# *Deep Learning Crash Course*

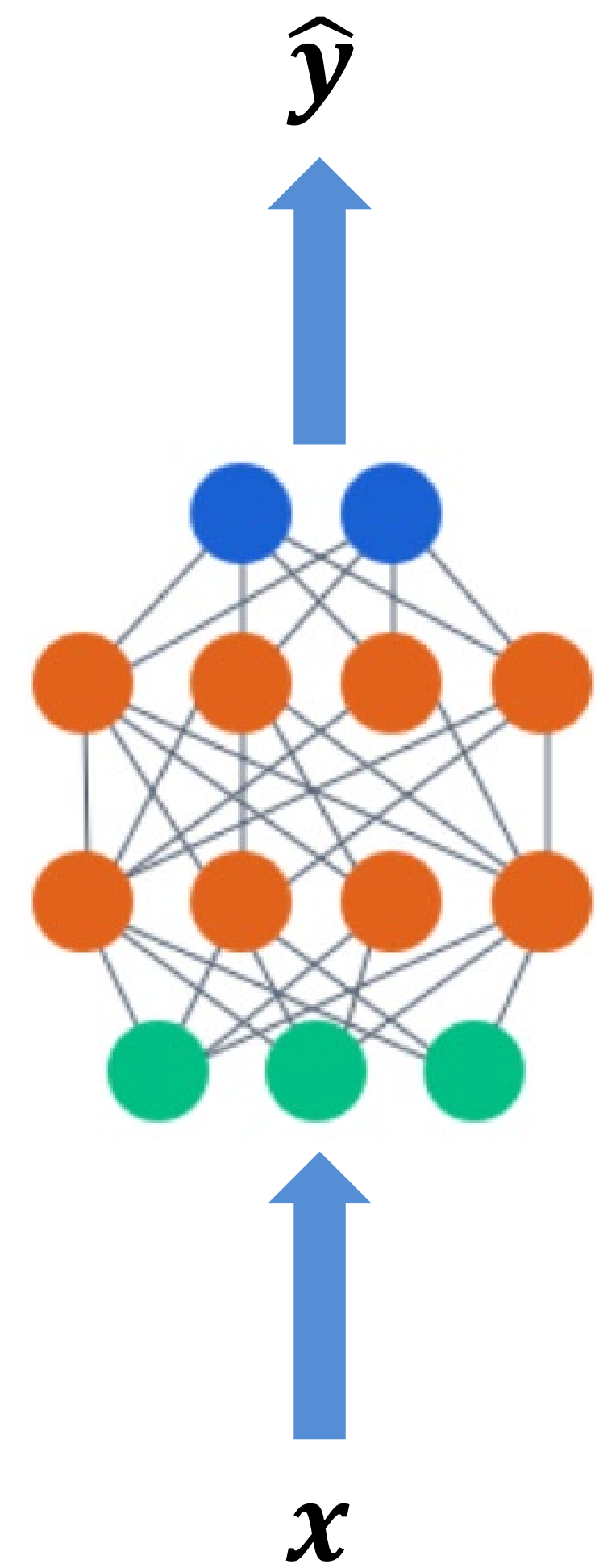**Hui Xue**

**Fall 2021**

# Outline

- Vanilla RNN and backprop through time

- Variation of RNNs

- LSTM and GRU

- Multi-layer RNN and bidirectional RNN

- Sequence pre-processing

Image Captioning

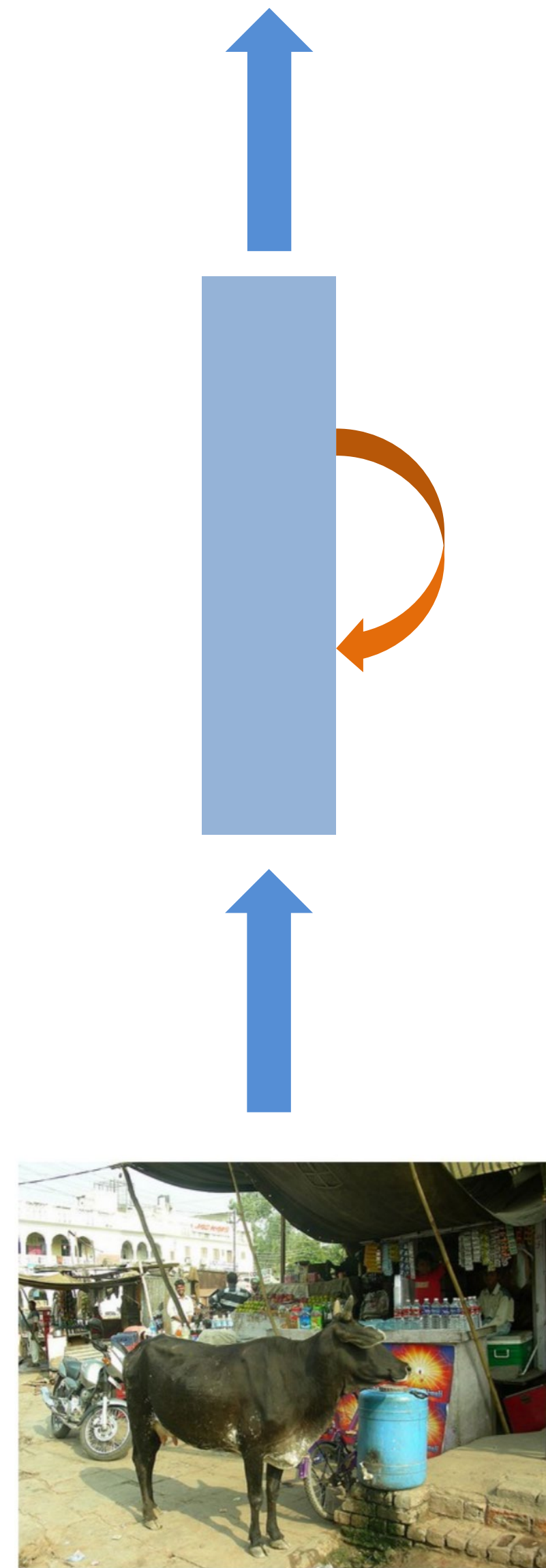$\hat{y}$



a cow is standing in the middle of a street

Sentimental classification

👎

Speech recognition

Today is nice. Let's go to the park.

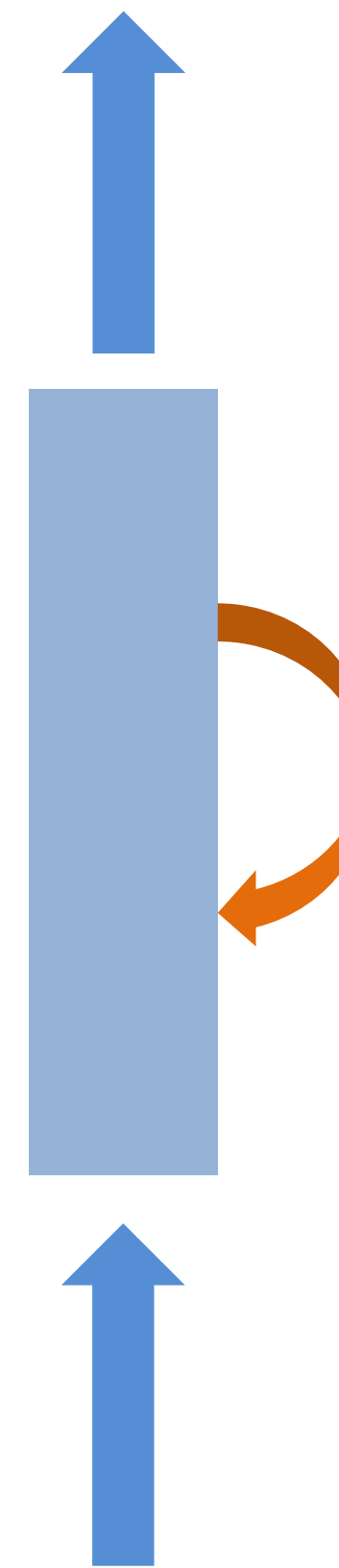Machine translation

今天天气很好。我们去公园吧。

$x$

- Image
- Video with fixed length
- Tabulate feature set

…
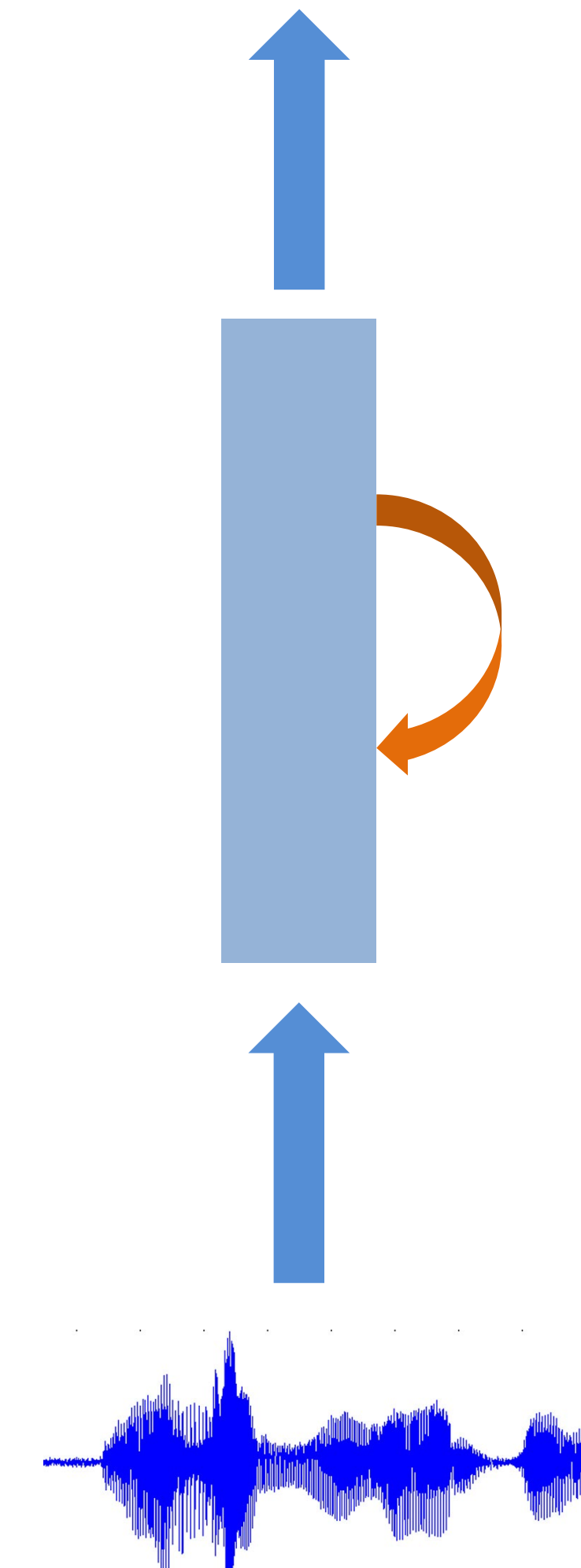
Three weeks ago, I bought this product. It worked fine at the beginning but stopped working after a week ..

Today is nice. Let's go to the park.
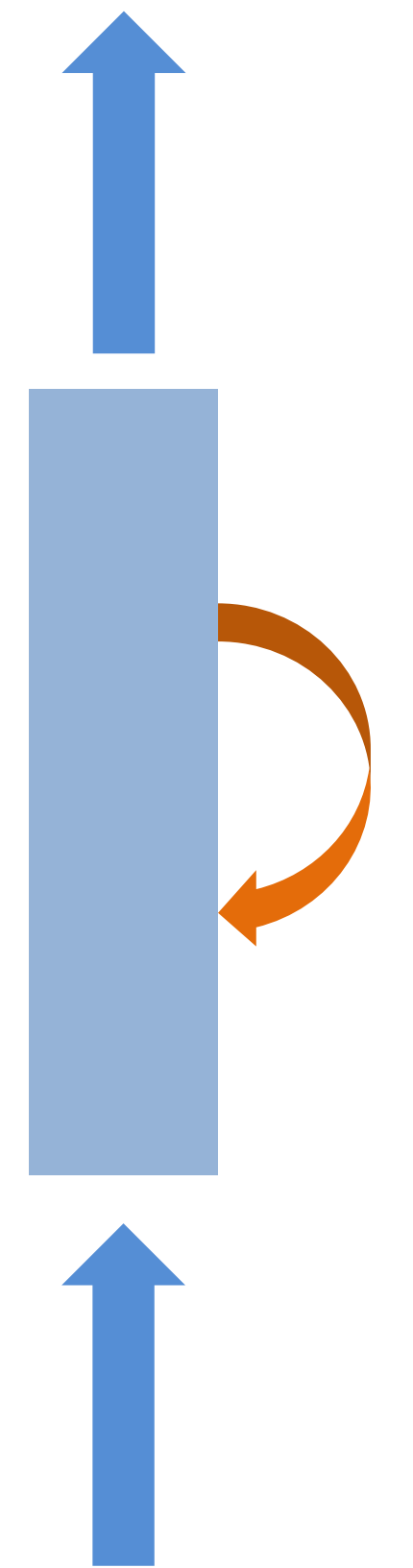
https://cs.stanford.edu/people/karpathy/sfmltalk.pdf

# Sequence model for regression or detection



Regression to estimate the number of sunspots

http://www.sidc.be/silso/dayssnplot



Detect R-wave trigger from ECG waveform

# Recurrent Neural Network

- Give a series of inputs $x^{<t>}, t = 1, \ldots, N$, a RNN is a model to receive every input and produce output $y^{<t>}$
- RNN has internal state $a^{<t>}$
- All time steps share the model parameters
- $a^{<0>}$ is often initialized as 0



unrolled along time

# Recurrent Neural Network

- input $x^{<t>}$, $t = 1, \dots, N$, output $y^{<t>}$
- internal state $a^{<t>}$

$$a^{<0>} = 0$$

$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + \textcolor{red}{W_{ax} \cdot x^{<t>}} + b_a)$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

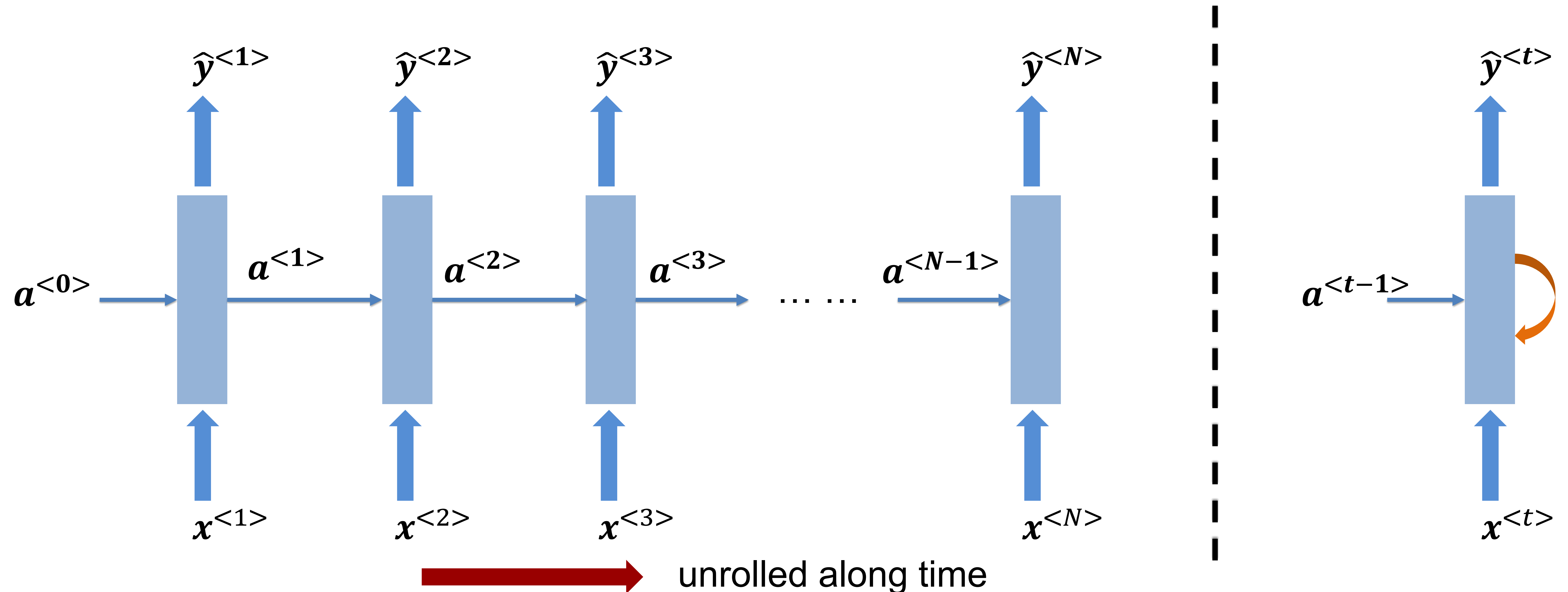$g_a$: nonlinear activation function for internal state, often tanh, ReLU

$g_y$: for output, sigmoid for binary classification, softmax for multi-class

$$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<3>} \quad \hat{y}^{<N>}$$

$$a^{<0>} \quad a^{<1>} \quad a^{<2>} \quad a^{<3>} \quad a^{<N-1>}$$

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<N>}$$

$$W_{aa}, W_{ax}, W_{ya}, b_a, b_y$$

"Elman RNN"

Elman, Jeffrey L. (1990). "Finding Structure in Time". *Cognitive Science*. **14** (2): 179–211. doi:10.1016/0364-0213(90)90002-E

6

$$\boldsymbol{Jordan\ RNN}: a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + \textcolor{red}{W_{ay} \cdot y^{<t-1>}} + b_a)$$

Jordan, Michael I. (1997-01-01). "Serial Order: A Parallel Distributed Processing Approach". *Neural-Network Models of Cognition - Biobehavioral Foundations. Advances in Psychology*. Neural-Network Models of Cognition. **121**. pp. 471–495.

# Recurrent Neural Network: Forward pass

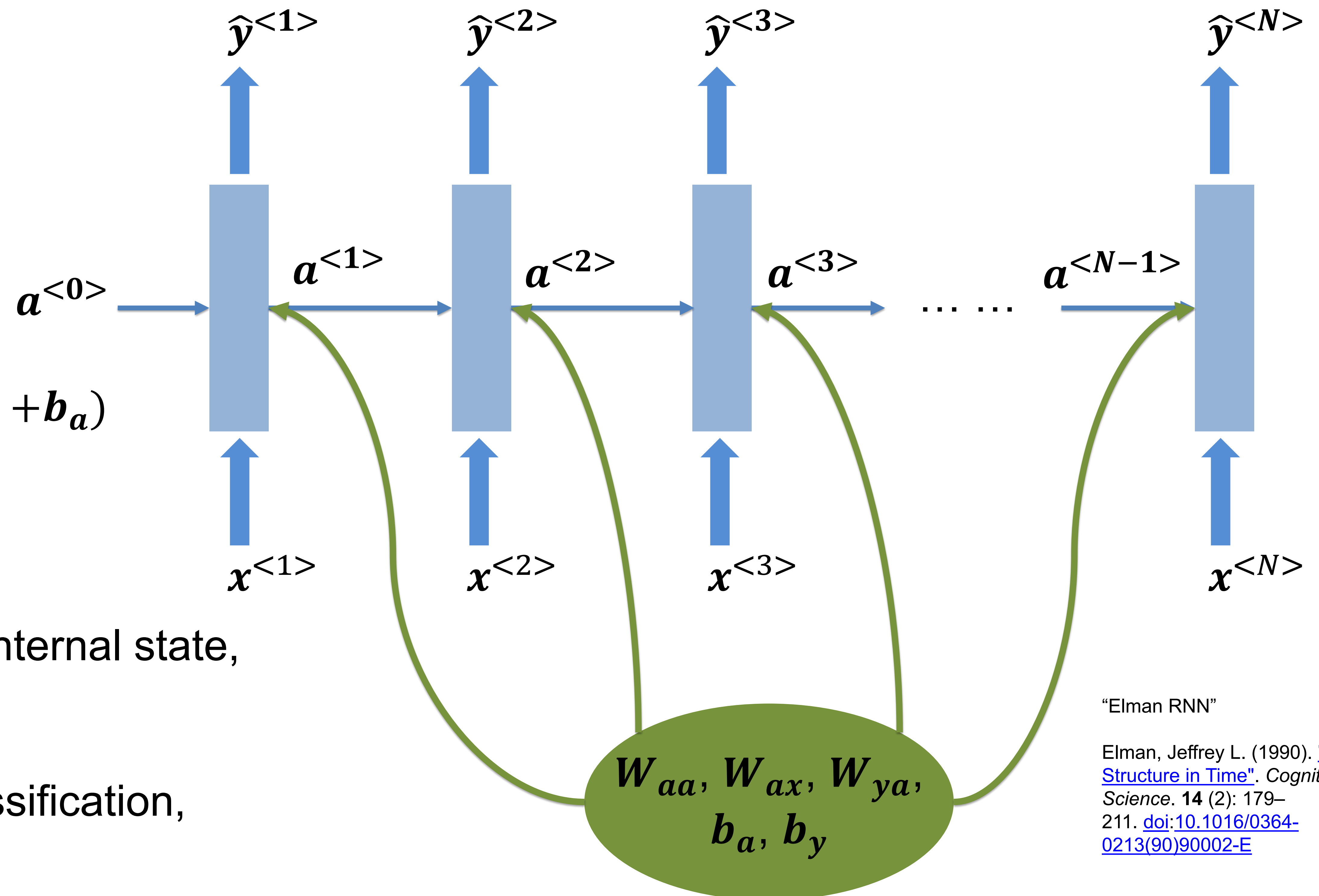- input $x^{<t>}, t = 1, ..., N$, output $y^{<t>}$
- internal state $a^{<t>}$

$a^{<0>} = 0$

$a^{<1>} = g_a(W_{aa} \cdot a^{<0>} + W_{ax} \cdot x^{<1>} + b_a)$

$y^{<1>} = g_y(W_{ya} \cdot a^{<1>} + b_y)$

$a^{<2>} = g_a(W_{aa} \cdot a^{<1>} + W_{ax} \cdot x^{<2>} + b_a)$

$y^{<2>} = g_y(W_{ya} \cdot a^{<2>} + b_y)$

$a^{<3>} = g_a(W_{aa} \cdot a^{<2>} + W_{ax} \cdot x^{<3>} + b_a)$

$y^{<3>} = g_y(W_{ya} \cdot a^{<3>} + b_y)$

$a^{<4>} = g_a(W_{aa} \cdot a^{<3>} + W_{ax} \cdot x^{<4>} + b_a)$

$y^{<4>} = g_y(W_{ya} \cdot a^{<4>} + b_y)$

- input $x^{<t>}, t = 1, \ldots, N$, output $y^{<t>}$
- internal state $a^{<t>}$

$$a^{<0>} = 0$$

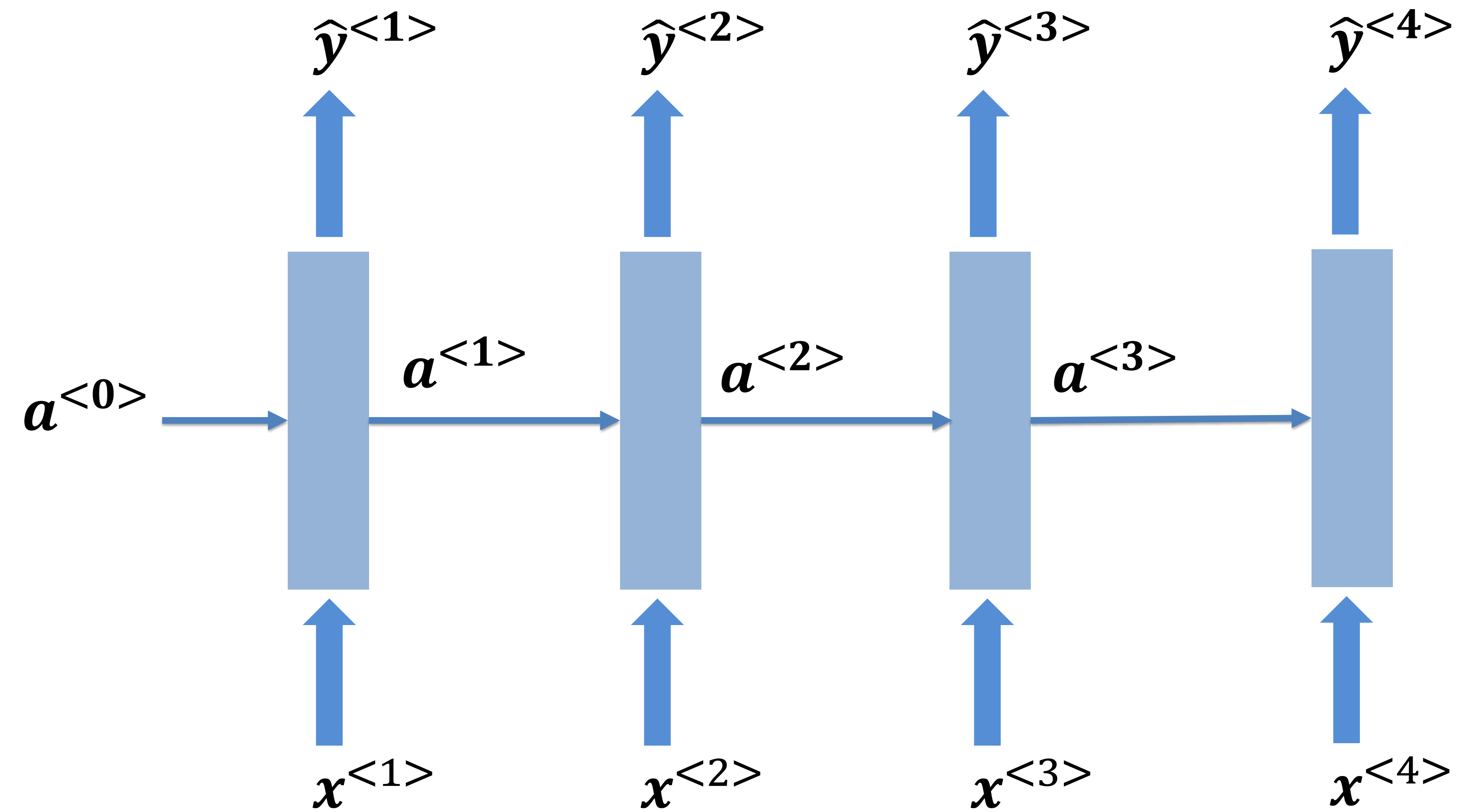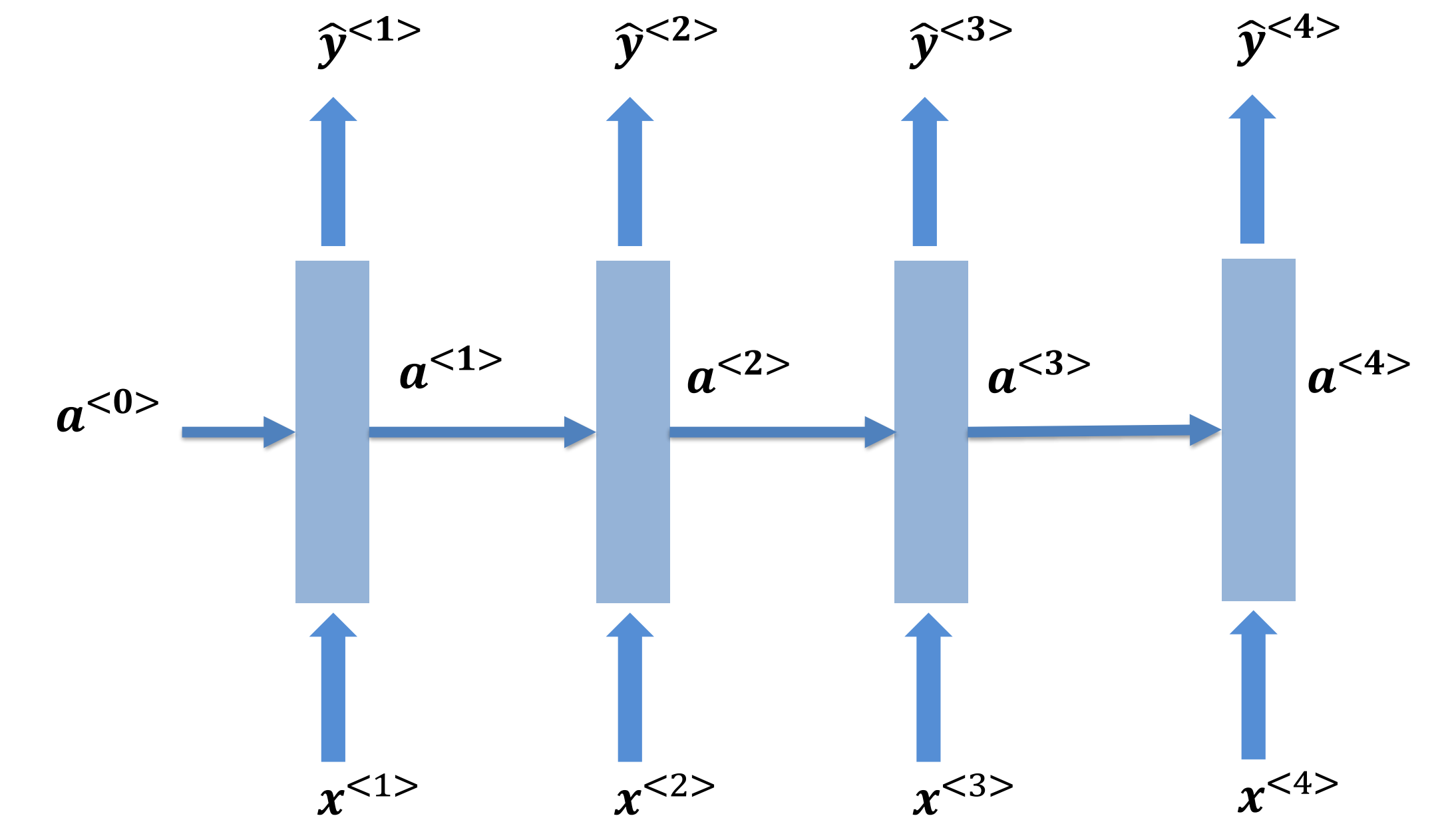$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$a^{<4>} = g_a(W_{aa} \cdot a^{<3>} + W_{ax} \cdot x^{<4>} + b_a)$$
$$= g_a(W_{aa} \cdot [g_a(W_{aa} \cdot a^{<2>} + W_{ax} \cdot x^{<3>} + b_a)] + W_{ax} \cdot x^{<4>} + b_a)$$
$$= g_a(W_{aa} \cdot [g_a(W_{aa} \cdot [g_a(W_{aa} \cdot a^{<1>} + W_{ax} \cdot x^{<2>} + b_a)] + W_{ax} \cdot x^{<3>} + b_a)] + W_{ax} \cdot x^{<4>} + b_a)$$
$$= g_a(W_{aa} \cdot [g_a(W_{aa} \cdot [g_a(W_{aa} \cdot [g_a(W_{aa} \cdot a^{<0>} + W_{ax} \cdot x^{<1>} + b_a)] + W_{ax} \cdot x^{<2>} + b_a)] + W_{ax} \cdot x^{<3>} + b_a)] + W_{ax} \cdot x^{<4>} + b_a)$$

$$y^{<4>} = g_y(W_{ya} \cdot a^{<4>} + b_y)$$

- Recurrently apply the model parameters to update internal state
- $y^{<t>}$ depends on all inputs on or before current time step

# Recurrent Neural Network: loss

- input $x^{<t>}, t = 1, \dots, N$, output $y^{<t>}$
- internal state $a^{<t>}$

Assume $y^{<t>}$ to take 0 or 1, e.g. like/dislike, binary classification

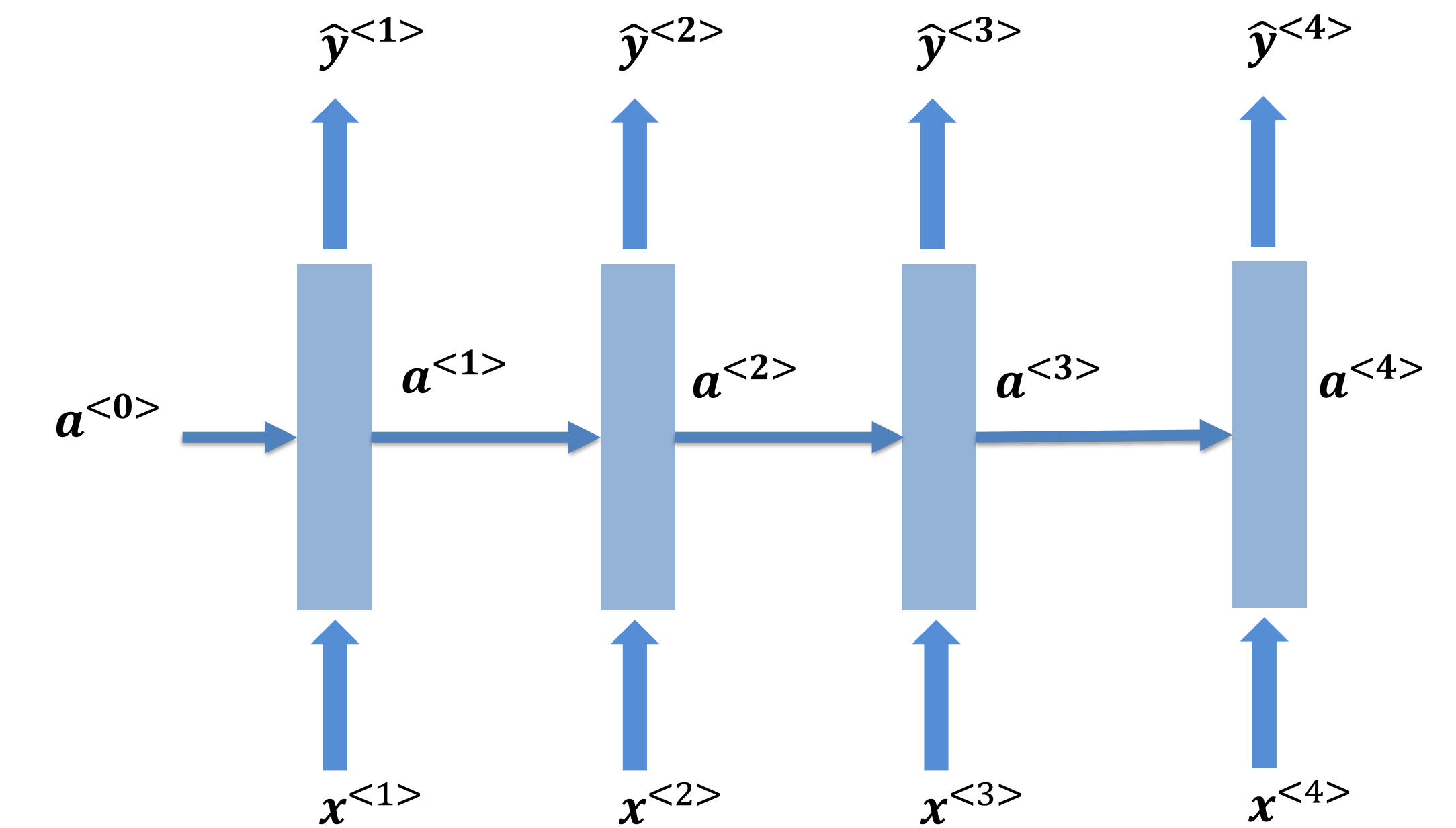$$\ell^{<t>}(y^{<t>}, \hat{y}^{<t>}) = -y^{<t>} log(\hat{y}^{<t>}) - (1 - y^{<t>}) log(1 - \hat{y}^{<t>})$$

$$\ell = \sum_{t=1}^{N} \ell^{<t>}(y^{<t>}, \hat{y}^{<t>})$$

$$\frac{\partial \ell}{\partial W_{aa}} = \sum_{t=1}^{N} \frac{\partial \ell^{<t>}(y^{<t>}, \hat{y}^{<t>})}{W_{aa}} \qquad \frac{\partial \ell}{\partial W_{ax}} = \sum_{t=1}^{N} \frac{\partial \ell^{<t>}(y^{<t>}, \hat{y}^{<t>})}{W_{ax}}$$

$$\frac{\partial \ell}{\partial W_{ya}} = \sum_{t=1}^{N} \frac{\partial \ell^{<t>}(y^{<t>}, \hat{y}^{<t>})}{W_{ya}} \qquad \dots \dots$$



$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

$$\frac{\partial \ell}{\partial W_{aa}} = \sum_{t=1}^{N} \frac{\partial \ell^{<t>}(y^{<t>}, \hat{y}^{<t>})}{W_{aa}}$$

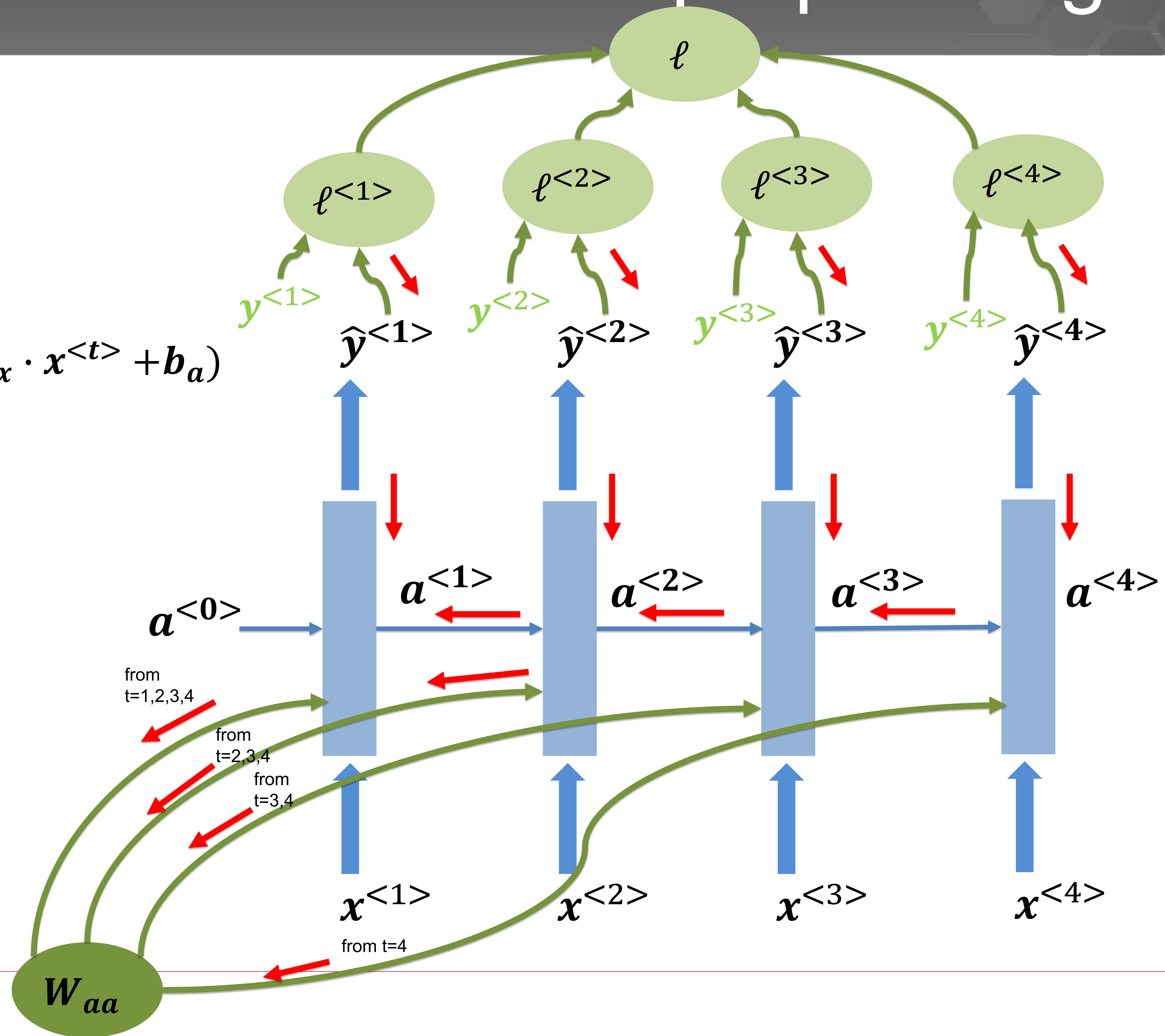$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$



10

# Recurrent Neural Network: variants
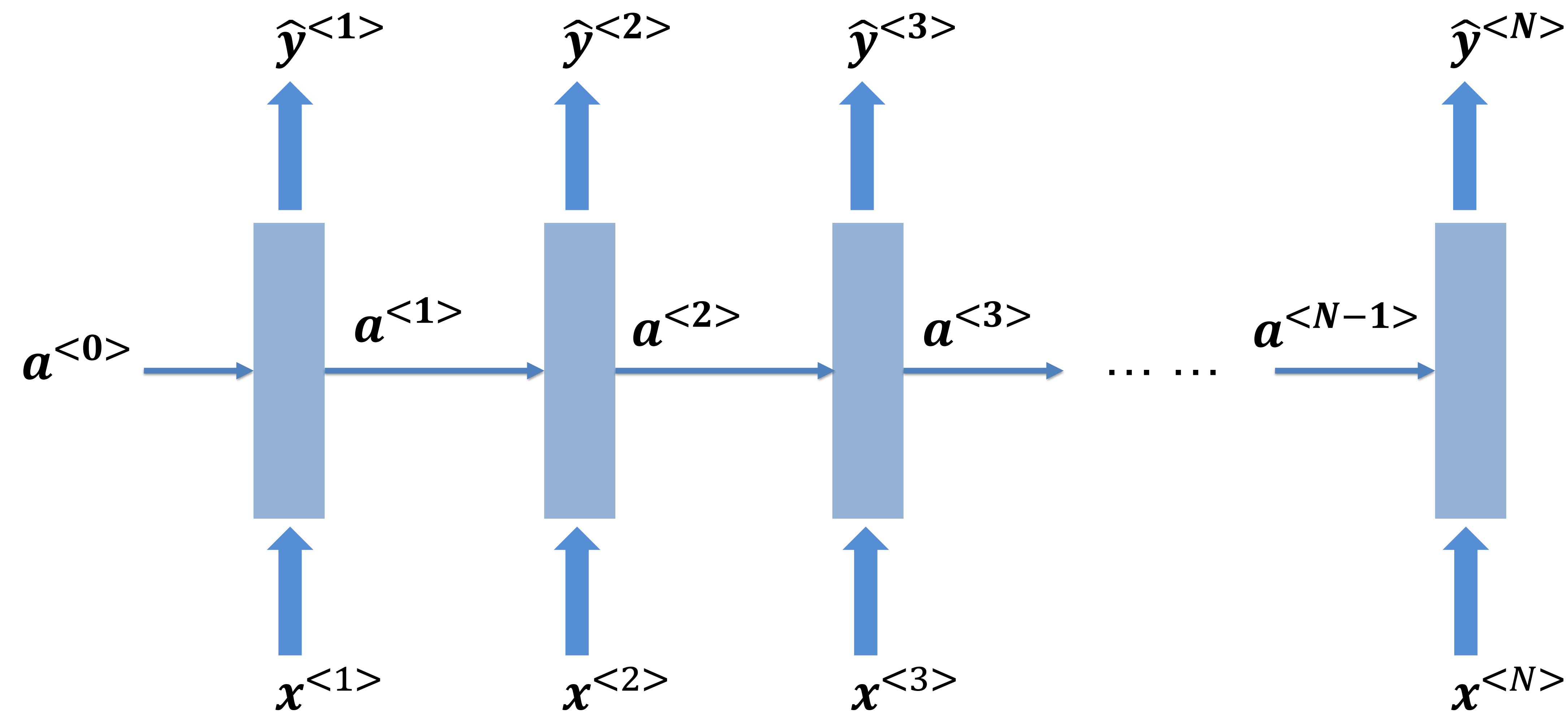
$$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<3>} \qquad\qquad \hat{y}^{<N>}$$

$a^{<0>}$ ⟶ [ ] $a^{<1>}$ ⟶ [ ] $a^{<2>}$ ⟶ [ ] $a^{<3>}$ ⟶ … … $a^{<N-1>}$ ⟶ [ ]

Many-to-many, same length

e.g. trigger word detection

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \qquad\qquad x^{<N>}$$

# Recurrent Neural Network: variants



$$\hat{y}$$

$a^{<0>} \rightarrow \boxed{\phantom{a}} \xrightarrow{a^{<1>}} \boxed{\phantom{a}} \xrightarrow{a^{<2>}} \boxed{\phantom{a}} \xrightarrow{a^{<3>}} \cdots \cdots \xrightarrow{a^{<N-1>}} \boxed{\phantom{a}}$

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad\quad\quad x^{<N>}$

Many-to-one, causal

e.g. video grading system

# Recurrent Neural Network: variants

Global pooling $\rightarrow \hat{\boldsymbol{y}}$

$\hat{\boldsymbol{y}}^{<1>}$ $\hat{\boldsymbol{y}}^{<2>}$ $\hat{\boldsymbol{y}}^{<3>}$ … … $\hat{\boldsymbol{y}}^{<N>}$

$\boldsymbol{a}^{<0>}$ $\boldsymbol{a}^{<1>}$ $\boldsymbol{a}^{<2>}$ $\boldsymbol{a}^{<3>}$ … … $\boldsymbol{a}^{<N-1>}$

$\boldsymbol{x}^{<1>}$ $\boldsymbol{x}^{<2>}$ $\boldsymbol{x}^{<3>}$ $\boldsymbol{x}^{<N>}$

Many-to-one, non-causal

e.g. video grading system

Better balanced for all time steps

$$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<3>} \quad \dots \dots \quad \hat{y}^{<N>}$$

$$a^{<0>} \quad a^{<1>} \quad a^{<2>} \quad a^{<3>} \quad \dots \dots \quad a^{<N-1>}$$
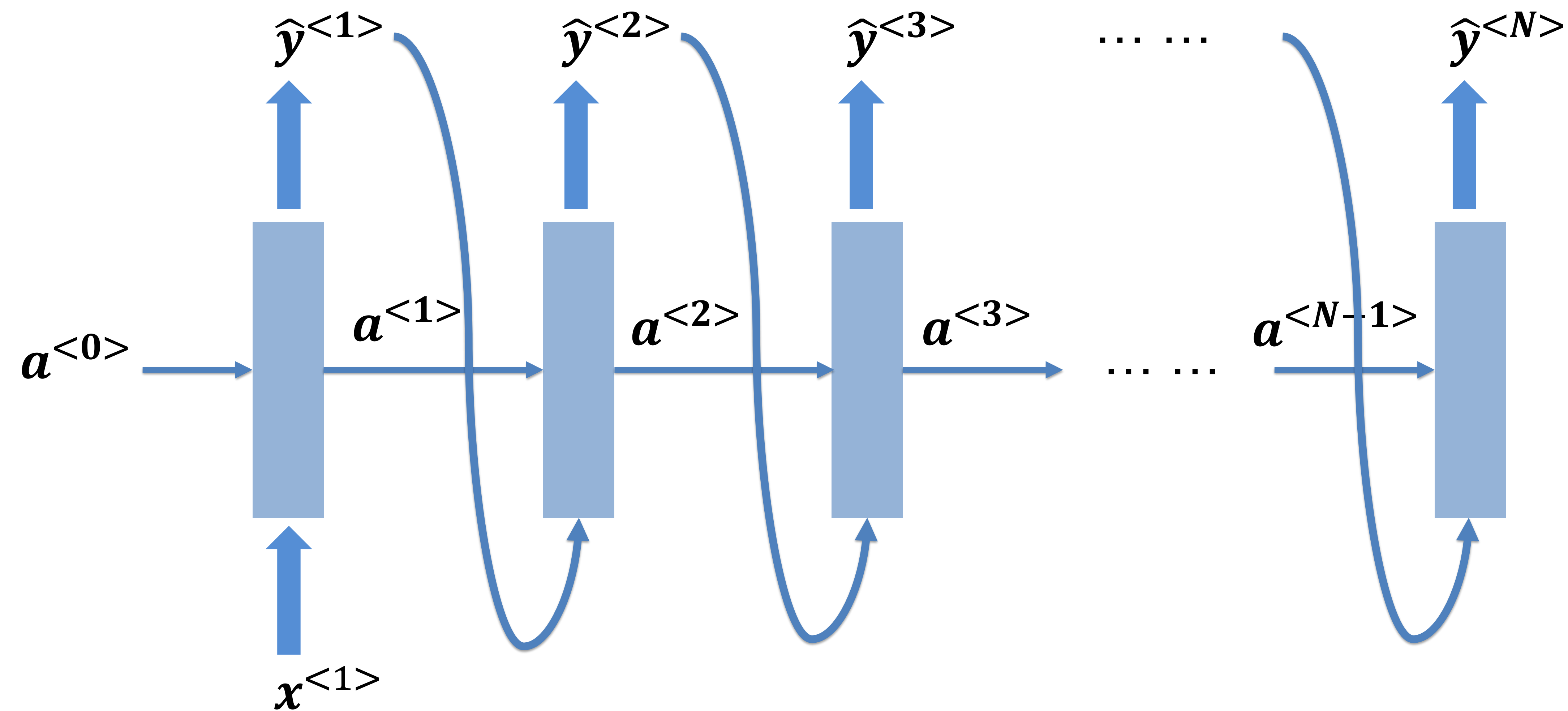
$$x^{<1>}$$

One-to-many

e.g. music or video generator, with input being a latent variable

Self-regression: Sample from distribution $\hat{y}^{<t>}$ and feed it as the next $x^{<t>}$, until <EOS> token is reached   (sample from a given distribution: np.random.choice(p=prob) )
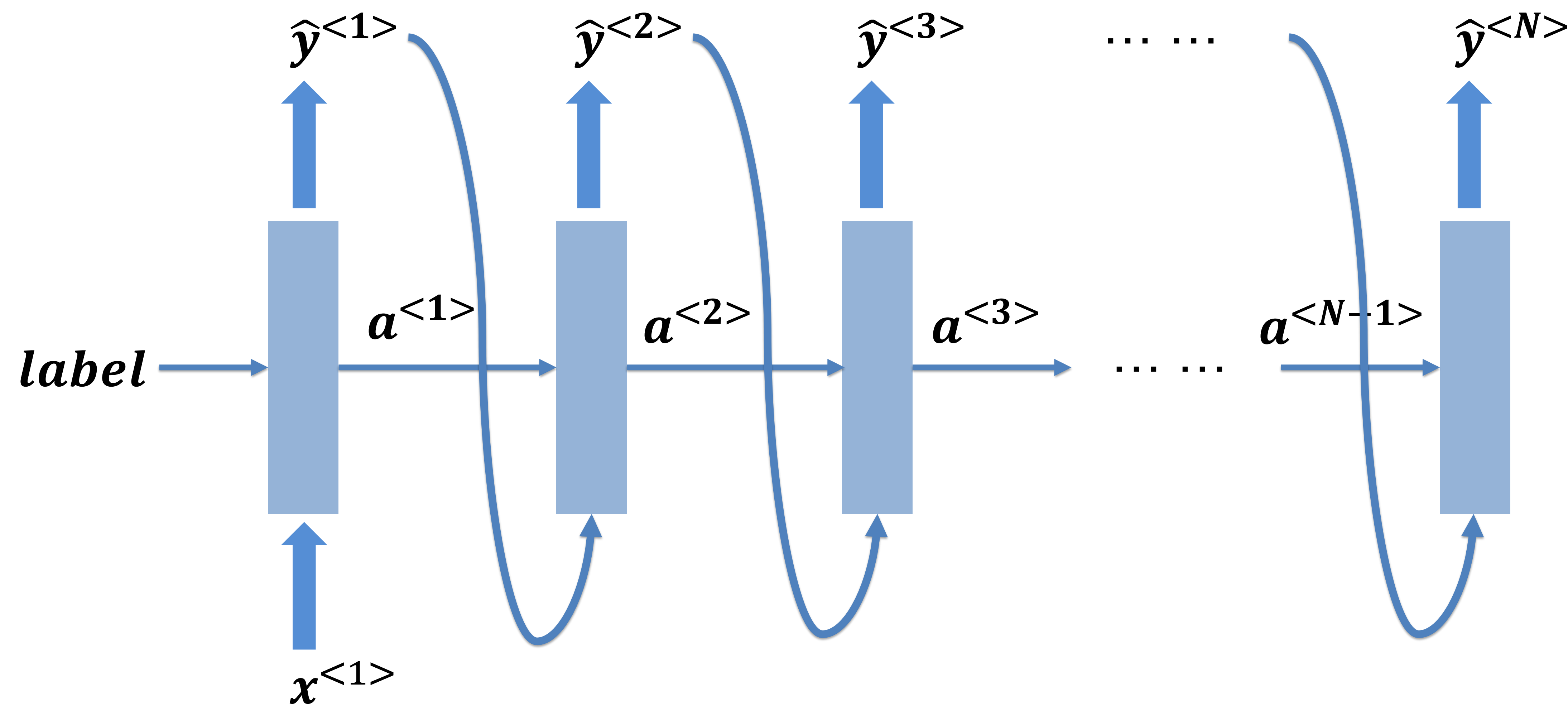
$$P(\hat{y}^{<1>}, \hat{y}^{<2>}, \hat{y}^{<3>}, ..., \hat{y}^{<N>})$$

$$= P\left(\hat{y}^{<1>} \middle| x^{<1>}\right) P\left(\hat{y}^{<2>} \middle| \hat{y}^{<1>}\right) P\left(\hat{y}^{<3>} \middle| \hat{y}^{<1>}, \hat{y}^{<2>}\right) ... P\left(\hat{y}^{<N>} \middle| \hat{y}^{<1>}, ..., \hat{y}^{<N-1>}\right)$$

Chain rule of joint probability

$$\hat{y}^{<1>} \qquad \hat{y}^{<2>} \qquad \hat{y}^{<3>} \qquad \dots \dots \qquad \hat{y}^{<N>}$$

$$label \qquad a^{<1>} \qquad a^{<2>} \qquad a^{<3>} \qquad \dots \dots \qquad a^{<N-1>}$$

$$x^{<1>}$$

One-to-many

e.g. music, with input being a latent variable and label being its style (classical, pop, …)

encoder

decoder

$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  ... ...  $\hat{y}^{<M>}$

$a^{<0>}$  $a^{<1>}$  $a^{<2>}$  $a^{<3>}$  $a^{<N-1>}$ ...  $a^{<N>}$

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<N>}$

Many-to-many, variable length    e.g. machine translation, video captioning

$$\frac{\partial \ell^{<N>}}{\partial W_{aa}}\bigg|_{t=1} = \frac{\partial \ell^{<N>}}{\partial \hat{y}^{<N>}} \frac{\partial \hat{y}^{<N>}}{\partial a^{<N>}} \frac{\partial a^{<N>}}{\partial a^{<N-1>}} \cdots \frac{\partial a^{<2>}}{\partial a^{<1>}} \frac{\partial a^{<1>}}{\partial W_{aa}}$$

$$a^{<N>} = g_a(W_{aa} \cdot a^{<N-1>} + W_{ax} \cdot x^{<N>} + b_a)$$

$$\frac{\partial \ell^{<N>}}{\partial W_{aa}} |_{t=1} = \frac{\partial \ell^{<N>}}{\partial \hat{y}^{<N>}} \frac{\partial \hat{y}^{<N>}}{\partial a^{<N>}} \frac{\partial a^{<N>}}{\partial a^{<N-1>}} \cdots \frac{\partial a^{<2>}}{\partial a^{<1>}} \frac{\partial a^{<1>}}{\partial W_{aa}}$$

$$a^{<N>} = g_a(W_{aa} \cdot a^{<N-1>} + W_{ax} \cdot x^{<N>} + b_a) \qquad \frac{\partial a^{<N>}}{\partial a^{<N-1>}} = W_{aa} \cdot g'_a (W_{aa} \cdot a^{<N-1>} + W_{ax} \cdot x^{<N>} + b_a)$$

$$\frac{\partial \ell^{<N>}}{\partial W_{aa}} |_{t=1} = \frac{\partial \ell^{<N>}}{\partial \hat{y}^{<N>}} \frac{\partial \hat{y}^{<N>}}{\partial a^{<N>}} \left[ \prod_{t=2}^{N} \frac{\partial a^{<t>}}{\partial a^{<t-1>}} \right] \frac{\partial a^{<1>}}{\partial W_{aa}}$$

$$= \frac{\partial \ell^{<N>}}{\partial \hat{y}^{<N>}} \frac{\partial \hat{y}^{<N>}}{\partial a^{<N>}} W_{aa}{}^{N-1} \left[ \prod_{t=2}^{N} g'_a (W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a) \right] \frac{\partial a^{<1>}}{\partial W_{aa}}$$
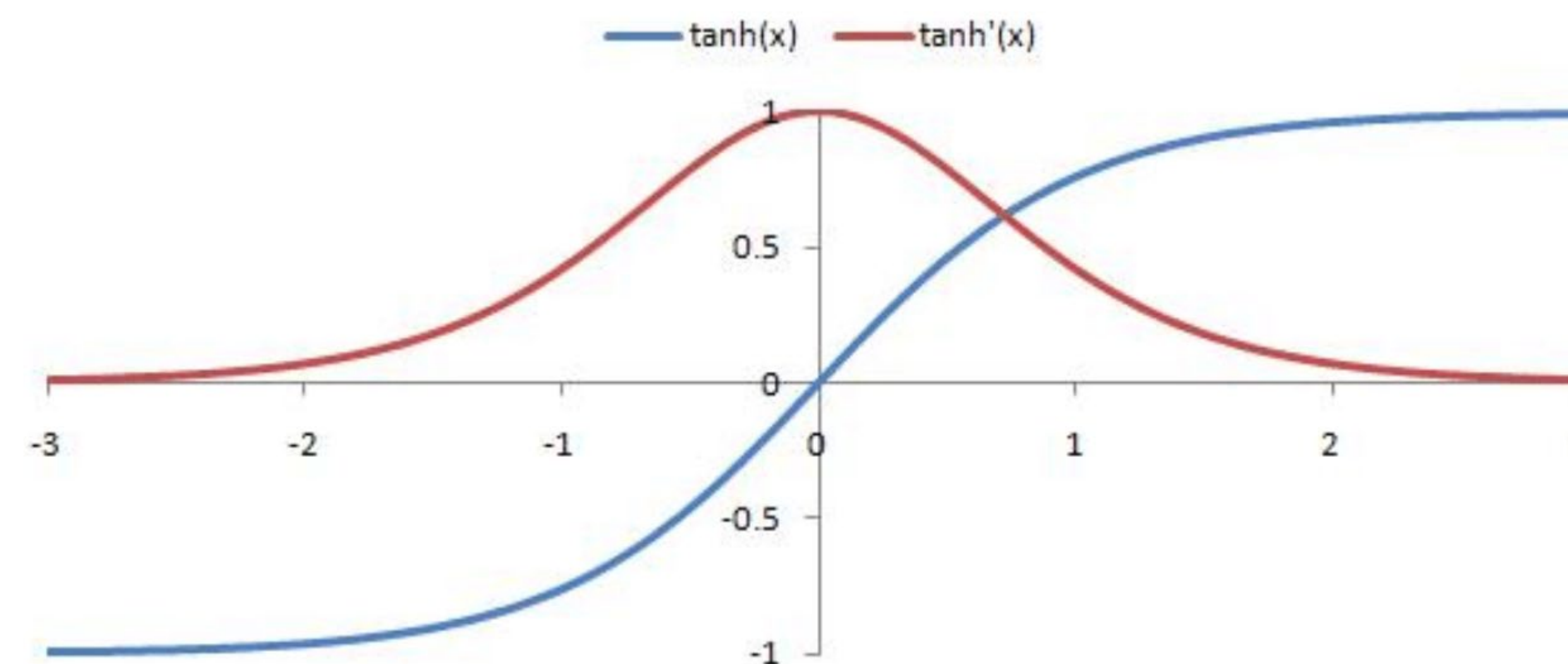
$$\frac{\partial \ell^{<N>}}{\partial W_{aa}}\Big|_{t=1} = \frac{\partial \ell^{<N>}}{\partial \hat{y}^{<N>}} \frac{\partial \hat{y}^{<N>}}{\partial a^{<N>}} W_{aa}^{N-1} \left[ \prod_{t=2}^{N} g_a' \left( W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a \right) \right] \frac{\partial a^{<1>}}{\partial W_{aa}}$$

If falling into the linear regime of tanh, then minimal eigen value of $W_{aa}$ comes into the play :
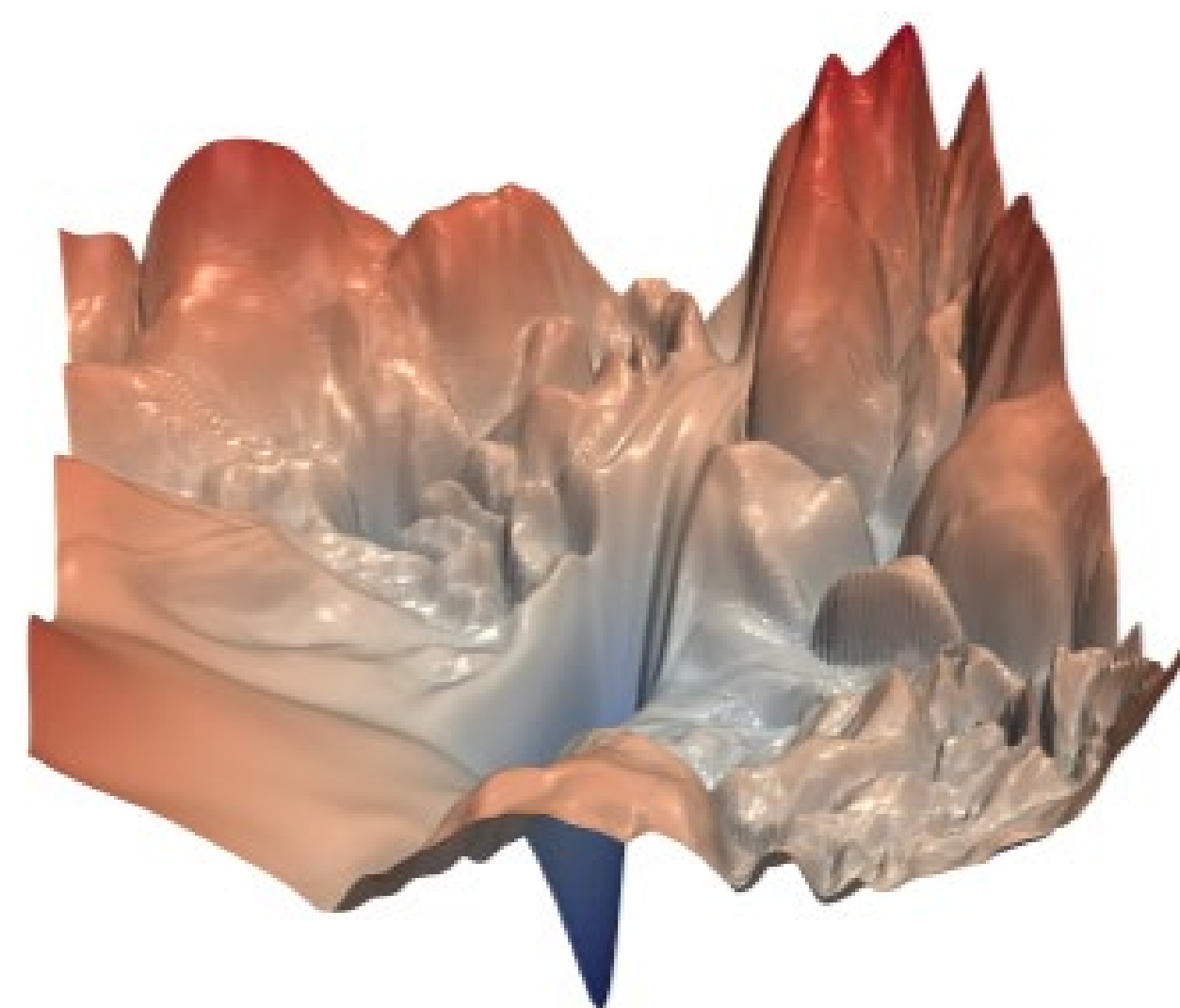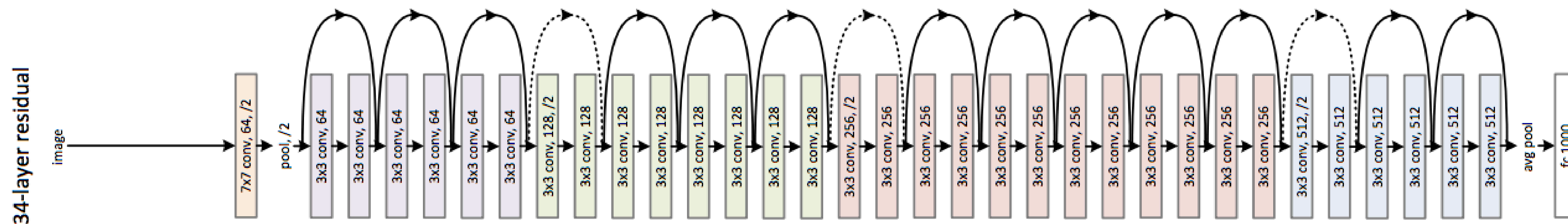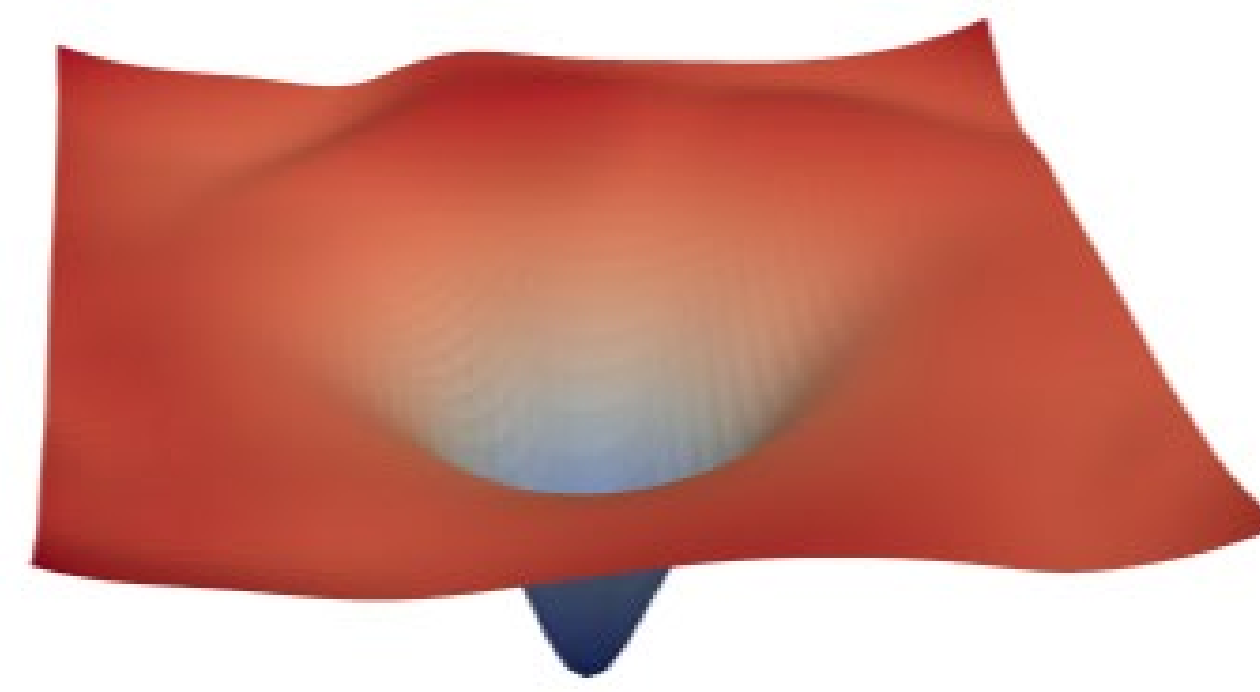
<1, vanishing gradient

>1 exploding gradient

Often cause vanishing gradient

# Fix the unstable gradient: skip connections



Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.
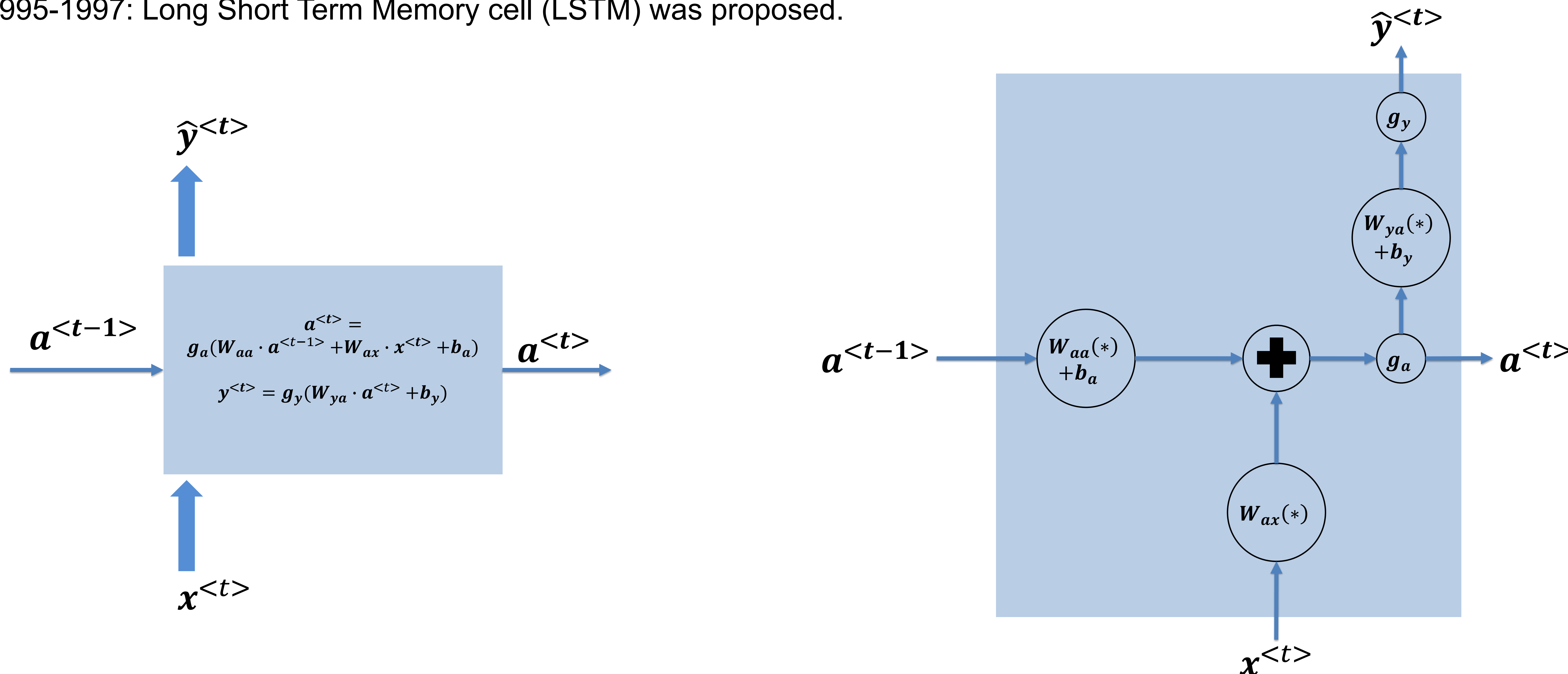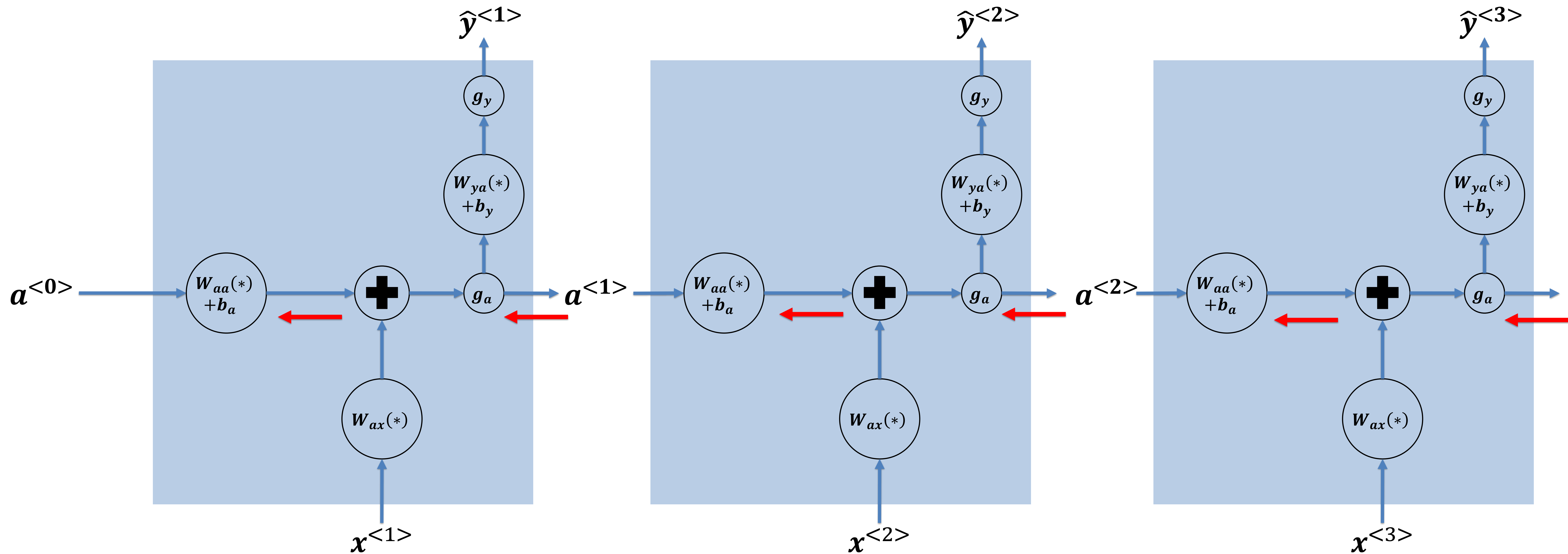
- Add skip-connections
- Allow gradient flow to pass back to earlier layers
- Network will learn to utilize these skip connections during training

Deep Residual Learning for Image Recognition. 2015. https://arxiv.org/abs/1512.03385
Visualizing the Loss Landscape of Neural Nets. https://arxiv.org/pdf/1712.09913.pdf

# Invented much earlier for RNN

1995-1997: Long Short Term Memory cell (LSTM) was proposed.



$$\hat{y}^{<t>}$$

$$a^{<t-1>}$$

$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

$$a^{<t>}$$

$$x^{<t>}$$

$$\hat{y}^{<t>}$$

$$g_y$$

$$W_{ya}(*) + b_y$$

$$a^{<t-1>}$$

$$W_{aa}(*) + b_a$$

$$g_a$$

$$a^{<t>}$$

$$W_{ax}(*)$$

$$x^{<t>}$$

LSTM can Solve Hard Long Time Lag Problems. Advances in Neural Information Processing Systems 9. 1997.

$$a^{<t>} = g_a(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

## Vanilla RNN

Introduce a memory cell $c^{<t>}$

1) Compute candidate $\tilde{c}^{<t>}$ to update memory cell, with learnable parameters to control what to remember, and what to forget

$$\tilde{c}^{<t>} = g_c(W_{ca} \cdot a^{<t-1>} + W_{cx} \cdot x^{<t>} + b_c)$$

2) Add two gates to control whether to update memory cell or keep its previous states

$$c^{<t>} = U \circ \tilde{c}^{<t>} + F \circ c^{<t-1>}$$

$$F = \sigma(W_{fa} \cdot a^{<t-1>} + W_{fx} \cdot x^{<t>} + b_f): \text{forget gate}$$

$$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + b_u): \text{update gate}$$

3) Add one output gate

$$a^{<t>} = O \circ c^{<t>}$$

$$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + b_o): \text{output gate}$$

4) Compute output

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

or

$$y^{<t>} = g_y(a^{<t>})$$

## LSTM

$$\tilde{c}^{<t>} = g_c(W_{ca} \cdot a^{<t-1>} + W_{cx} \cdot x^{<t>} + b_c)$$

$$c^{<t>} = U \circ \tilde{c}^{<t>} + F \circ c^{<t-1>}$$

$$F = \sigma(W_{fa} \cdot a^{<t-1>} + W_{fx} \cdot x^{<t>} + b_f): \text{forget gate}$$

$$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + b_u): \text{update gate}$$

$$a^{<t>} = O \circ c^{<t>}$$

$$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + b_o): \text{output gate}$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

$F$, $U$, $O$ are vectors, same size as memory cell
$U \circ \tilde{c}^{<t>}$, $F \circ c^{<t-1>}$, $O \circ c^{<t>}$ : element-wise multiplication

## LSTM

$$\tilde{c}^{<t>} = g_c(W_{ca} \cdot a^{<t-1>} + W_{cx} \cdot x^{<t>} + b_c)$$

$$c^{<t>} = U \circ \tilde{c}^{<t>} + F \circ c^{<t-1>}$$

$$F = \sigma(W_{fa} \cdot a^{<t-1>} + W_{fx} \cdot x^{<t>} + b_f): \text{forget gate}$$

$$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + b_u): \text{update gate}$$

$$a^{<t>} = O \circ c^{<t>}$$

$$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + b_o): \text{output gate}$$

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

LSTM can Solve Hard Long Time Lag Problems. Advances in Neural Information Processing Systems 9. 1997.

26

# LSTM: peephole connection

Add the memory cell state to the computation of gates

$$F = \sigma(W_{fa} \cdot a^{<t-1>} + W_{fx} \cdot x^{<t>} + b_f): \text{forget gate}$$

$$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + b_u): \text{update gate}$$

$$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + b_o): \text{output gate}$$

$$F = \sigma(W_{fa} \cdot a^{<t-1>} + W_{fx} \cdot x^{<t>} + \textcolor{red}{W_{fc} \cdot c^{<t-1>}} + b_f): \text{forget gate}$$

$$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + \textcolor{red}{W_{uc} \cdot c^{<t-1>}} + b_u): \text{update gate}$$

$$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + \textcolor{red}{W_{oc} \cdot c^{<t-1>}} + b_o): \text{output gate}$$

27

Replace linear matrix multiplication with the convolution and flatten/reshape

$$F = \sigma(W_{fa} * a^{<t-1>} + W_{fx} * x^{<t>} + W_{fc} \cdot c^{<t-1>} + b_f)\text{: forget gate}$$

$$U = \sigma(W_{ua} * a^{<t-1>} + W_{ux} * x^{<t>} + W_{uc} \cdot c^{<t-1>} + b_u)\text{: update gate}$$

$$O = \sigma(W_{oa} * a^{<t-1>} + W_{ox} * x^{<t>} + W_{oc} \cdot c^{<t-1>} + b_o)\text{: output gate}$$

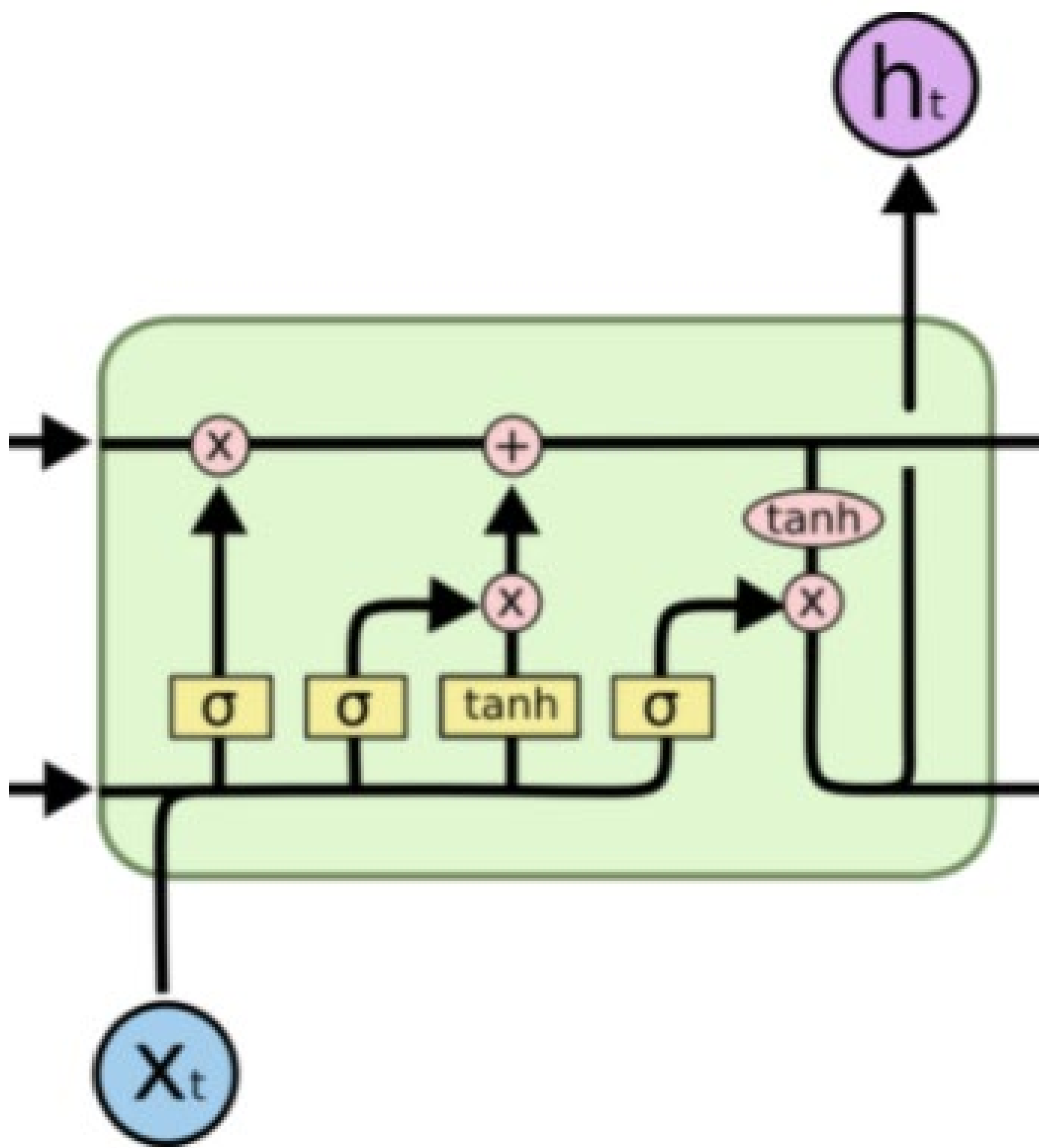$$\tilde{c}^{<t>} = g_c(W_{ca} * a^{<t-1>} + W_{cx} * x^{<t>} + b_c)$$

$$c^{<t>} = U \circ \tilde{c}^{<t>} + F \circ c^{<t-1>}$$
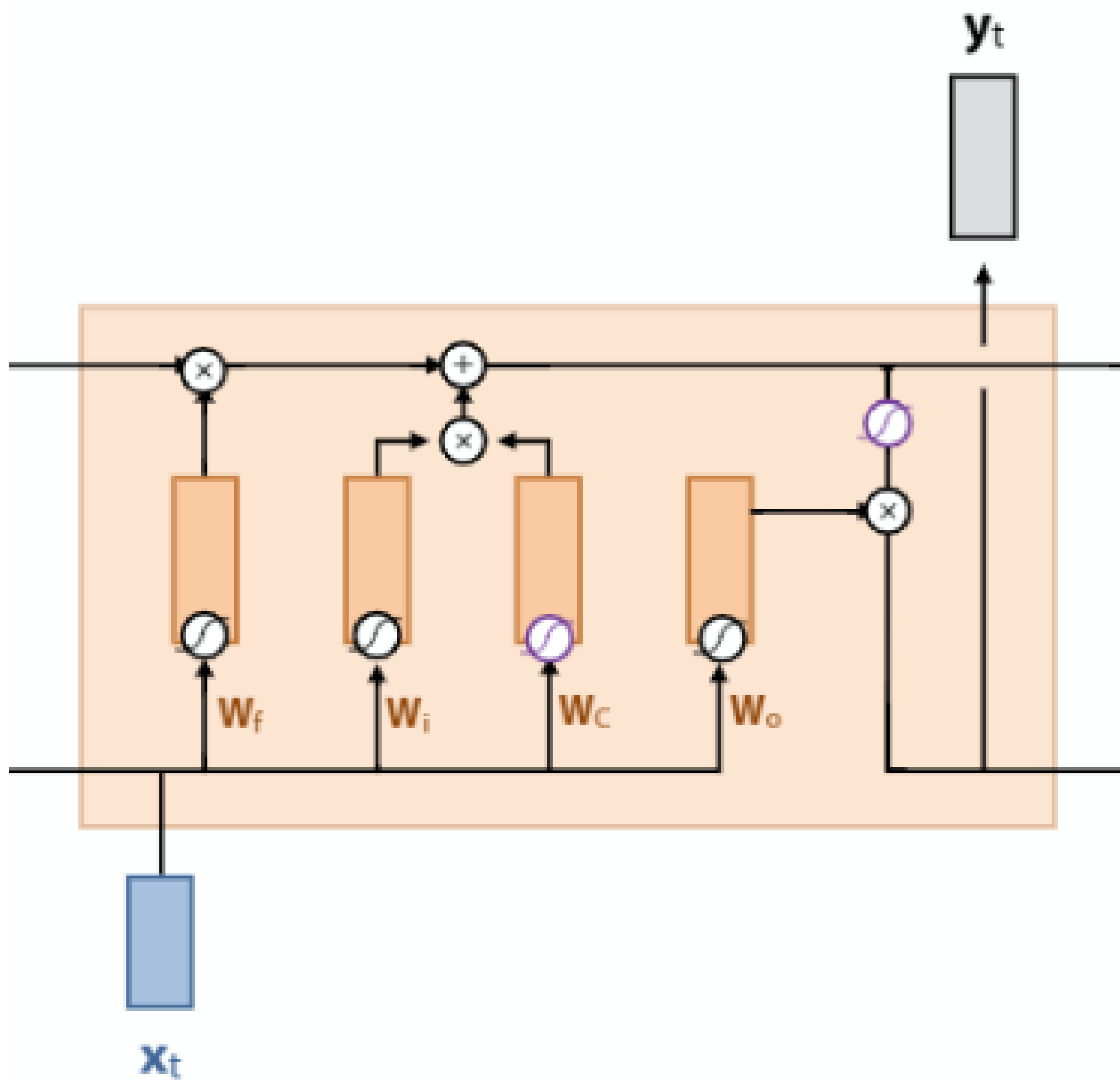
$a^{<t>}, x^{<t>}$ are images
$*$ : image convolution
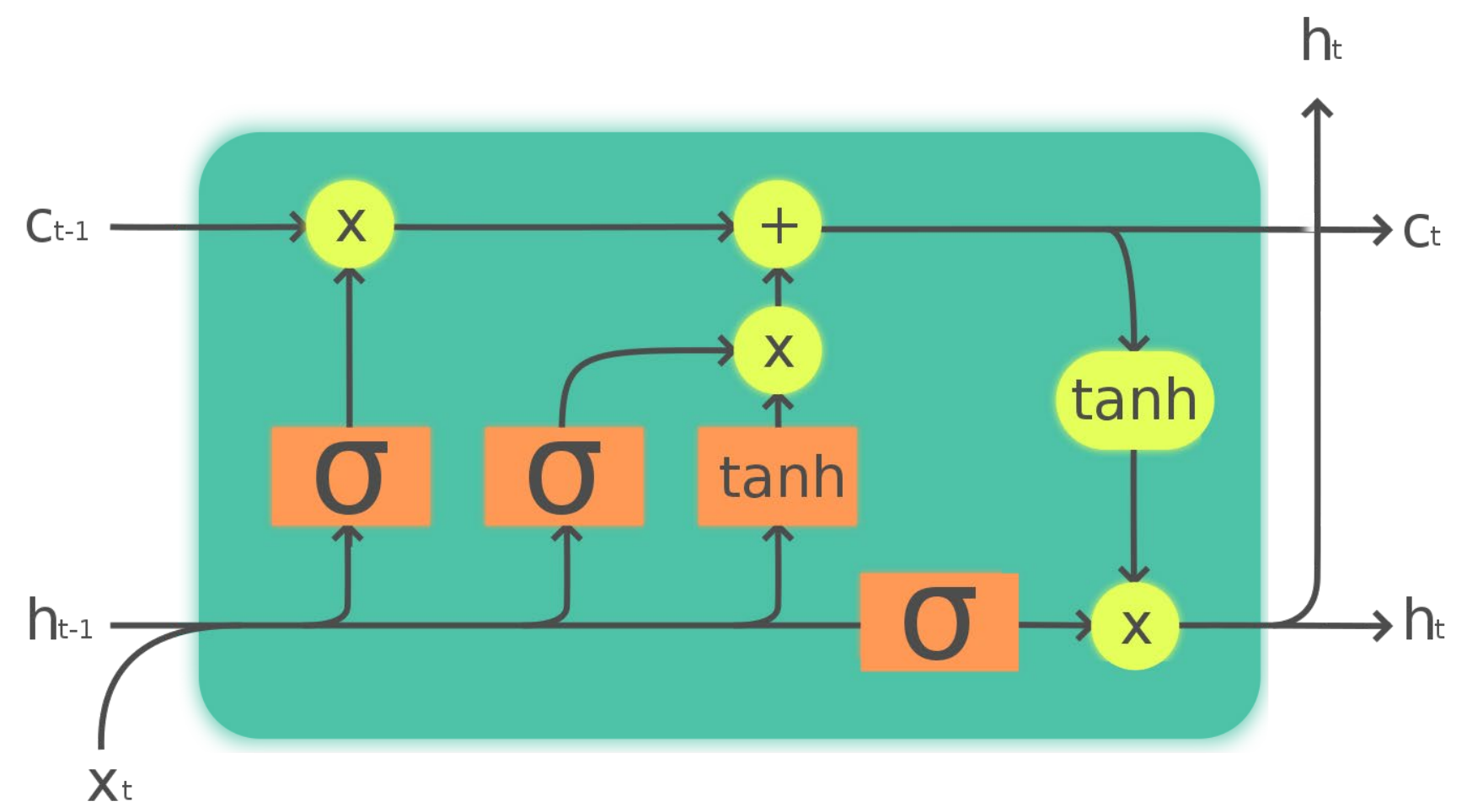After convolution, results are flattened to be a vector

https://colah.github.io/posts/2015-08-Understanding-LSTMs/

https://github.com/dlvu/dlvu.github.io/blob/main/slides/dlvu.lecture05.pdf

https://en.wikipedia.org/wiki/Long_short-term_memory#cite_note-43

## GRU

$$\tilde{c}^{<t>} = g_c(W_{ca} \cdot R \circ a^{<t-1>} + W_{cx} \cdot x^{<t>} + b_c)$$

$$c^{<t>} = U \circ \tilde{c}^{<t>} + (1 - U) \circ c^{<t-1>}$$

$F = 1 - U$: forget gate

$U = \sigma(W_{ua} \cdot a^{<t-1>} + W_{ux} \cdot x^{<t>} + b_u)$: update gate

$$a^{<t>} = c^{<t>}$$

~~$O = \sigma(W_{oa} \cdot a^{<t-1>} + W_{ox} \cdot x^{<t>} + b_o)$: output gate~~

$R = \sigma(W_{ra} \cdot a^{<t-1>} + W_{rx} \cdot x^{<t>} + b_r)$: relevant gate

$$y^{<t>} = g_y(W_{ya} \cdot a^{<t>} + b_y)$$

Learning phrase representations using RNN encoder-decoder for statistical machine translation. 2014. https://arxiv.org/abs/1406.1078

# Other variants

**MUT1:**

$$z = \text{sigm}(W_{xz}x_t + b_z)$$
$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$
$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z$$
$$+ \; h_t \odot (1 - z)$$

**MUT2:**

$$z = \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z)$$
$$r = \text{sigm}(x_t + W_{hr}h_t + b_r)$$
$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z$$
$$+ \; h_t \odot (1 - z)$$

**MUT3:**

$$z = \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z)$$
$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$
$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z$$
$$+ \; h_t \odot (1 - z)$$

- More than ten thousand architectures
- Different sequence tasks
- "We have evaluated a variety of recurrent neural network architectures in order to find an architecture that reliably outperforms the LSTM. **Though there were architectures that outperformed the LSTM on some problems, we were unable to find an architecture that consistently beat the LSTM and the GRU in all experimental conditions**."

An Empirical Exploration of Recurrent Network Architectures. 2015. http://proceedings.mlr.press/v37/jozefowicz15.pdf
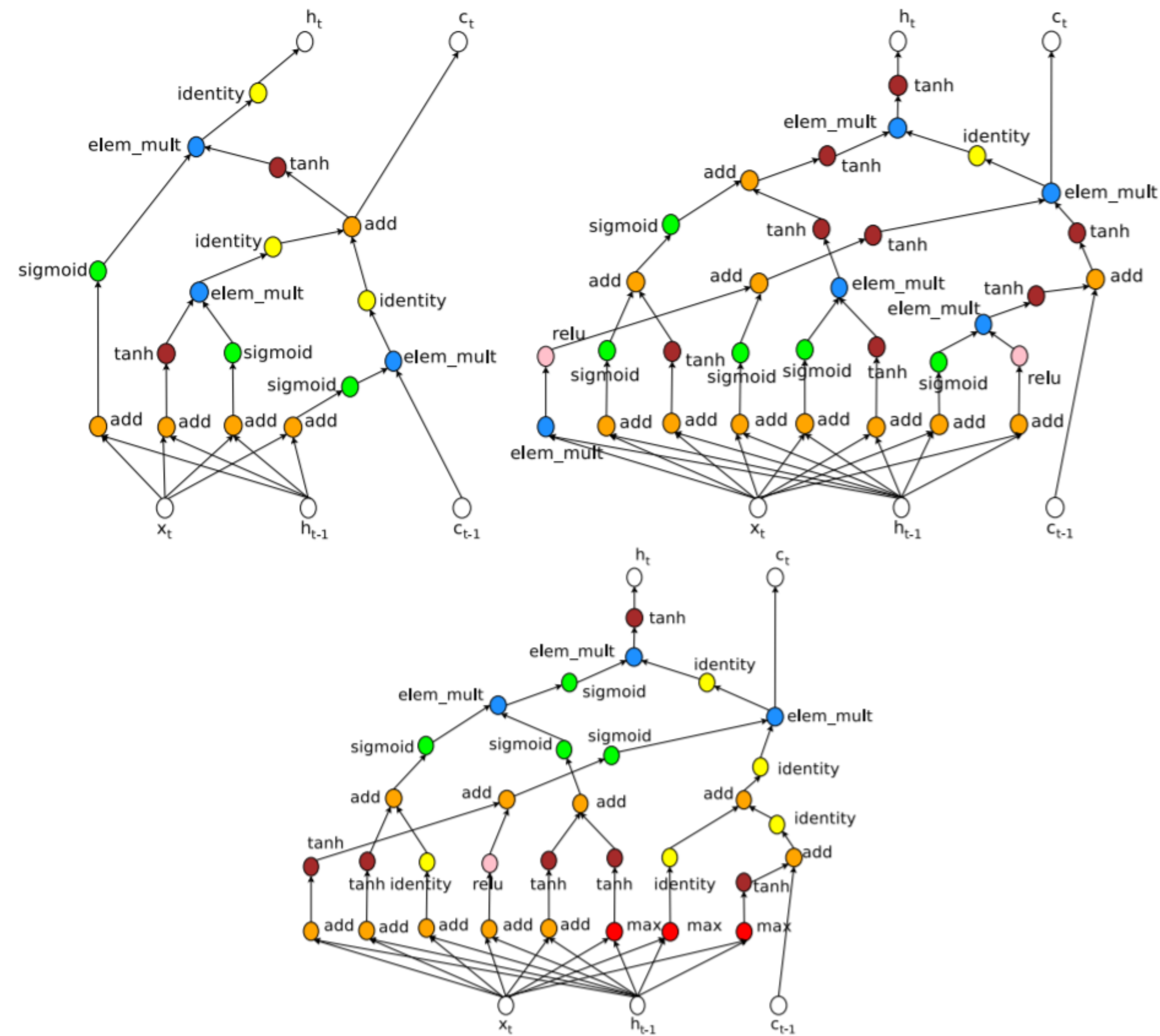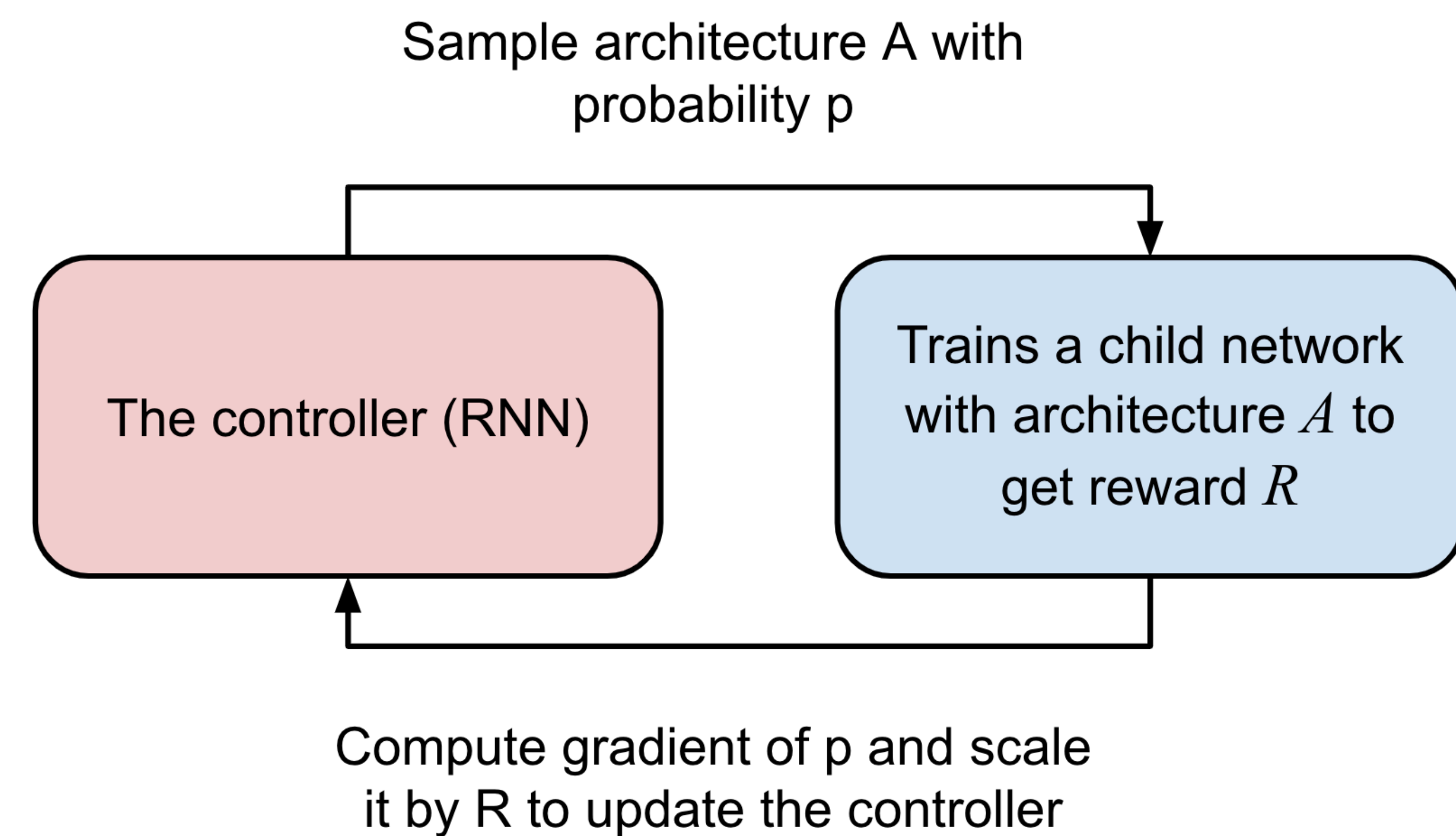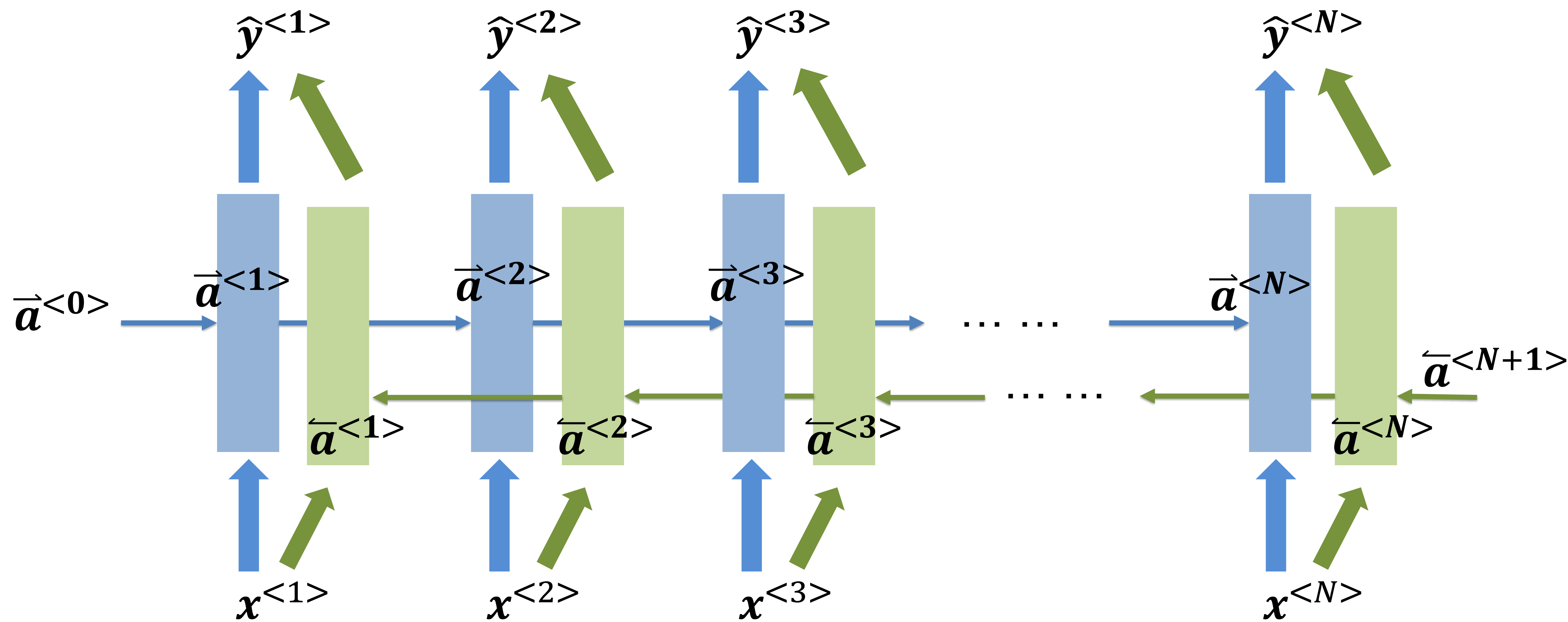
Figure 8: A comparison of the original LSTM cell vs. two good cells our model found. Top left: LSTM cell. Top right: Cell found by our model when the search space does not include *max* and *sin*. Bottom: Cell found by our model when the search space includes *max* and *sin* (the controller did not choose to use the *sin* function).

Neural Architecture Search with Reinforcement Learning. 2017.

- Require large computing resource
- Can find better architecture on given tasks



Sample architecture A with probability p

The controller (RNN)

Trains a child network with architecture $A$ to get reward $R$

Compute gradient of p and scale it by R to update the controller

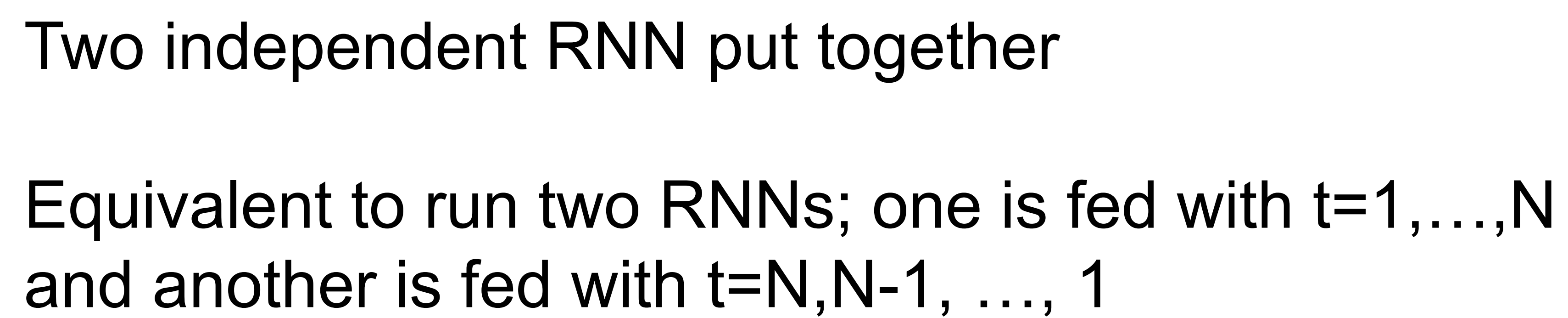https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Bidirectional RNN (BRNN)



- Some applications provide the entire series, e.g. machine translation

- Or, the deployment does not require real-time processing

- Allow utilizing information from the past and the future

$$y^{<t>} = g_y(W_{ya} \cdot [\overrightarrow{a}^{<t>}; \overleftarrow{a}^{<t>}] + b_y) \quad \text{or} \quad y^{<t>} = g_y(W_{ya} \cdot (\overrightarrow{a}^{<t>} + \overleftarrow{a}^{<t>}) + b_y)$$

# Bidirectional RNN (BRNN)

- Still acyclic graph
- Can run forward-pass and backprop
- Supported by PyTorch or Keras

$$\overleftarrow{a}^{<1>} \quad x^{<1>} \quad \overrightarrow{a}^{<1>} \quad \overrightarrow{a}^{<2>}$$

$$\overleftarrow{a}^{<2>} \quad \overrightarrow{a}^{<3>}$$

$$\overleftarrow{a}^{<3>} \quad x^{<4>} \quad \overrightarrow{a}^{<4>}$$

$$\overleftarrow{a}^{<4>}$$

Two independent RNN put together

Equivalent to run two RNNs; one is fed with t=1,…,N and another is fed with t=N,N-1, …, 1

CLASS torch.nn.GRU(*args, **kwargs) [SOURCE]

Applies a multi-layer gated recurrent unit (GRU) RNN to an input sequence.

For each element in the input sequence, each layer computes the following function:

$$
\begin{aligned}
r_t &= \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr}) \\
z_t &= \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz}) \\
n_t &= \tanh(W_{in}x_t + b_{in} + r_t * (W_{hn}h_{(t-1)} + b_{hn})) \\
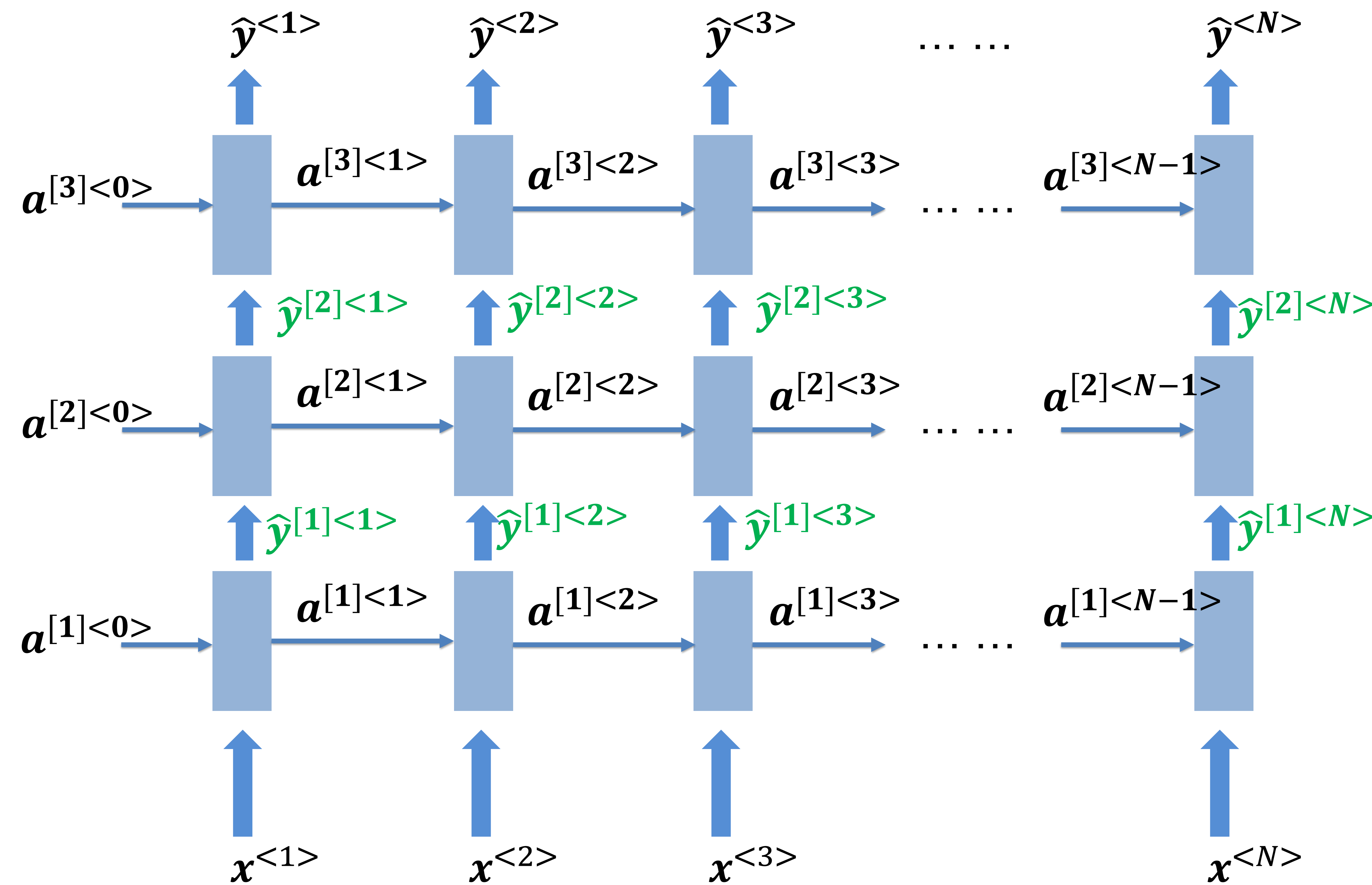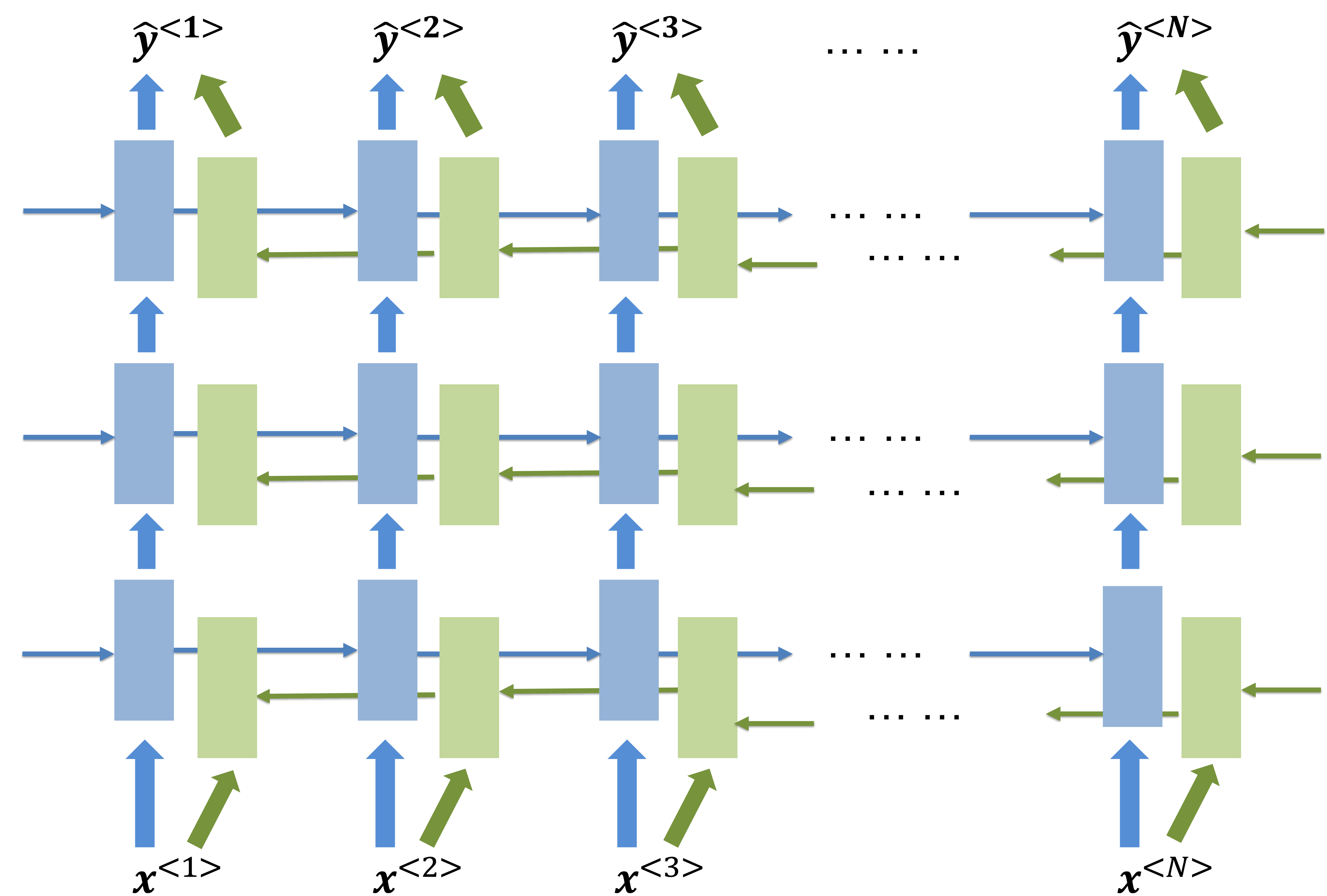h_t &= (1 - z_t) * n_t + z_t * h_{(t-1)}
\end{aligned}
$$

where $h_t$ is the hidden state at time $t$, $x_t$ is the input at time $t$, $h_{(t-1)}$ is the hidden state of the layer at time $t$-1 or the initial hidden state at time $o$, and $r_t$, $z_t$, $n_t$ are the reset, update, and new gates, respectively. $\sigma$ is the sigmoid function, and $*$ is the Hadamard product.

In a multilayer GRU, the input $x_t^{(l)}$ of the $l$-th layer ($l >= 2$) is the hidden state $h_t^{(l-1)}$ of the previous layer multiplied by dropout $\delta_t^{(l-1)}$ where each $\delta_t^{(l-1)}$ is a Bernoulli random variable which is $0$ with probability dropout.

**Parameters**

- **input_size** – The number of expected features in the input $x$
- **hidden_size** – The number of features in the hidden state $h$
- **num_layers** – Number of recurrent layers. E.g., setting num_layers=2 would mean stacking two GRUs together to form a *stacked GRU*, with the second GRU taking in outputs of the first GRU and computing the final results. Default: 1
- **bias** – If False, then the layer does not use bias weights $b\_ih$ and $b\_hh$. Default: True
- **batch_first** – If True, then the input and output tensors are provided as (batch, seq, feature). Default: False
- **dropout** – If non-zero, introduces a *Dropout* layer on the outputs of each GRU layer except the last layer, with dropout probability equal to dropout. Default: 0
- **bidirectional** – If True, becomes a bidirectional GRU. Default: False

$$a^{[l]<t>}$$
$$= g_a(W_{aa}^{[l]} \cdot a^{[l]<t-1>} + W_{ax}^{[l]} \cdot \hat{y}^{[l-1]<t>} + b_a^{[l]})$$

$$y^{[l]<t>} = g_y(W_{ya}^{[l]} \cdot a^{[l]<t>} + b_y^{[l]})$$

CLASS `torch.nn.GRU(*args, **kwargs)`                    [SOURCE]

Applies a multi-layer gated recurrent unit (GRU) RNN to an input sequence.

For each element in the input sequence, each layer computes the following function:

$$r_t = \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr})$$
$$z_t = \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz})$$
$$n_t = \tanh(W_{in}x_t + b_{in} + r_t * (W_{hn}h_{(t-1)} + b_{hn}))$$
$$h_t = (1 - z_t) * n_t + z_t * h_{(t-1)}$$

where $h_t$ is the hidden state at time $t$, $x_t$ is the input at time $t$, $h_{(t-1)}$ is the hidden state of the layer at time $t$-$1$ or the initial hidden state at time $o$, and $r_t$, $z_t$, $n_t$ are the reset, update, and new gates, respectively. $\sigma$ is the sigmoid function, and $*$ is the Hadamard product.

In a multilayer GRU, the input $x_t^{(l)}$ of the $l$ -th layer ($l >= 2$) is the hidden state $h_t^{(l-1)}$ of the previous layer multiplied by dropout $\delta_t^{(l-1)}$ where each $\delta_t^{(l-1)}$ is a Bernoulli random variable which is $0$ with probability `dropout`.

Parameters

- **input_size** – The number of expected features in the input $x$
- **hidden_size** – The number of features in the hidden state $h$
- **num_layers** – Number of recurrent layers. E.g., setting `num_layers=2` would mean stacking two GRUs together to form a *stacked GRU*, with the second GRU taking in outputs of the first GRU and computing the final results. Default: 1
- **bias** – If `False`, then the layer does not use bias weights $b\_ih$ and $b\_hh$. Default: `True`
- **batch_first** – If `True`, then the input and output tensors are provided as (batch, seq, feature). Default: `False`
- **dropout** – If non-zero, introduces a *Dropout* layer on the outputs of each GRU layer except the last layer, with dropout probability equal to `dropout`. Default: 0
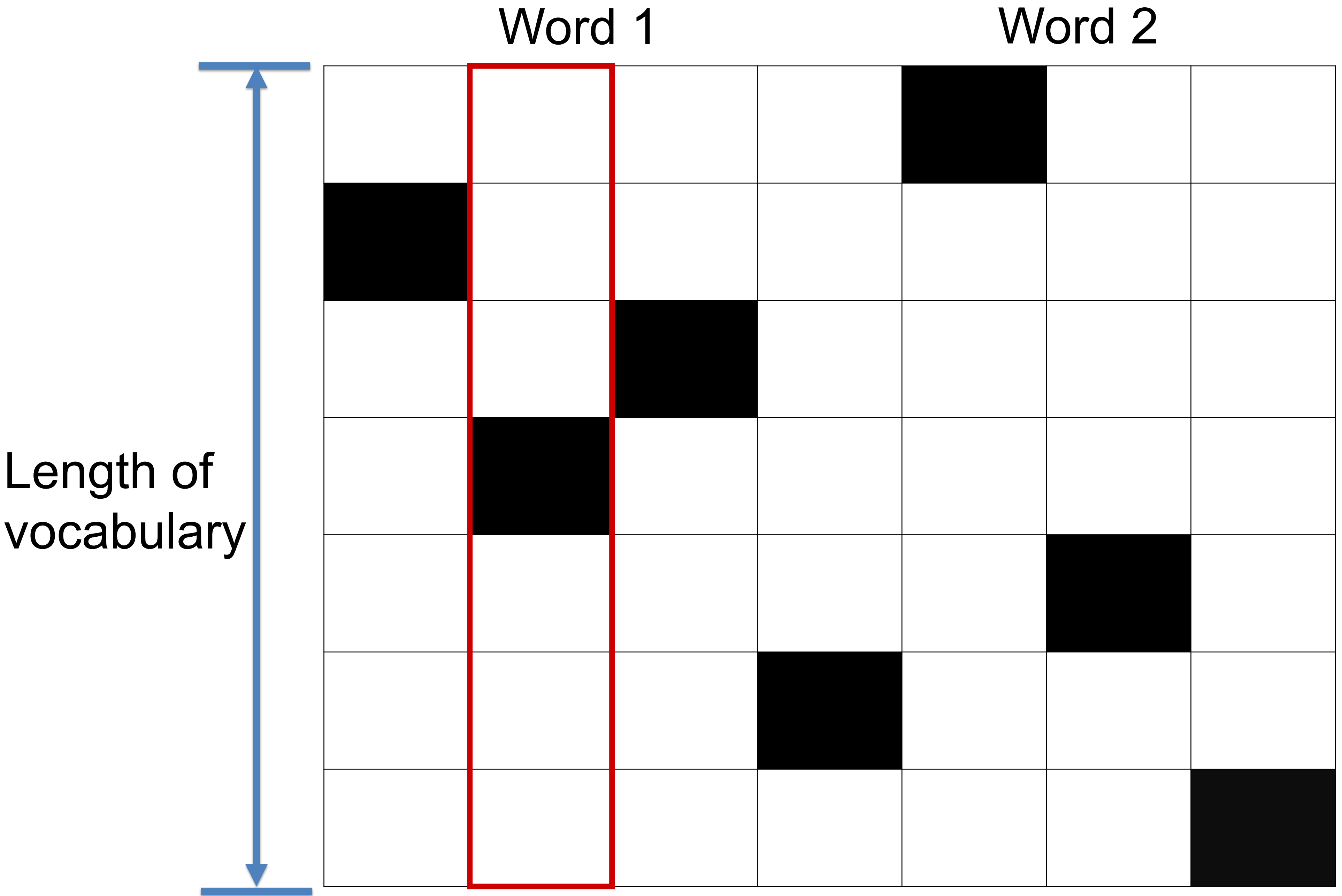- **bidirectional** – If `True`, becomes a bidirectional GRU. Default: `False`

**This is a picture showing kids are playing.**

Word level vocabulary:

Use one-hot encoding for words

[

| | |
|---|---|
| a | 0 |
| aa | 1 |
| … | |
| air | |
| … | |
| is | 2876 |
| … | |
| grass | |
| … | |
| this | |
| … | |
| playing | 27890 |
| … | |
| zinc | |

]

Word 1                    Word 2

Length of vocabulary

- 10K-30K entries in English vocabulary
- Assume all words are orthogonal to each other
- Does not utilize any "pre-training"

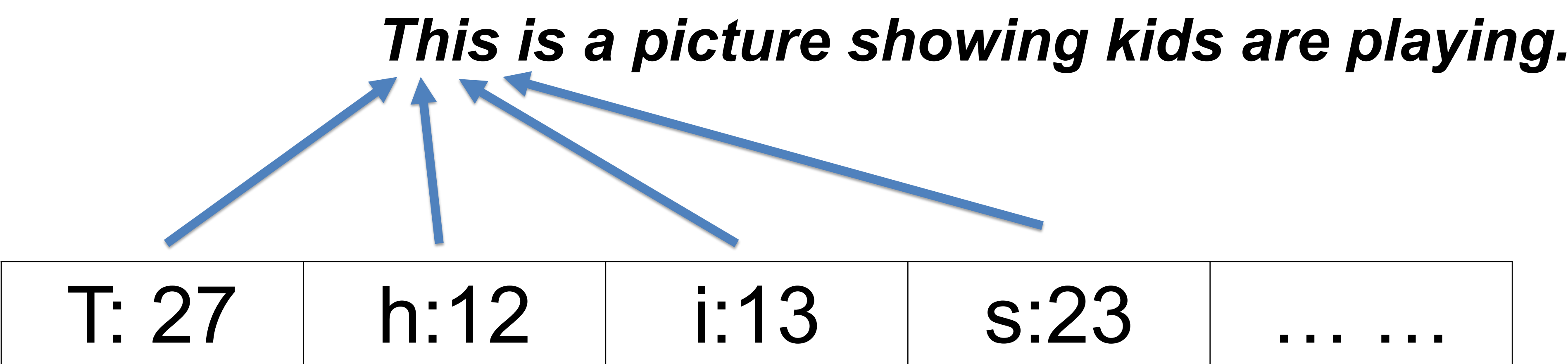Character level vocabulary:

Use one-hot encoding for characters

Often use compact representation:

| ASCII printable characters | | | | | |
|---|---|---|---|---|---|
| 32 | space | 64 | @ | 96 | ` |
| 33 | ! | 65 | A | 97 | a |
| 34 | " | 66 | B | 98 | b |
| 35 | # | 67 | C | 99 | c |
| 36 | $ | 68 | D | 100 | d |
| 37 | % | 69 | E | 101 | e |
| 38 | & | 70 | F | 102 | f |
| 39 | ' | 71 | G | 103 | g |
| 40 | ( | 72 | H | 104 | h |
| 41 | ) | 73 | I | 105 | i |
| 42 | * | 74 | J | 106 | j |
| 43 | + | 75 | K | 107 | k |
| 44 | , | 76 | L | 108 | l |
| 45 | - | 77 | M | 109 | m |
| 46 | . | 78 | N | 110 | n |
| 47 | / | 79 | O | 111 | o |
| 48 | 0 | 80 | P | 112 | p |
| 49 | 1 | 81 | Q | 113 | q |
| 50 | 2 | 82 | R | 114 | r |
| 51 | 3 | 83 | S | 115 | s |
| 52 | 4 | 84 | T | 116 | t |
| 53 | 5 | 85 | U | 117 | u |
| 54 | 6 | 86 | V | 118 | v |
| 55 | 7 | 87 | W | 119 | w |
| 56 | 8 | 88 | X | 120 | x |
| 57 | 9 | 89 | Y | 121 | y |
| 58 | : | 90 | Z | 122 | z |
| 59 | ; | 91 | [ | 123 | { |
| 60 | < | 92 | \ | 124 | | |
| 61 | = | 93 | ] | 125 | } |
| 62 | > | 94 | ^ | 126 | ~ |
| 63 | ? | 95 | _ | | |

***This is a picture showing kids are playing.***

| T: 27 | h:12 | i:13 | s:23 | … … |
|---|---|---|---|---|

*Hello!*
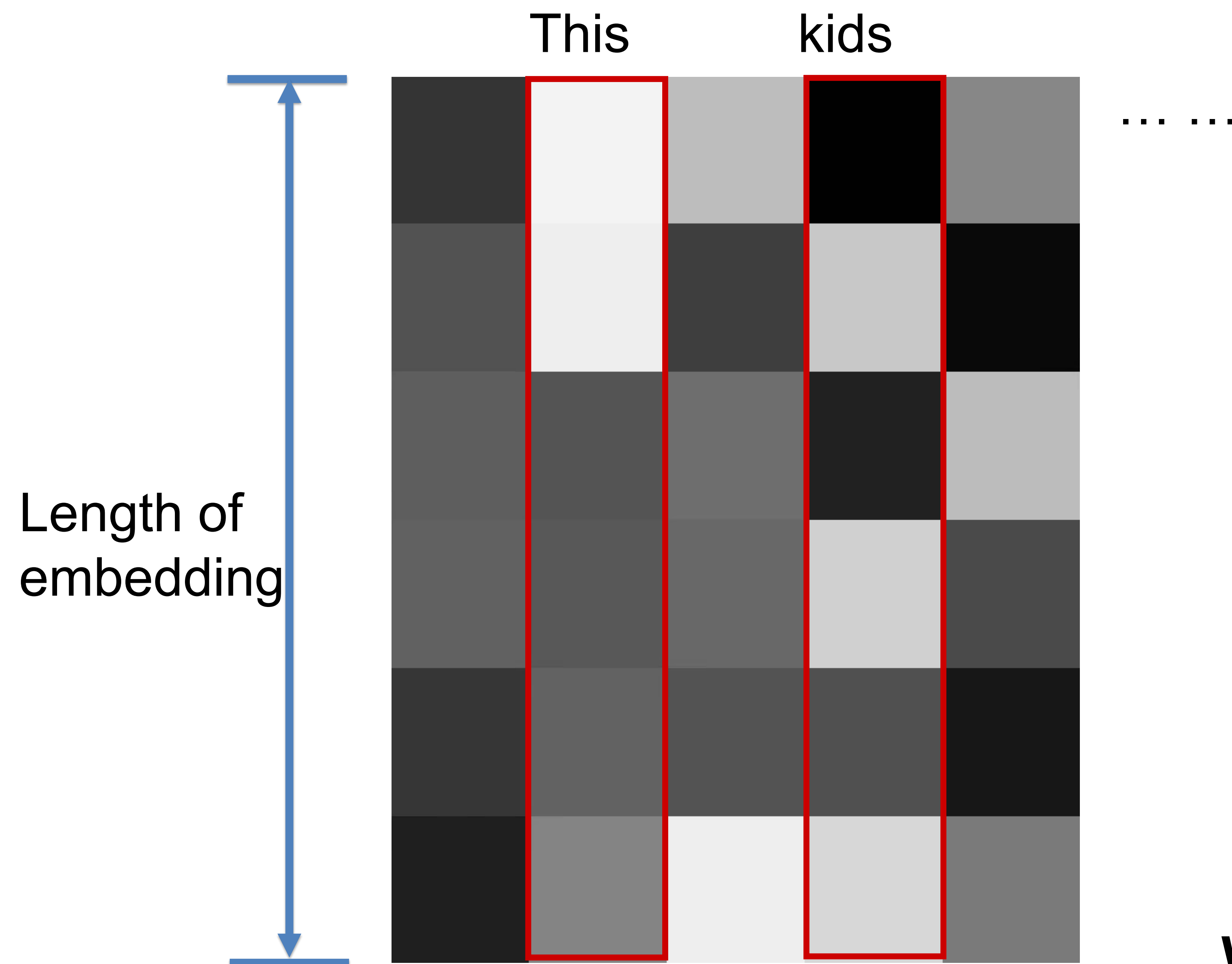


- Casual training
- Finish with <EOS> token

# Word embedding

***This is a picture showing kids are playing.***
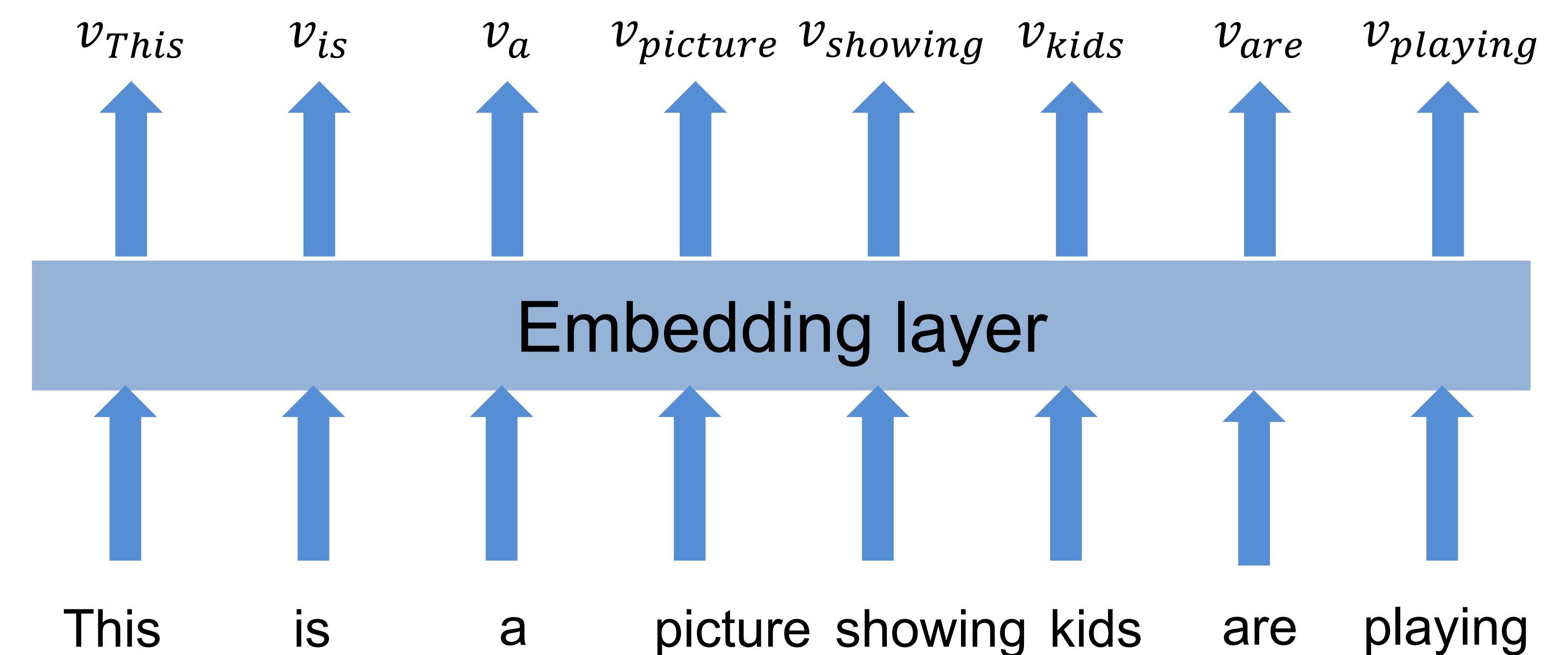
Use a real-value vector to represent words

This          kids



… …

Length of embedding

EMBEDDING 🔗

CLASS  `torch.nn.Embedding(`*`num_embeddings, embedding_dim, padding_idx=None,`* [SOURCE]
`      `*`max_norm=None, norm_type=2.0, scale_grad_by_freq=False,`*
`      `*`sparse=False, _weight=None`*`)`

A simple lookup table that stores embeddings of a fixed dictionary and size.

$$v = W \times one\_hot\_vector \quad W: \text{embedding\_dim x num\_embedding}$$

$v_{This}$   $v_{is}$   $v_a$   $v_{picture}$   $v_{showing}$   $v_{kids}$   $v_{are}$   $v_{playing}$

Embedding layer

This      is      a      picture   showing   kids      are      playing

**W** can be learned or pre-trained.

**This is a _picture_ showing kids are playing.**

$x$

Skip-grams

Self-supervised learning

X : center word

Y: neighboring word

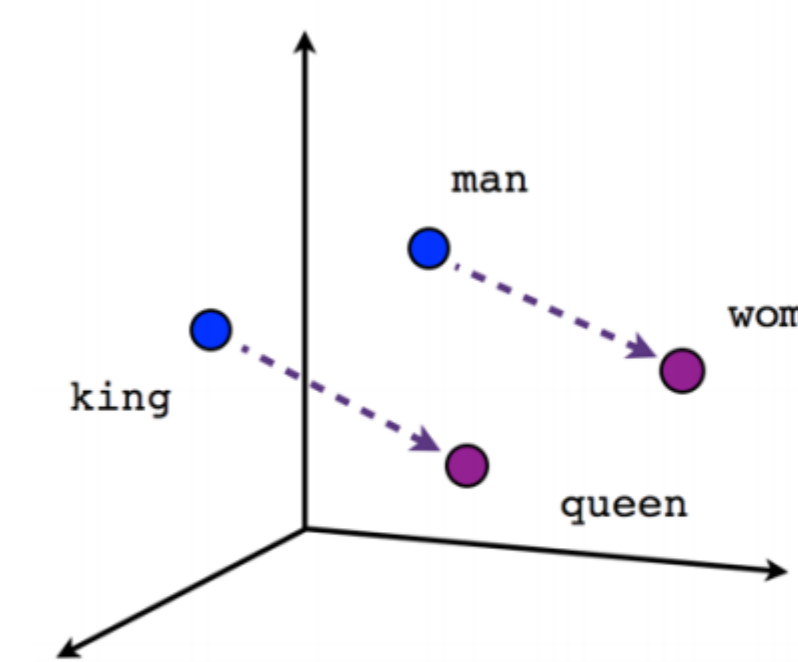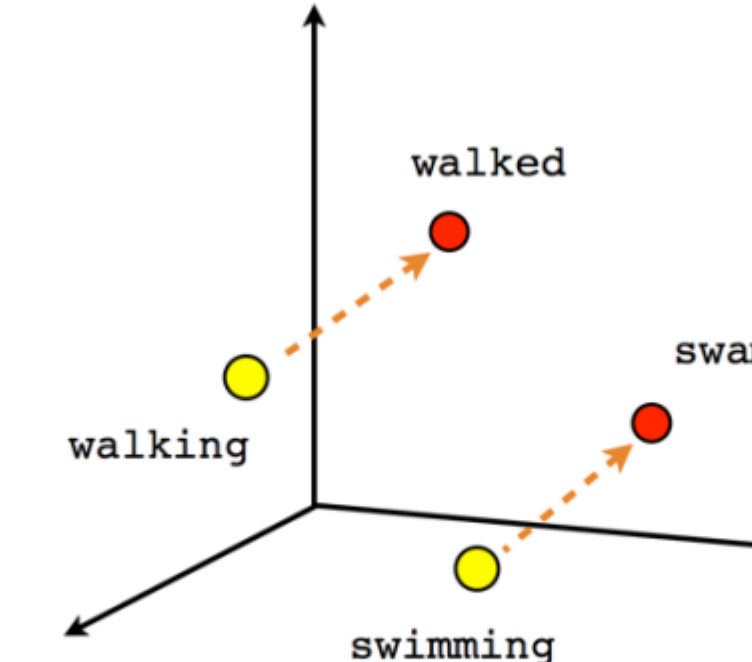| Center | Context |
|--------|---------|
| _picture_ | showing |
| _picture_ | a |
| _picture_ | is |
| _picture_ | This |

**showing**

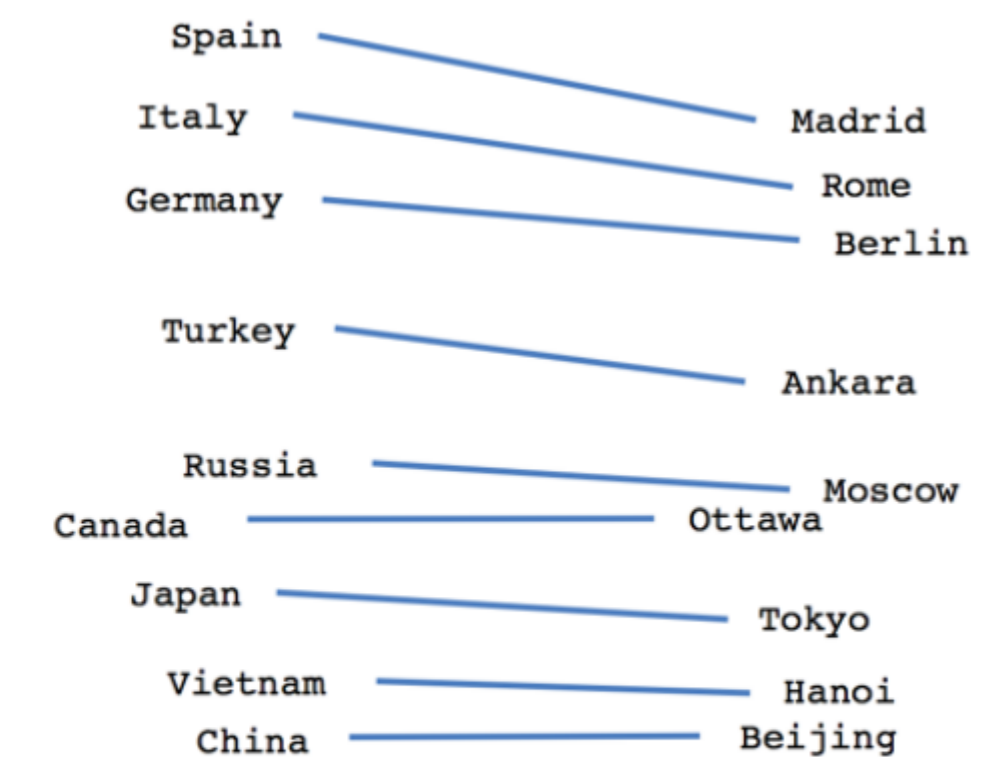$$SoftMax$$

$$W \times (\cdot) + b$$

Embedding layer

**_picture_**

Male-Female     Verb tense     Country-Capital

- Train on a large corpse of texts
- Map words to N-dimensional space
- Learned semantic relationship between words

$$e_{king} + (e_{woman} - e_{man}) \approx e_{queen}$$

https://www.tensorflow.org/tutorials/representation/word2vec

Efficient Estimation of Word Representations in Vector Space.2013. https://arxiv.org/abs/1301.3781

# Word2Vec : carry the accumulated "bias" of the world



he (128)

she (72)

http://wordbias.umiacs.umd.edu/
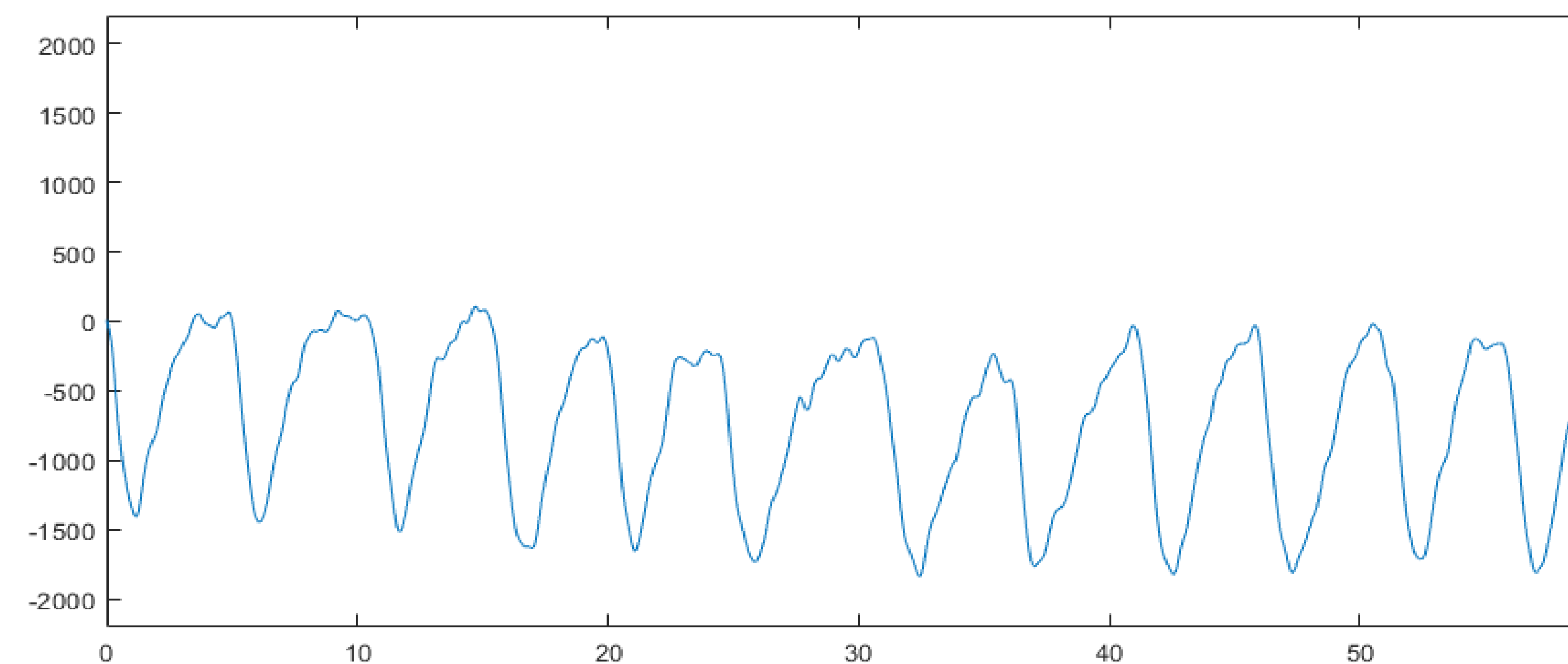
# Sequence pre-processing

Multiple instances of sequences:

**This is a picture showing kids are playing** ➡ 8x1024 matrix
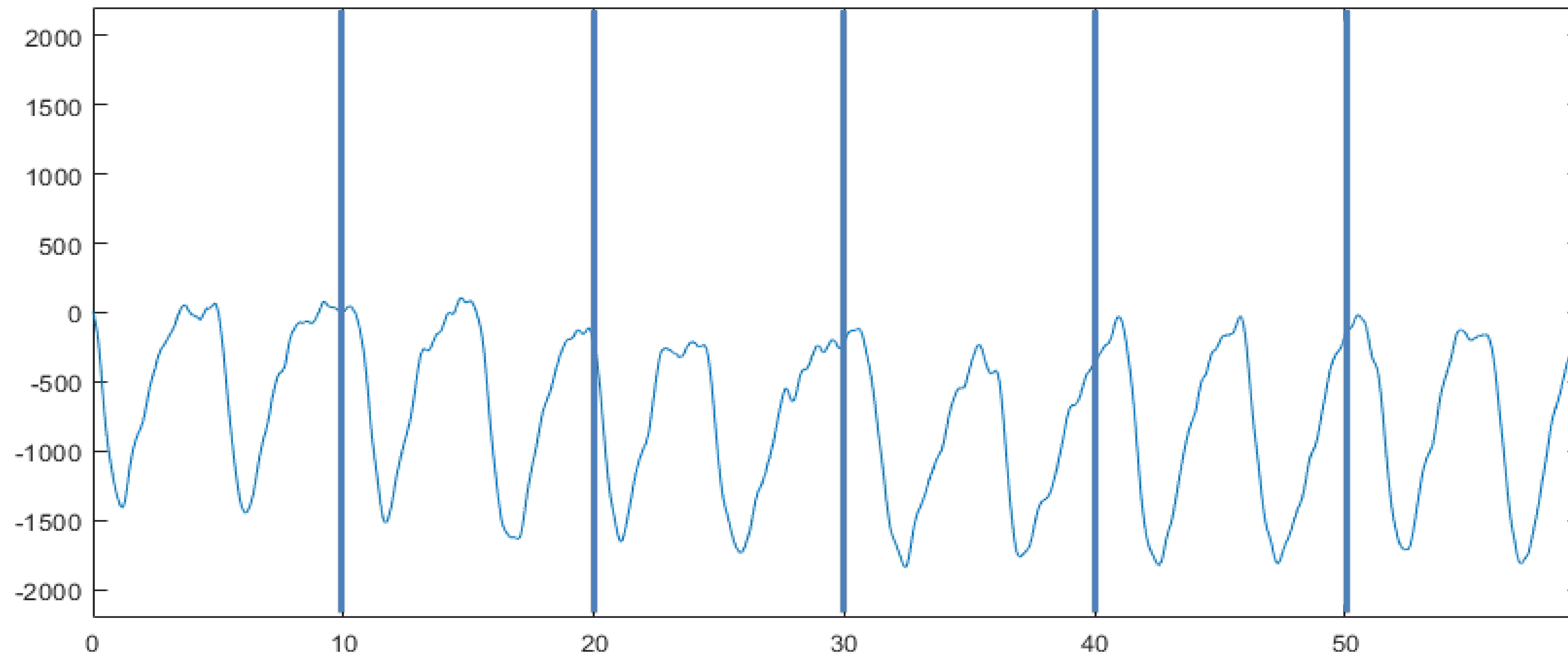
**Today is a good day** ➡ 5x1024 matrix

**The boy is reading a book under the tree** ➡ 9x1024 matrix

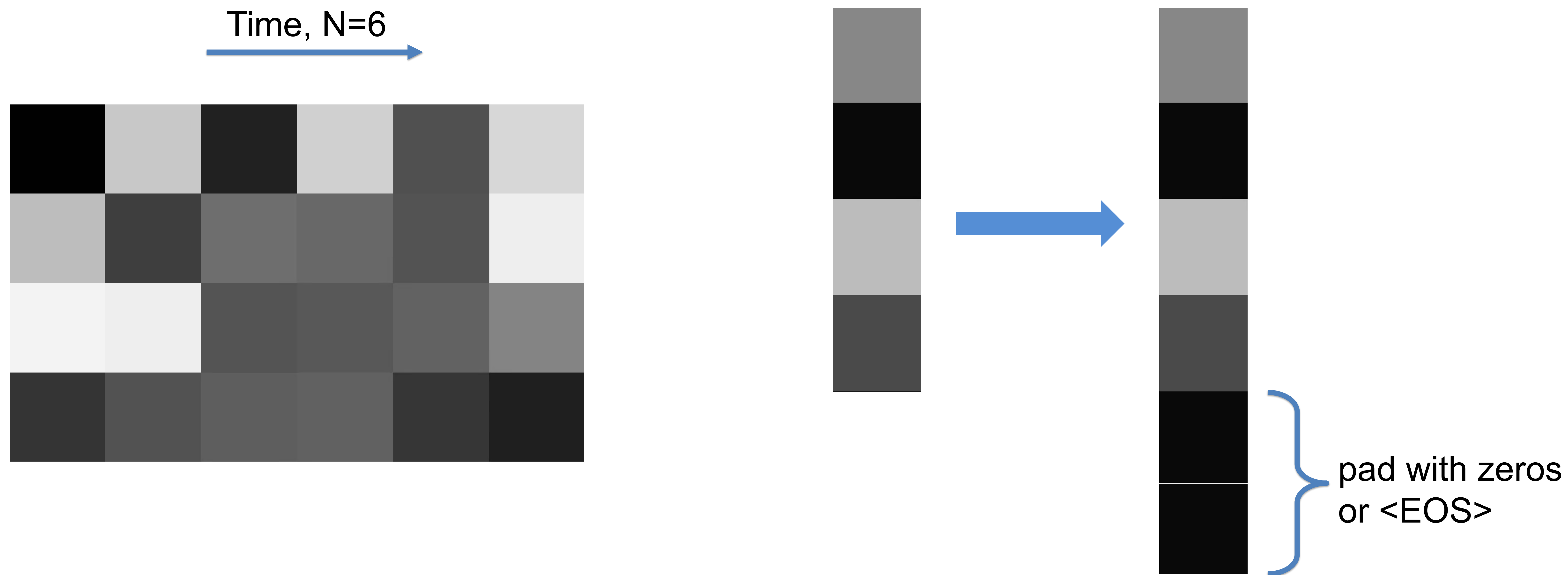One instance of long sequence:

# Sequence pre-processing

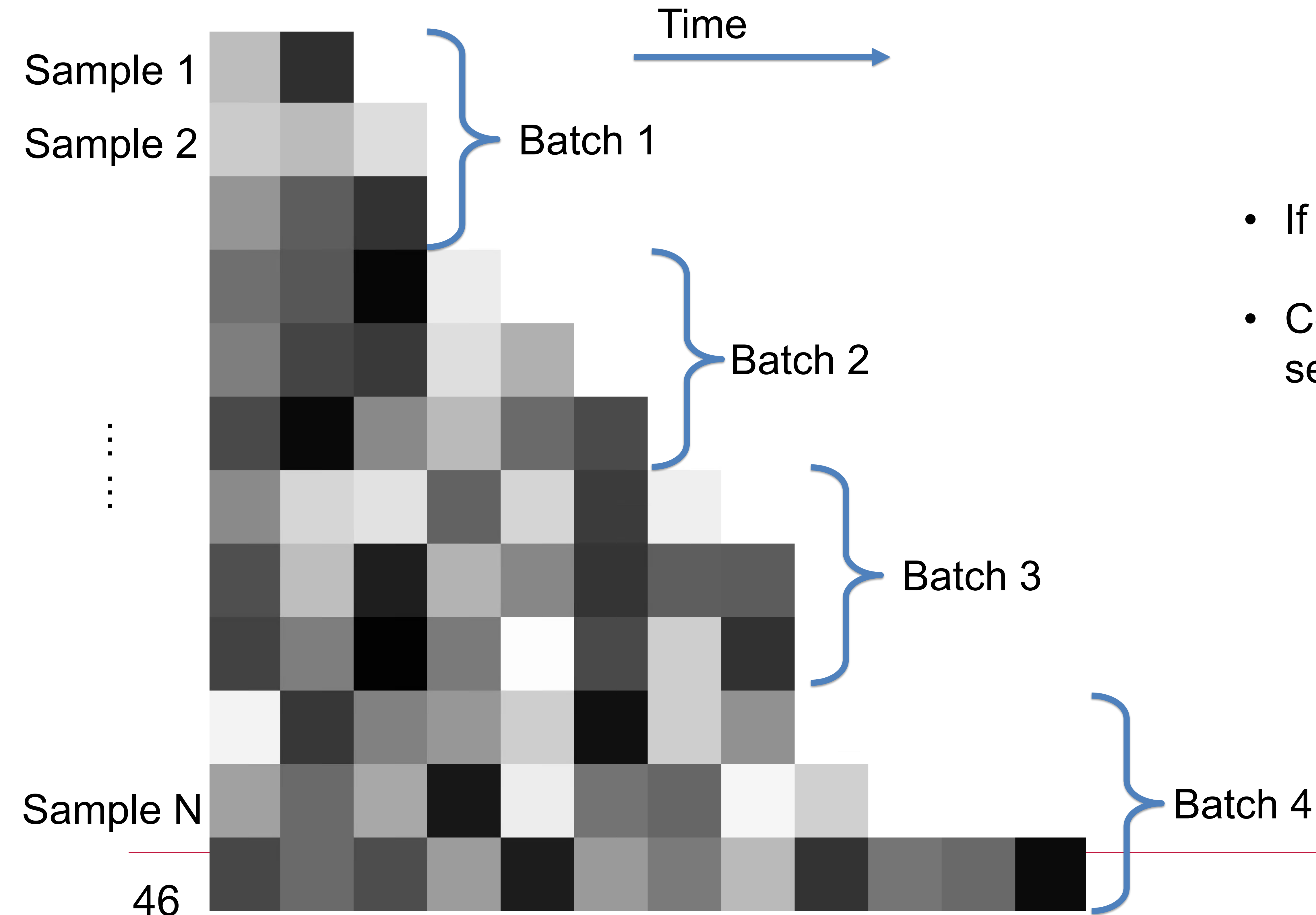One instance of long sequence:



Split to equal length chunk

# Sequence pre-processing

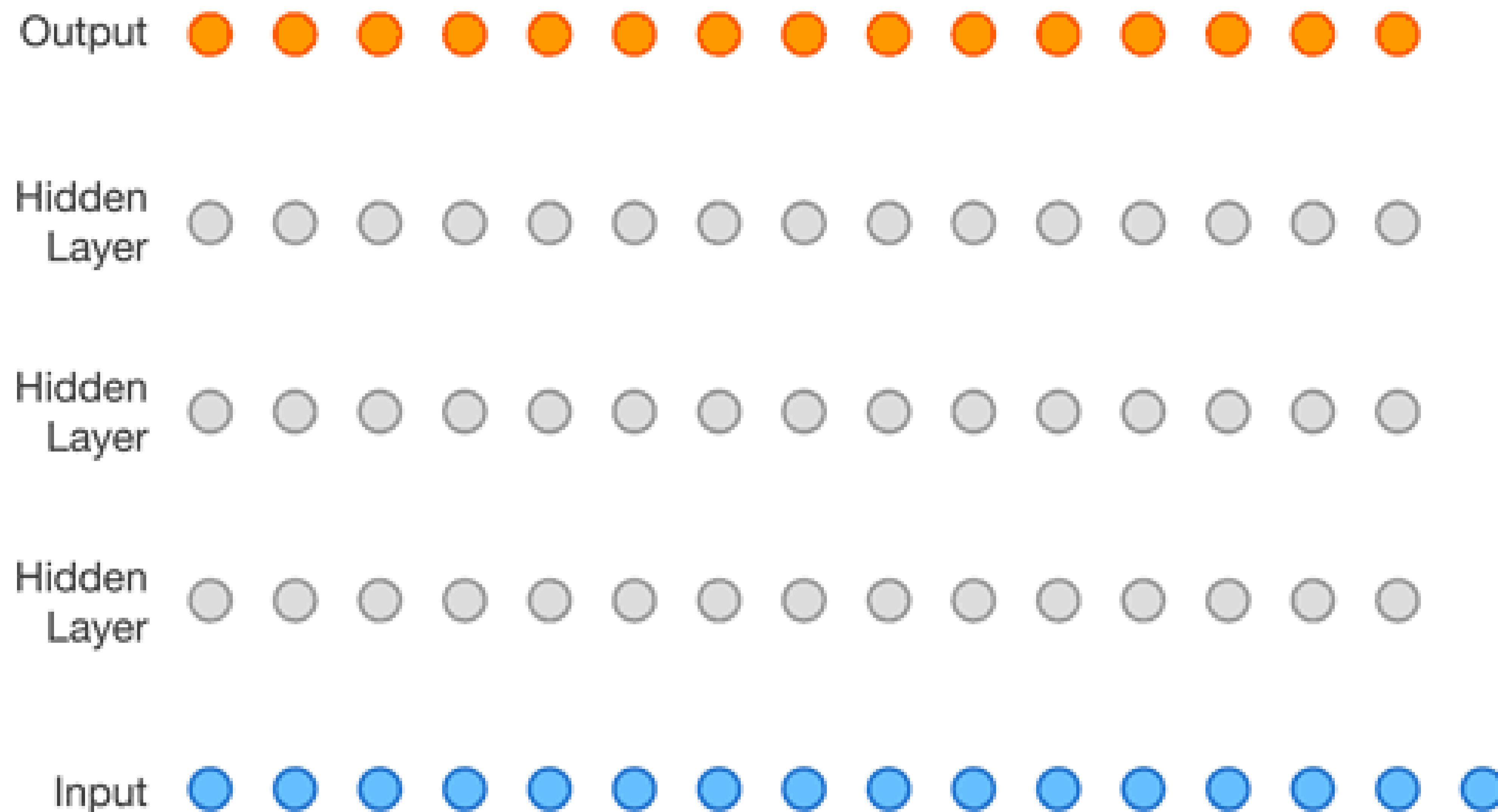Multiple instances of sequences : padding to equal length

Time, N=6

pad with zeros
or <EOS>

# Sequence pre-processing

Multiple instances of sequences : padding to equal length, with sorting

Time

Sample 1

Sample 2

Batch 1

- If data distribution is not biased by its length

- Consider to reduce batch size for longer sequences

Batch 2

Batch 3

Sample N

Batch 4

- Use convolution to process sequence data
- Dilated + Casual convolution

$$Receptive\_length\_in\_time = 2^L(kernel_{size} - 1)$$

Dilated convolution



dilation=1    dilation=2    dilation=3

Output

Hidden Layer

Hidden Layer

Hidden Layer
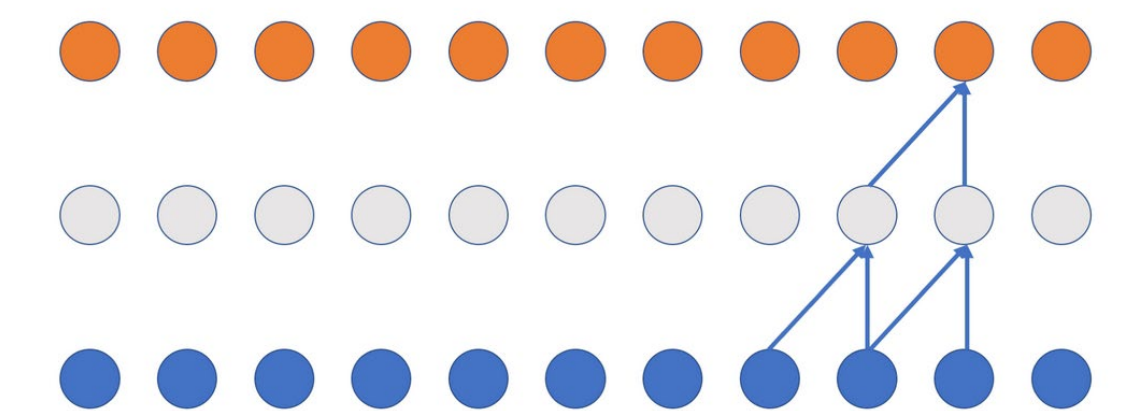
Input

Normal convolution    Casual convolution



- Fast training, fast inference
- Google Wavenet – text to audio generation, used in e.g. Amazon Alexa

https://deepmind.com/blog/article/wavenet-generative-model-raw-audio