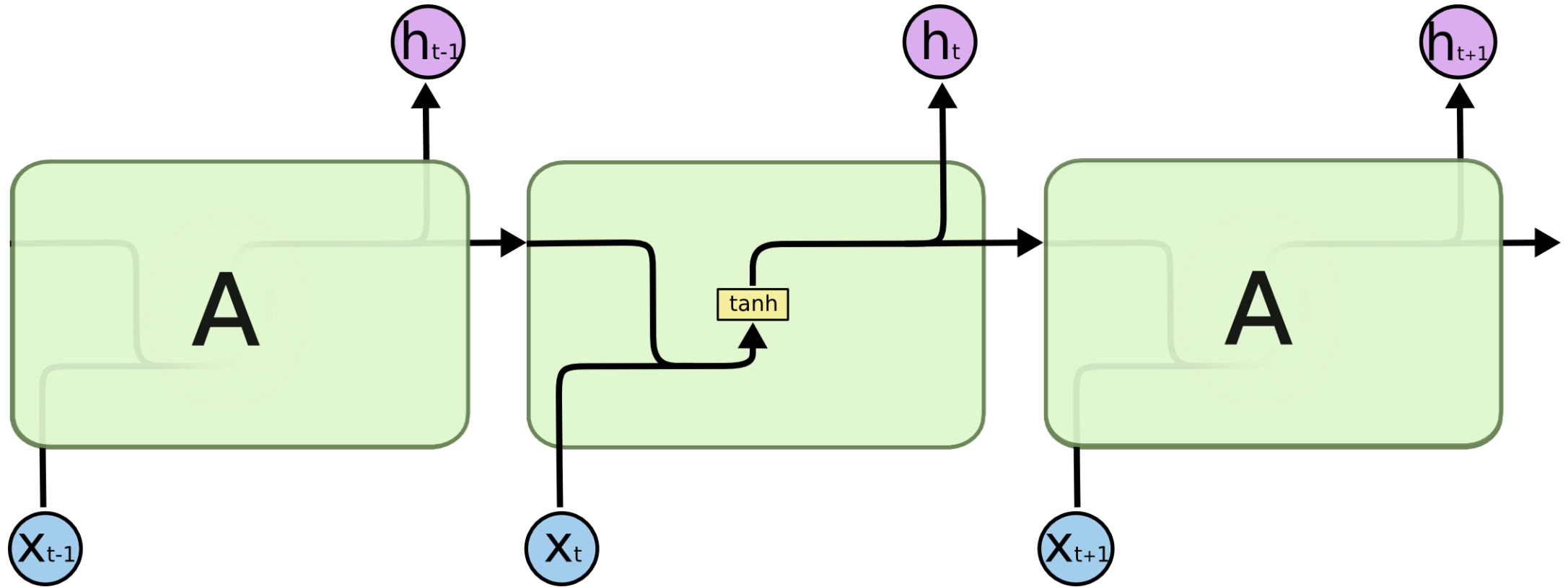


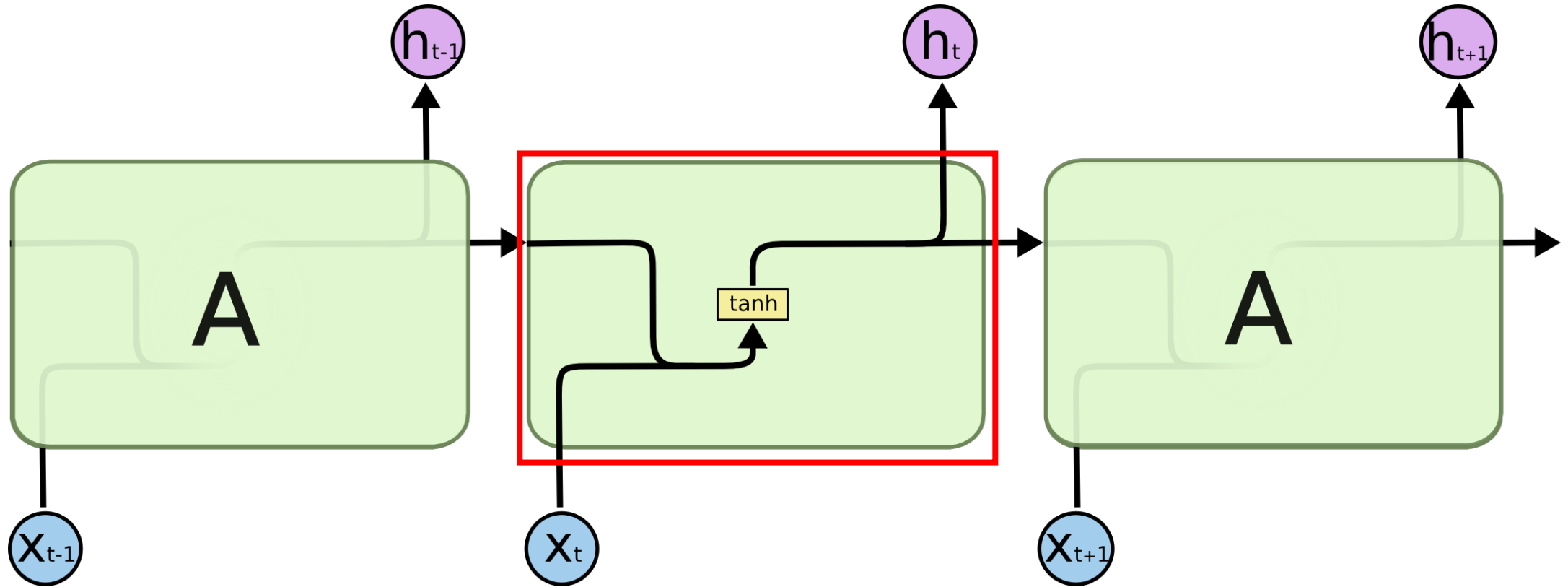
Long-Short Term Memory

Neural Networks Design And Application

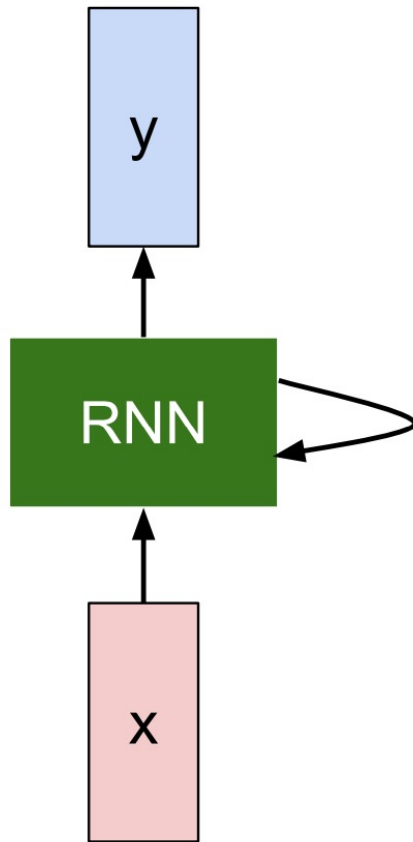
A conventional RNN



A conventional RNN



Recurrent neural networks



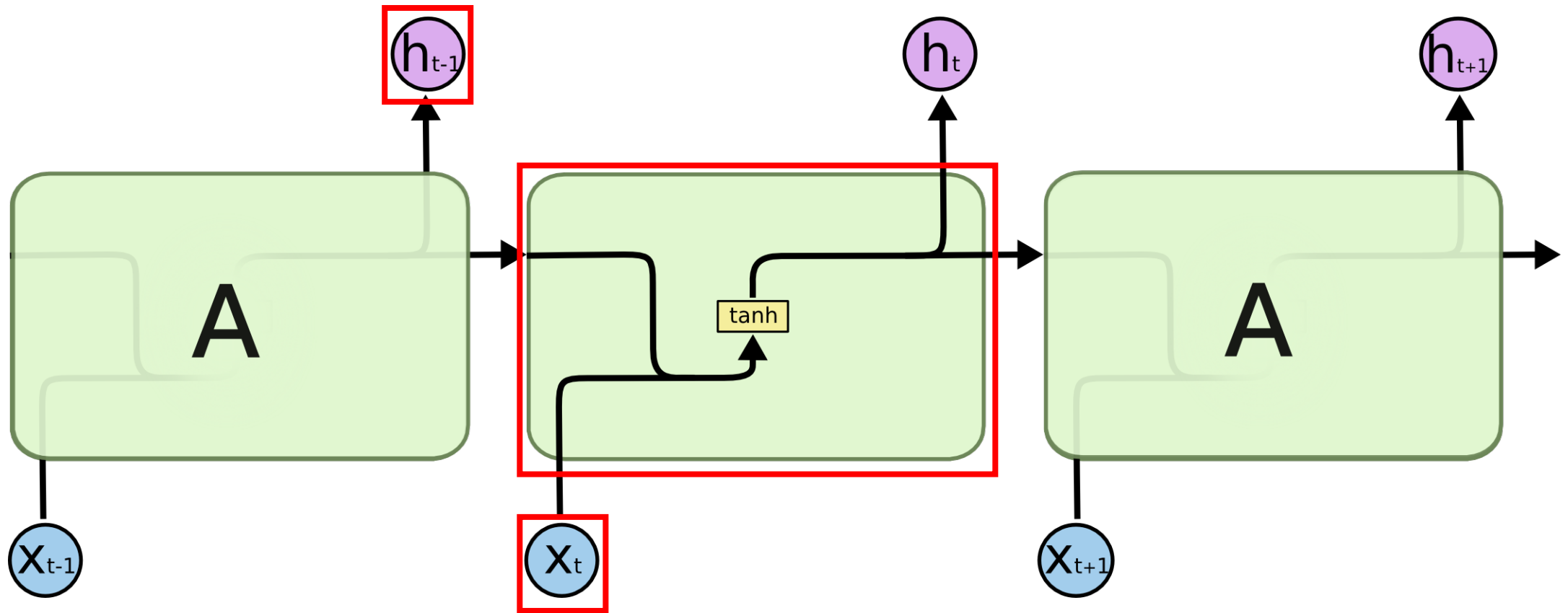
$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right]$$

$$f = \tanh(\cdot)$$

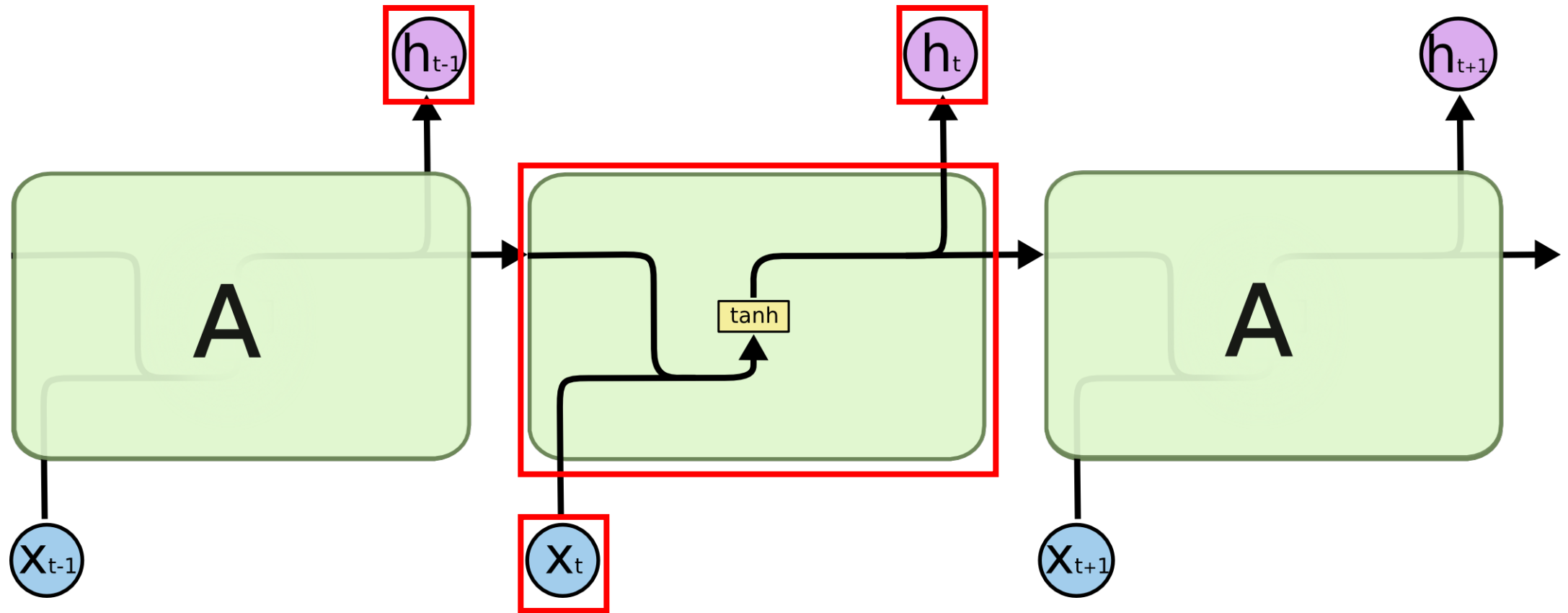
$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state some function with parameters W old state input vector at some time step

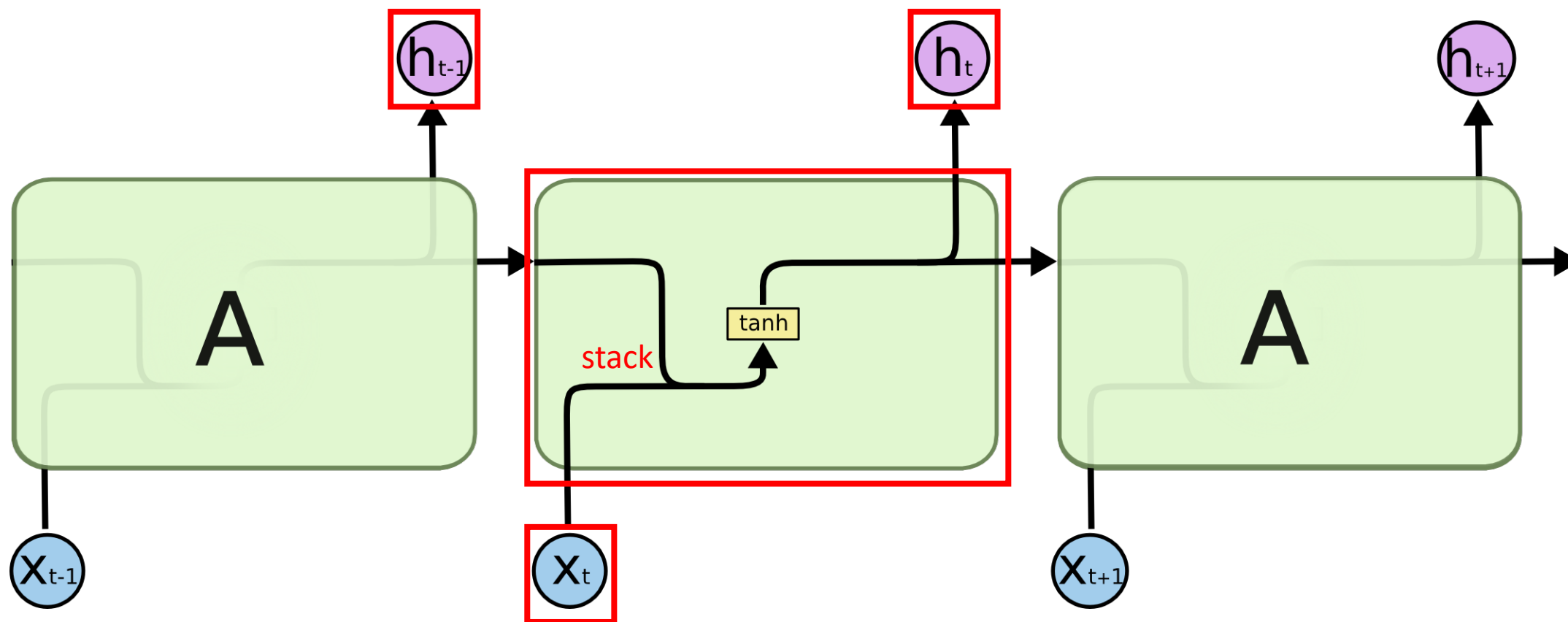
A conventional RNN



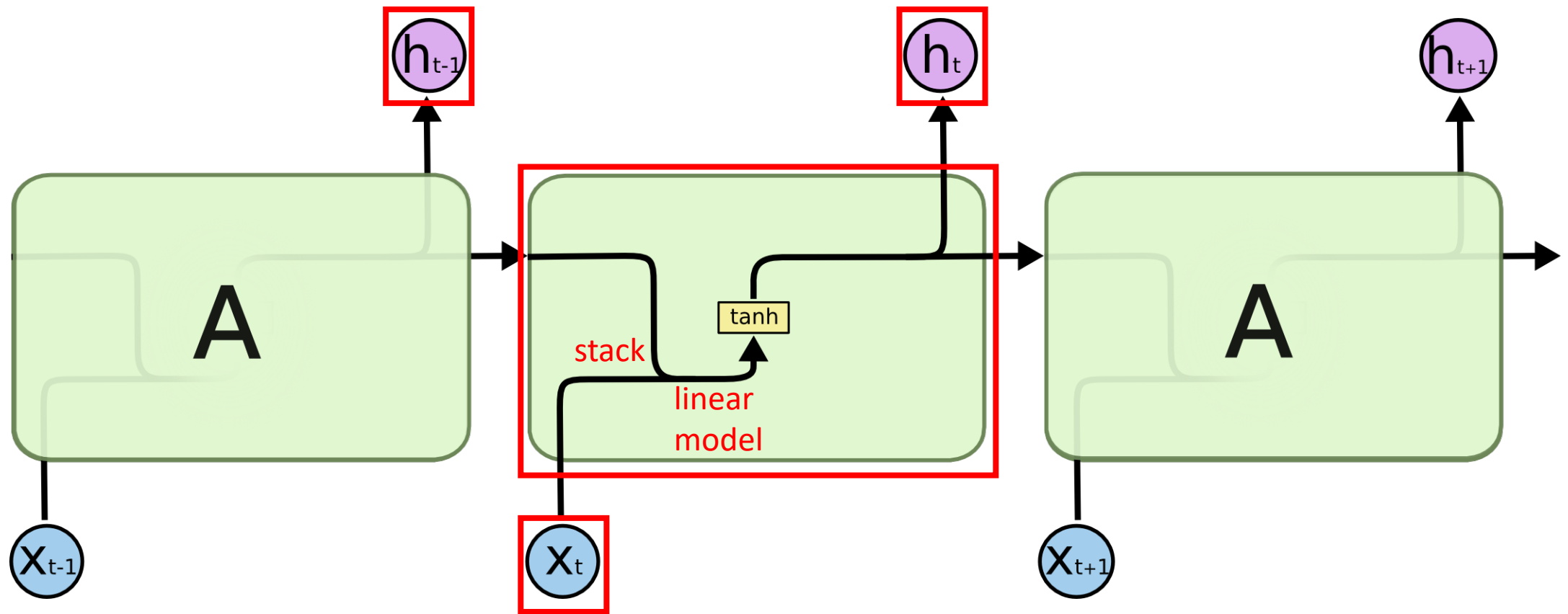
A conventional RNN



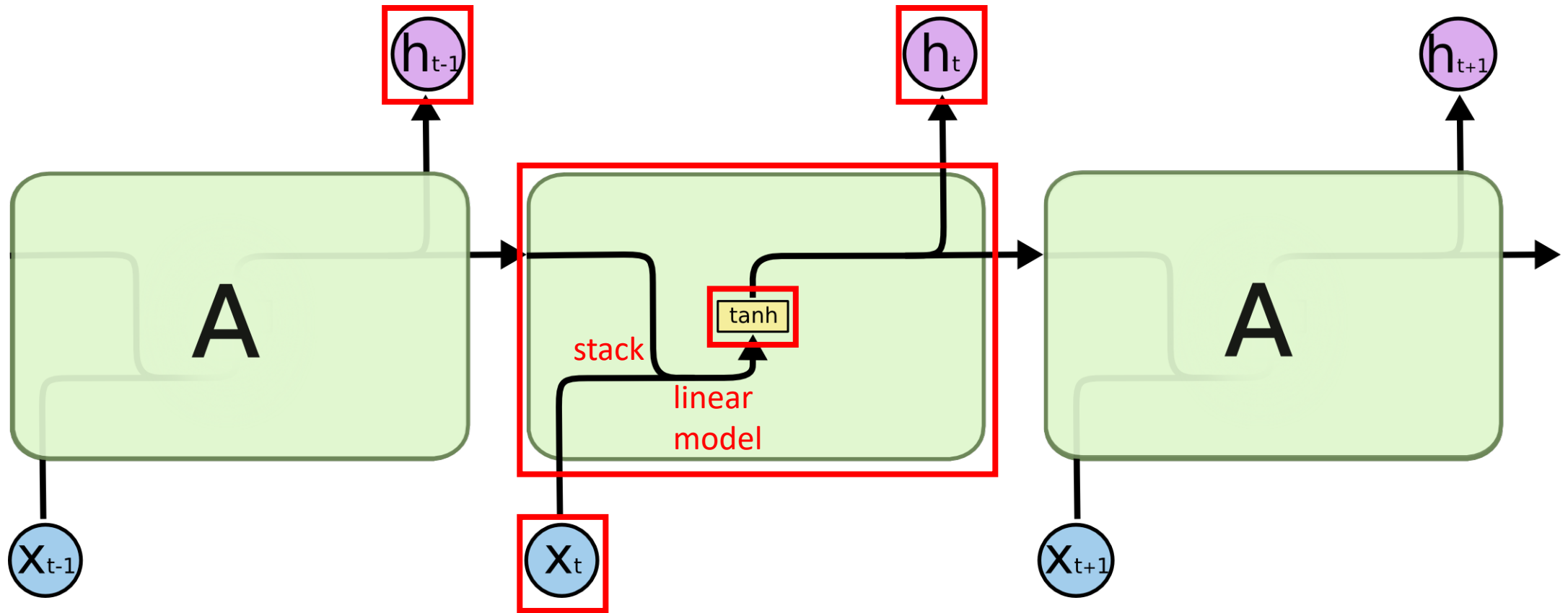
A conventional RNN



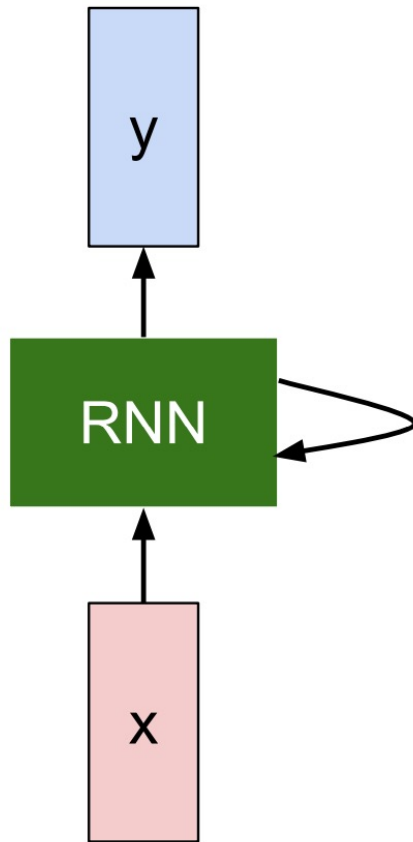
A conventional RNN



A conventional RNN



Recurrent neural networks



$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right]$$

$$f = \tanh(\cdot)$$

$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

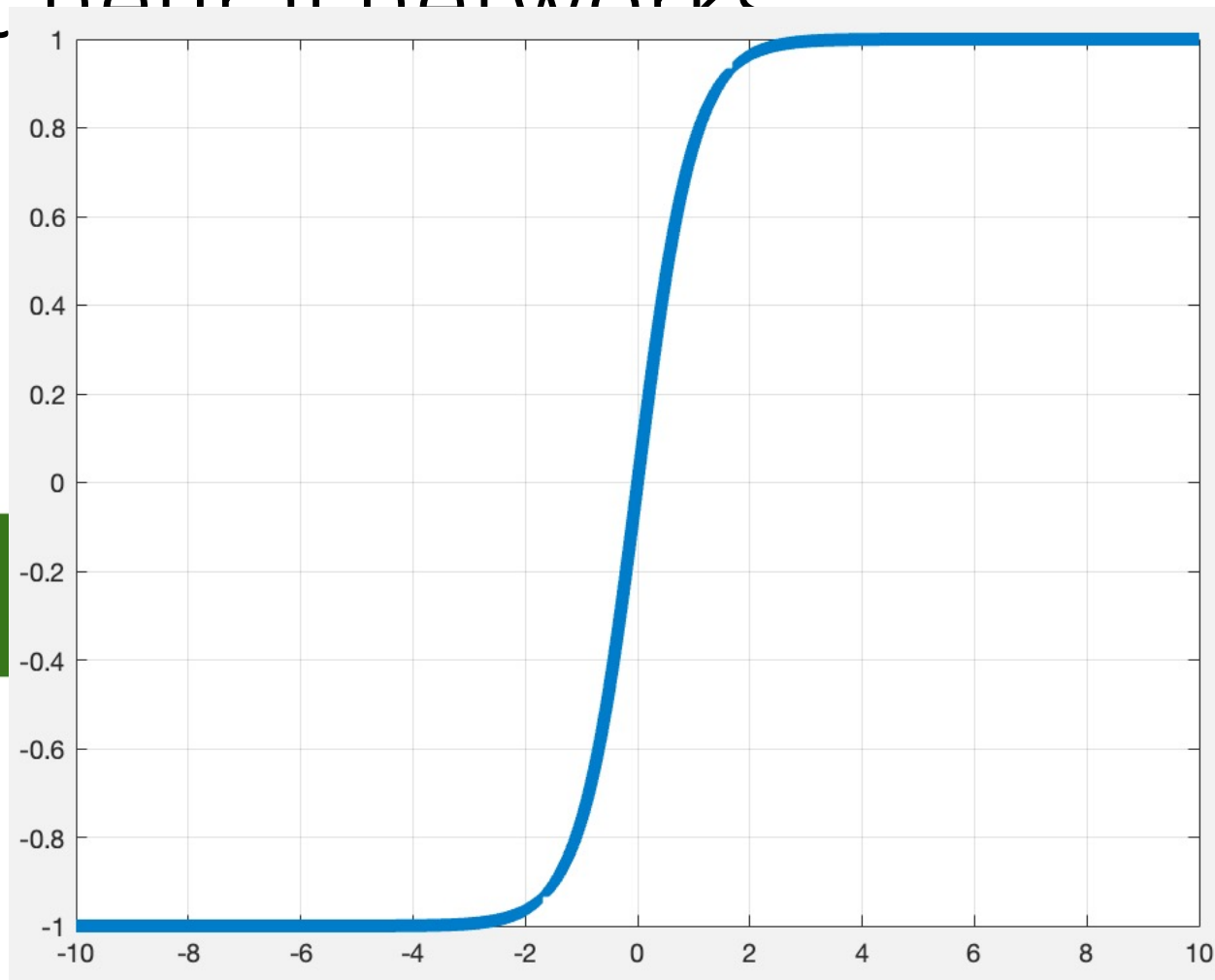
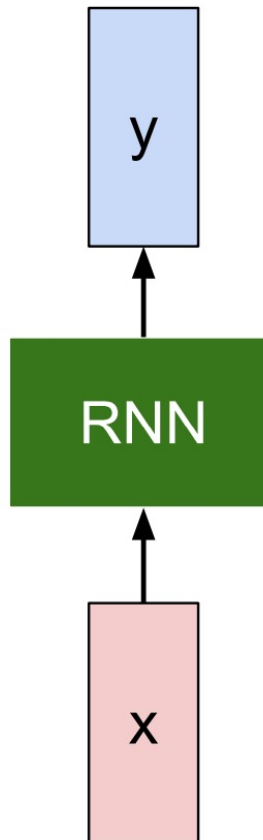
new state

some function with parameters W

old state

input vector at some time step

Recurrent neural networks



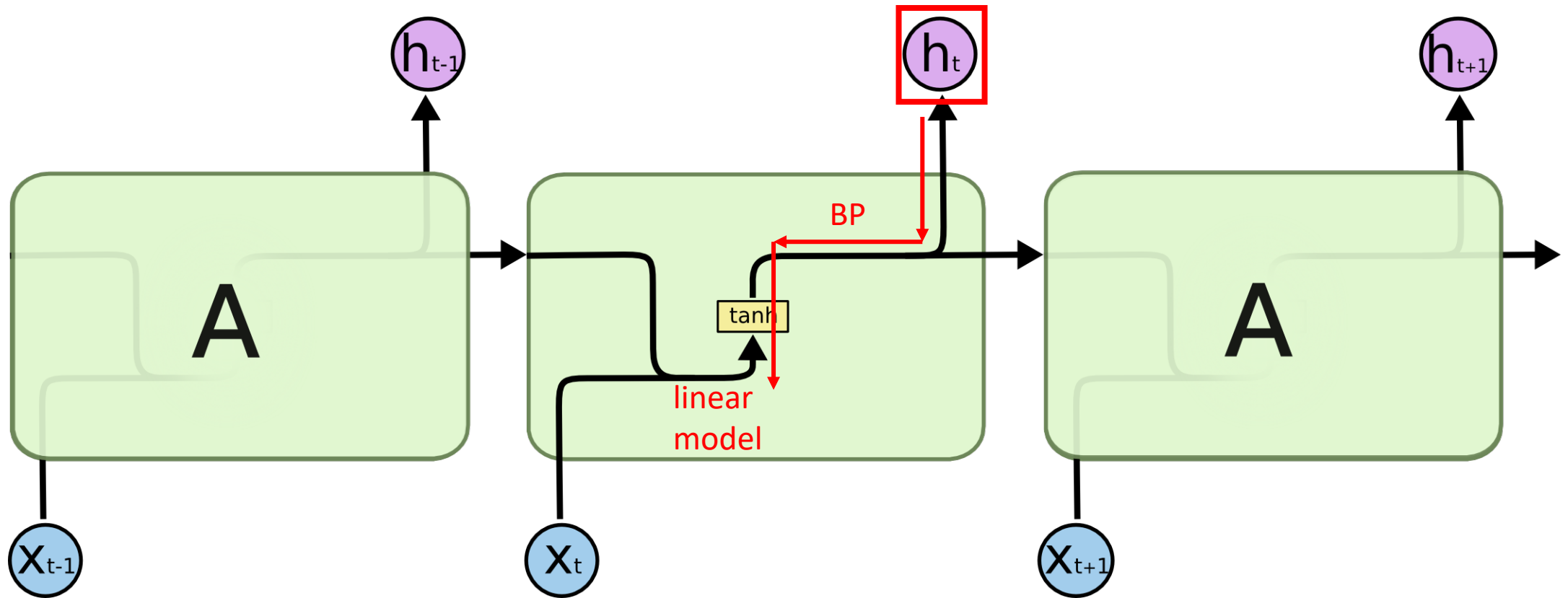
tanh function

$$f = \tanh(\cdot)$$

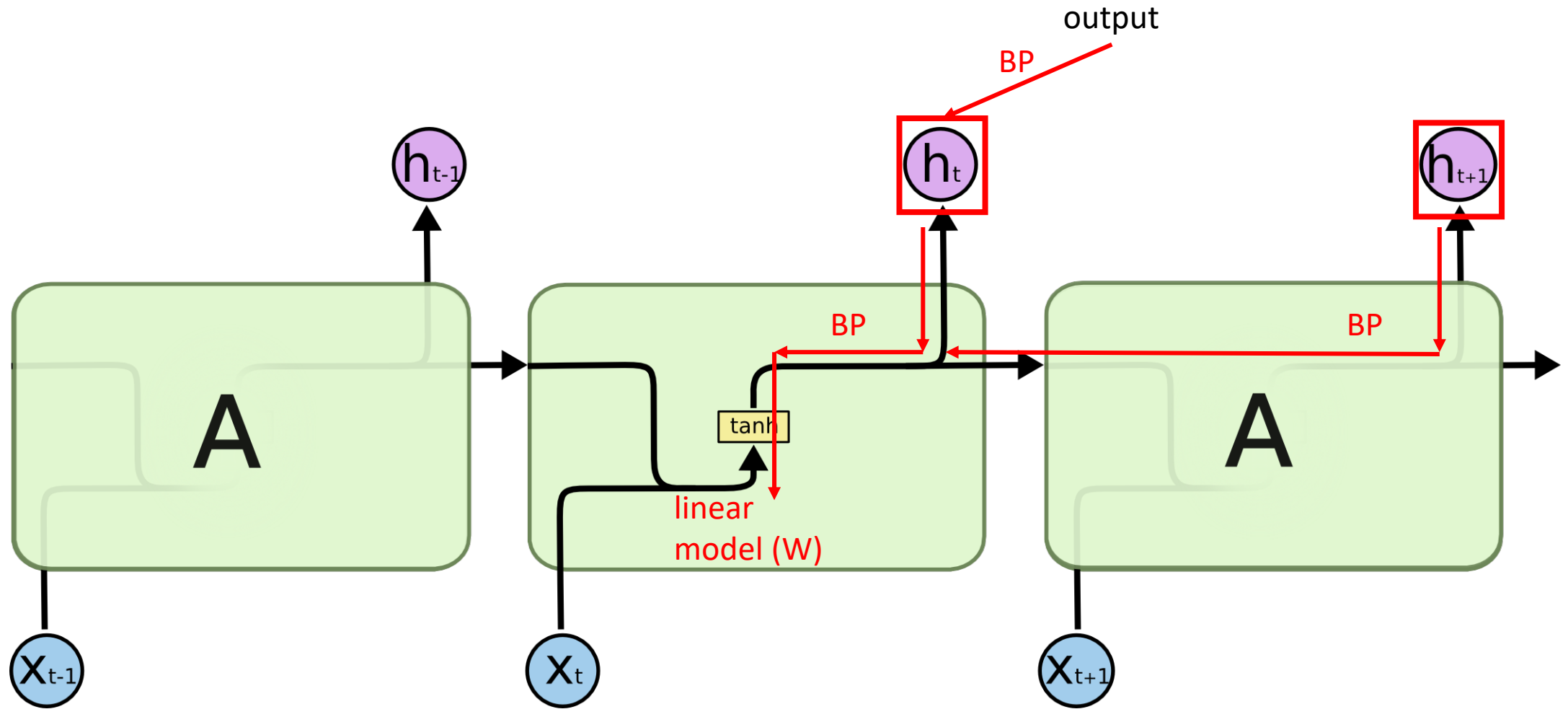
c_t

put vector at
some time step

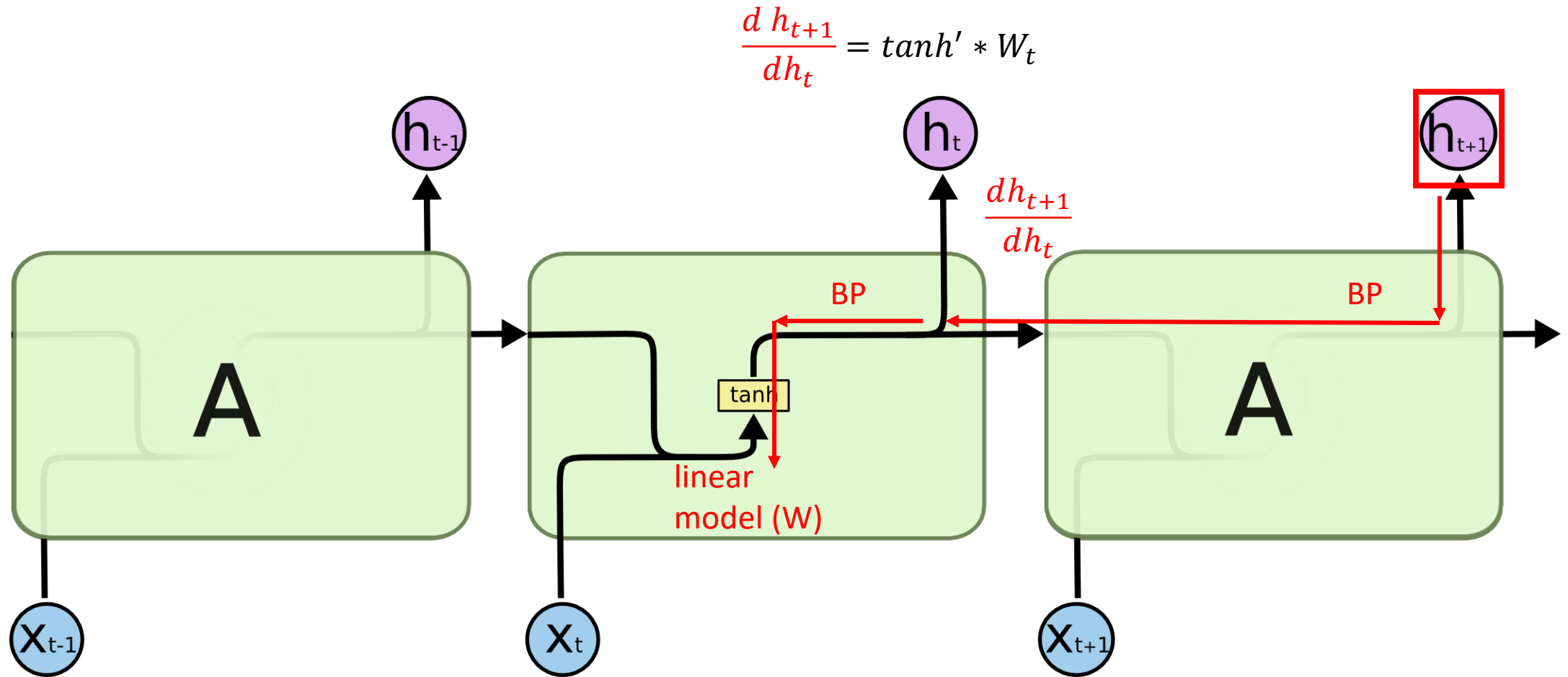
A conventional RNN



A conventional RNN

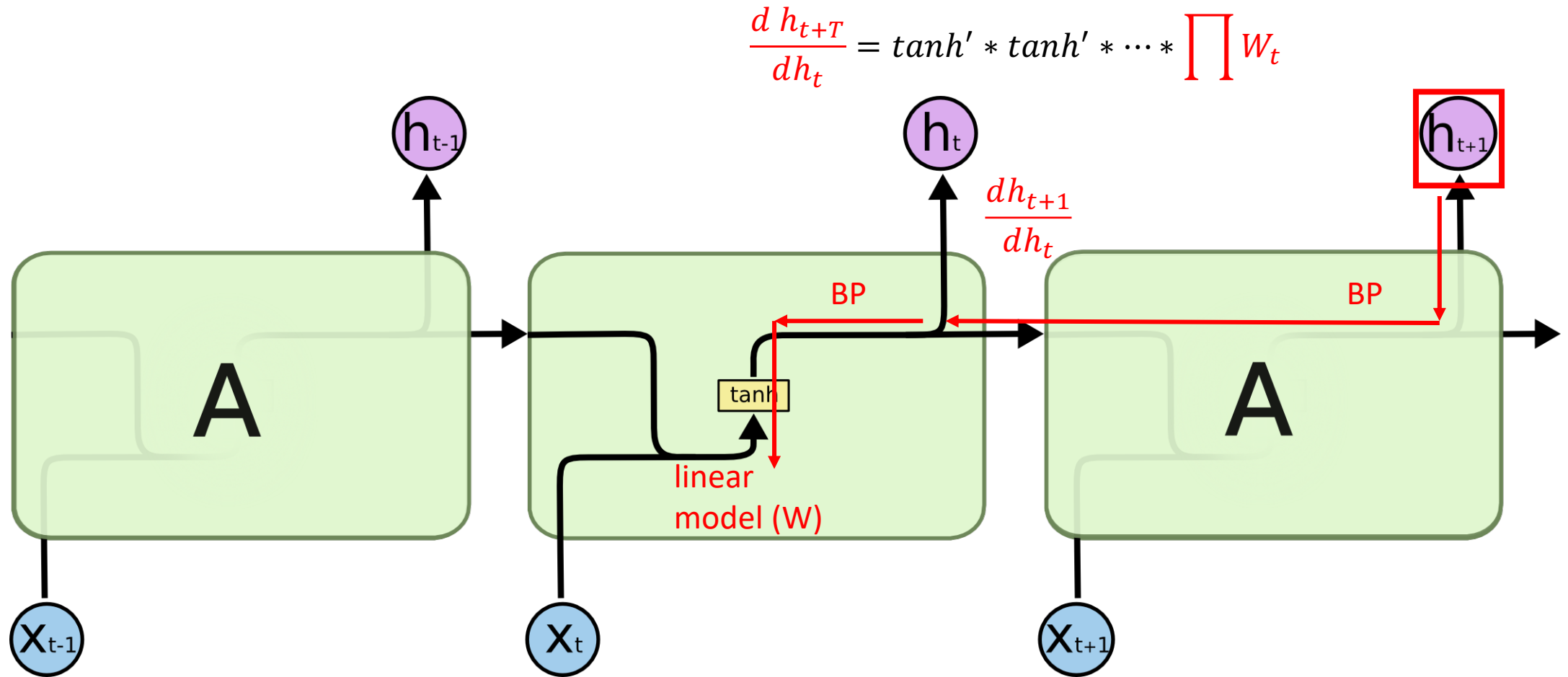


A conventional RNN



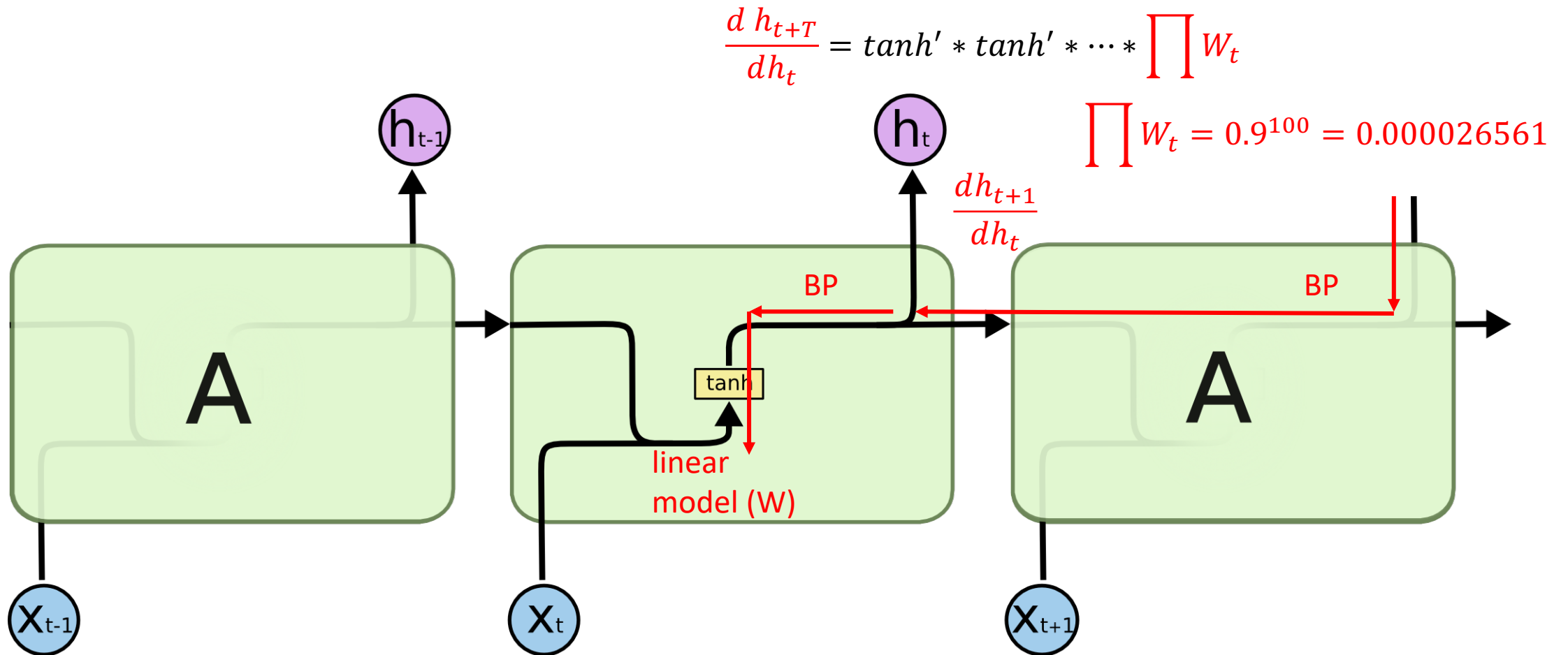
A conventional RNN

Q: what issue will we have?



A conventional RNN

Q: what issue will we have?



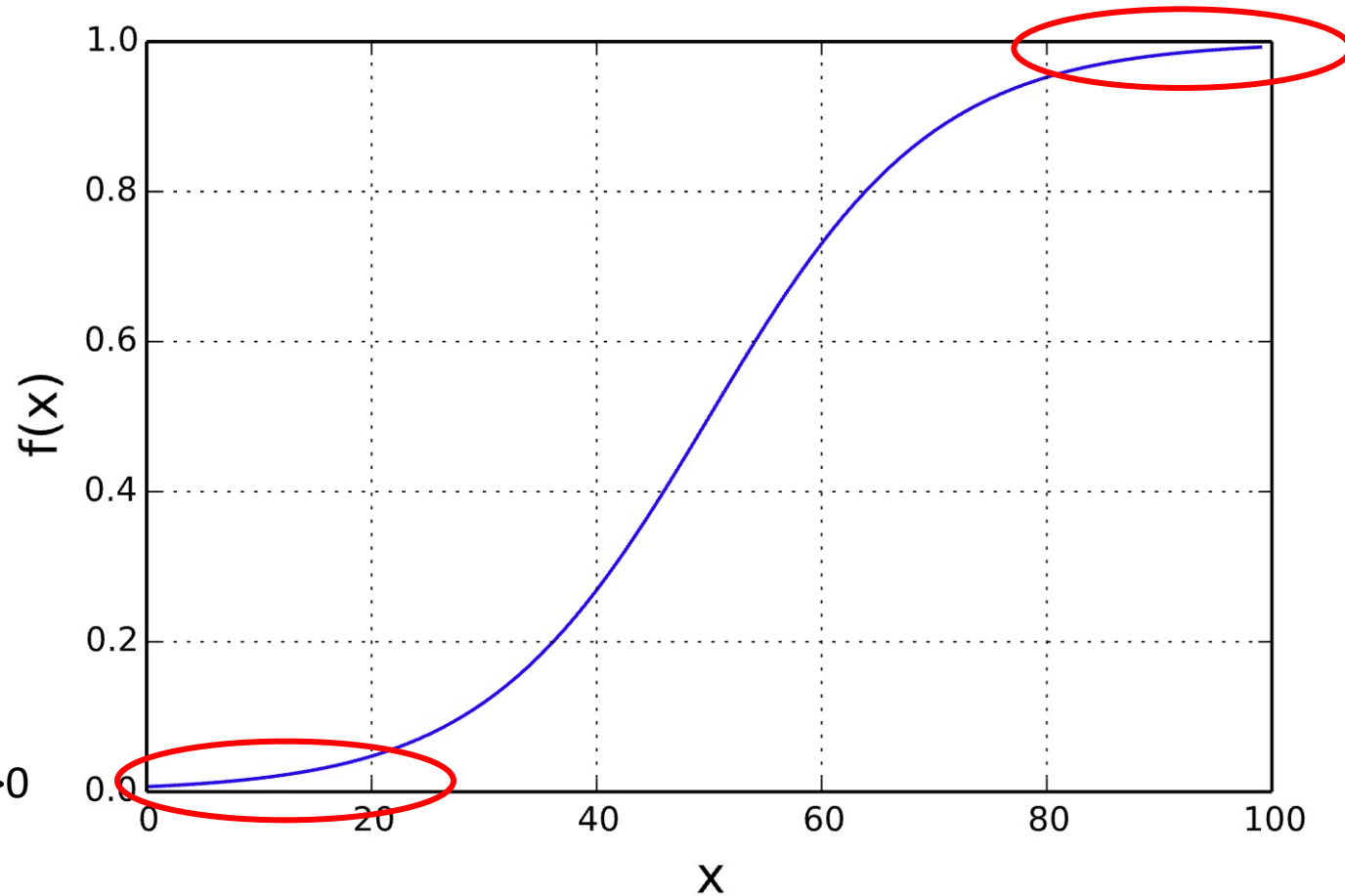
Gradient vanish

$$f_n \left(\dots \left(f_2(f_1(x)) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

gradients $\rightarrow 0$



gradients $\rightarrow 0$

Sigmoid function

ResNet: shortcut connection

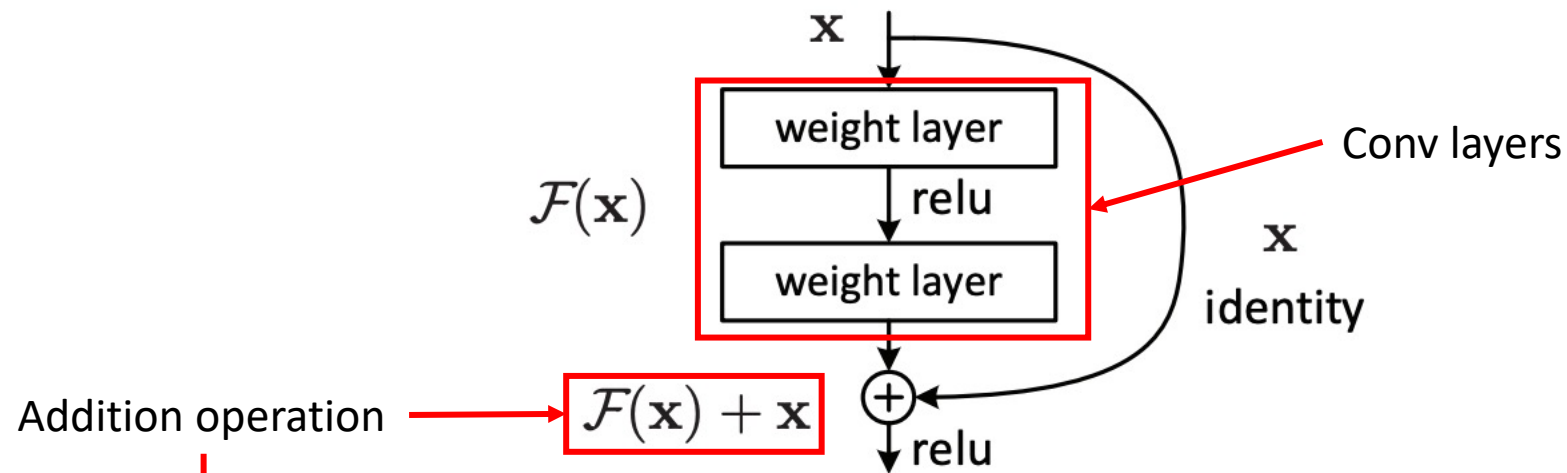
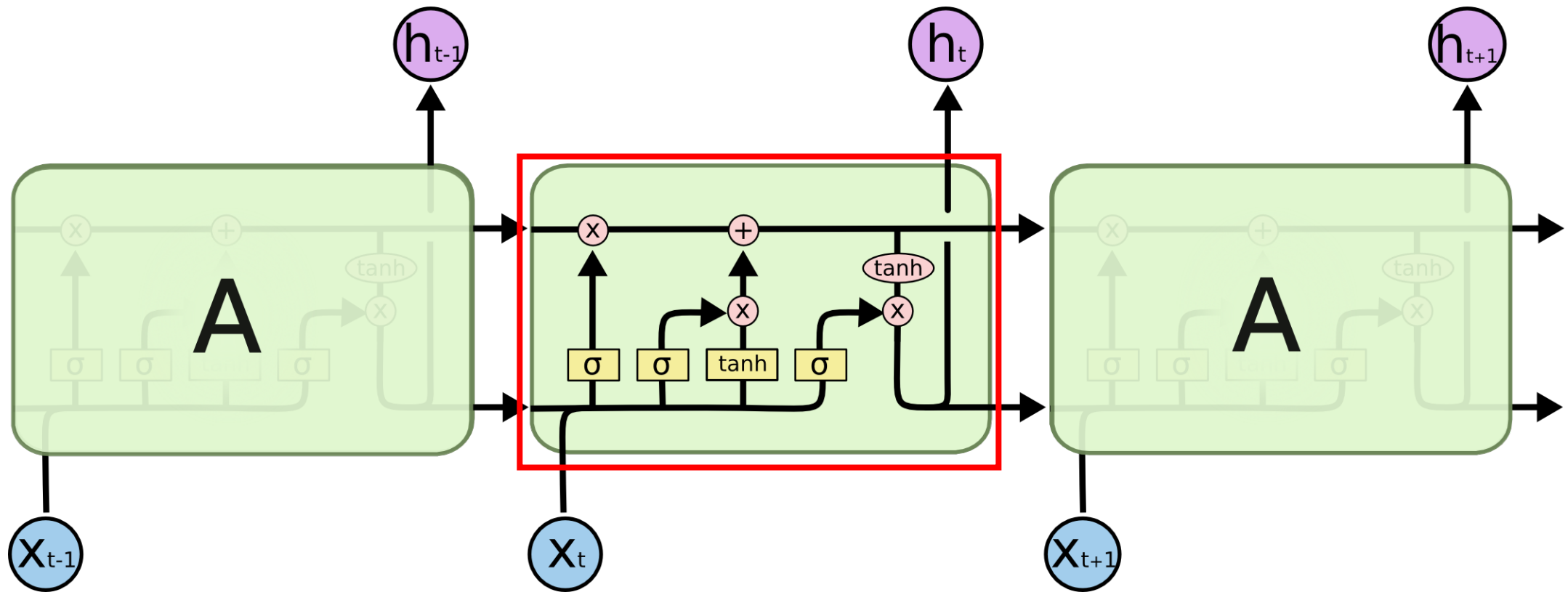


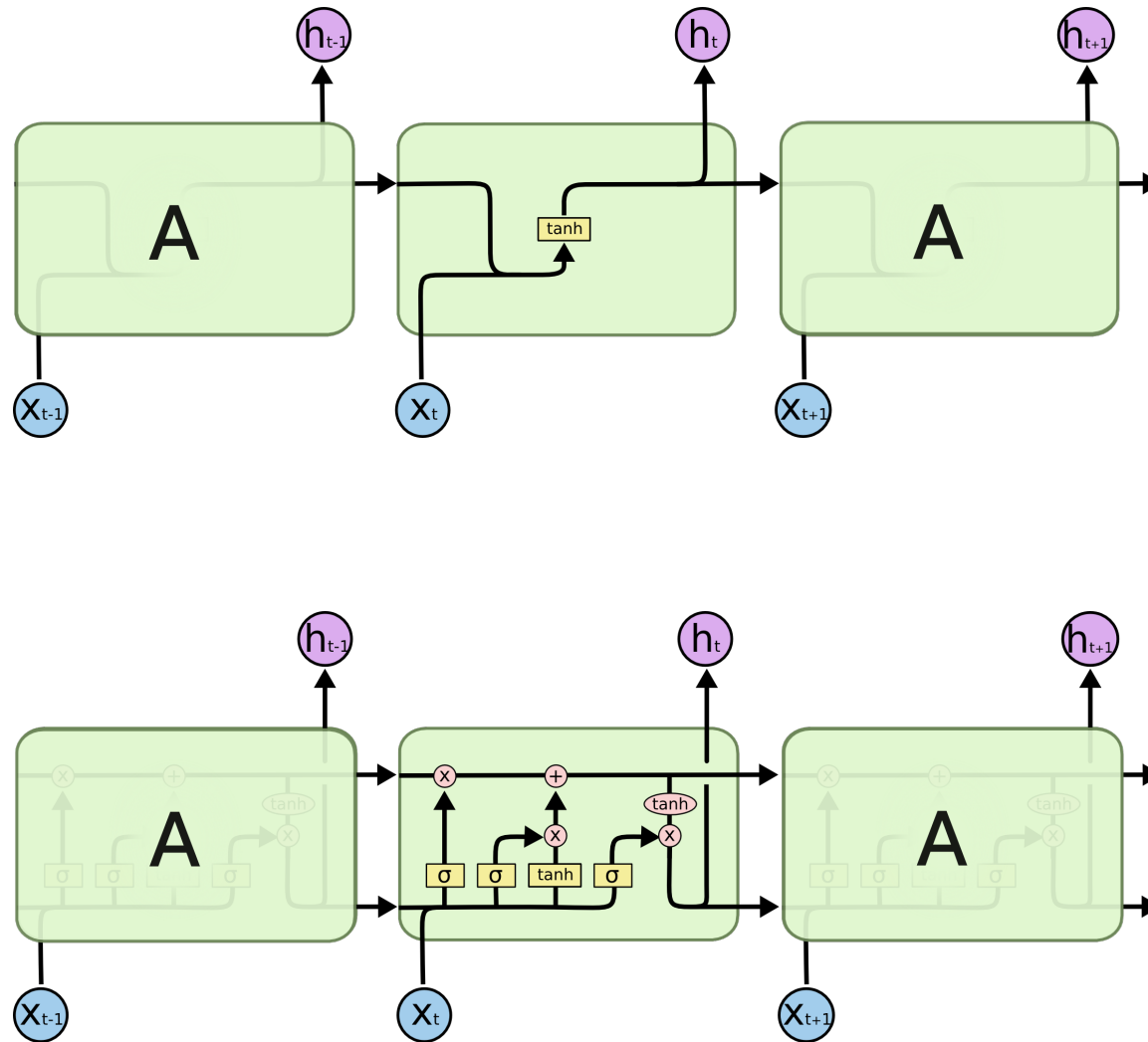
Figure 2. Residual learning: a building block.

implication: same dimension

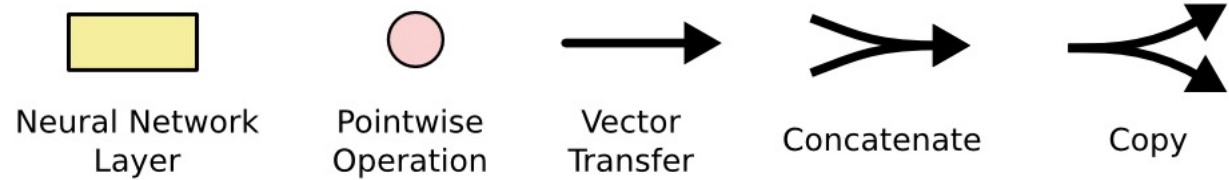
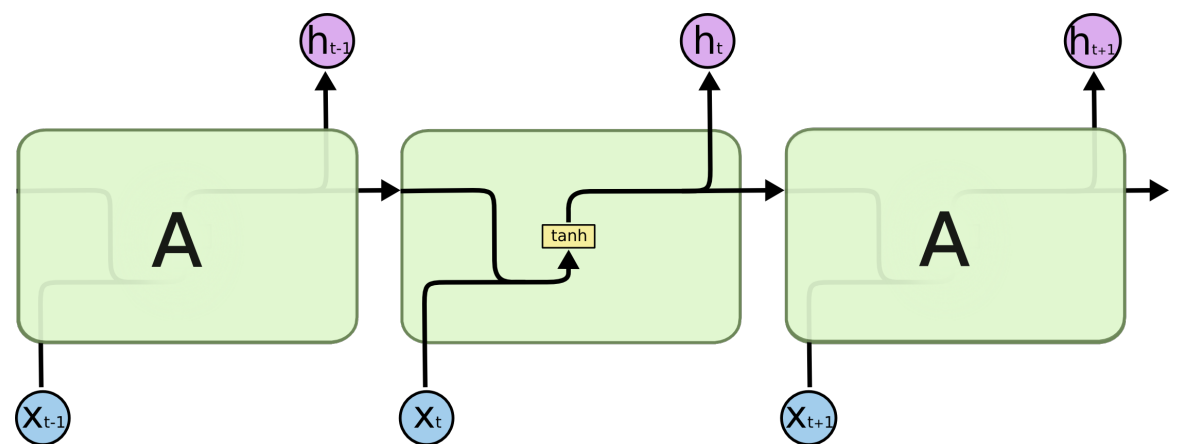
Long-short term memory (LSTM) networks



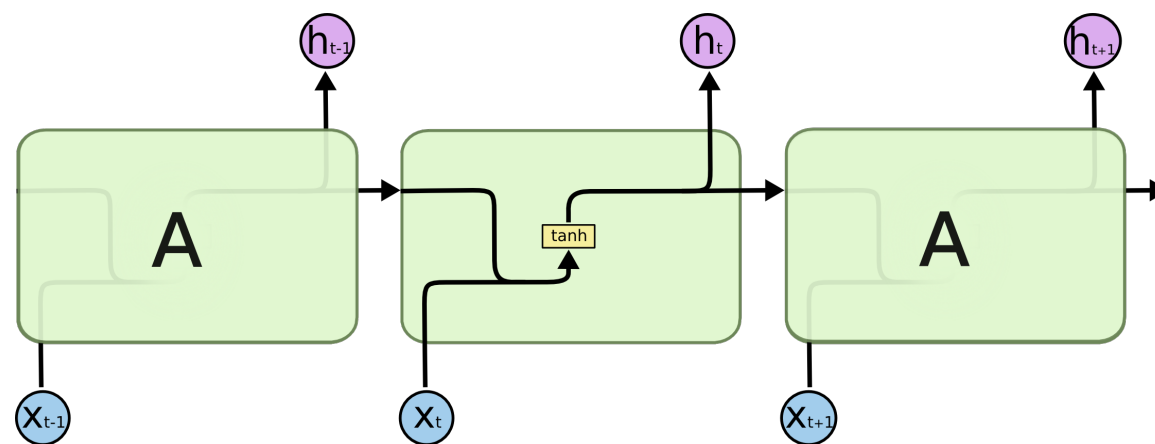
LSTM




LSTM

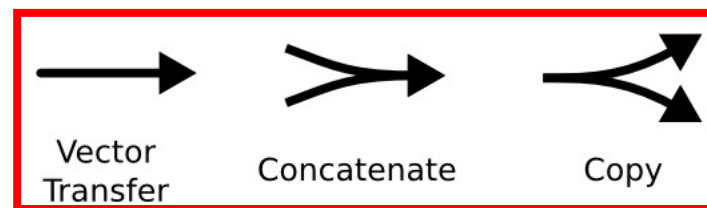


LSTM

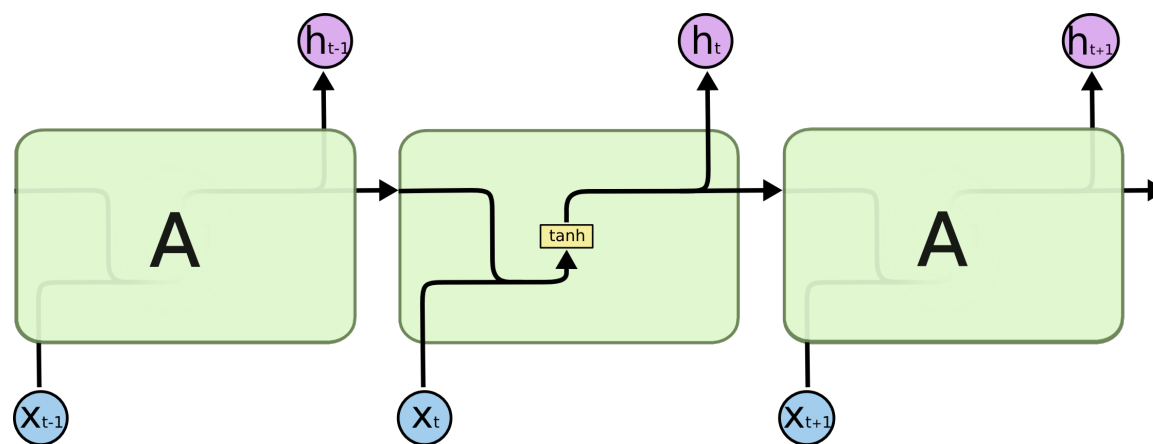



Neural Network
Layer



Pointwise
Operation



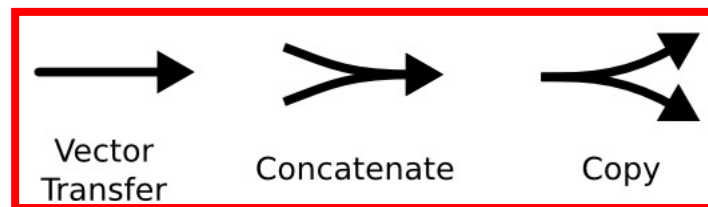
LSTM



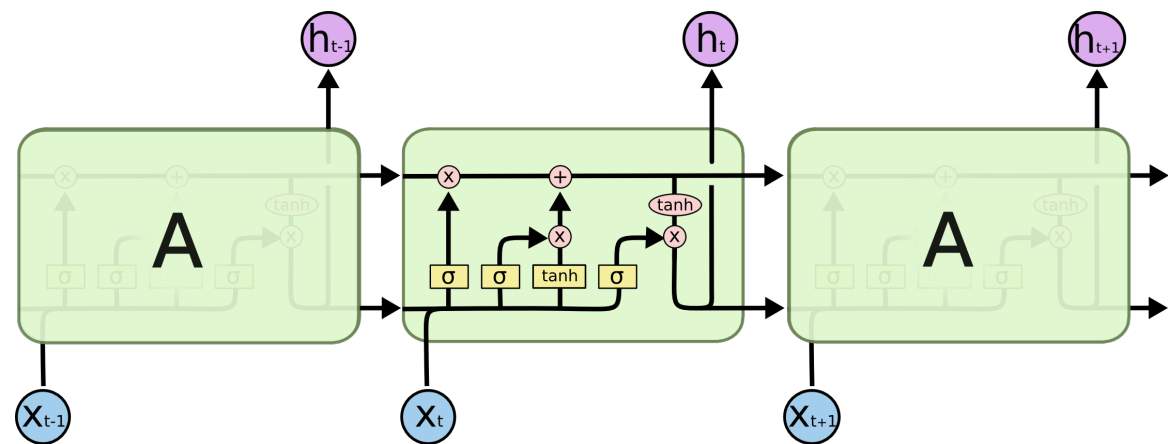
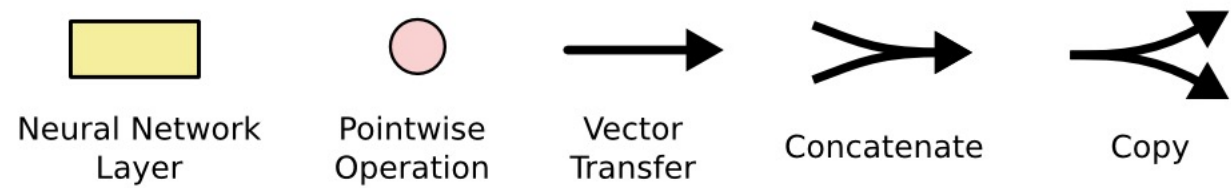
Composition of all $\{h_t\}$


Neural Network
Layer

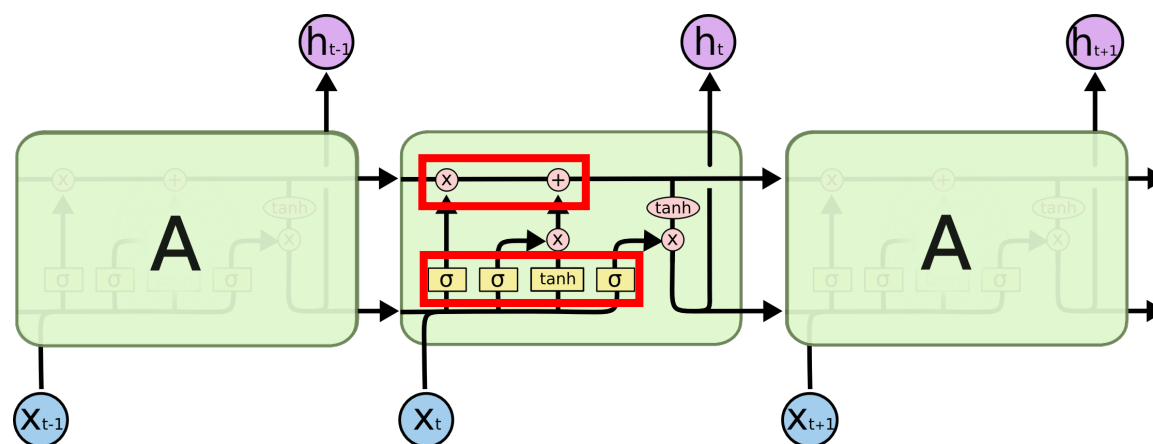
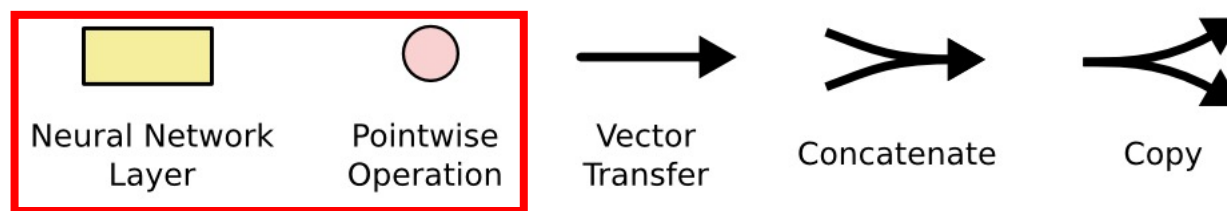

Pointwise
Operation



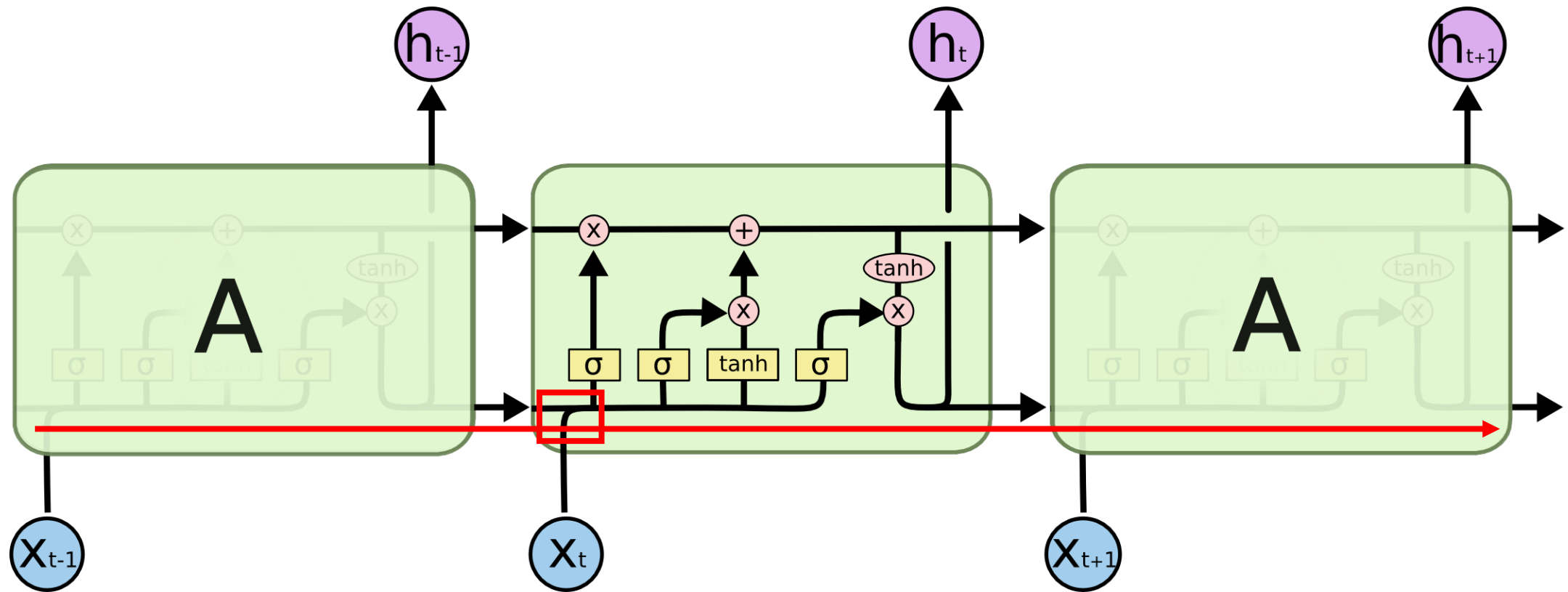
LSTM



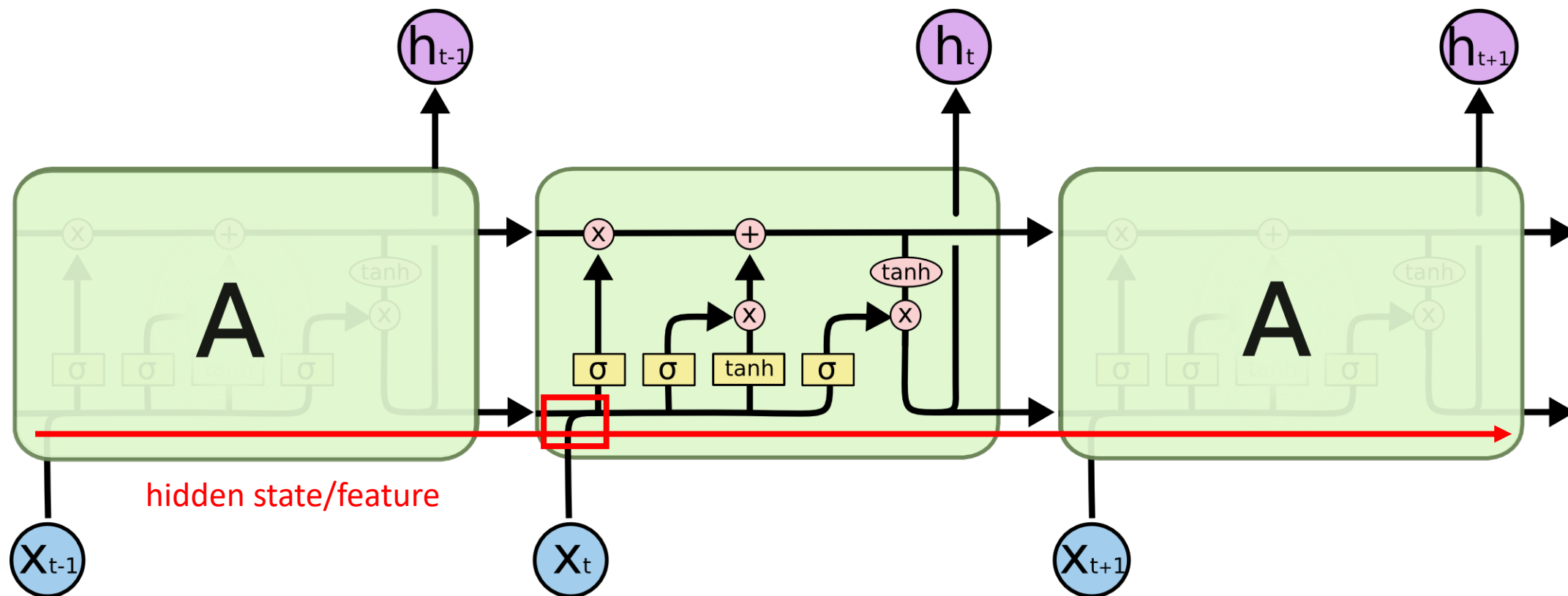
LSTM



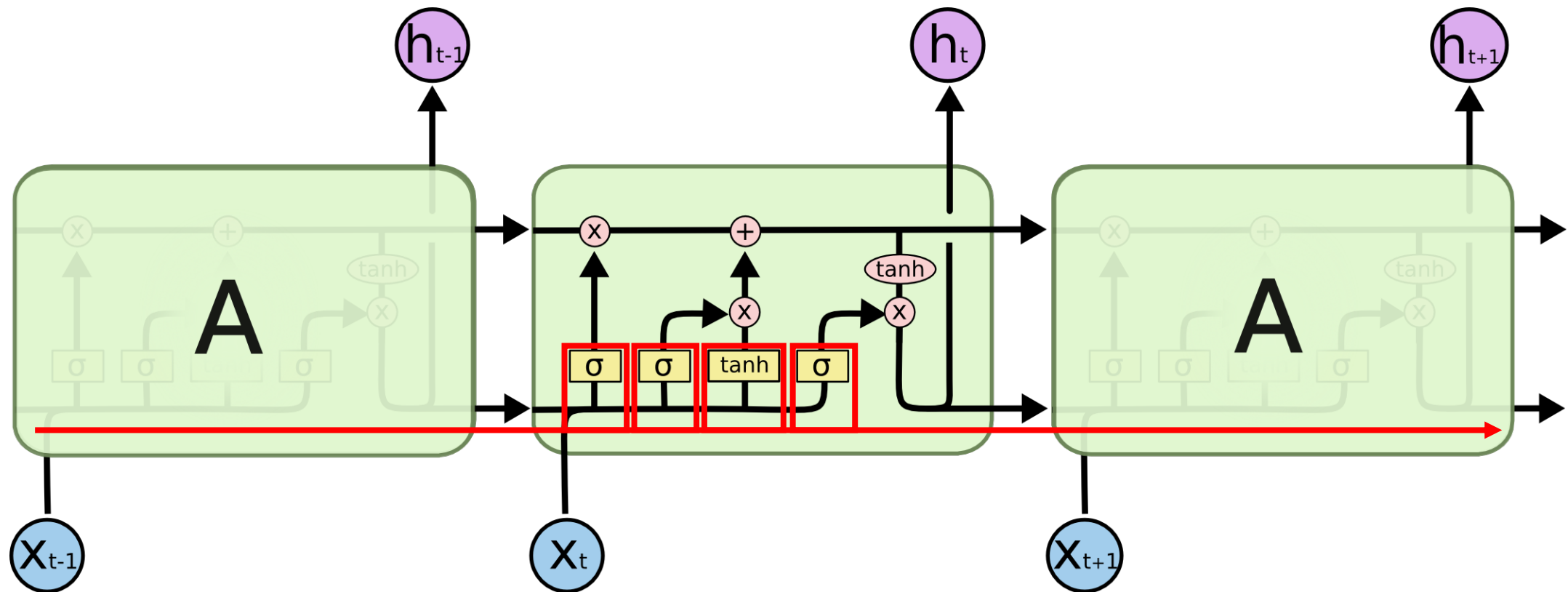
LSTM



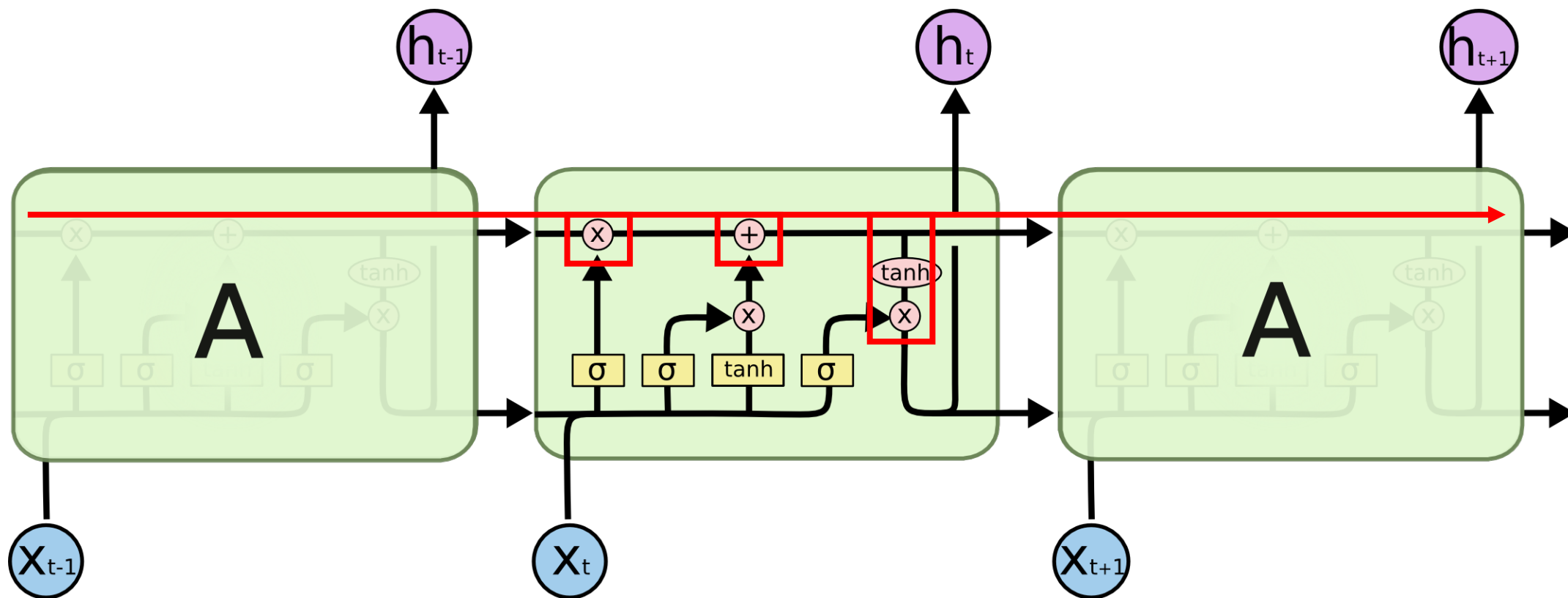
LSTM



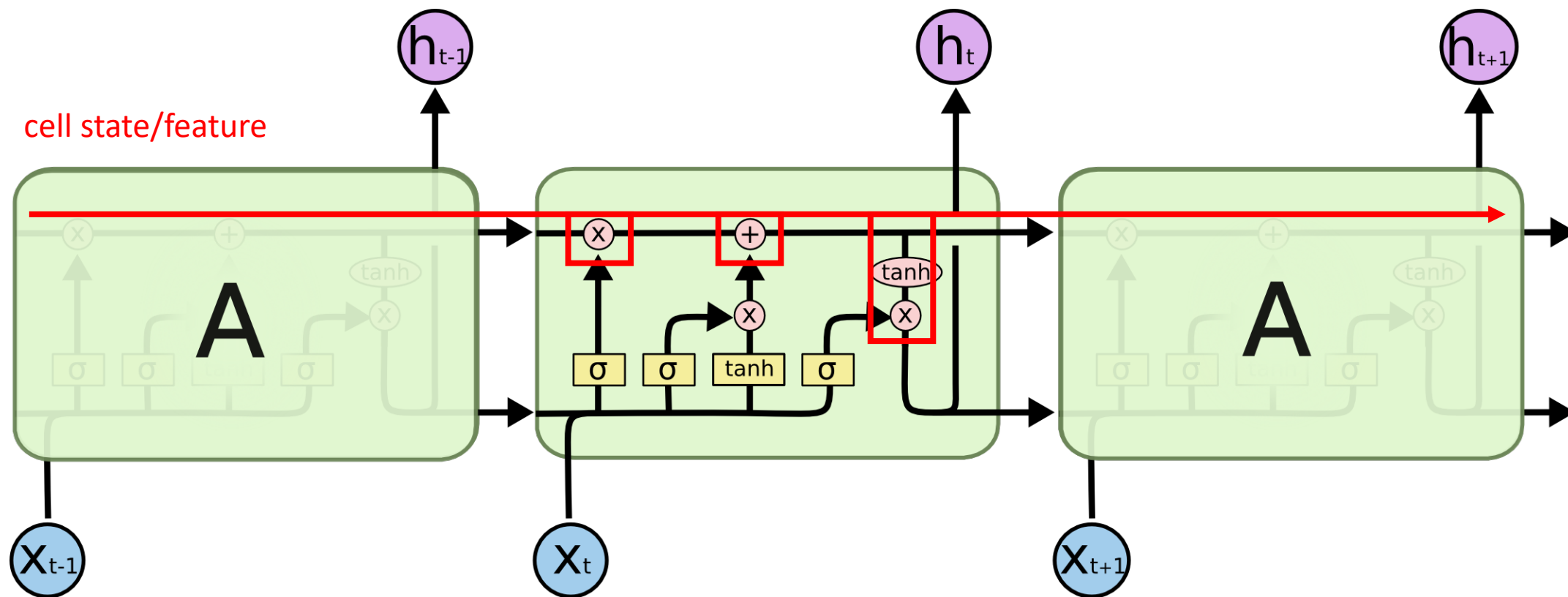
LSTM



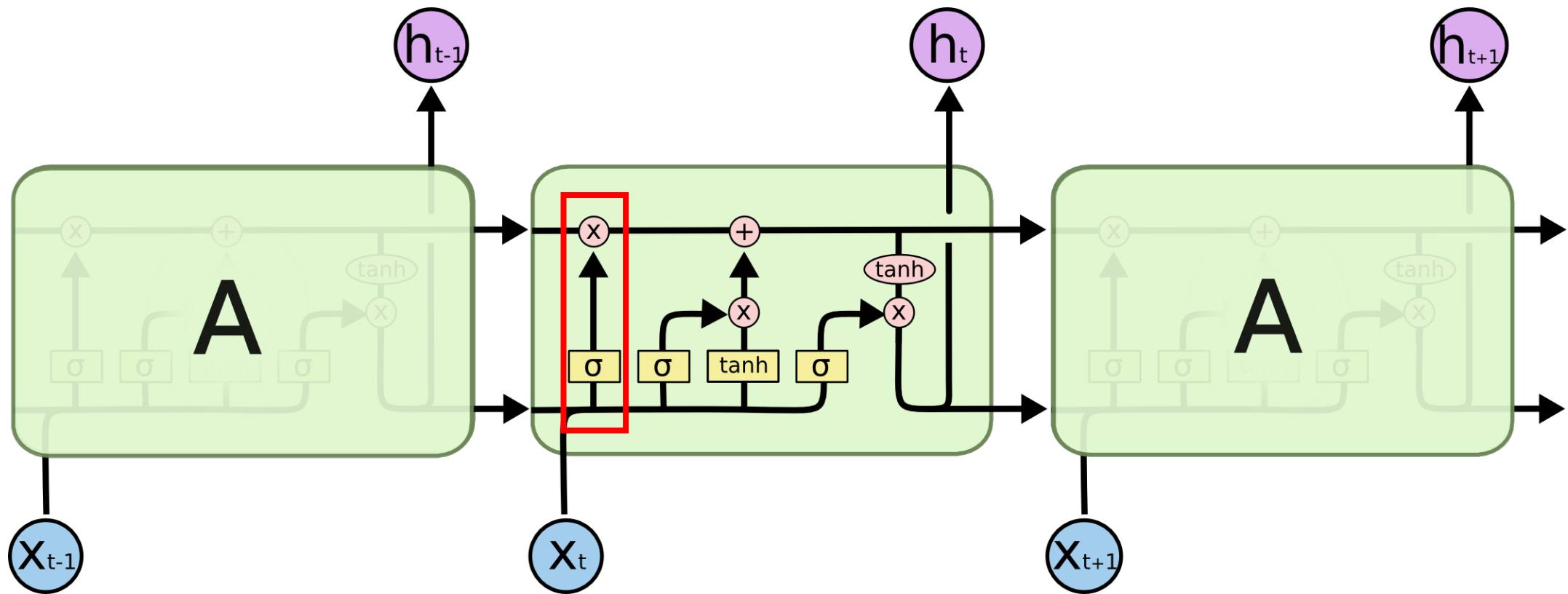
LSTM



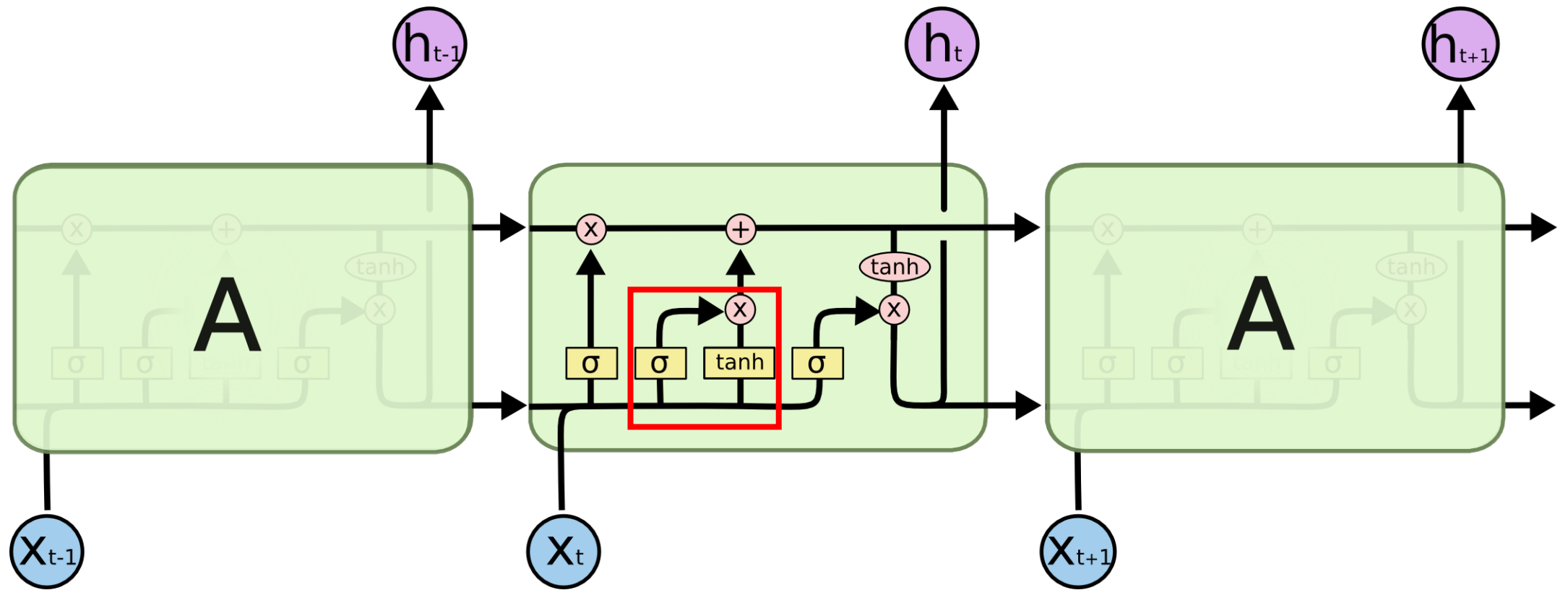
LSTM



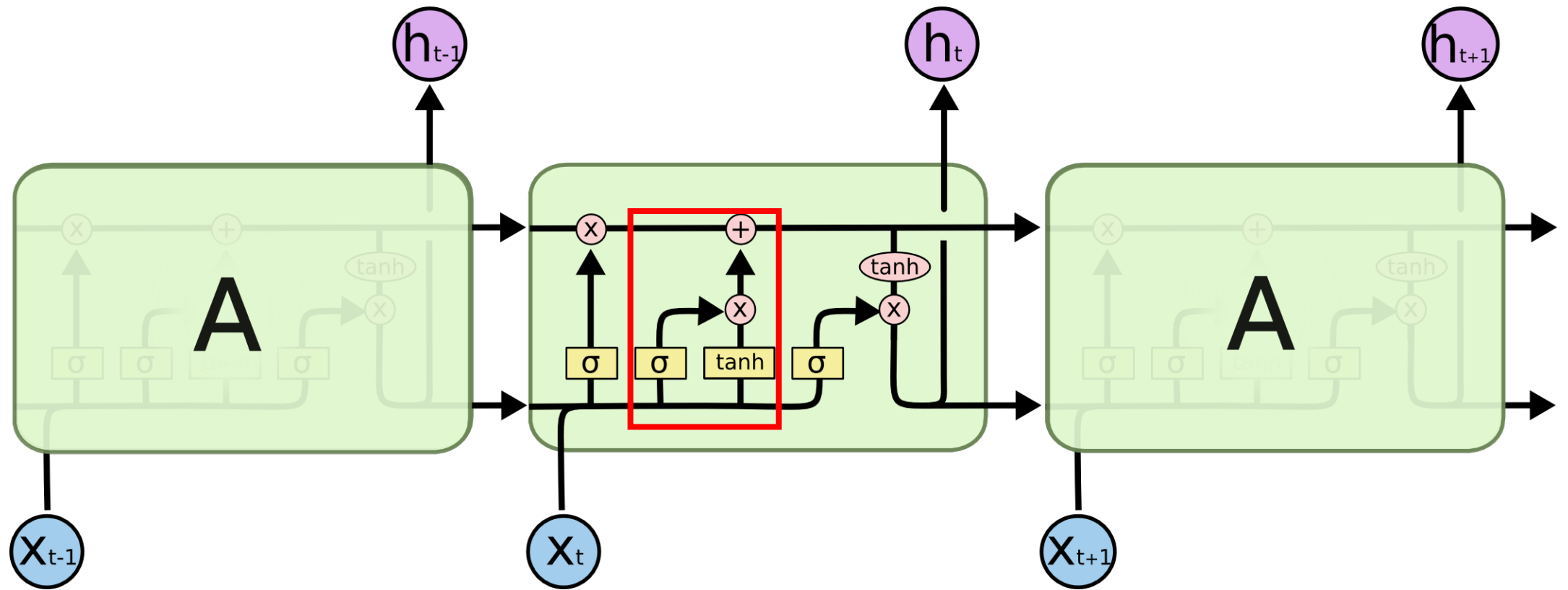
LSTM



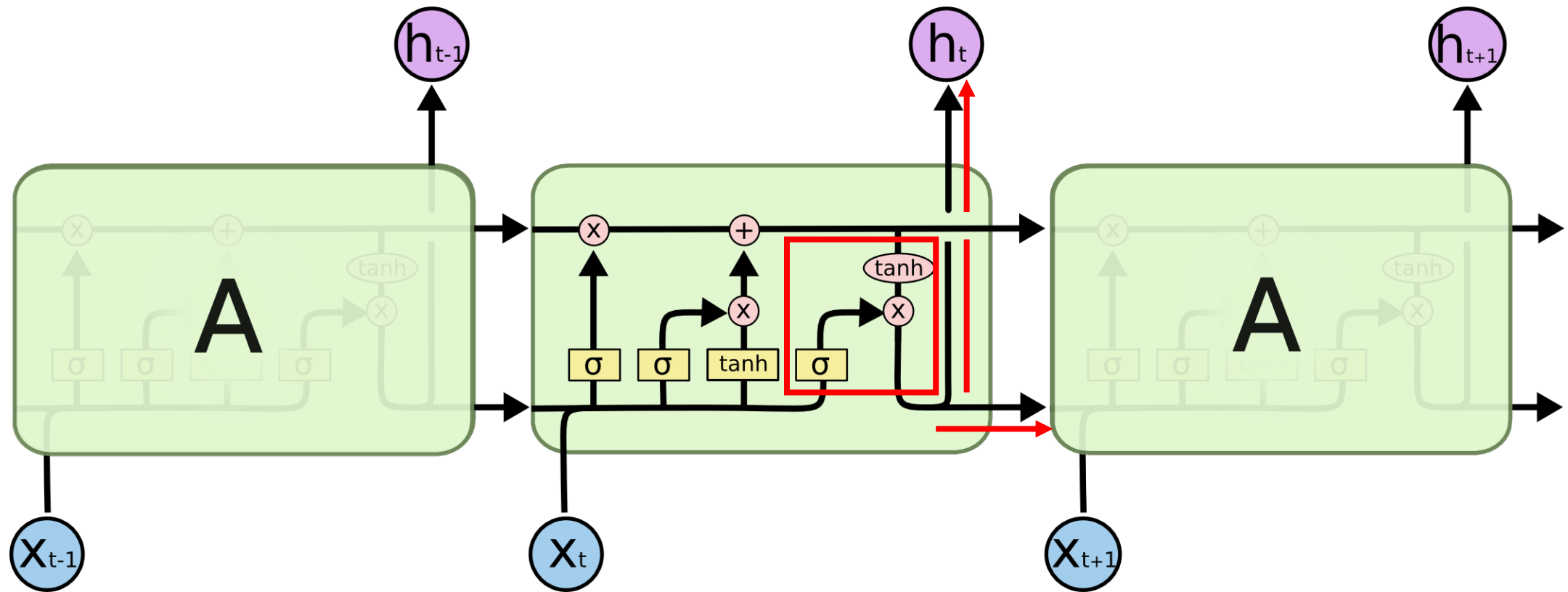
LSTM



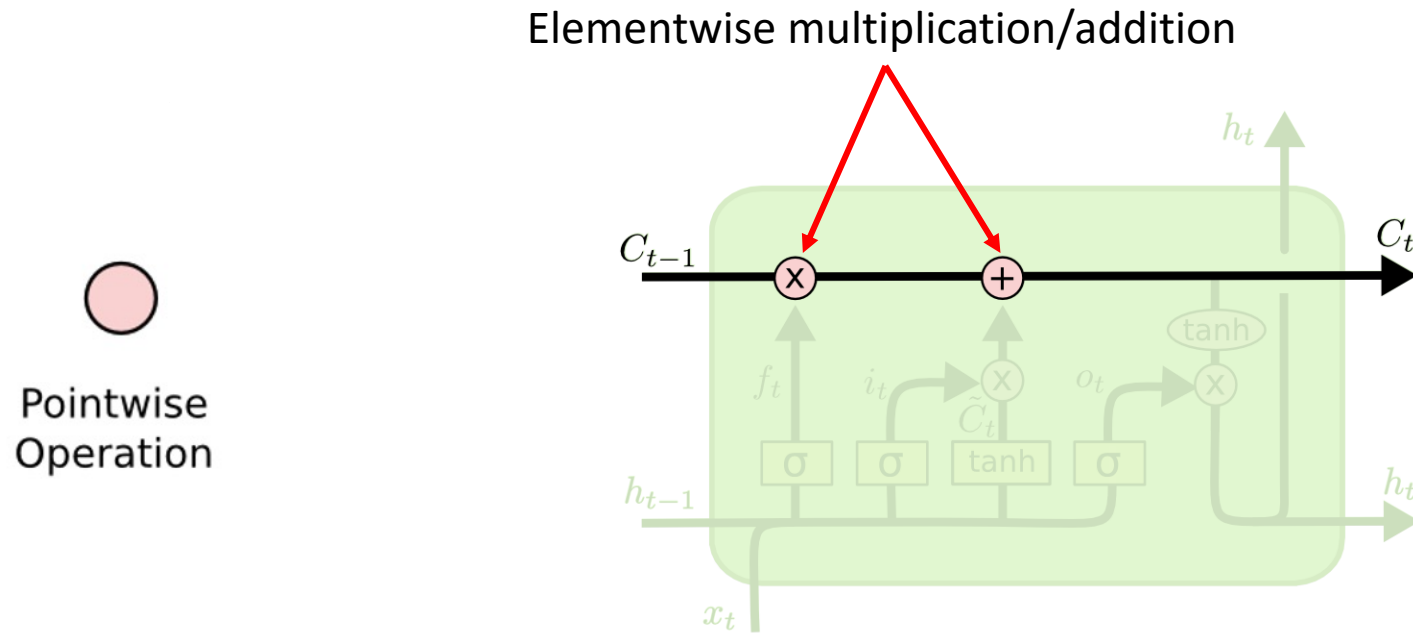
LSTM



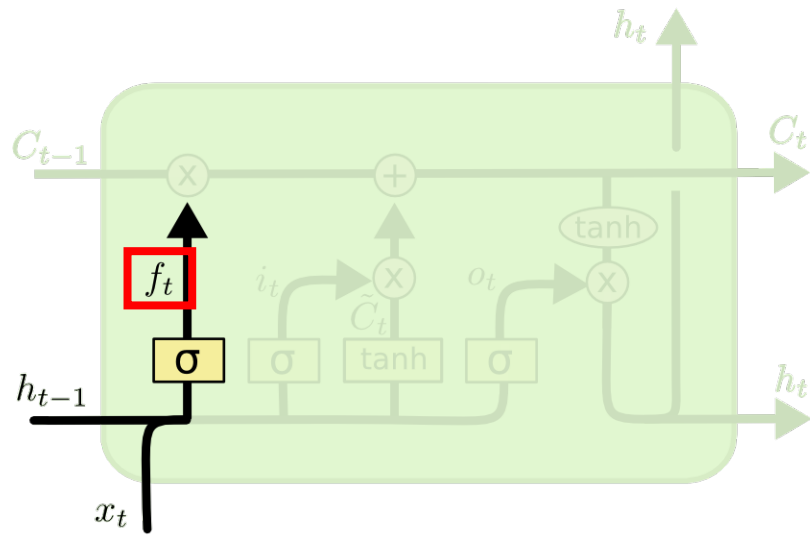
LSTM



LSTM

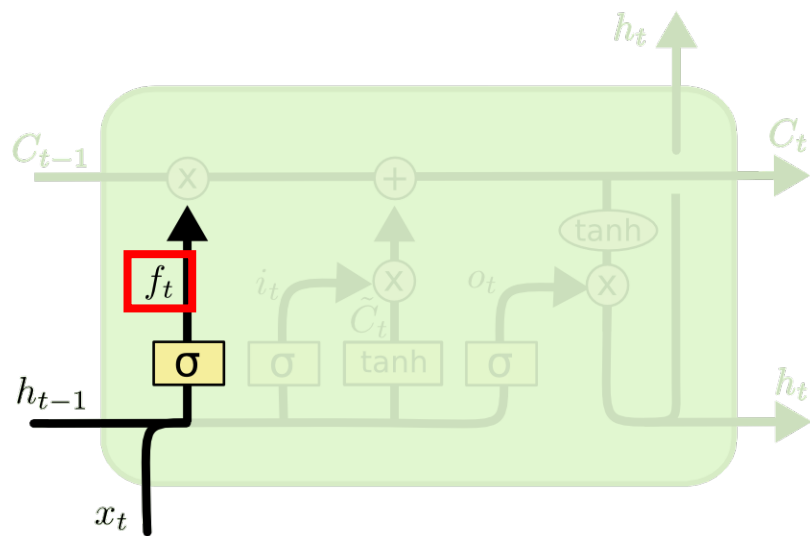


LSTM

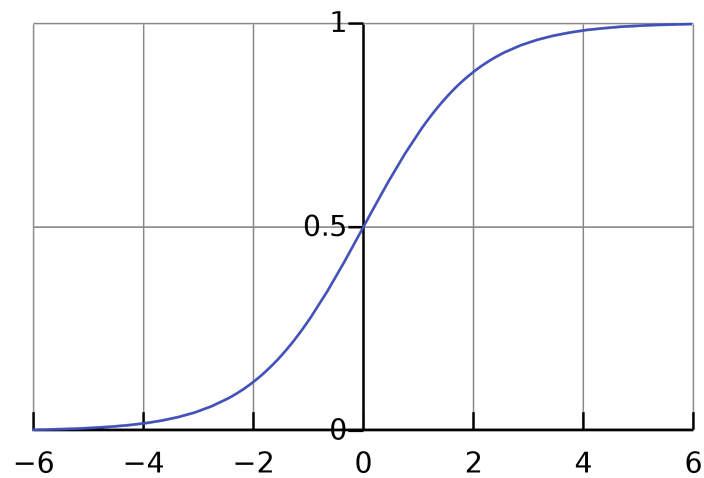


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

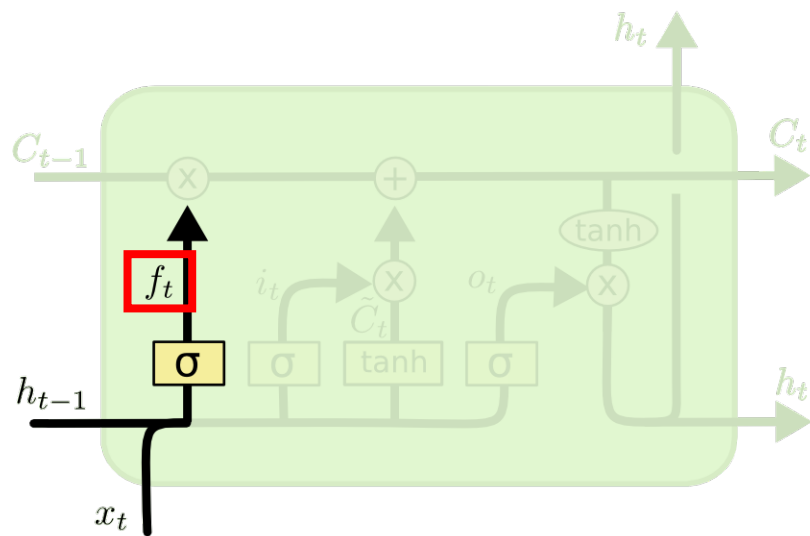
LSTM



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



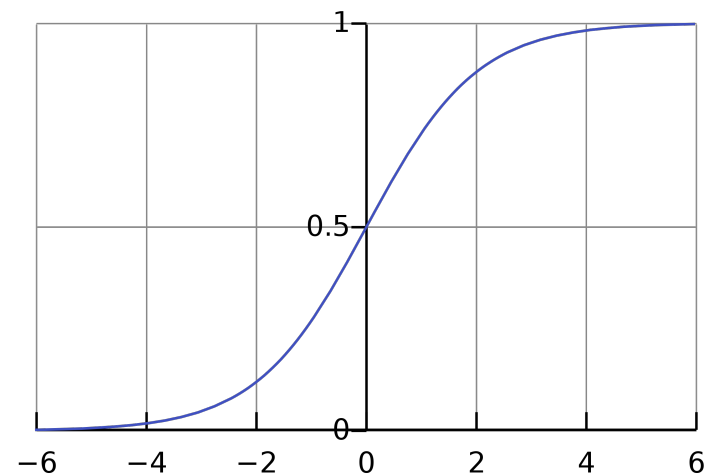
LSTM



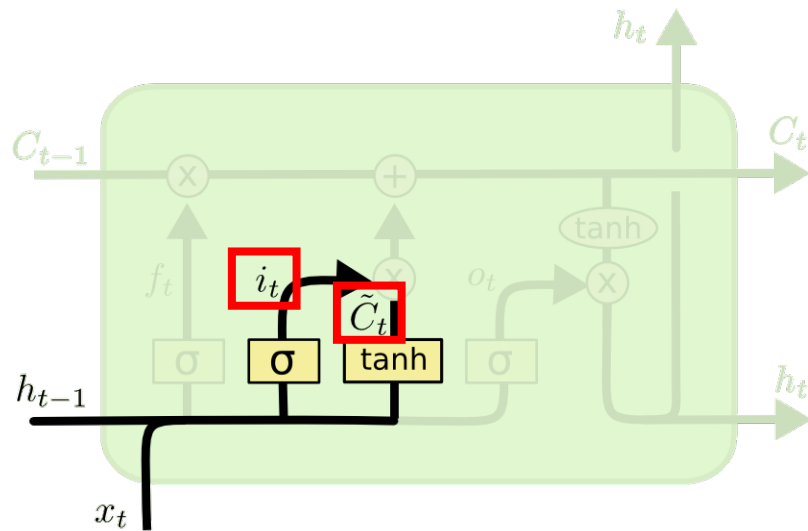
Forget gate:

Whether to erase cell

$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$



LSTM



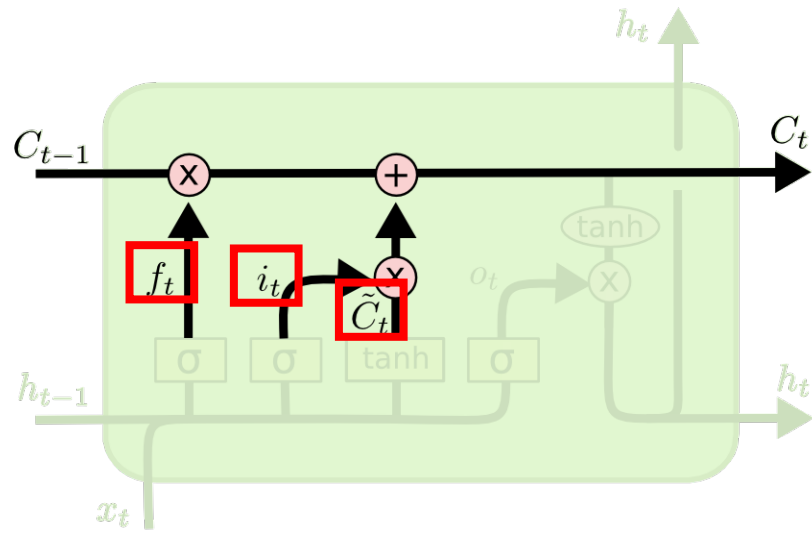
Input gate

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

cell input activation vector

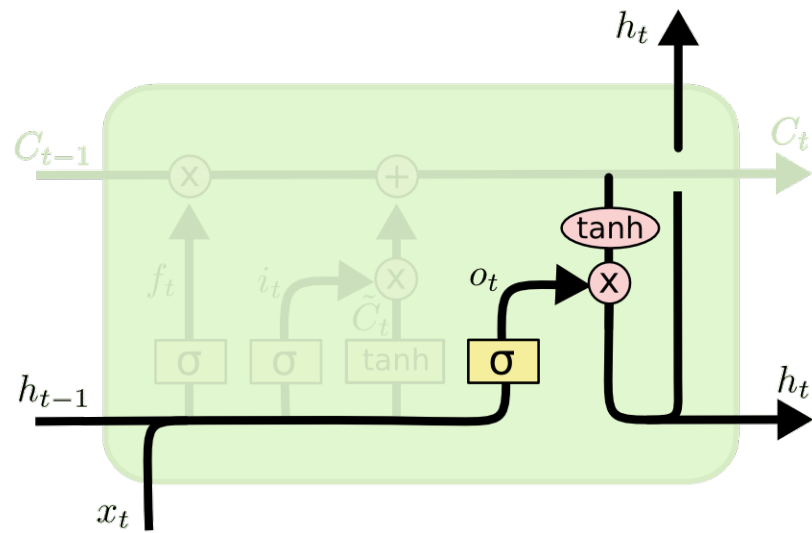
LSTM



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Cell state

LSTM

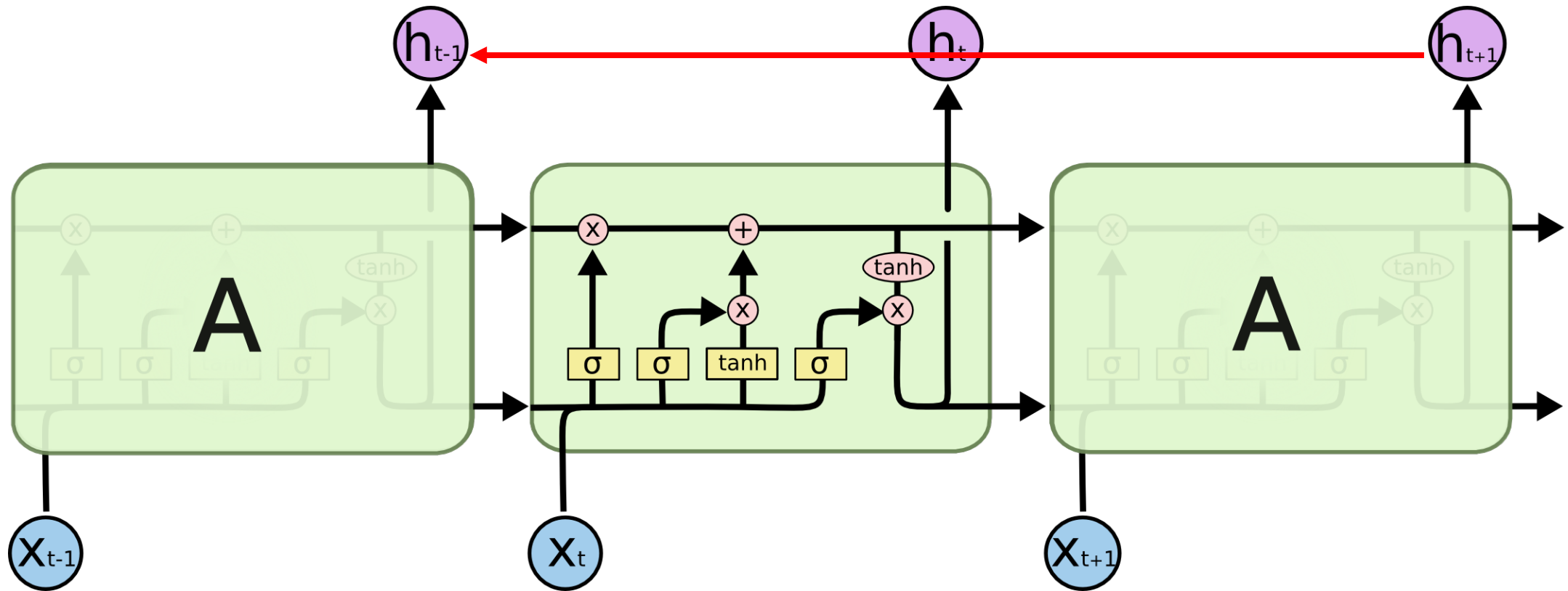


Output gate

$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

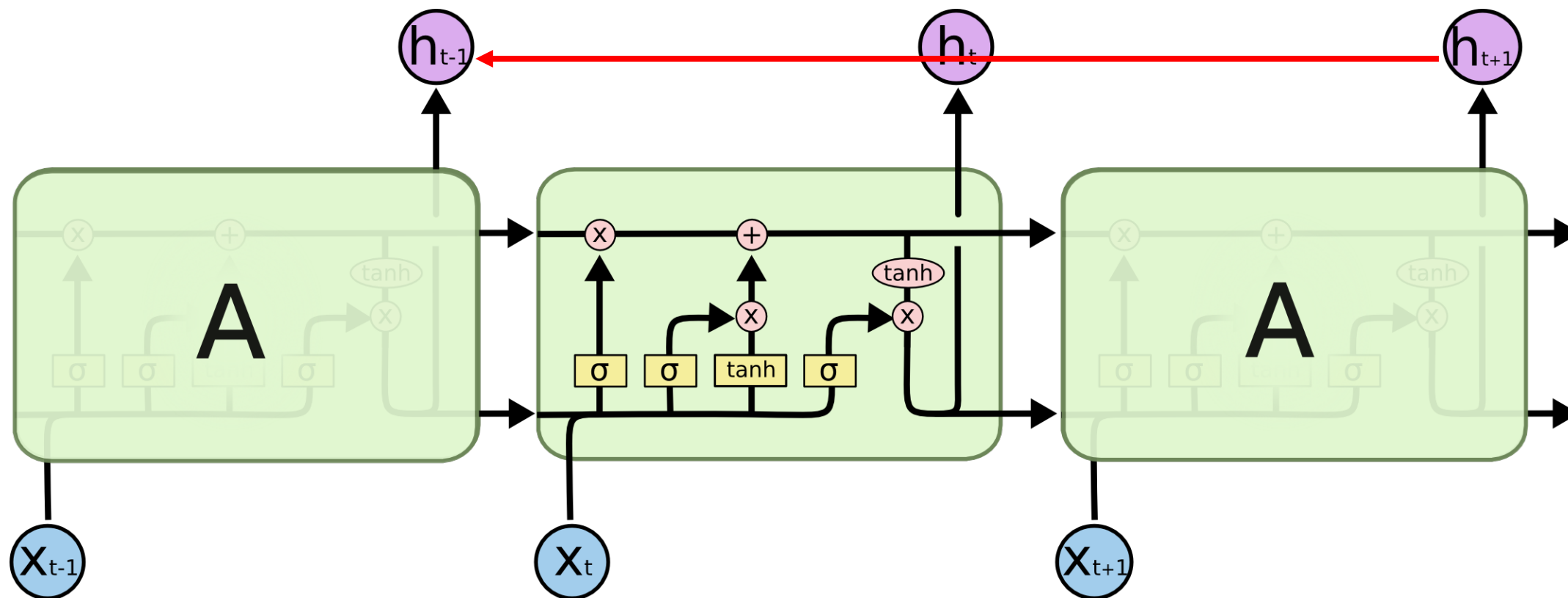
$$h_t = o_t * \tanh (C_t)$$

LSTM



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

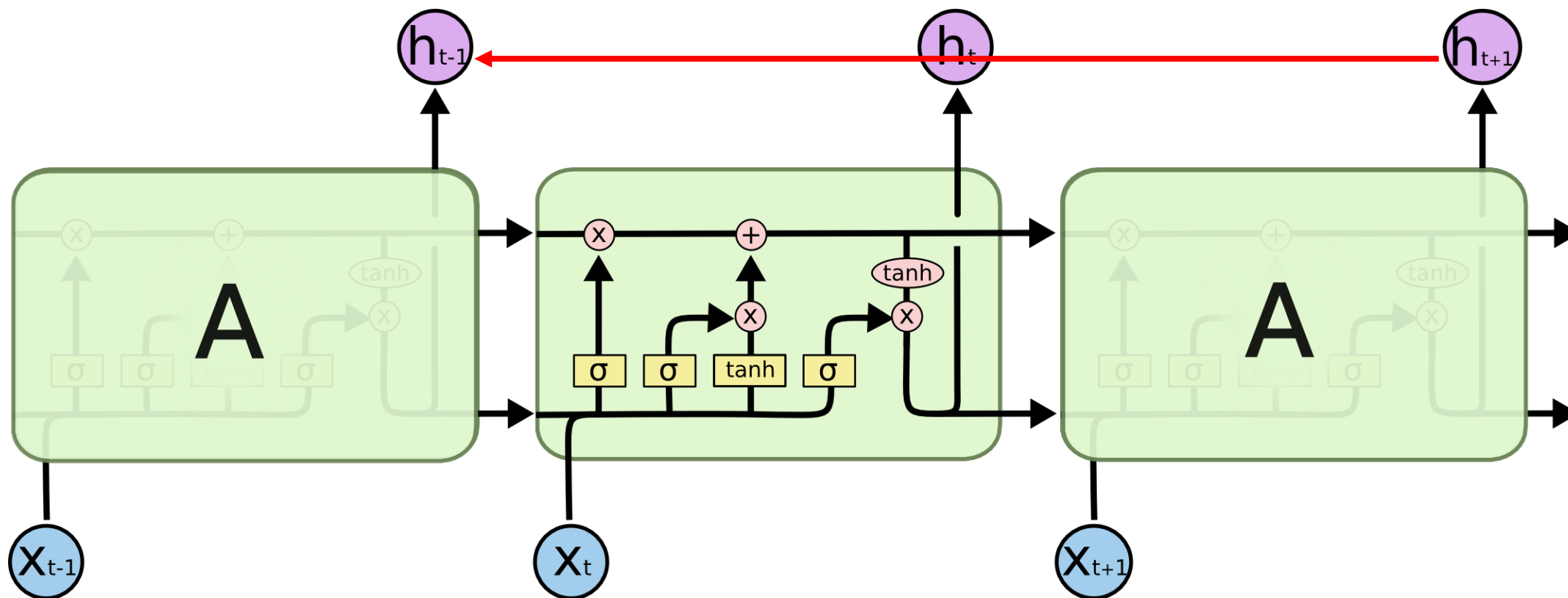
LSTM



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t = f_t * f_{t-1} * C_{t-2} + others$$

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTM



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t = f_t * f_{t-1} * C_{t-2} + \text{others}$$

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad \text{Not multiplication of weights } \prod W_t$$

Deep RNN

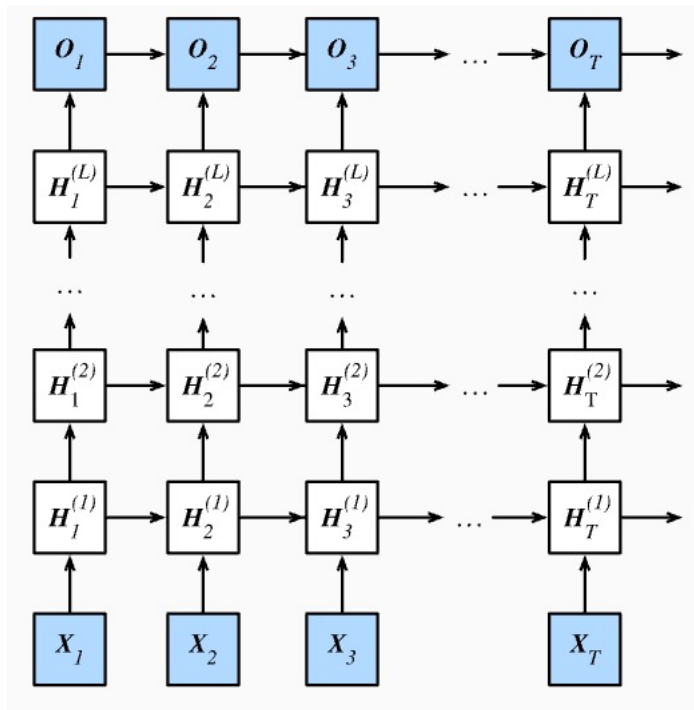


Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Deep RNN

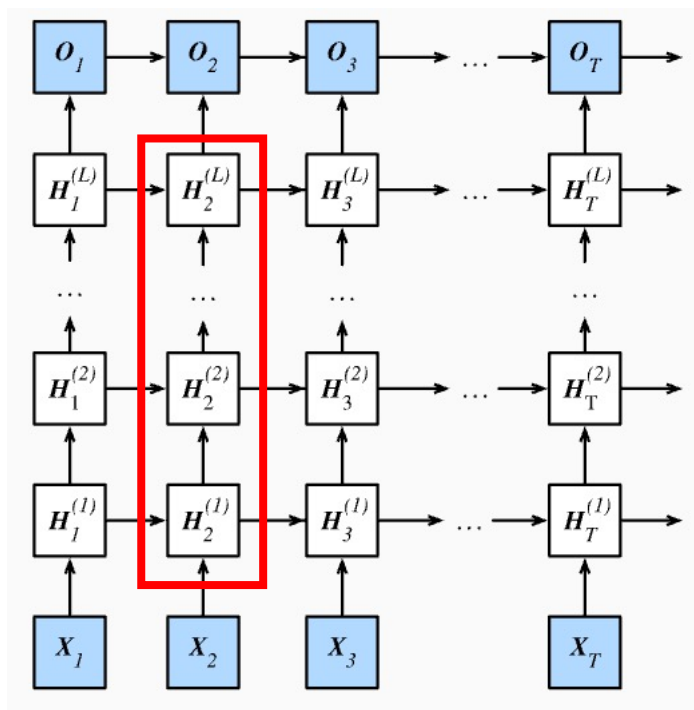
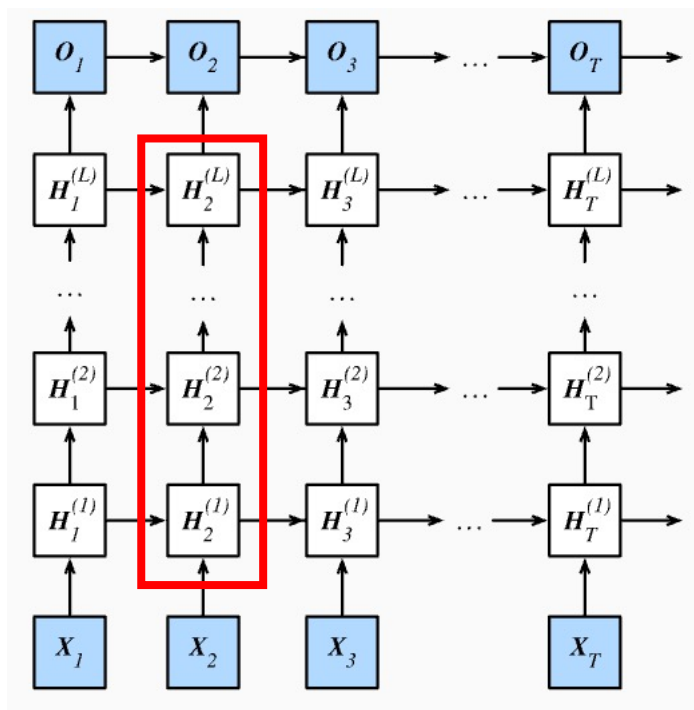


Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Deep RNN



For each unit: a mapping

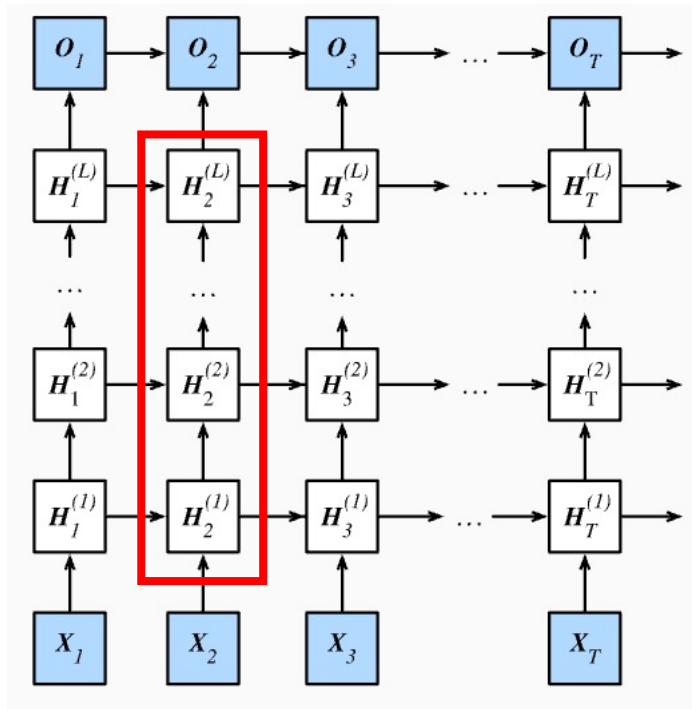
$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l$$

$$h_t^l \in \mathbb{R}^n$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

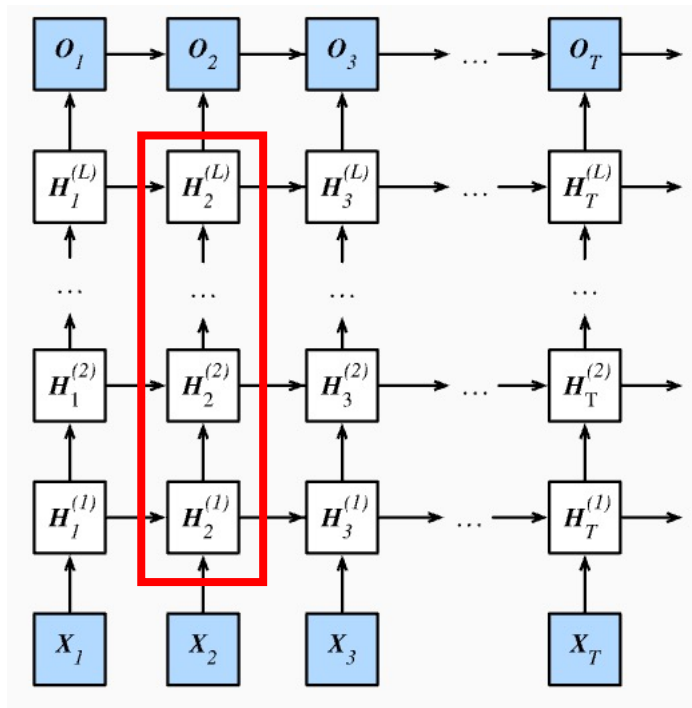
$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

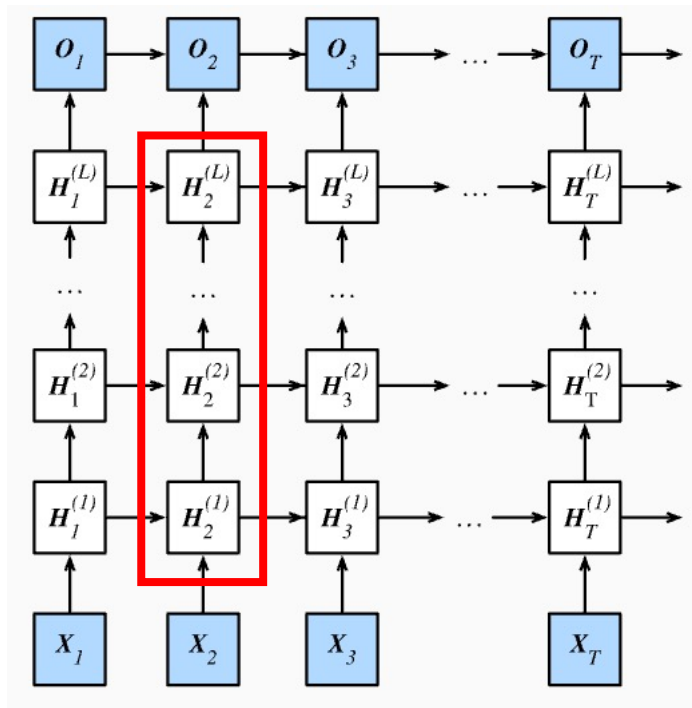
$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

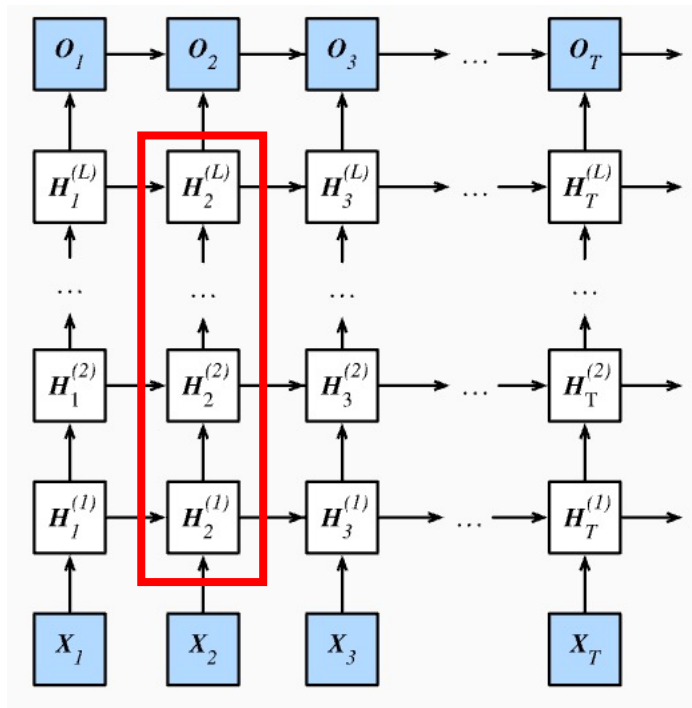
$$h_t^l = f(T_{n,n}h_t^{l-1} + T_{n,n}h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n}h_t^{l-1} + T_{n,n}h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

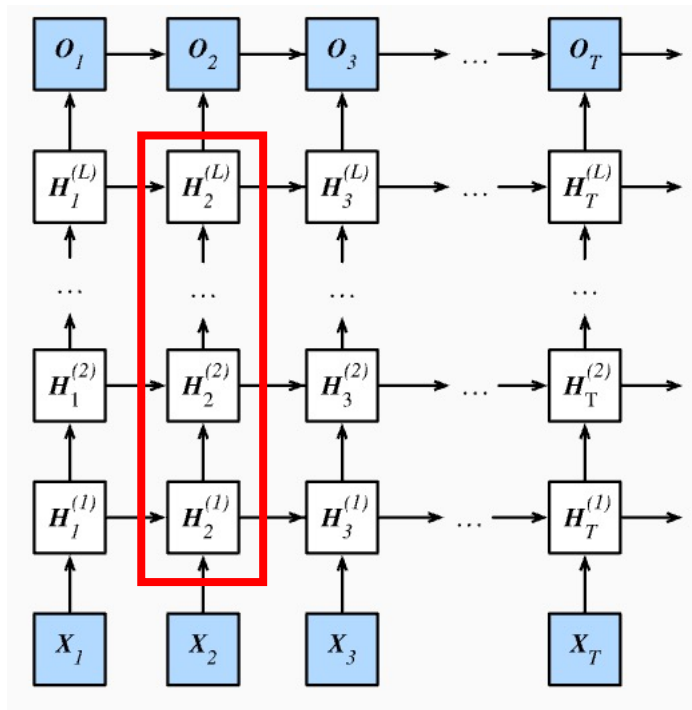
$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

$$Wx + b$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

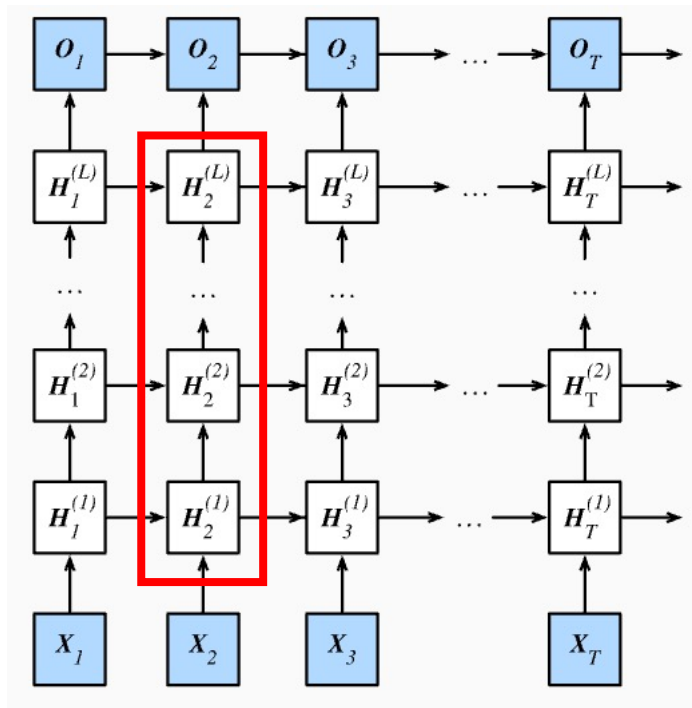
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

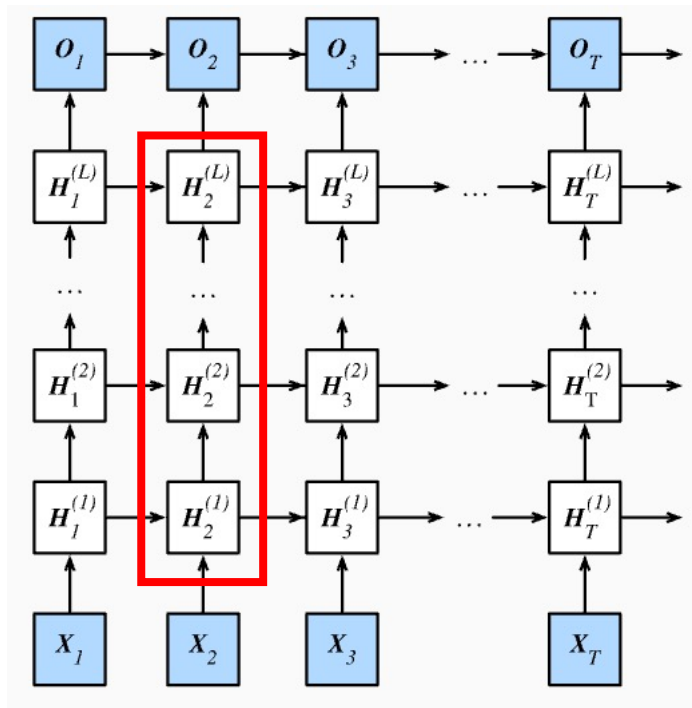
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

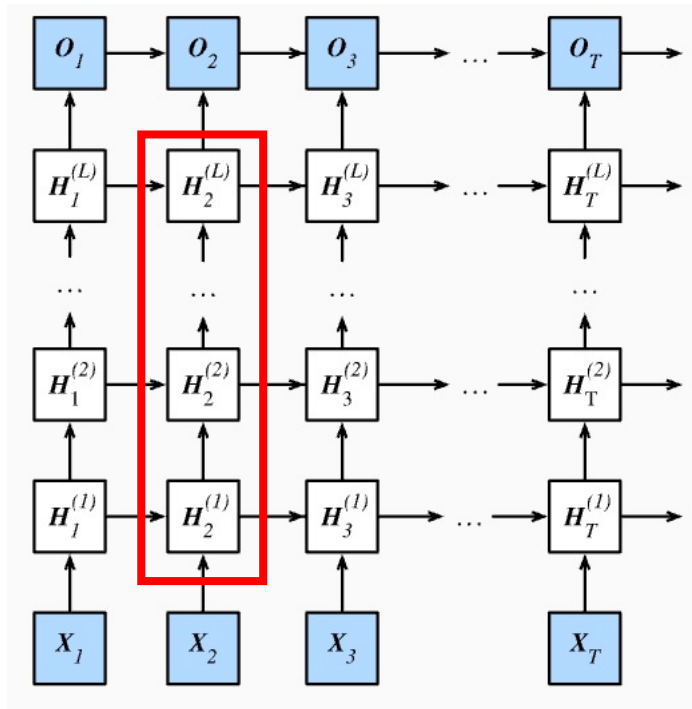
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right]$$

For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \text{tanh}\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

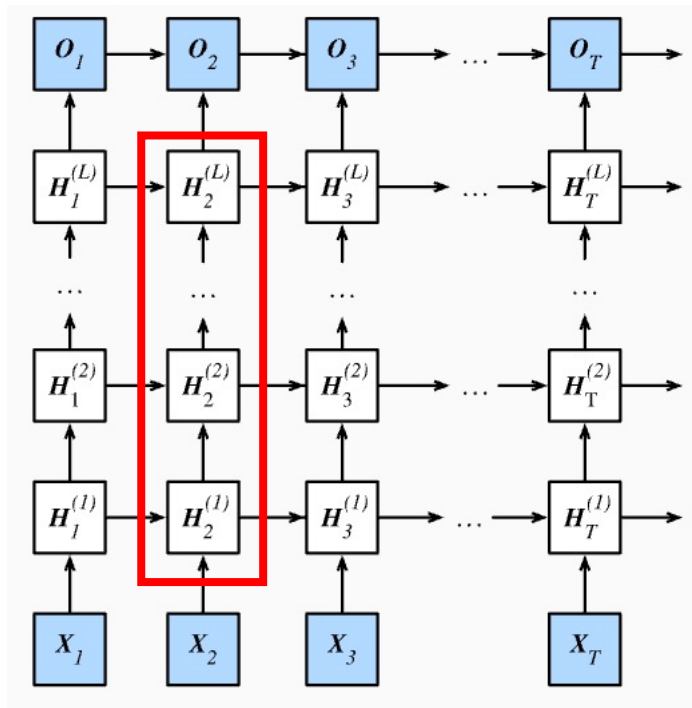
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right] \quad f = \tanh(\cdot)$$

For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \tanh\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

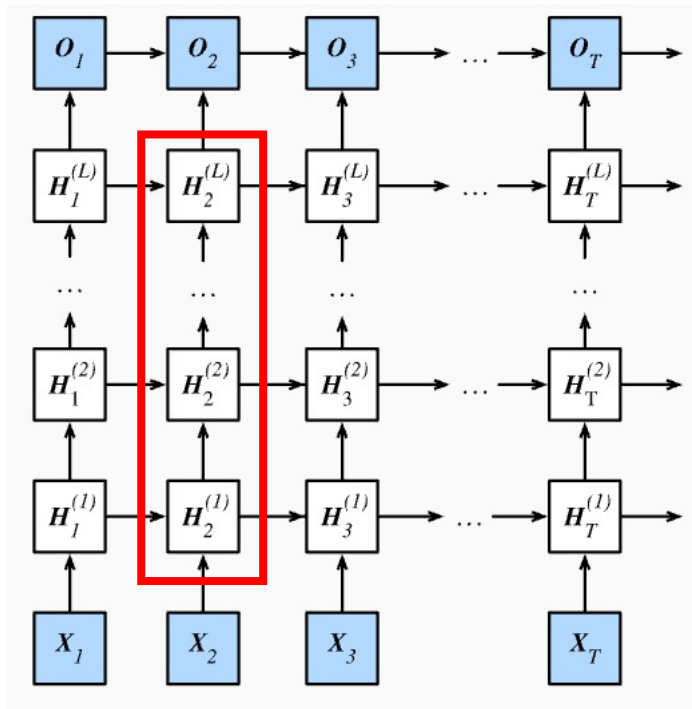
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



Vector feature

$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right] \quad f = \tanh(\cdot)$$

For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \tanh\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

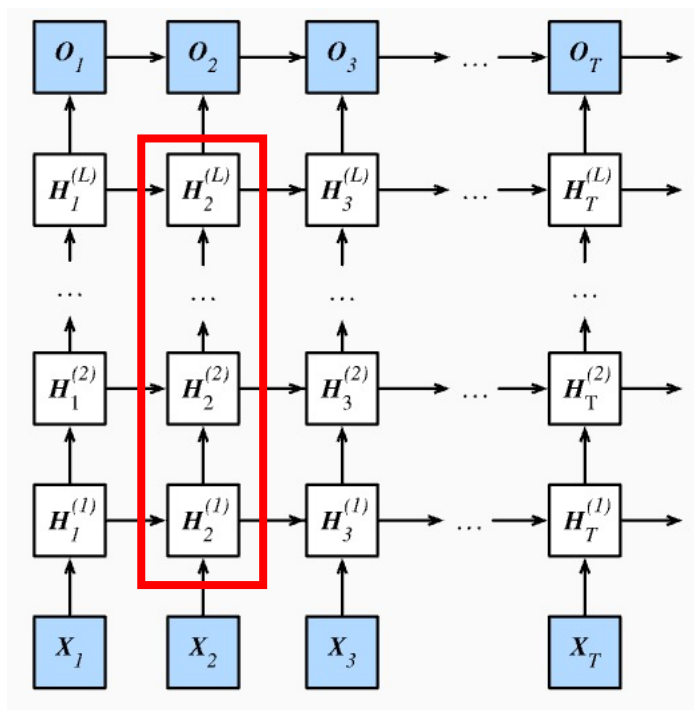
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep RNN



elementwise operation

Vector feature

$$h_t = f \left[W * \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right]$$

$f = \tanh(\cdot)$

For each unit: a mapping

$$\text{RNN} : h_t^{l-1}, h_{t-1}^l \rightarrow h_t^l \quad h_t^l \in \mathbb{R}^n$$

$$h_t^l = f(T_{n,n} h_t^{l-1} + T_{n,n} h_{t-1}^l), \text{ where } f \in \{\text{sigm}, \tanh\}$$

$$T_{n,m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ (a mapping function)}$$

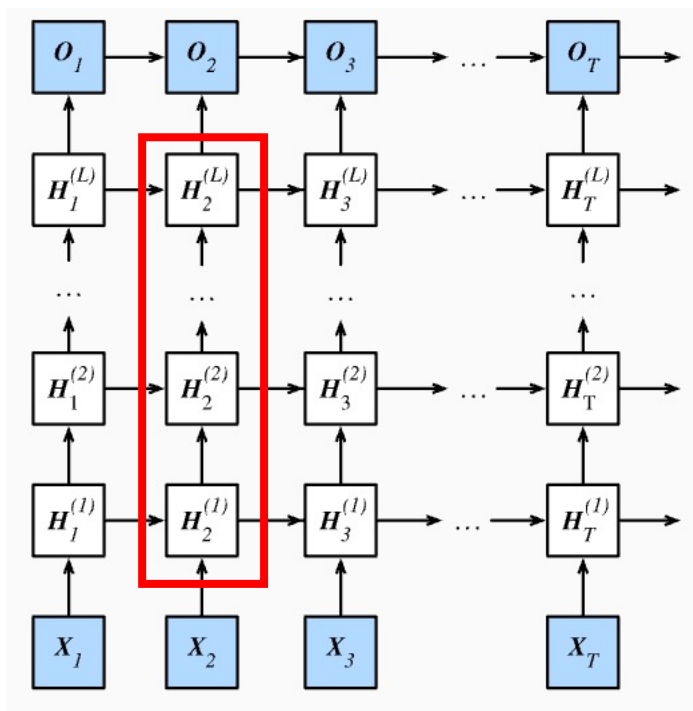
$$Wx + b$$

$$x \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep LSTM

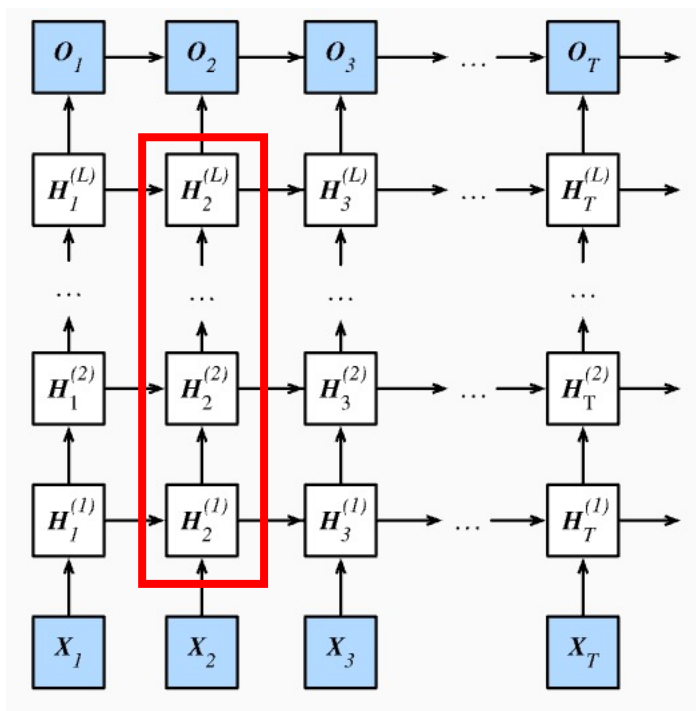


$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow h_t^l, c_t^l$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep LSTM



$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow h_t^l, c_t^l$$

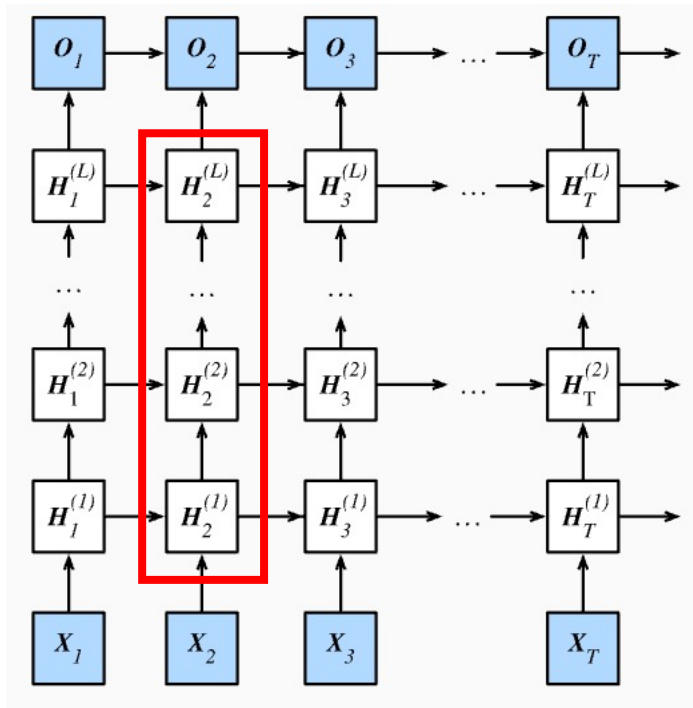
$$c_t^l \in \mathbb{R}^n$$

$$h_t^l \in \mathbb{R}^n$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep LSTM



$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow h_t^l, c_t^l$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} T_{2n,4n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

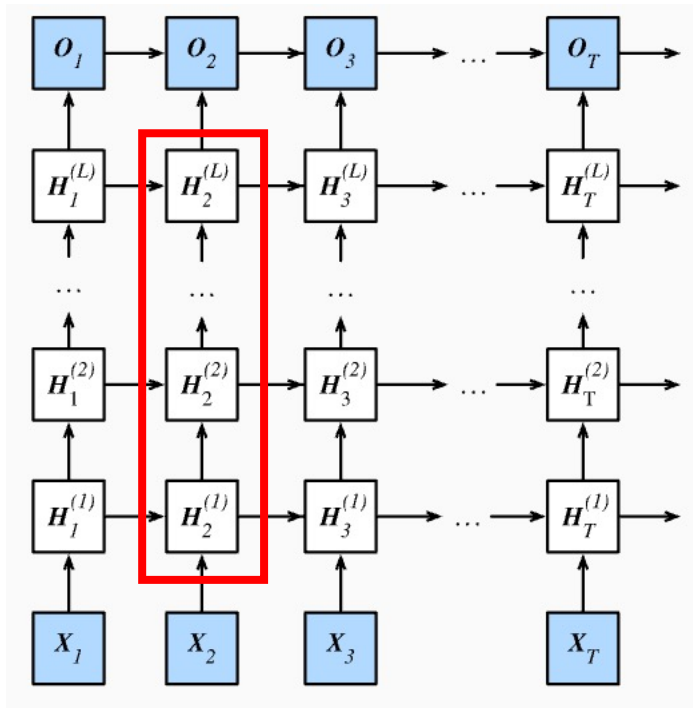
$$c_t^l \in \mathbb{R}^n$$

$$h_t^l \in \mathbb{R}^n$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep LSTM



$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow \boxed{h_t^l, c_t^l}$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} T_{2n,4n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$\boxed{c_t^l} = f \odot c_{t-1}^l + i \odot g$$

$$\boxed{h_t^l} = o \odot \tanh(c_t^l)$$

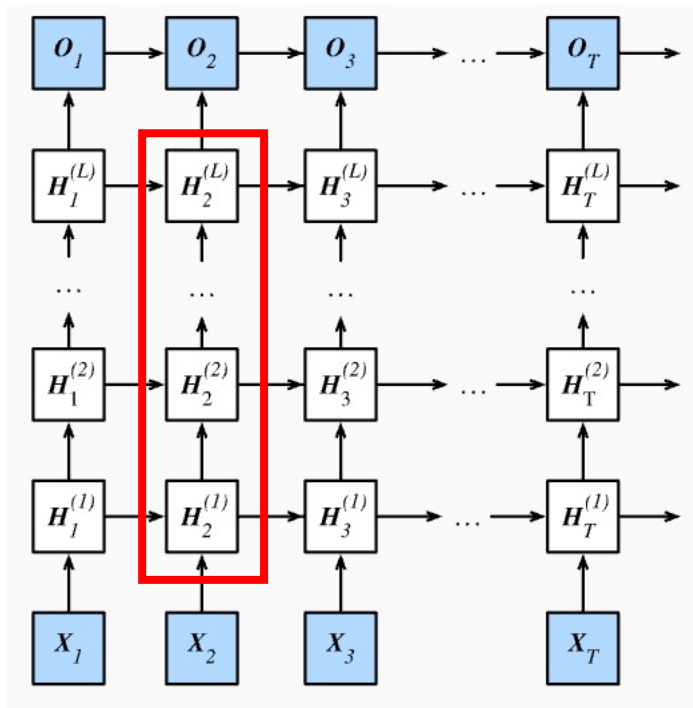
$$c_t^l \in \mathbb{R}^n$$

$$h_t^l \in \mathbb{R}^n$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Deep LTSM



$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow \boxed{h_t^l, c_t^l}$$

$$\begin{pmatrix} \boxed{i} \\ \boxed{f} \\ \boxed{o} \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} T_{2n,4n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$\boxed{c_t^l} = \boxed{f} \odot c_{t-1}^l + \boxed{i} \odot g$$

$$\boxed{h_t^l} = \boxed{o} \odot \tanh(\boxed{c_t^l})$$

weights

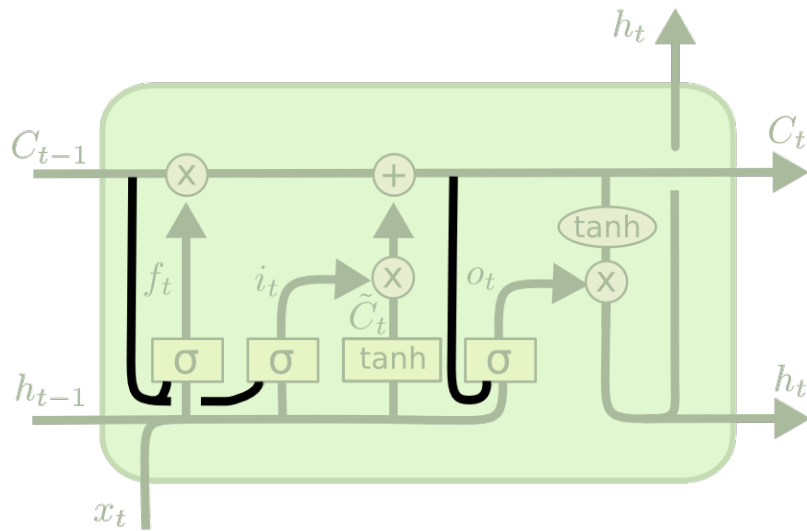
$$c_t^l \in \mathbb{R}^n$$

$$h_t^l \in \mathbb{R}^n$$

Image from Fig 8.10.1 of *Dive into Deep Learning* at https://classic.d2l.ai/chapter_recurrent-neural-networks/deep-rnn.html

Formulations from: Zaremba, Wojciech, Ilya Sutskever, and Oriol Vinyals. "Recurrent neural network regularization." *arXiv preprint arXiv:1409.2329* (2014).

Variants: peephole connections

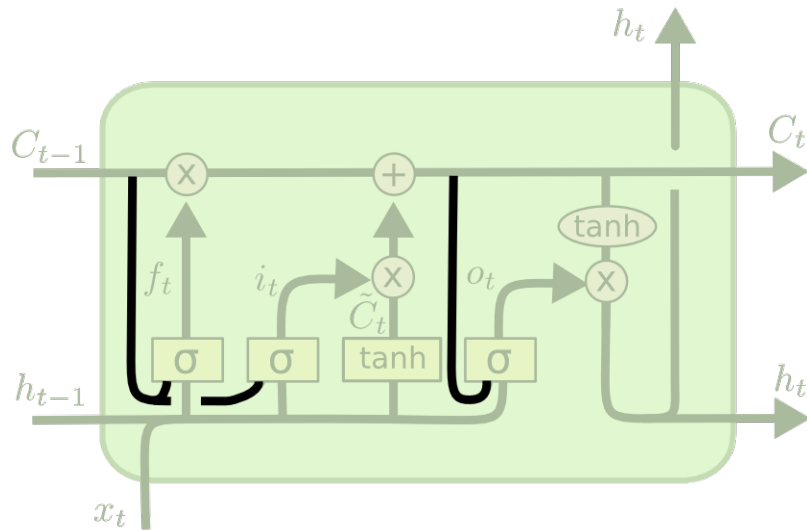


$$f_t = \sigma (W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma (W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma (W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

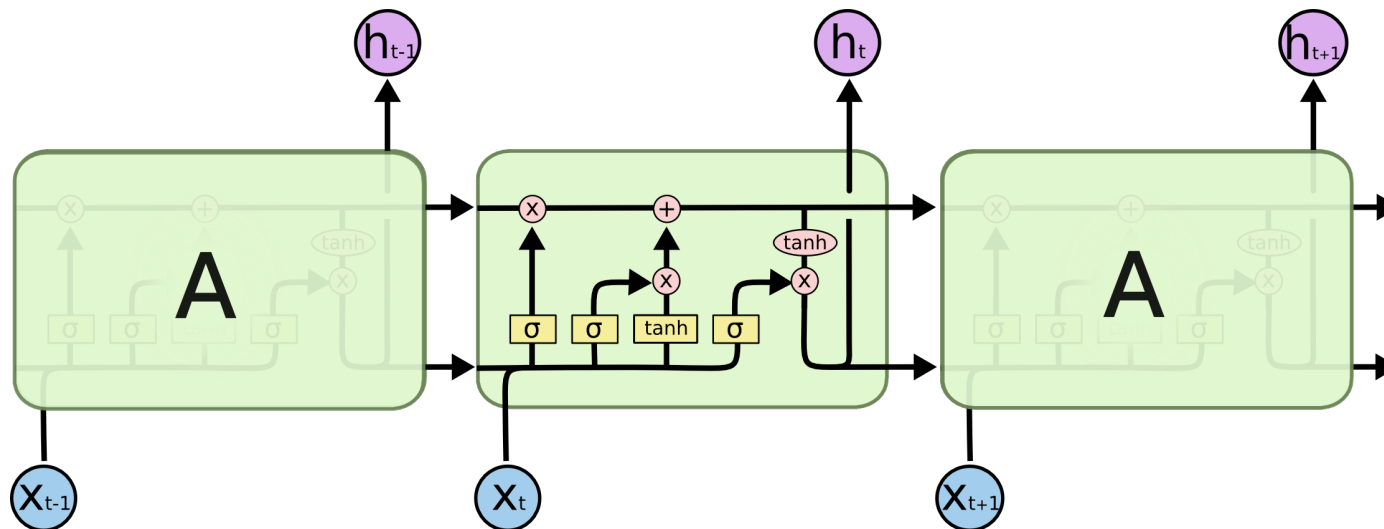
Variants: peephole connections



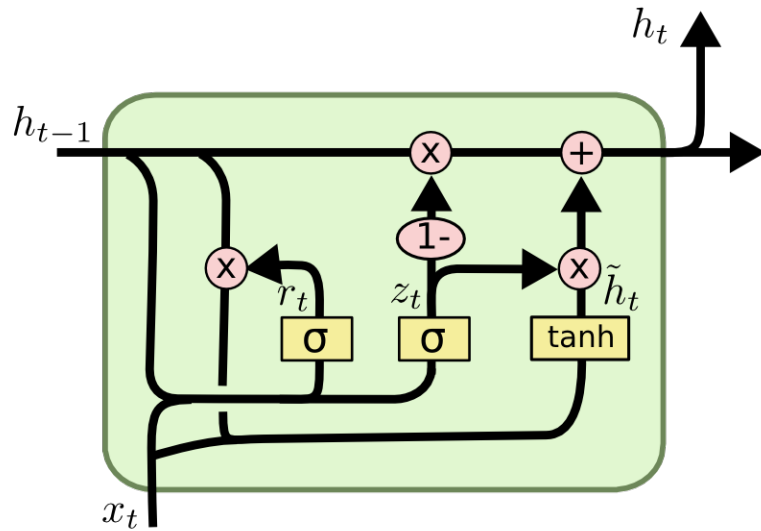
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$



Variant: gated recurrent unit (GRU)



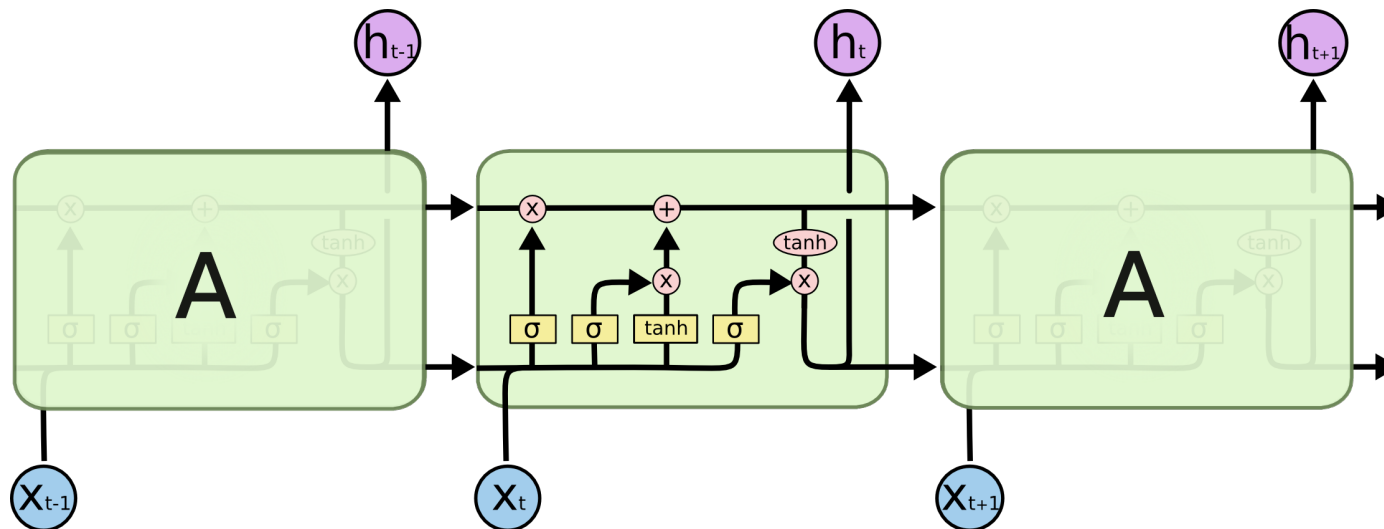
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

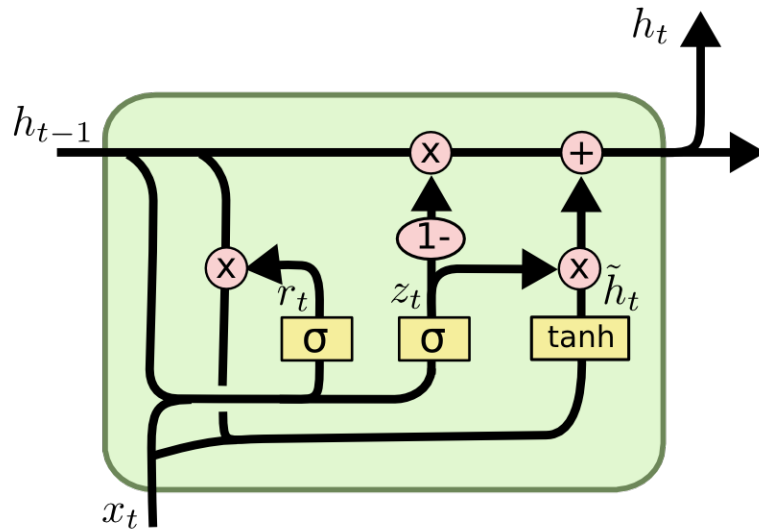
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Weighted average



Variant: gated recurrent unit (GRU)



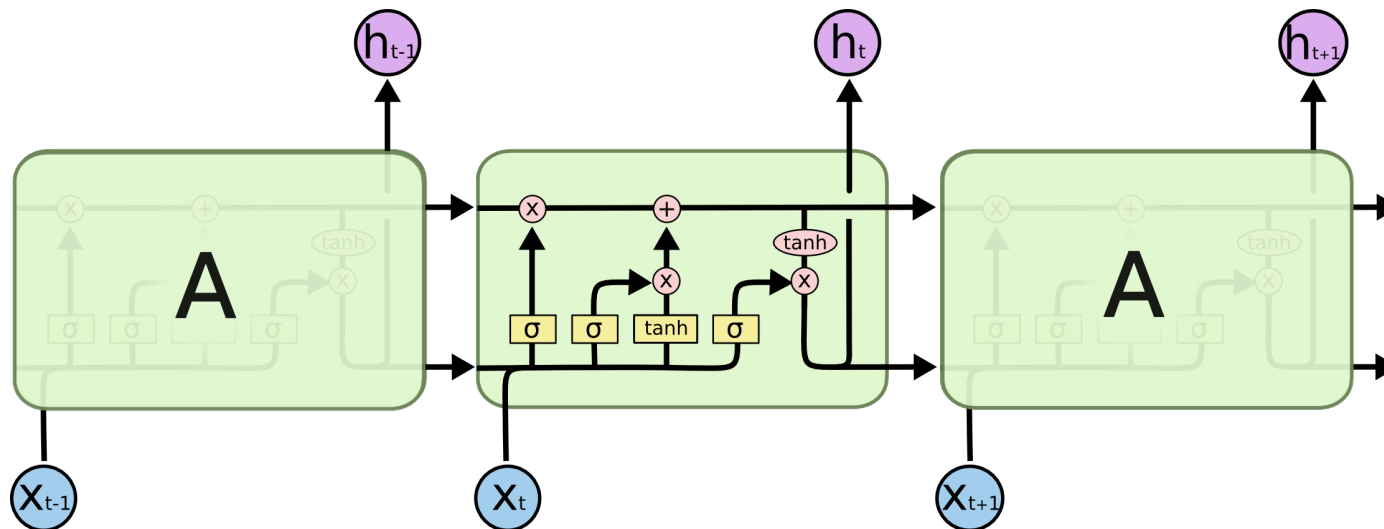
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Weighted average



Reading

- Deep learning book
 - Section 10.2.2, 10.7, 10.10
- Blog: understanding LSTM
 - <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>