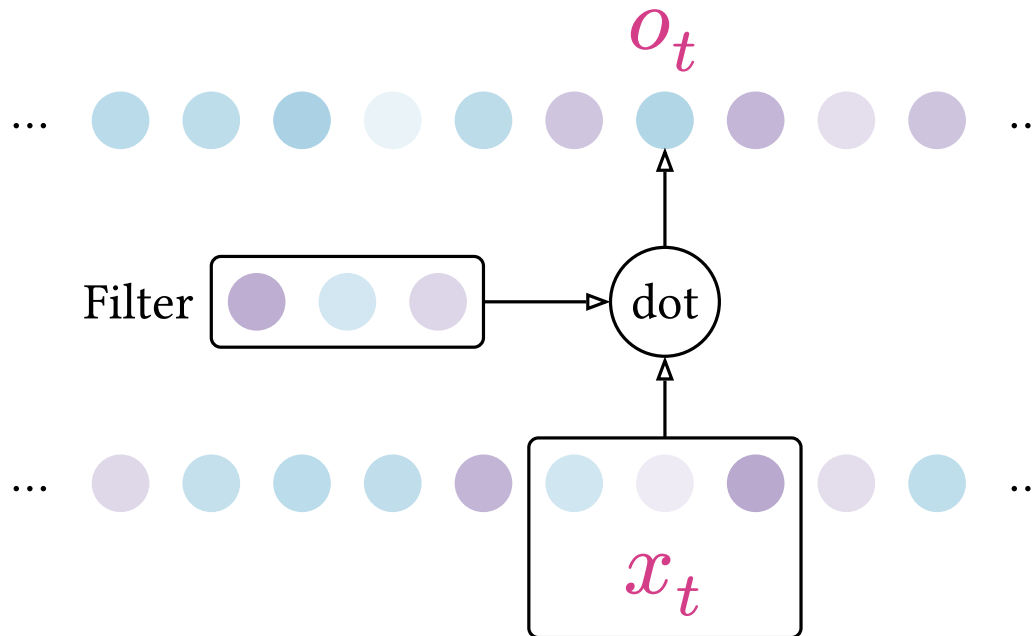# Deep Learning for Time Series
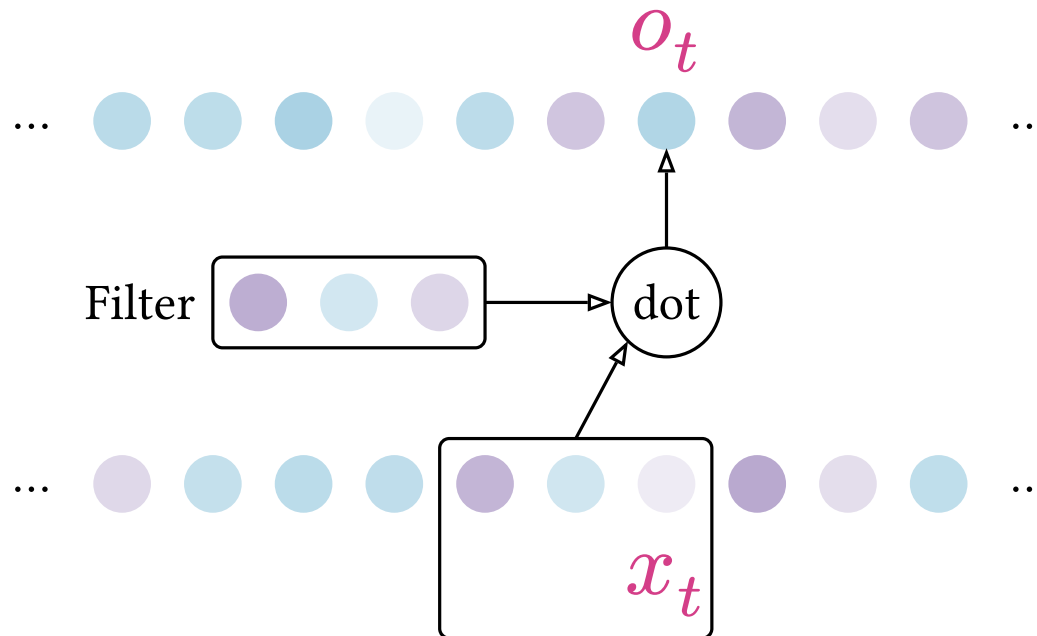
## Session 2: ConvNets and Recurrent architectures

Romain Tavenard

# Convolutional architectures

- Basic time series processing: 1d convolutions (over time)
- Limited receptive field: co-localization matters

$o_t$

... ● ● ● ● ● ● ● ● ● ● ...

Filter [ ● ● ● ] → (dot)
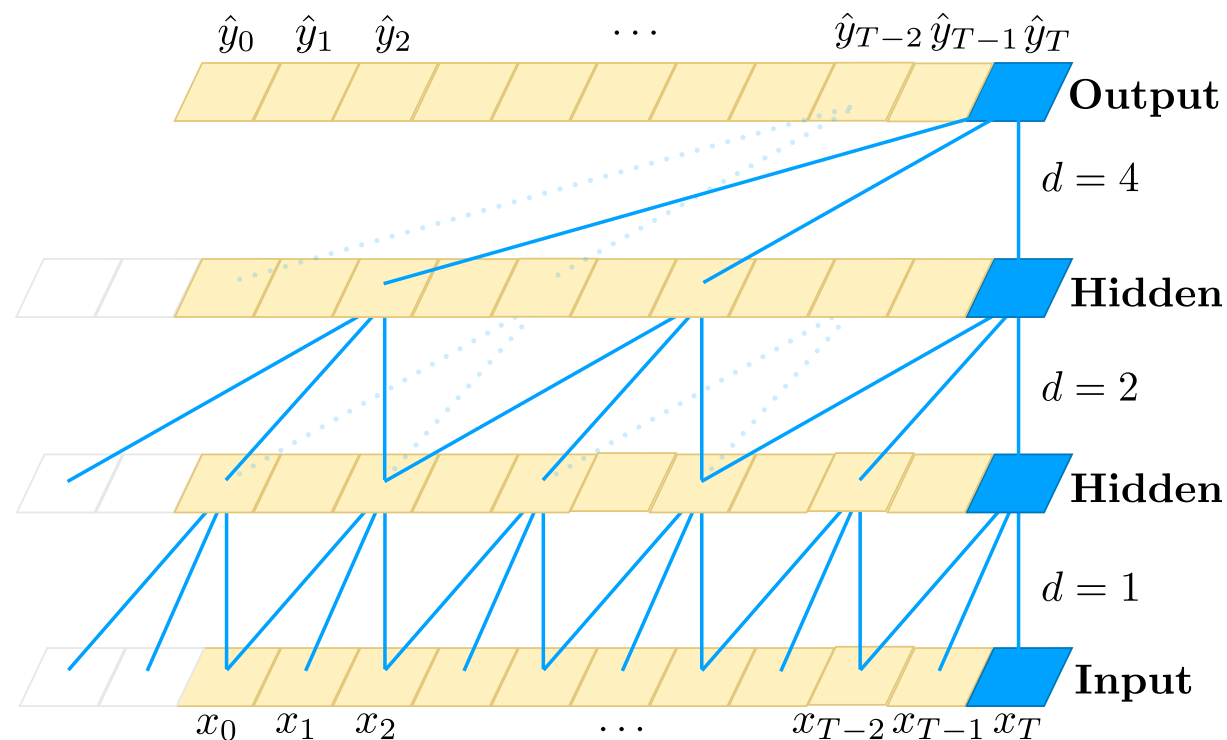
... ● ● ● ● ● ● ● ● ● ● ...

$x_t$

- Forecasting tasks: cannot access the future
- Causal convolution: convolve on past information alone (asymmetric window)
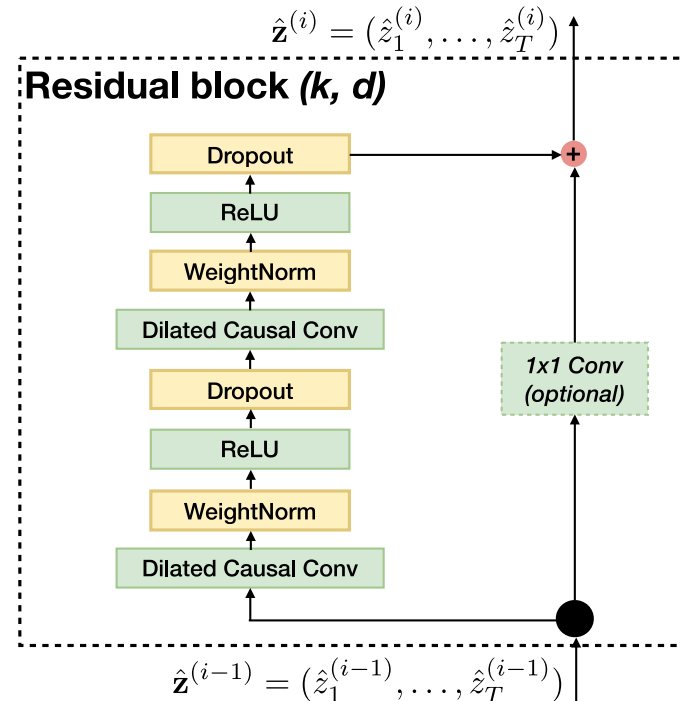
# Temporal Convolution Network (TCN)

- Main idea: cascade dilated causal convolutions

  $\Rightarrow$ Larger receptive field

$$\hat{y}_0 \quad \hat{y}_1 \quad \hat{y}_2 \quad \cdots \quad \hat{y}_{T-2}\ \hat{y}_{T-1}\ \hat{y}_T$$

**Output**

$d = 4$

**Hidden**

$d = 2$

**Hidden**

$d = 1$

**Input**

$$x_0 \quad x_1 \quad x_2 \quad \cdots \quad x_{T-2}\ x_{T-1}\ x_T$$

Source: "An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling", Bai et al.
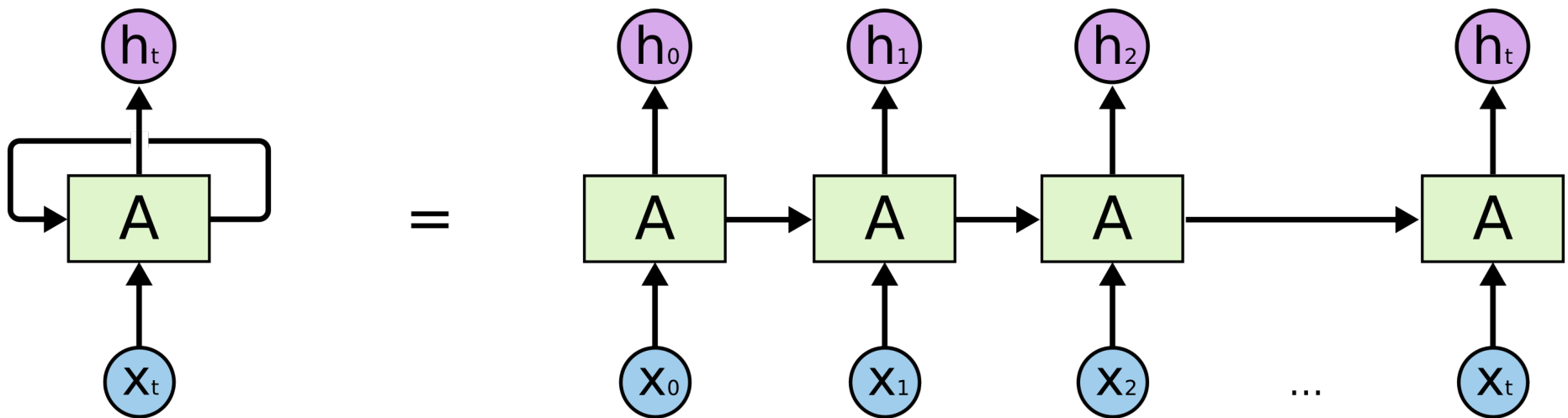
# Temporal Convolution Network (TCN)

- Additional improvements:
  - ▸ Residual connections
    $\Rightarrow$ Multi-resolution analysis
  - ▸ Normalization+Dropout layers



Source: "An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling", Bai et al.

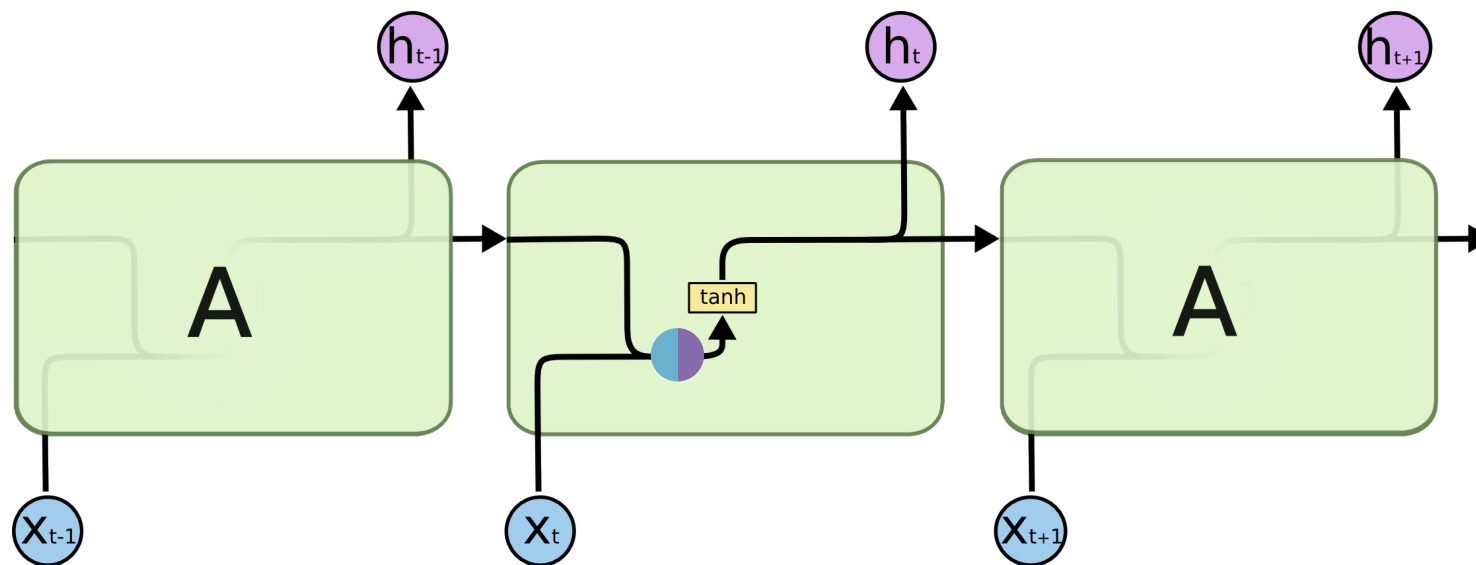# Recurrent architectures

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



Source: Christopher Olah's blog

- Hidden state is computed as:

$$h_t = \varphi(\;\bullet\;)$$



Source: Christopher Olah's blog

$\bullet\; x_t \quad \bullet\; h_{t-1} \quad \bullet\;$ Linear mixing of $x_t$ and $h_{t-1}$

- 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) or exploding gradients (if $\varphi$ is unbounded)

- At each time step, keep only part of the information
  - ‣ Through **gating mechanism**
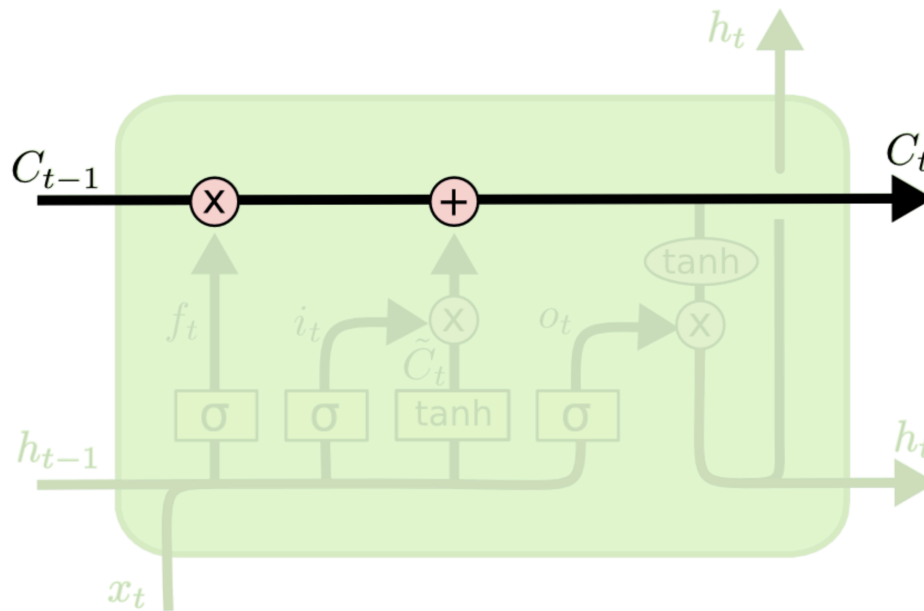


$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right)$$

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

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

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

Source: Christopher Olah's blog
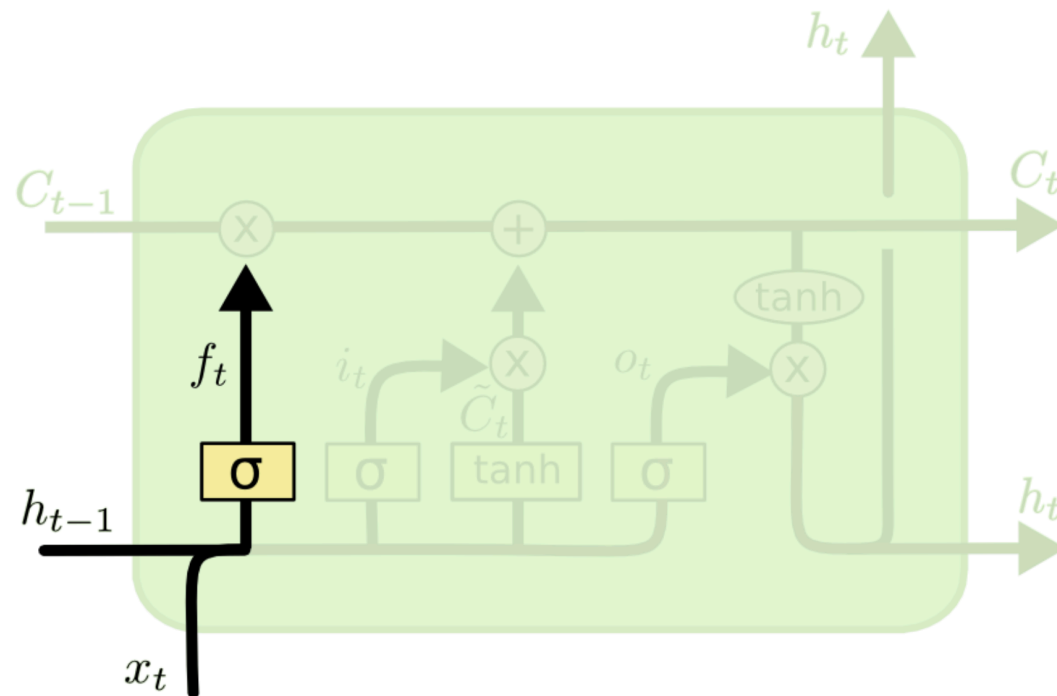
# Long Short Term Memory (LSTM)

- Similar ideas as in GRUs, but:
  - ‣ an additional *cell state* $C_t$



Source: Christopher Olah's blog

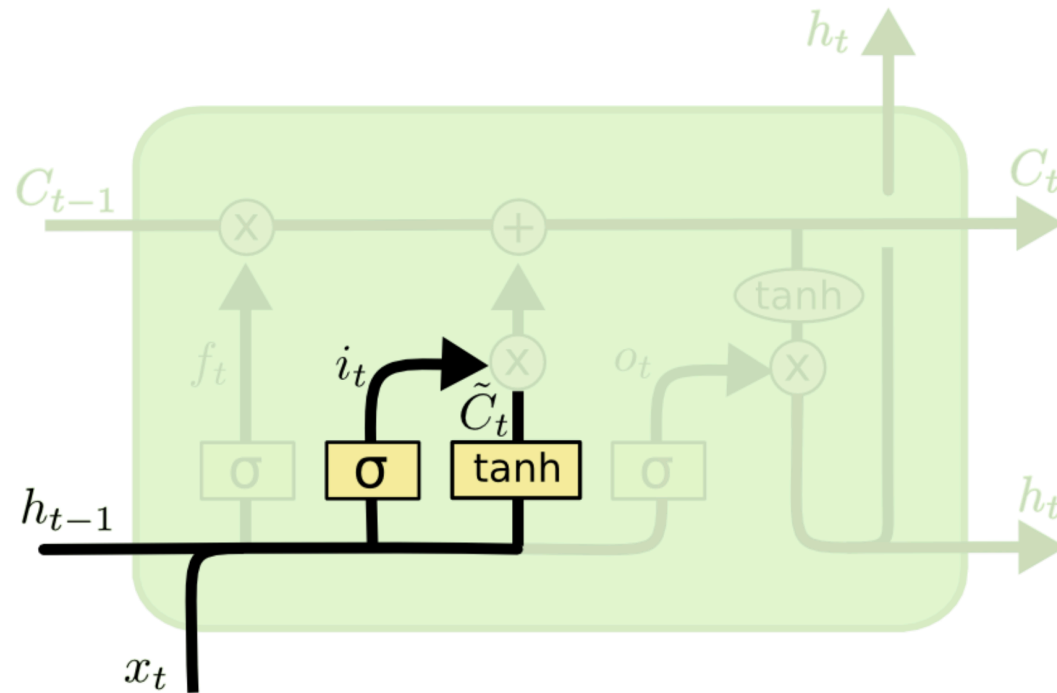  - ‣ input and forget gates are made independent (in place of $z_t$ in GRU)

# Long Short Term Memory (LSTM)

- **Forget gate**: $f_t = \sigma(\text{●}\!\!\text{◗})$



Source: Christopher Olah's blog

$\text{●}\, x_t \quad \text{●}\, h_{t-1} \quad \text{◗}\,$ Linear mixing of $x_t$ and $h_{t-1}$

- **Input gate**: $i_t = \sigma(\text{⬤})$
- **Suggested $C_t$ update**: $\tilde{C}_t = \varphi(\text{⬤})$
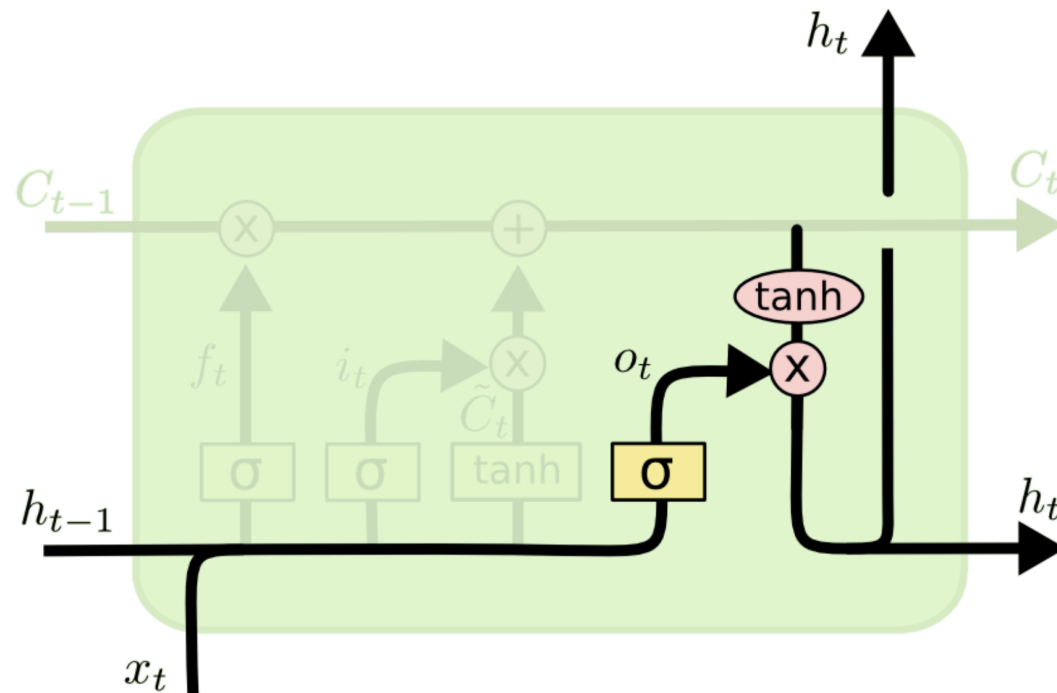


Source: Christopher Olah's blog

⬤ $x_t$ ⬤ $h_{t-1}$ ⬤ Linear mixing of $x_t$ and $h_{t-1}$

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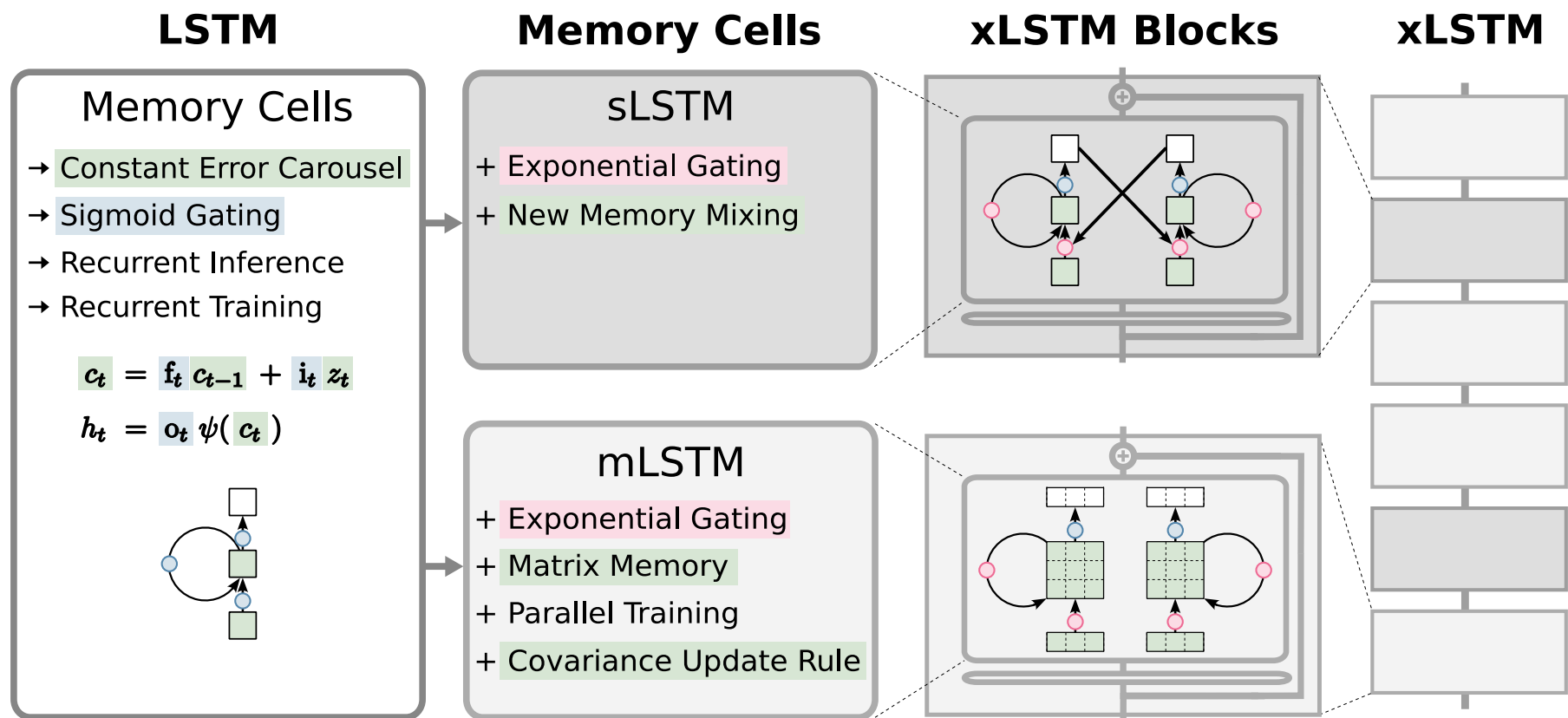


Source: Christopher Olah's blog

# Long Short Term Memory (LSTM)

- **Output gate**: $o_t = \sigma(\bullet)$
- **Hidden state update rule**: $h_t = o_t \odot \varphi(C_t)$



Source: Christopher Olah's blog

$\bullet \, x_t \quad \bullet \, h_{t-1} \quad \bullet$ Linear mixing of $x_t$ and $h_{t-1}$

- A "modern" LSTM variant
  - ▸ Made of sLTSM and mLSTM layers
  - ▸ Embedded in blocks with normalization layers, residual connections, *à la* Transformer



Source: "xLSTM: Extended Long Short-Term Memory" by Beck et al., NeurIPS 2024

- What's "new"?
  - ‣ In both sLSTM and mLSTM layers:
    - – Exponential activation (to face vanishing gradients)
  - ‣ In sLTSM only:
    - – Multi-head
  - ‣ In mLSTM only:
    - – Novel memory store
    - – Drop recurrence for gate computations: better parallelism

- Exponential activation for input and forget gates:

$$i_t = \exp(\bullet)$$

$$f_t = \max(\exp(\bullet), \sigma(\bullet))$$

- Multi-head: keep separate linear mixings per head

$\bullet$ $x_t$  $\bullet$ $h_{t-1}$  $\bullet$ Linear mixing of $x_t$ and $h_{t-1}$

# xLSTM: focus on mLSTM layers

- Exponential activation as in sLSTM
- Memory store

$$C_t = f_t \odot C_{t-1} + i_t \odot v_t k_t^\top$$

$$\tilde{h}_t = C_t q_t \qquad \text{(up to normalization)}$$

▸ Simplified case (no gate):
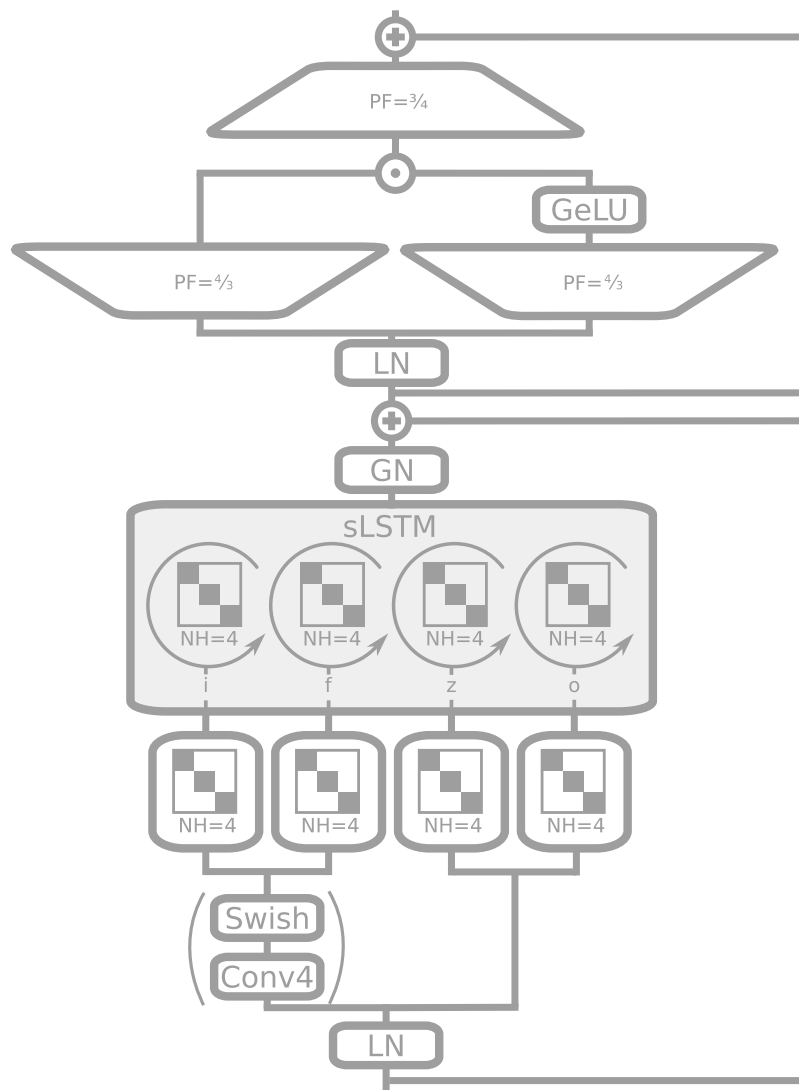
$\Rightarrow$ Mimicks QKV behaviour in self-attention

- Drop recurrence for gate computations: better parallelism
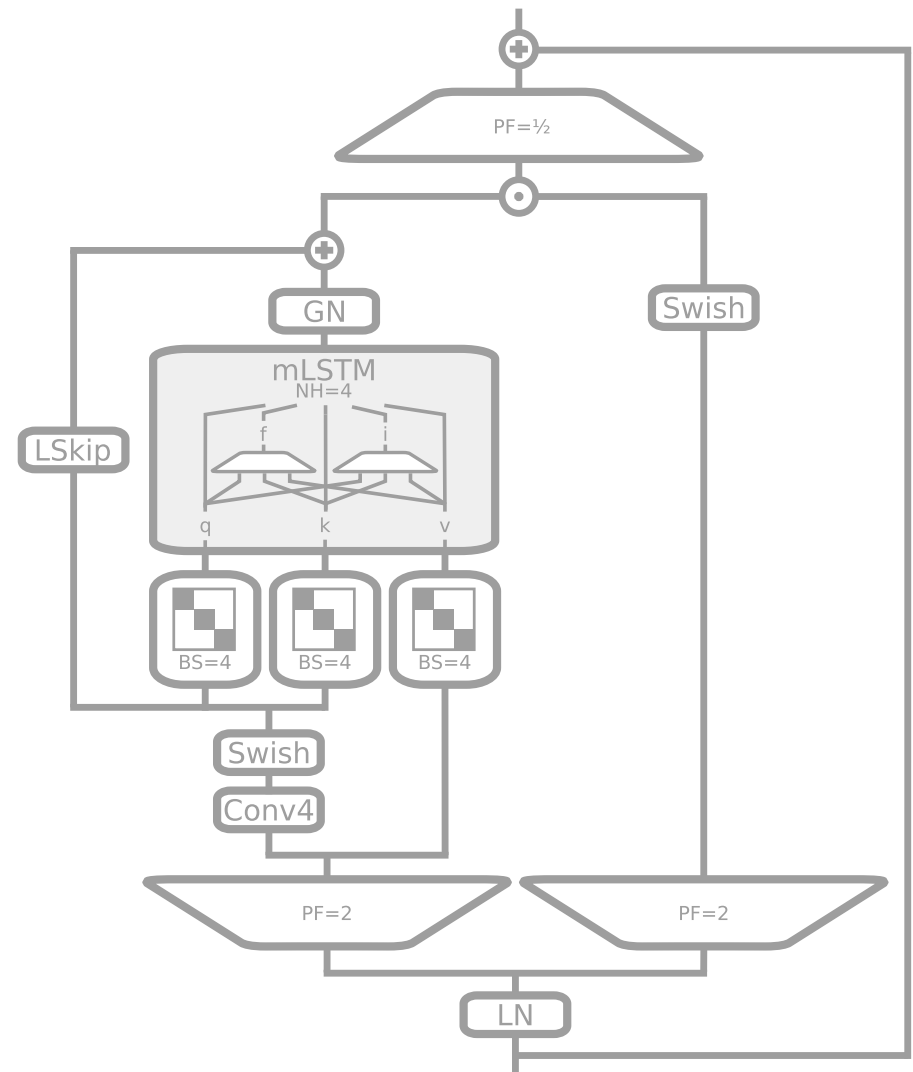
$$i_t = \exp(\bullet)$$

$$f_t = \max(\exp(\bullet), \sigma(\bullet))$$

$$o_t = \sigma(\bullet)$$

$\bullet\ x_t$   $\bullet\ h_{t-1}$   $\bullet$ Linear mixing of $x_t$ and $h_{t-1}$

An sLSTM block

An mLSTM block