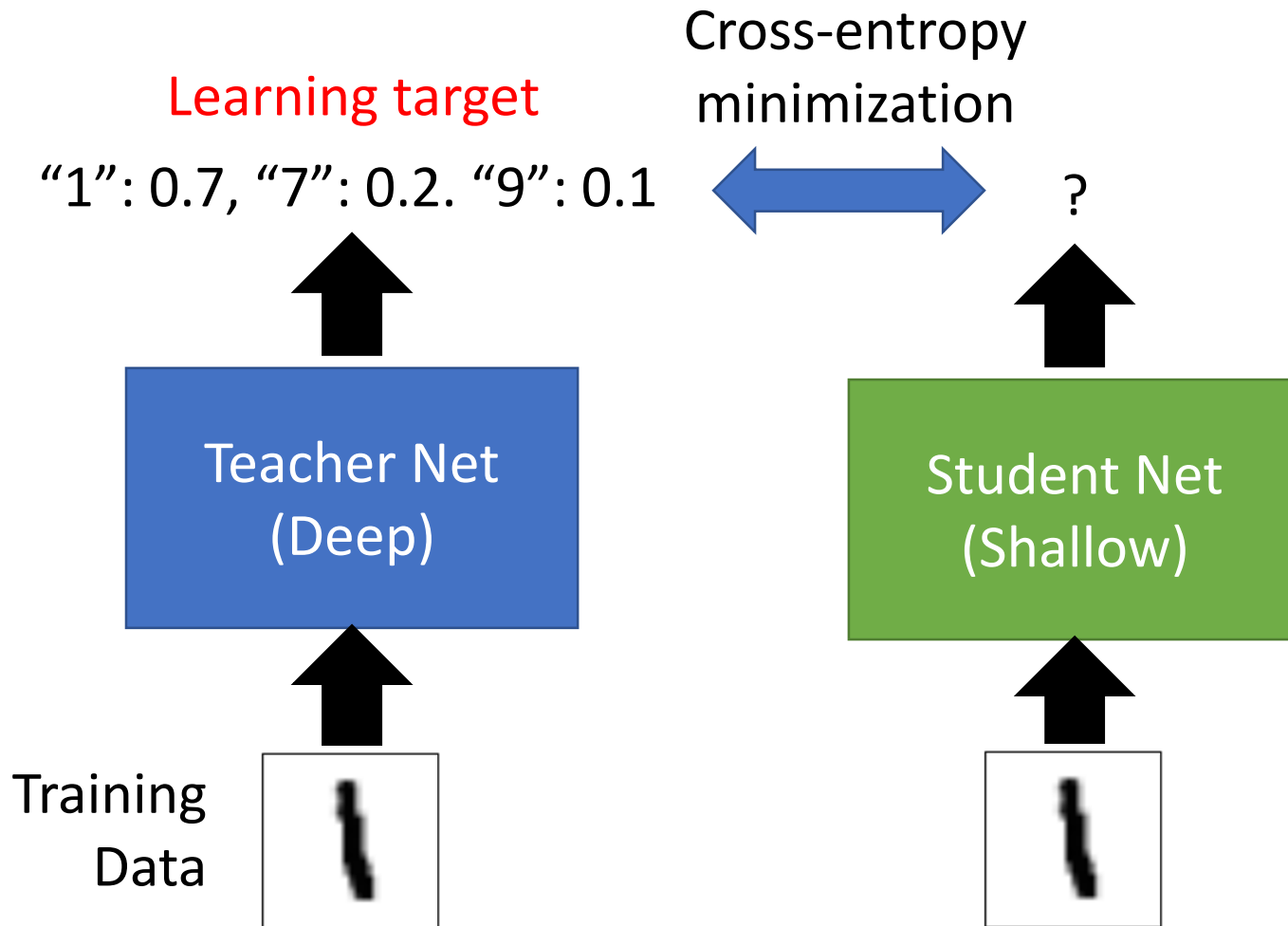
- Simple pattern first, then memorize exception

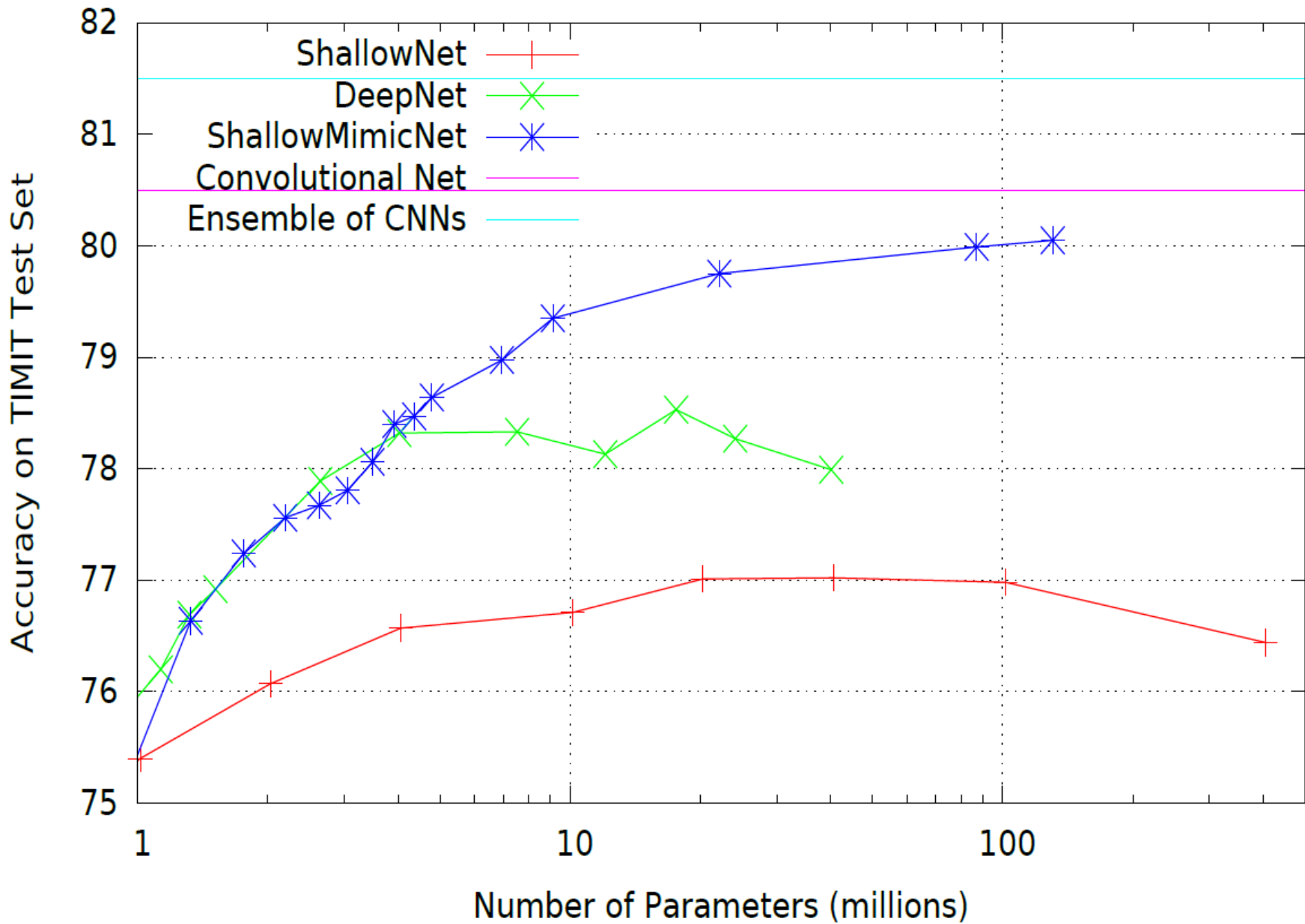


(b) Noise added on classification labels.

Knowledge Distillation

Knowledge Distillation
<https://arxiv.org/pdf/1503.02531.pdf>
Do Deep Nets Really Need to be Deep?
<https://arxiv.org/pdf/1312.6184.pdf>





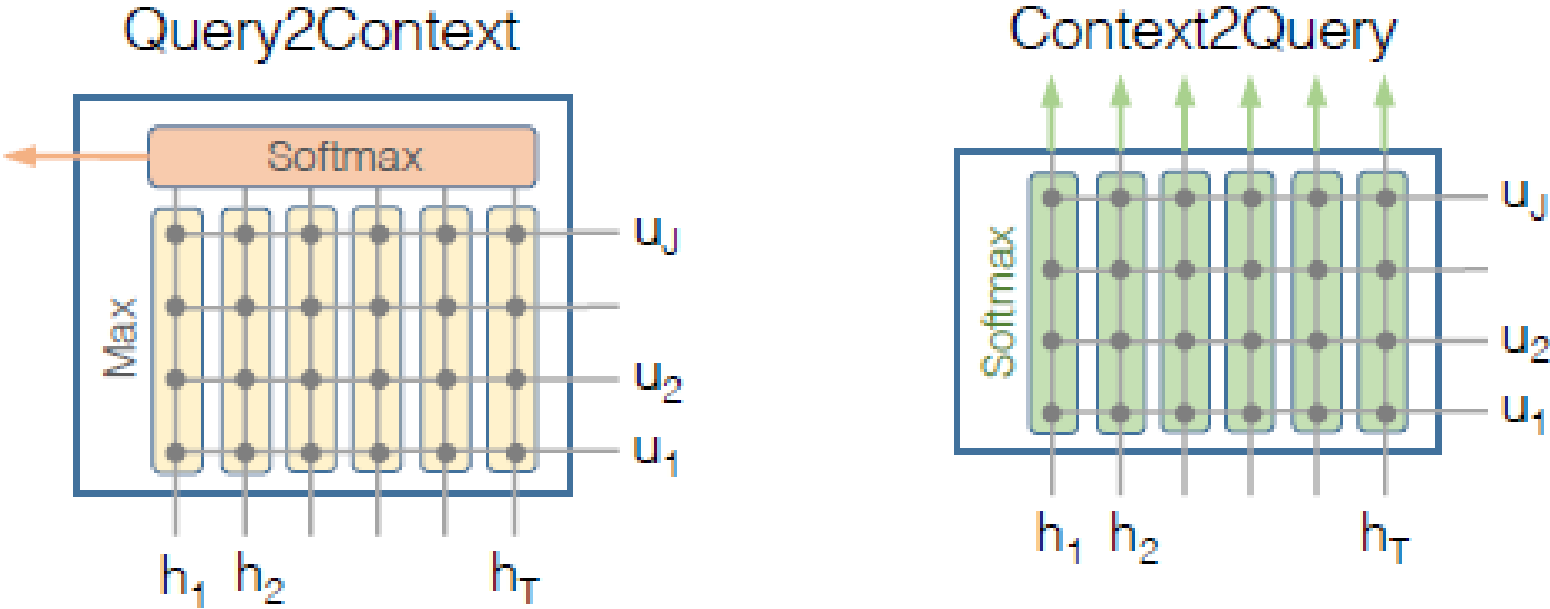
Deep Learning for Question Answering

Question Answering

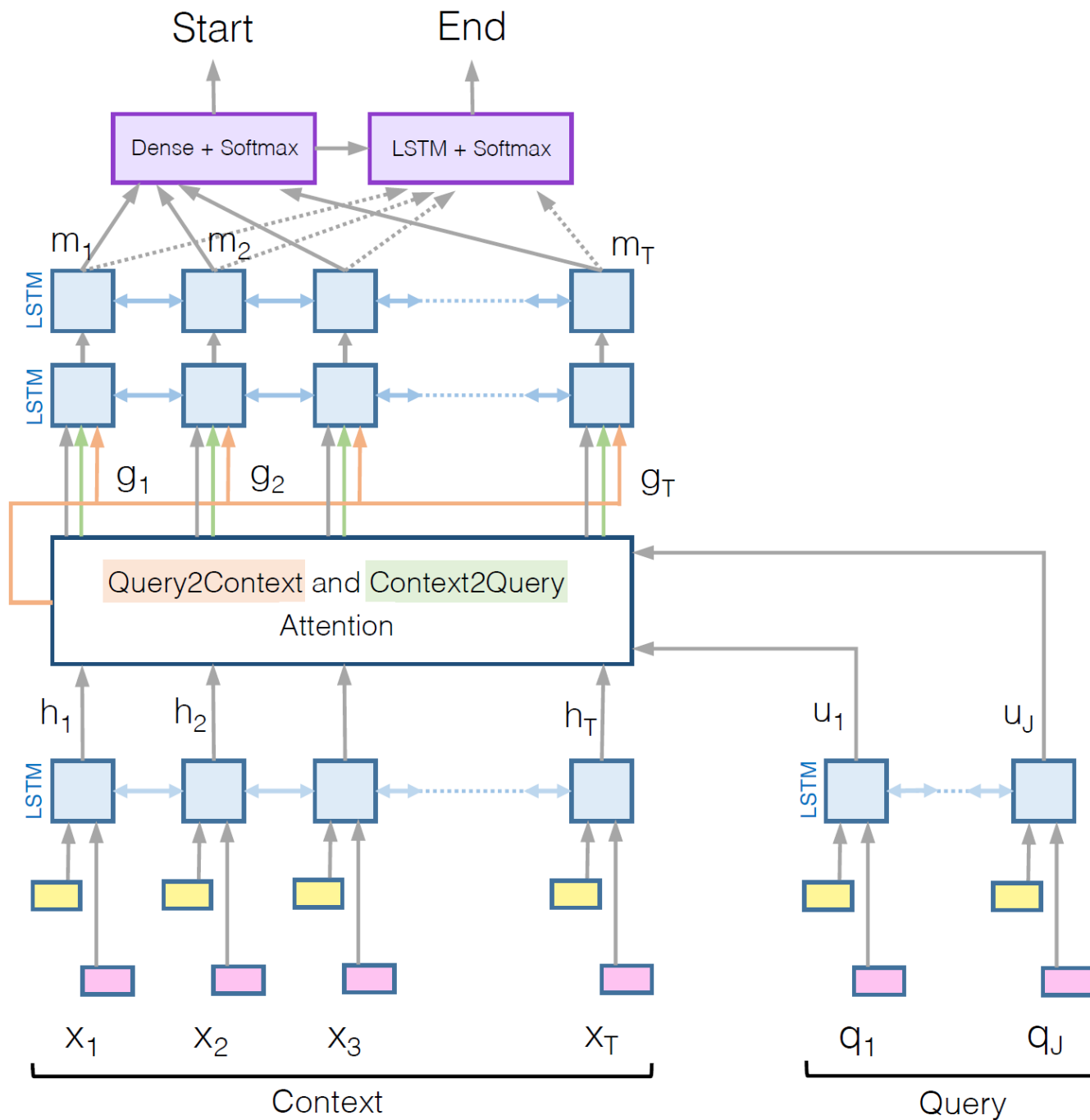
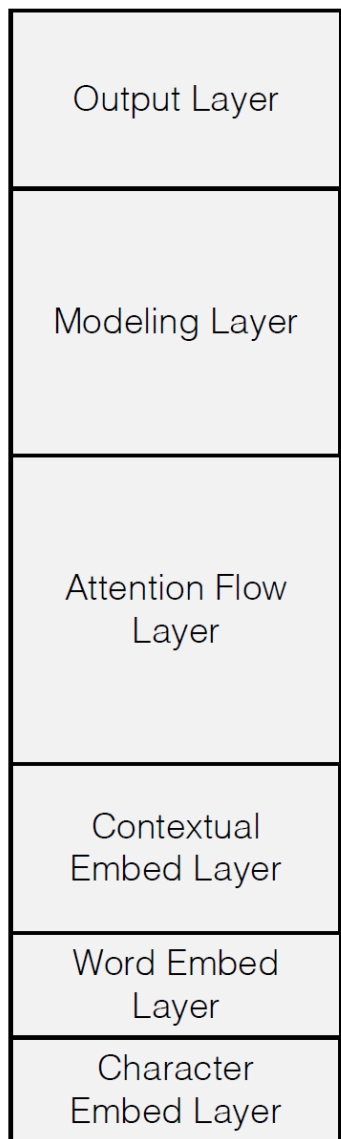
- Given a document and a query, output an answer
- bAbl: the answer is a word
 - <https://research.fb.com/downloads/babi/>
- SQuAD: the answer is a sequence of words (in the input document)
 - <https://rajpurkar.github.io/SQuAD-explorer/>
- MS MARCO: the answer is a sequence of words
 - <http://www.msmarco.org>
- MovieQA: Multiple choice question (output a number)
 - <http://movieqa.cs.toronto.edu/home/>
- More: <https://github.com/dapurv5/awesome-question-answering>

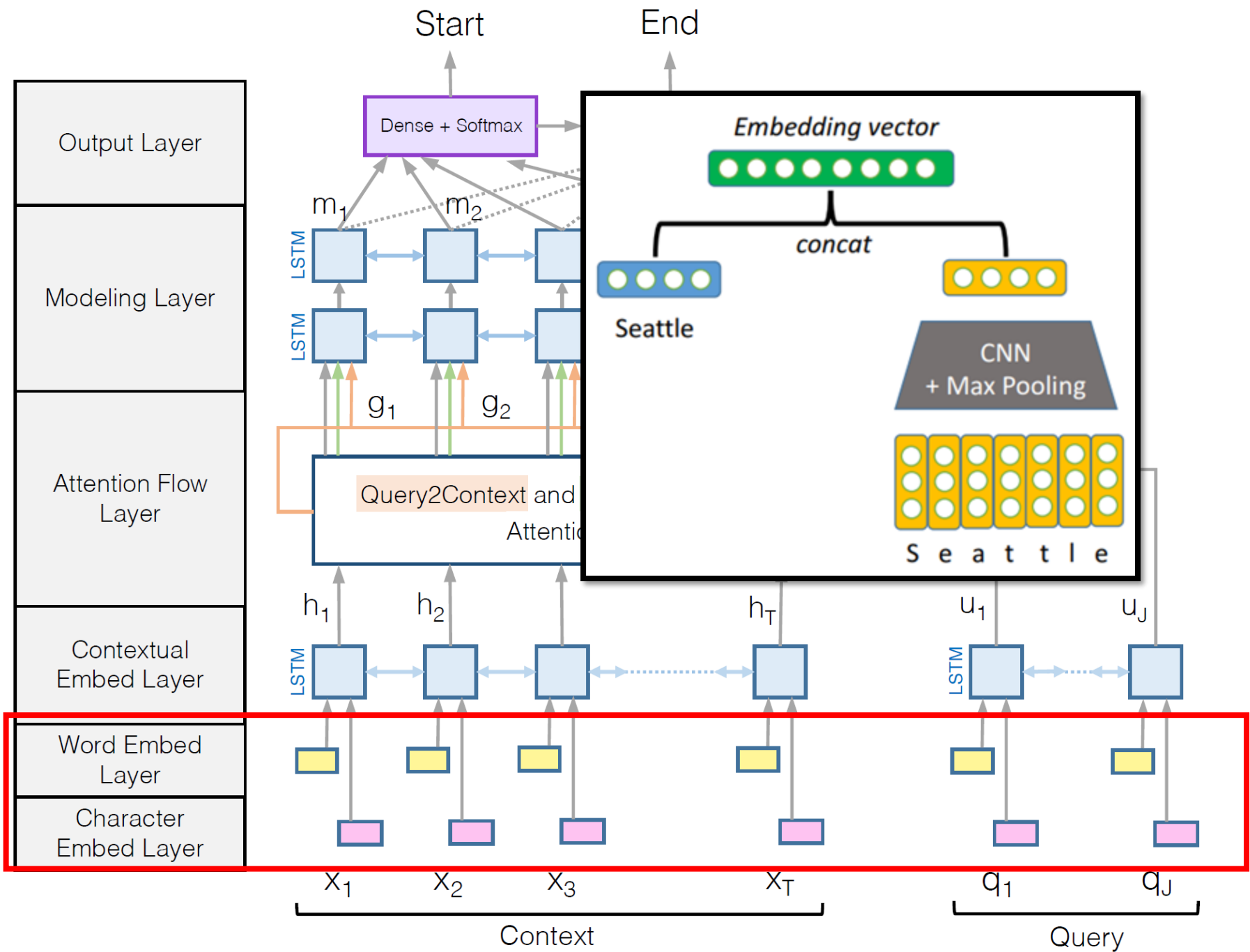
Rank	Model	EM	F1
1 Oct 17, 2017	Interactive AoA Reader+ (ensemble) <i>Joint Laboratory of HIT and iFLYTEK</i>	79.083	86.450
2 Oct 24, 2017	FusionNet (ensemble) <i>Microsoft Business AI Solutions Team</i>	78.978	86.016
3 Nov 03, 2017	BiDAF + Self Attention + ELMo (single model) <i>Allen Institute for Artificial Intelligence</i>	78.580	85.833
3 Oct 12, 2017	r-net (ensemble) <i>Microsoft Research Asia</i> http://aka.ms/rnet	78.926	85.722
3 Oct 22, 2017	DCN+ (ensemble) <i>Salesforce Research</i>	78.852	85.996
4 Oct 22, 2017	BiDAF + Self Attention + ELMo (single model) <i>Allen Institute for Artificial Intelligence</i>	77.856	85.344
5 Jul 25, 2017	Interactive AoA Reader (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	77.845	85.297
6 Aug 21, 2017	Reinforced Mnemonic Reader (ensemble) <i>NUDT and Fudan University</i> https://arxiv.org/abs/1705.02798	77.678	84.888

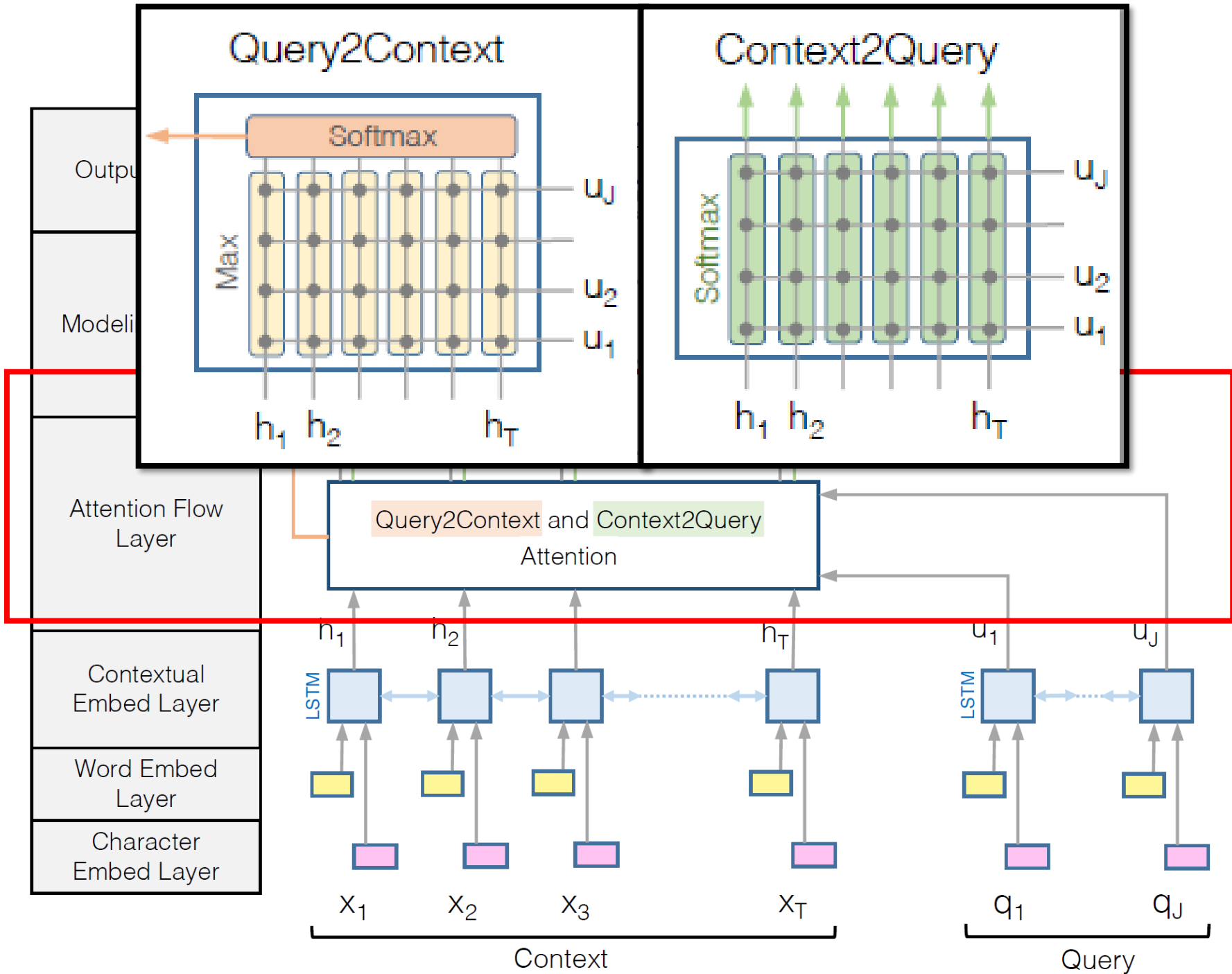
Bi-directional Attention Flow

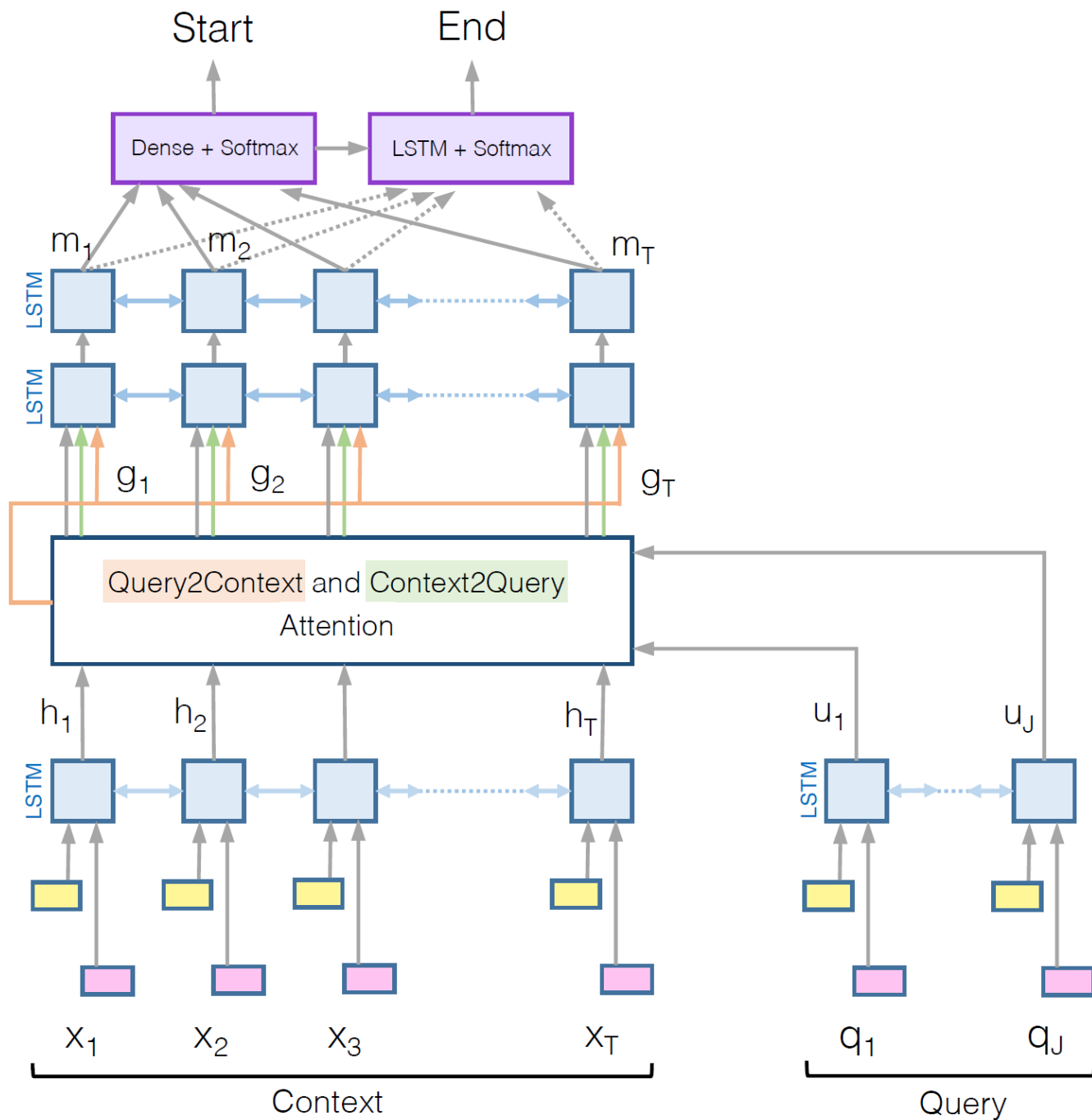
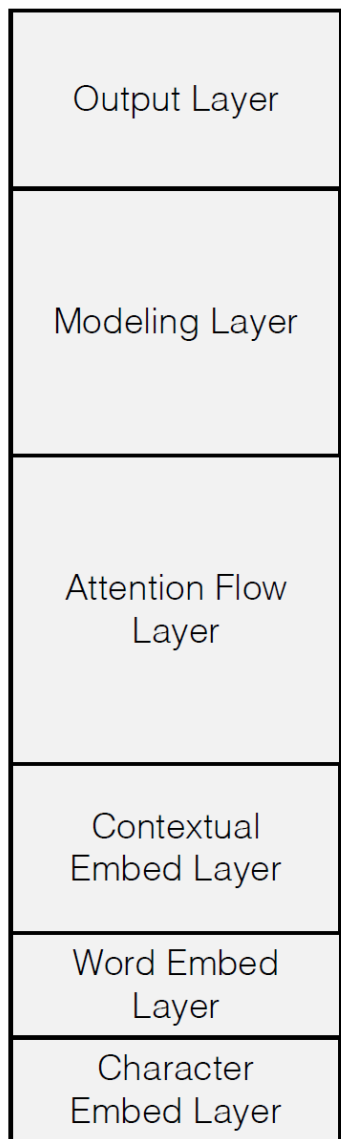


3	BiDAF + Self Attention + ELMo (single model)	78.580	85.833
Nov 03, 2017	Allen Institute for Artificial Intelligence		









Dynamic Coattention Networks

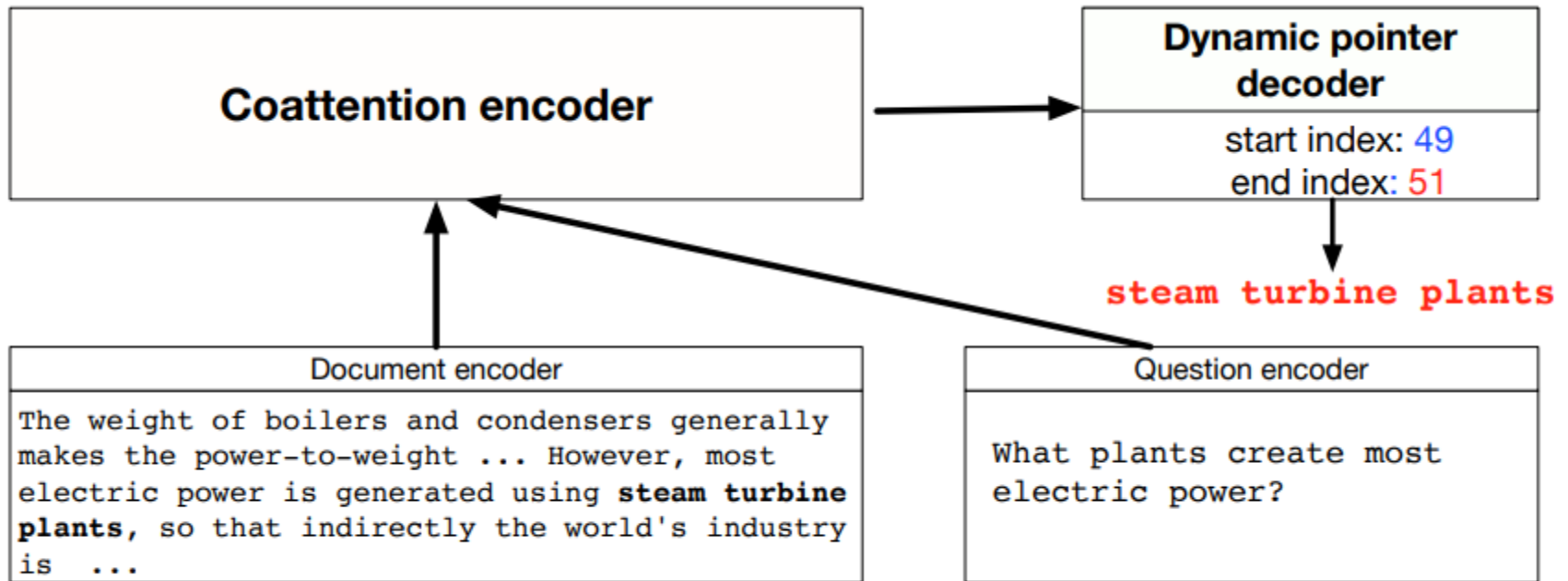
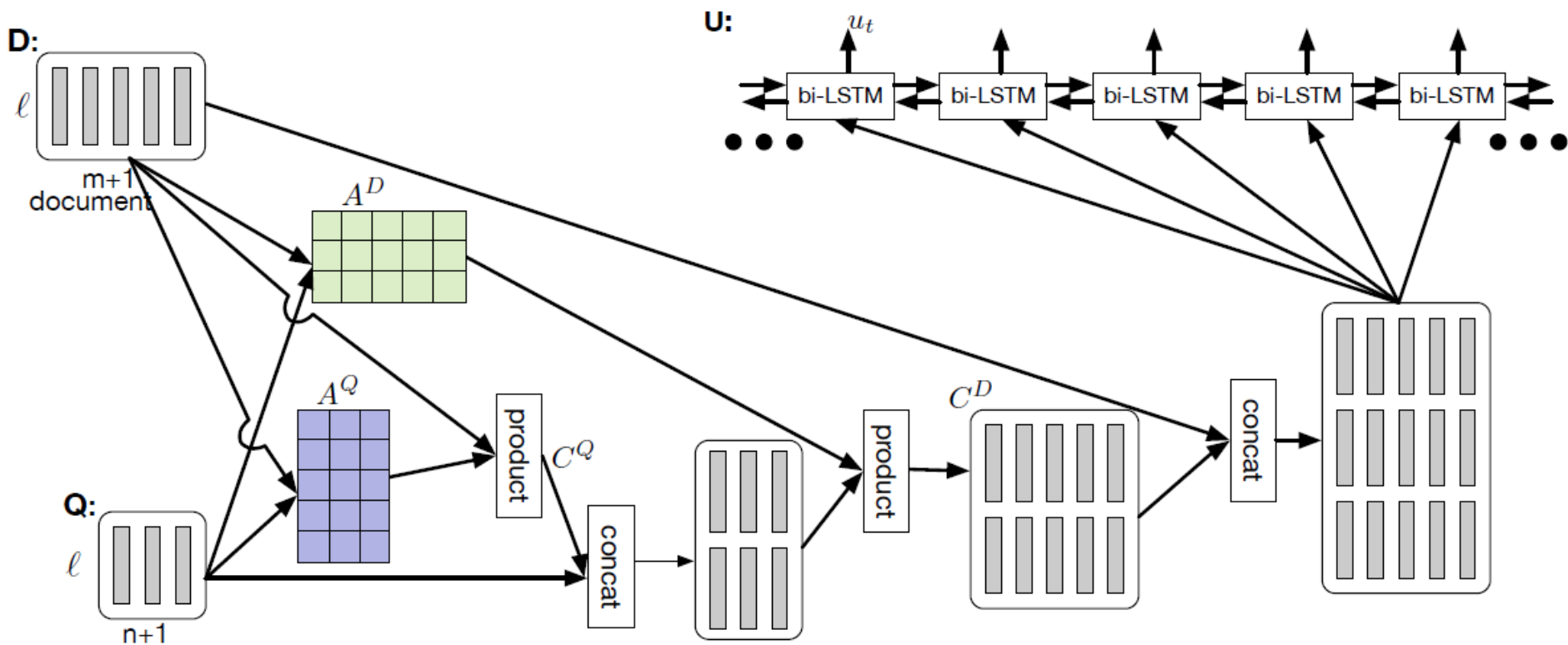


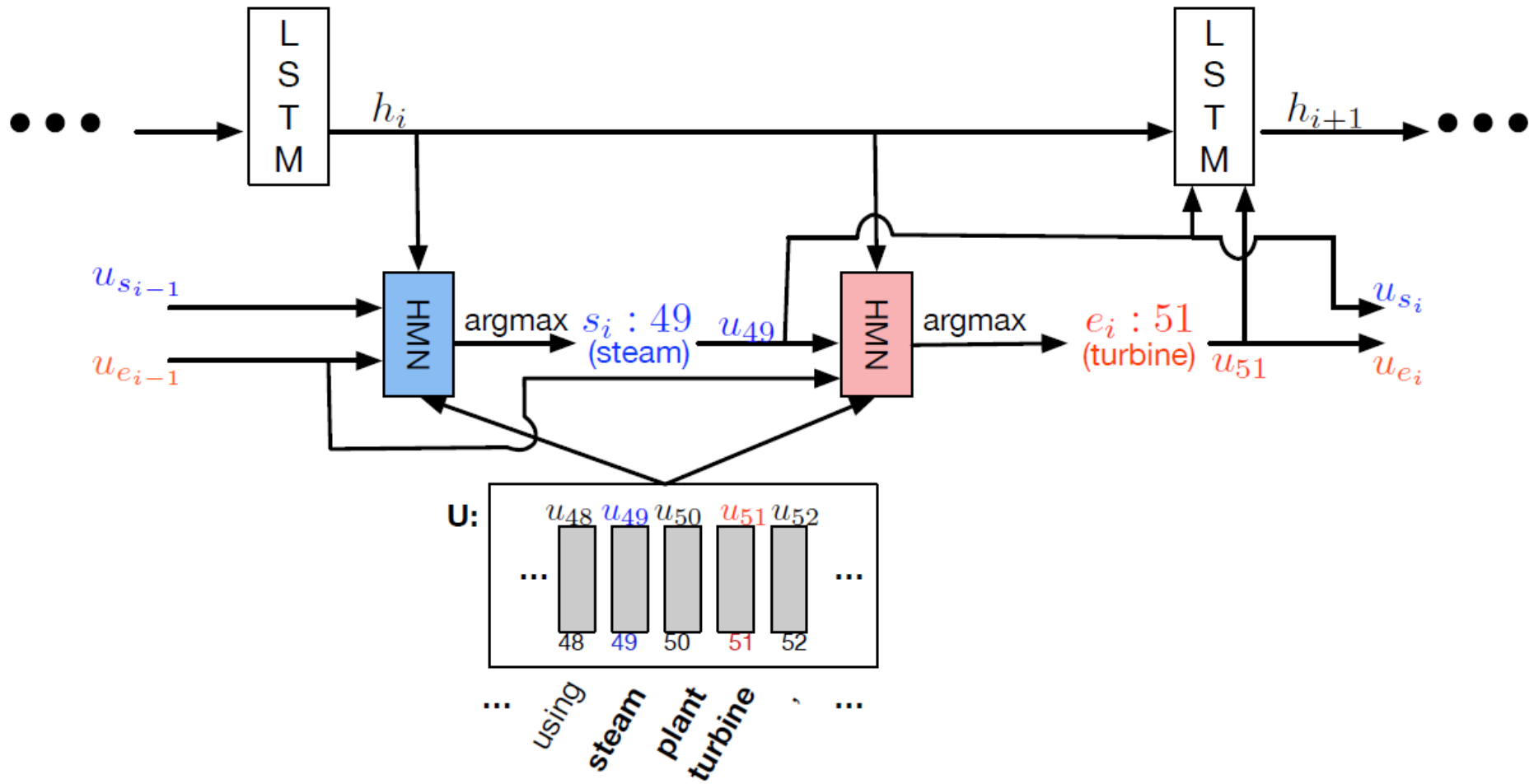
Figure 1: Overview of the Dynamic Coattention Network.

3	DCN+ (ensemble)	78.852	85.996
Oct 22, 2017	Salesforce Research		

Dynamic Coattention Networks



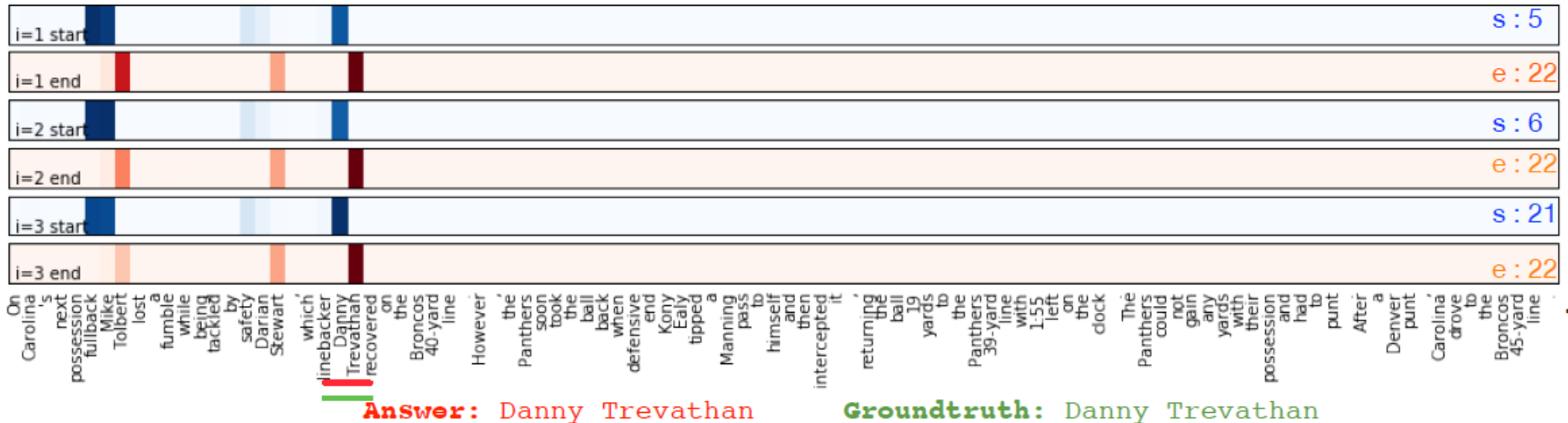
Dynamic Coattention Networks



Dynamic Coattention Networks

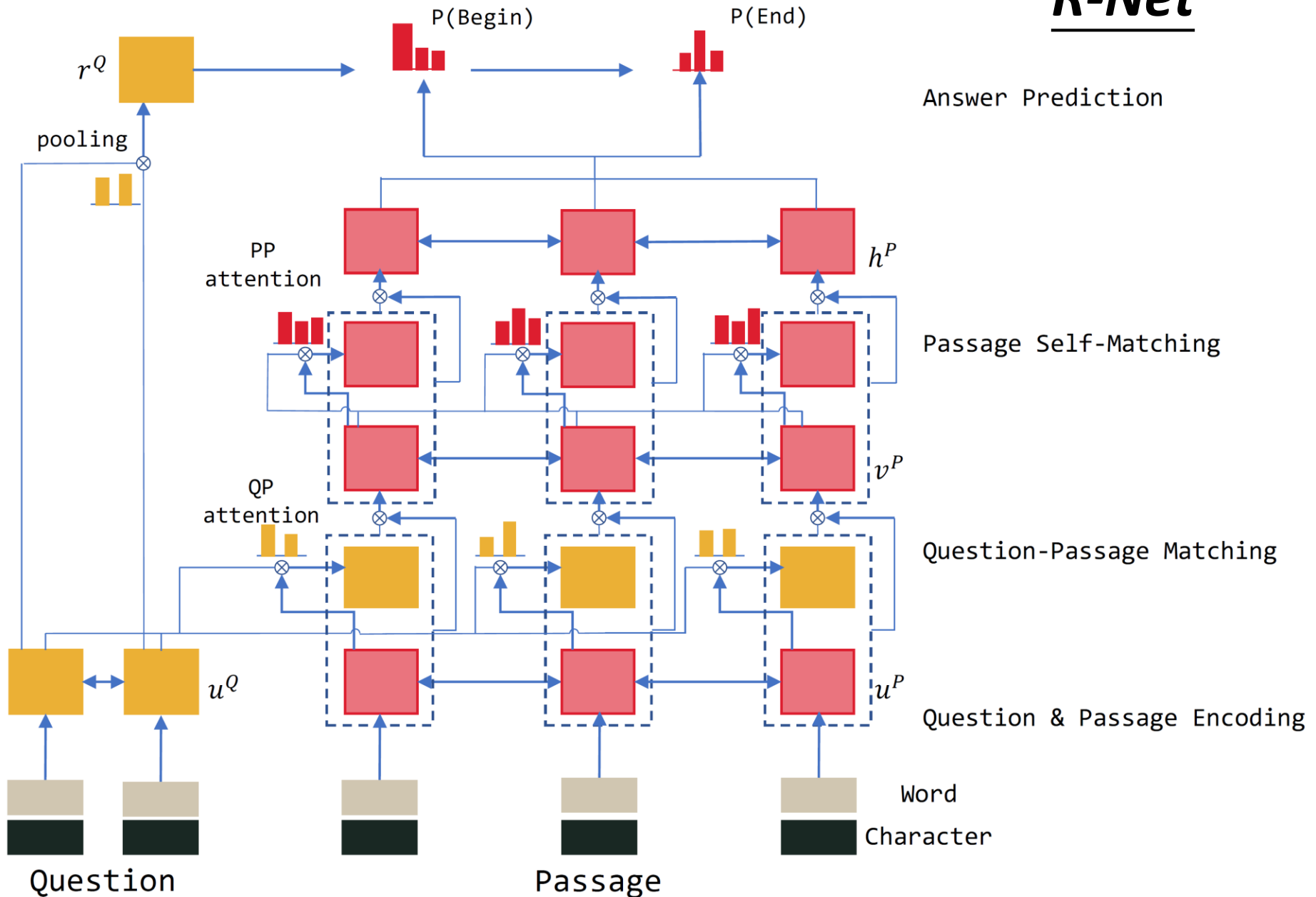
- Experimental Results

Question 1: Who recovered Tolbert's fumble?



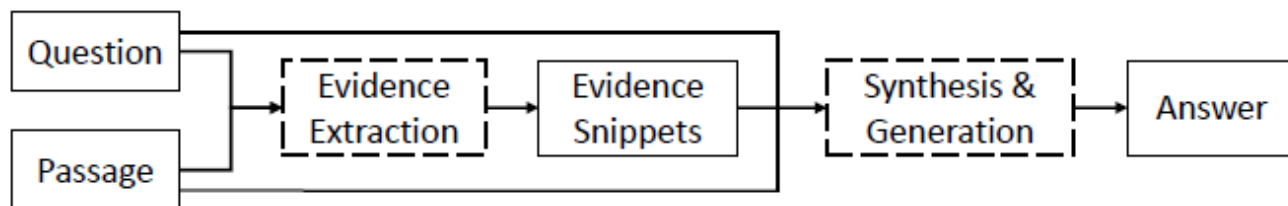
DCN+: <https://arxiv.org/pdf/1711.00106.pdf>

R-Net

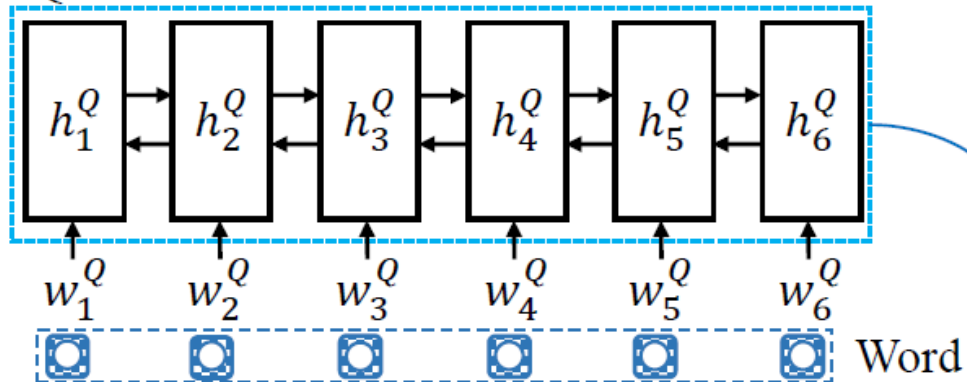


S-net

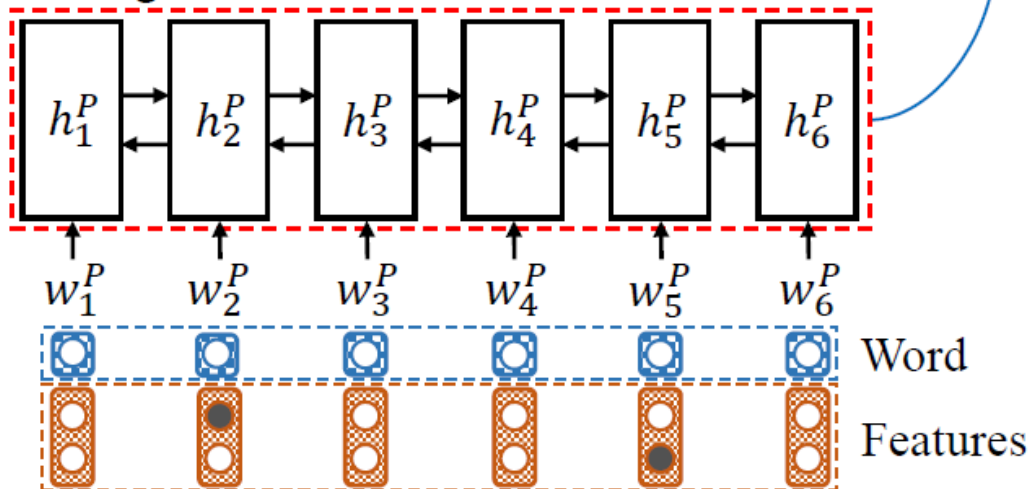
MS MARCO



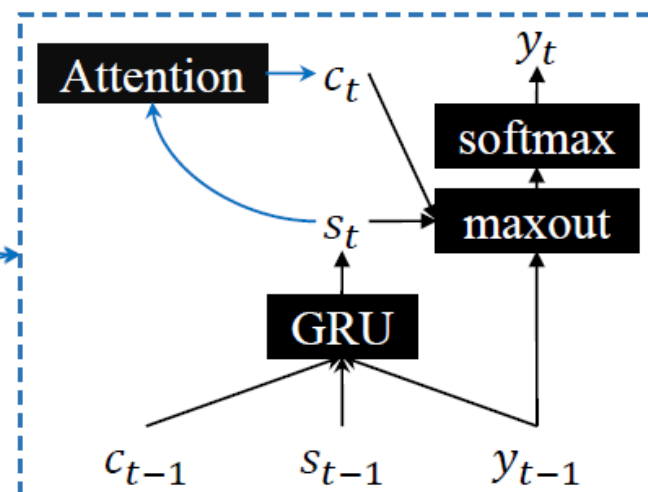
Question Encoder



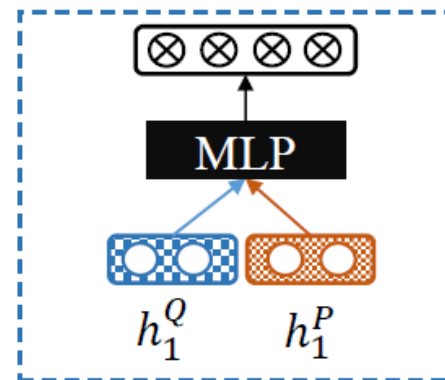
Passage Encoder



Decoder



Decoder Initialization



Attention-over-Attention (AoA)

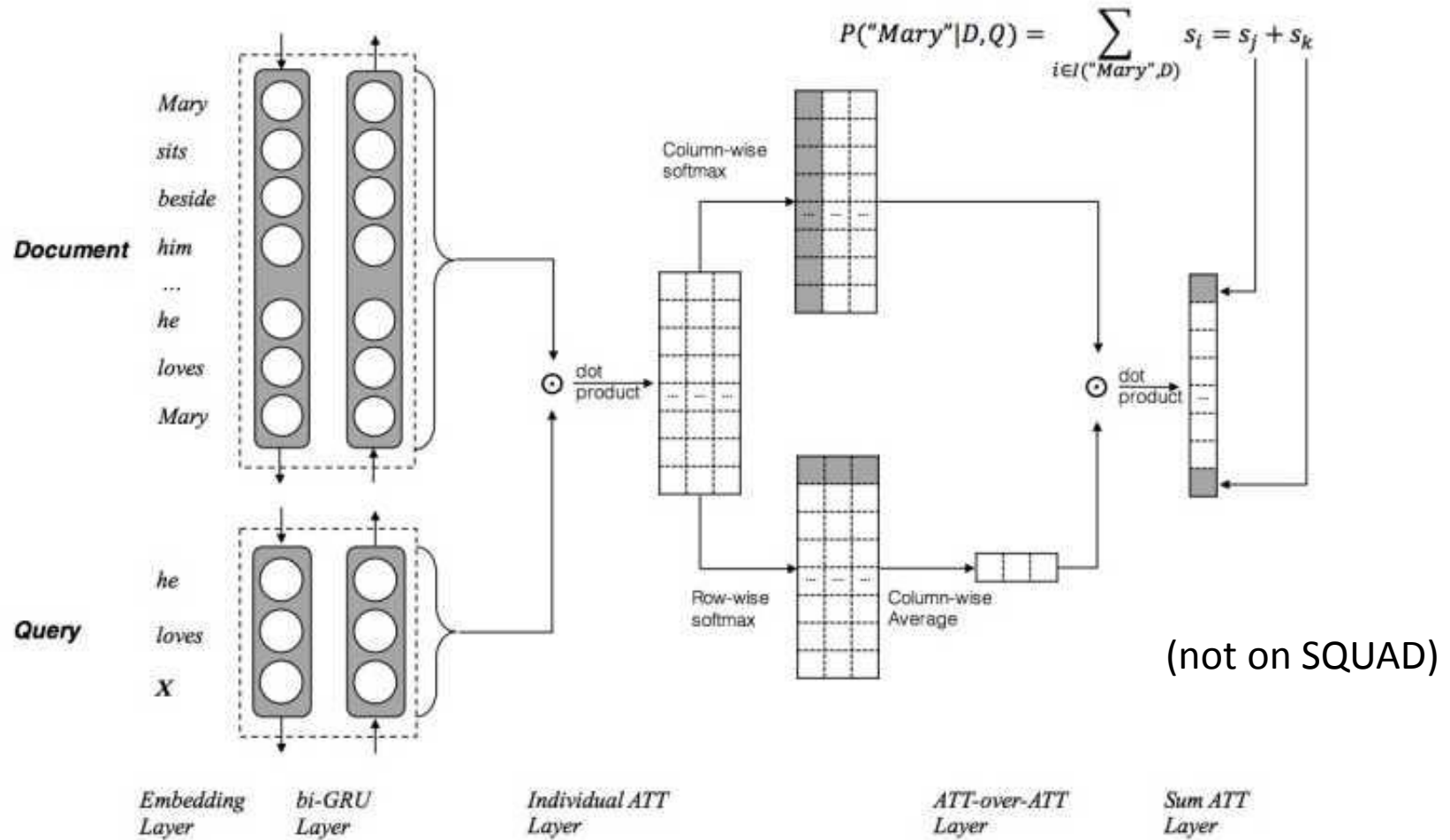


Figure 1: Neural network architecture of the proposed Attention-over-Attention Reader (AoA Reader).

1	Interactive AoA Reader+ (ensemble)	79.083	86.450
Oct 17, 2017	Joint Laboratory of HIT and iFLYTEK		

Reinforced Mnemonic Reader

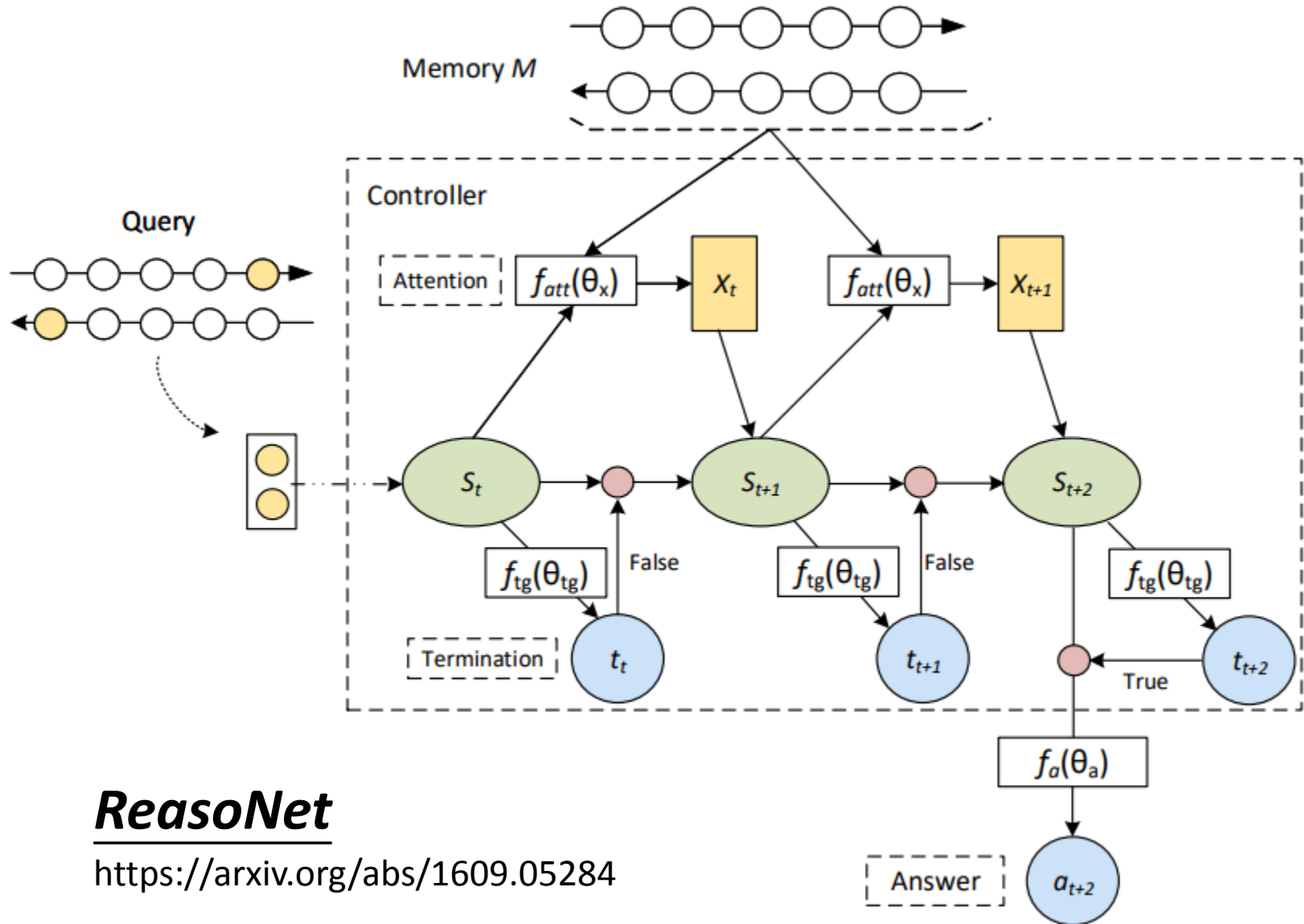
Reinforcement Learning for Machine Comprehension

One way to tackle this problem is to directly optimizing the F1 score with reinforcement learning. The F1 score measures the overlap between the predicted answer and the ground-truth answer, serving as a “soft” metric compared to the “hard” EM. Taking the F1 score as reward, we use the REINFORCE algorithm (Williams 1992) to maximize the model’s expected reward. For each sampled answer \hat{A} , we define the loss as:

$$J_{RL}(\theta) = -\mathbb{E}_{\hat{A} \sim p_{\theta}(A|C, Q)}[R(\hat{A}, A^*)] \quad (10)$$

where p_{θ} is the policy to be learned, and $R(\hat{A}, A^*)$ is the reward function for a sampled answer, computed as the F1 score with the ground-truth answer A^* . \hat{A} is obtained by sampling from the predicted probability distribution $p_{\theta}(A|C, Q)$.

Multiple-hop



ReasonNet

<https://arxiv.org/abs/1609.05284>

FusionNet

Architectures	(1)	(2)	(2')	(3)	(3')
Match-LSTM (Wang & Jiang, 2016)		✓			
DCN (Xiong et al., 2017)		✓			✓
FastQA (Weissenborn et al., 2017)	✓				
FastQAExt (Weissenborn et al., 2017)	✓	✓		✓	
BiDAF (Seo et al., 2017)		✓			✓
RaSoR (Lee et al., 2016)	✓		✓		
DrQA (Chen et al., 2017)	✓				
MPCM (Wang et al., 2016)	✓	✓			
Mnemonic Reader (Hu et al., 2017)	✓	✓		✓	
R-net (Wang et al., 2017)		✓		✓	

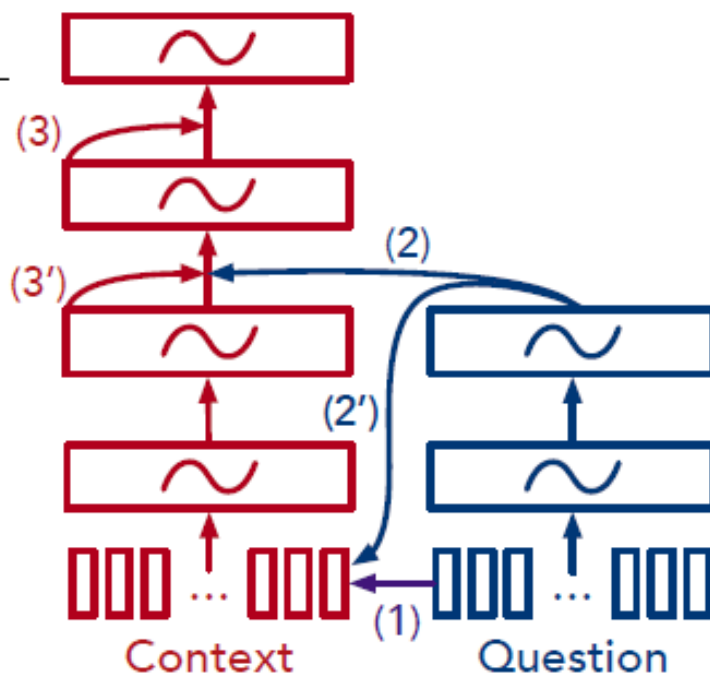
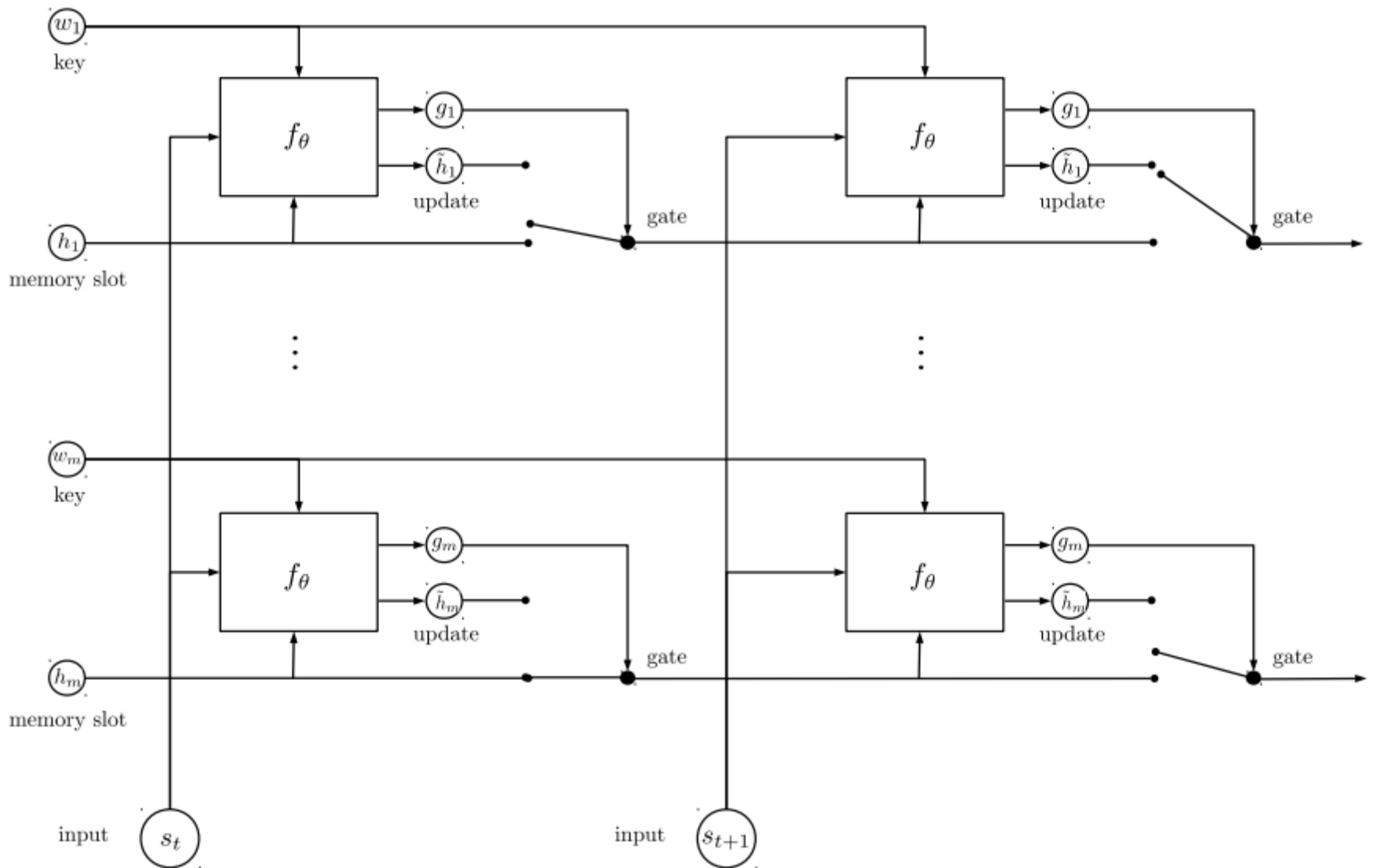


Table 1: A summarized view on the fusion processes used in several state-of-the-art architectures.

Figure 2: A conceptual architecture illustrating recent advances in MRC.

1	FusionNet (ensemble)	78.978	86.016
Oct 24, 2017	Microsoft Business AI Solutions Team		

Recurrent Entity Networks



Task	NTM	D-NTM	MemN2N	DNC	DMN+	EntNet
1: 1 supporting fact	31.5	4.4	0	0	0	0
2: 2 supporting facts	54.5	27.5	0.3	0.4	0.3	0.1
3: 3 supporting facts	43.9	71.3	2.1	1.8	1.1	4.1
4: 2 argument relations	0	0	0	0	0	0
5: 3 argument relations	0.8	1.7	0.8	0.8	0.5	0.3
6: yes/no questions	17.1	1.5	0.1	0	0	0.2
7: counting	17.8	6.0	2.0	0.6	2.4	0
8: lists/sets	13.8	1.7	0.9	0.3	0.0	0.5
9: simple negation	16.4	0.6	0.3	0.2	0.0	0.1
10: indefinite knowledge	16.6	19.8	0	0.2	0	0.6
11: basic coreference	15.2	0	0.0	0	0.0	0.3
12: conjunction	8.9	6.2	0	0	0.2	0
13: compound coreference	7.4	7.5	0	0	0	1.3
14: time reasoning	24.2	17.5	0.2	0.4	0.2	0
15: basic deduction	47.0	0	0	0	0	0
16: basic induction	53.6	49.6	51.8	55.1	45.3	0.2
17: positional reasoning	25.5	1.2	18.6	12.0	4.2	0.5
18: size reasoning	2.2	0.2	5.3	0.8	2.1	0.3
19: path finding	4.3	39.5	2.3	3.9	0.0	2.3
20: agent's motivation	1.5	0	0	0	0	0

Failed Tasks (> 5% error):

16

9

3

2

1

0

Mean Error:

20.1

12.8

4.2

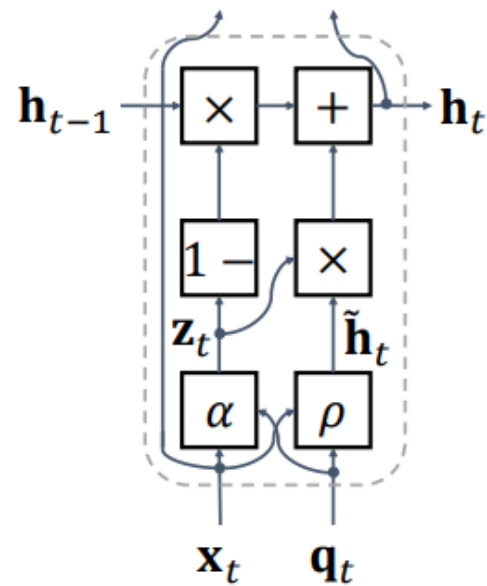
3.8

2.8

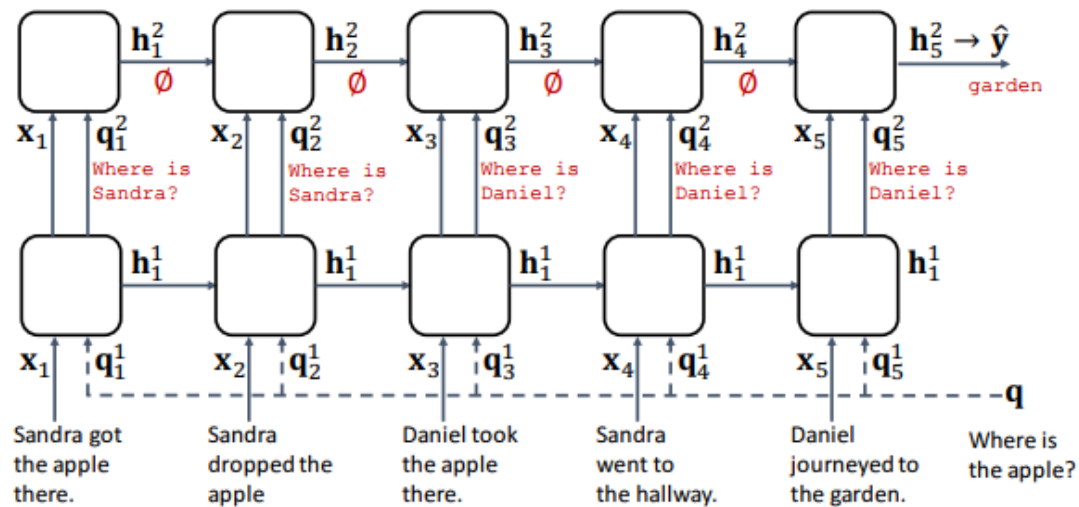
0.5

Query-Reduction Networks for Question Answering

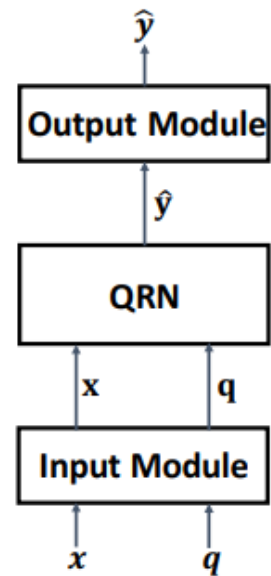
- <https://arxiv.org/pdf/1606.04582.pdf>



(a) QRN unit



(b) 2-layer QRN



(c) Overview

Query-Reduction Networks for Question Answering

- <https://arxiv.org/pdf/1606.04582.pdf>

Task	1k										10k						
	Previous works				QRN						Previous works			QRN			
	LSTM	N2N	DMN+	GMemN2N	1r	2	2r	3r	6r	6r200*	N2N	DMN+	GMemN2N	2r	2rv	3r	6r200
1: Single supporting fact	50.0	0.1	1.3	0.0	0.0	0.0	0.0	0.0	0.0	13.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2: Two supporting facts	80.0	18.8	72.3	8.1	65.7	1.2	0.7	0.5	1.5	15.3	0.3	0.3	0.0	0.4	0.8	0.4	0.0
3: Three supporting facts	80.0	31.7	73.3	38.7	68.2	17.5	5.7	1.2	15.3	13.8	2.1	1.1	4.5	0.4	1.4	0.0	0.0
4: Two arg relations	39.0	17.5	26.9	0.4	0.0	0.0	0.0	0.7	9.0	13.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5: Three arg relations	30.0	12.9	25.6	1.0	1.0	1.1	1.1	1.2	1.3	12.5	0.8	0.5	0.2	0.5	0.2	0.3	0.0
6: Yes/no questions	52.0	2.0	28.5	8.4	0.1	0.0	0.9	1.2	50.6	15.5	0.1	0.0	0.0	0.0	0.0	0.0	0.0
7: Counting	51.0	10.1	21.9	17.8	10.9	11.1	9.6	9.4	13.1	15.3	2.0	2.4	1.8	1.0	0.7	0.7	0.0
8: Lists/sets	55.0	6.1	21.9	12.5	6.8	5.7	5.6	3.7	7.8	15.1	0.9	0.0	0.3	1.4	0.6	0.8	0.4
9: Simple negation	36.0	1.5	42.9	10.7	0.0	0.6	0.0	0.0	32.7	13.0	0.3	0.0	0.0	0.0	0.0	0.0	0.0
10: Indefinite knowledge	56.0	2.6	23.1	16.5	0.8	0.6	0.0	0.0	3.5	12.9	0.0	0.0	0.2	0.0	0.0	0.0	0.0
11: Basic coreference	38.0	3.3	4.3	0.0	11.3	0.5	0.0	0.0	0.9	14.7	0.1	0.0	0.0	0.0	0.0	0.0	0.0
12: Conjunction	26.0	0.0	3.5	0.0	0.0	0.0	0.0	0.0	0.0	15.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0
13: Compound coreference	6.0	0.5	7.8	0.0	5.3	5.5	0.0	0.3	8.9	13.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0
14: Time reasoning	73.0	2.0	61.9	1.2	20.2	1.3	0.8	3.8	18.2	14.5	0.1	0.0	0.0	0.2	0.0	0.0	0.1
15: Basic deduction	79.0	1.8	47.6	0.0	39.4	0.0	0.0	0.0	0.1	14.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0
16: Basic induction	77.0	51.0	54.4	0.1	50.6	54.8	53.0	53.4	53.5	15.5	51.8	45.3	0.0	49.4	50.4	49.1	0.0
17: Positional reasoning	49.0	42.6	44.1	41.7	40.6	36.5	34.4	51.8	52.0	13.0	18.6	4.2	27.8	0.9	0.0	5.8	4.1
18: Size reasoning	48.0	9.2	9.1	9.2	8.2	8.6	7.9	8.8	47.5	14.9	5.3	2.1	8.5	1.6	8.4	1.8	0.7
19: Path finding	92.0	90.6	90.8	88.5	88.8	89.8	78.7	90.7	88.6	13.6	2.3	0.0	31.0	36.1	1.0	27.9	0.1
20: Agents motivations	9.0	0.2	2.2	0.0	0.0	0.0	0.2	0.3	5.5	14.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0
# Failed	20	10	16	10	12	8	7	5	13	20	3	1	3	2	2	3	0
Average error rates (%)	51.3	15.2	33.2	12.7	20.1	11.7	9.9	11.3	20.5	14.2	4.2	2.8	3.7	4.6	3.2	4.3	0.3

Acknowledgement

- 感謝 曹燁文 同學發現投影片上的錯字