

# Natural Language Processing with PyTorch

**Week 3** Neural Language Processing

# Week 2

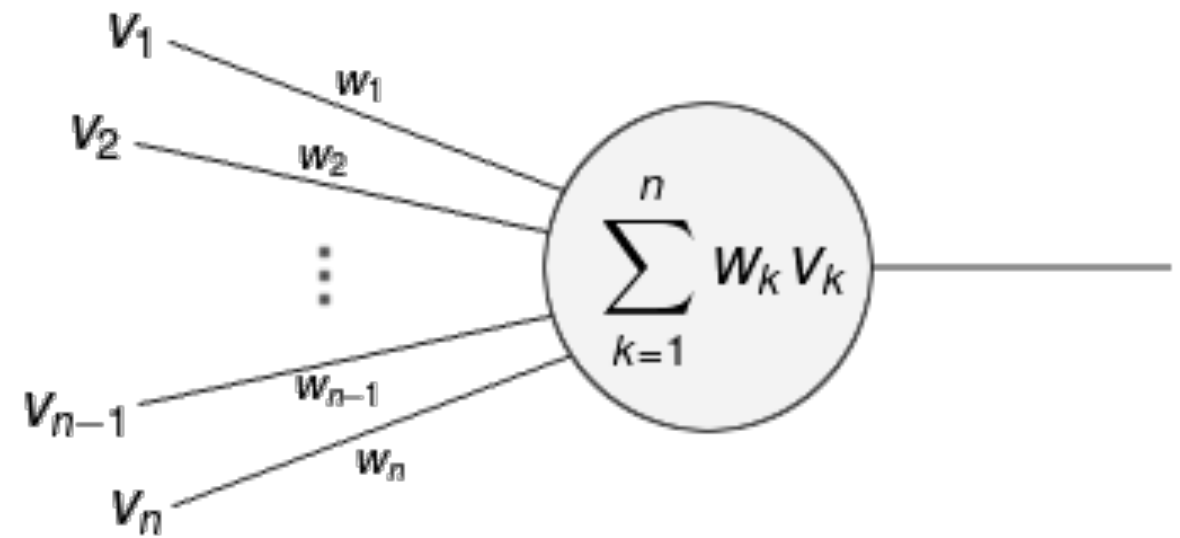
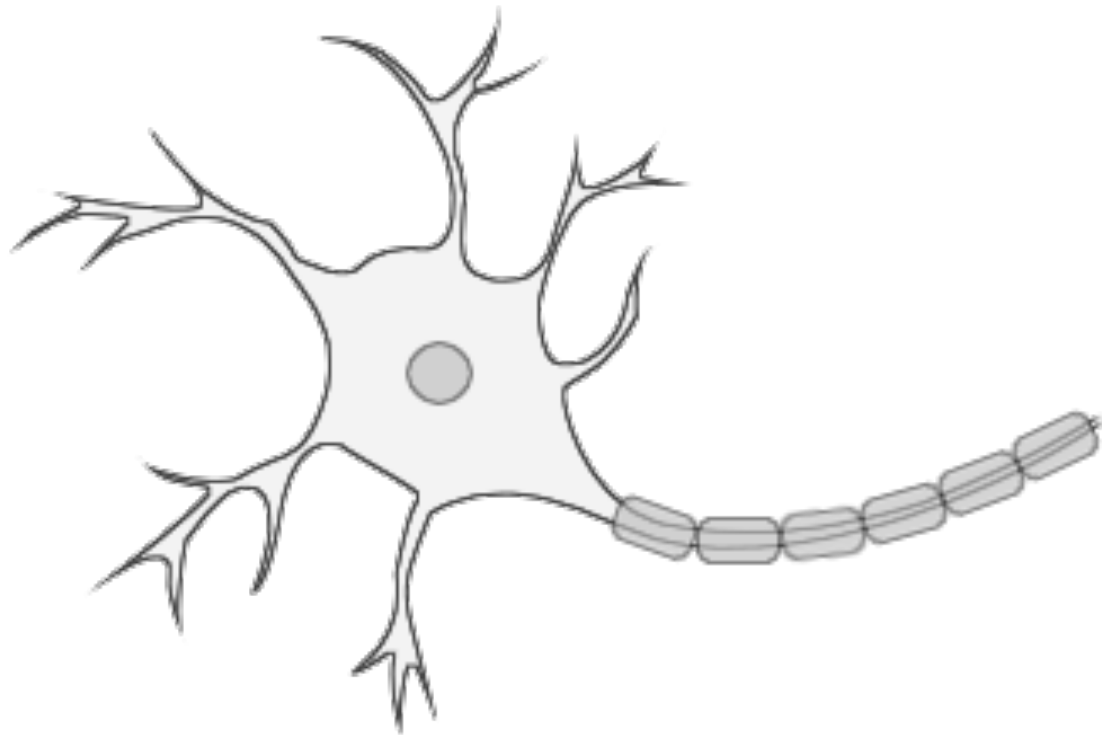
# Review &

# Warranty

The story of CNNs

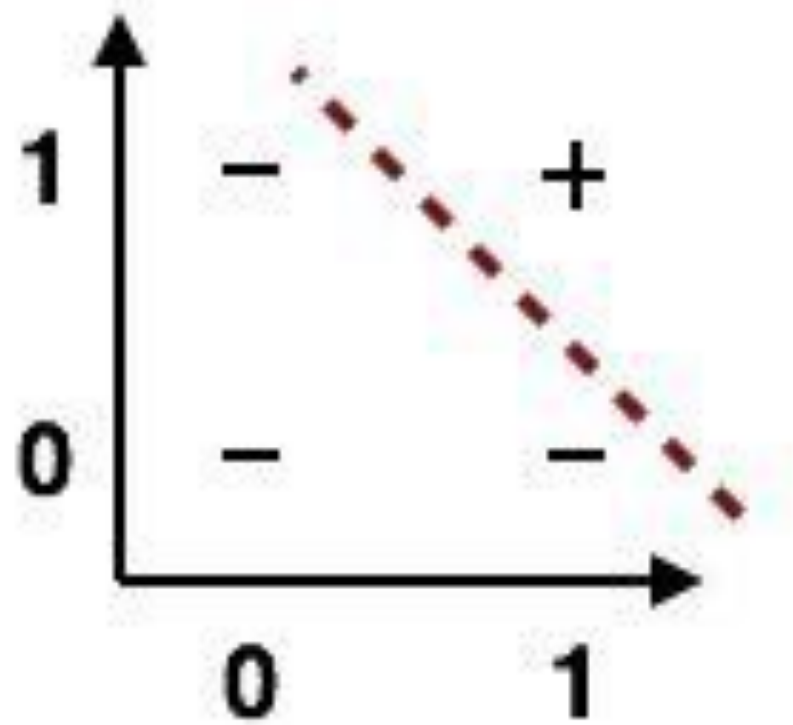
The story of RNNs

The story of residual  
connections

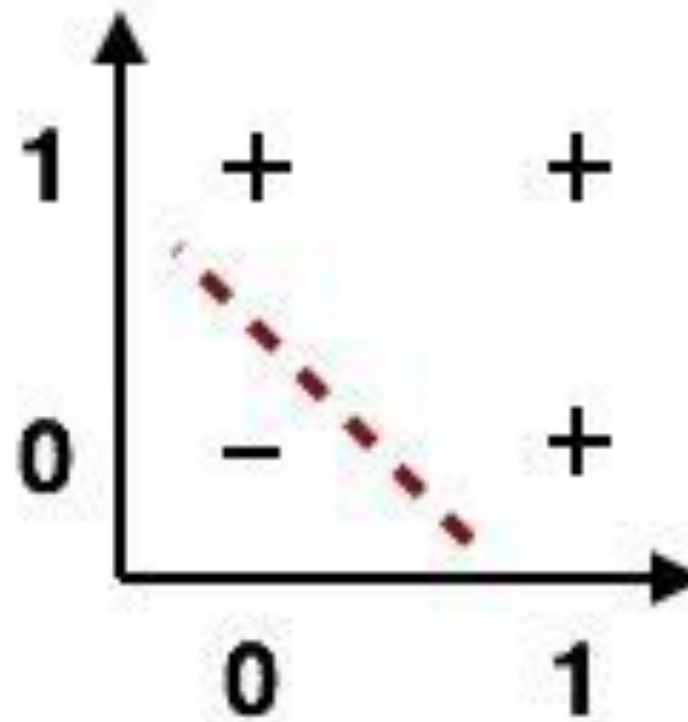


A biological neuron and its mathematical modelling

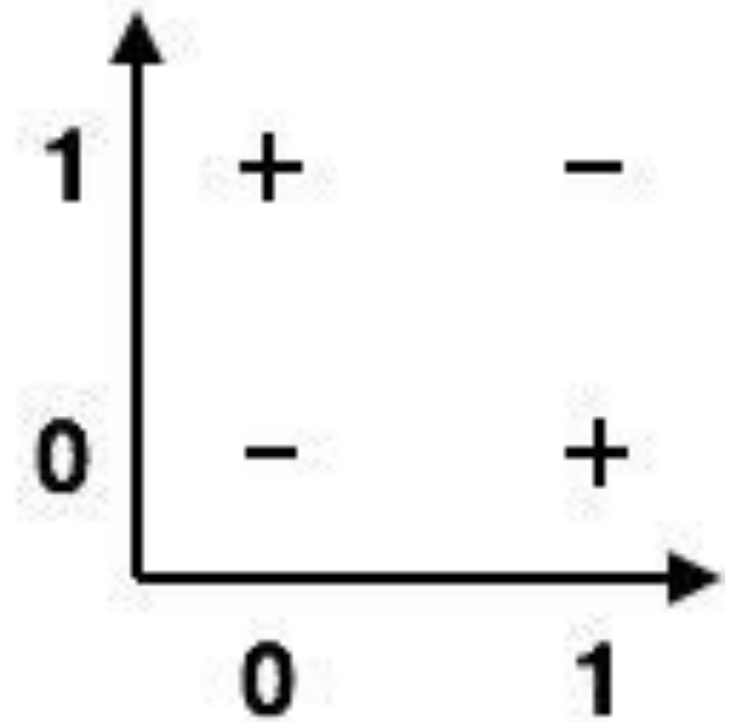
**and**



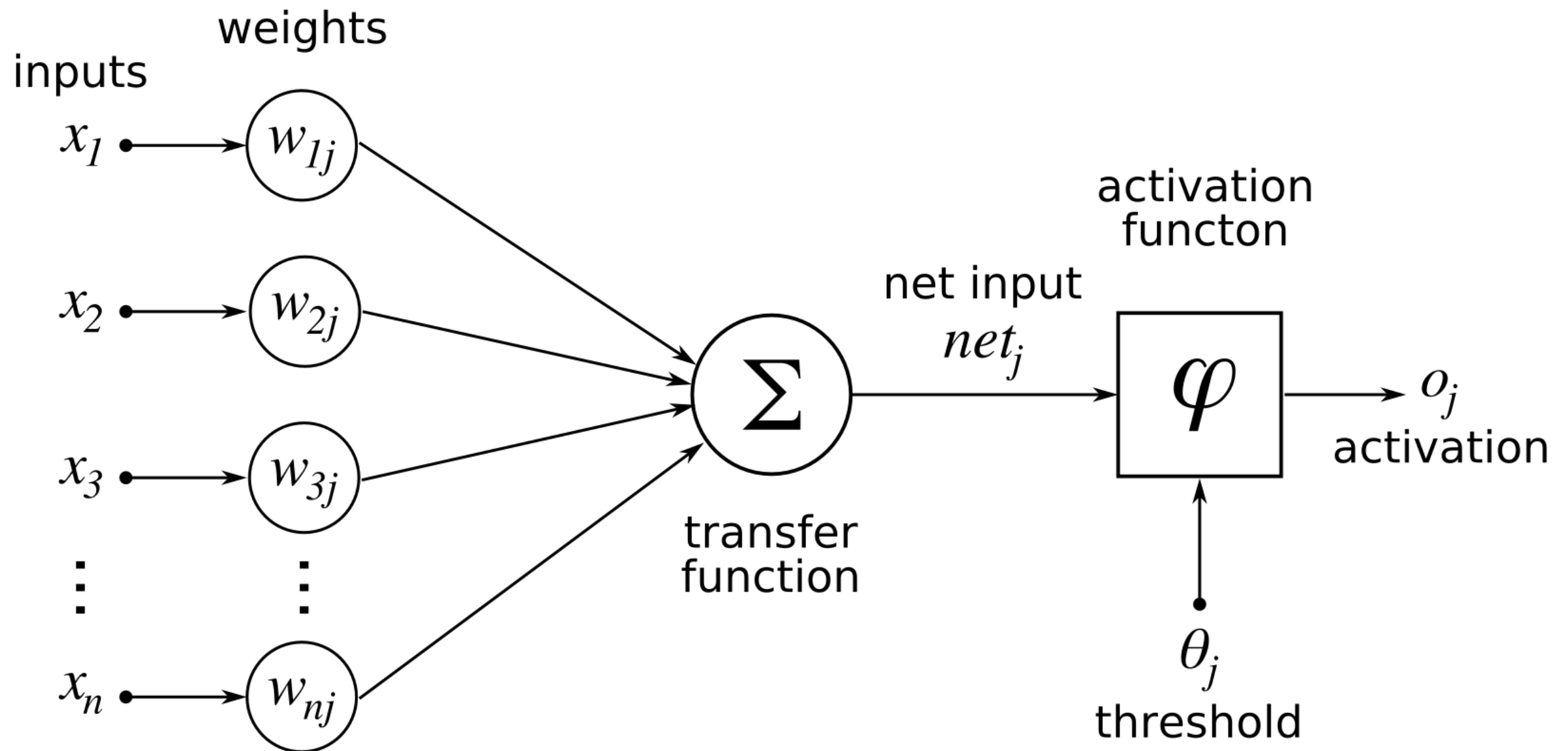
**or**



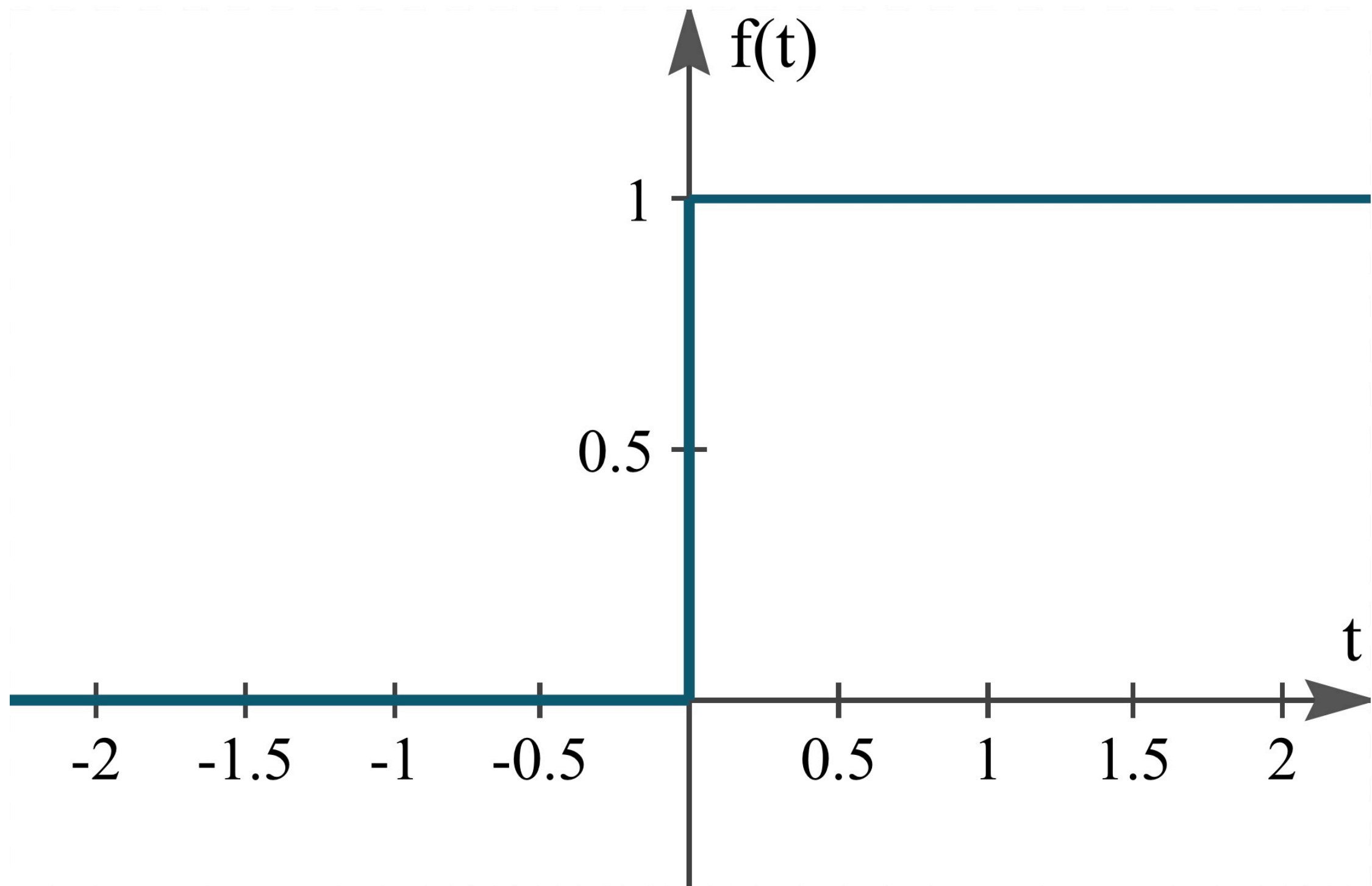
**xor**



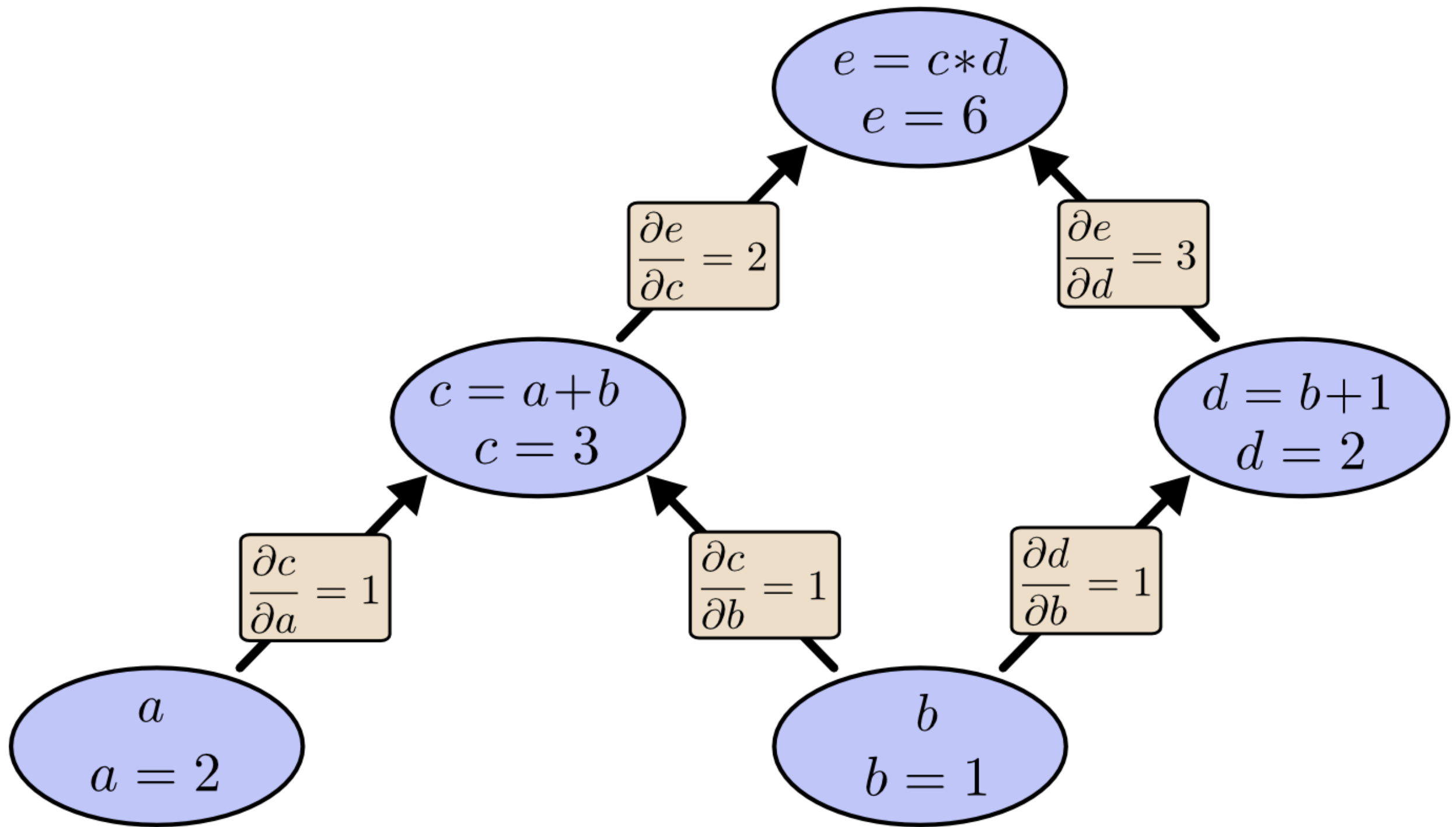
The XOR problem



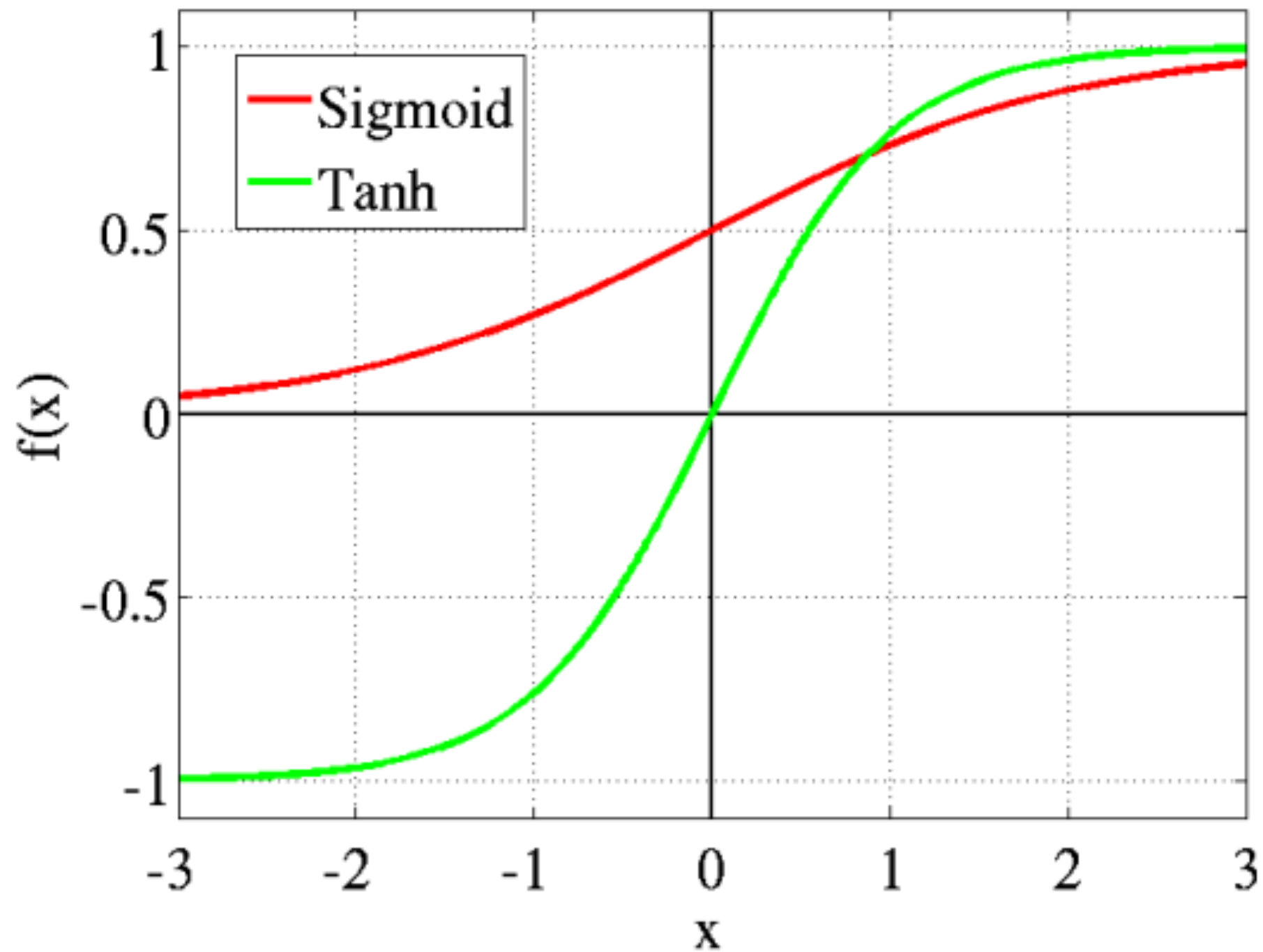
A perceptron with an activation function



The step function

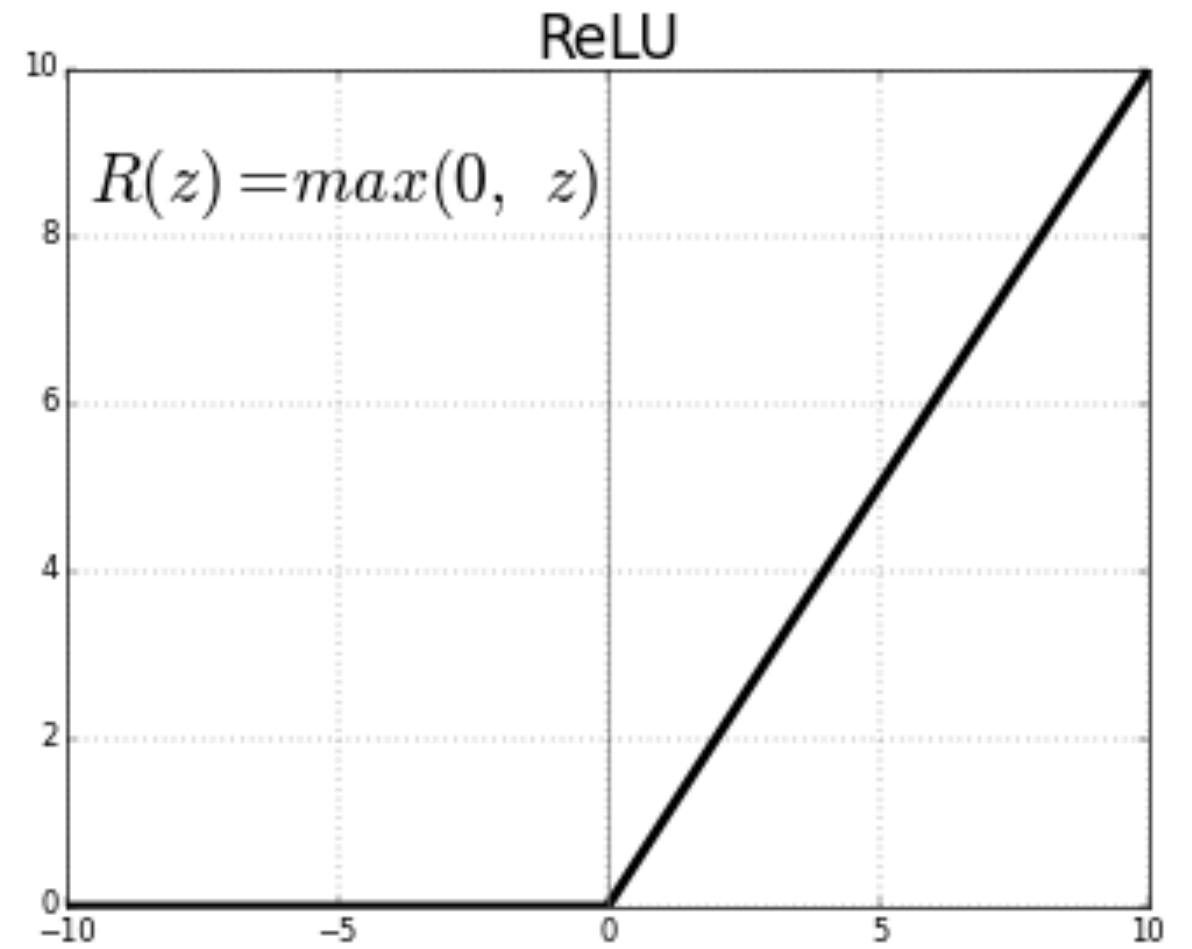
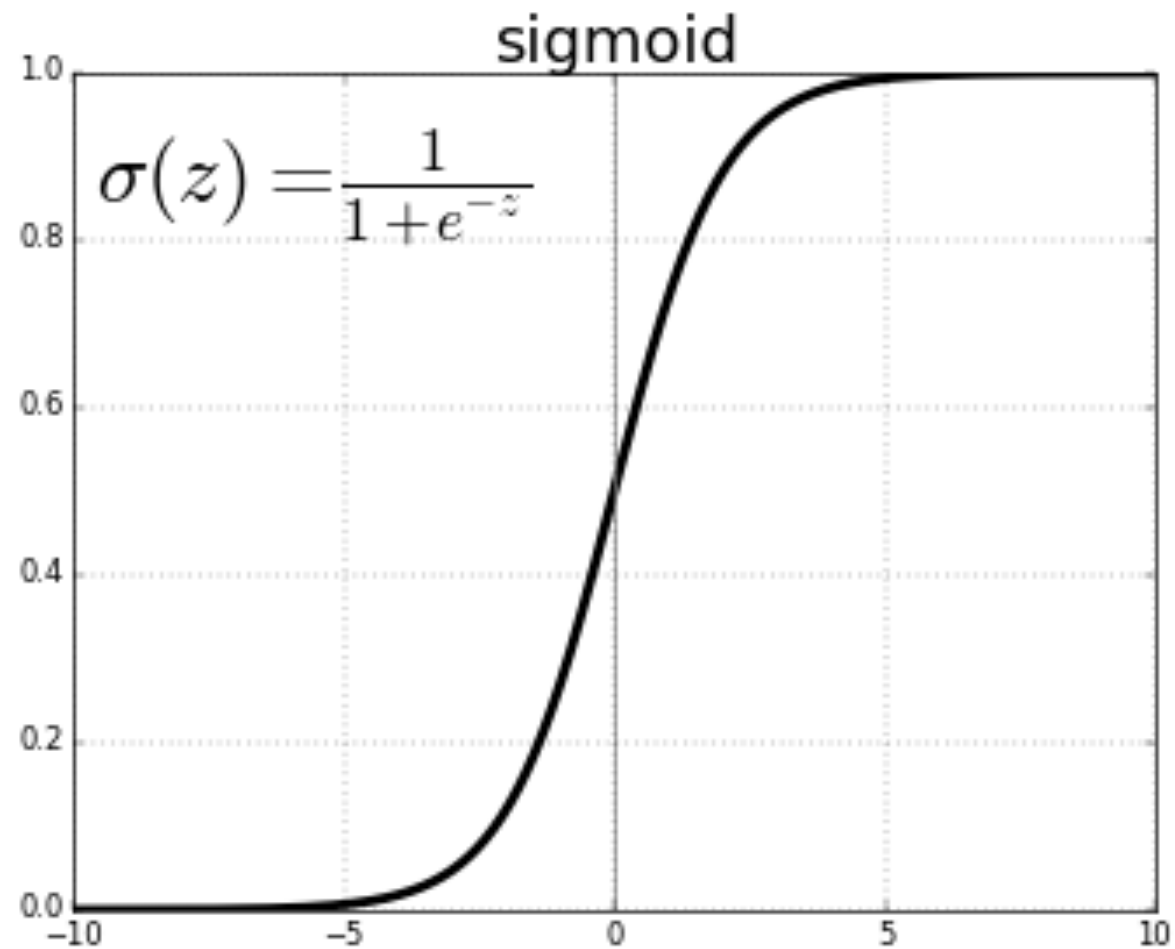


Backpropagation

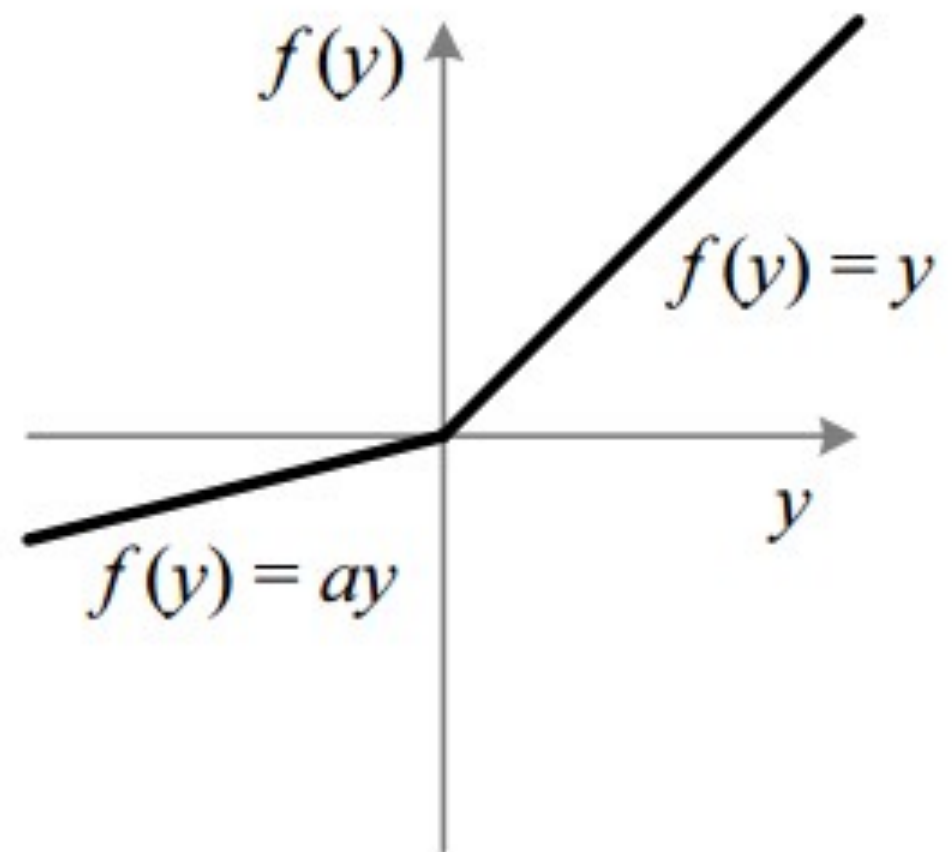
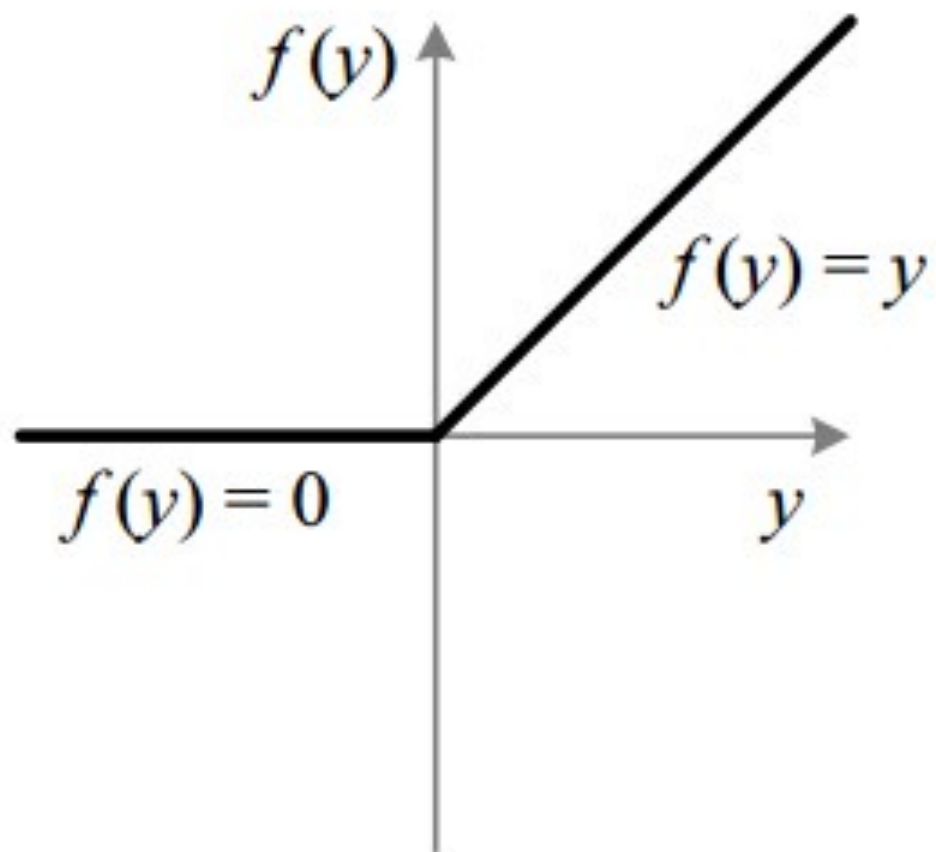


Sigmoid and Tanh





## Sigmoid and ReLU



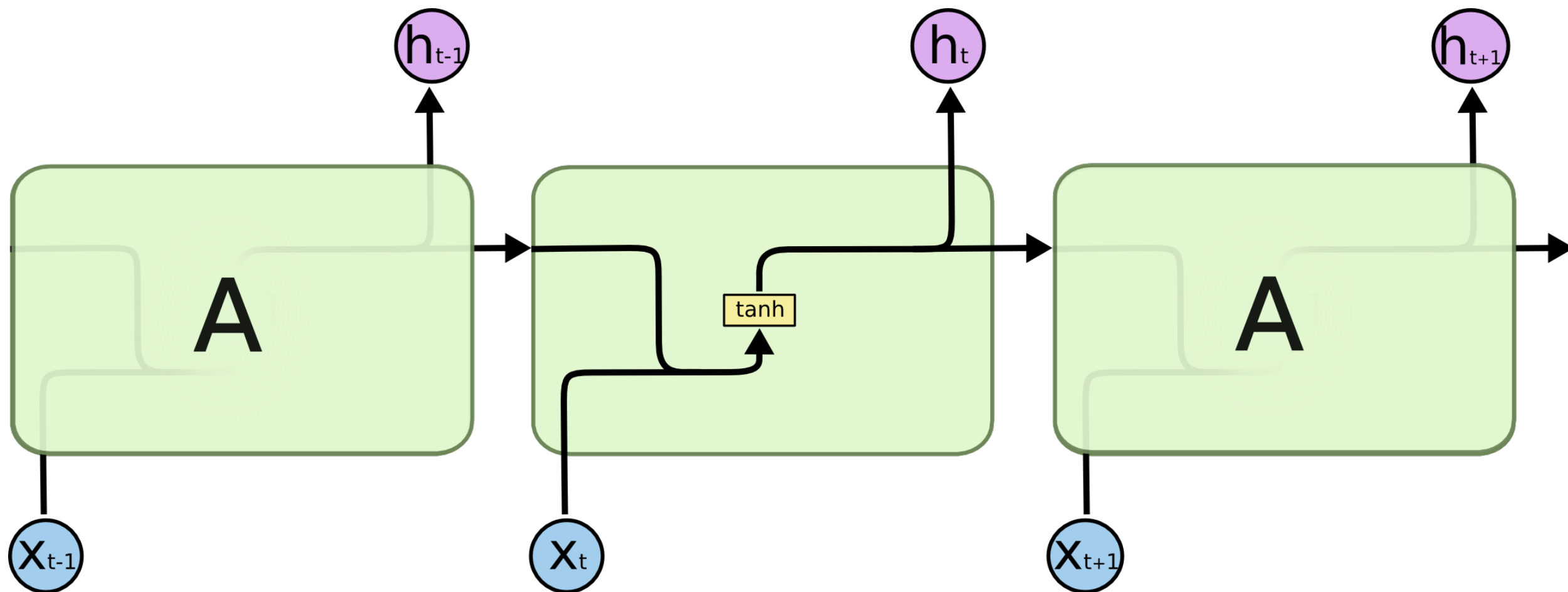
## ReLU and Leaky ReLU

# The story of CNNs

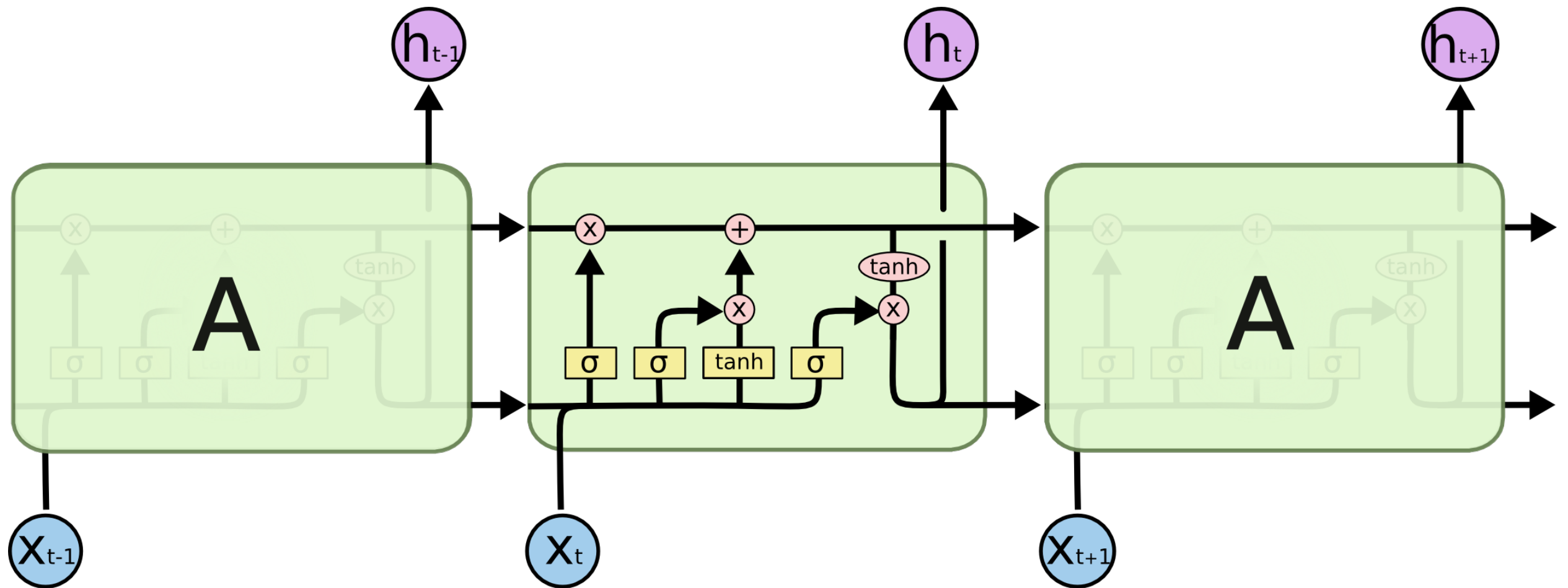
- **Perceptrons!**
- A single perceptron cannot solve **the XOR problem**; and for perceptrons to go multi-layer, some form of **nonlinearity** is needed
- The **step function** as a nonlinear activation function!
- For **backpropagation**, a **differentiable** activation function is needed
- **Sigmoid** and **Tanh** functions!

# The story of CNNs

- Sigmoid and tanh activations suffer from **gradient saturation** and **gradient vanishing**
- **ReLU!**
- ReLU suffers from **dying ReLUs**.
- **Leaky ReLU!**



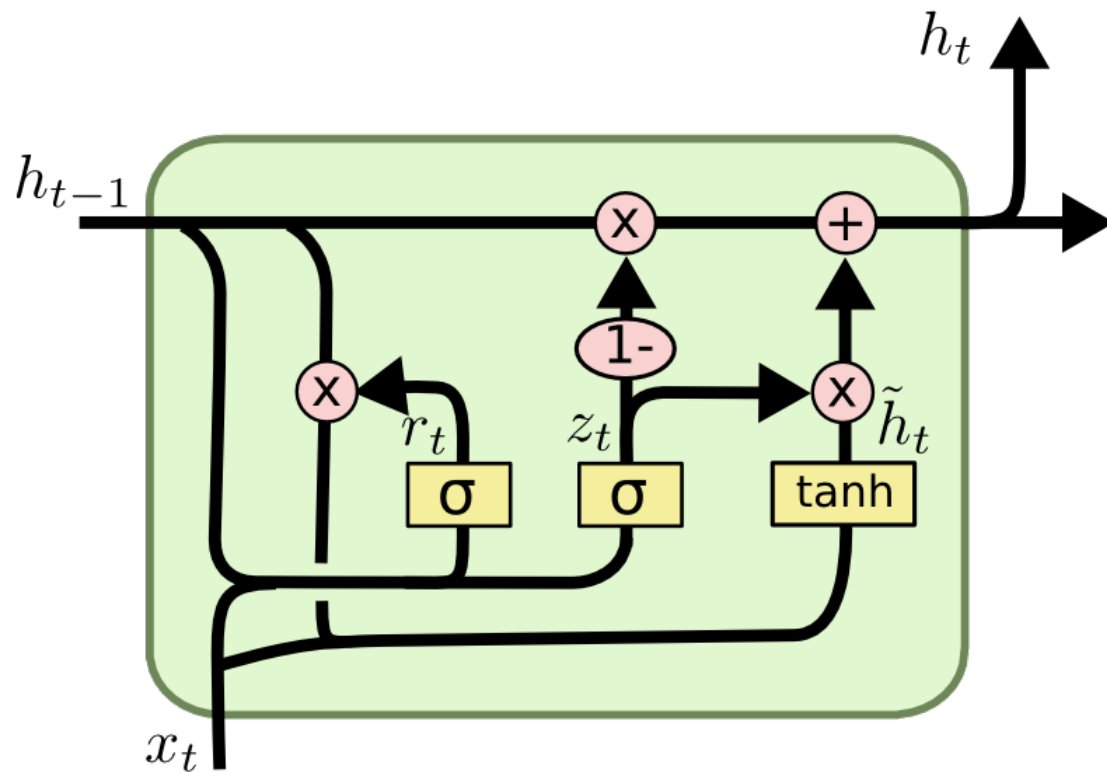
RNN



LSTM

$$\mu_t = \alpha \cdot \mu_{t-1} + (1 - \alpha) \cdot v_t$$

Leaky units, or Exponential Moving Average



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

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

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

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

GRU



$$\nabla [\tanh(\mathbf{X}) \otimes \sigma(\mathbf{X})] = \tanh'(\mathbf{X}) \nabla \mathbf{X} \otimes \sigma(\mathbf{X}) + \sigma'(\mathbf{X}) \nabla \mathbf{X} \otimes \tanh(\mathbf{X}).$$

## Gated Tanh Units

$$\nabla[\mathbf{X} \otimes \sigma(\mathbf{X})] = \nabla \mathbf{X} \otimes \sigma(\mathbf{X}) + \mathbf{X} \otimes \sigma'(\mathbf{X}) \nabla \mathbf{X}$$

## Gated Linear Units

# The story of RNNs

- **RNNs** for sequential, variable-length data!
- As a cell shares weight for all timesteps, its easy for **gradient vanishing and exploding** to happen
- **LSTMs** with gradient super highway!
- LSTMs are unnecessarily complex
- **GRUs** for decreased time resolution!

# The story

- But sigmoid and tanh keeps causing **gradient saturation**, and ReLU cannot be used because of gradient exploding
- **Gated Linear Units!**

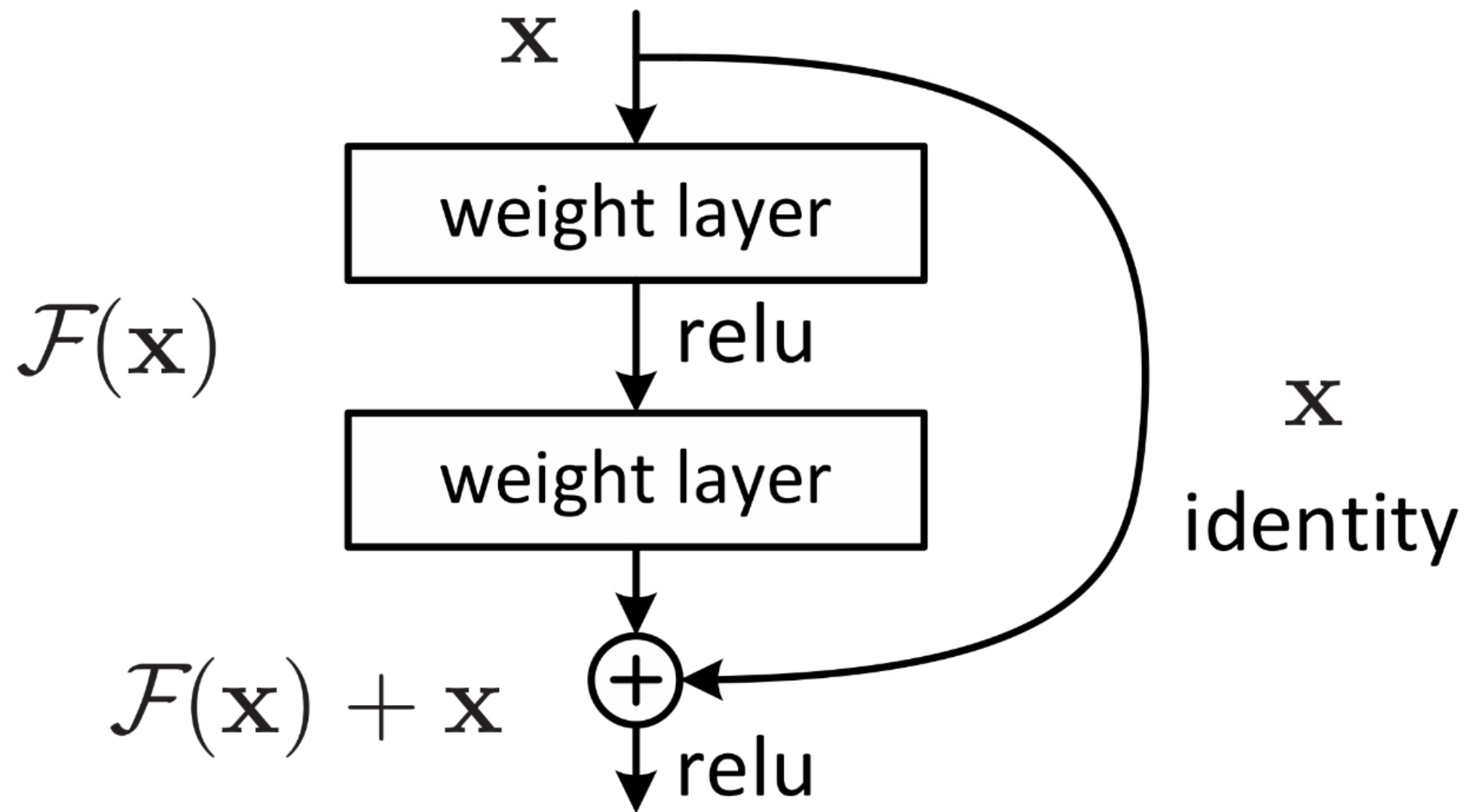
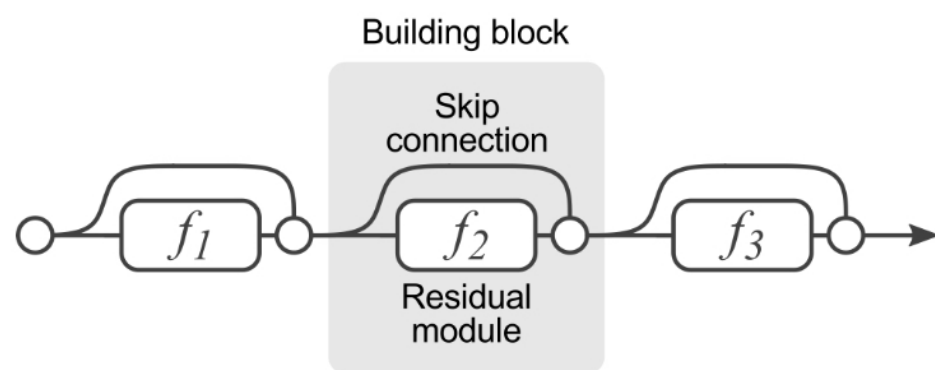


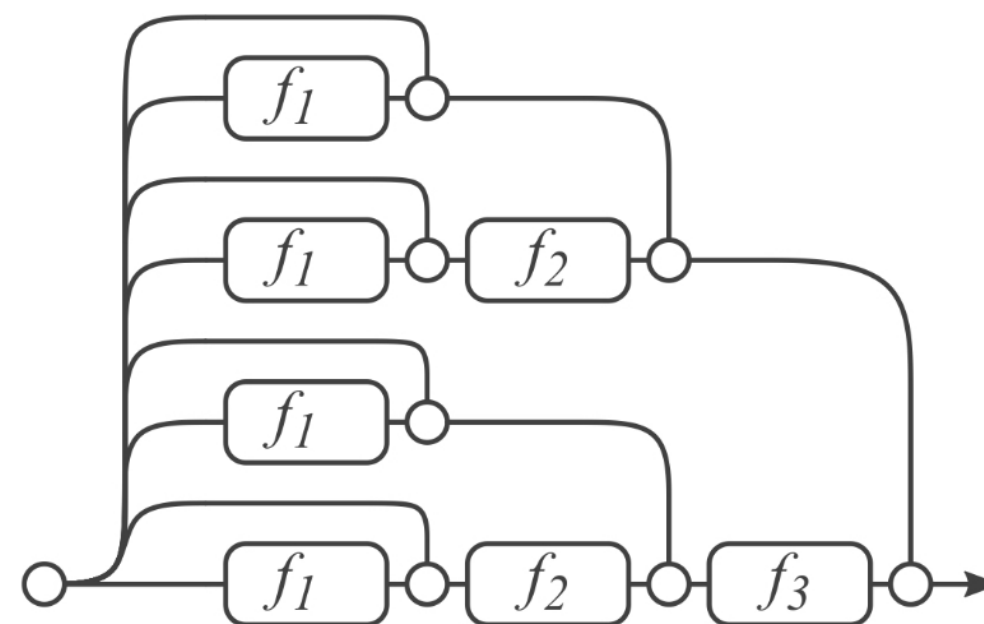
Figure 2. Residual learning: a building block.

Residual connections in CNNs



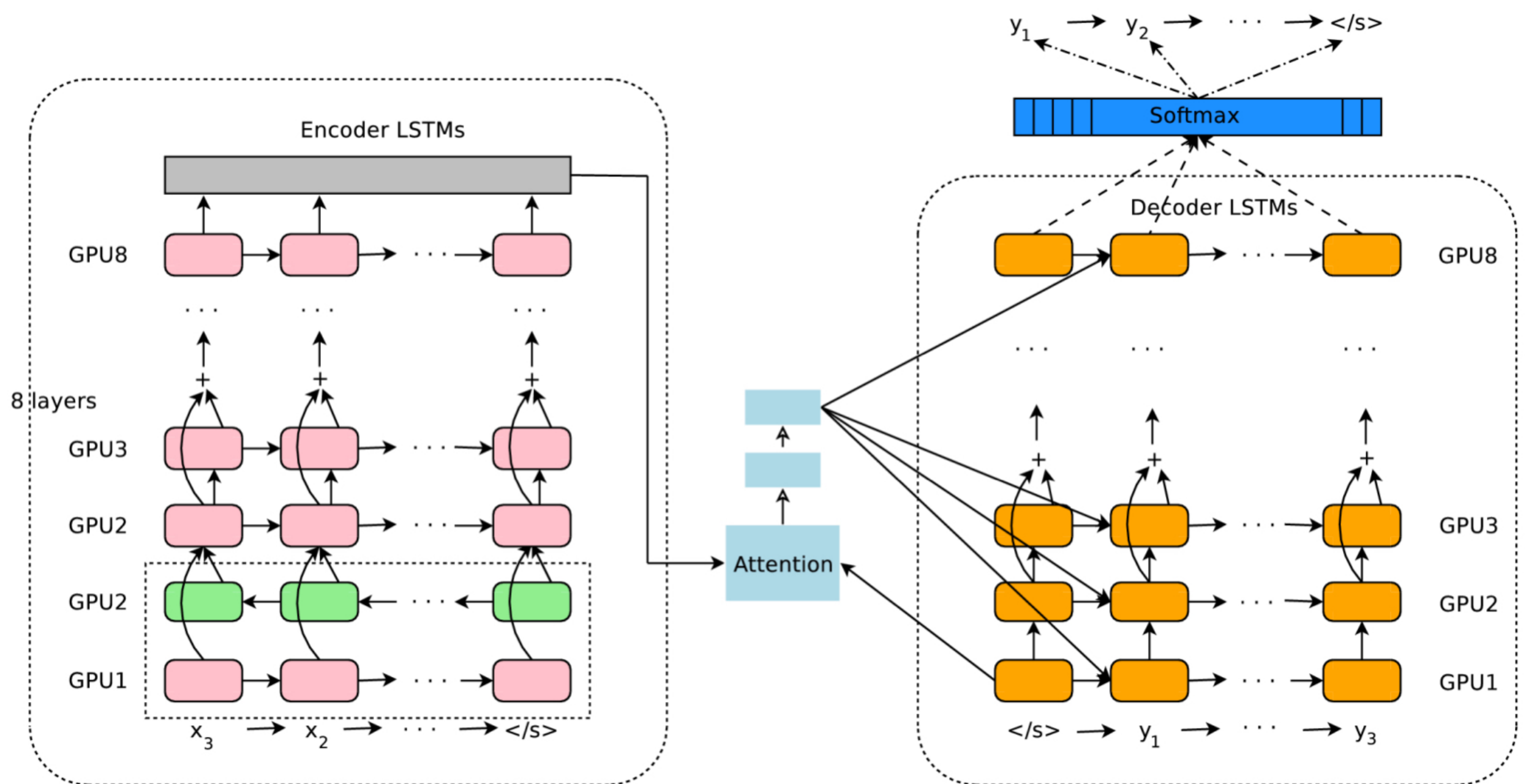
(a) Conventional 3-block residual network

=



(b) Unraveled view of (a)

# Residual networks allowing stochastic layer depths



## Stack-wise residual connections in RNNs

In **zoneout**, the values of the hidden state and memory cell randomly either maintain their previous value or are updated as usual. This introduces stochastic identity connections between subsequent time steps:

$$c_t = d_t^c \odot c_{t-1} + (1 - d_t^c) \odot (f_t \odot c_{t-1} + i_t \odot g_t)$$
$$h_t = d_t^h \odot h_{t-1} + (1 - d_t^h) \odot (o_t \odot \tanh(f_t \odot c_{t-1} + i_t \odot g_t))$$

Time-wise residual connections in RNNs (Zoneout)



# The story of residual connections

- Deep neural networks are difficult to train
- **ResNet!**
- Residual connections in RNNs?
  - Stack-wise residual connections: Google NMT (2016)
  - Time-wise residual connections: Zoneout (2017)