

Lecture 7: RNNs and Distributed Word Reps

Word Representations

Consider a simple classification task on text

"I love this movie.
I've seen it many times
and it's still awesome."



ML model



"This movie is bad.
I don't like it at all.
It's terrible."



ML model



How to represent text as input to ML models

Word Representations

for vocabulary size $D=10$, the one-hot vector of word ID $w=4$ is

$$\mathbf{e}(w) = [0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$$

The major problem with the one-hot representation is that it is very high-dimensional

- ▶ the dimensionality of $\mathbf{e}(w)$ is the size of the vocabulary
- ▶ a typical vocabulary size is $\approx 100\ 000$
- ▶ a window of 10 words would correspond to an input vector of at least 1 000 000 units!

Distributed Word Representations

- A common approach is to associate each discrete object with a “dense” continuous vector

Word	w	$C(w)$
“the”	1	[0.6762, -0.9607, 0.3626, -0.2410, 0.6636]
“a”	2	[0.6859, -0.9266, 0.3777, -0.2140, 0.6711]
“have”	3	[0.1656, -0.1530, 0.0310, -0.3321, -0.1342]
“be”	4	[0.1760, -0.1340, 0.0702, -0.2981, -0.1111]
“cat”	5	[0.5896, 0.9137, 0.0452, 0.7603, -0.6541]
“dog”	6	[0.5965, 0.9143, 0.0899, 0.7702, -0.6392]
“car”	7	[-0.0069, 0.7995, 0.6433, 0.2898, 0.6359]
...

Slide Credit Hugo Larochelle

- We can now try to make meaningful distances in vector space by learning $C(w)$

Distributed vs Local

- Distributed representations for symbol -> using dense vs local (sparse)

"the"	1	[0.6762, -0.9607, 0.3626, -0.2410, 0.6636]
-------	---	--

vs

[0 0 0 1 0 0 0 0 0 0]

- Note: the word “distributed” is sometimes used in slightly different context in deep learning literature

Pytorch Embedding

```
import torch.nn as nn
import torch

words = 5
dim = 100

embed = nn.Embedding(words, dim)

embed(torch.LongTensor([1]))
```



```
tensor([[-0.6355, -0.2671,  0.356,  0.2689,  0.2781, -1.3727,  0.7757,  2.0445,
        1.4461,  0.3756, -0.0162,  1.1468, -1.1845,  0.8442, -1.3146,  0.171,
        0.5611, -0.8307,  1.1716, -1.3098, -2.9567,  0.4804,  1.4355, -0.3178,
        0.2934,  1.1778,  1.2912,  0.1803, -1.9347,  1.4485,  0.0195, -0.8828,
        1.2395,  0.2929, -1.9109, -0.6097, -1.4886, -0.530,  0.793,  1.2733,
       -0.4736,  0.020, -0.0703,  0.1745, -2.8354,  1.3762, -0.7413, -0.4601,
       -0.4509, -0.7949,  0.0967, -0.6718, -0.4257,  0.0081, -0.6909, -0.211,
       -0.334, -0.2481,  0.1715,  0.4038, -1.9862,  1.5716,  0.8283, -0.3971,
       -0.293,  1.4967, -0.0121, -0.3886,  0.2584, -0.1783,  0.2082,  2.0378,
       -1.5603, -0.9889, -0.7937, -0.4811,  0.515, -0.6355,  0.0782,  0.6339,
        0.6144,  0.4604,  1.6851,  0.6491, -1.2013, -0.385, -0.7349,  1.4328,
       -0.9501, -0.5964,  0.091,  0.6858, -0.985, -2.502, -0.6703,  0.8197,
        0.9445,  0.2693, -0.5862, -0.935], grad_fn=<EmbeddingBackward>)
```

```
embed.weight.shape
```

```
torch.Size([5, 100])
```

- Embedding class is essentially a weight matrix where the forward indexes the rows
- Can be optimized for large vocabulary

Language Modeling

- Language modelling is concerned with capturing the statistical properties of text
- A common objective used for learning representations of words and sentences is to predict the distribution for the next word or the missing word

I love this movie and I've seen it many _____

$$p(w_t \mid w_{t-(n-1)}, \dots, w_{t-1})$$

$$P(w_t = k \mid w_{t-n+1}, \dots, w_{t-1}) = \frac{e^{a_k}}{\sum_{l=1}^N e^{a_l}}$$

I love this _____ and I've seen it many times

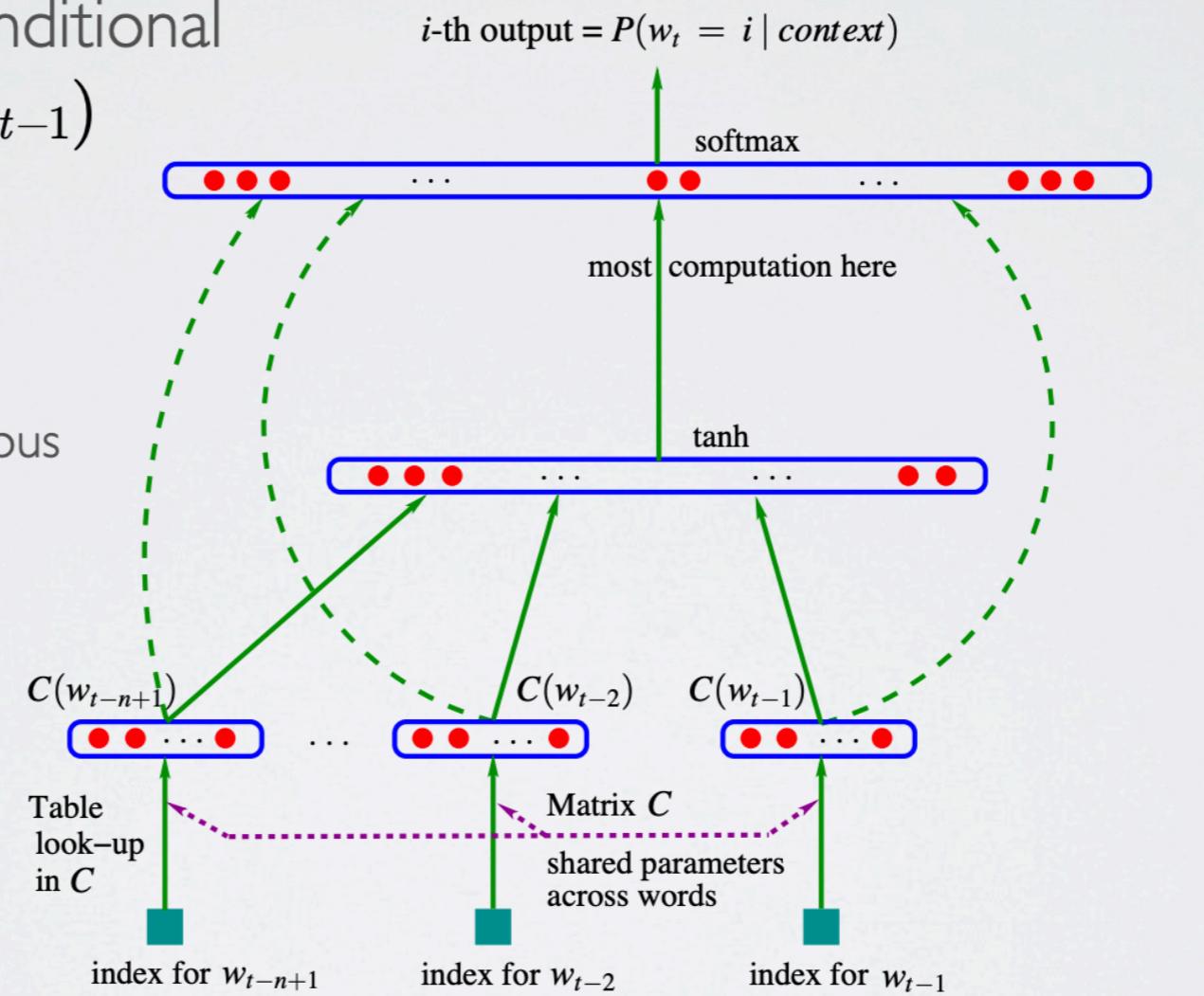
$$L(\theta) = \sum_t \log P(w_t \mid w_{t-n+1}, \dots, w_{t-1}).$$

Simple Neural Language Model

Solution: model the conditional
 $p(w_t | w_{t-(n-1)}, \dots, w_{t-1})$
with a neural network

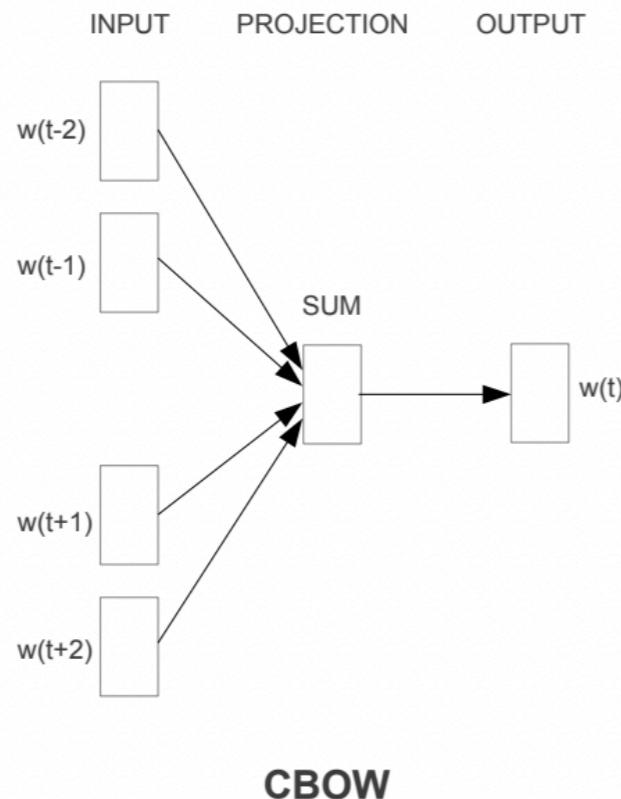
- learn word representations to allow transfer to n -grams not observed in training corpus

Bengio, Ducharme,
Vincent and Jauvin, 2003



Slide Credit Hugo Larochelle

Word2vec



- Learns useful features using a simple linear model

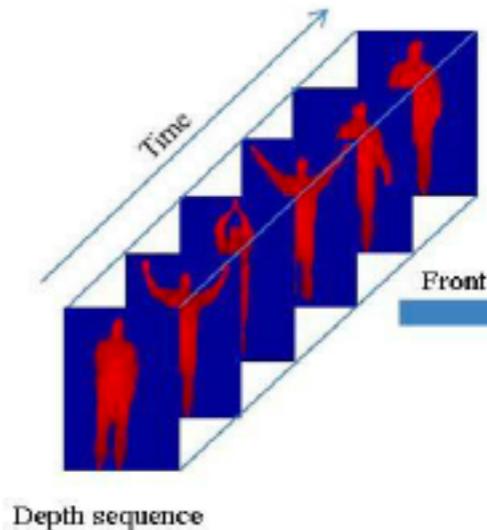
What is missing in these models (Word2vec and Bengio et al) ?

Sequence Modeling

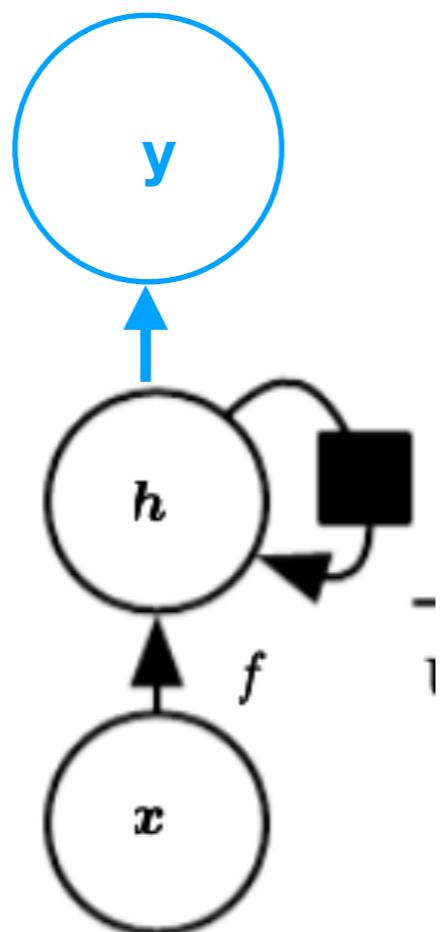
- Audio
- Video
- Text
- Genomes

Target Genome

ATTTGCGCAGAGACCTAAGGCATTAGCTTGGCCCTAAAG



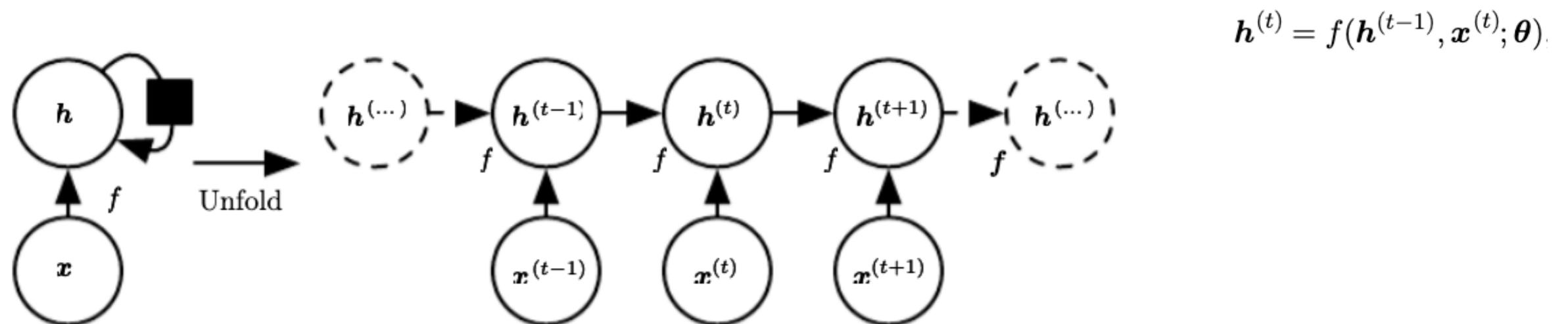
Recurrent Neural Networks



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}),$$

Recurrent Neural Networks

- Recurrent networks have inductive bias for sequential data
 - Have a state
 - Take input at each time-step of the sequence
 - Can be variable length - infinitely deep network
- Often helpful to think about an unrolled sequence of steps with same operations
 - Becomes analogous to feedforward network with repeated function applied



Example Tasks

- Output at the end e.g. Sentence classification
- Output at each time step e.g. Next word prediction

Sentence Classification Example

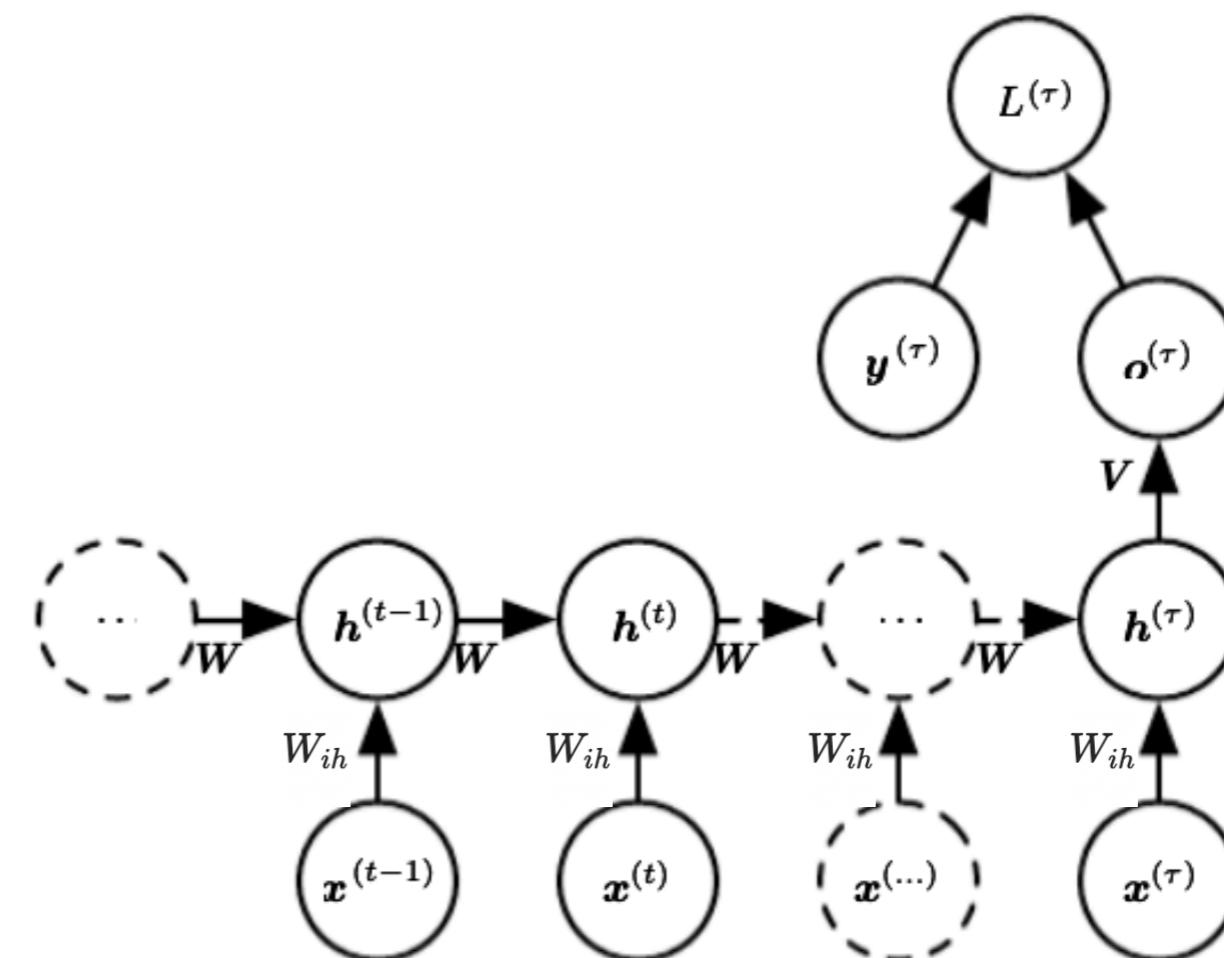
"I love this movie.
I've seen it many times
and it's still awesome."



"This movie is bad.
I don't like it at all.
It's terrible."



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta})$$

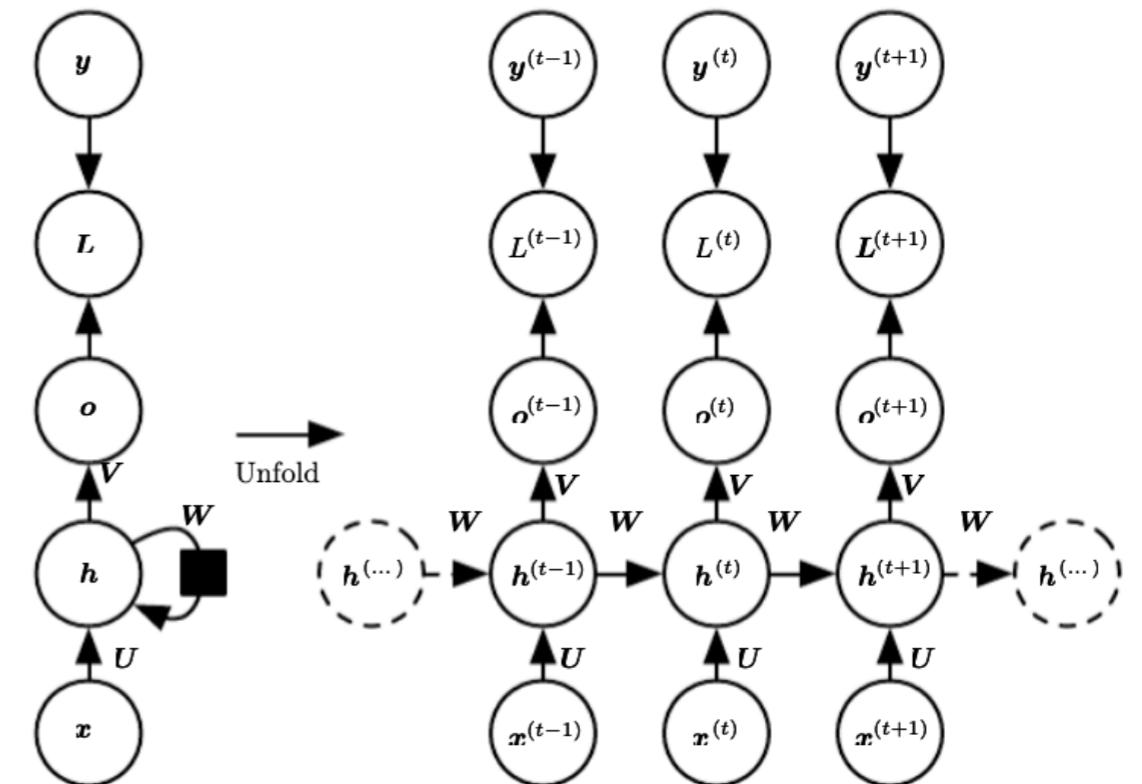


$$h_t = \tanh(W_{ih}x_t + b_{ih} + W_{hh}h_{(t-1)} + b_{hh})$$

RNN with Output at Each Step

- Can have output at each time step besides a state

$$\begin{aligned}\mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}),\end{aligned}$$

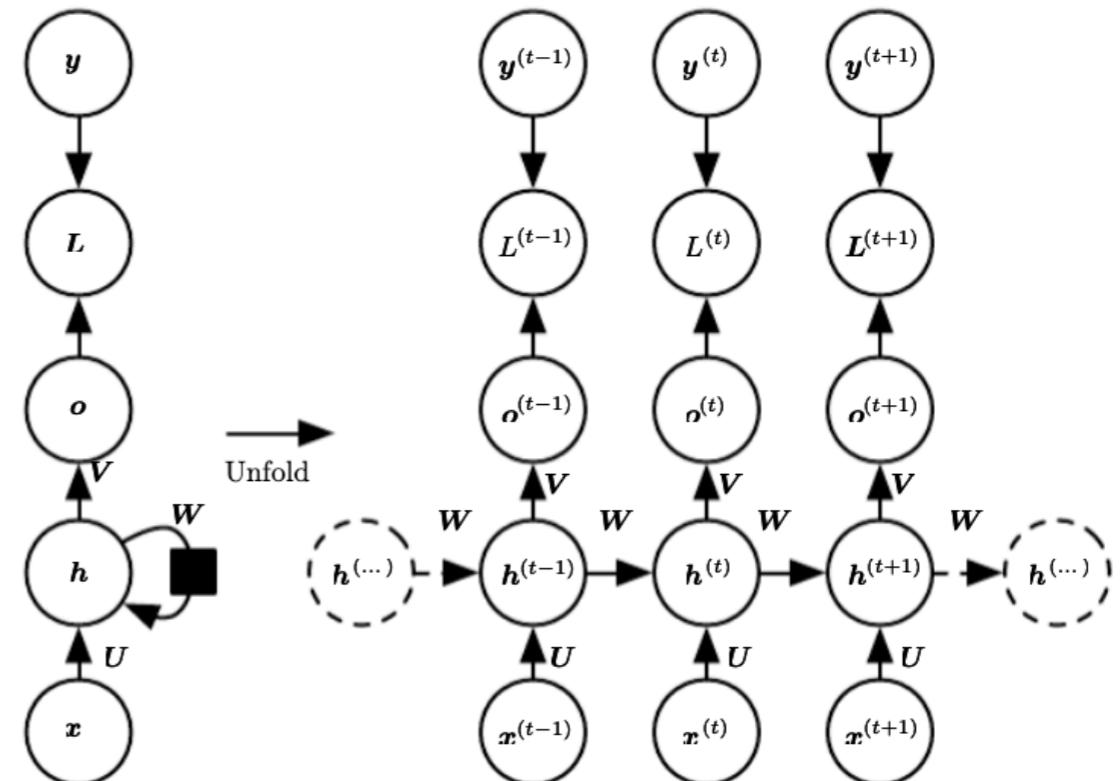


- “Deeper” RNN refers to more layers after h

Next Char/Word Prediction

- Can have output at each time step besides a state

$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \end{aligned}$$



Loss is summed over t

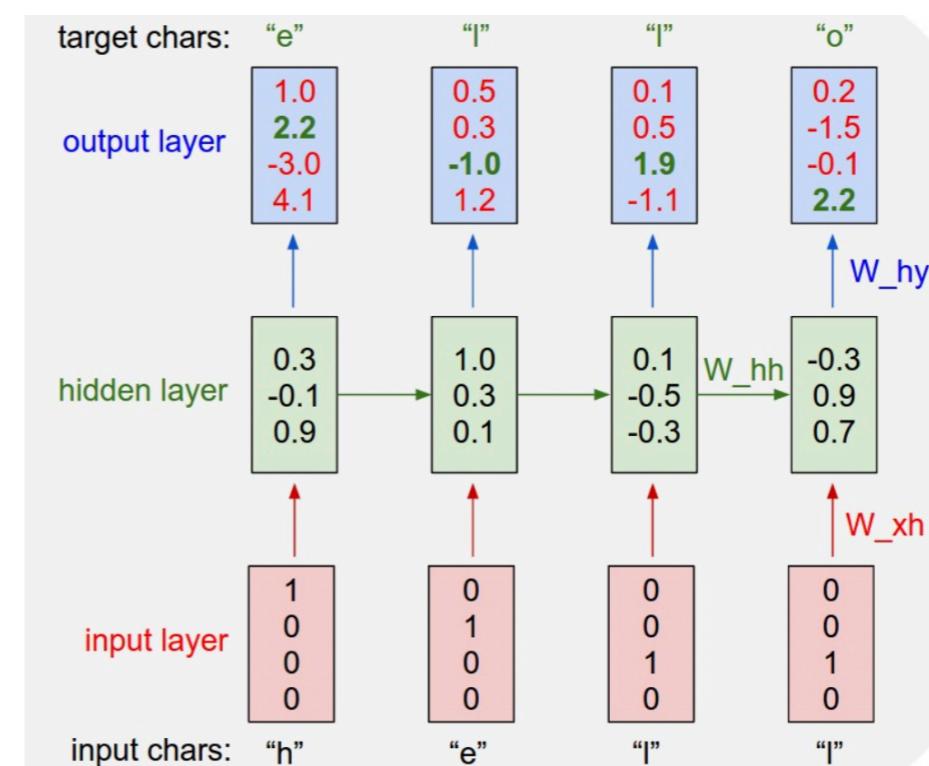
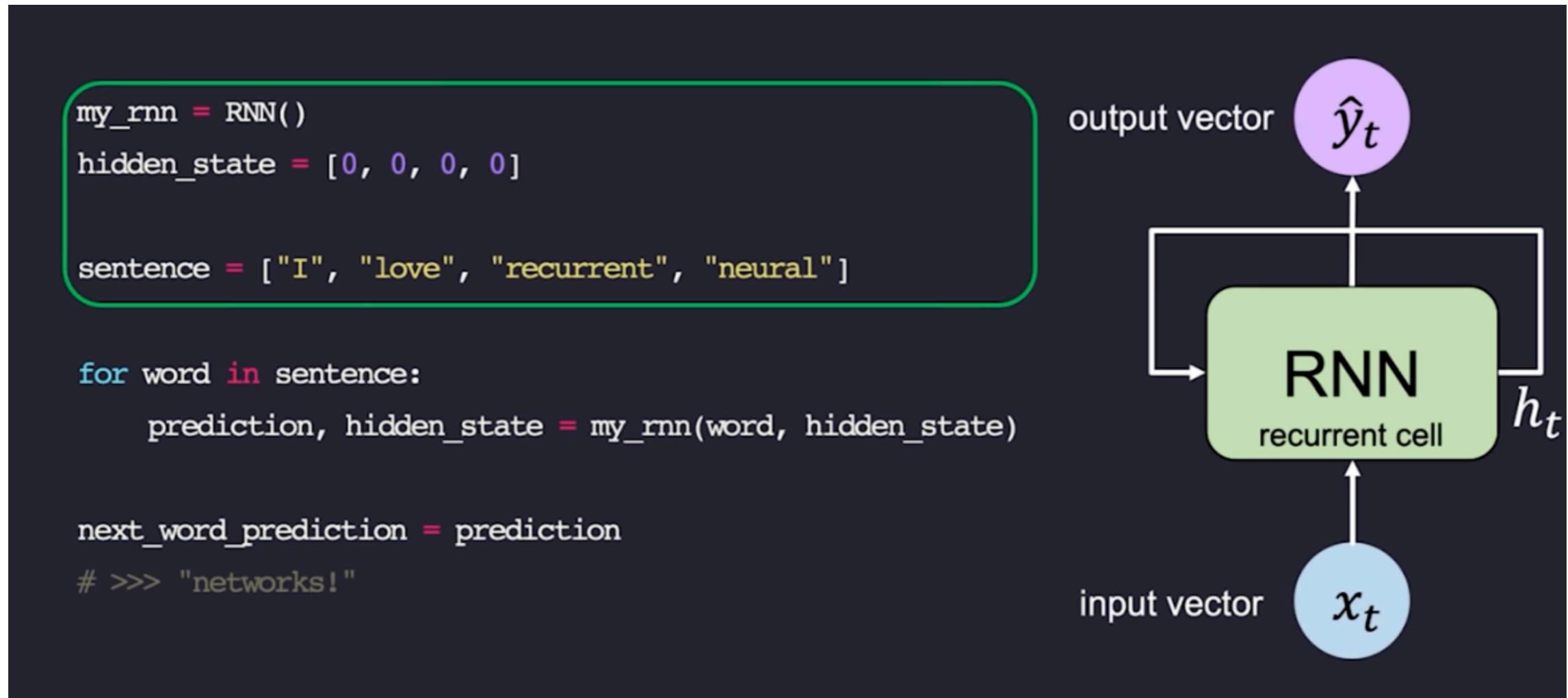


Image Credit: Karpathy

Next Word Prediction Example



Slide Credit: Ava Soleimany

RNN as Generative Model

- Given a start character generate an entire sequence

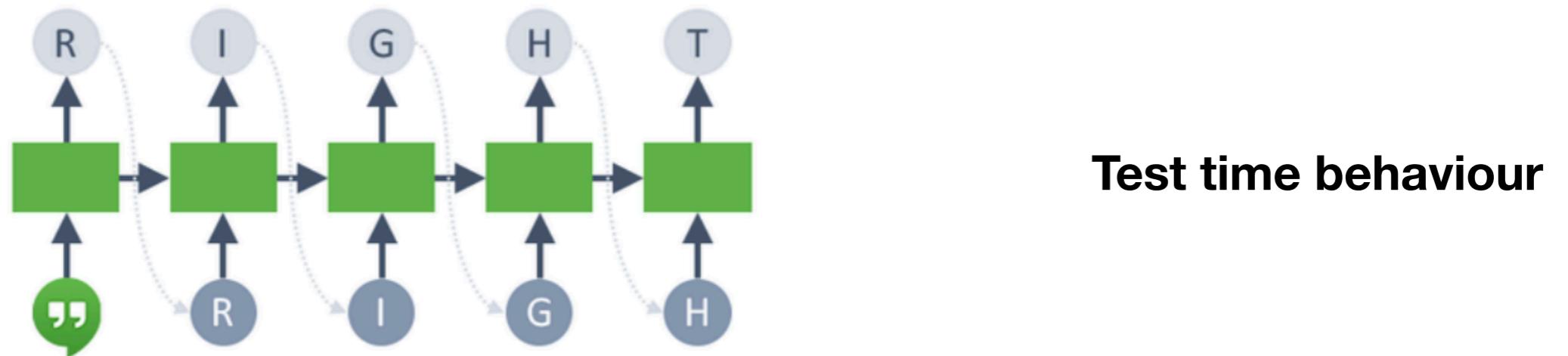
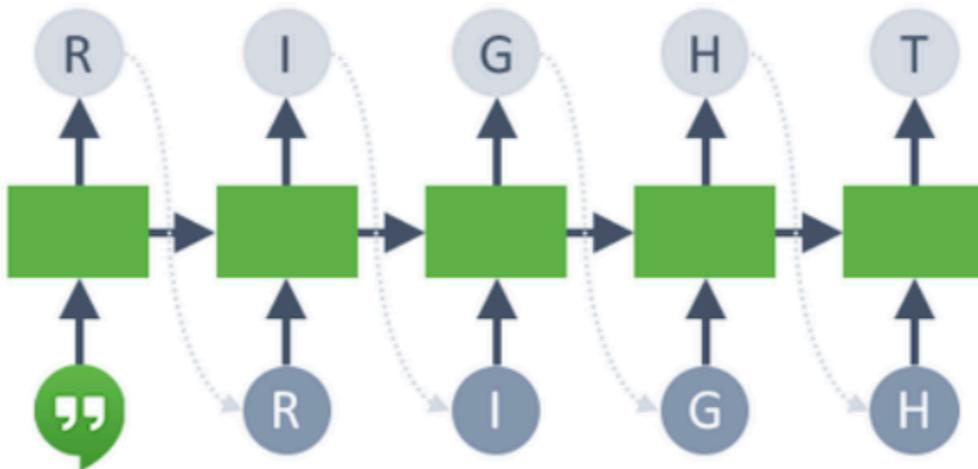


Image credit: Lisa Zhang

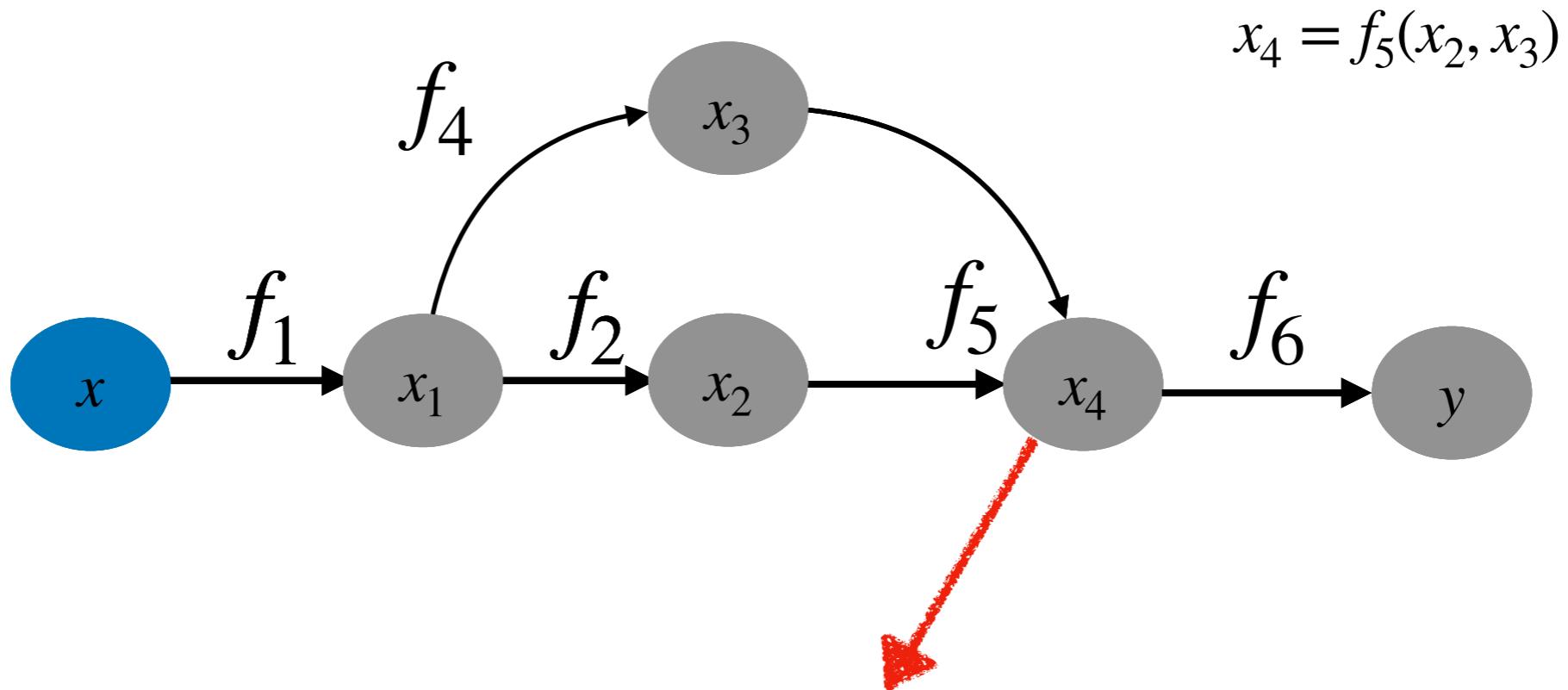
- We sample from the output distribution at each time step

Teacher Forcing



- At training time we feed GT characters (teacher forcing)
- This can however cause an issue with training and test time distributions being different
 - MIXER, SEARN, Scheduled sampling Bengio et al
 - Similar problem arises in “imitation learning behavioral cloning” e.g. for autonomous driving

Reminder

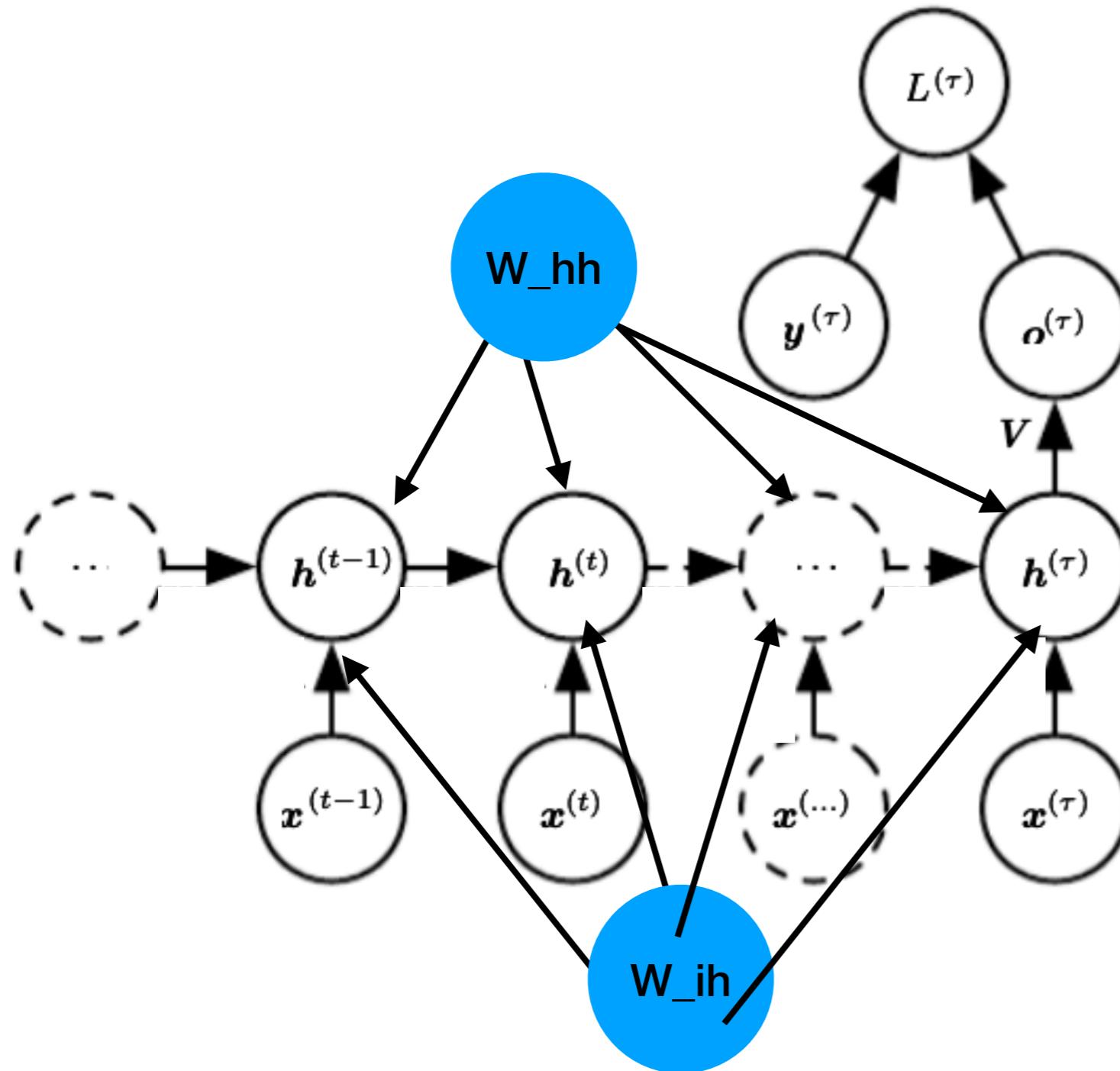


$$\frac{\partial y}{\partial \mathbf{x}_0} = \frac{\partial y}{\partial \mathbf{x}_4} \left(\frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1} + \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \mathbf{x}_1} \right) \frac{\partial \mathbf{x}_1}{\partial \mathbf{x}_0}$$

$$\frac{\partial y}{\partial \mathbf{x}_j} = \sum_{s \in Child(j)} \frac{\partial y}{\partial \mathbf{x}_s} \frac{\partial \mathbf{x}_s}{\partial \mathbf{x}_j}$$

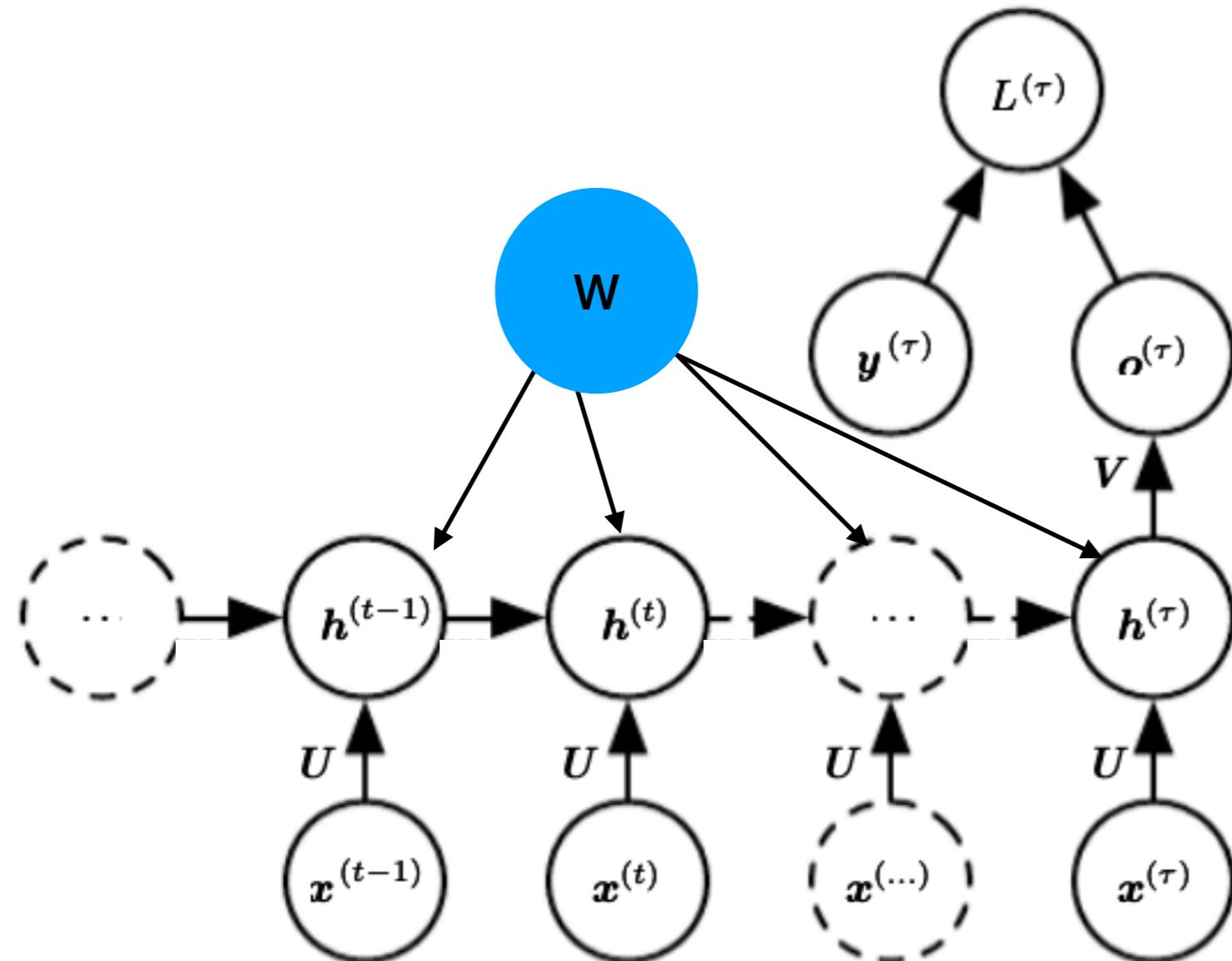
Backpropagation Through Time

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



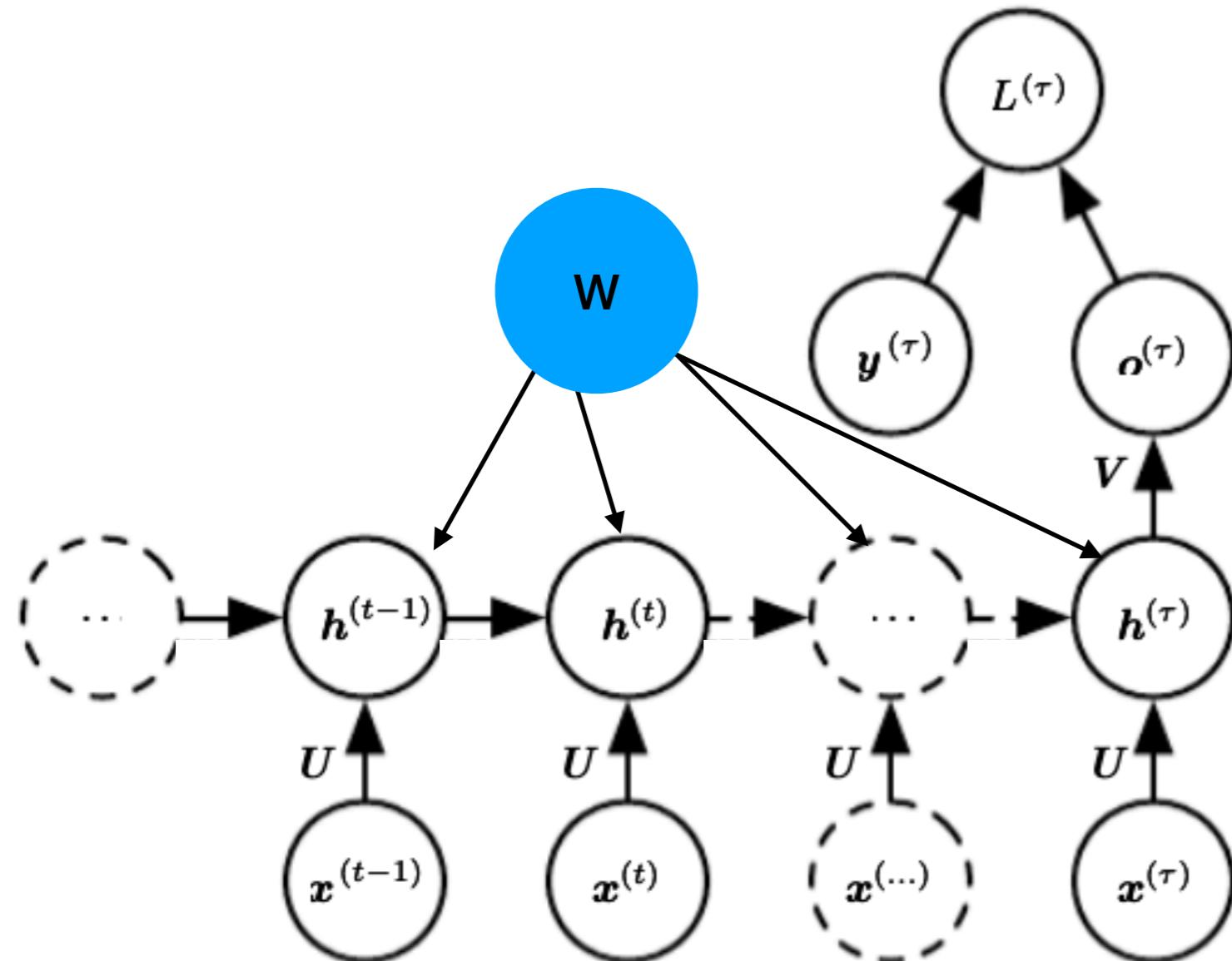
- BPTT – backprop through unrolled RNN

Backprop and Truncated BPTT



- For long sequence memory and update cost for a single instance can be expensive
- Truncated BPTT is a technique that simply stops backprop before reaching $t=0$

Backprop and Truncated BPTT



- For long sequence memory and update cost for a single instance can be expensive
- Truncated BPTT is a technique that simply stops backprop before reaching $t=0$

Vanishing and Exploding Gradients

$$h_t = W_{hh}\rho(h_{t-1}) + W_{ih}x_t \quad L_T = l(W_T h_T, y)$$

$$\frac{\partial L_T}{\partial \mathbf{W}_{hh}} = \frac{\partial L_T}{\partial h_T} \prod_{i=T}^1 \frac{\partial h_i}{\partial h_{i-1}} = \frac{\partial L_T}{\partial h_T} \prod_{T=t..1} W_{hh}^T diag(\rho'(h_{t-1}))$$

$$\frac{\partial h_t}{\partial h_{t-1}} = W_{hh}^T diag(\rho'(h_{t-1}))$$

Long Term Dependency

- Due to vanishing gradients it is difficult to learn long term dependencies
- Popular variants of RNNs aim to address this issue they are the LSTM and variants
- Gradient clipping also helps to avoid exploding gradients



LONG SHORT-TERM MEMORY

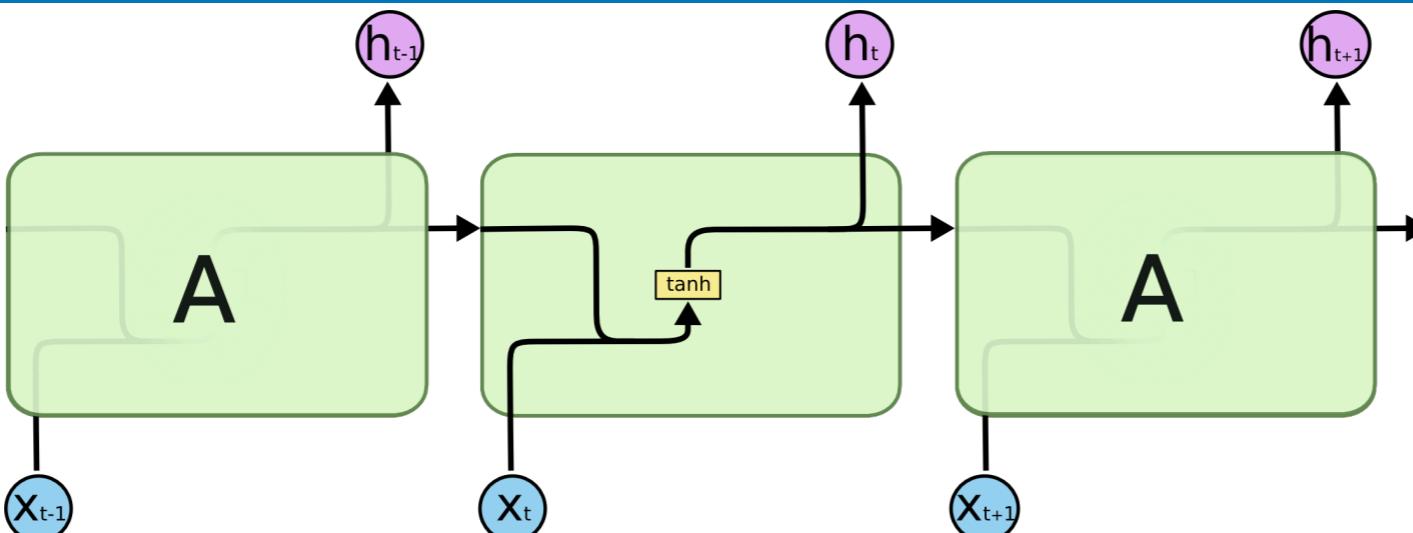
NEURAL COMPUTATION 9(8):1735–1780, 1997

Sepp Hochreiter
Fakultät für Informatik
Technische Universität München
80290 München, Germany
hochreit@informatik.tu-muenchen.de
<http://www7.informatik.tu-muenchen.de/~hochreit>

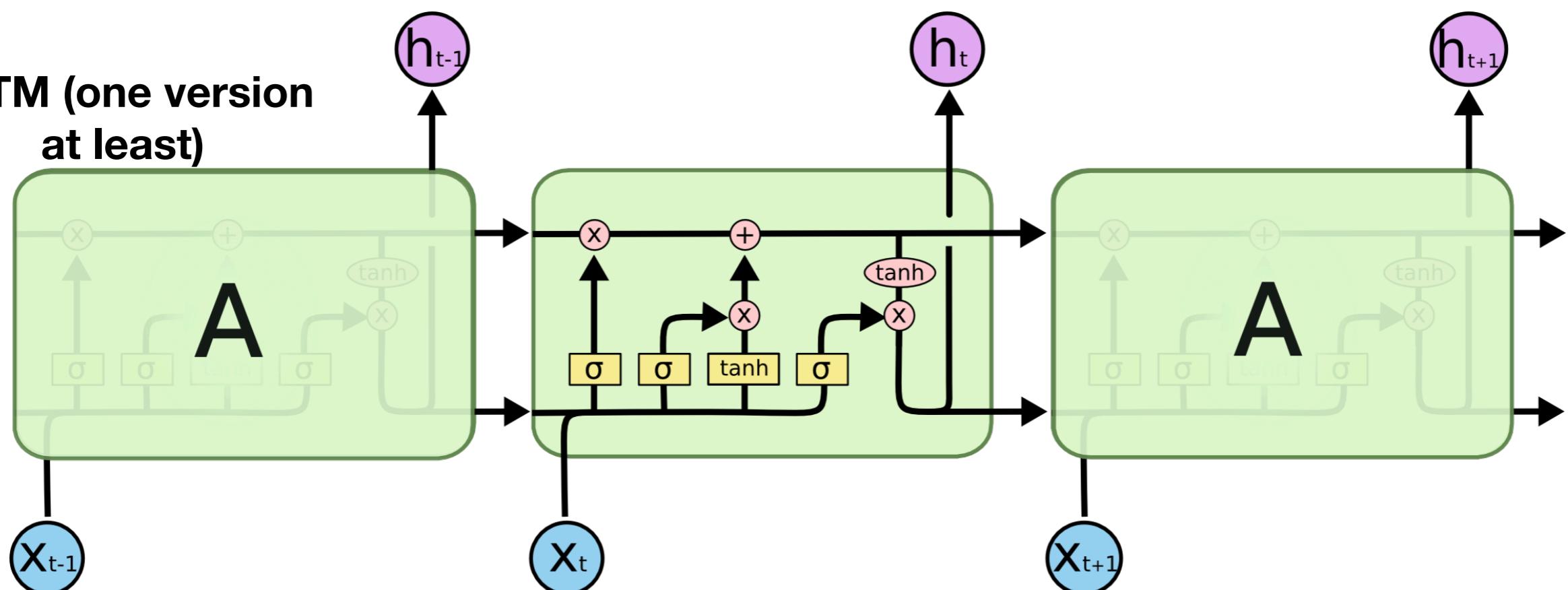
Jürgen Schmidhuber
IDSIA
Corso Elvezia 36
6900 Lugano, Switzerland
juergen@idsia.ch
<http://www.idsia.ch/~juergen>

LSTMs

Vanilla RNN



LSTM (one version
at least)



Neural Network Layer



Pointwise Operation



Vector Transfer

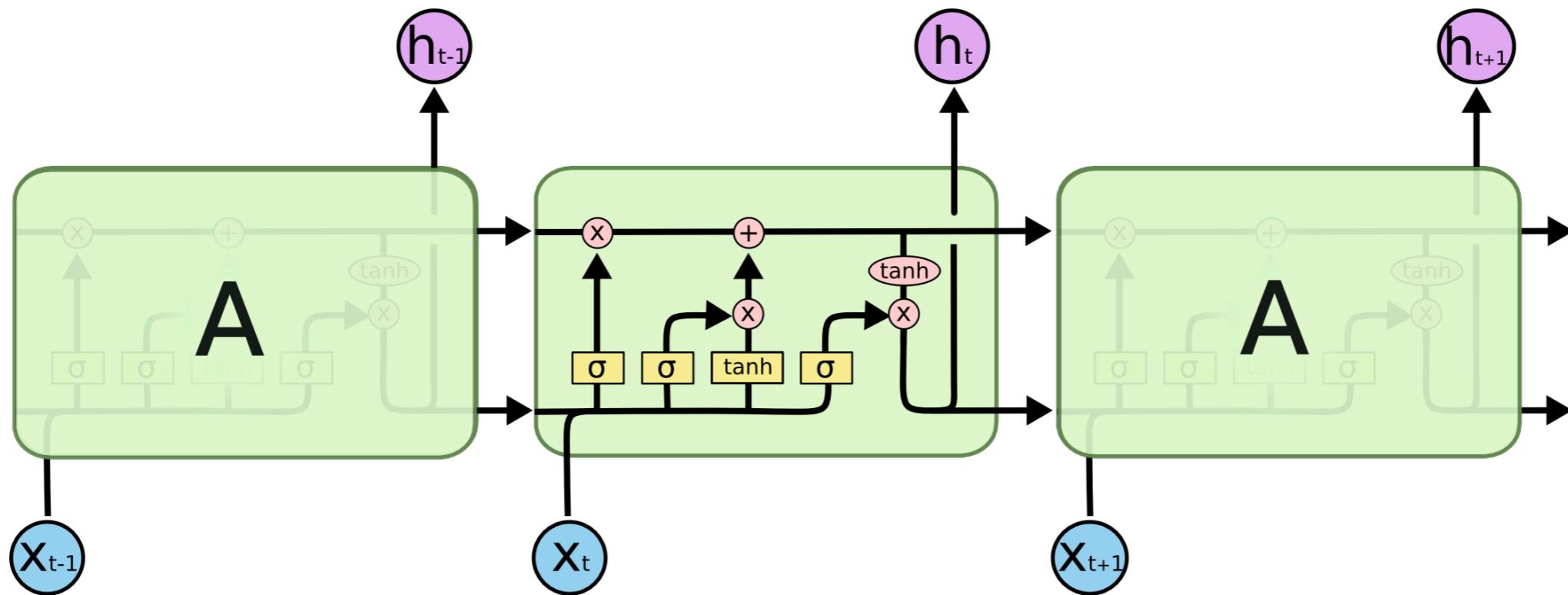


Concatenate

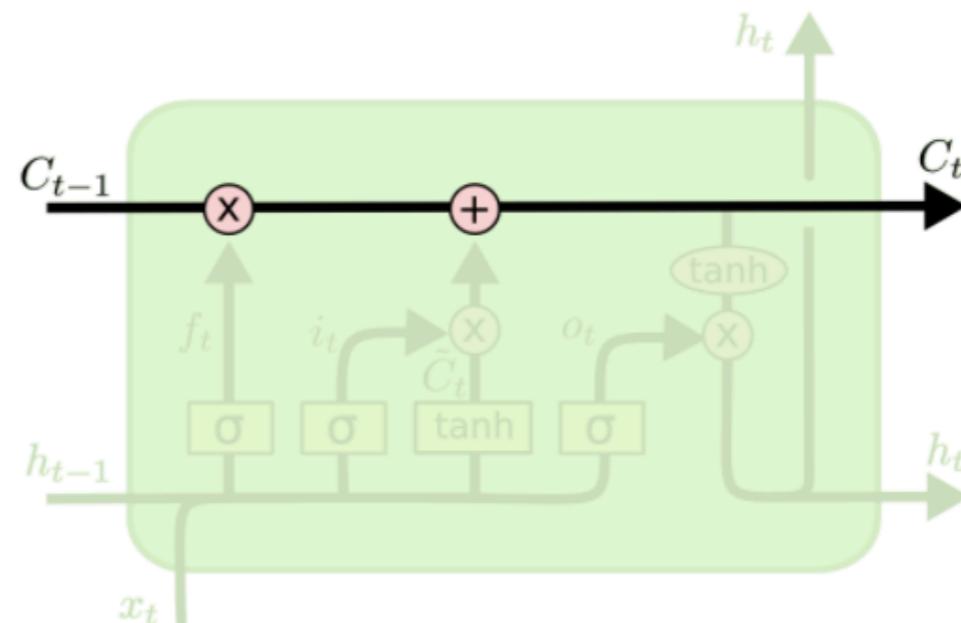


Copy

LSTMs

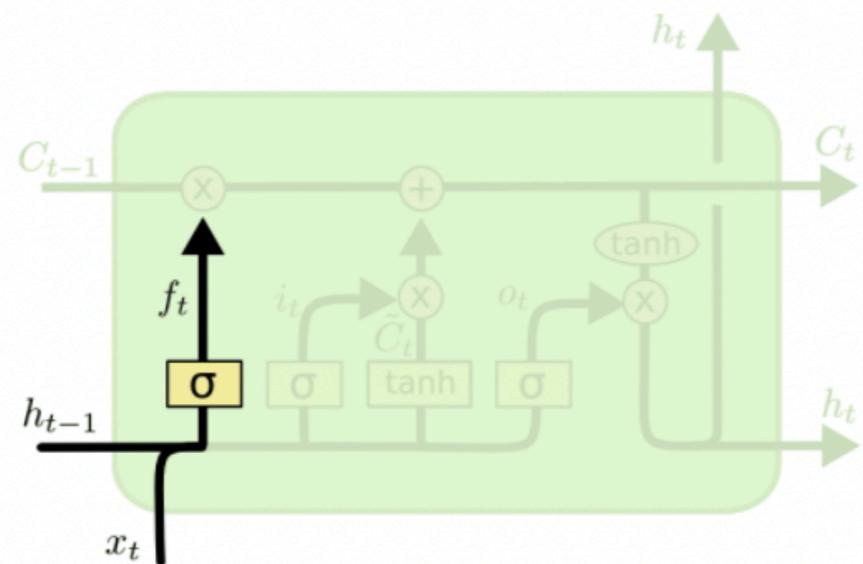


- Cell state allows information to flow across many time steps modulated by gates and hidden state
- Can draw some analogies to Skip connections in ResNets
- “RNNs overwrite the hidden state, LSTM adds to it” - Ilya Sutskever



LSTMs

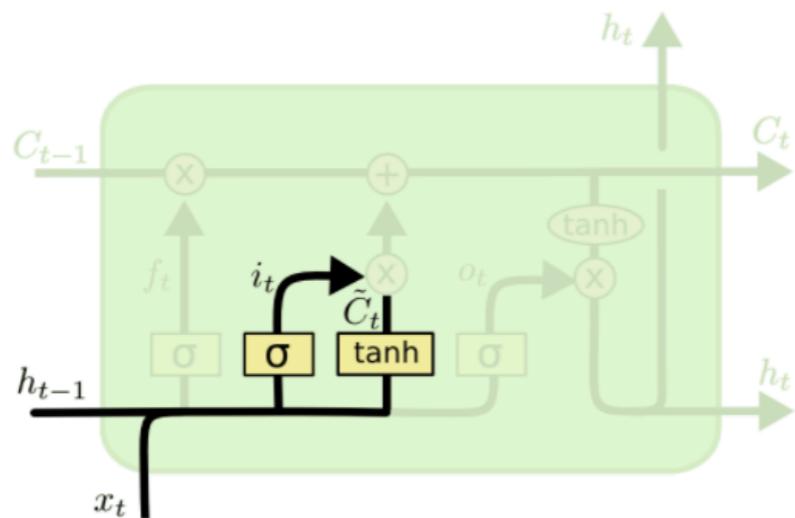
Forget gate



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

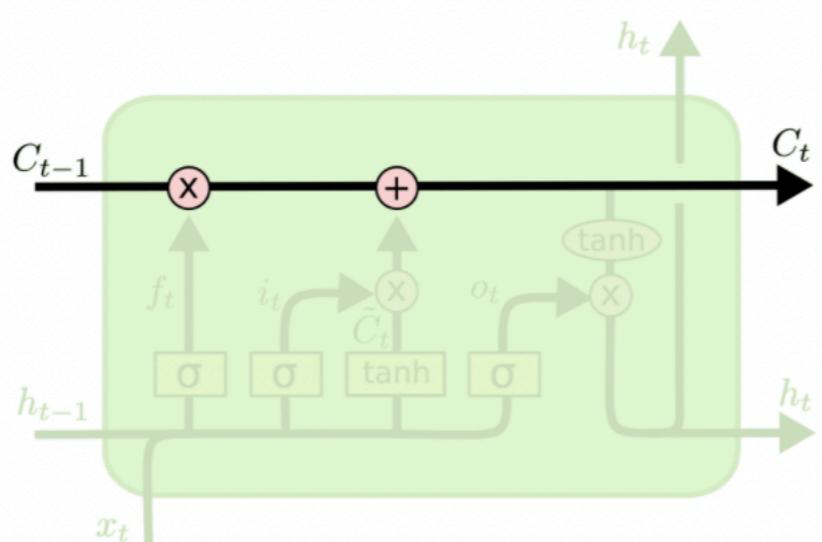
LSTMs

Input gate



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

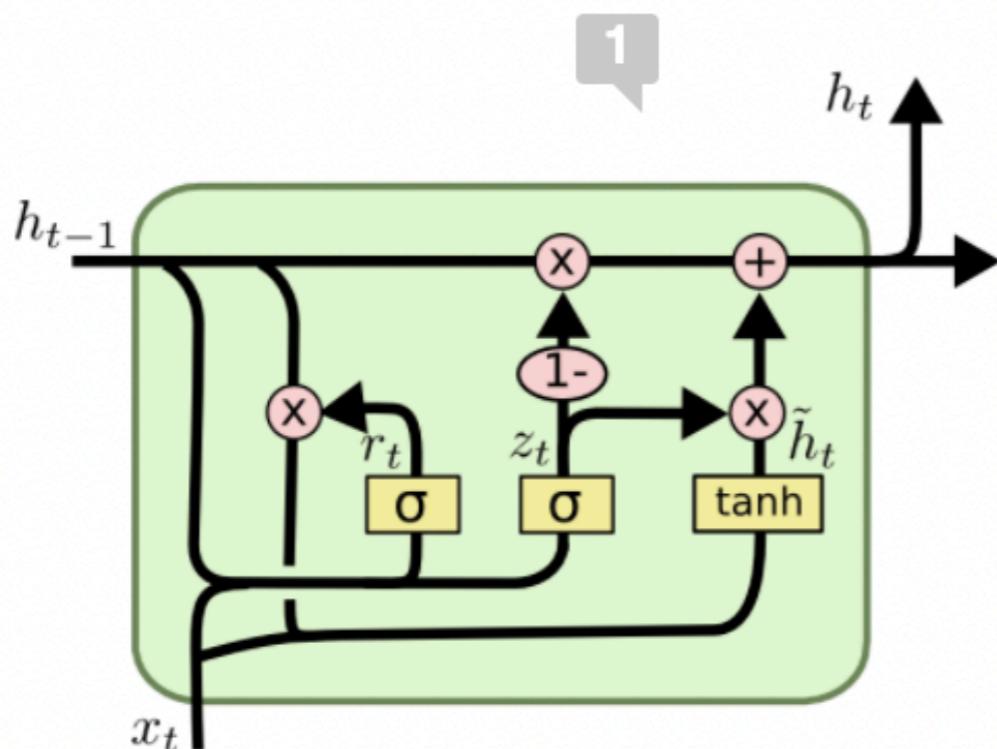
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



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

GRU

No cell state we just modulate hidden state



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

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

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

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

Issues with RNNs

- Very long term dependencies (e.g. across paragraph length) are still not well tackled by these models
 - Improved on by Memory based network and Attention based models
- Sequential processing makes them inherently slow

RNNs are more general than FeedForward Networks

- Unlike feedforward they are Turing complete
- Can implement while and for loops
- Allow having a concept of memory
- Allow us to think about building “differentiable computers”

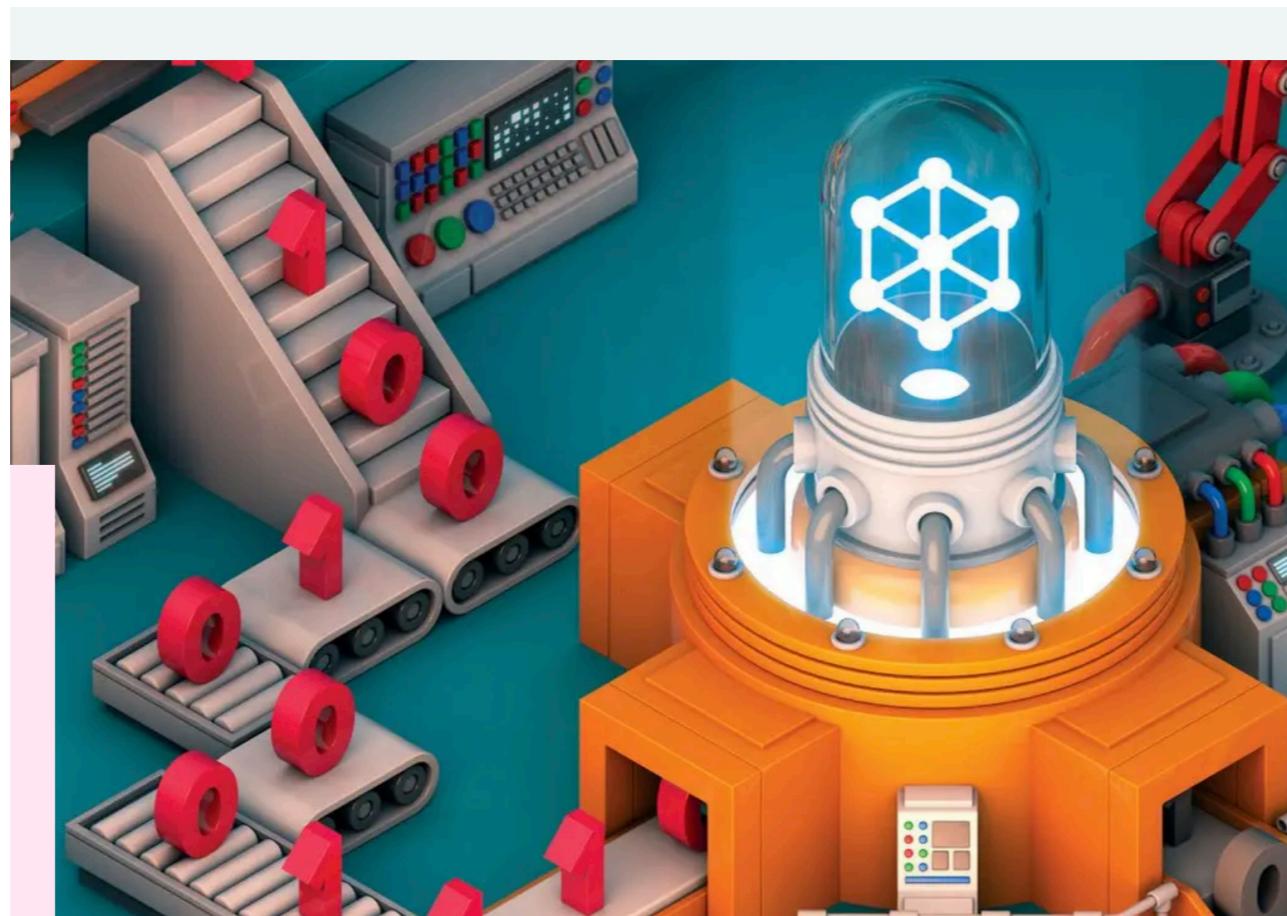
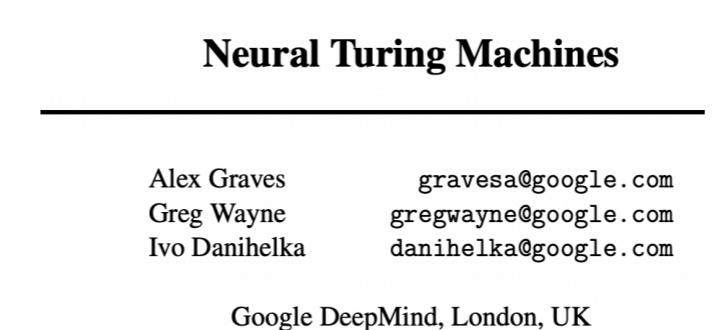


Image Credit: Deepmind

From RNNs to Differentiable Computers

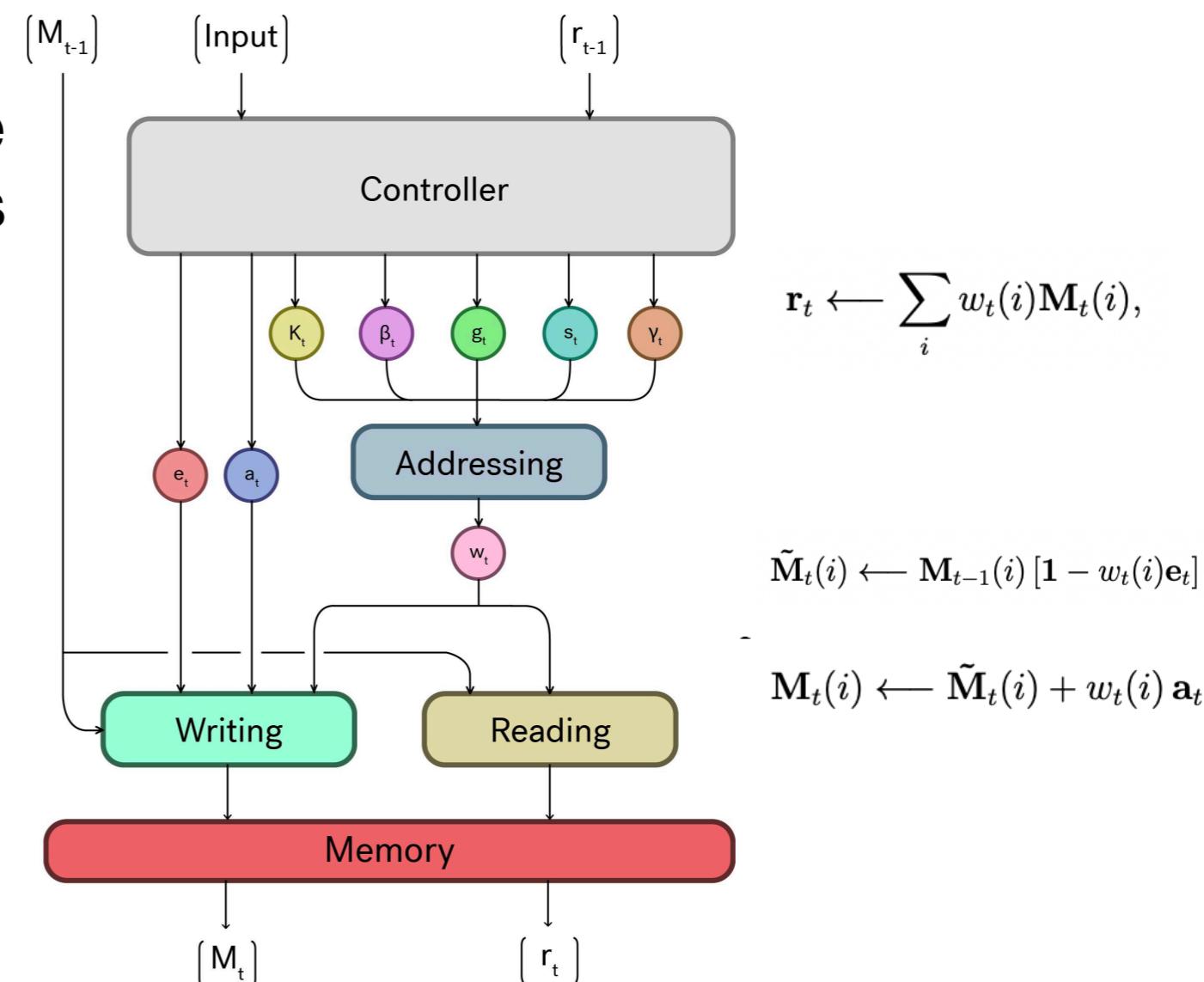
- Convnets and RNNs are just one example of how we can use network architectures to add inductive bias to our models



MEMORY NETWORKS

Jason Weston, Sumit Chopra & Antoine Bordes
Facebook AI Research
770 Broadway
New York, USA
{jase, spchopra, abordes}@fb.com

- A common approach in deep learning architecture design for various problems is to parametrize components of existing algorithms
- Important to make any discrete operations differentiable
- We can take the LSTM gating idea to an extreme by creating differentiable read and write operations into Memory



Applications of RNNs

- Speech Recognition
- Machine Translation (before 2018)
- Video
- Many others!

Sequence to Sequence

- Popularized in Machine Translation
- Way of thinking of encoding sequence and decoding sequence
- Efficiently trained RNNs + sequence to sequence model led to some breakthrough results in MT (circa 2014)

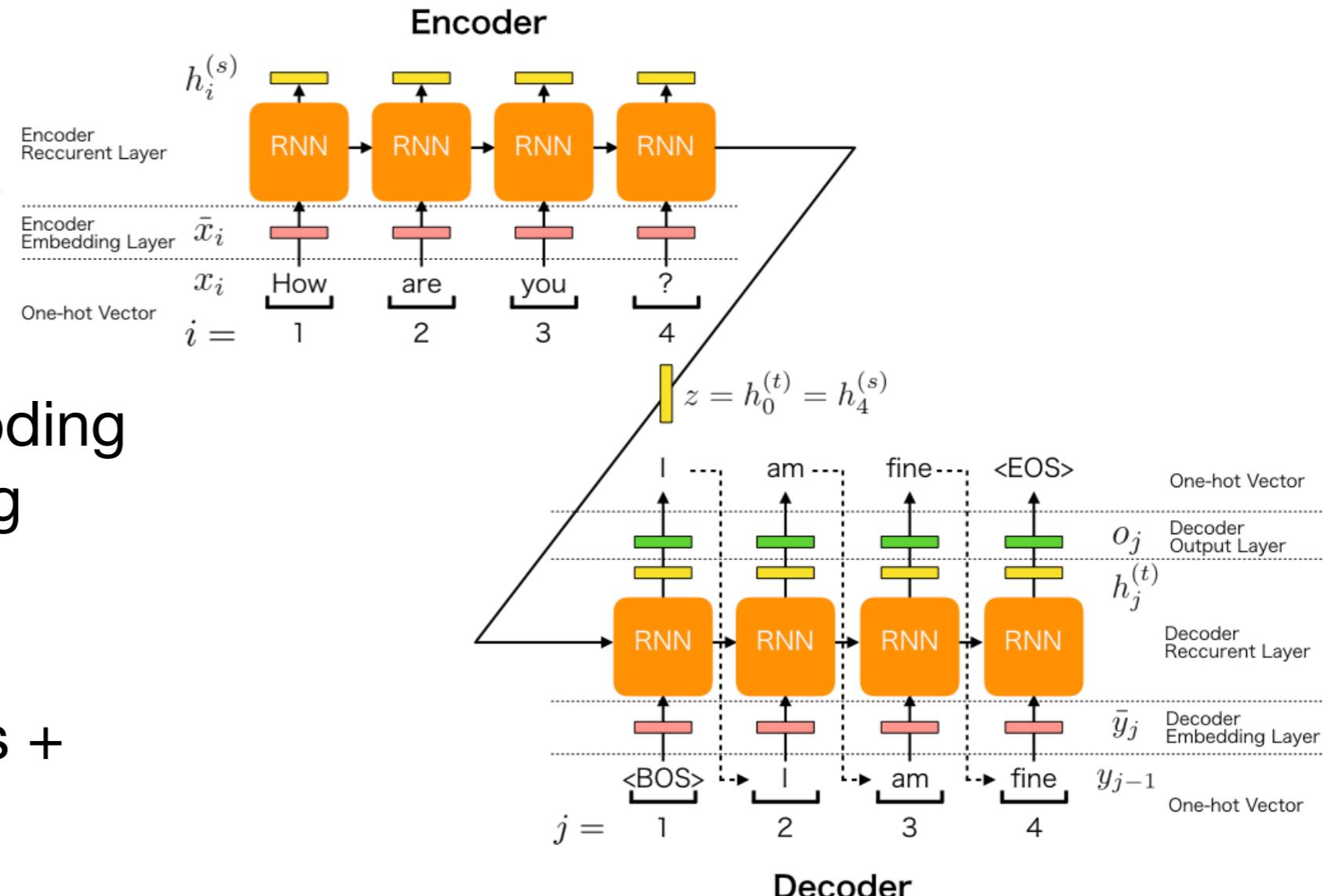
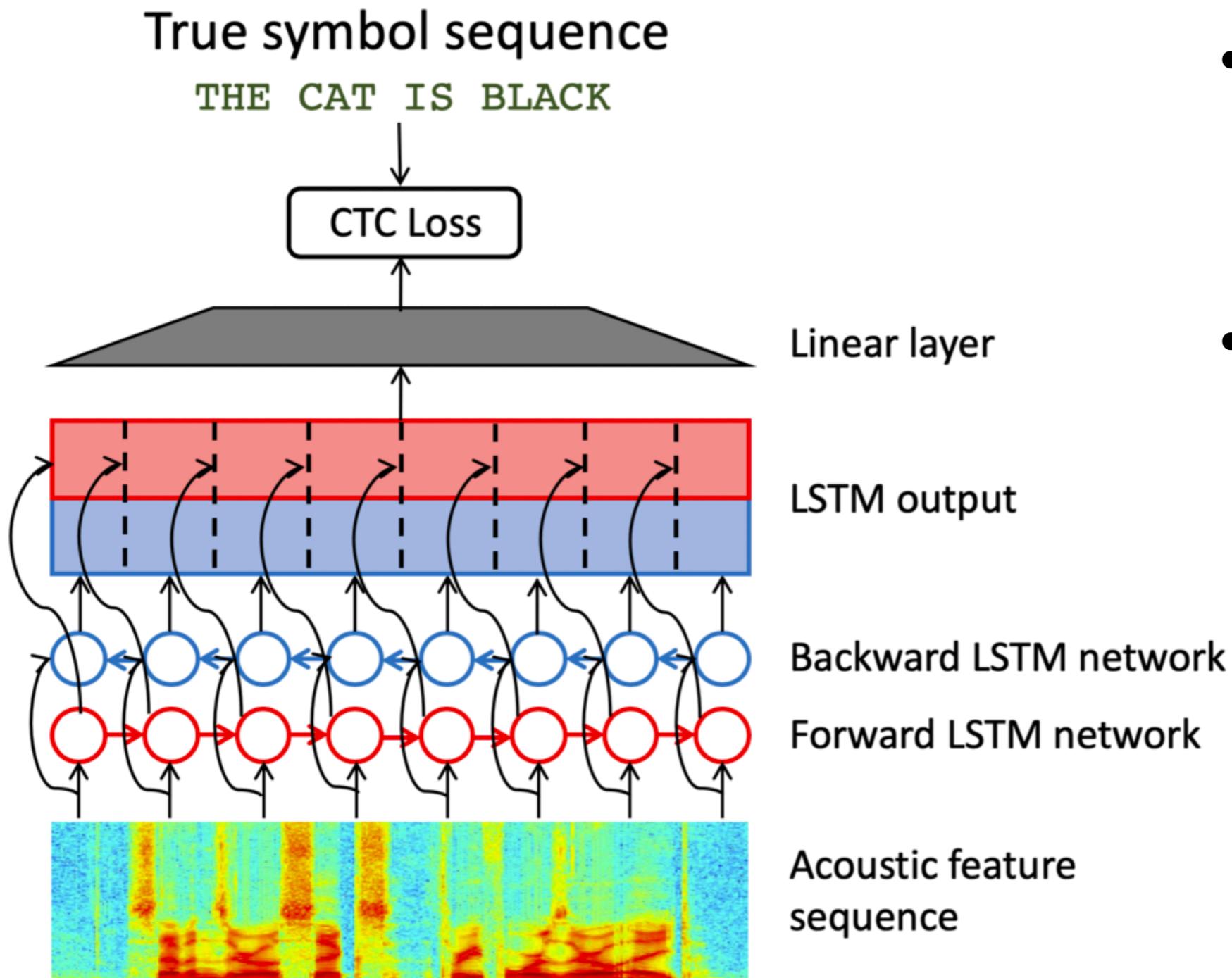


Image source: chainer.org

- Limitations on long sequences

Sequence to Sequence Learning
with Neural Networks

Speech Recognition



- Bi-directional LSTM gives extra context at each step
- Losses can go beyond sum of likelihood at each step e.g. CTC

Caption Generation

Caption Generation

