

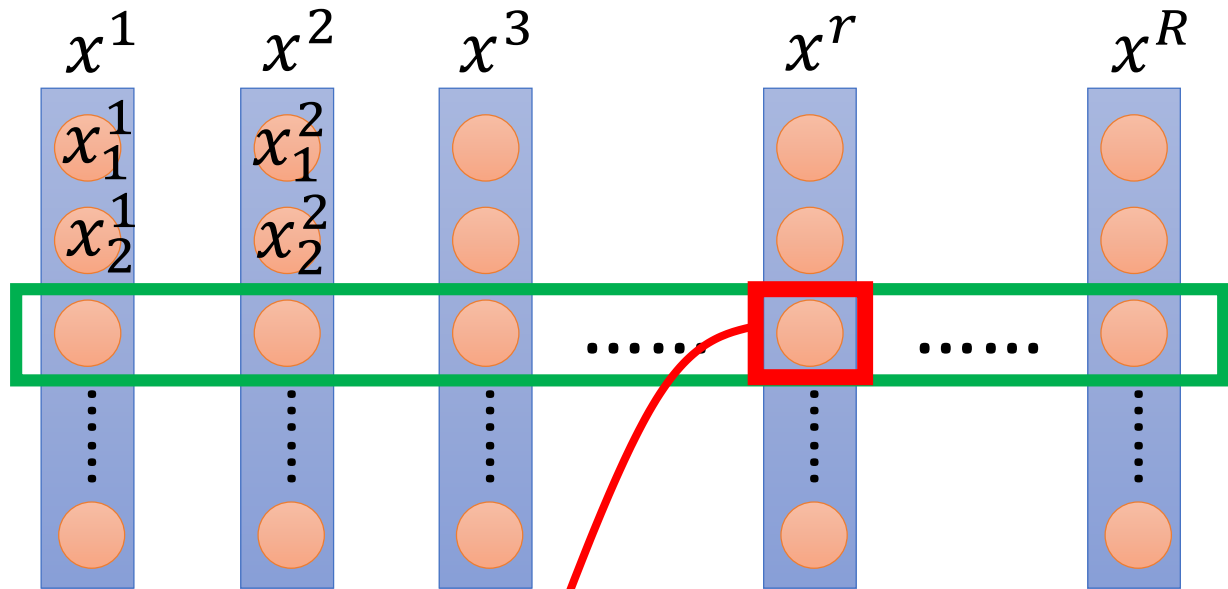
Tips for Training Deep Network

Output

- Training Strategy: Batch Normalization
- Activation Function: SELU
- Network Structure: Highway Network

Batch Normalization

Feature Scaling



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

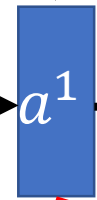
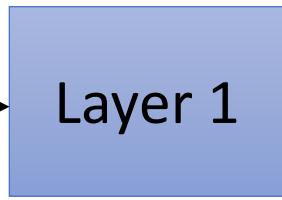
$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

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

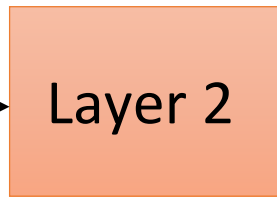
In general, gradient descent converges much faster with feature scaling than without it.

How about Hidden Layer?

Feature Scaling



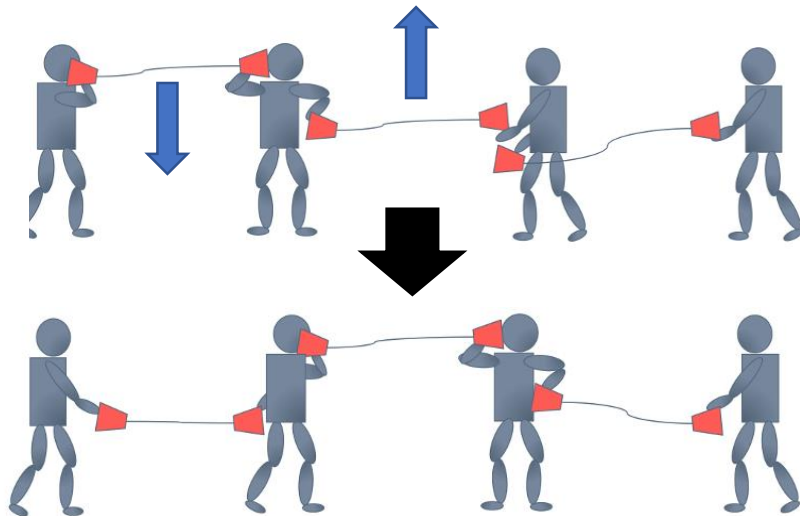
Feature Scaling ?



Feature Scaling ?



.....



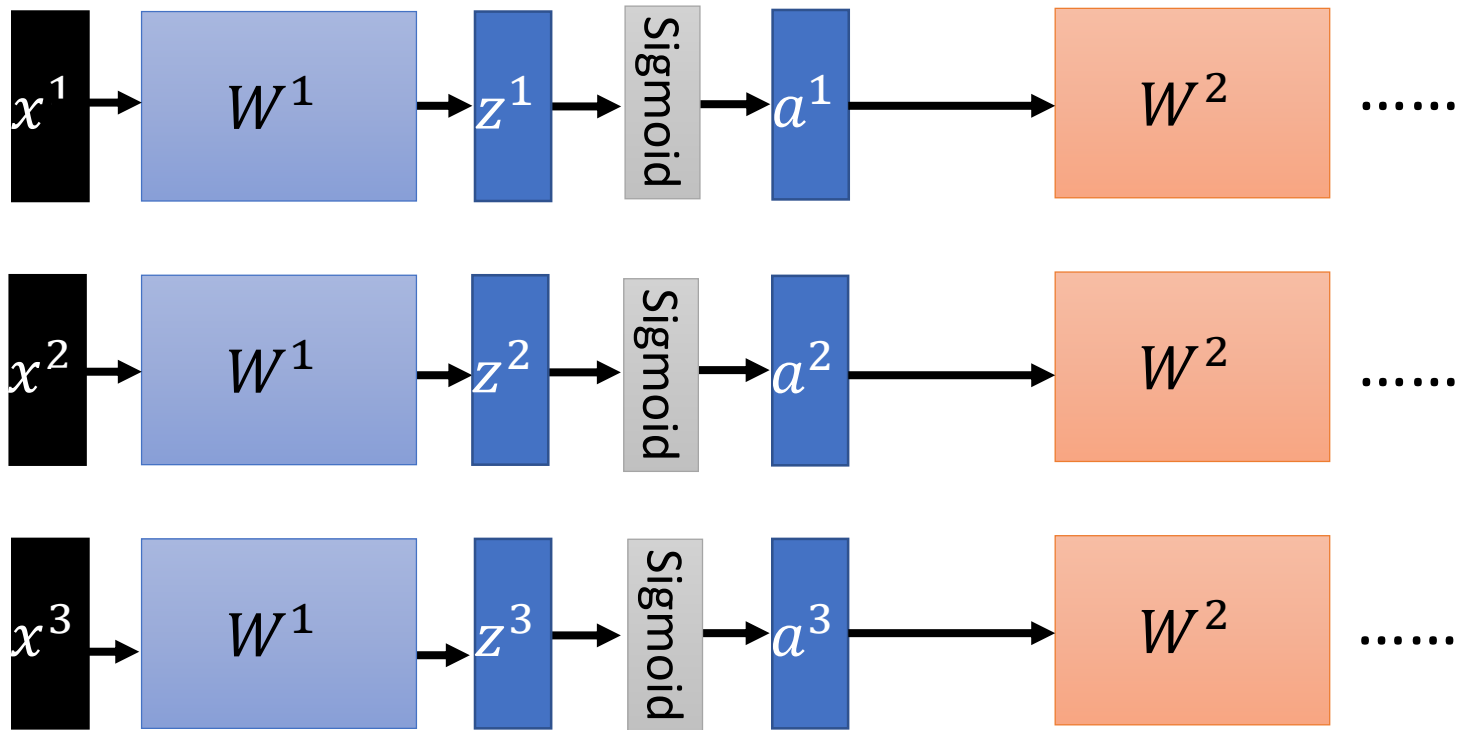
Internal Covariate Shift

Difficulty: their statistics change during the training ...

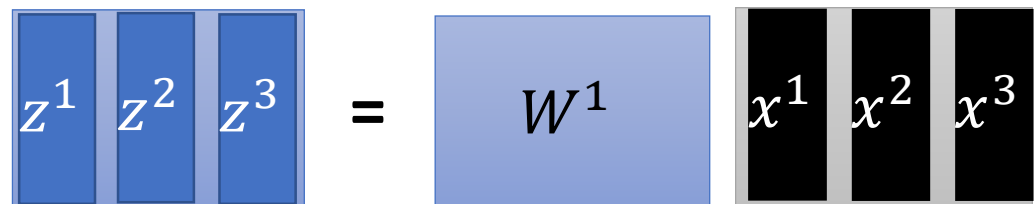
➔ **Batch normalization**

Smaller learning rate can be helpful, but the training would be slower.

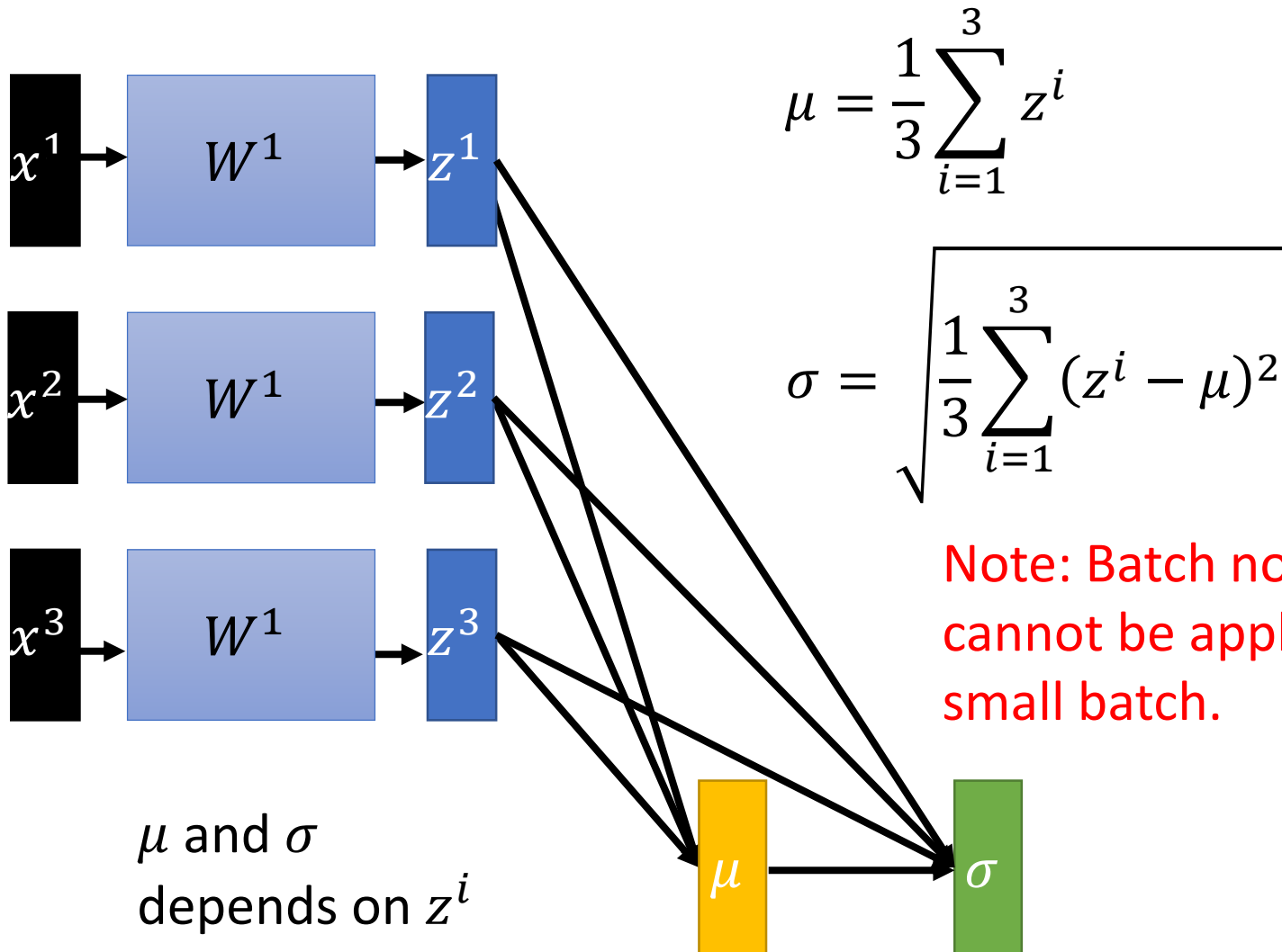
Batch



Batch

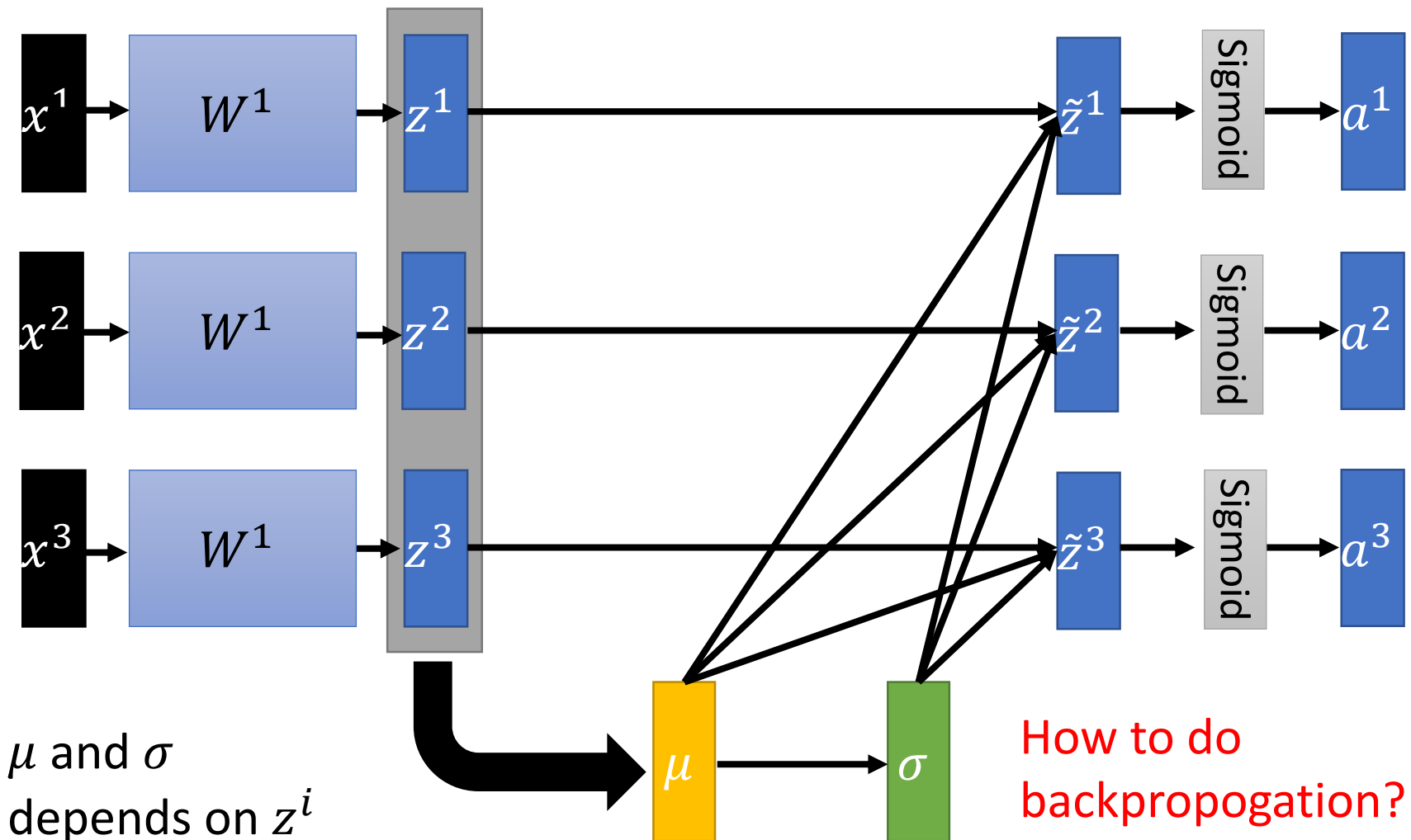


Batch normalization



Batch normalization

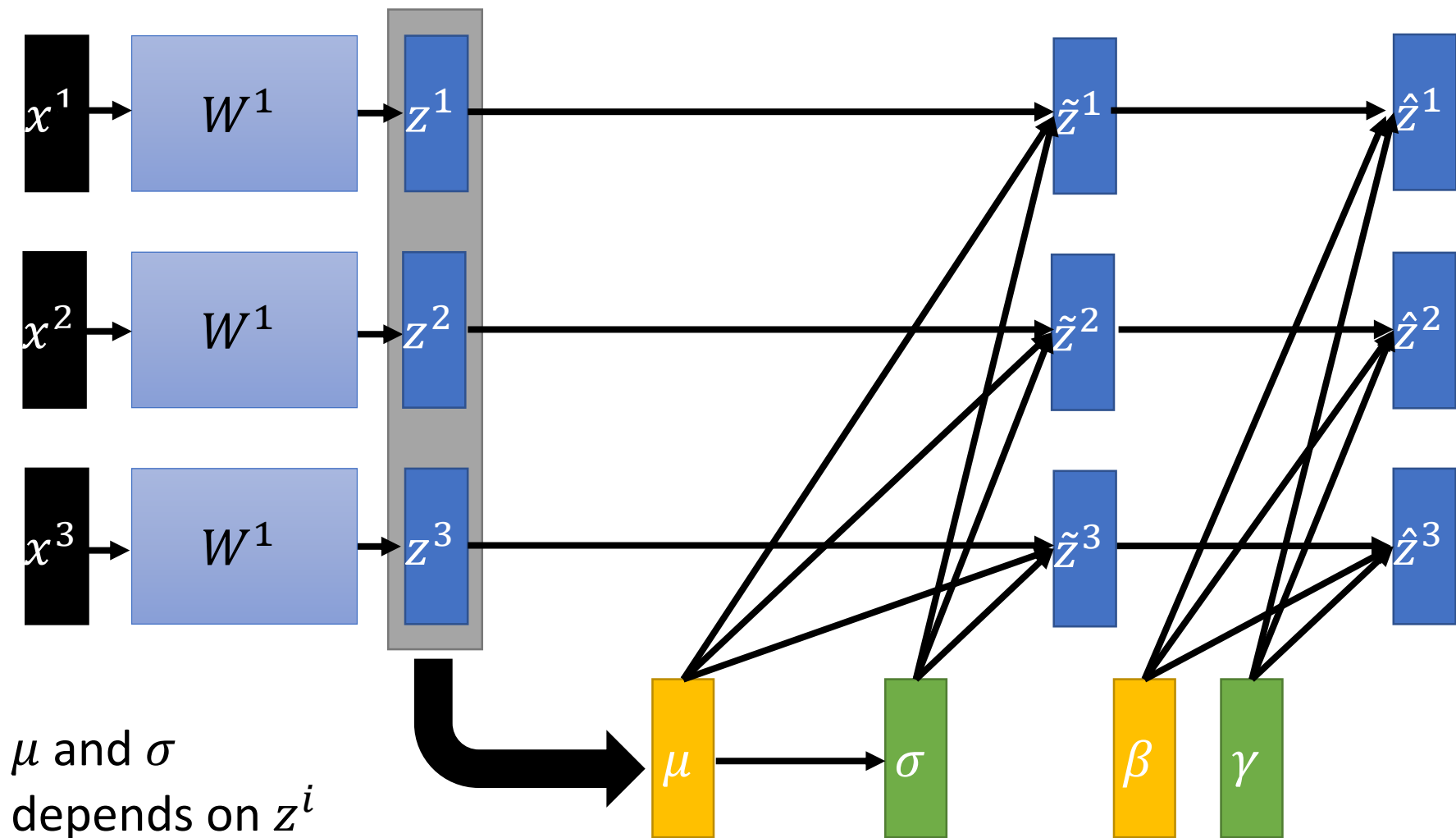
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$



Batch normalization

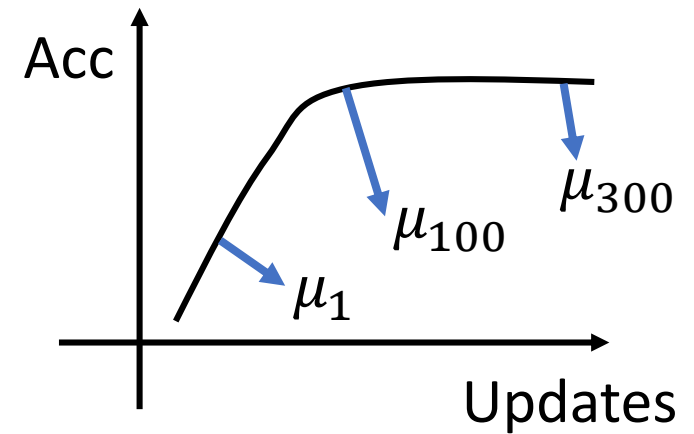
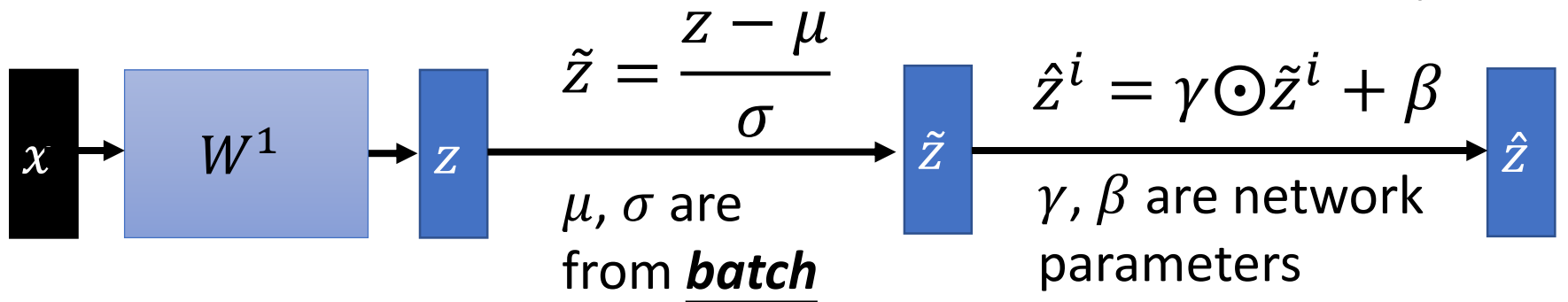
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

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



Batch normalization

- At testing stage:



We do not have batch at testing stage.

Ideal solution:

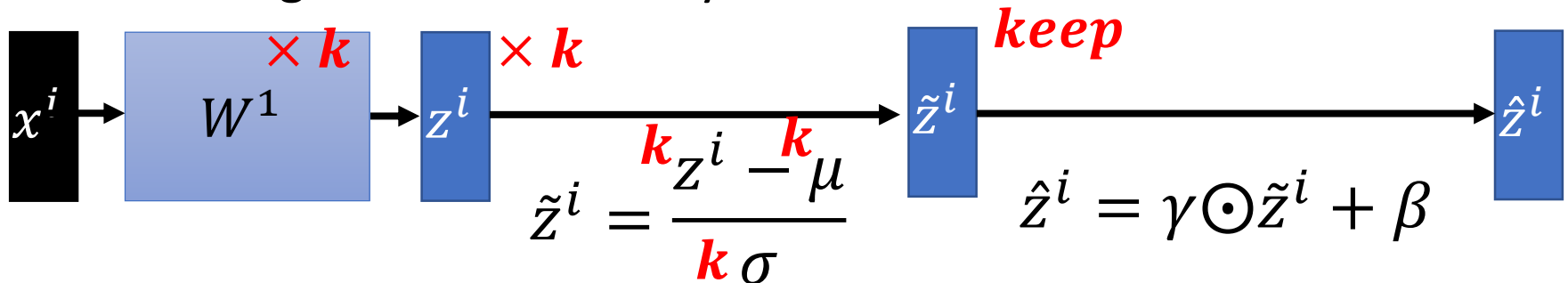
Computing μ and σ using the whole training dataset.

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.

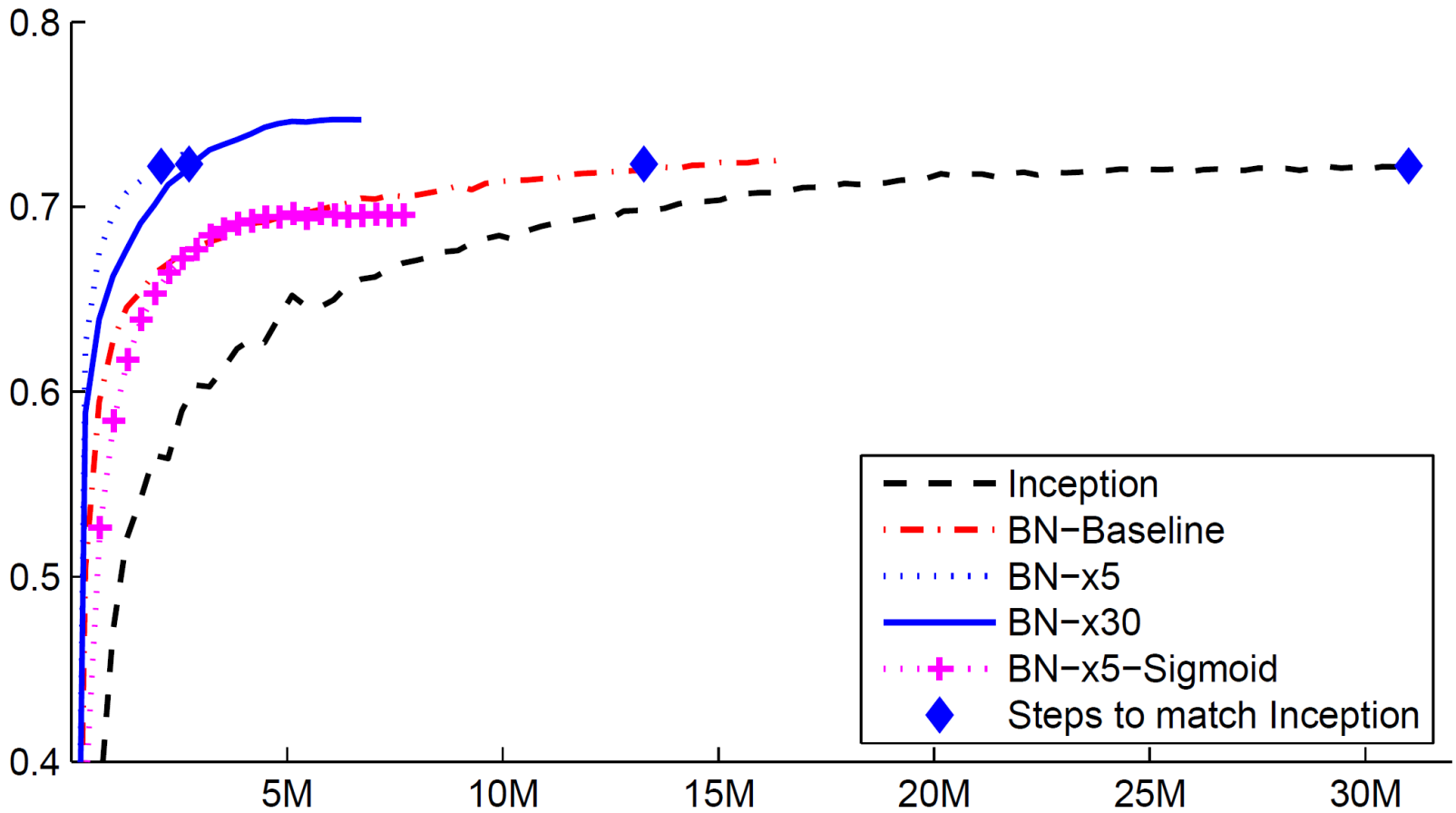


Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*

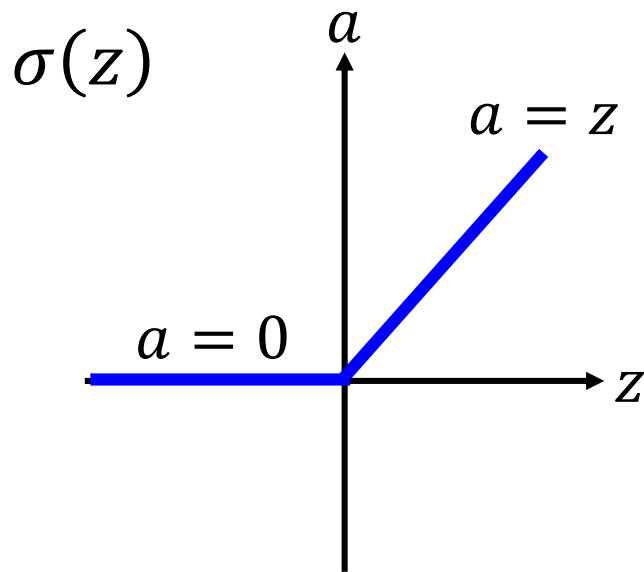
To learn more

- Batch Renormalization
- Layer Normalization
- Instance Normalization
- Weight Normalization
- Spectrum Normalization

Activation Function: SELU

ReLU

- Rectified Linear Unit (ReLU)

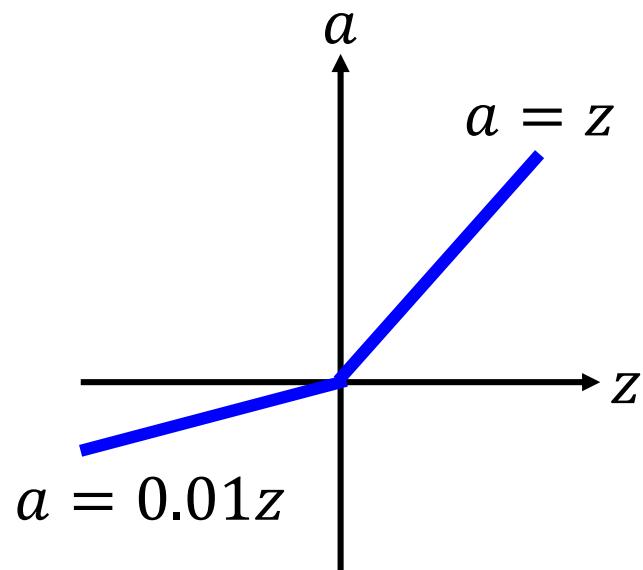


Reason:

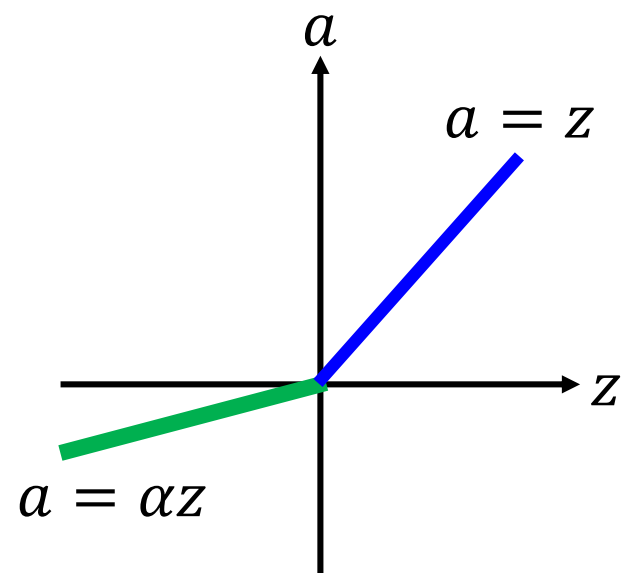
1. Fast to compute
2. Biological reason
3. Infinite sigmoid with different biases
4. Vanishing gradient problem

ReLU - variant

Leaky ReLU



Parametric ReLU



α also learned by
gradient descent

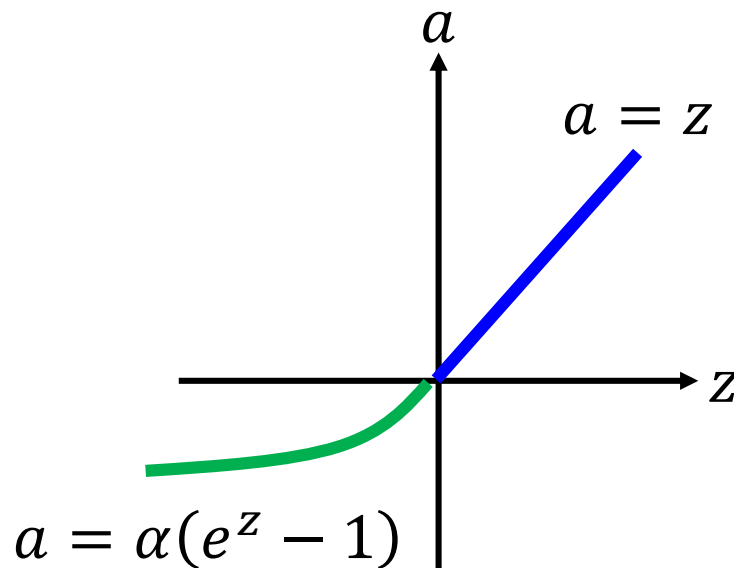
(1) Definition of scaled exponential linear units (SELUs)

In [3]:

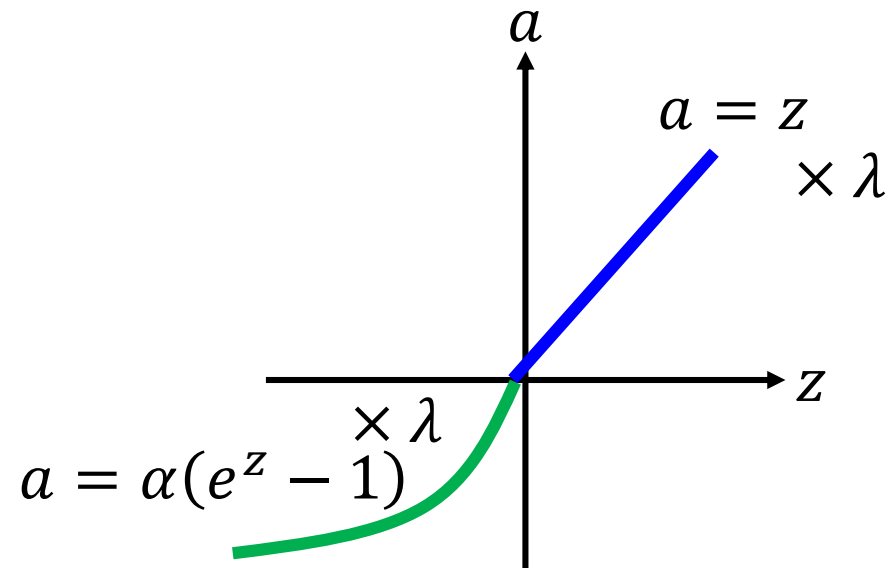
```
def selu(x):  
    with ops.name_scope('elu') as scope:  
        alpha = 1.6732632423543772848170429916717  
        scale = 1.0507009873554804934193349852946  
        return scale*tf.where(x>=0.0, x, alpha*tf.nn.elu(x))
```

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

Exponential Linear Unit (ELU)



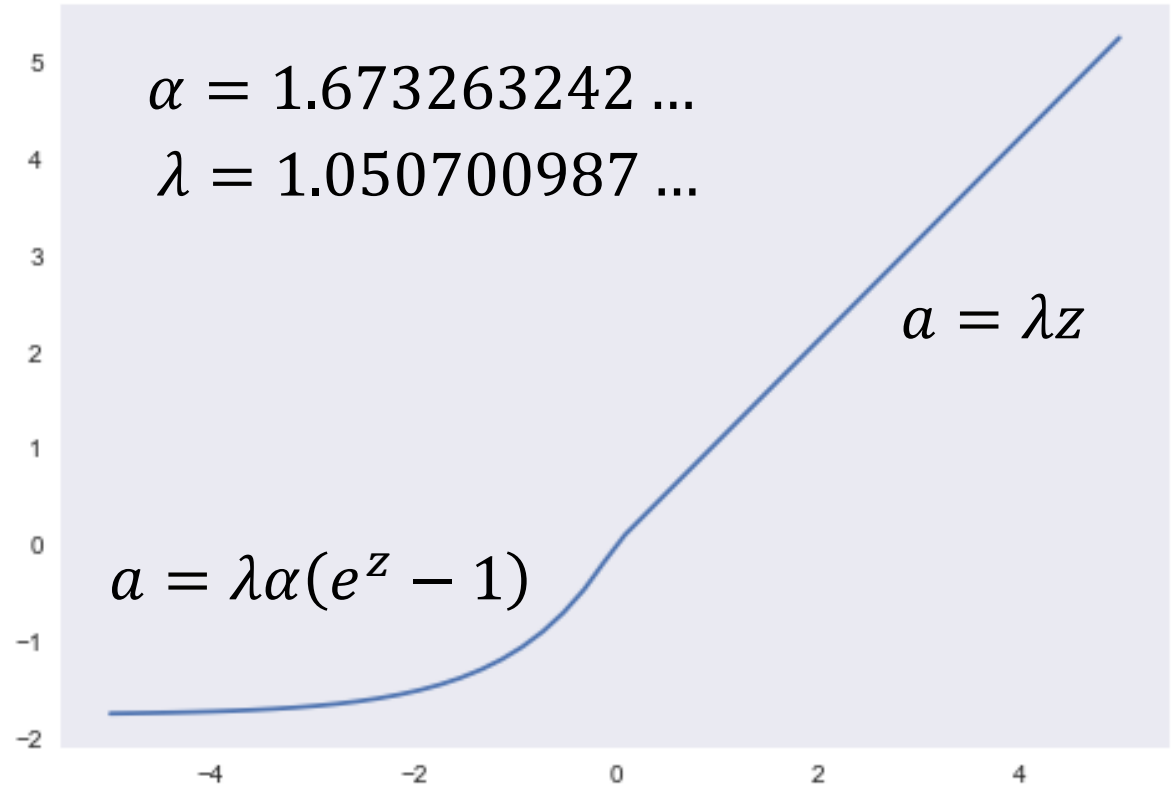
Scaled ELU (SELU)



$$\alpha = 1.6732632423543772848170429916717$$

$$\lambda = 1.0507009873554804934193349852946$$

SELU



Positive and negative values

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

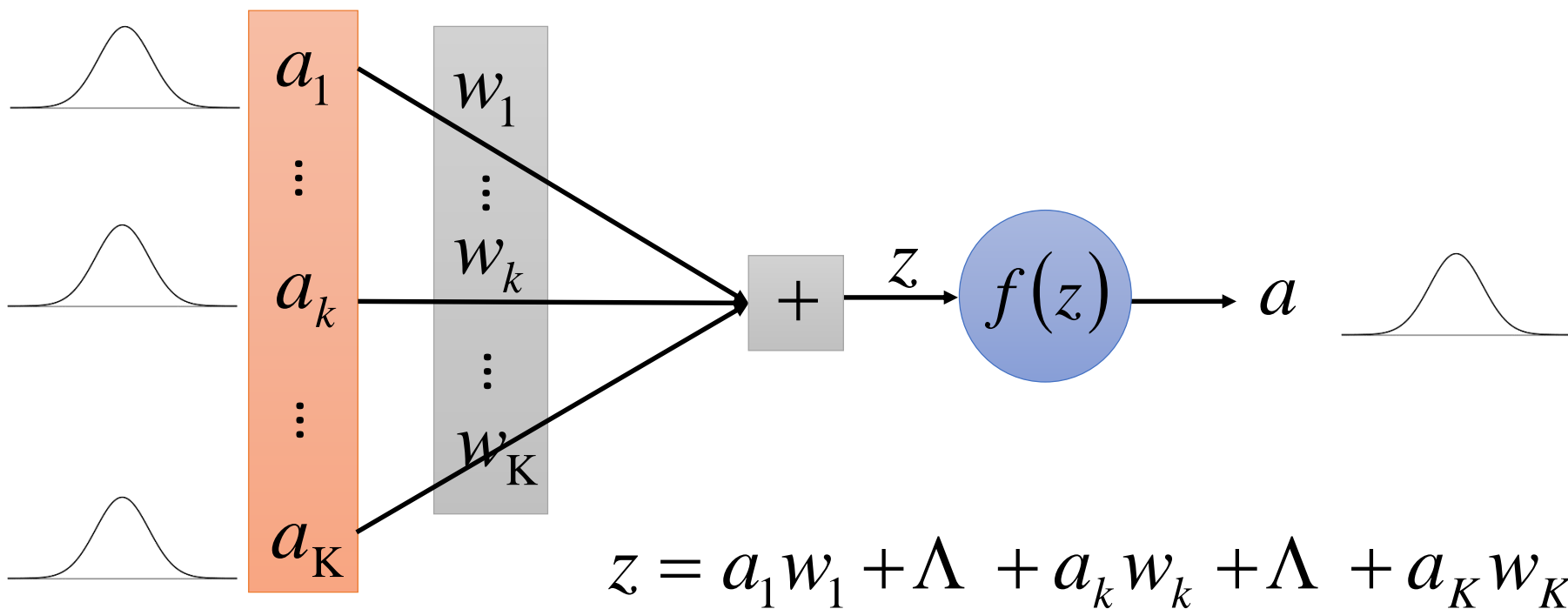
Saturation region ➡ ELU also has this property

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

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

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

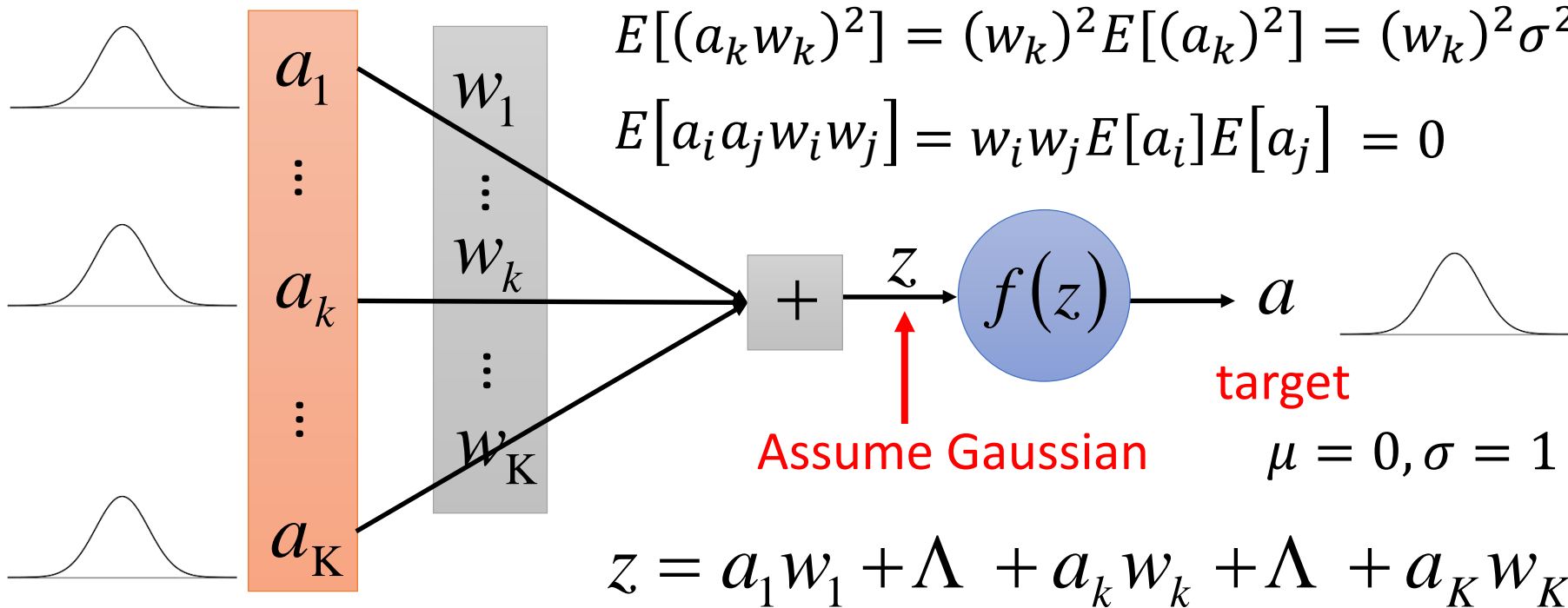
$$= E[(a_1 w_1 + a_2 w_2 + \dots)^2]$$

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

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

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

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



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

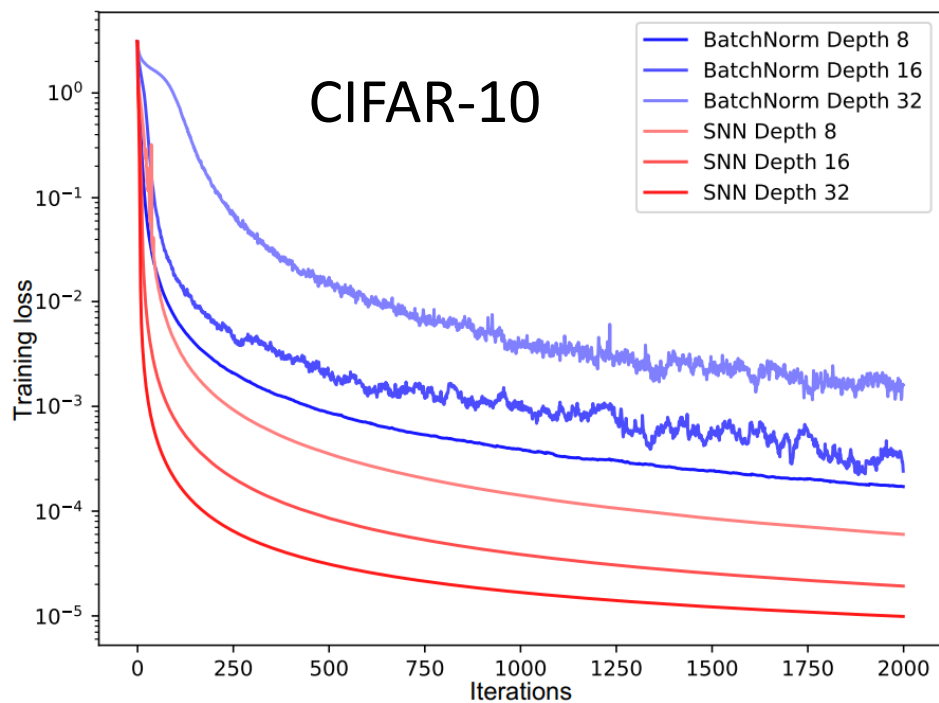
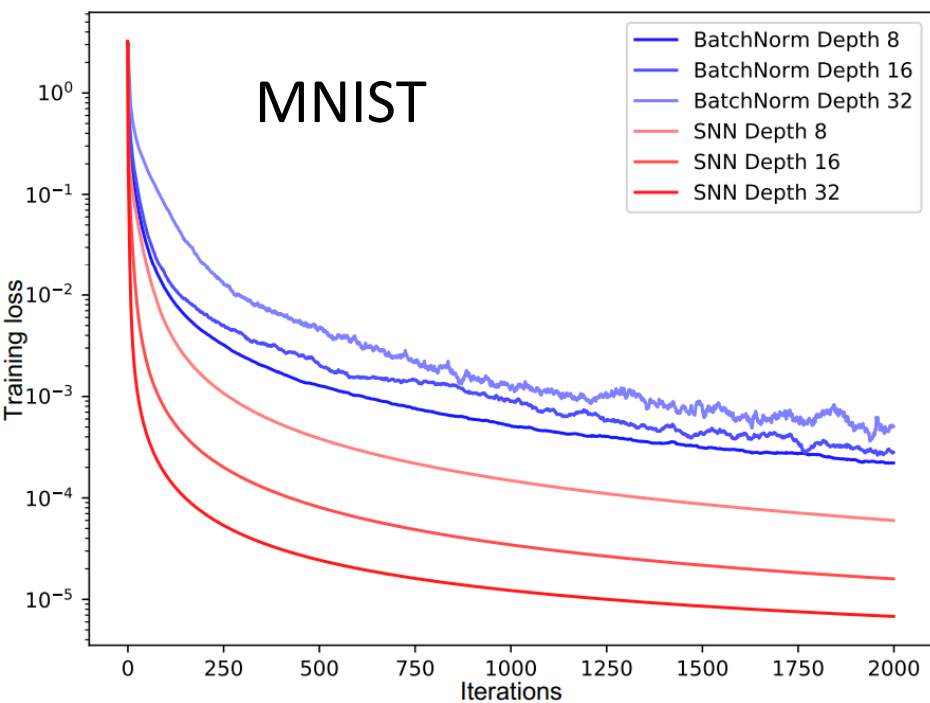


4



41





FNN method comparison

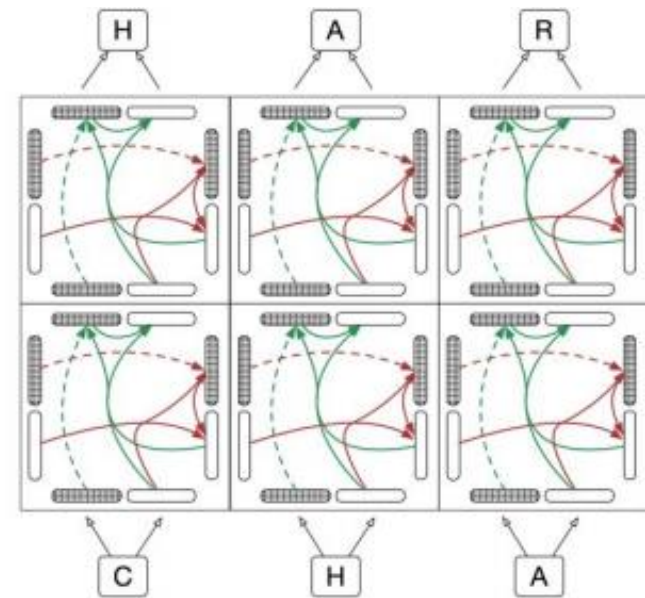
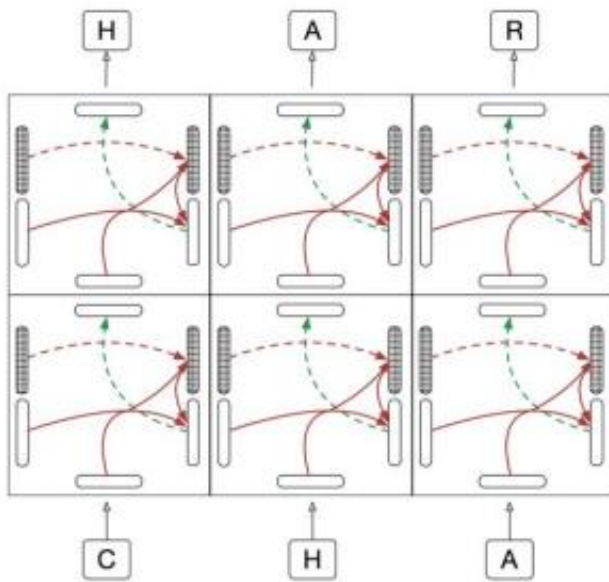
Method	avg. rank diff.	<i>p</i> -value
SNN	-0.756	
MSRAinit	-0.240*	2.7e-02
LayerNorm	-0.198*	1.5e-02
Highway	0.021*	1.9e-03
ResNet	0.273*	5.4e-04
WeightNorm	0.397*	7.8e-07
BatchNorm	0.504*	3.5e-06

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
...

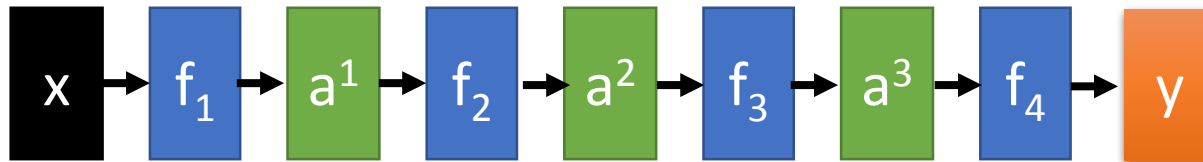
Demo

Highway Network & Grid LSTM



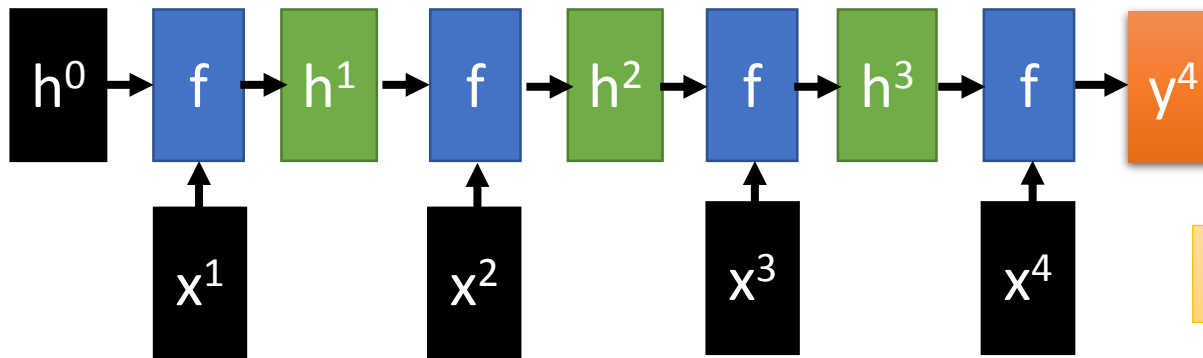
Feedforward v.s. Recurrent

1. Feedforward network does not have input at each step
2. Feedforward network has different parameters for each layer



$$a^t = f_l(a^{t-1}) = \sigma(W^t a^{t-1} + b^t)$$

t is layer



$$h^t = f(h^{t-1}, x^t) = \sigma(W^h h^{t-1} + W^i x^t + b^i)$$

t is time step

Applying gated structure in feedforward network

GRU \rightarrow Highway Network

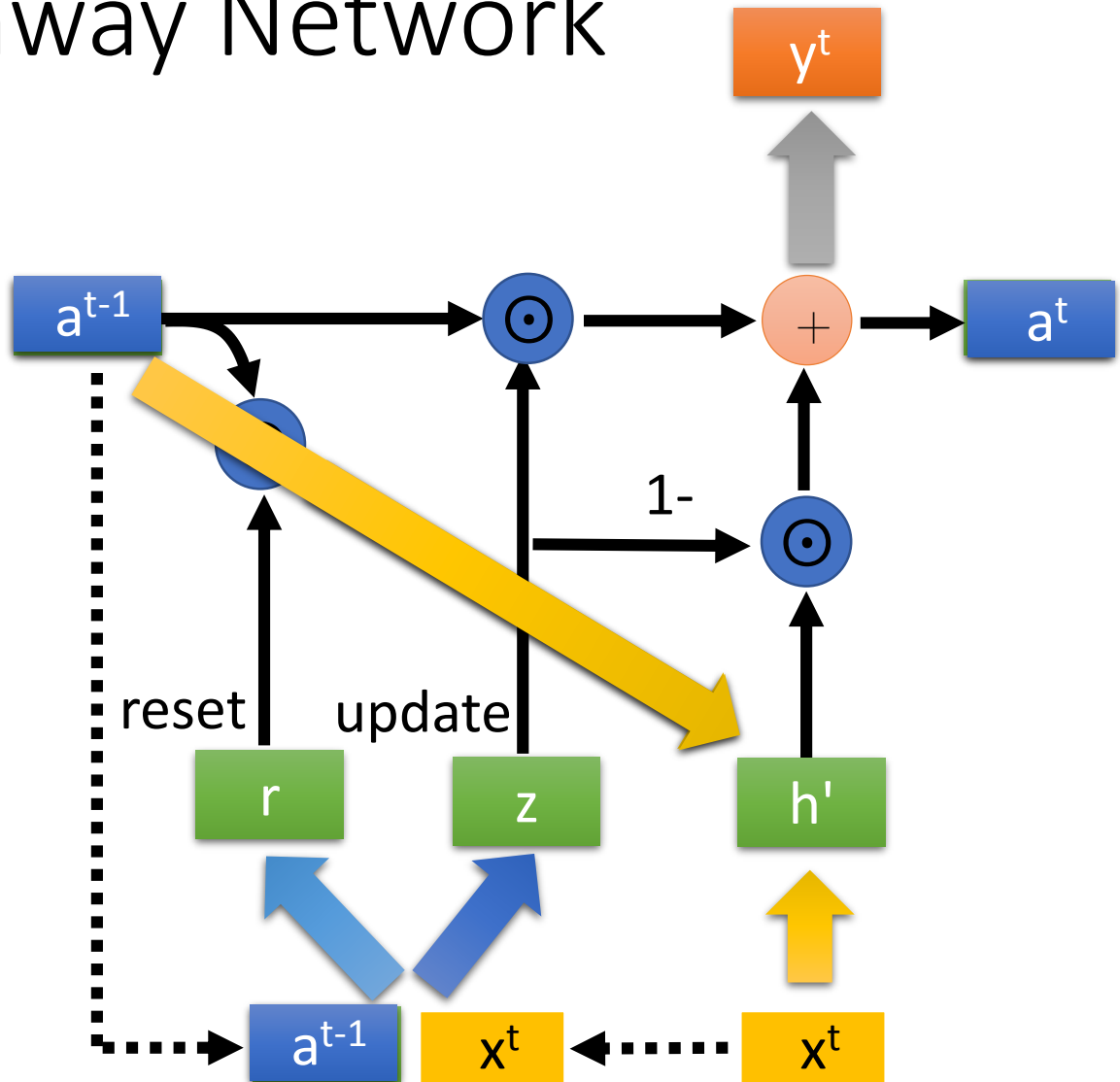
No input x^t at each step

No output y^t at each step

a^{t-1} is the output of the (t-1)-th layer

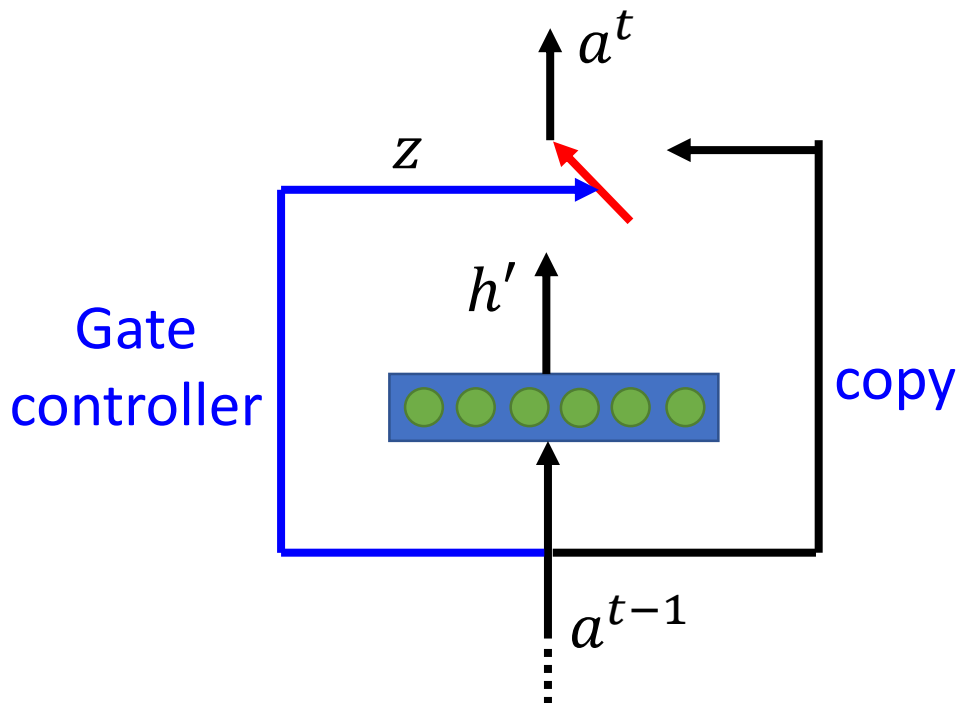
a^t is the output of the t-th layer

No reset gate



Highway Network

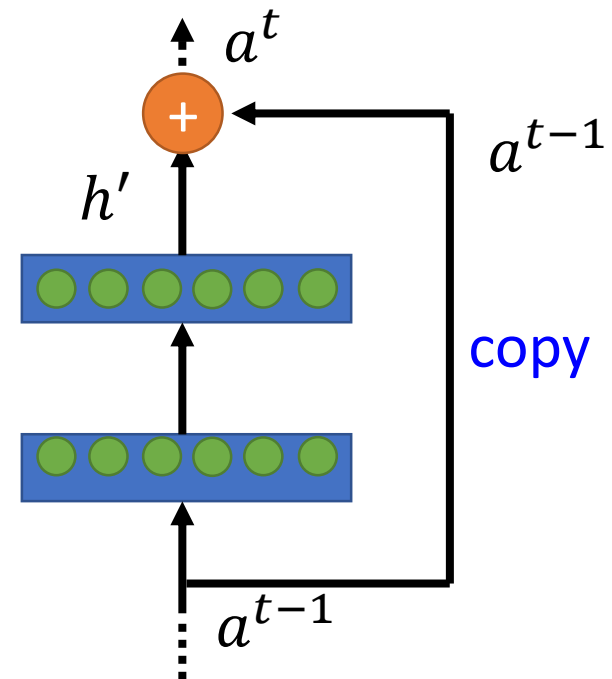
- **Highway Network**



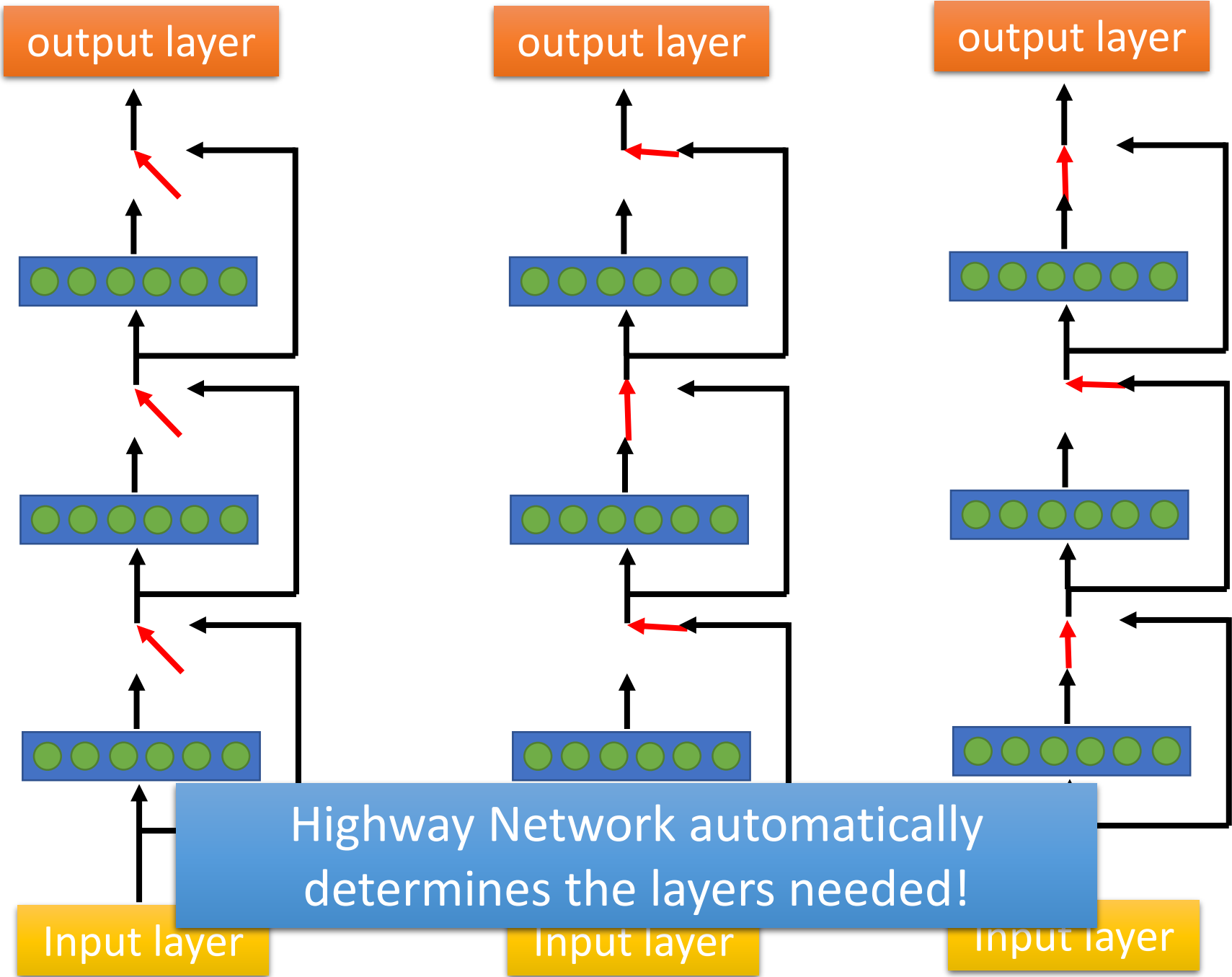
Training Very Deep Networks
<https://arxiv.org/pdf/1507.06228v2.pdf>

$$h' = \sigma(Wa^{t-1})$$
$$z = \sigma(W'a^{t-1})$$
$$a^t = z \odot a^{t-1} + (1 - z) \odot h$$

- **Residual Network**



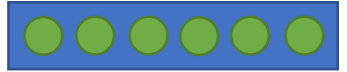
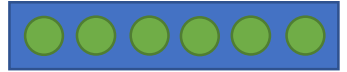
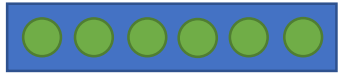
Deep Residual Learning for Image Recognition
<http://arxiv.org/abs/1512.03385>



output layer

output layer

output layer



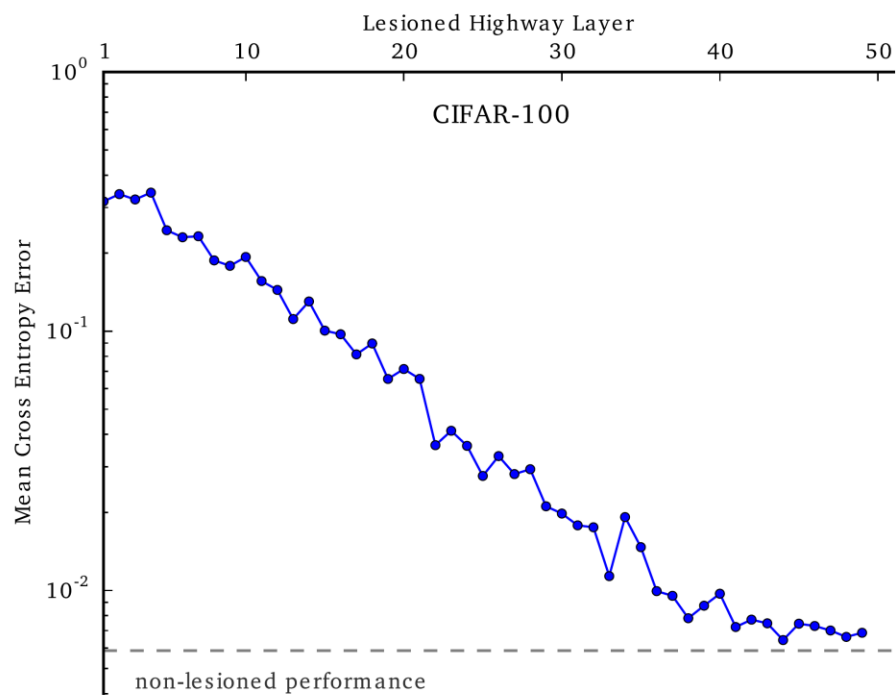
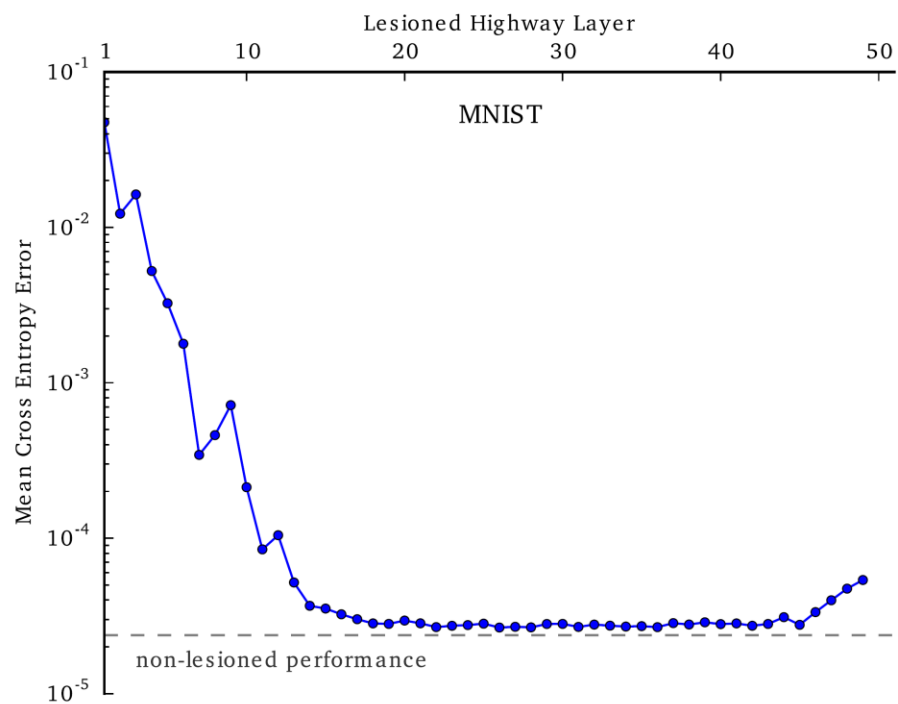
Highway Network automatically determines the layers needed!

Input layer

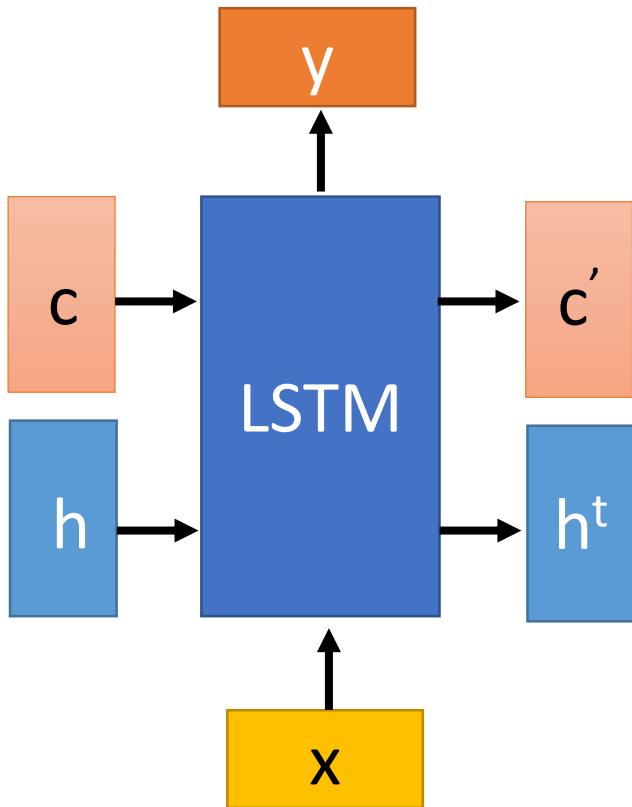
Input layer

Input layer

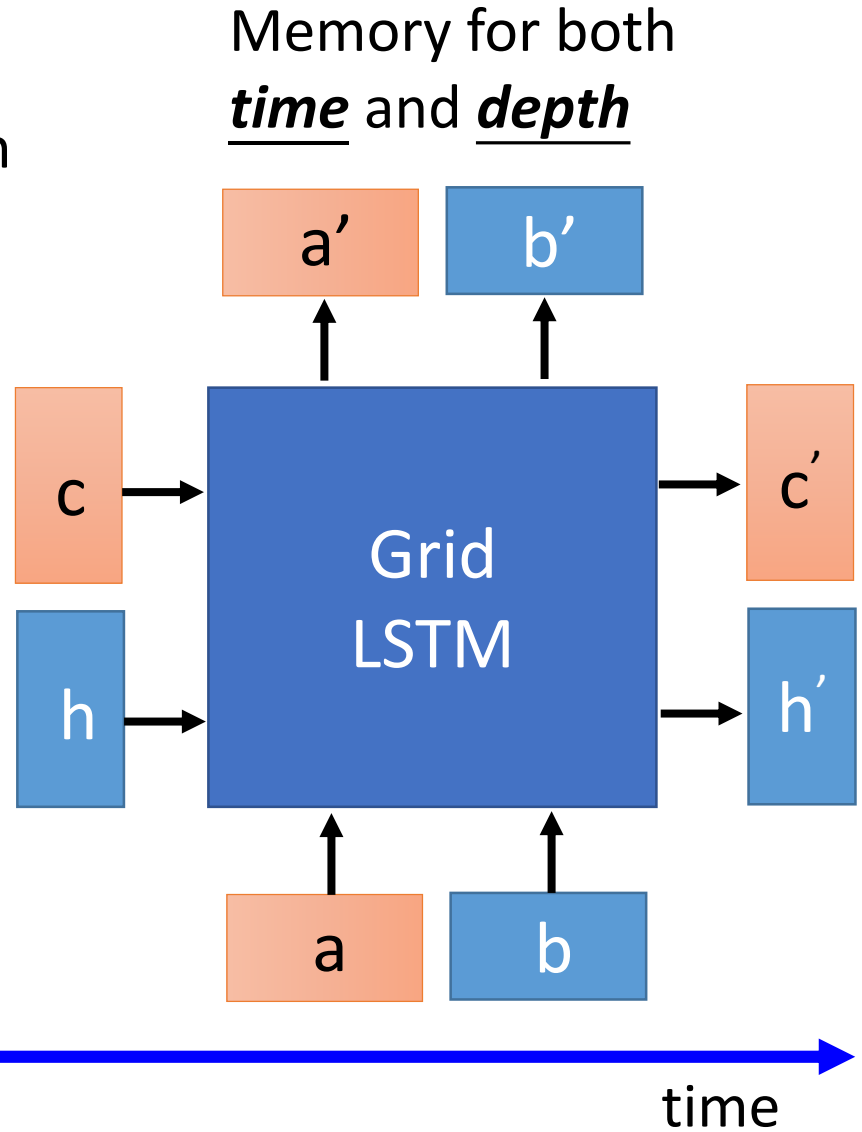
Highway Network

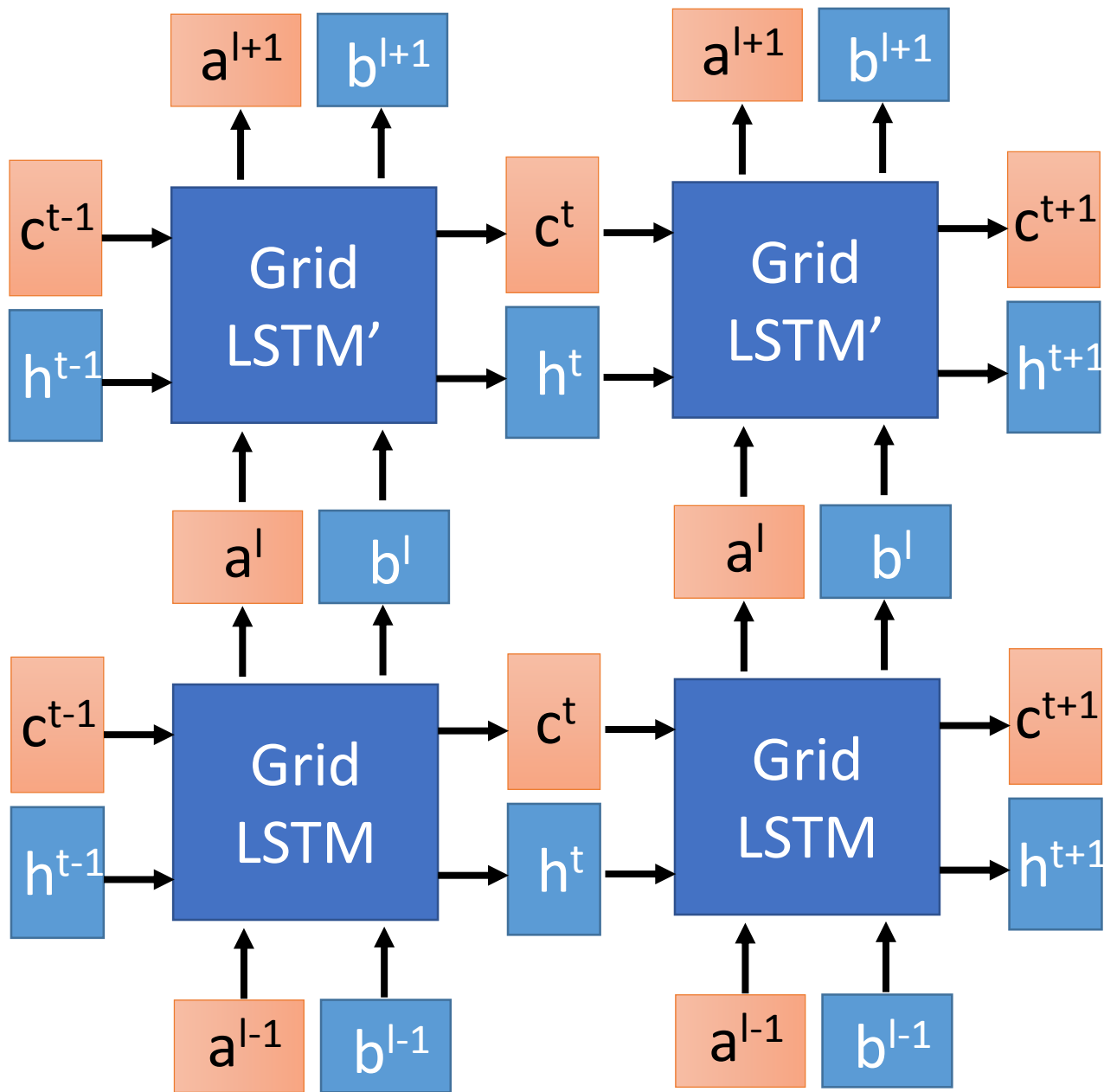


Grid LSTM

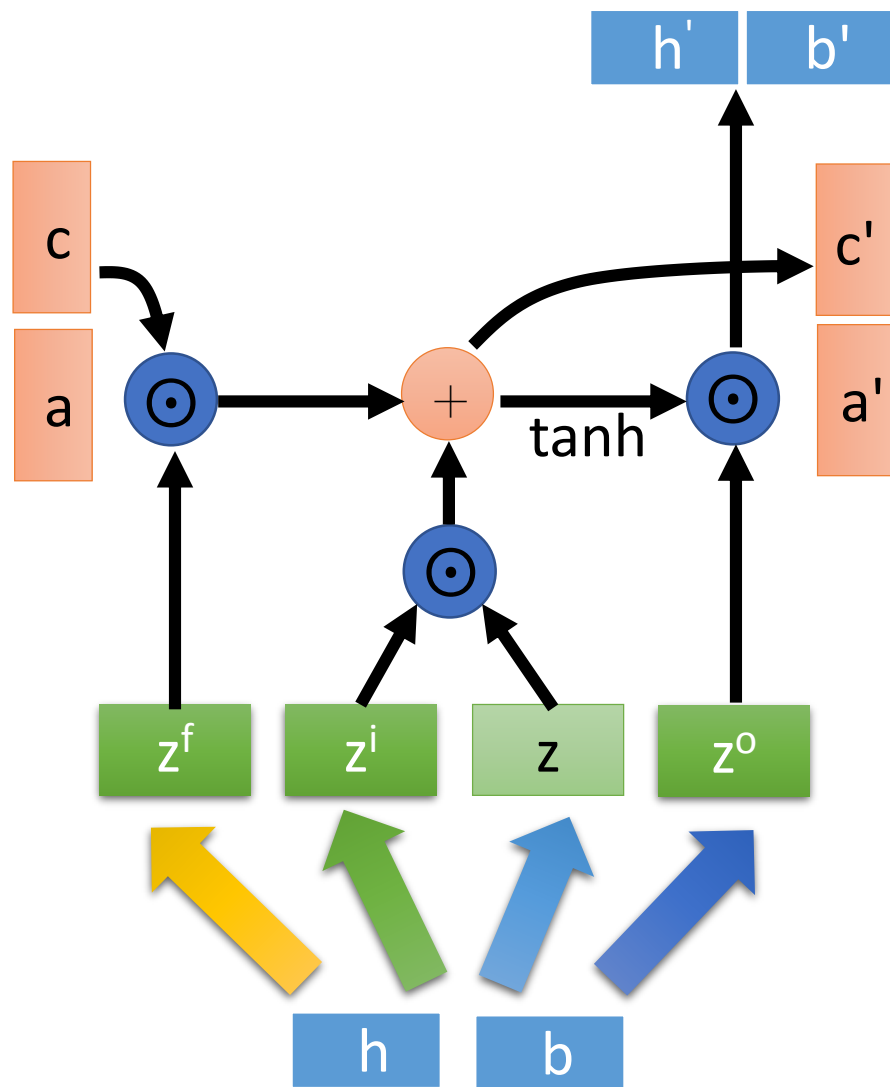
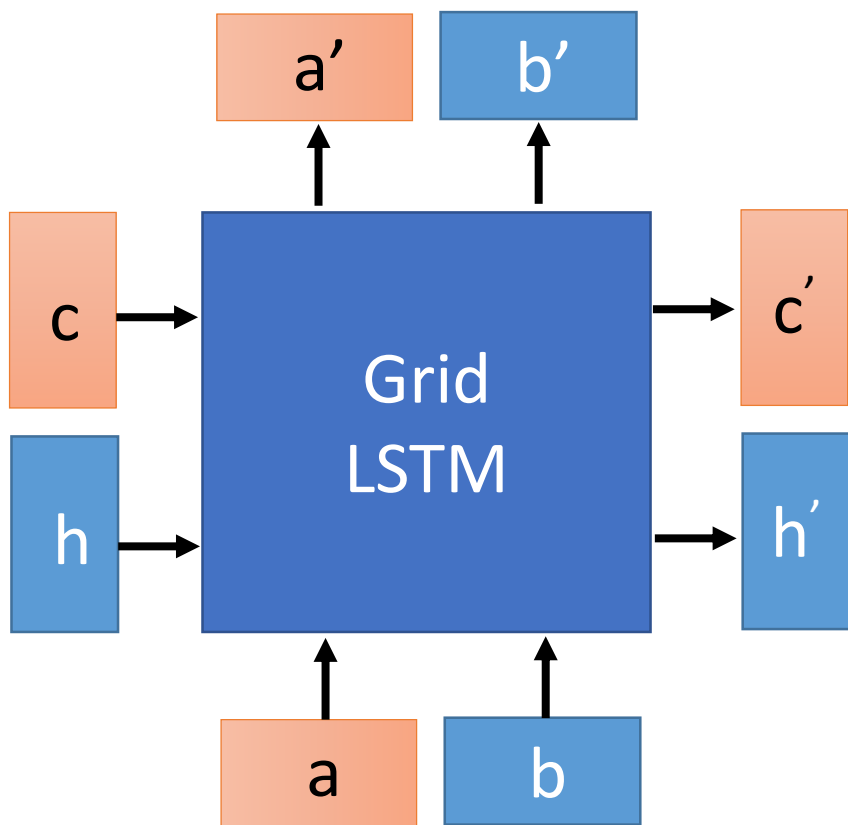


depth

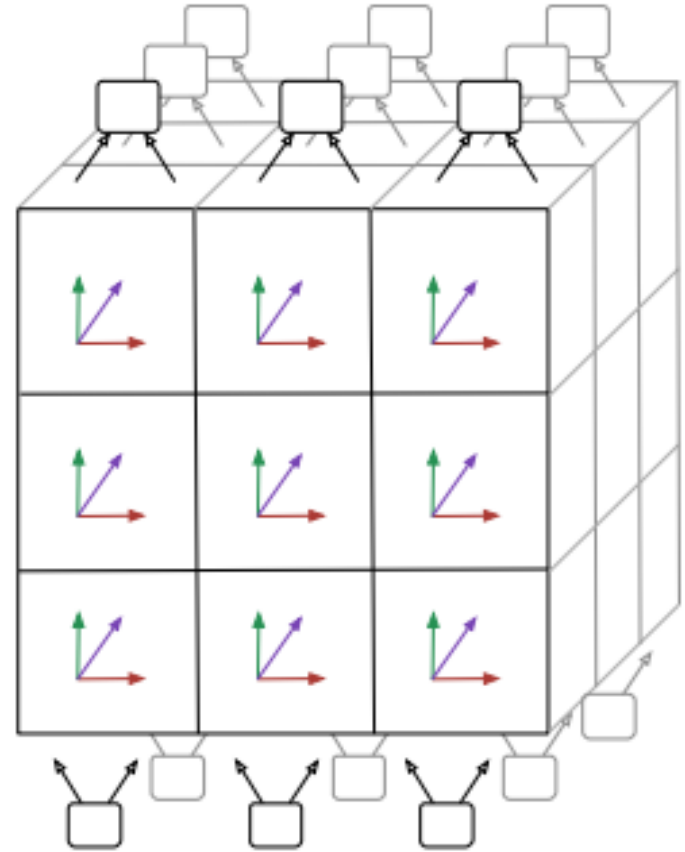
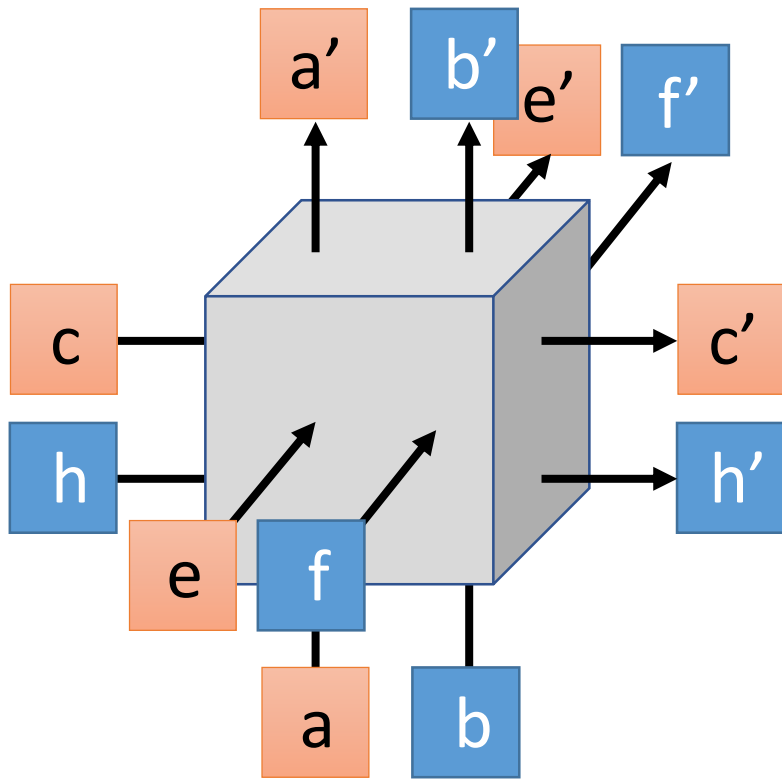




Grid LSTM

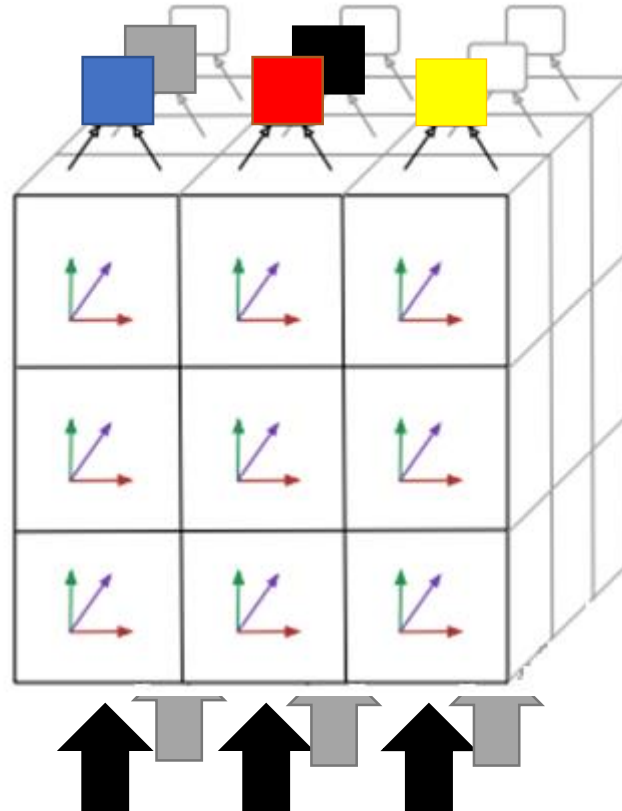
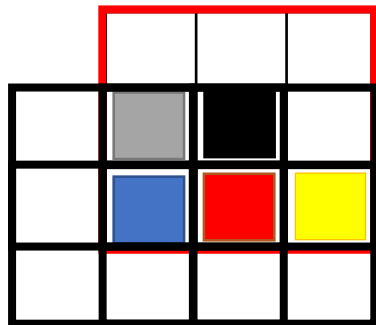
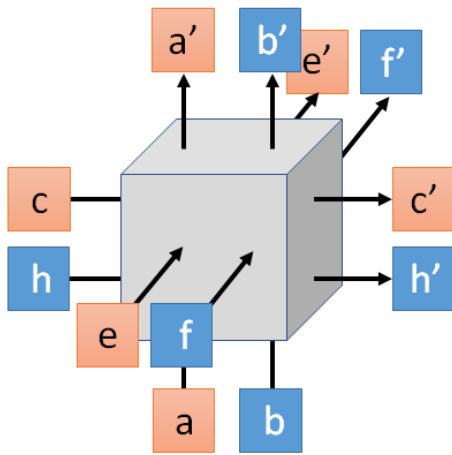


3D Grid LSTM



3D Grid LSTM

- Images are composed of pixels



3 x 3 images

