

Introduction to recurrent neural networks

Deep Learning problems with sequential data

- Speech recognition
- Sentiment classification
- Machine translation
- Video activity recognition

Input



“The quick brown fox jumped over
the lazy dog.”

“There is nothing to like in
this movie.”



Voulez-vous chanter
avec moi?



Do you want to sing with me?



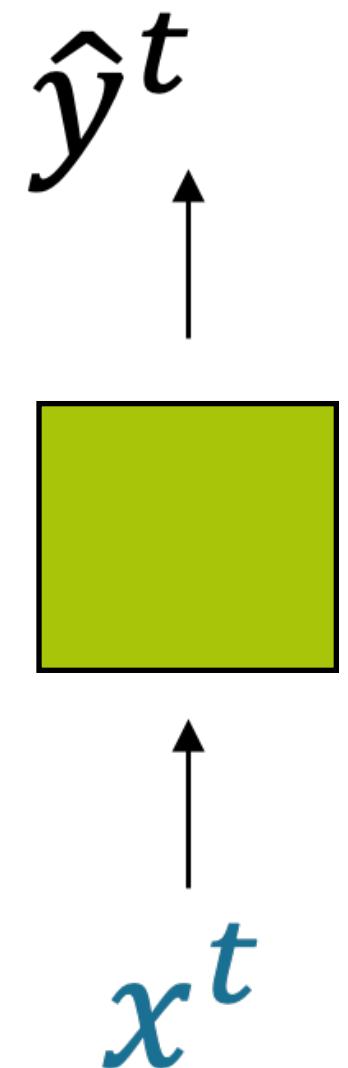
Running

Idea 1: use the previous word

Глаза Степана Аркадьича весело заблестели, и он задумался, улыбаясь

- Features - one hot encoding for a dictionary of all words:
- The number of "features" - hundreds of thousands

[0, 0, 0, 1, 0, ..., 0]



Issue 1: Long-Term Dependencies

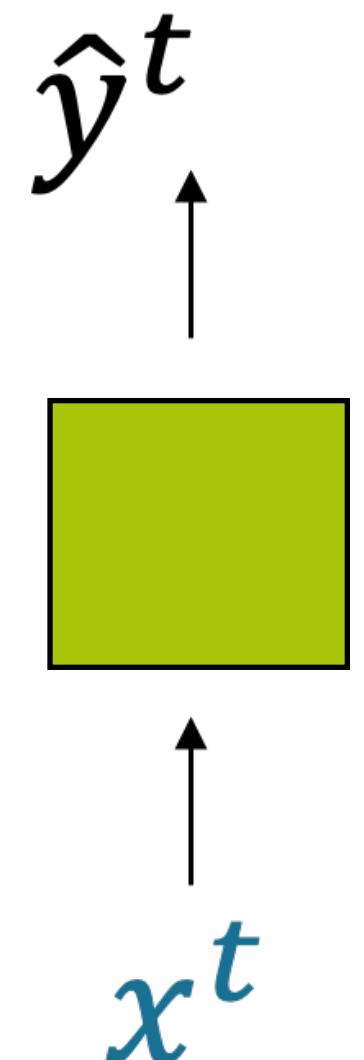
“France is where I grew up, but I now live in Boston. I speak fluent ____.”

Idea 2: bag of words model

Глаза Степана Аркадьича весело заблестели, и он задумался, улыбаясь

- Features - counting the number of words for a dictionary of all words:
- ***Bag of Words Model*** - Order Doesn't Matter

[0, 1, 0, 1, 0, ..., 1]



Issue 2: Use of word order information

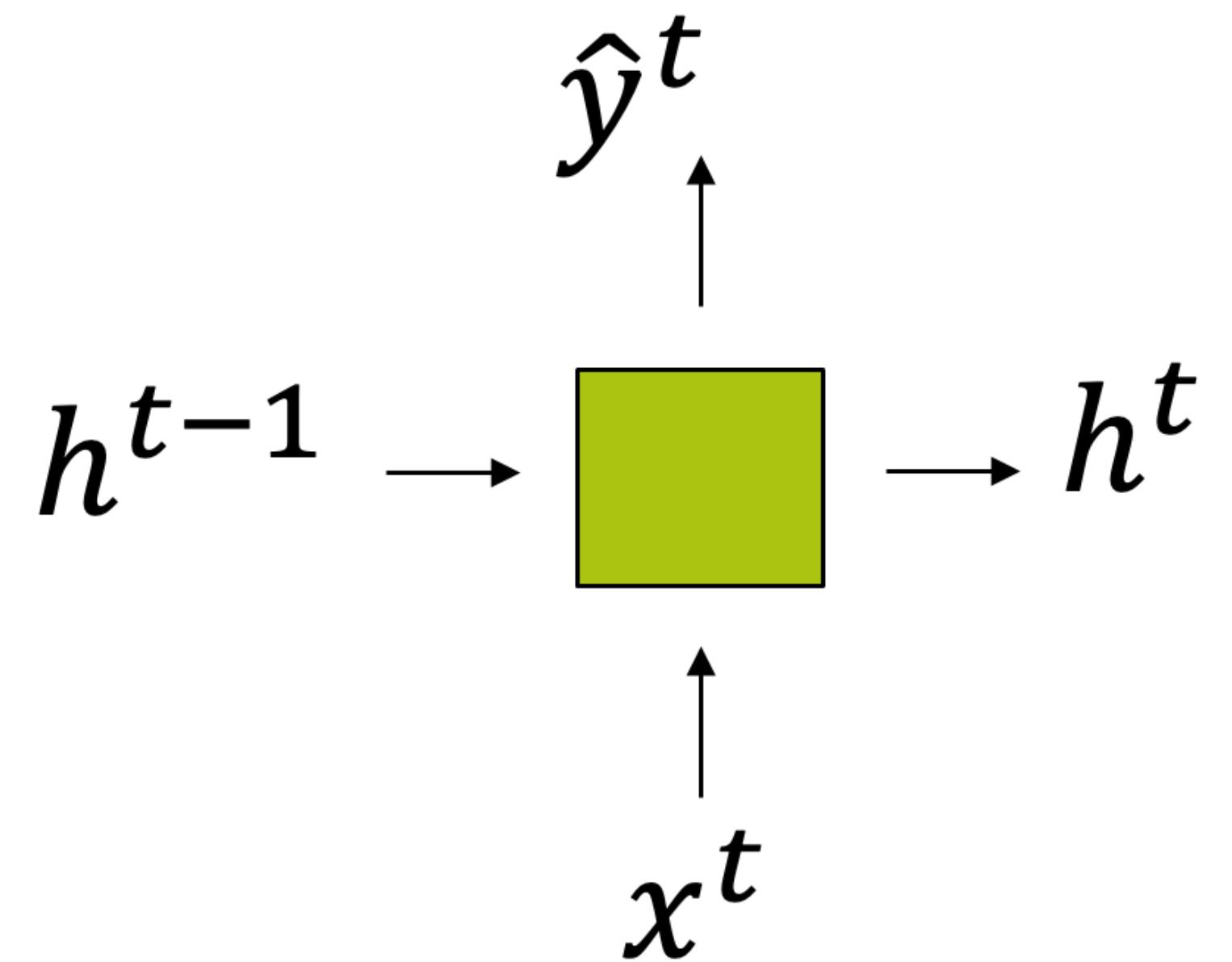
The food was good, not bad at all.

vs.

The food was bad, not good at all.

Model Design Criteria

- Long term memory
- Order Information
- Natural preprocessing
- Handling Variable Length Sequences



Solution: Recurrent neural networks (RNNs)

$$\mathbf{h}^t = f_h(\mathbf{x}^t, \mathbf{h}^{t-1})$$

$$\mathbf{h}^t = \tanh(V\mathbf{x}^t + W\mathbf{h}^{t-1} + b_h)$$

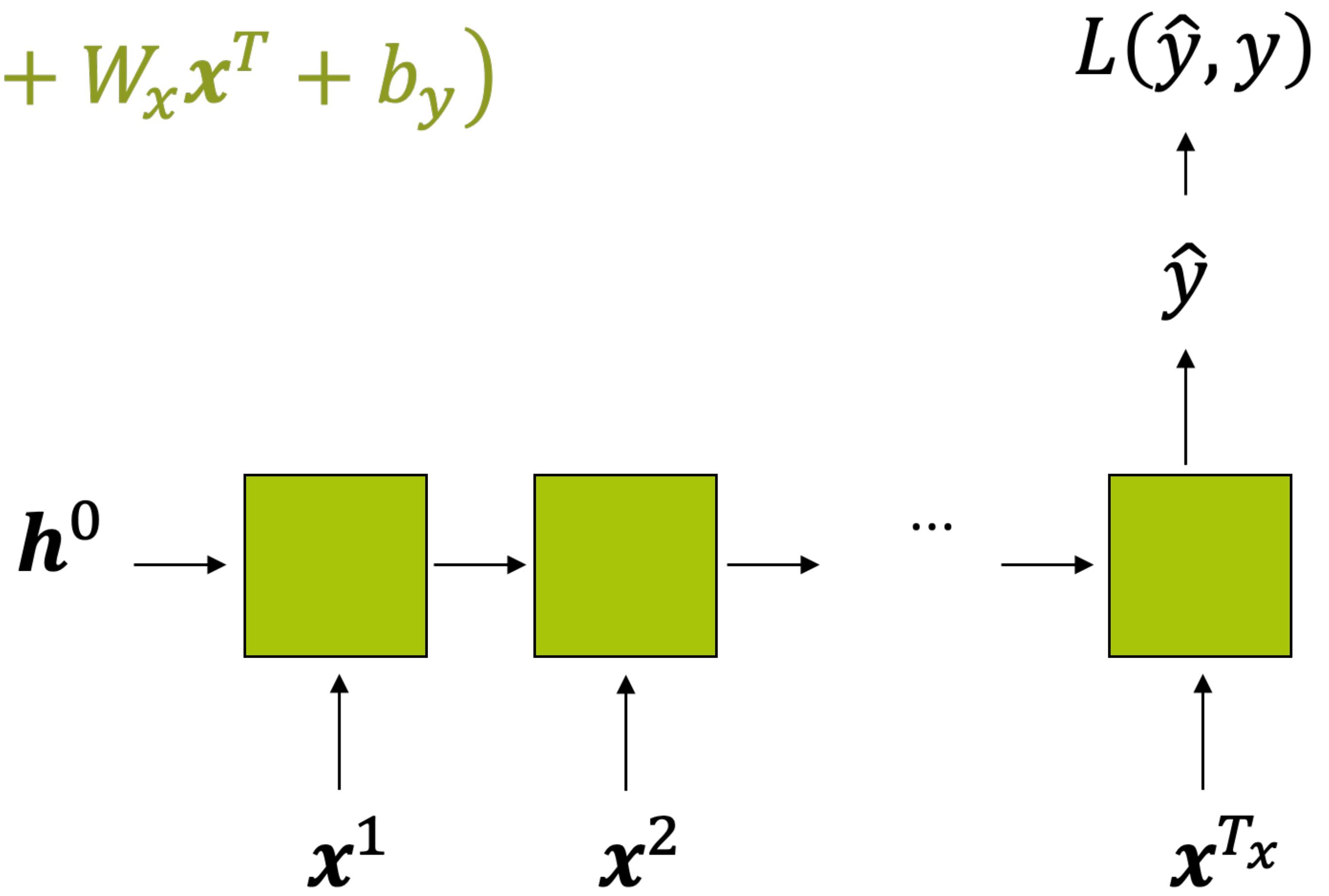
$$\hat{y}^t = f_y(\mathbf{h}^t)$$

$$\hat{y}^t = \text{softmax}(U\mathbf{h}^t + b_y)$$

Model learning: back propagation

$$\hat{y}^t = \text{softmax}(W_h h^t + W_x x^T + b_y)$$

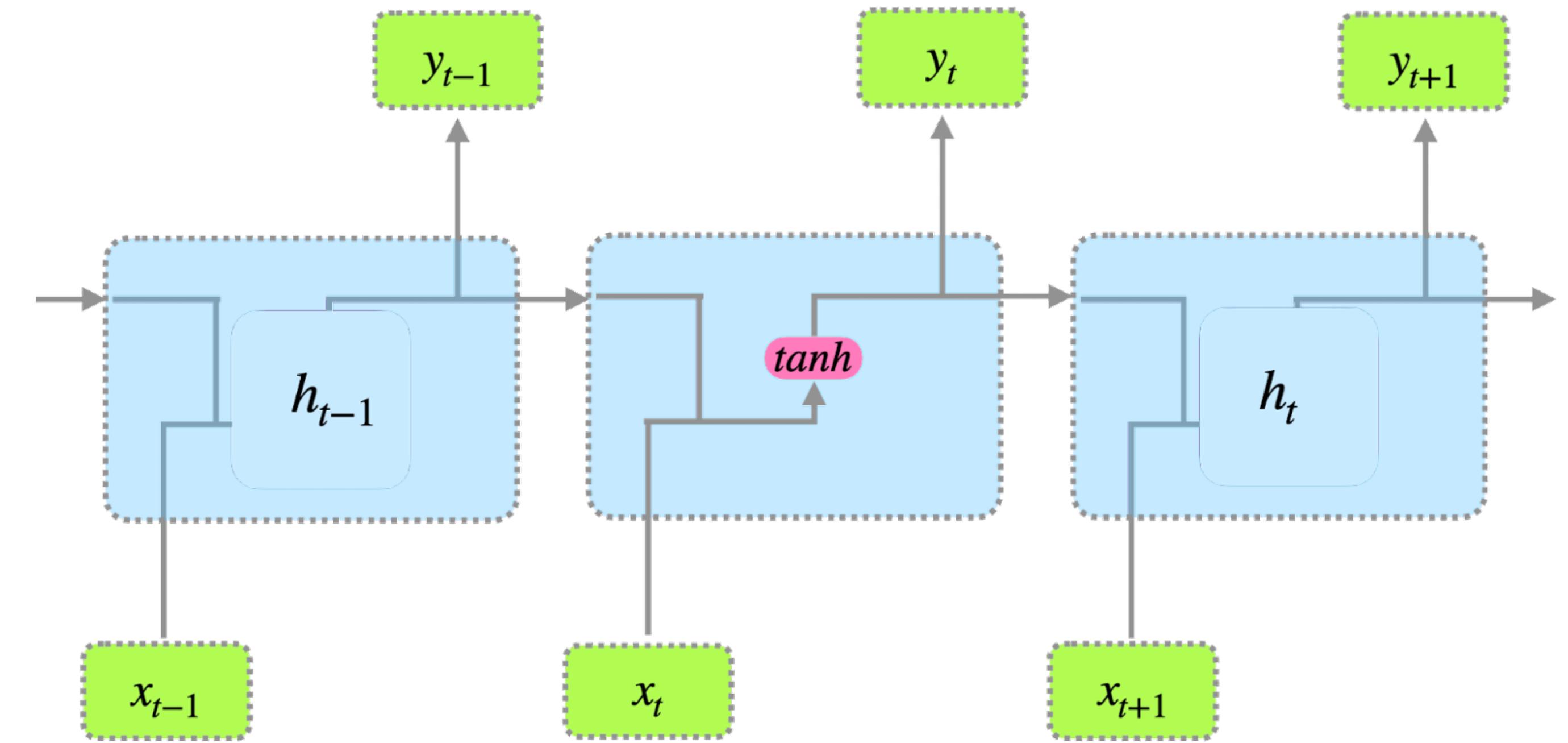
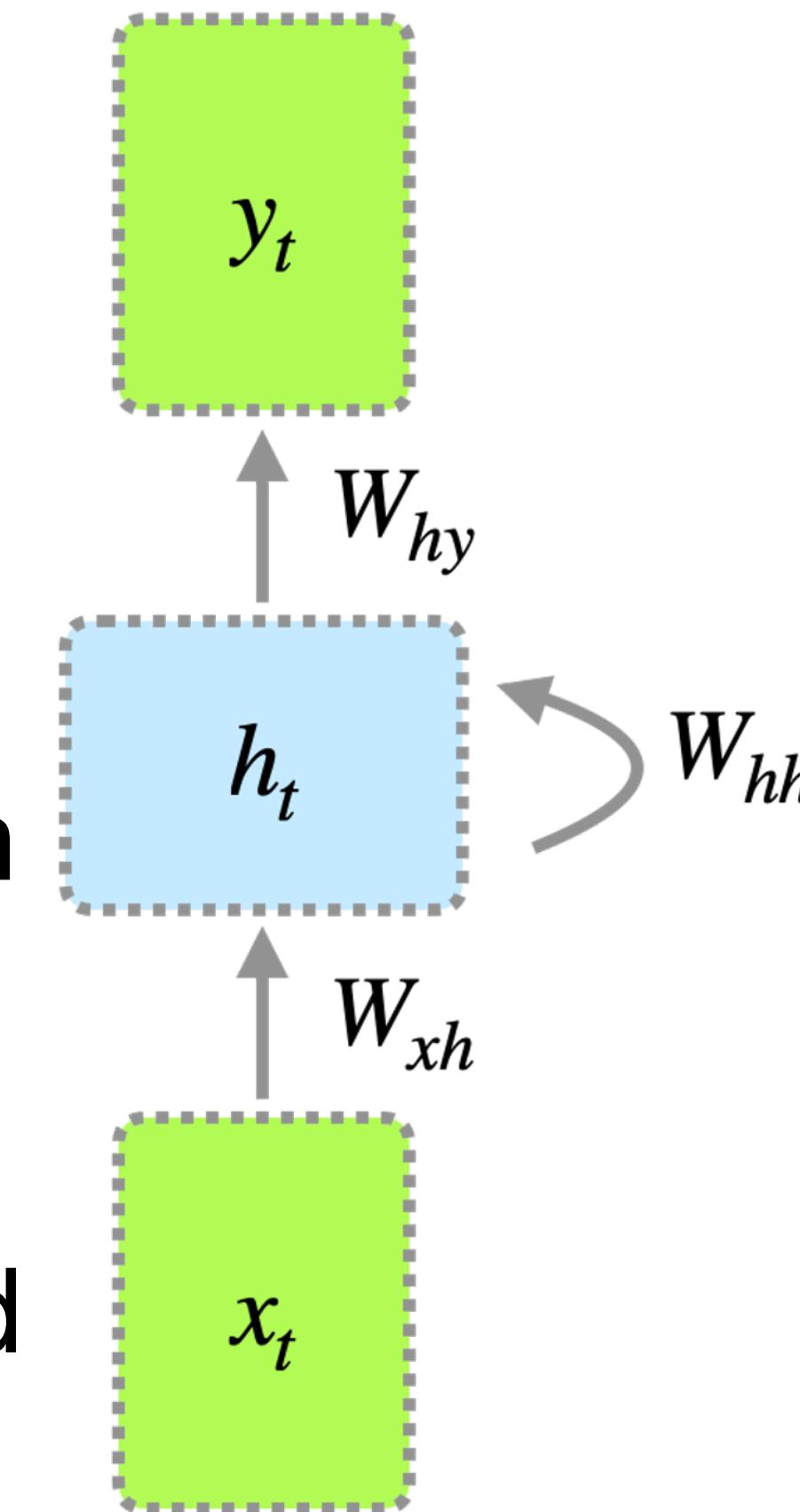
$$\frac{\partial L}{\partial W_h} = \frac{\partial L}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial W_h}$$



RNN intuition

Model Design

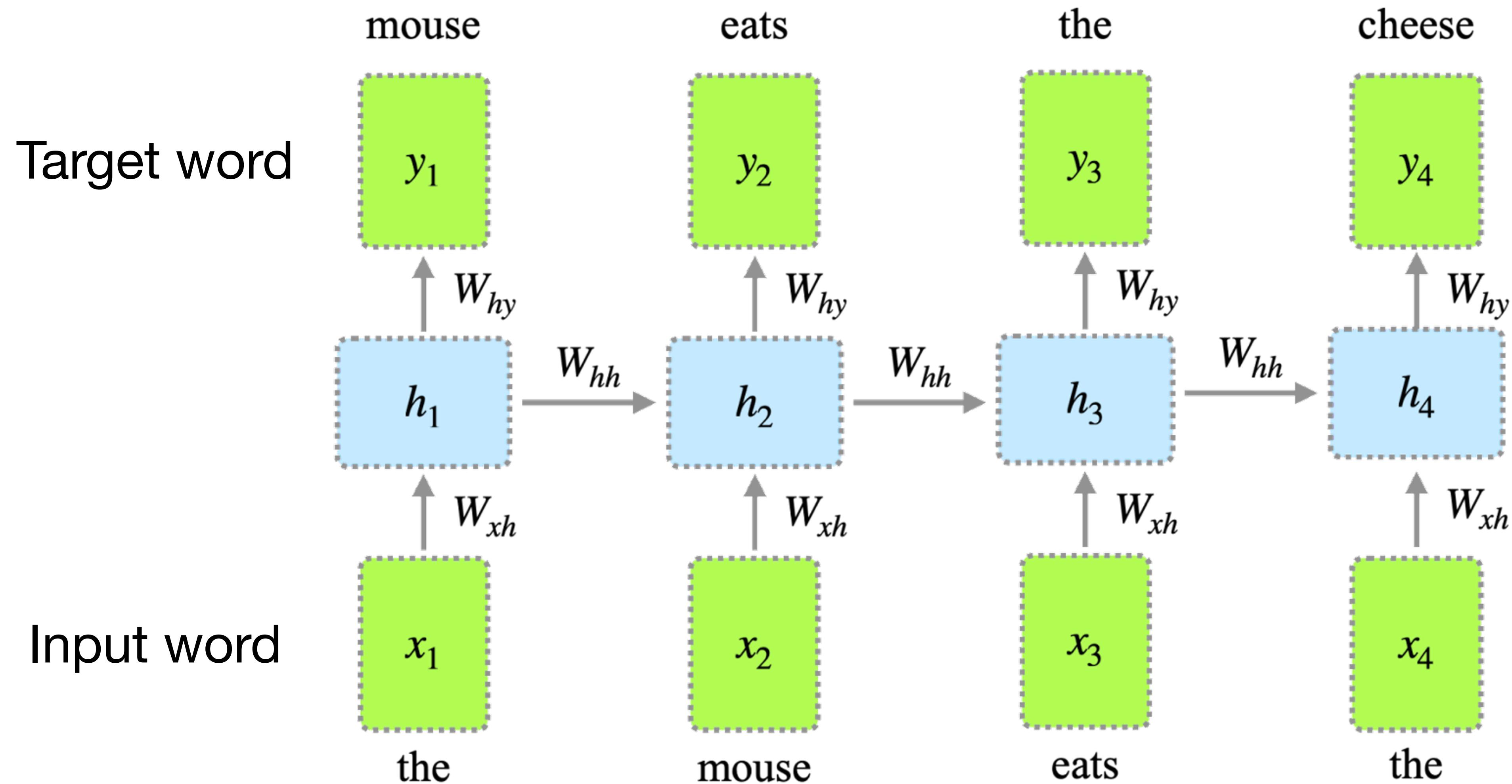
Predicted likelihood
Hidden representation
Encoded word



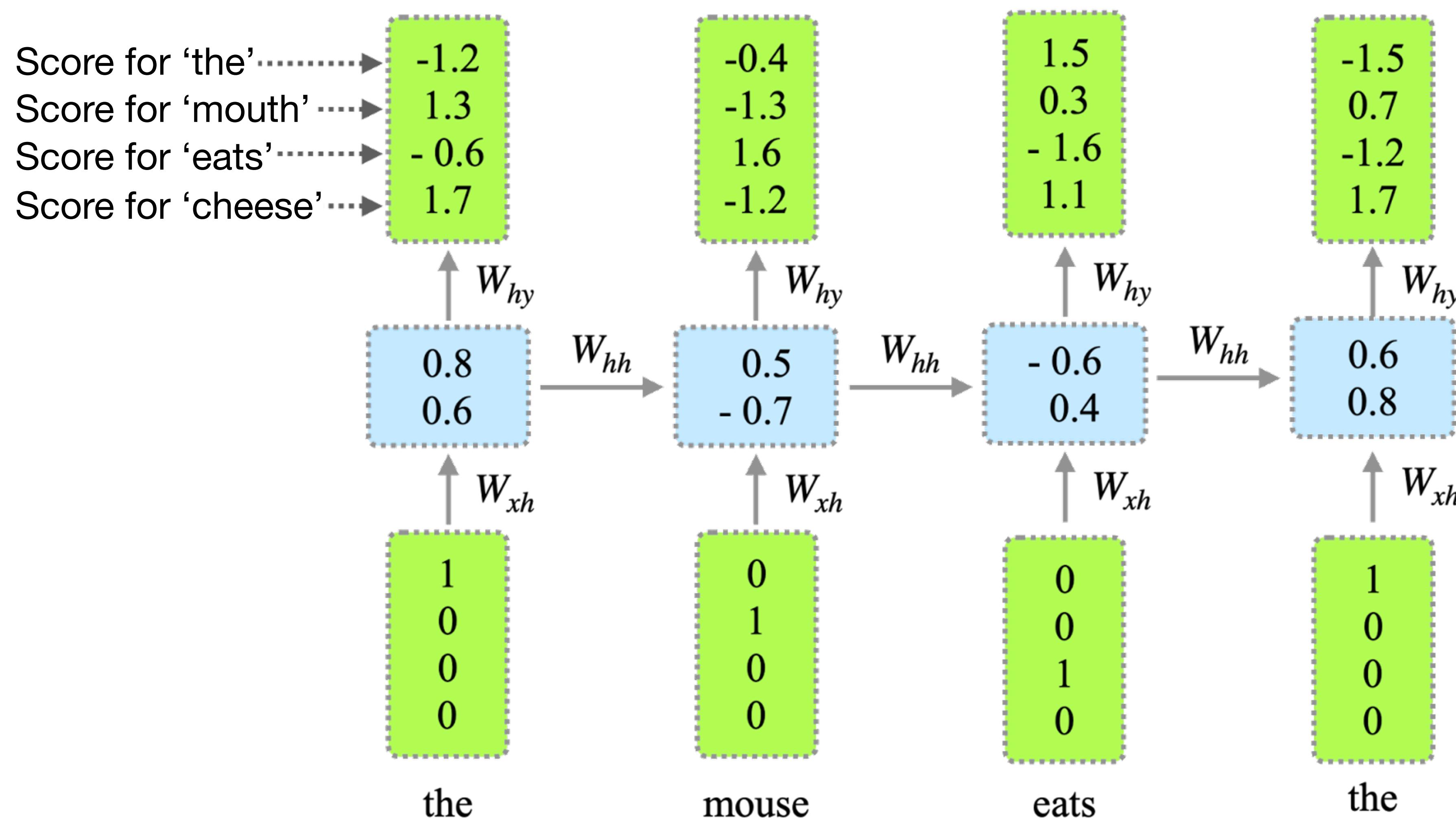
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

x_t and y_t — n-dimensional for dictionary of size t

Model Design Example

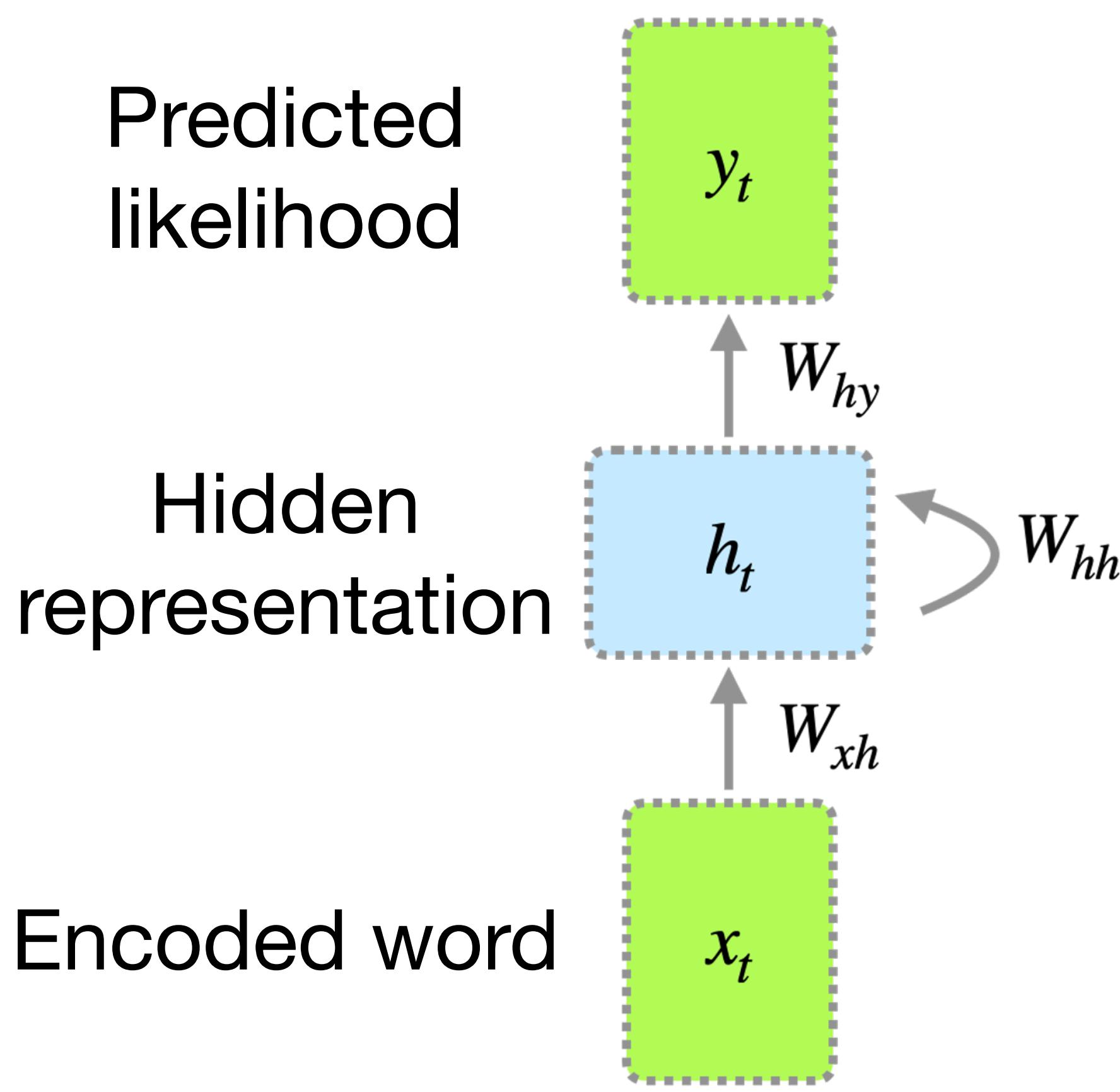


Model Design Example

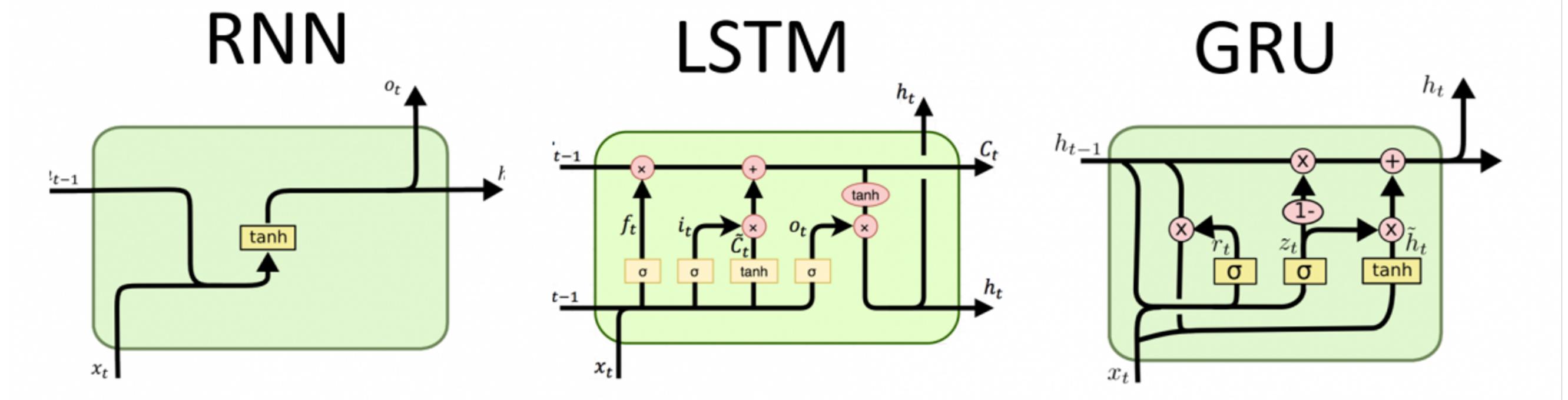


LSTM and GRU

Proper configuration of the recurrent block ensures long-term dependencies



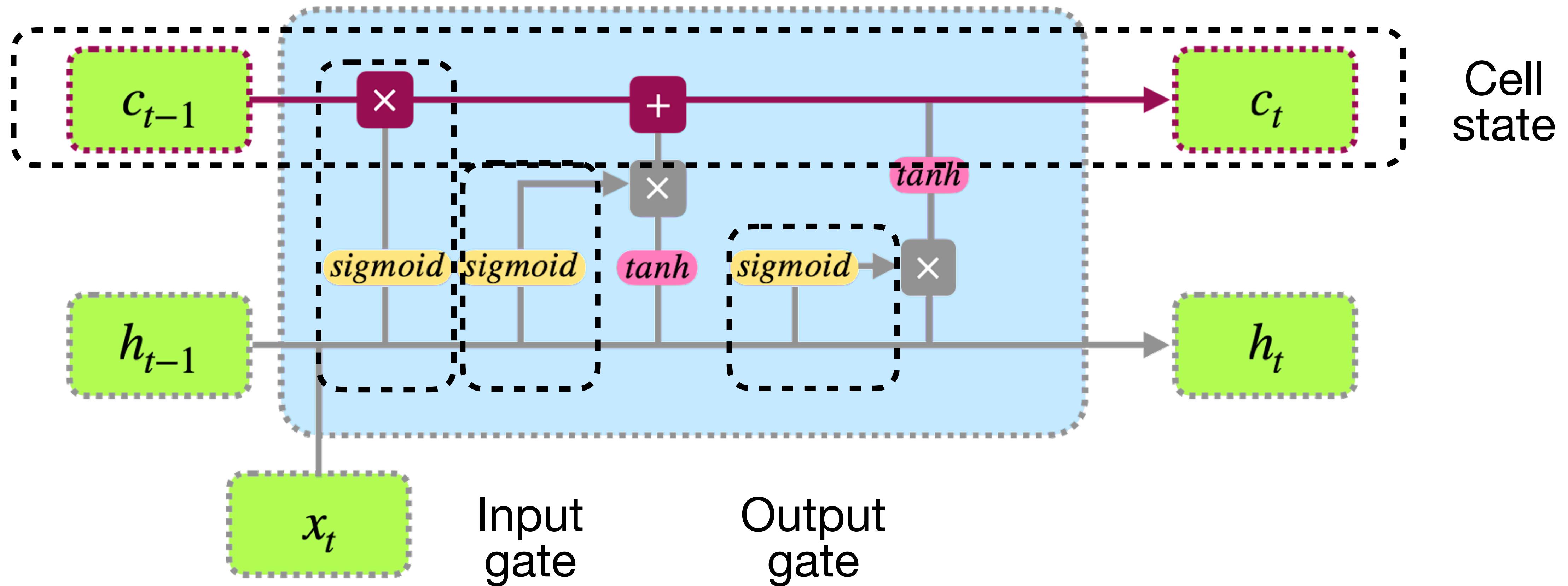
- Plain RNN is not a good choice
- LSTM and GRU cells work better
- LSTM is slightly better
- GRU is a little faster



Long Short-Term Memory (LSTM)

Forget gate

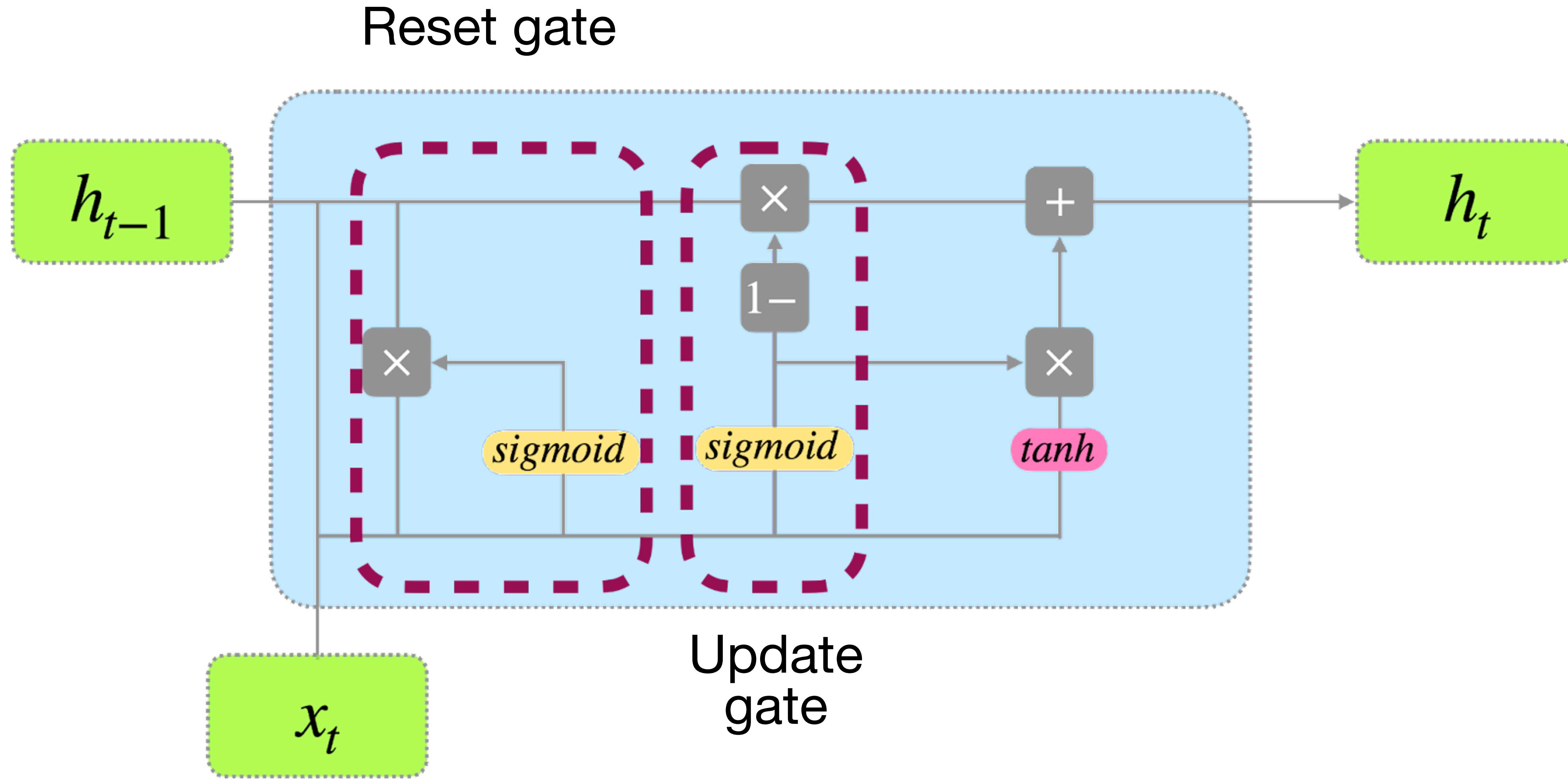
Skoltech



LSTM gates

- **Forget gate** decides what information should be thrown away or kept
- **Input Gate:** update the cell state
- **Cell state:**
 - 1) Cell state gets pointwise multiplied by the forget vector: drop values in the cell state if it gets multiplied by values near 0
 - 2) Take the output from the input gate and do a pointwise addition which updates the cell state to new values that the NN finds relevant
 - 3) Get a new cell state
- **Output Gate** decides what the next hidden state should be

Gated Recurrent Unit (GRU)



GRU gates

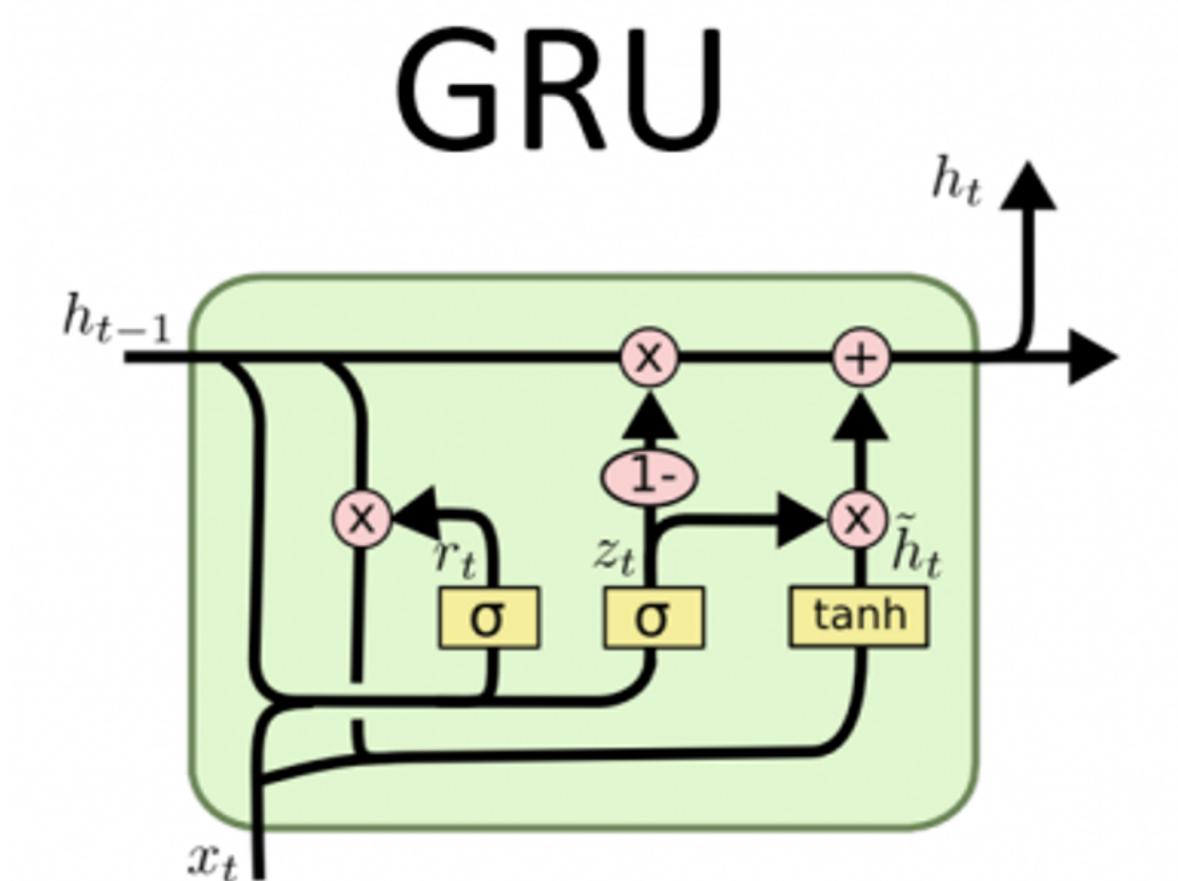
- **Reset gate** is used to decide how much past information to forget
- **Update Gate** acts similar to the forget and input gate of an LSTM

$$\mathbf{r}^t = \sigma(W_{xr}\mathbf{x}^t + W_{hr}\mathbf{h}^{t-1} + b_r)$$

$$\mathbf{z}^t = \sigma(W_{xz}\mathbf{x}^t + W_{hz}\mathbf{h}^{t-1} + b_z)$$

$$\tilde{\mathbf{h}}^t = \tanh(W_{xh}\mathbf{x}^t + W_{hr}(\mathbf{r}^t \odot \mathbf{h}^{t-1}) + b_h)$$

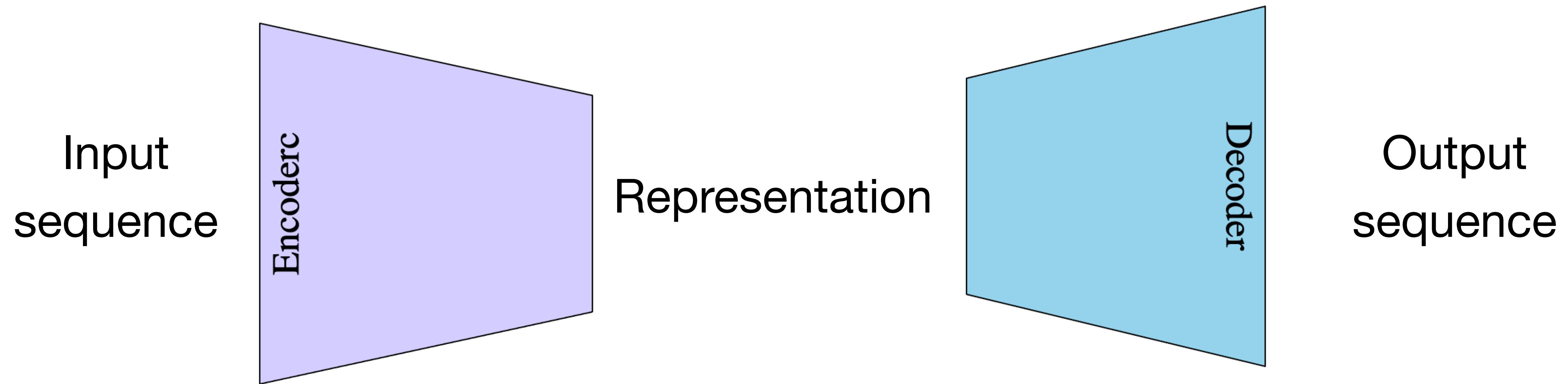
$$\mathbf{h}^t = \mathbf{z}^t \odot \mathbf{h}^{t-1} + (1 - \mathbf{z}^t) \odot \tilde{\mathbf{h}}^t$$



Machine translation: application example

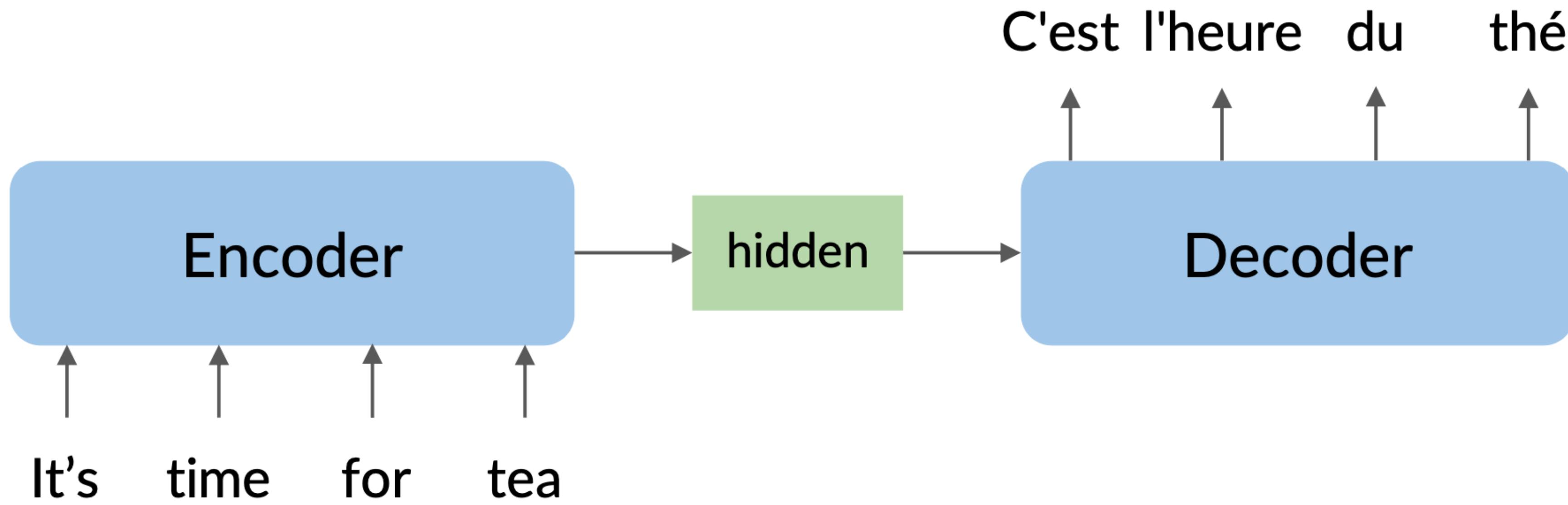
Neural network for machine translation

- Let's bring two RNNs
- Let's make **seq2seq (sequence to sequence)** architecture out of them

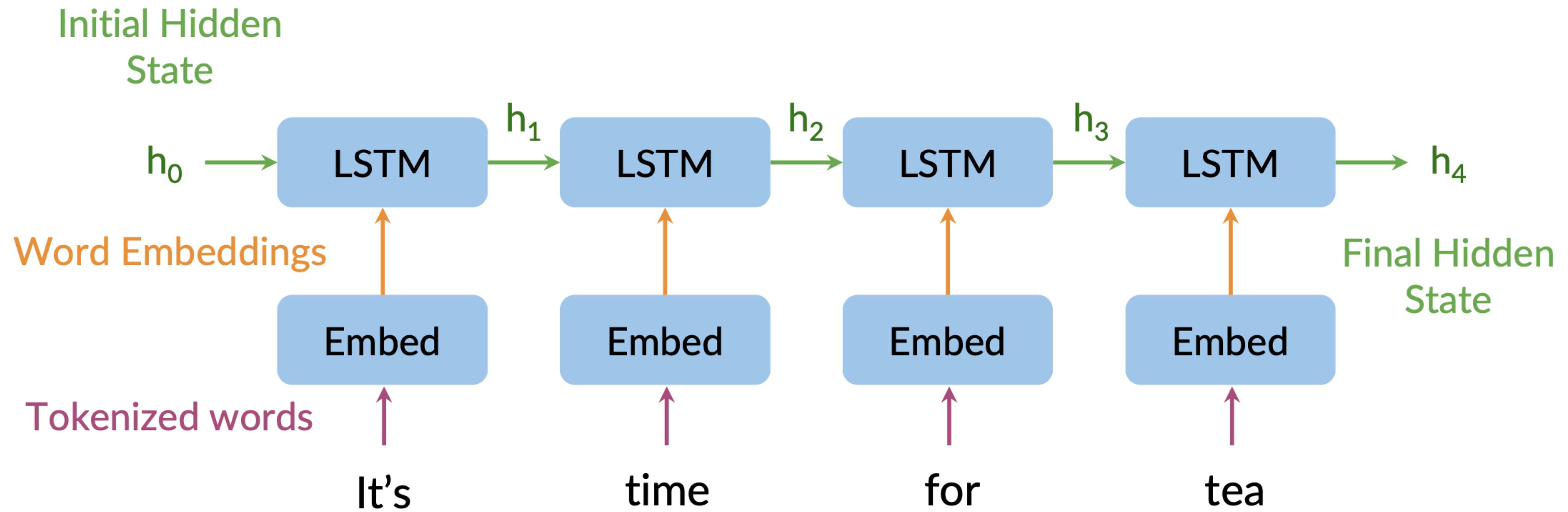


Neural network for machine translation

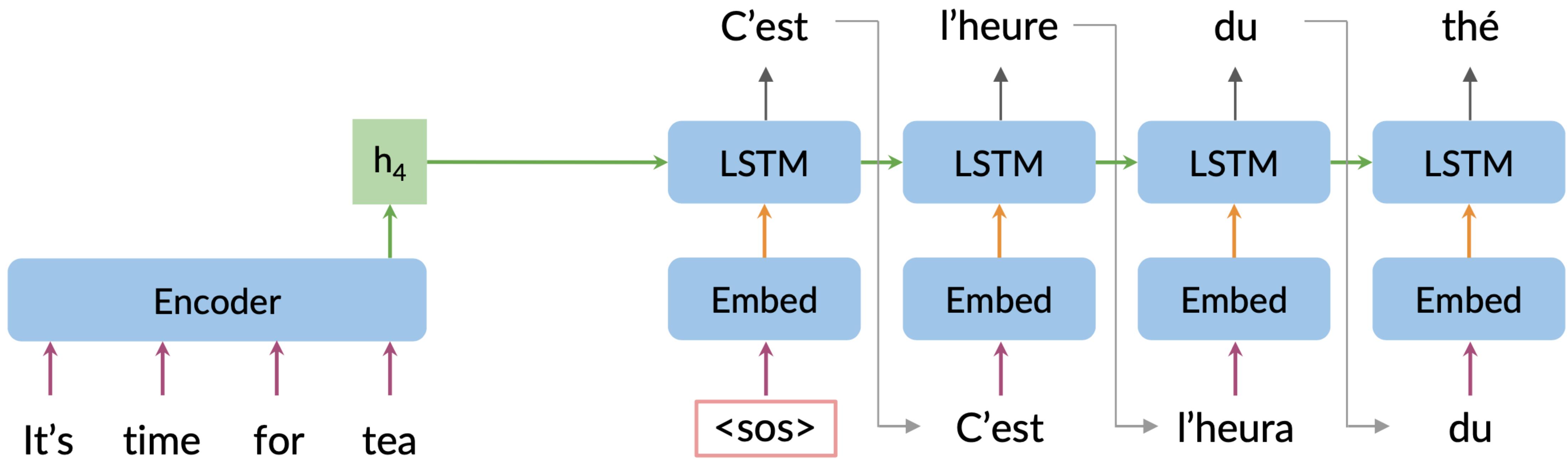
- Let's bring two RNNs
- Let's make **seq2seq (sequence to sequence)** architecture out of them



Encoder



Decoder



Main issue: bottleneck, the network will have time to forget everything

Conclusion

- RNN's are good for processing sequence data for predictions but suffers from short-term memory
- LSTM's and GRU's were created as a method to mitigate short-term memory using mechanisms called gates
- RNNs struggle dealing with long sequences

Next lecture: Models with attention