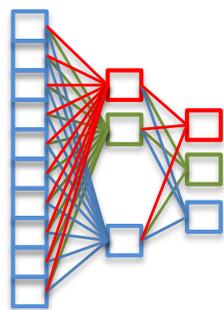


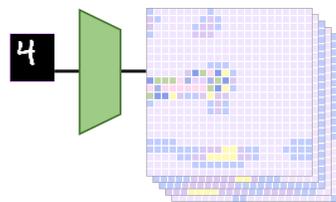
AI and Data Analytics

4.2 Recurrent Neural Networks

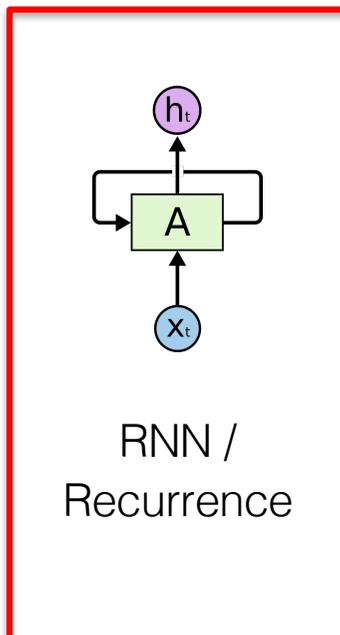
The Deep Toolbox



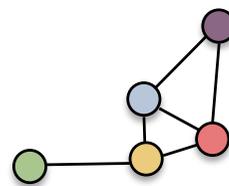
MLP



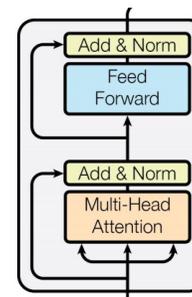
CNN /
Convolutions



RNN /
Recurrence



GN, GCN /
Graphs, geometry



Transformers /
Self-attention

What do I know about the data and the task?

*Nothing
(vector space)*

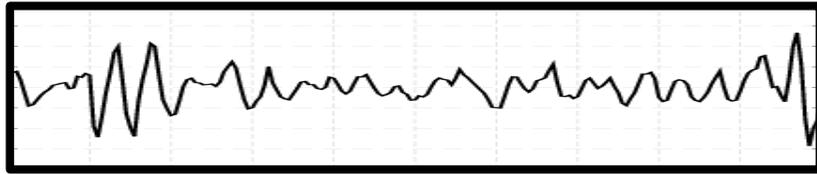
*Translation
equivariance*

*Sequential data,
Markov property*

*Graph structured
data*

*Permutation
equivariance*

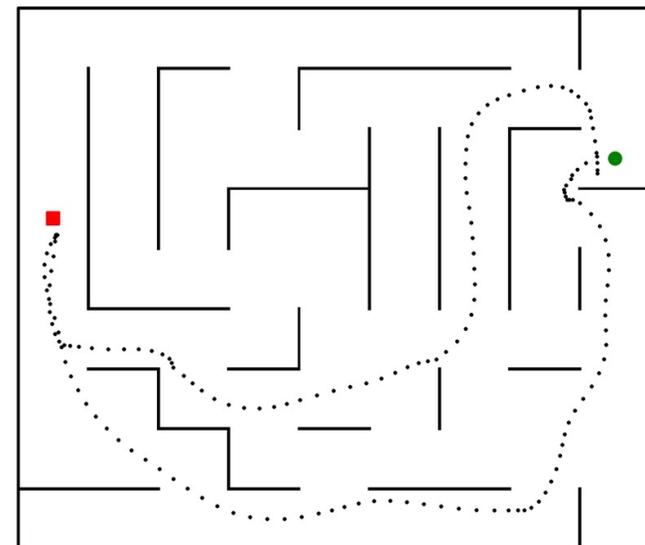
Problem: dealing with sequences



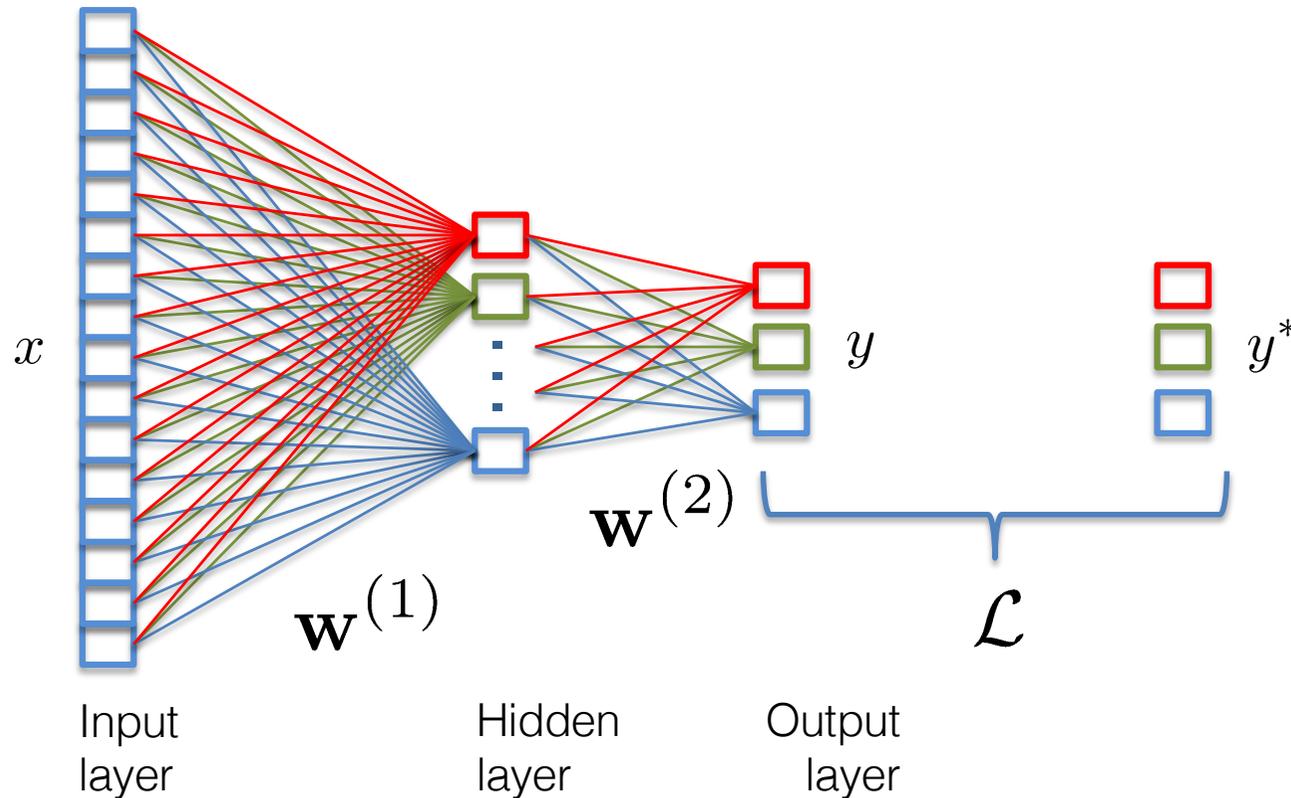
Toutes les familles heureuses se ressemblent. Chaque famille malheureuse, au contraire, l'est à sa façon.



Happy families are all alike. Every unhappy family is unhappy in its own way.

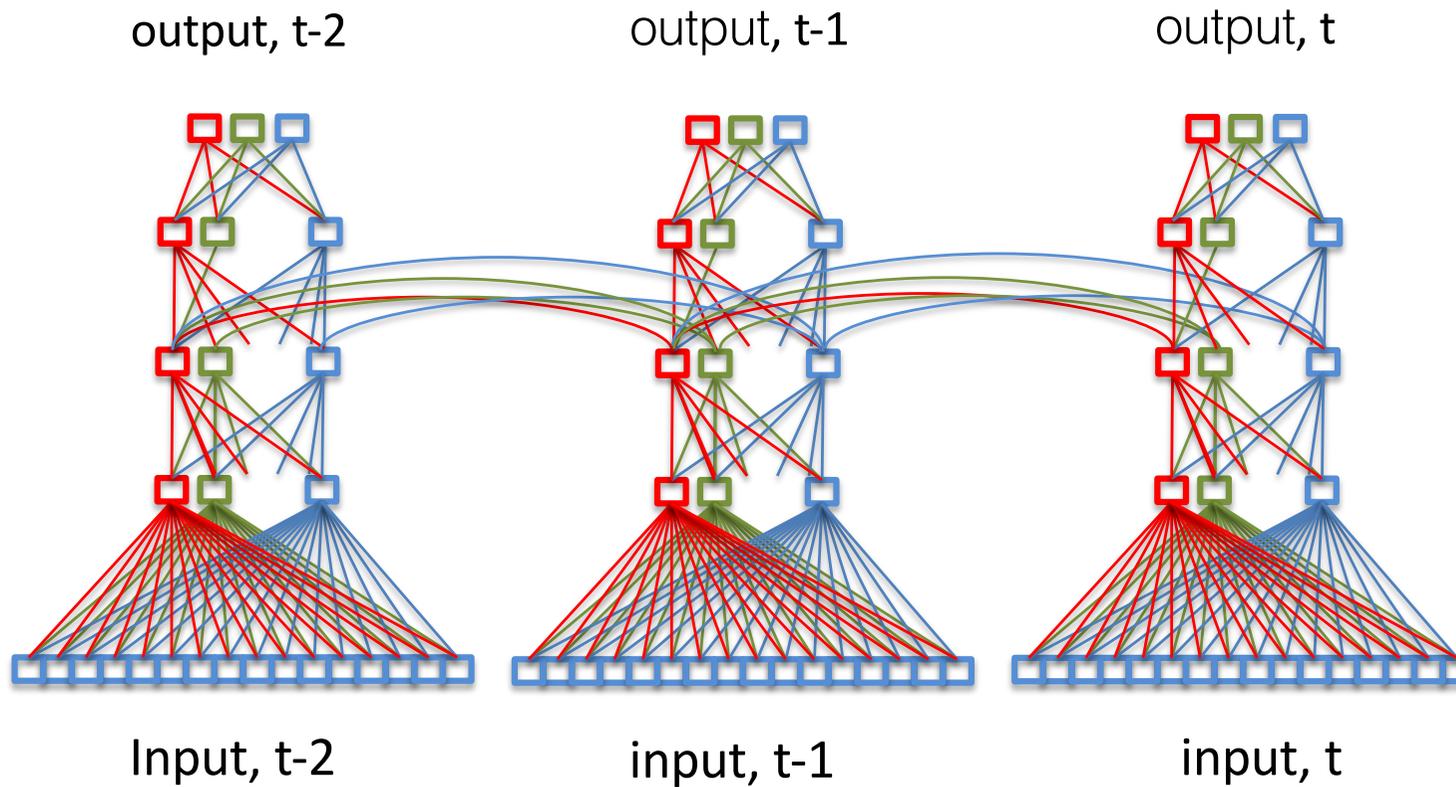


Feed Forward Networks



Prediction: feed-forward computation in a DAG (directed asyclic graph).

Recurrent neural networks



A shout-out ...

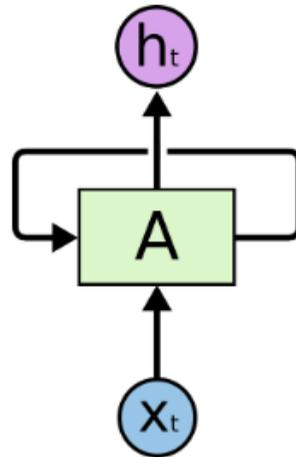
... to Chris Olah's excellent blog post on RNN and LSTM networks, which completely dominates lectures on this topic. The following couple of slides are based on his excellent drawings.

Chris Olah, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Recurrent Neural Networks (RNNs)

As feed-forward networks, Recurrent Neural Networks (RNNs) predict some output from a given input.

However, they also pass information over time, from instant $(t-1)$ to (t) :

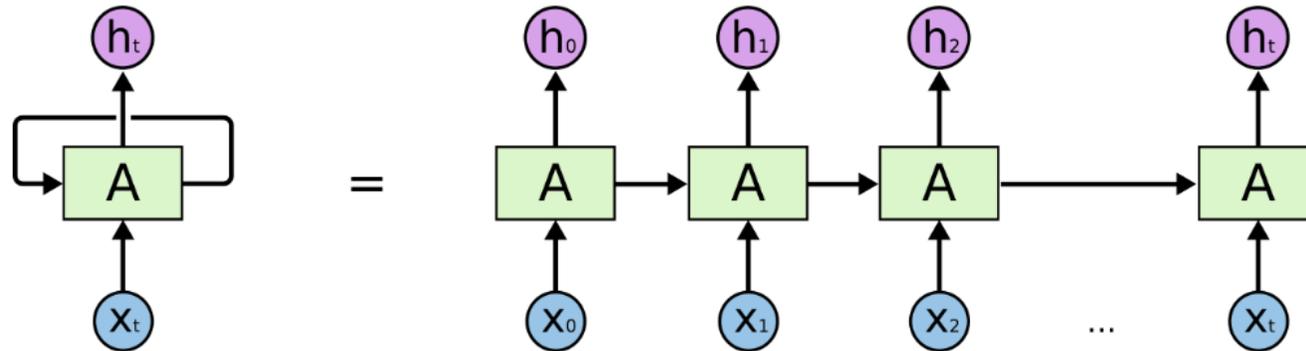


Here, we write h_t for the output, since these networks can be stacked into multiple layers, i.e. h_t is input into a new layer.

Chris Olah, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Recurrent Neural Networks (RNNs)

A more intuitive view of an RNN is to unroll it over time:



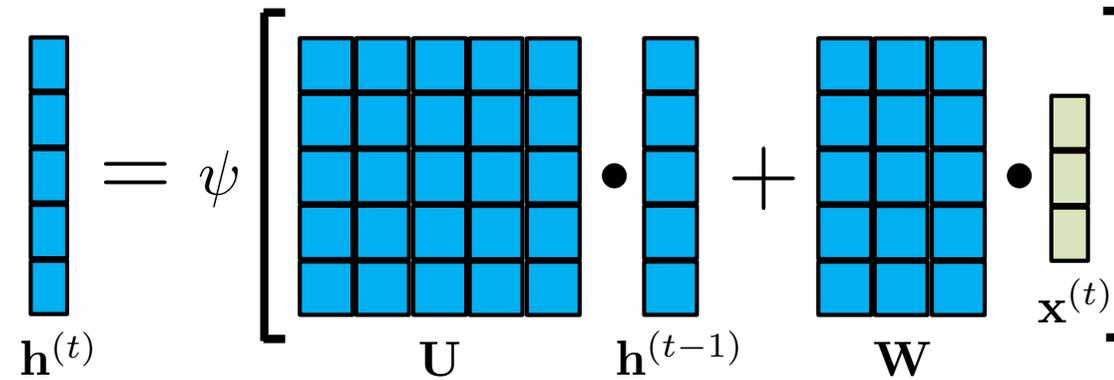
The update equations are:

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

\Rightarrow two weight matrices: one for the recurrent connections over time, one for the input connections.

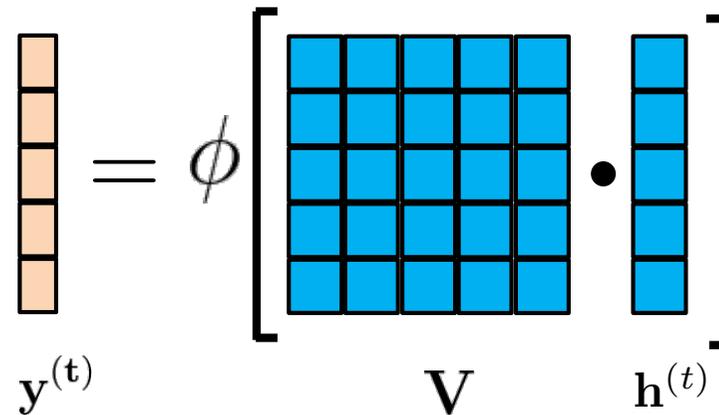
σ is an activation function.

From state h to output y



The diagram illustrates the state transition equation. On the left, a vertical blue vector $h^{(t)}$ is shown. To its right is an equals sign followed by the Greek letter ψ . This is followed by a large square bracket containing two terms. The first term is a 5x5 blue grid labeled U multiplied by a vertical blue vector $h^{(t-1)}$. The second term is a 5x3 blue grid labeled W multiplied by a vertical light green vector $x^{(t)}$. The two terms are separated by a plus sign. The entire expression is enclosed in a large square bracket.

$$h^{(t)} = \psi \left[U \cdot h^{(t-1)} + W \cdot x^{(t)} \right]$$



The diagram illustrates the output equation. On the left, a vertical orange vector $y^{(t)}$ is shown. To its right is an equals sign followed by the Greek letter ϕ . This is followed by a large square bracket containing a 5x5 blue grid labeled V multiplied by a vertical blue vector $h^{(t)}$. The entire expression is enclosed in a large square bracket.

$$y^{(t)} = \phi \left[V \cdot h^{(t)} \right]$$

Toy example with handcrafted parameters

In a sequential problem, we surveil a farm and watch for the appearance of objects. At each instant t we observe a vector \mathbf{v} which can indicate the appearance of

- a wolf 
- a farmer 

The objective is to output an estimate of danger, i.e.

- presence of the wolf w/o the farmer, or
- presence of both, arrival of the wolf before the farmer.

Toy example with handcrafted parameters

$$\mathbf{h}^{(t)} = \psi \left[\mathbf{U} \cdot \mathbf{h}^{(t-1)} + \dots \right]$$

\mathbf{U} matrix:

1	0	0	0
0	1	0	0
1	0	1	0
0	1	0	1

$\mathbf{h}^{(t-1)}$ components:

- Wolf arrived *
- Farmer arrived *
- Time since wolf arrived
- Time since farmer arrived

$$\left[\begin{array}{l} \text{Appearance template "wolf"} \\ \text{Appearance template "farmer"} \end{array} + \mathbf{W} \cdot \mathbf{x}^{(t)} \right]$$

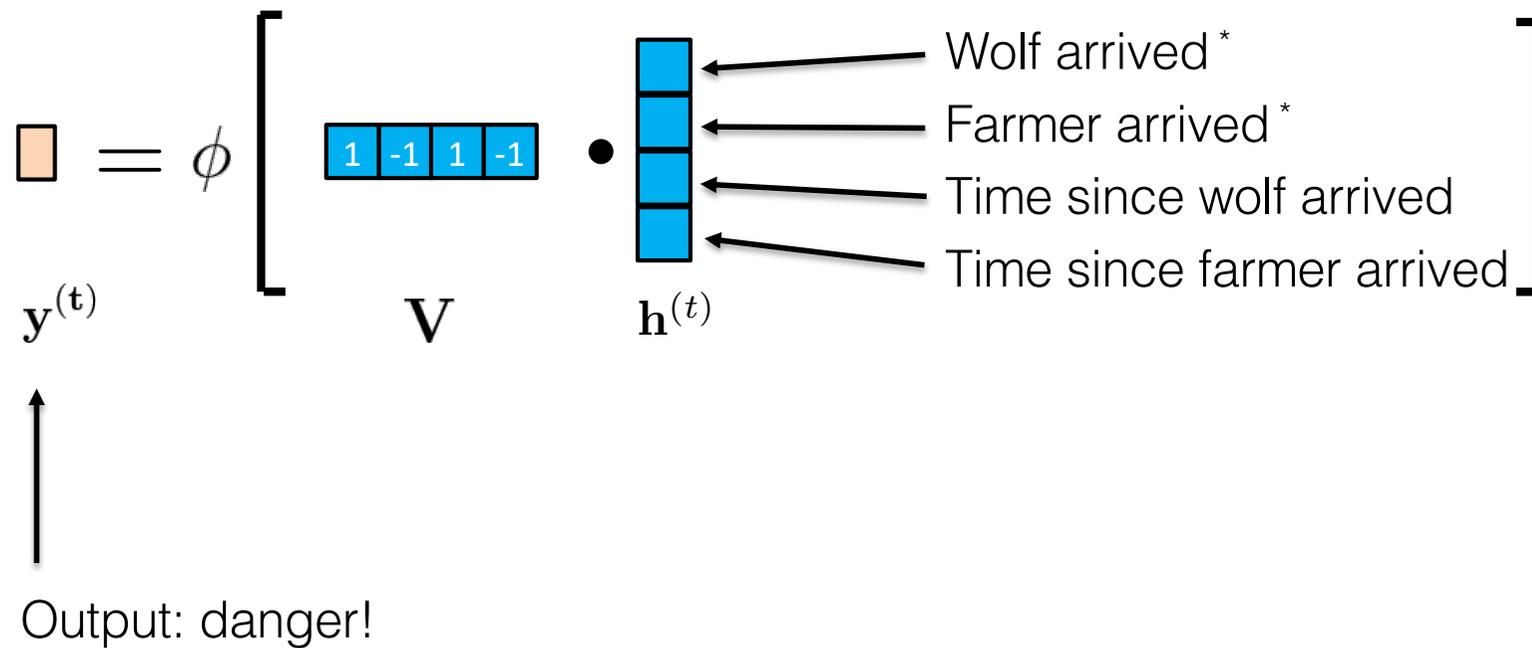
\mathbf{W} matrix:

0	0	0
0	0	0

$\mathbf{x}^{(t)}$ vector:

*Appropriate activation functions required to normalize state values

Toy example with handcrafted parameters



RNN training & problems

RNNs are trained with backpropagation through time (BPTT): the graph is unrolled, and the loss derived w.r.t. the parameters of all different time instants.

Standard vanilla RNNs are difficult to train and suffer from shortcomings:

- Vanishing gradients: small gradients vanish over long time ranges.
- Exploding gradients: high gradients explode over long ranges.
- Lack of long-term dependency handling: short-term updates between individual time instances dominate.

Solution: gating mechanisms (LSTM and GRU networks).

Gating

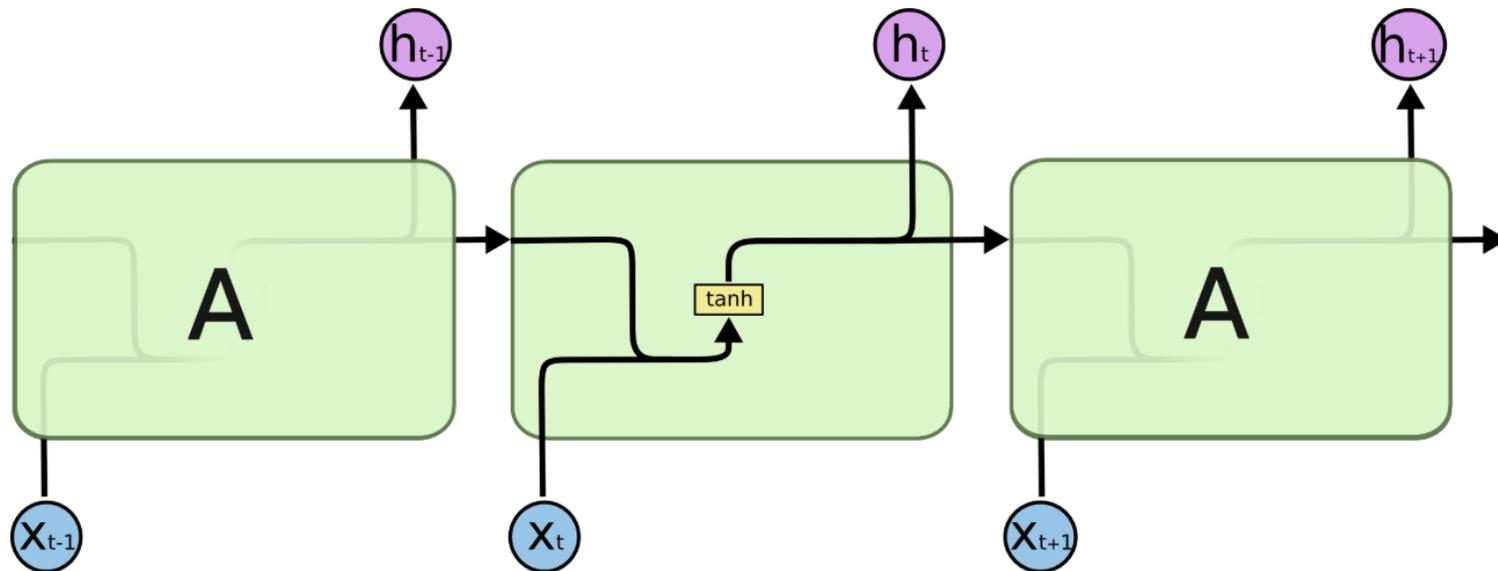
Problem: hidden state is updated at each time step

Solution: at each instant, for each state value, decide how much to

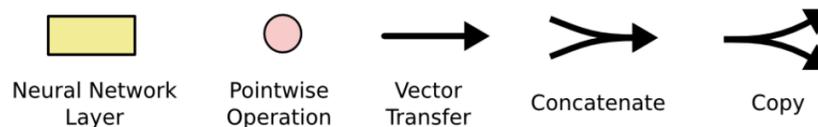
- input
 - forget
 - output
- These decisions are also learned

LSTM Networks: walk through

We start by illustrating a vanilla RNN in a new way:

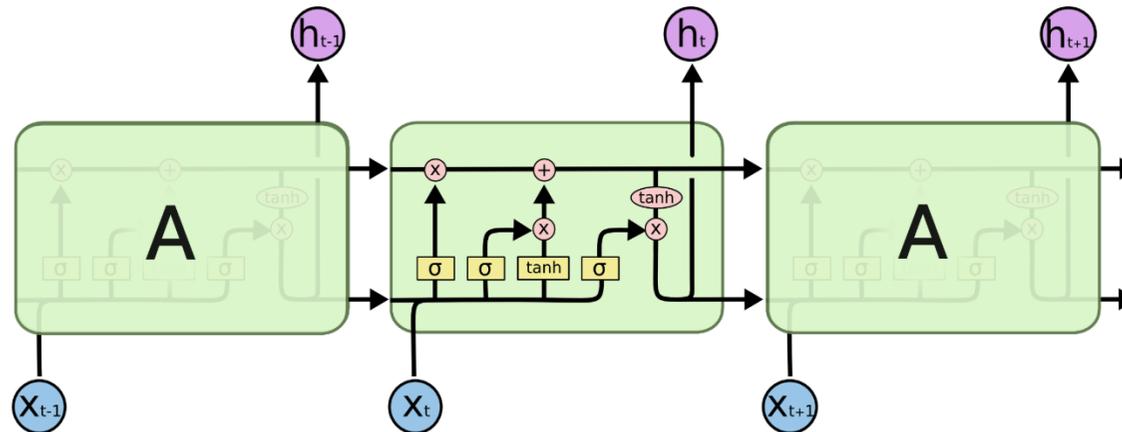


where:



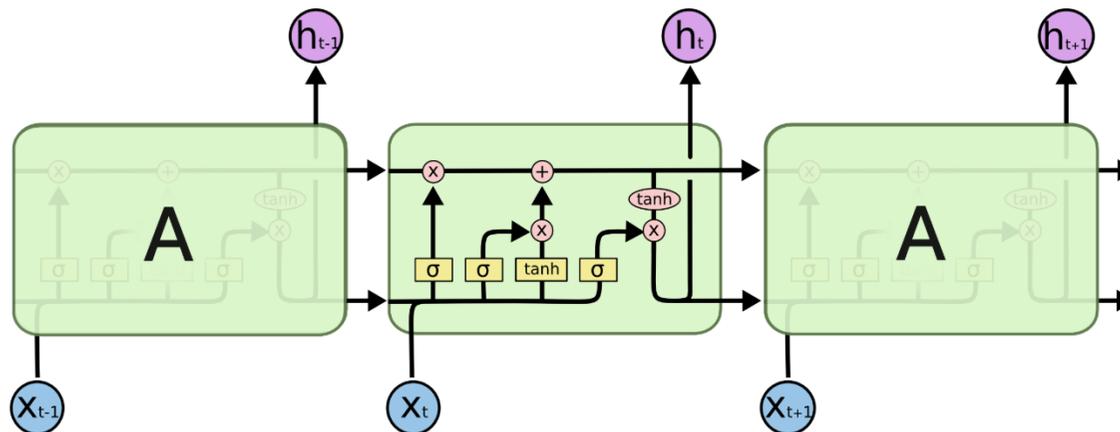
LSTM Networks

LSTM (=Long-short term Memory) networks use gating mechanisms, which handle information flow in a fully trainable way.



LSTM Networks

LSTM (=Long-short term Memory) networks use gating mechanisms, which handle information flow in a fully trainable way.



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

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

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

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

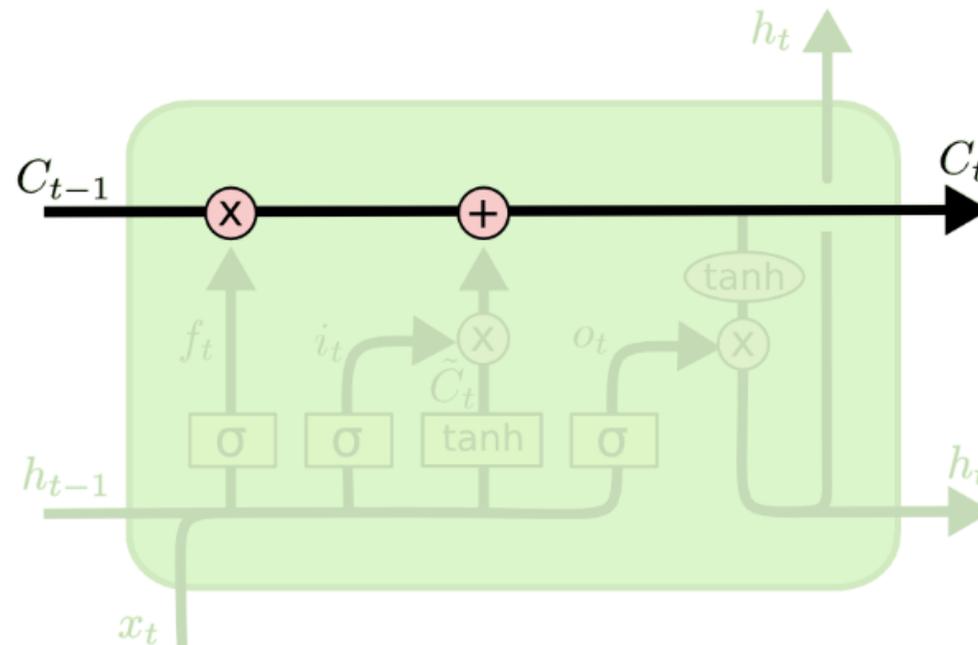
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

LSTM Networks: walk through

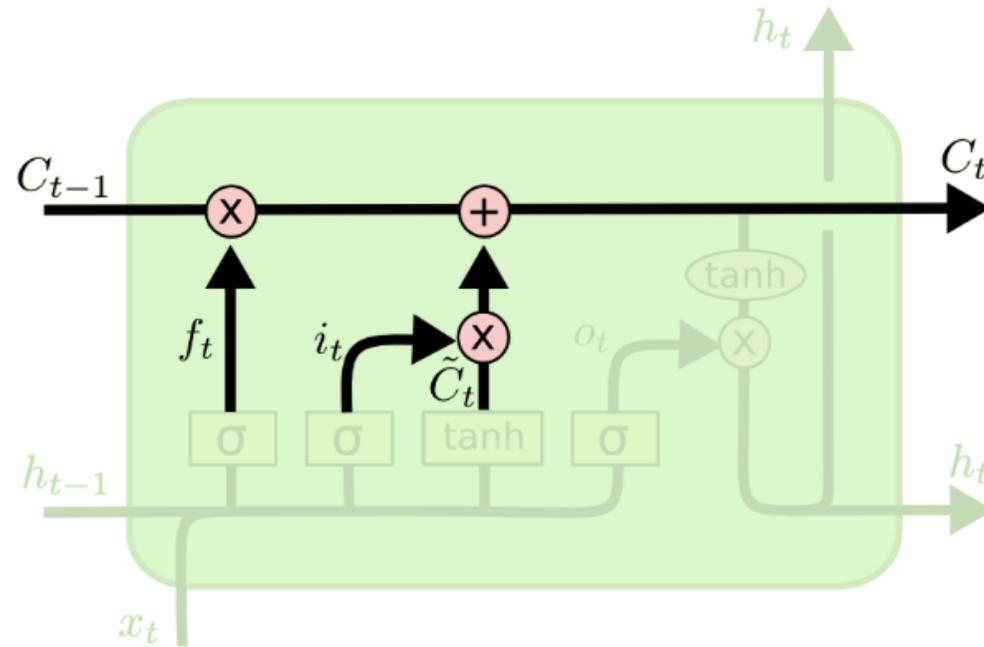
An LSTM has two different memory representations:

- A cell state C ,
- the hidden state h .



LSTM Networks: walk through

The new cell state C_t is a linear combination of the old cell state C_{t-1} and some new updated information \tilde{C}_t (described later):

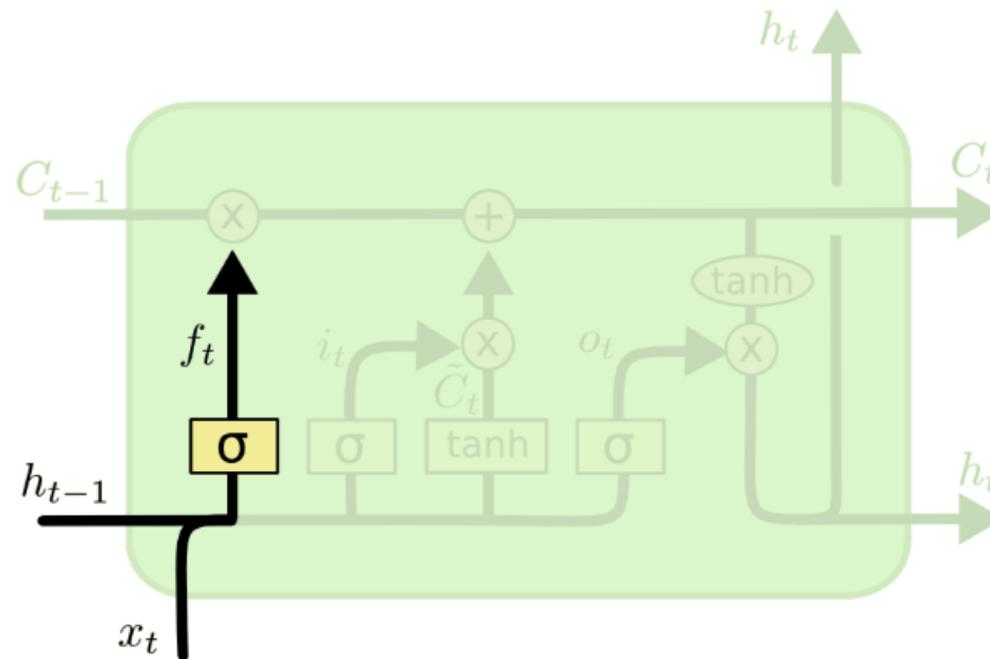


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

f_t and i_t are “**gates**” (trainable functions), which govern the information flow (based on the hidden state h_t).

LSTM Networks: walk through

The forget gate controls how much of the old cell state is forgotten or passed through:

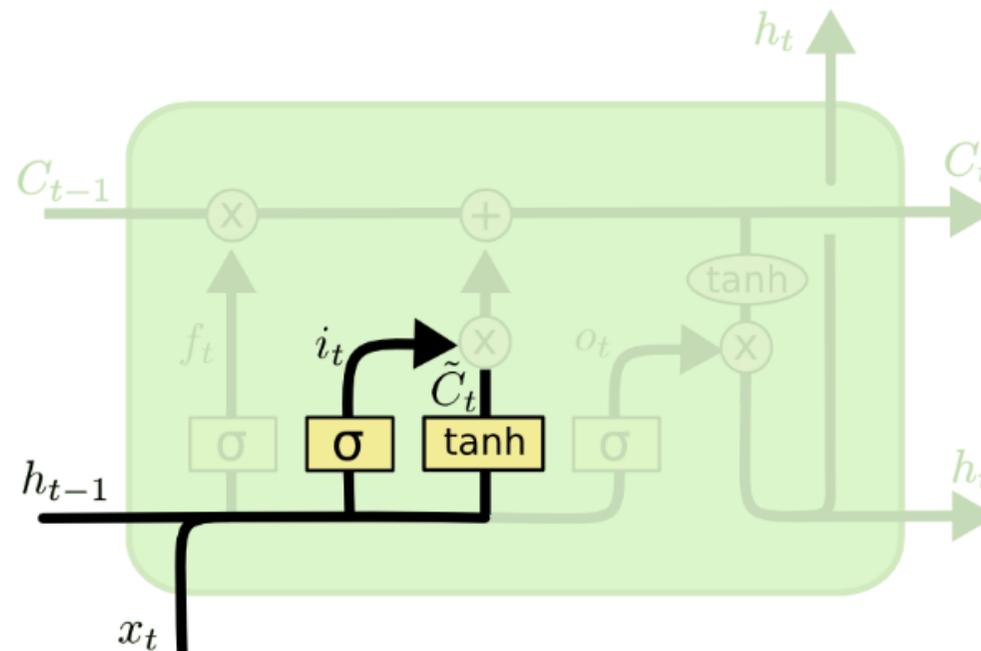


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTM Networks: walk through

The input gate controls how much new information is passed over to the cell state.

The new information is predicted from the hidden state h :

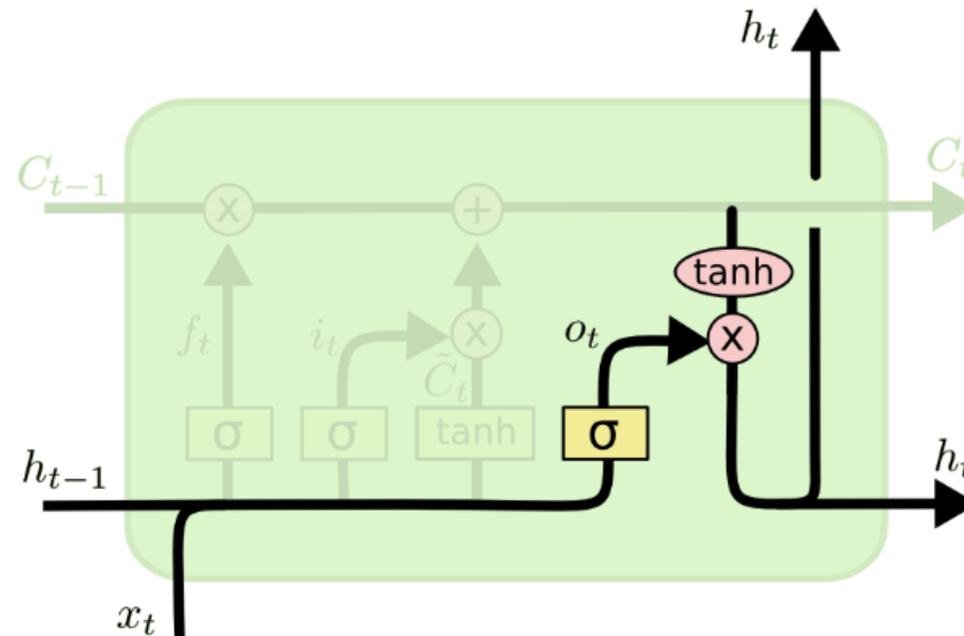


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

LSTM Networks: walk through

The output gate controls, how the information from the cell state is translated into the hidden state h :

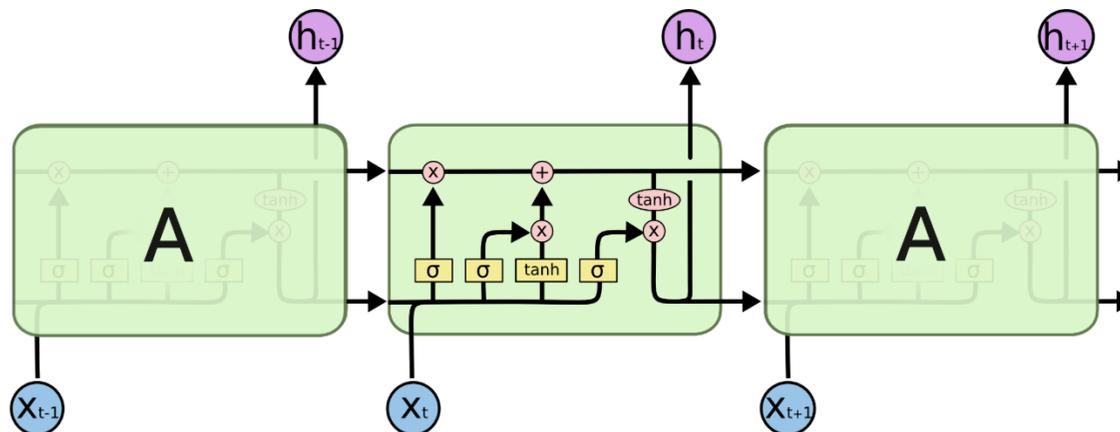


$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

LSTM Networks

The full model, again:



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

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

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

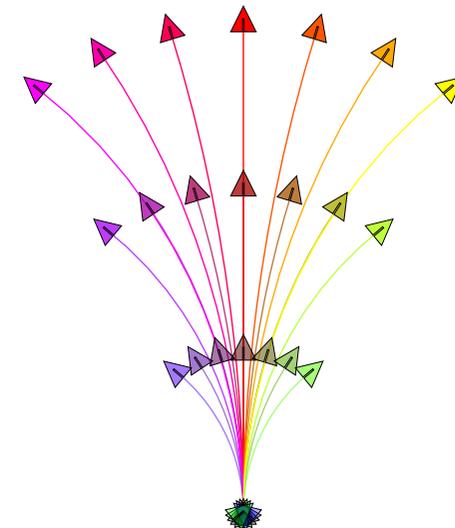
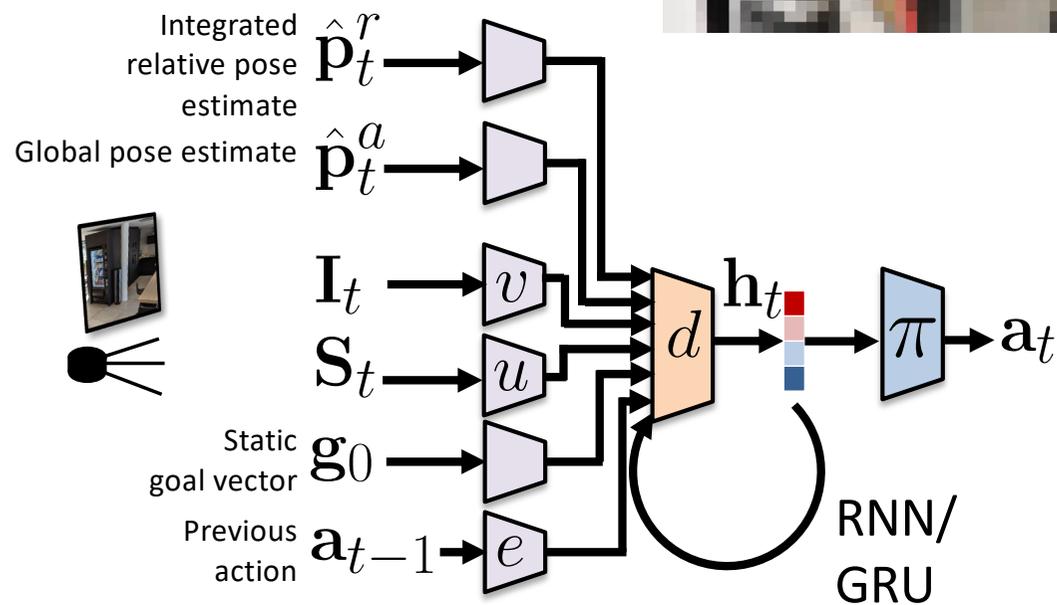
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

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

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

Example applications: Robotics

Navigation with recurrent memory



Input: RGB, Lidar, pose estimates from onboard sensors.
The agent does **not** have access to a map.

28 discrete actions
Pairs of linear+angular velocities