

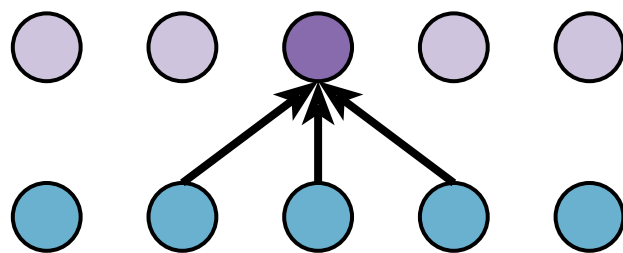
Deep Learning

5. Time series & sequences

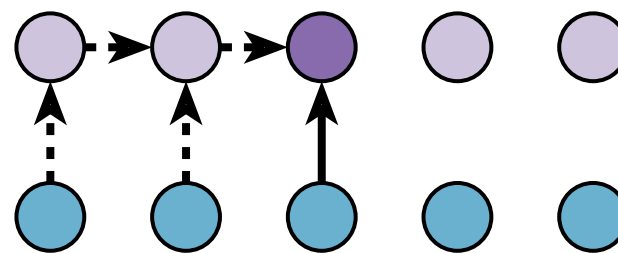
A course @EDHEC
by Romain Tavenard (Prof. @Univ. Rennes 2)

Neural network architectures and inductive biases

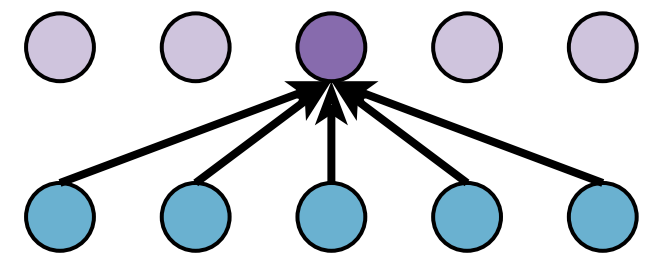
1D Conv.



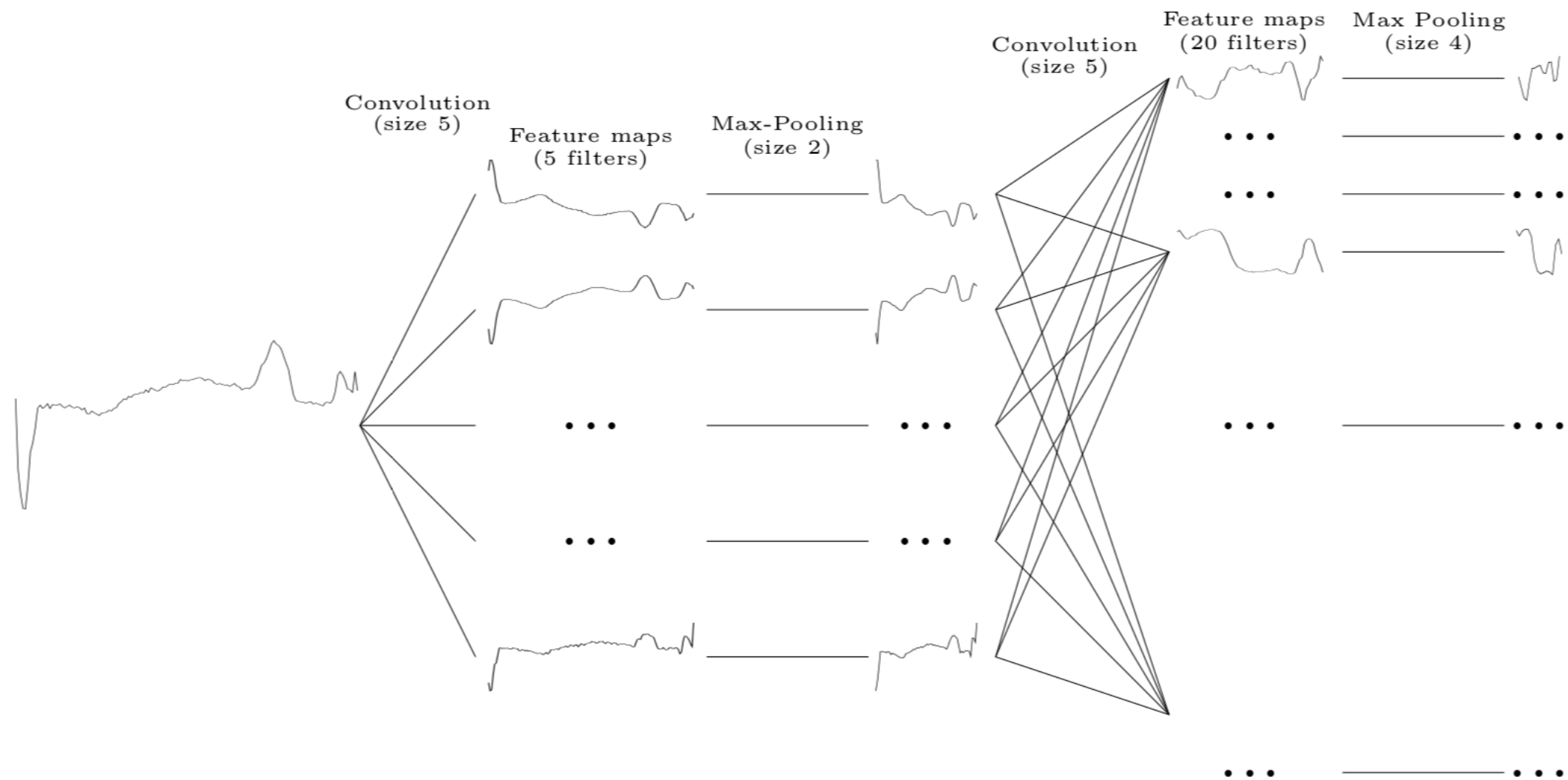
RNN



Self-Attention



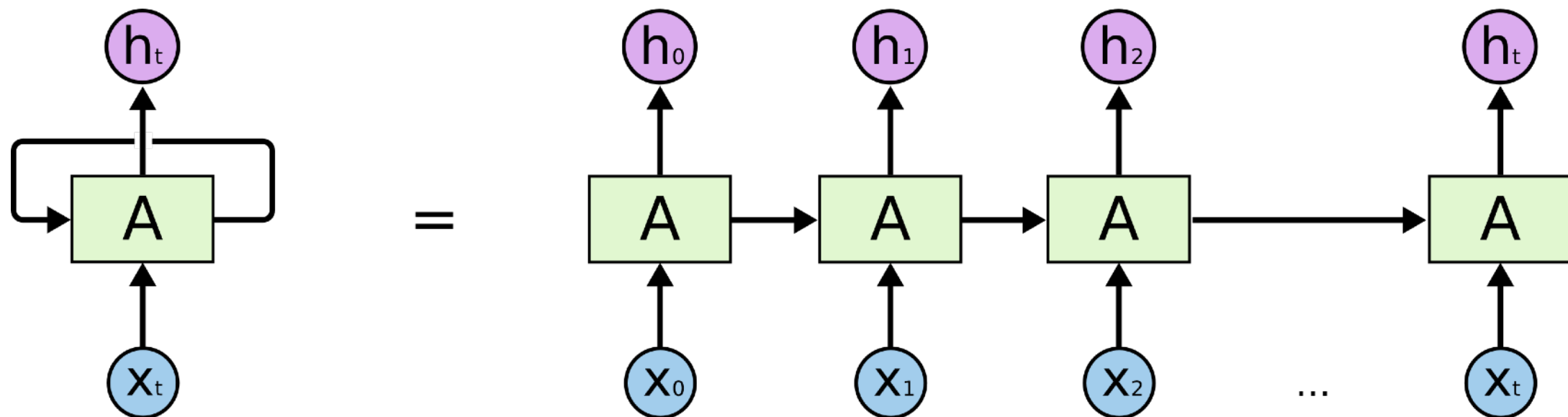
Convolutional neural nets for time series



Source: [Le Guennec *et al.*, 2014]

Recurrent neural nets

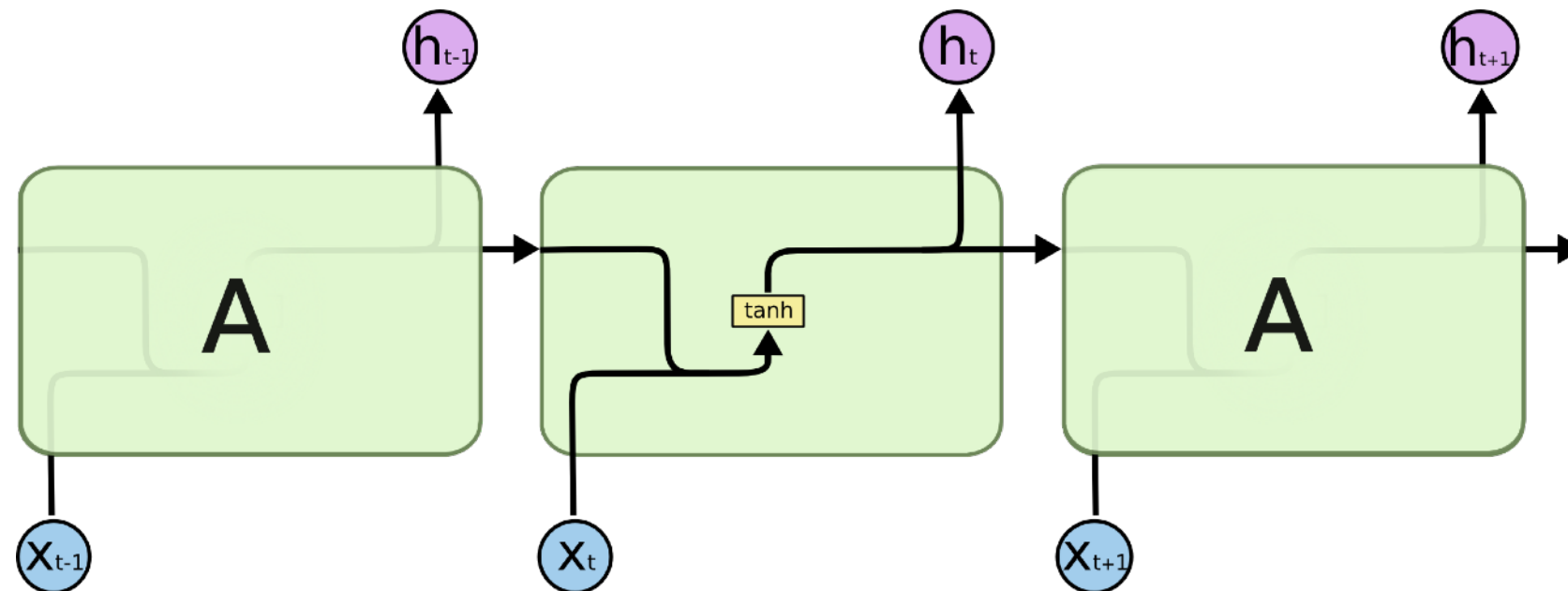
- Very flexible model (any length, let the model learn its memory needs, ...)



Source: [Christopher Olah's blog](#)

Recurrent neural nets

- "Vanilla" RNN in more details



$$h_t = \tanh(W_h h_{t-1} + W_x x_t + b)$$

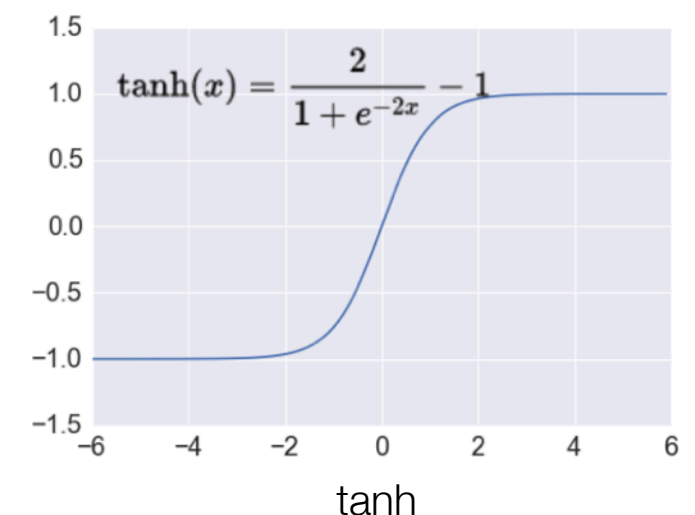
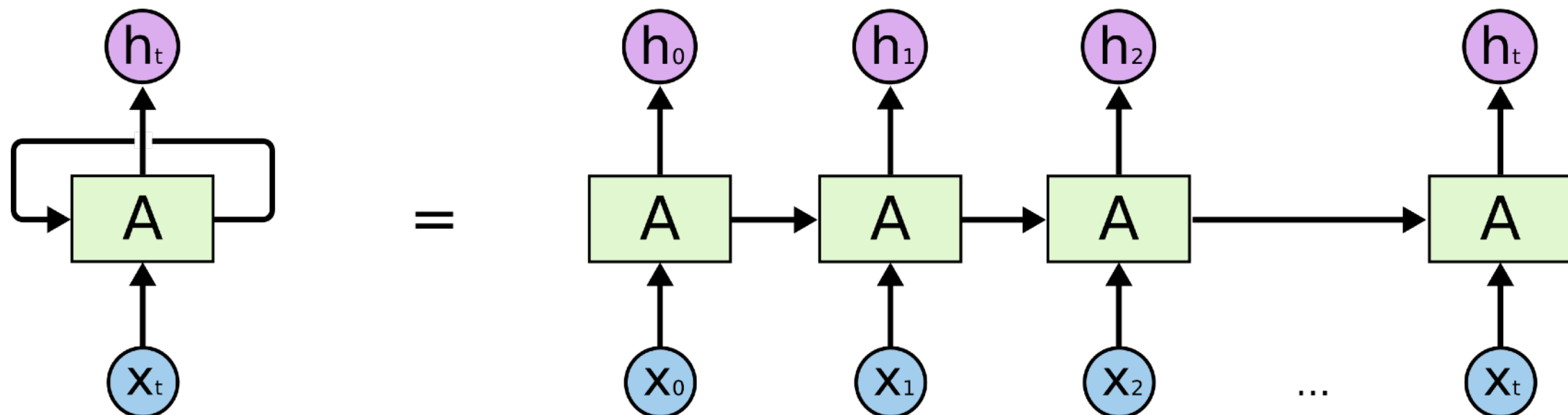


Illustration: RNN cell, source: [Christopher Olah's blog](#)

Recurrent neural nets

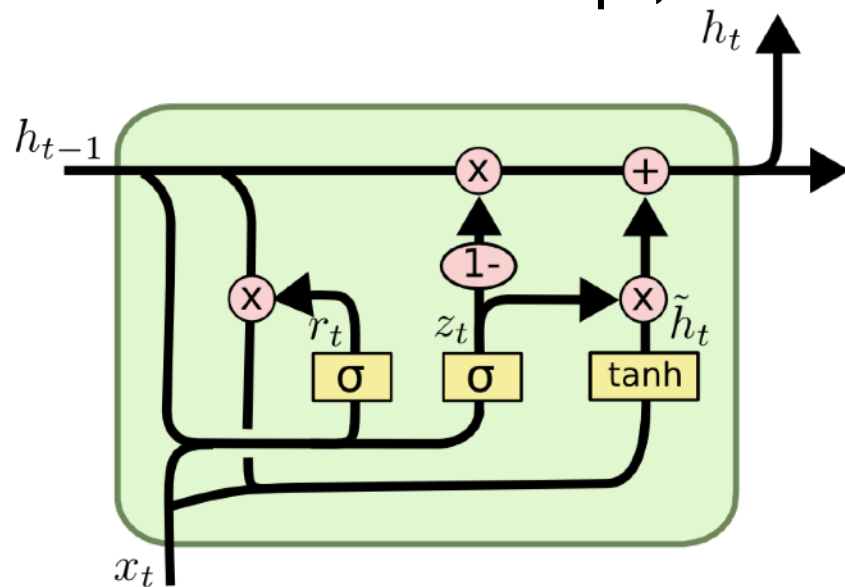
- Very flexible model (any length, let the model learn its memory needs, ...)
- Difficult to learn in practice
 - Slow (lack of parallelism)
 - Vanishing gradients (hard to learn long-term dependencies)



Source: [Christopher Olah's blog](#)

Recurrent neural nets

- Gated Recurrent Unit (GRU)
- Principle
 - At each time step, keep only part of the information



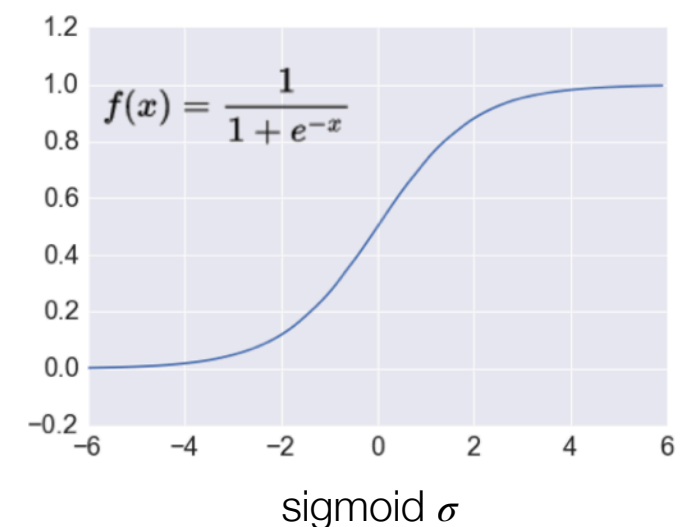
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

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

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

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

Illustration: GRU cell, source: [Christopher Olah's blog](#)



Recurrent neural nets

- Long Short-Term Memory (LSTM)
- Principle: similar to GRU

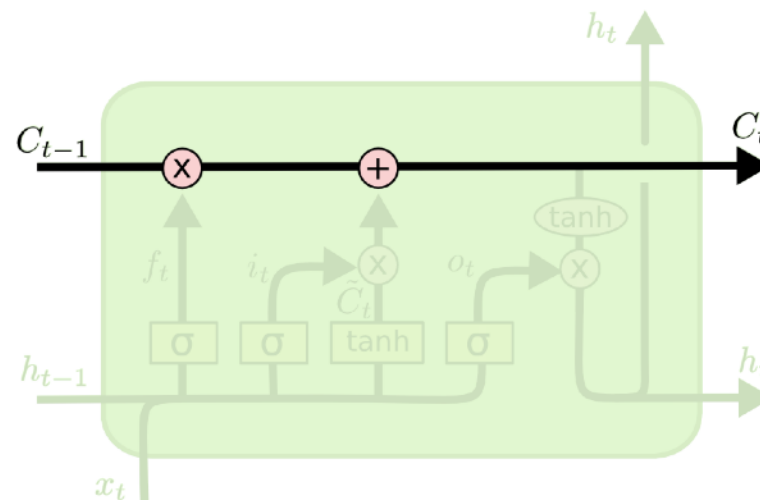
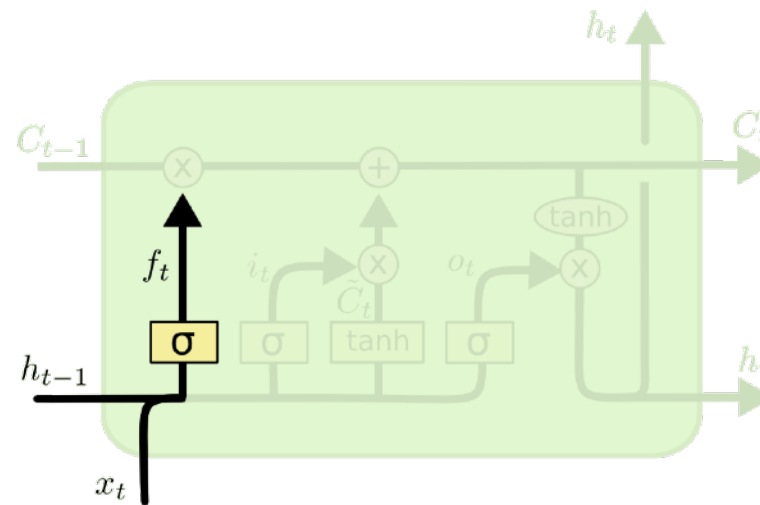


Illustration: LSTM cell, source: [Christopher Olah's blog](#)

Recurrent neural nets

- Long Short-Term Memory (LSTM)
- Principle: similar to GRU

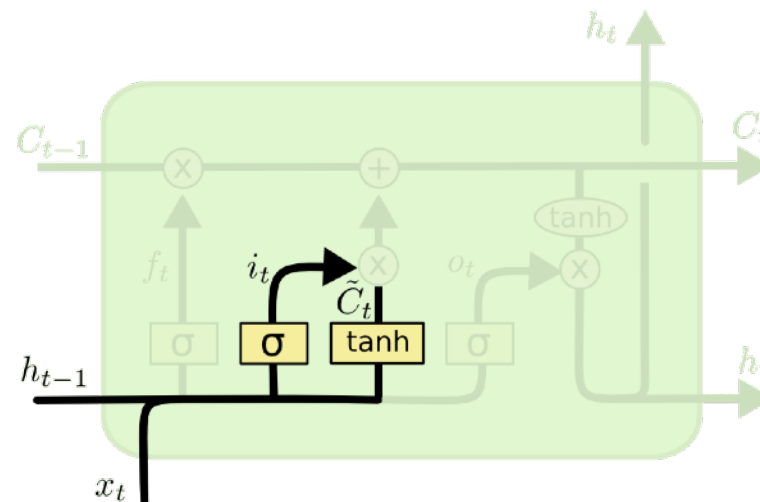


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Illustration: LSTM cell, source: [Christopher Olah's blog](#)

Recurrent neural nets

- Long Short-Term Memory (LSTM)
- Principle: similar to GRU

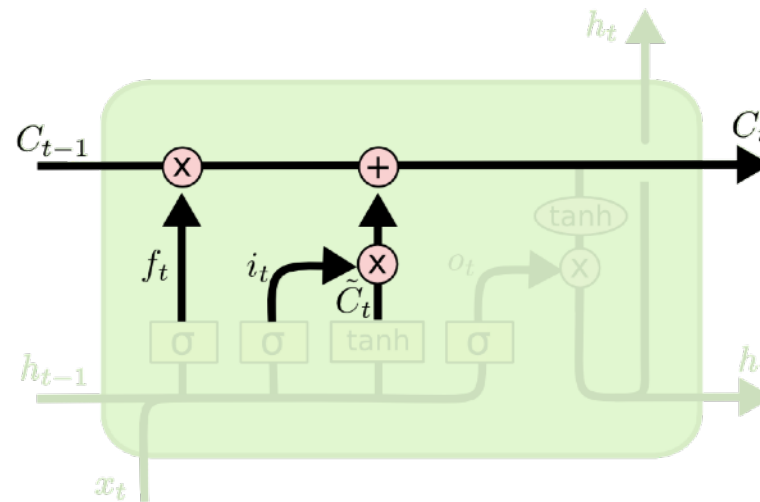


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Illustration: LSTM cell, source: [Christopher Olah's blog](#)

Recurrent neural nets

- Long Short-Term Memory (LSTM)
- Principle: similar to GRU

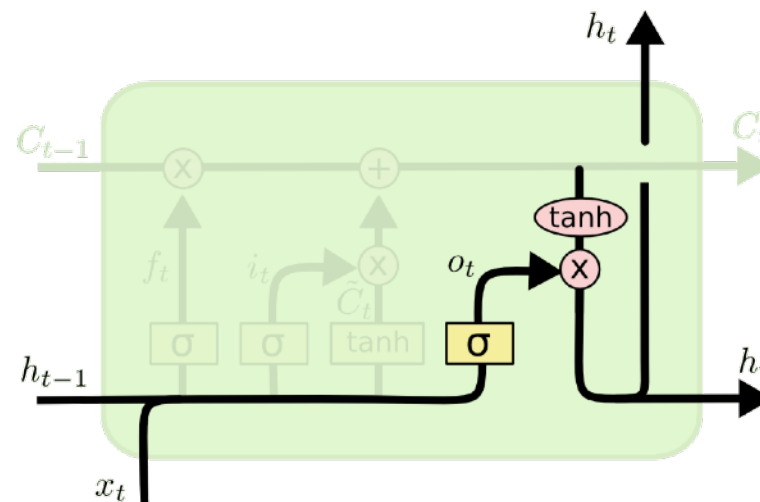


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Illustration: LSTM cell, source: [Christopher Olah's blog](#)

Recurrent neural nets

- Long Short-Term Memory (LSTM)
- Principle: similar to GRU



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

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

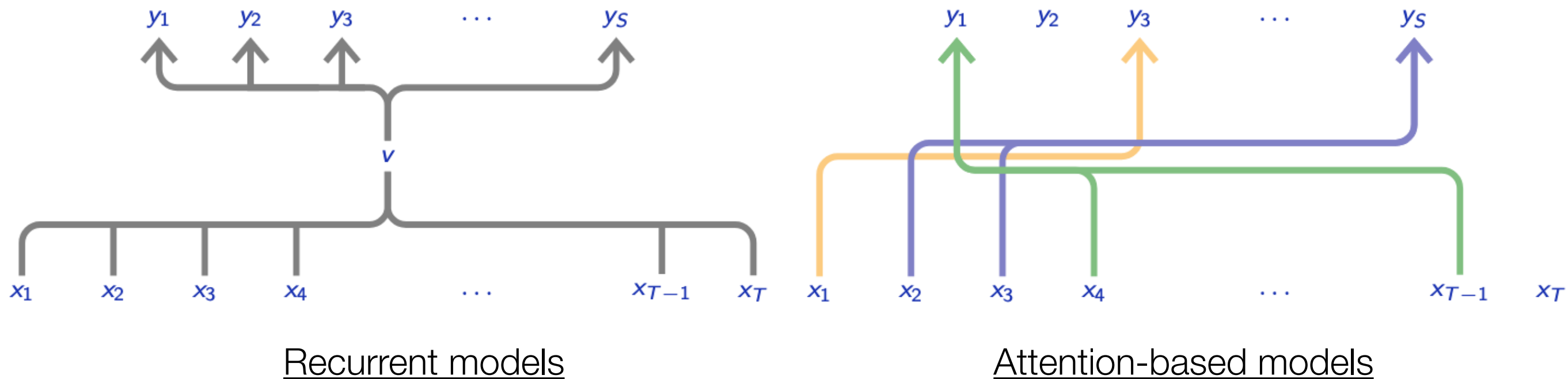
Illustration: LSTM cell, source: [Christopher Olah's blog](#)

Attention-based models

Motivating example

- Consider the following translation task:
 - From English: "**An apple** that had been on the tree in the garden for weeks had finally been **picked up**."
 - To French: "**Une pomme** qui était sur l'arbre du jardin depuis des semaines avait finalement été **ramassée**."
- Both recurrent and convolutional architectures fail at modelling long-range dependencies (for different reasons)
 - Attention aims at tackling this limitation

Attention for Sequence-to-Sequence (seq2seq) tasks



Recurrent models

- Compute a bottleneck representation
- Generate output sequence from this bottleneck

Attention-based models

- For each token in the output sequence, aggregate input features
- Aggregation is importance-based
- Importance depends on the features, not their localization

Assessing importance: Queries, Keys & Values

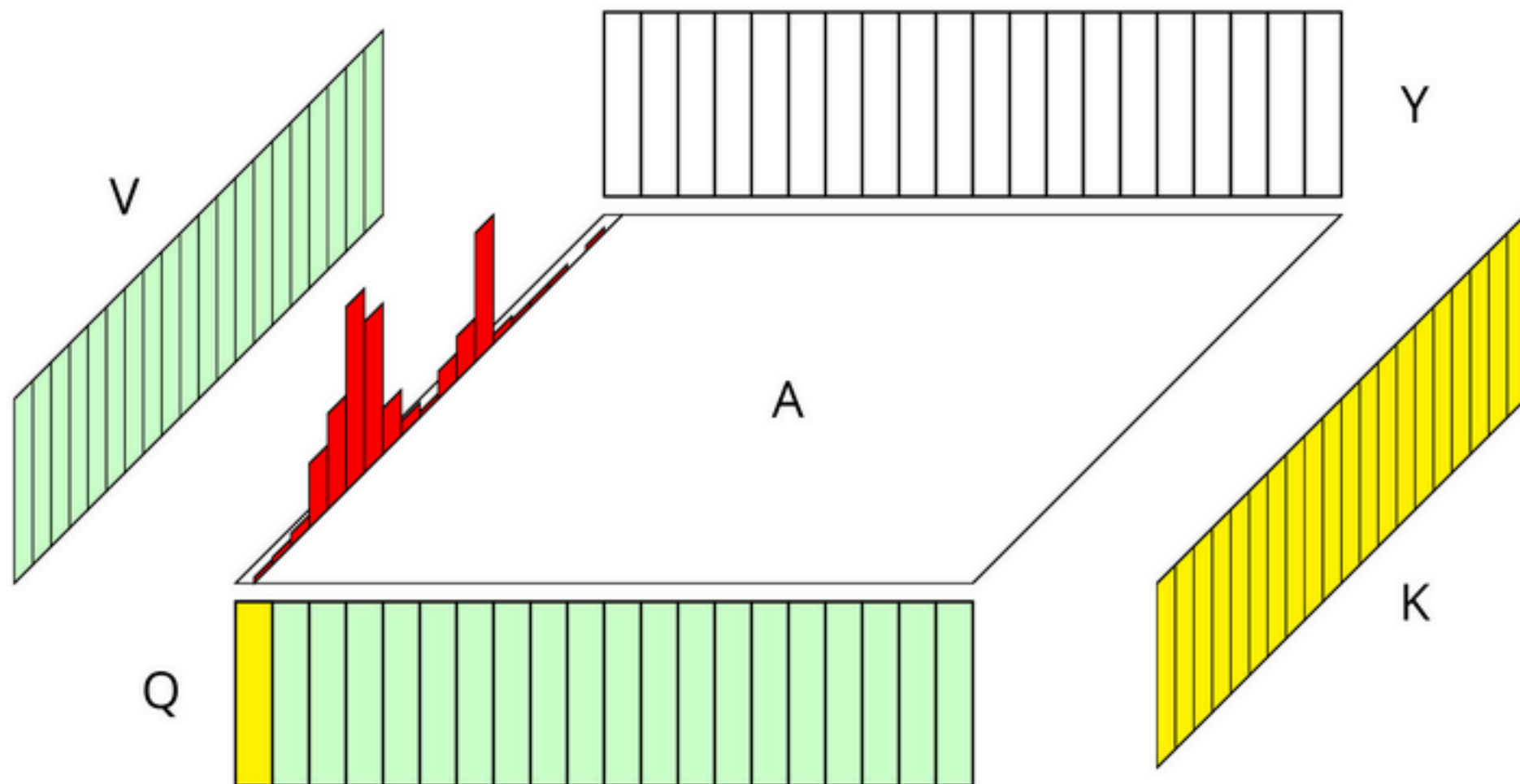
- The **Query** / **Key** / **Value** metaphor: Python dictionaries

```
d = {"a": 12, "b": 7}
print(d["a"])
```

- Queries, Keys & Values: the continuous case
 - Look for keys that are **similar** to the query (not equal)
 - Retrieved value is a weighted average of values associated to keys that are close to the query
 - Could be seen as weighted k-nearest neighbors

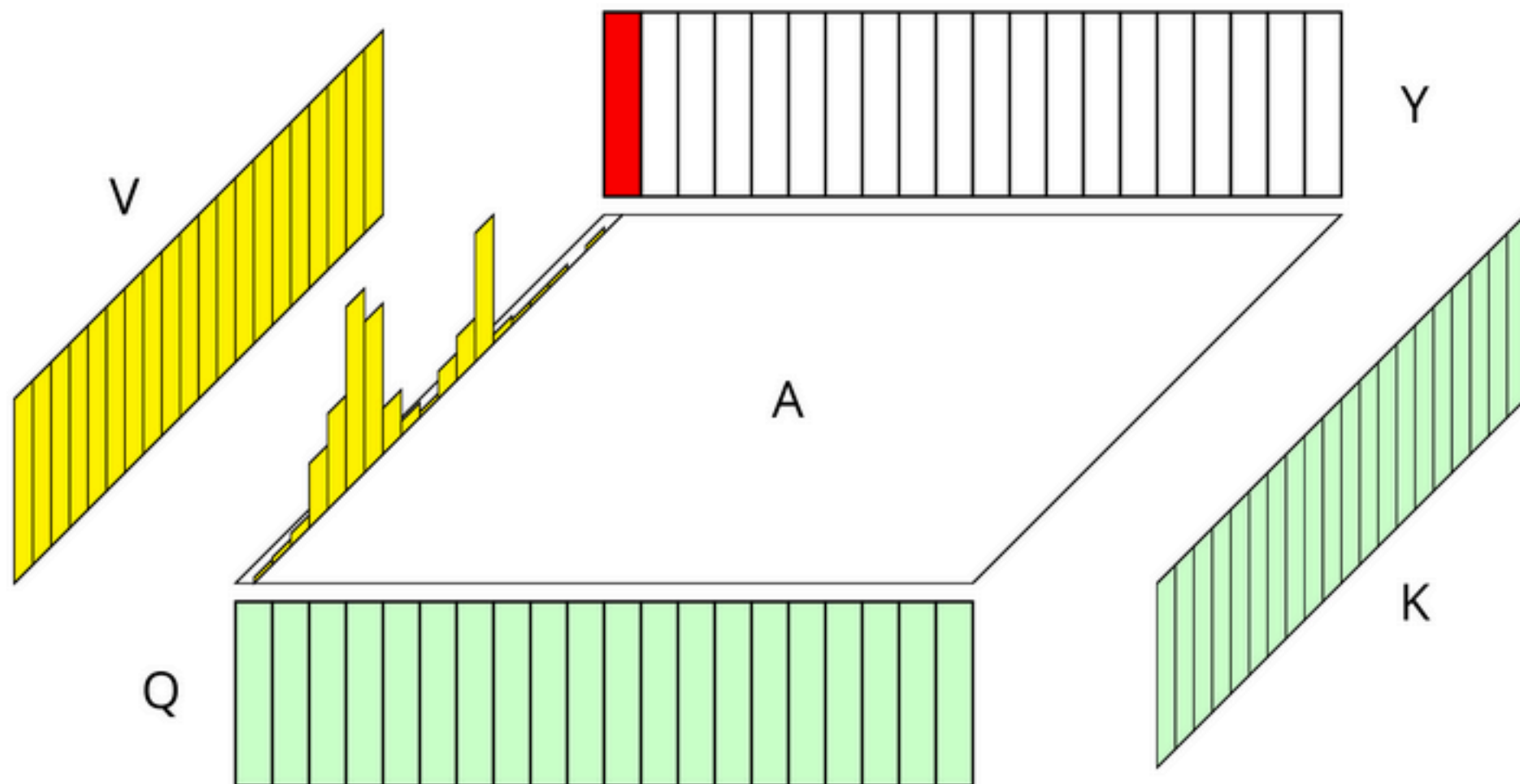
Visual explanation of the Attention mechanism

$$A_{i,j} = \text{softmax}(Q_i \cdot K_j)$$



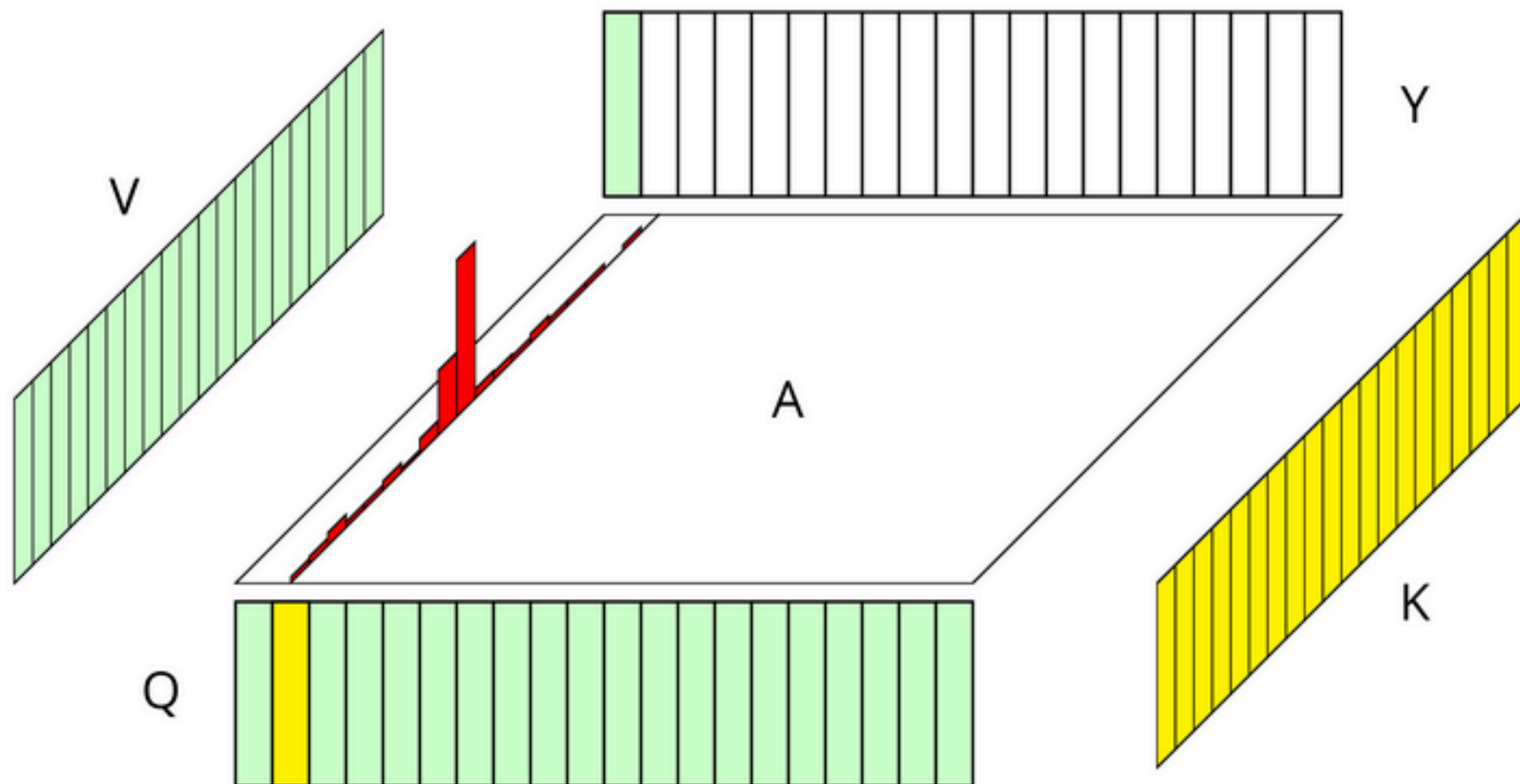
Visual explanation of the Attention mechanism

$$Y_i = \sum_j A_{i,j} V_j$$



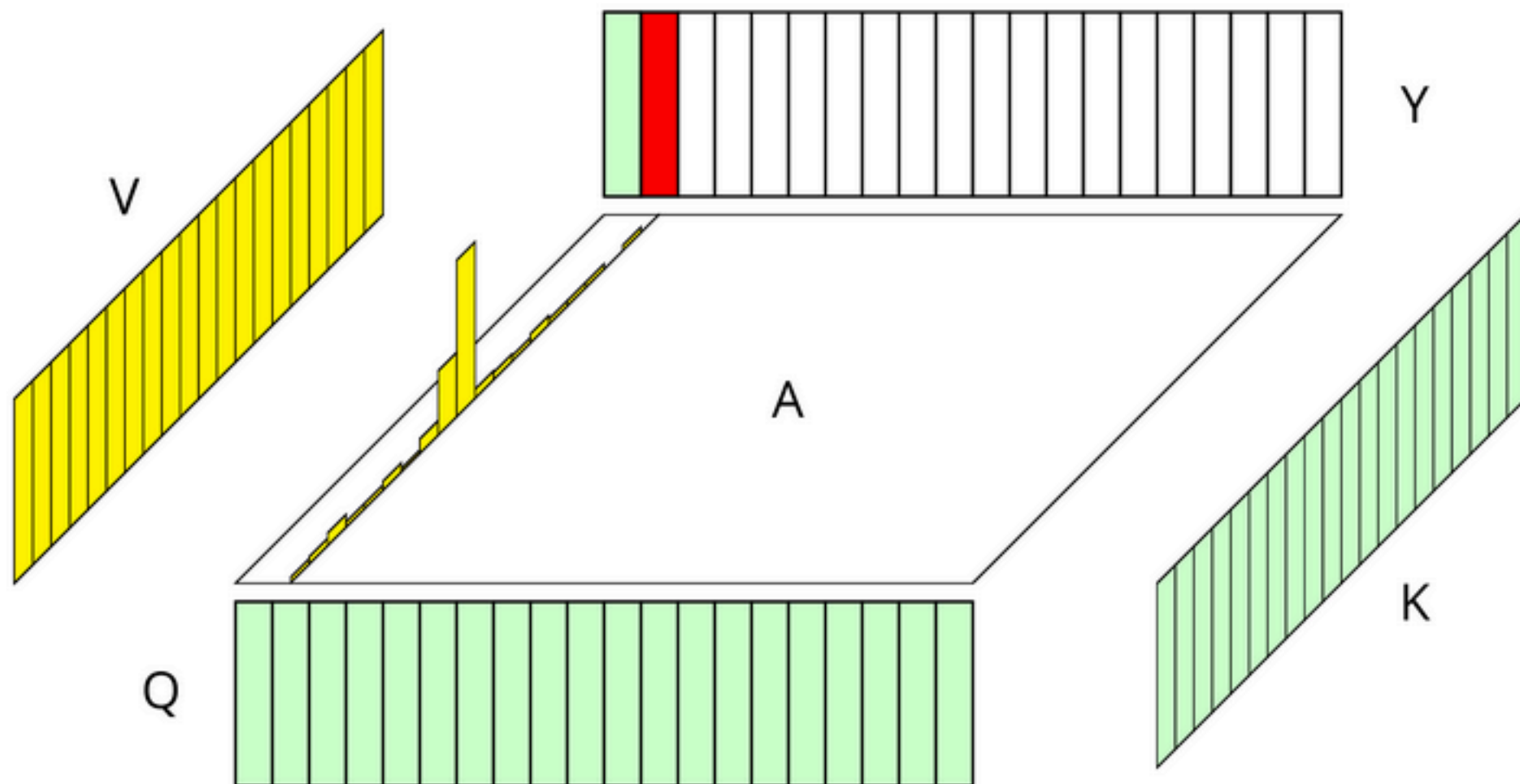
Visual explanation of the Attention mechanism

$$A_{i,j} = \text{softmax}(Q_i \cdot K_j)$$



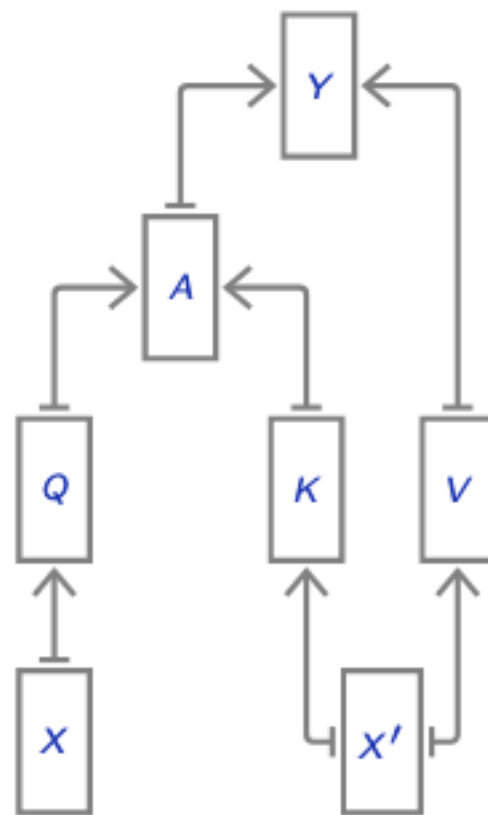
Visual explanation of the Attention mechanism

$$Y_i = \sum_j A_{i,j} V_j$$



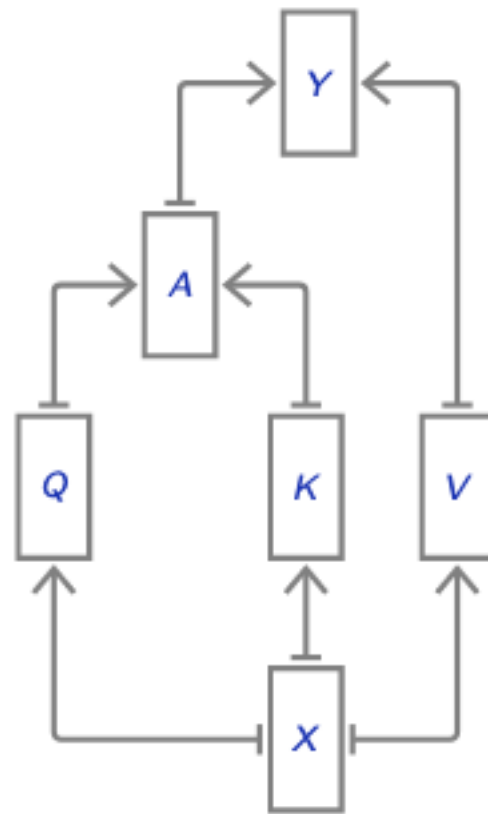
Standard attention layers

- Standard Attention layers
 - take as inputs 2 sequences X and X'
 - output a sequence Y



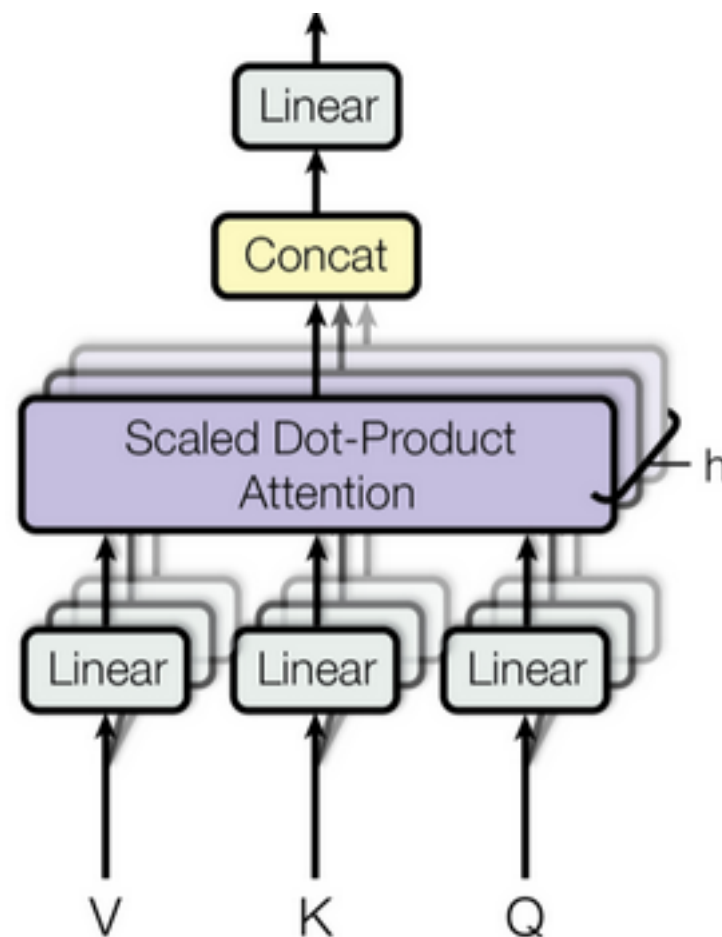
Standard attention layers

- In self-attention layers, we have $X=X'$

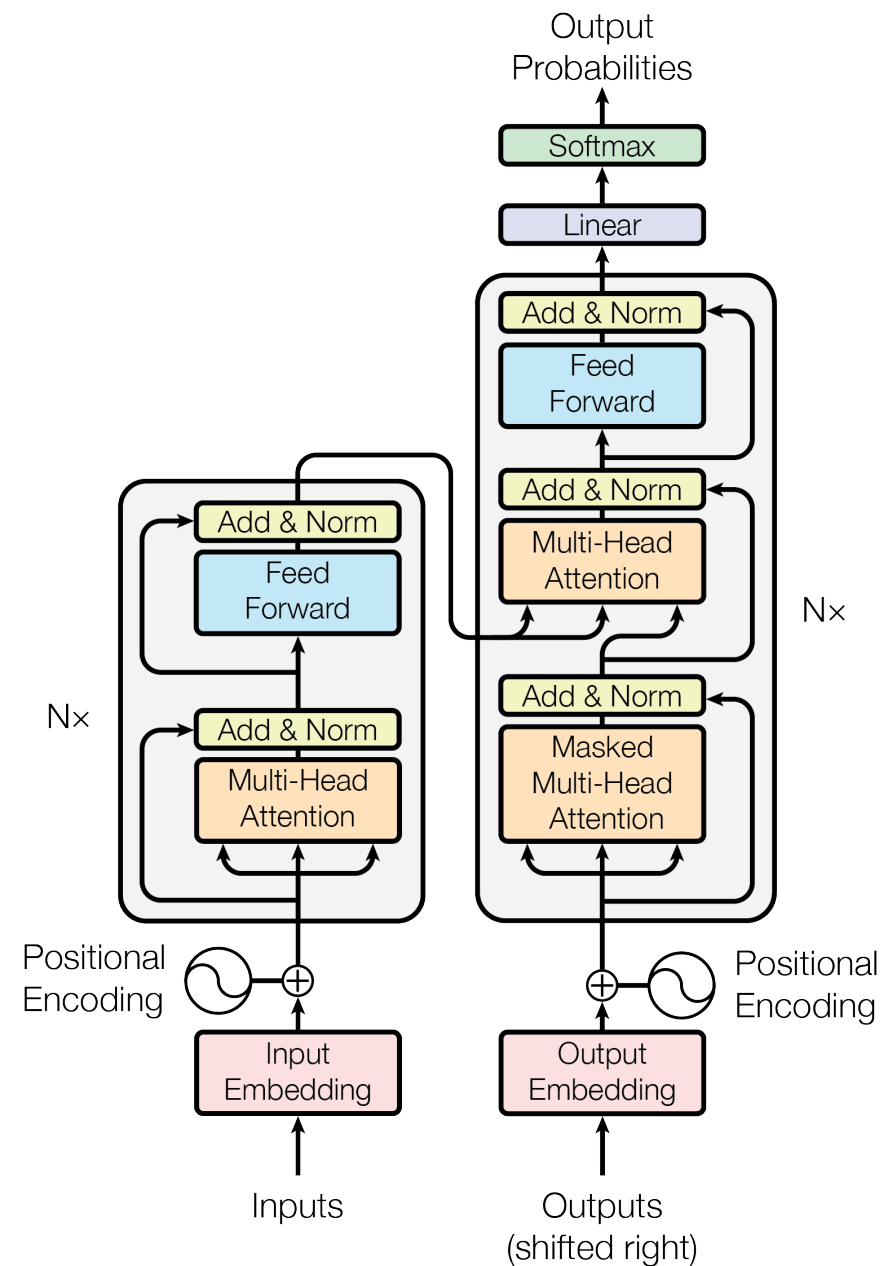


Standard attention layers

- In multi-head attention layers,
 - several (h here) such blocks operate in parallel
 - Their output is concatenated in feature dimension



The Transformer: The typical attention-based architecture



Summary

- 1d-CNN, RNN and Attention-based models can be used
 - depends on the context
 - slightly different underlying assumptions
 - locality (ConvNets)
 - sequentiality (RNNs)
 - relationships between any set of items in the sequence (Attention-based models)
- 1d-CNN are faster to train