# Natural Language Processing with PyTorch

**Week 2** Deep Neural Networks in PyTorch

**fast campus**

https://github.com/juneoh/fastcampus-pytorch-nlp

# Review & Warranty

# Contents

1. Convolutional Neural Networks

   ○ The convolutional layer

   ○ Batch normalization

   ○ Residual connections

   ○ Variations: dilated convolution, deconvolution, separable convolution

**fast campus**

# Contents

2. Recurrent Neural Networks

  ○ The recurrent layer

  ○ Gradient vanishing and exploding

  ○ Gradient clipping

  ○ Variations: Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU)

**Fast campus**

# Contents

3. Cryptocurrency price prediction using CNN and RNN

- ○ Cryptocurrency 101
- ○ Obtaining and preprocessing the data
- ○ Building our first CNN model
- ○ Building our first RNN model
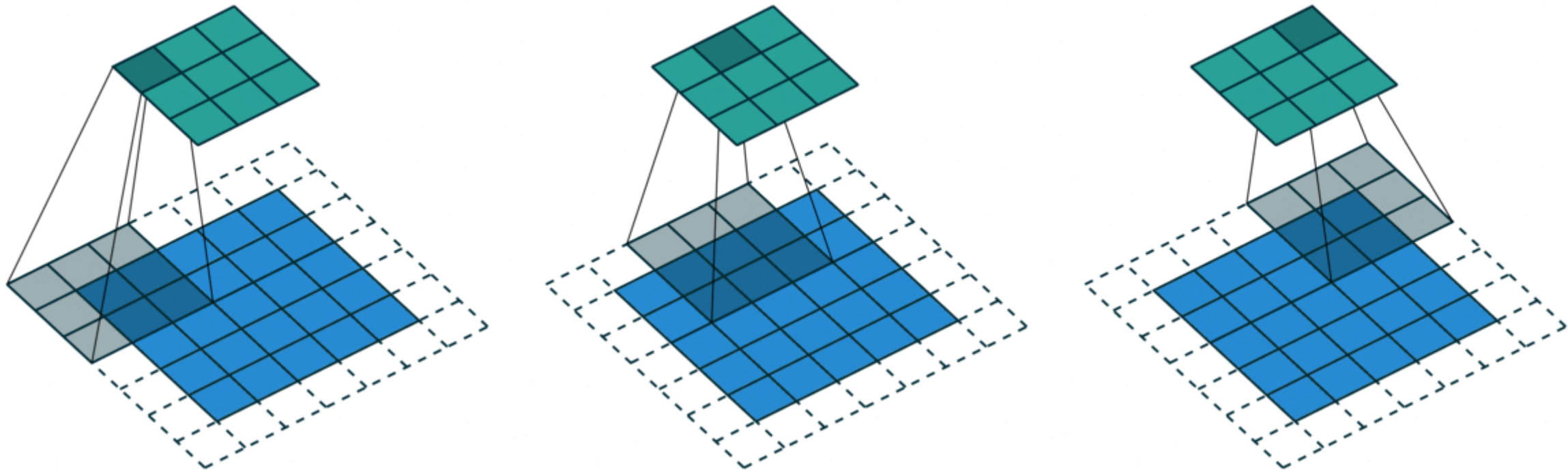- ○ Running live

**fast campus**

# Convolutional Neural Networks

# Contents

- The convolutional layer

- Batch normalization

- Residual connections

- Variations: dilated convolution, deconvolution, separable convolution

**fast campus**

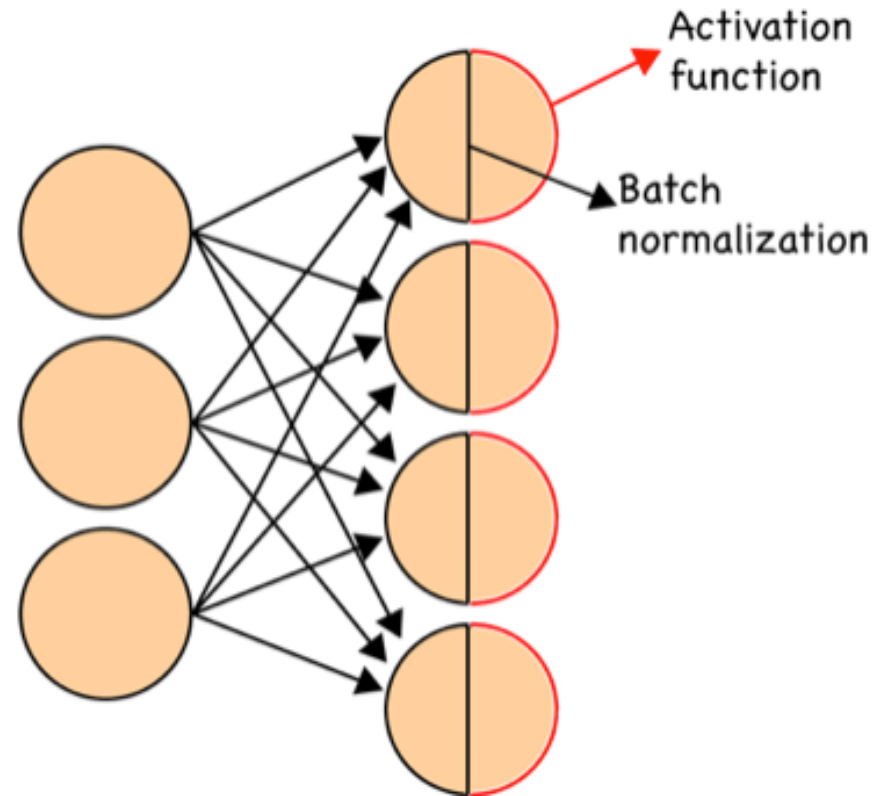# The convolutional layer



```
torch.nn.Conv2d(in_channels=1, out_channels=1, kernel=(3, 3),
stride=2, padding=1, dilation=1, bias=False)
```

fast campus

# Batch normalization

Training Deep Neural Networks is complicated by the fact that **the distribution of each layer's inputs changes during training**, as the parameters of the previous layers change. This slows down the training by requiring lower learning rates and careful parameter initialization, and makes it notoriously hard to train models with saturating nonlinearities. We refer to this phenomenon as **internal covariate shift**, and address the problem by **normalizing layer inputs**.

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift (2015)

**fast campus**

# Batch normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \qquad \text{// scale and shift}$$

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

# Residual connections

## ResNet



Figure 5. A deeper residual function $\mathcal{F}$ for ImageNet. Left: a building block (on 56×56 feature maps) as in Fig. 3 for ResNet-34. Right: a "bottleneck" building block for ResNet-50/101/152.
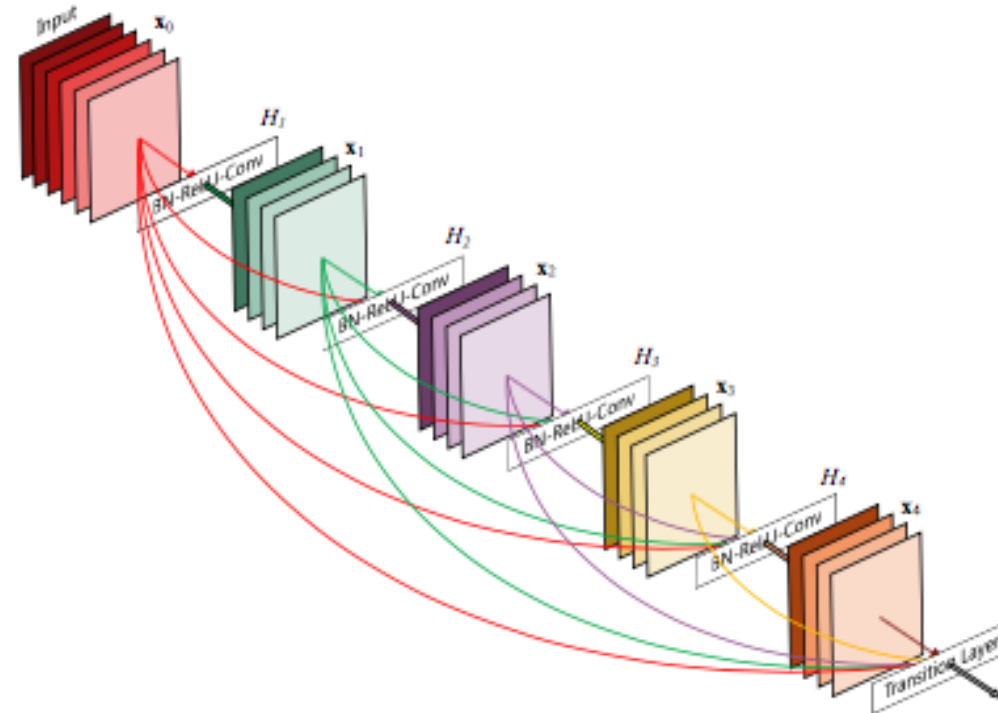
# Residual connections

## DenseNet



**Figure 1:** A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.
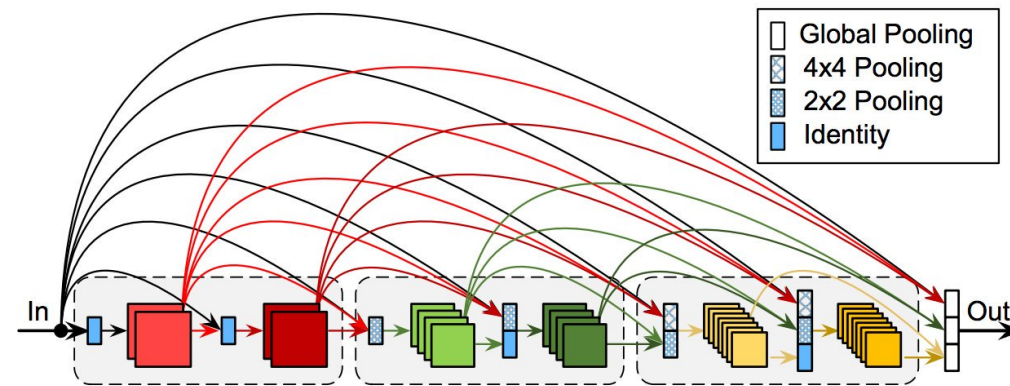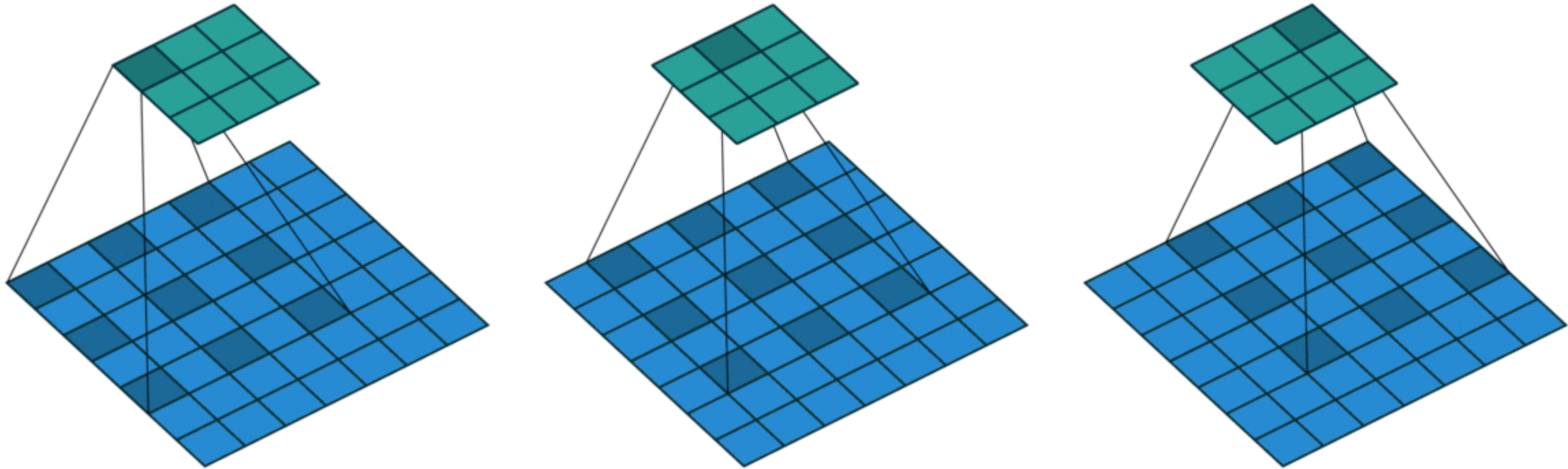
# Residual connections

## CondenseNet



Figure 5. The proposed DenseNet variant. It differs from the original DenseNet in two ways: (1) layers with different resolution feature maps are also directly connected; (2) the growth rate doubles whenever the feature map size shrinks (far more features are generated in the third, yellow, dense block than in the first).

# Variations

- Dilated convolution (a.k.a. atrous convolution

- Deconvolution (a.k.a. transposed convolution, fractionally strided convolution, upconvolution)

- Separable convolution

# Dilated convolution



```
torch.nn.Conv2d(in_channels=1, out_channels=1, kernel=(3, 3),
stride=2, padding=1, dilation=2, bias=False)
```

# Deconvolution



```
torch.nn.ConvTranspose2d(in_channels=1, out_channels=1, kernel=(3,
3), stride=2, padding=1, dilation=1, bias=False)
```

https://github.com/vdumoulin/conv_arithmetic

# Separable convolution

- = Depthwise convolution + Pointwise convolution



Figure 3. Left: Standard convolutional layer with batchnorm and ReLU. Right: Depthwise Separable convolutions with Depthwise and Pointwise layers followed by batchnorm and ReLU.
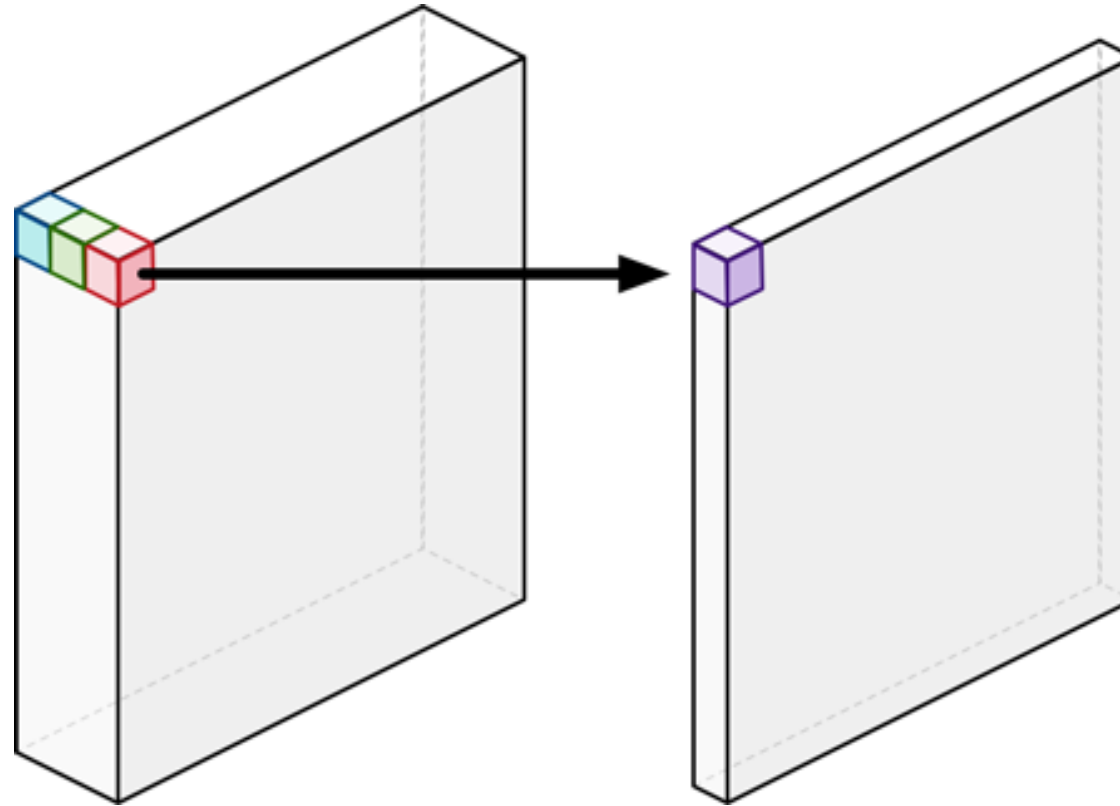
MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications

# Separable convolution



```
regular = torch.nn.Conv2d(in_channels=in_channels,
out_channels=out_channels, kernel_size=3, padding=1)
```

# Separable convolution



```
depthwise = torch.nn.Conv2d(in_channels=in_channels,
out_channels=in_channels, kernel_size=3, padding=1,
groups=in_channels)
```

# Separable convolution



```
pointwise = torch.nn.Conv2d(in_channels=in_channels,
out_channels=out_channels, kernel_size=1)
```

# Separable convolution

- = Depthwise convolution + Pointwise convolution

```python
depthwise = torch.nn.Conv2d(in_channels=in_channels,
out_channels=in_channels, kernel_size=3, padding=1,
groups=in_channels)
pointwise = torch.nn.Conv2d(in_channels=in_channels,
out_channels=out_channels, kernel_size=1)
```

# Recurrent Neural Networks

# Contents

- The recurrent layer

- Gradient vanishing and exploding

- Gradient clipping
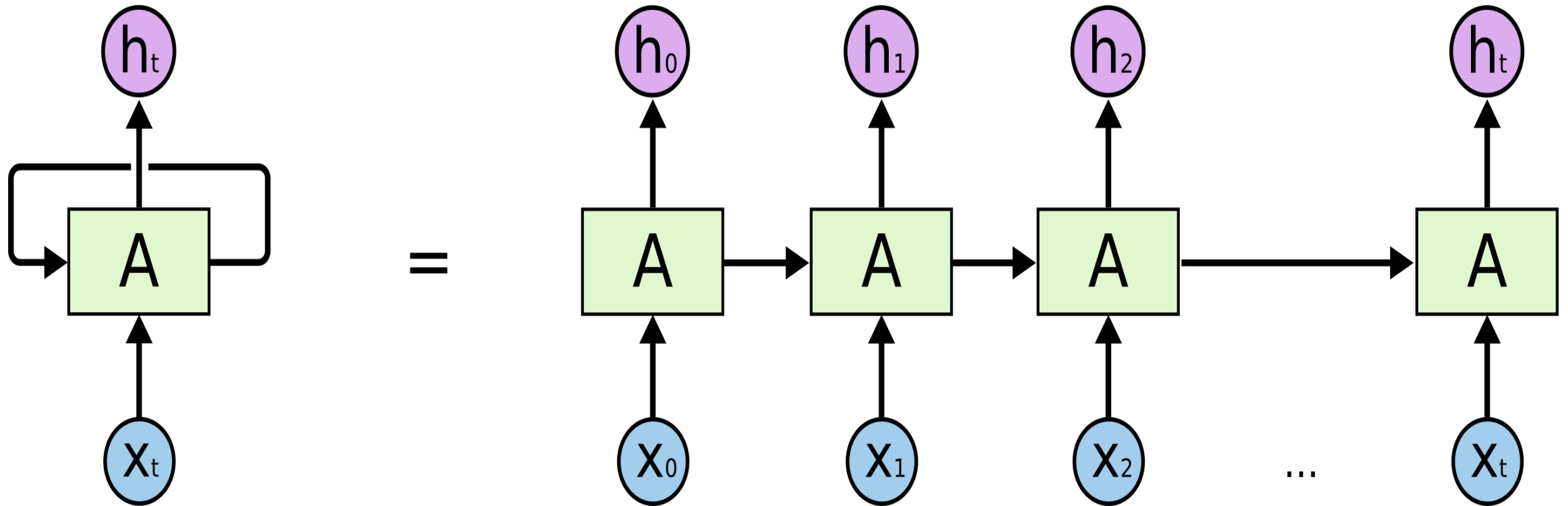
- Variations: Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU)

# The recurrent layer

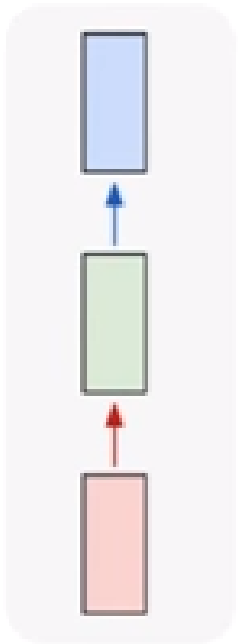A single RNN cell:

# The recurrent layer
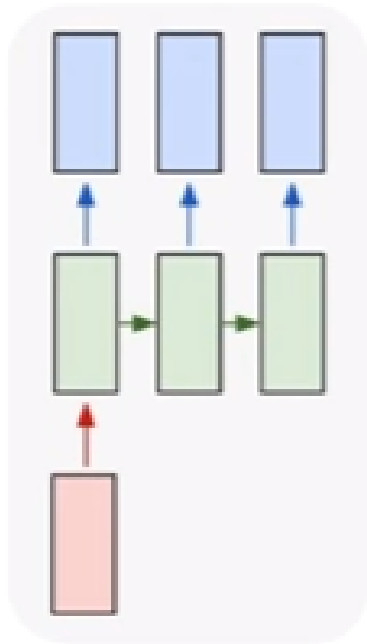
The RNN cell, unrolled:

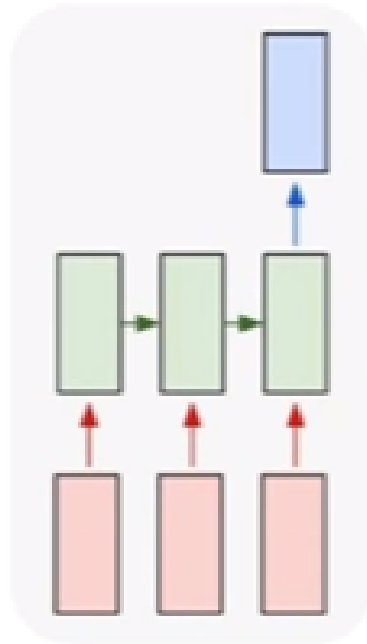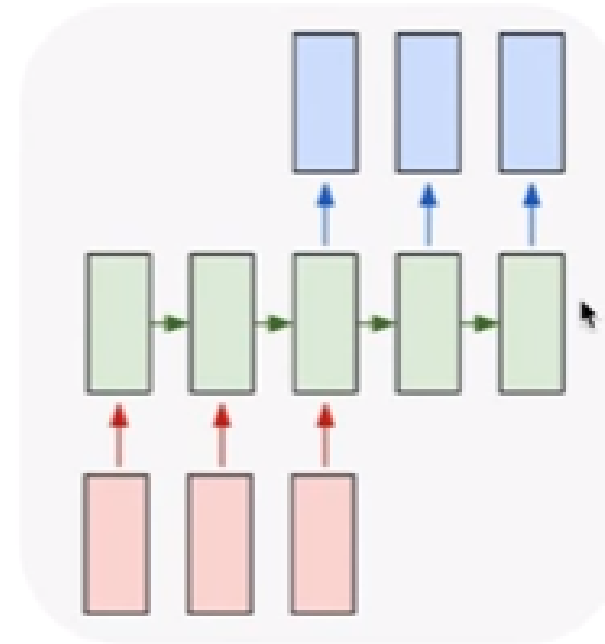# The recurrent layer
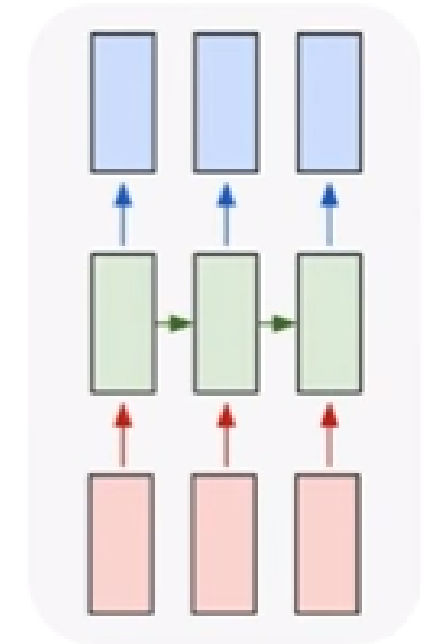
Possibilities:
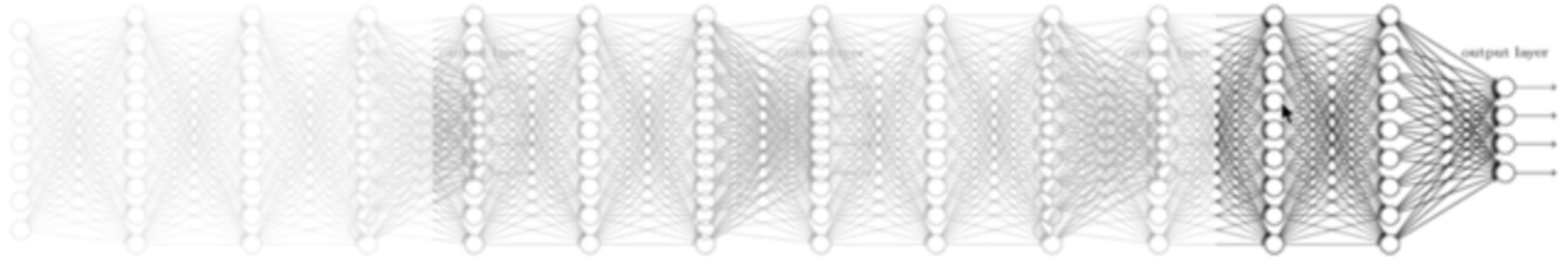


one to one     one to many     many to one     many to many     many to many

# Gradient vanishing and exploding
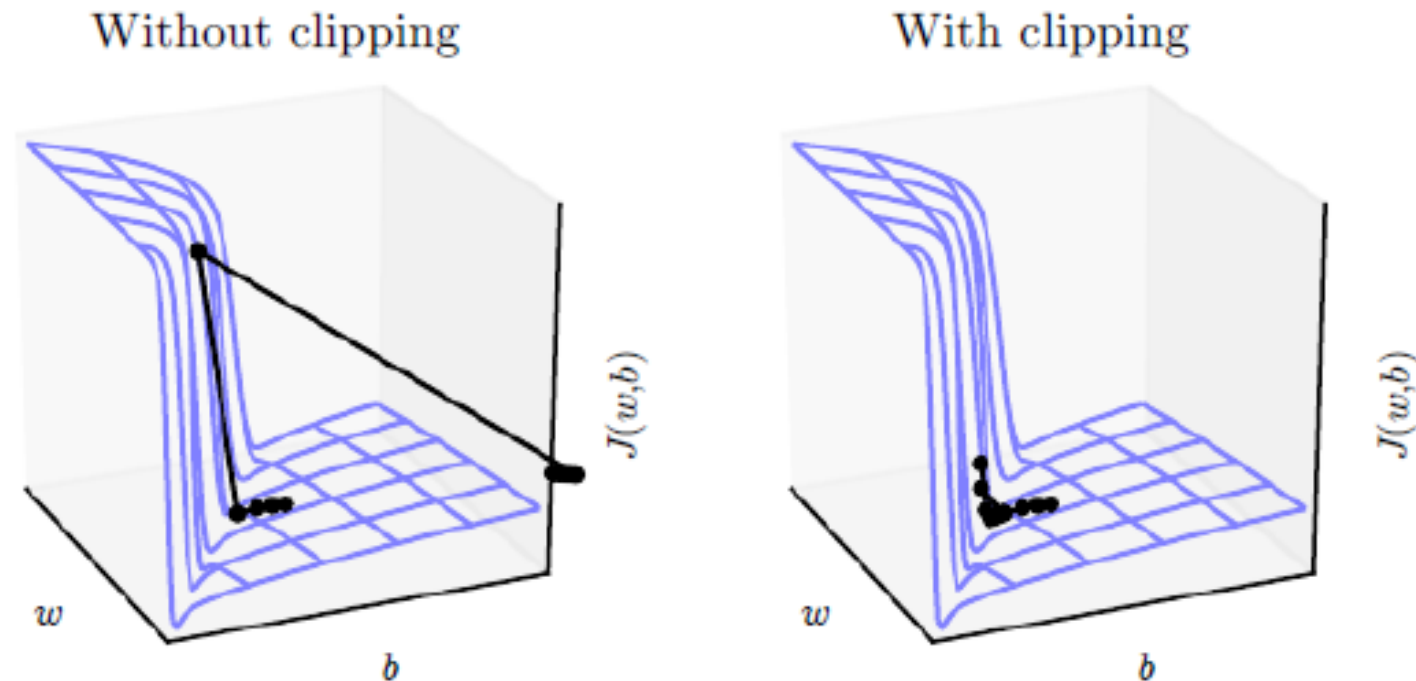
## Gradient vanishing

- Symptom: error signals fail to reach the beginning.
- Solution: let error signals skip layers! (e.g. ResNet, LSTM)

# Gradient vanishing and exploding

**Gradient exploding**

- Symptom: error signals explode on "gradient cliffs".

# Gradient vanishing and exploding

## Gradient exploding

- Solution: set limits on gradients! (e.g. gradient clipping)

**Algorithm 1** Pseudo-code for norm clipping

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$
**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**
$$\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$
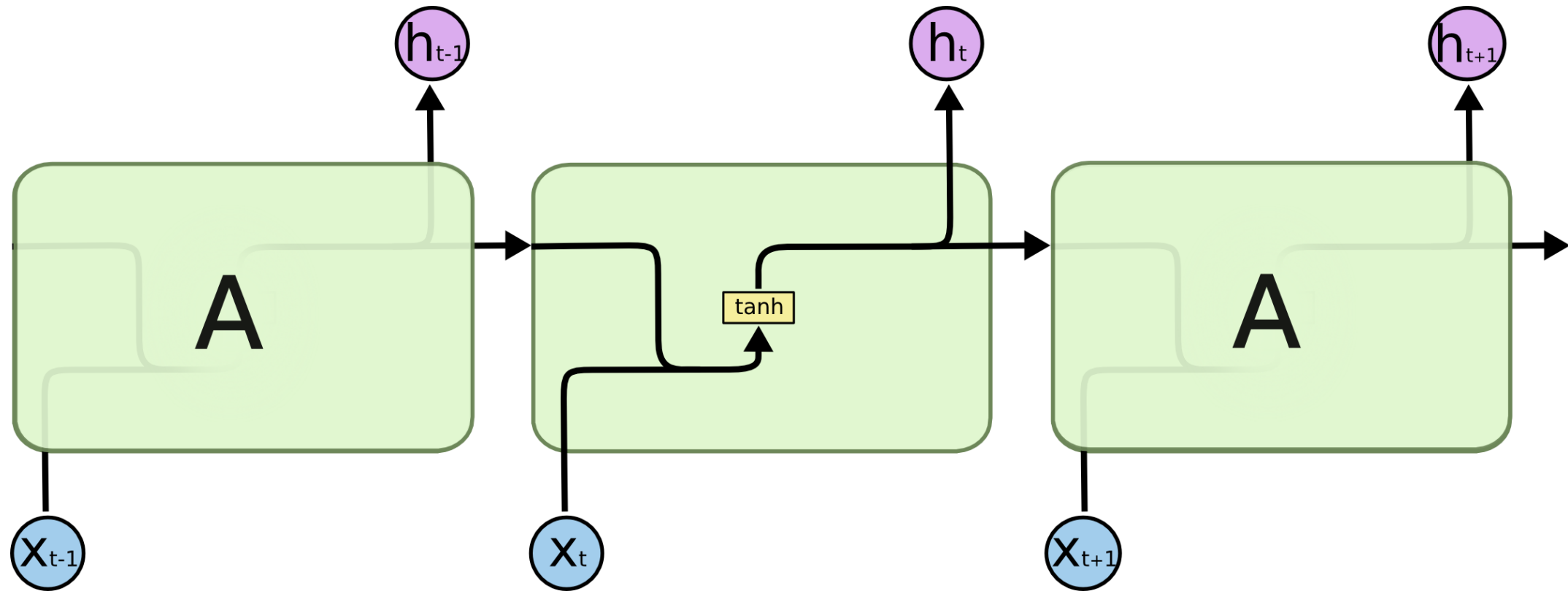**end if**

# Variations of RNN

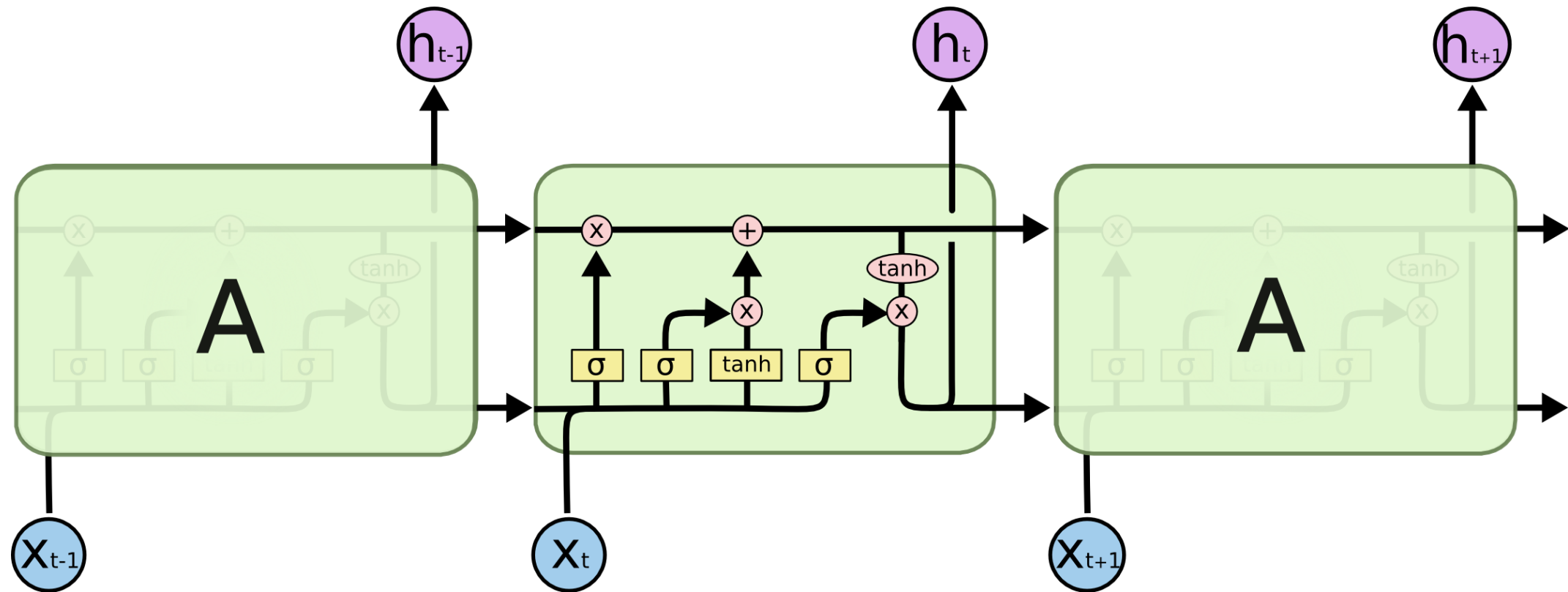- Long Short-Term Memory (LSTM)
- Gated Recurrent Unit (GRU)

# Long Short-Term Memory (LSTM)
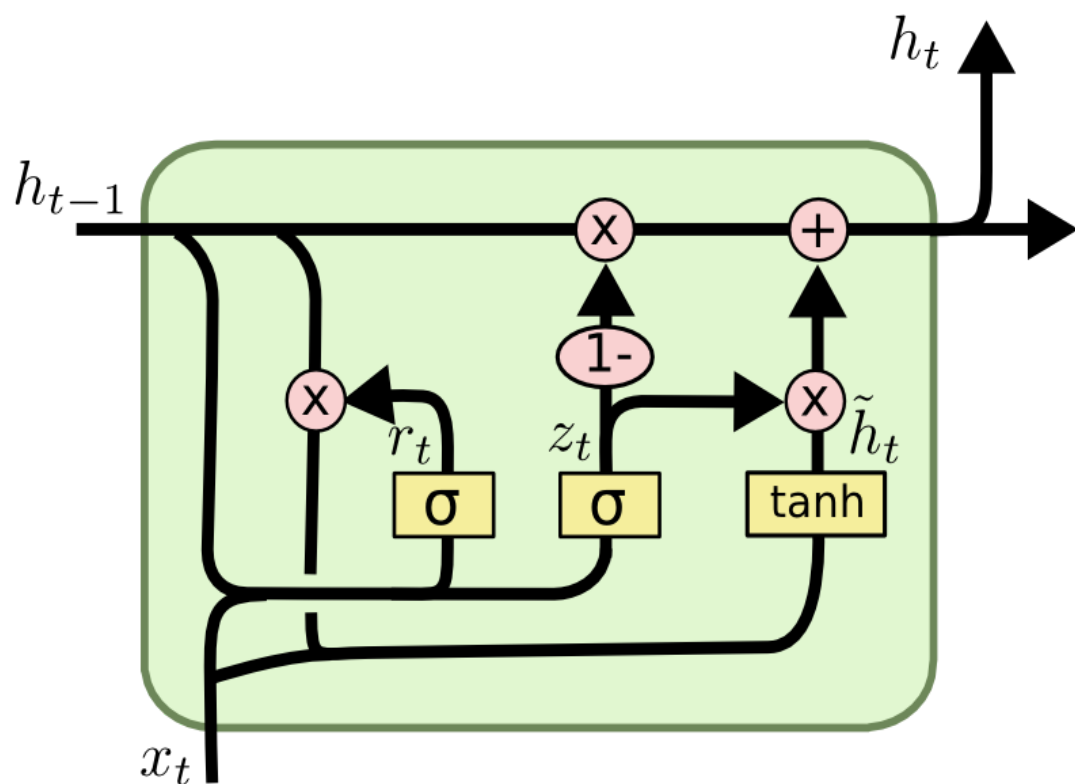
A simple RNN:

# Long Short-Term Memory (LSTM)

LSTM:

# Gated Recurrent Unit (GRU)

- Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation (2014)



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

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

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

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

# Cryptocurrency price prediction

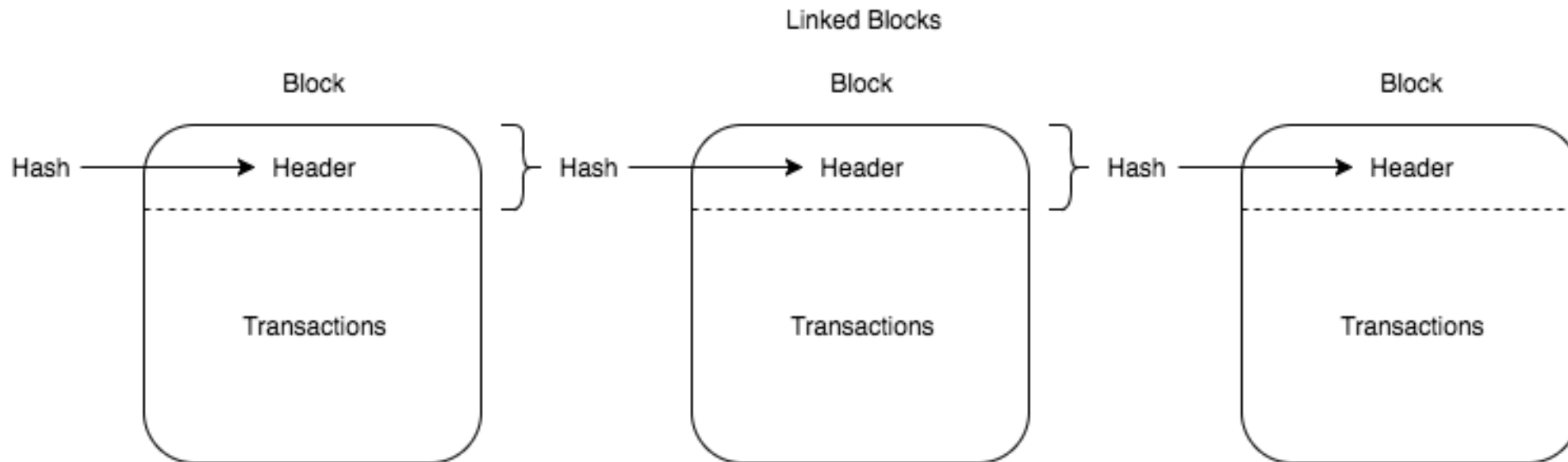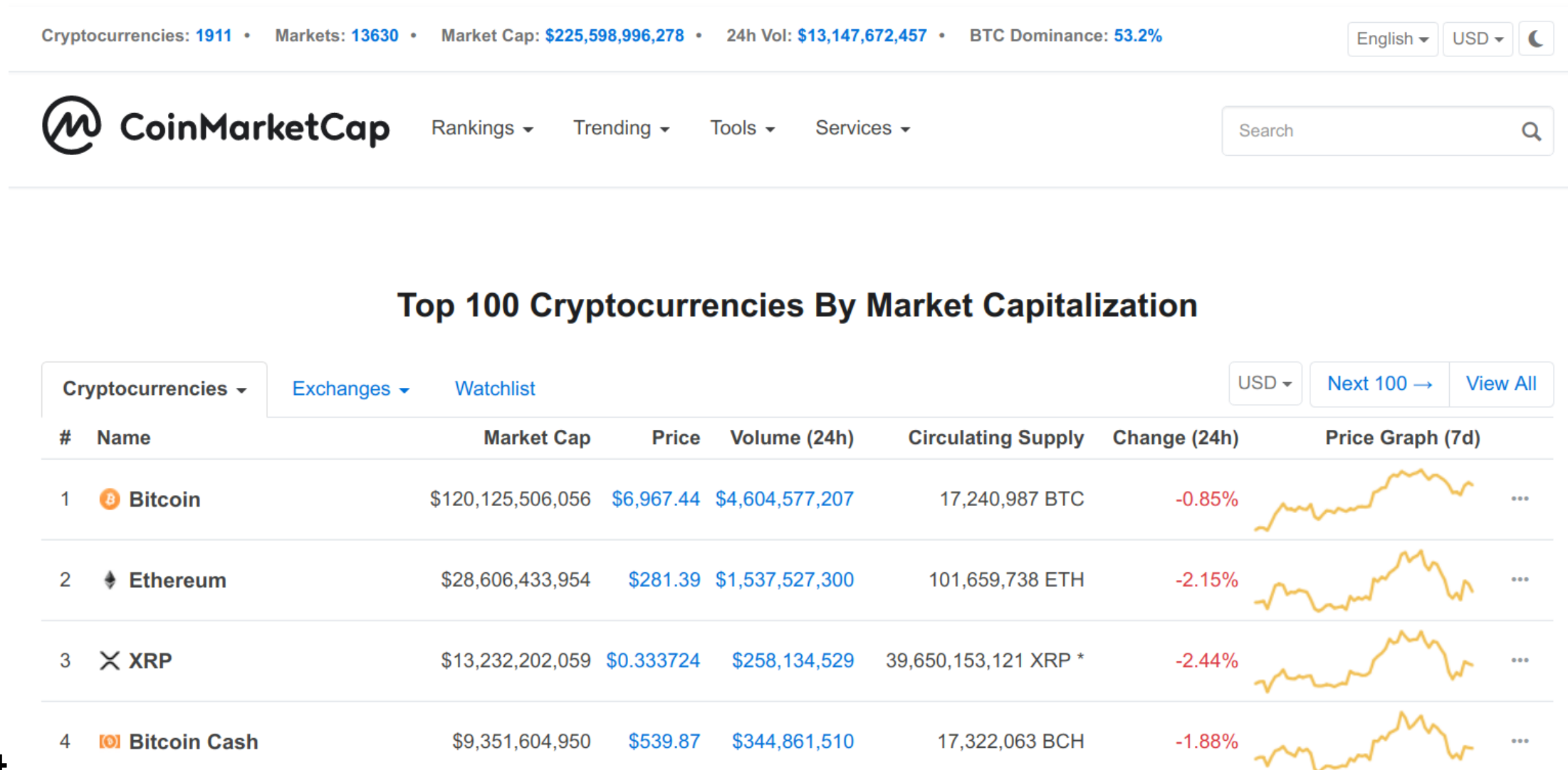https://github.com/juneoh/cryptocurrency_price_prediction

# Contents

- Cryptocurrency 101
- Obtaining and preprocessing the data
- Building our first CNN model
- Building our first RNN model
- Running live

# Cryptocurrency 101

- Cryptocurrency $\neq$ blockchain
- Blockchain-based cryptocurrencies
  - Bitcoin, Ethereum, Bitcoin Cash, …

Linked Blocks

| | Block | | Block | | Block |
|---|---|---|---|---|---|
| Hash → | Header | Hash → | Header | Hash → | Header |
| | Transactions | | Transactions | | Transactions |

**Fast**

# Obtaining and preprocessing the data

# Building our first CNN model

`cnn.py`

# Building our first CNN model

`rnn.py`

# Running live

`python run.py`

# Thank you!