

CS 4803 / 7643: Deep Learning

Topics:

- Recurrent Neural Networks (RNNs)
 - Beam Search
 - LSTMs
- Reinforcement Learning

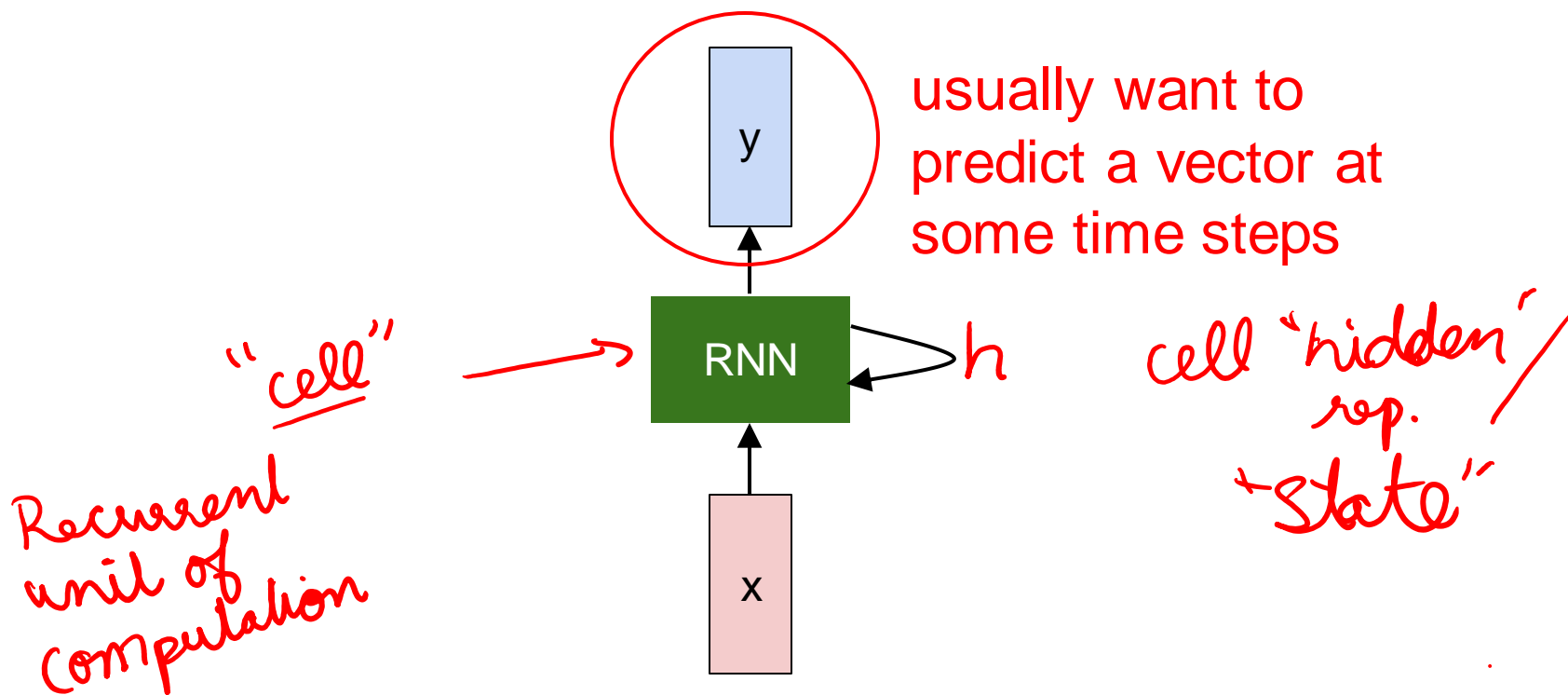
Dhruv Batra
Georgia Tech

Administrativa

- HW4
 - Due: 10/21, 11:59pm
 - RNNs, LSTMs, Image Captioning, Transformers!

Recap from last time

Recurrent Neural Network

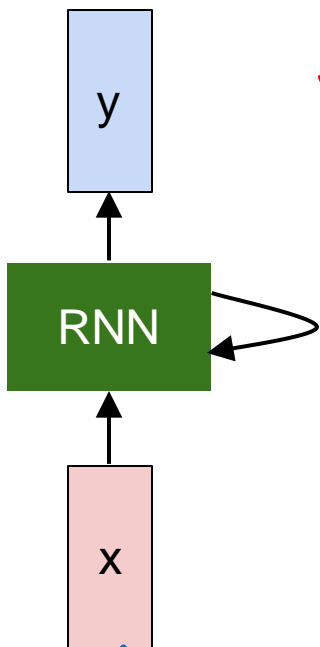


$$P(y_t | x_1 \dots x_t) \approx P(y_t | h_t)$$

(Vanilla) Recurrent Neural Network

The state consists of a single "hidden" vector h :

"Scores" $\vec{y}_t = \text{Softmax}(W_{hy} \vec{h}_t + \vec{b}_y)$



$h_t \in \mathbb{R}^{d_h}$
 $x \in \mathbb{R}^{d_x}$

$$\vec{h}_t = f_W(\vec{h}_{t-1}, \vec{x}_t)$$

$\frac{512}{\mathbb{R}} \Rightarrow \vec{h}_t = \tanh(W_{hh} \vec{h}_{t-1} + W_{xh} x_t + \vec{b}_h)$

$(\begin{bmatrix} W_{hh} & W_{xh} \end{bmatrix} \begin{bmatrix} \vec{h}_{t-1} \\ x_t \end{bmatrix} + \vec{b}_h)$

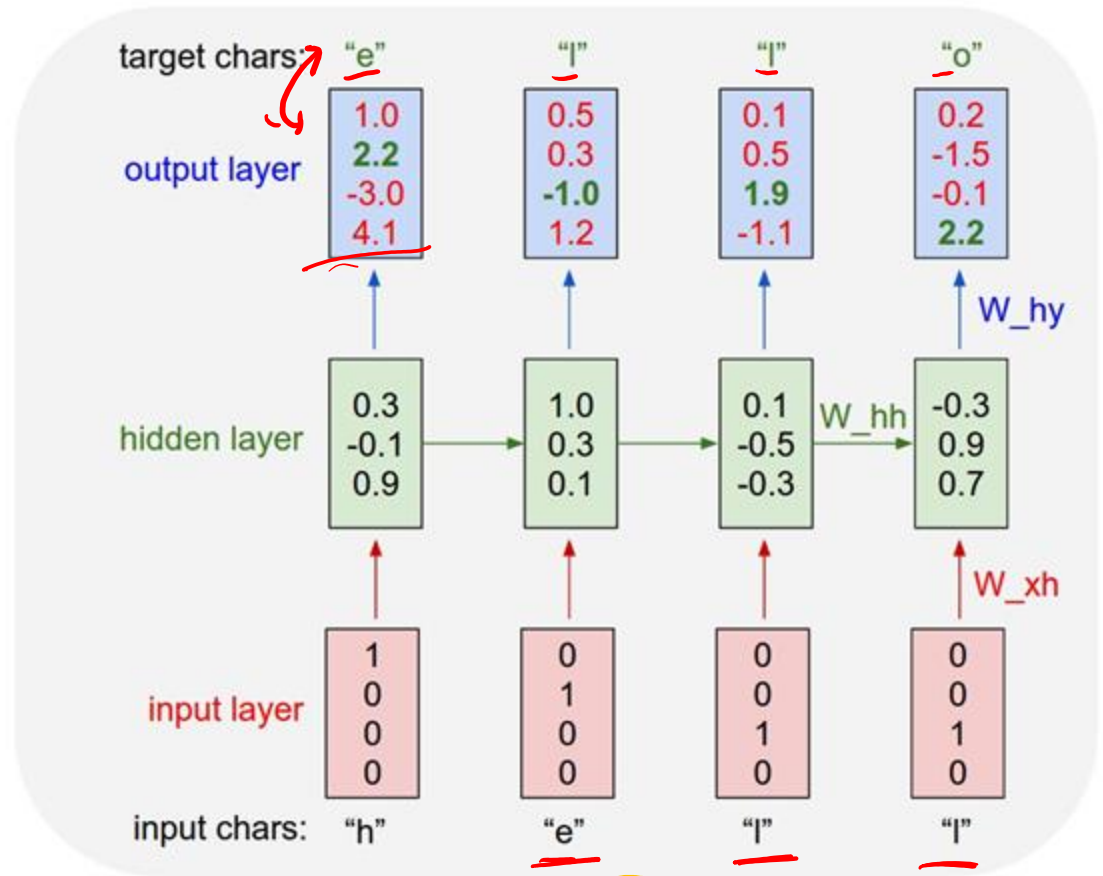
Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

$$\min_{\vec{w}} \sum_t -\log P(\underline{y}_t^{gt} | h_t, \vec{w}) = \mathcal{L}_1 \quad \mathcal{L}_2 \quad \mathcal{L}_3 \quad \mathcal{L}_4$$

Example:
Character-level
Language Model

Vocabulary:
 [h,e,l,o]

Example training
 sequence:
 "hello"



$$x_2 = y_1^{gt} \quad x_3 = y_2^{gt} \quad \dots$$

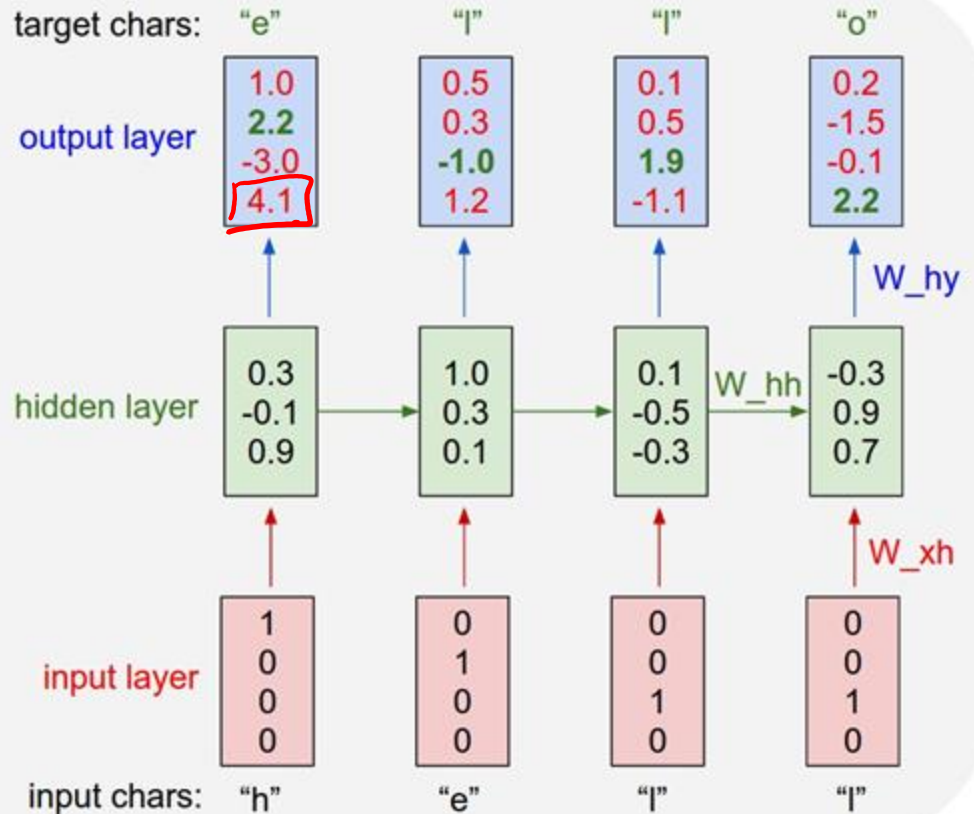
Training Time: MLE / “Teacher Forcing”

$$P(y_1 \dots y_T) = \prod_t P(y_t | y_{1:t-1})$$

Example:
Character-level
Language Model

Vocabulary:
 [h,e,l,o]

Example training
 sequence:
 “hello”



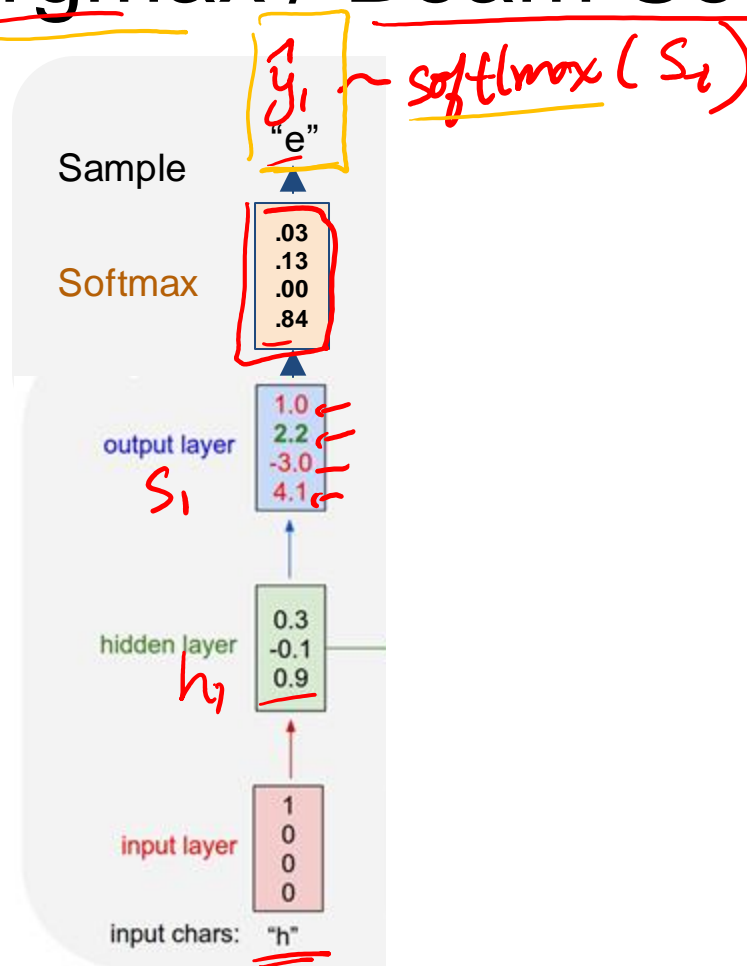
$$x_2 = y_1^{36}$$

Test Time: Sample / Argmax / Beam Search

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a
time, feed back to
model

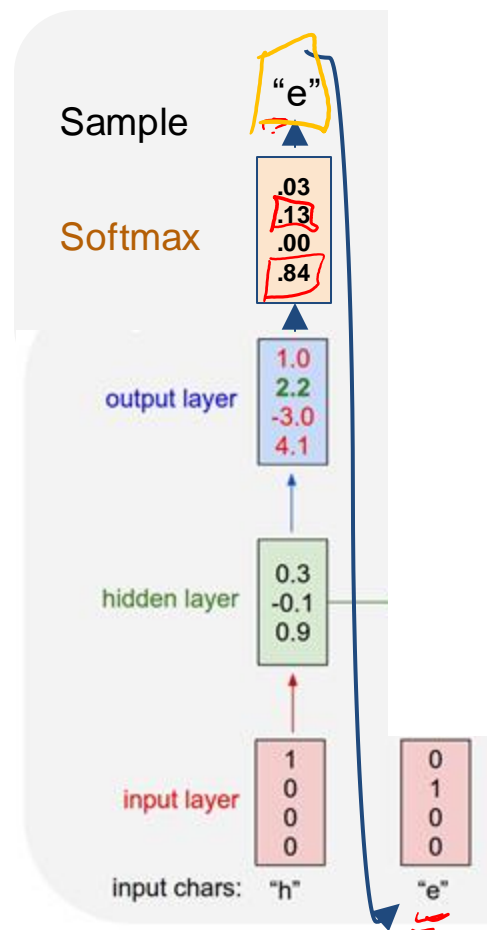


Test Time: Sample / Argmax / Beam Search

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a
time, feed back to
model

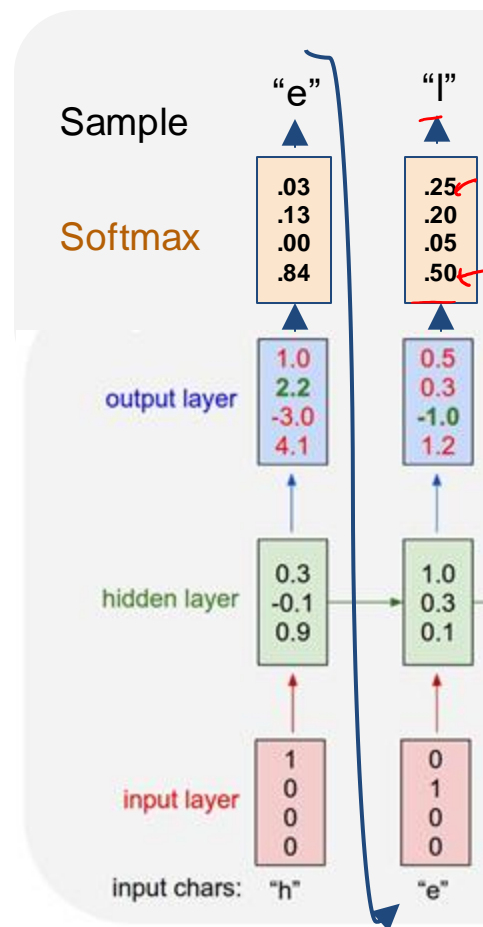


Test Time: Sample / Argmax / Beam Search

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a
time, feed back to
model

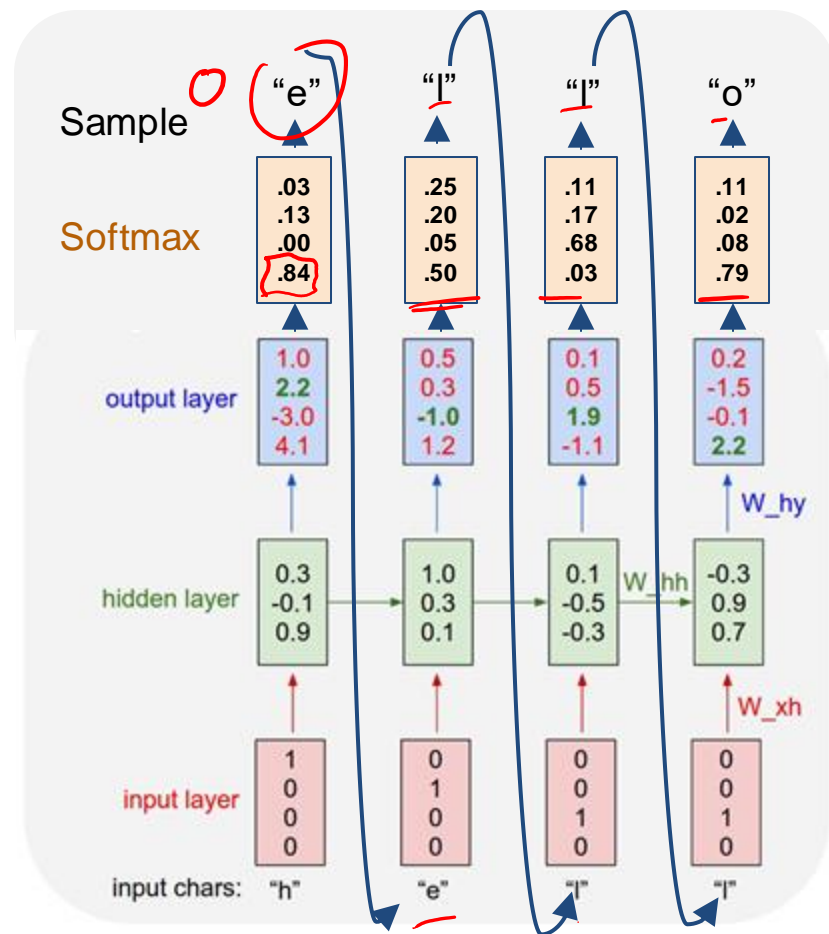


Test Time: Sample / Argmax / Beam Search

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a
time, feed back to
model



h o

Plan for Today

- Recurrent Neural Networks (RNNs)

- Inference: Beam Search

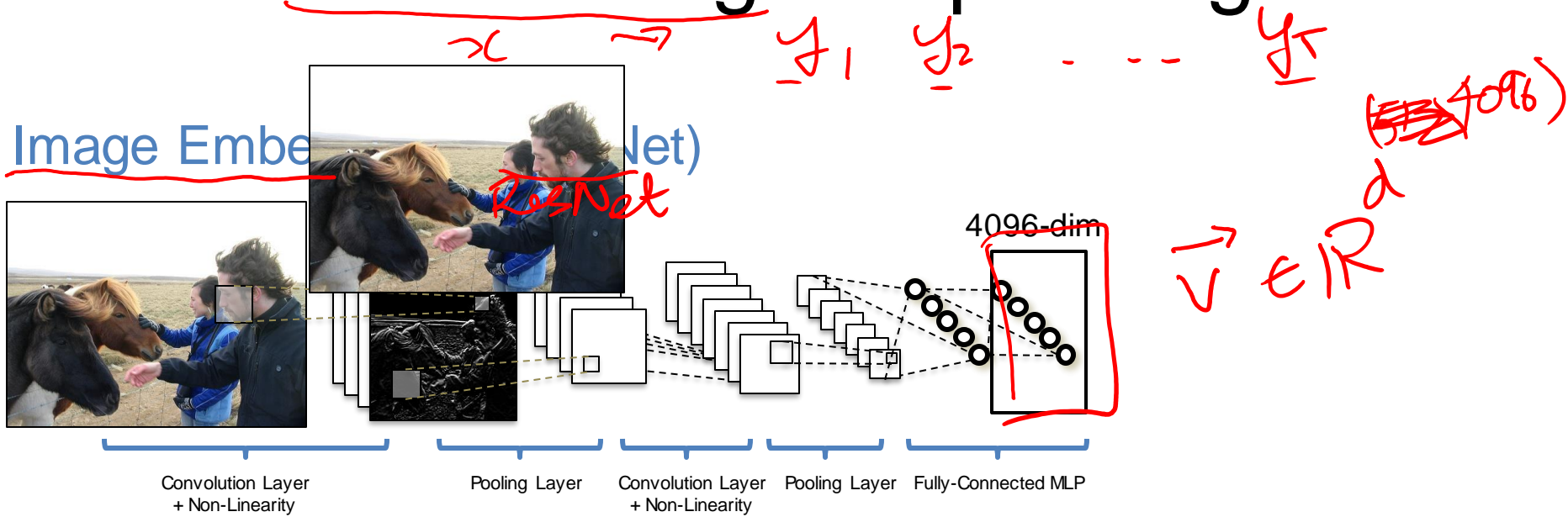
- Example: Image Captioning

- Multilayer RNNs

- Problems with gradients in “vanilla” RNNs

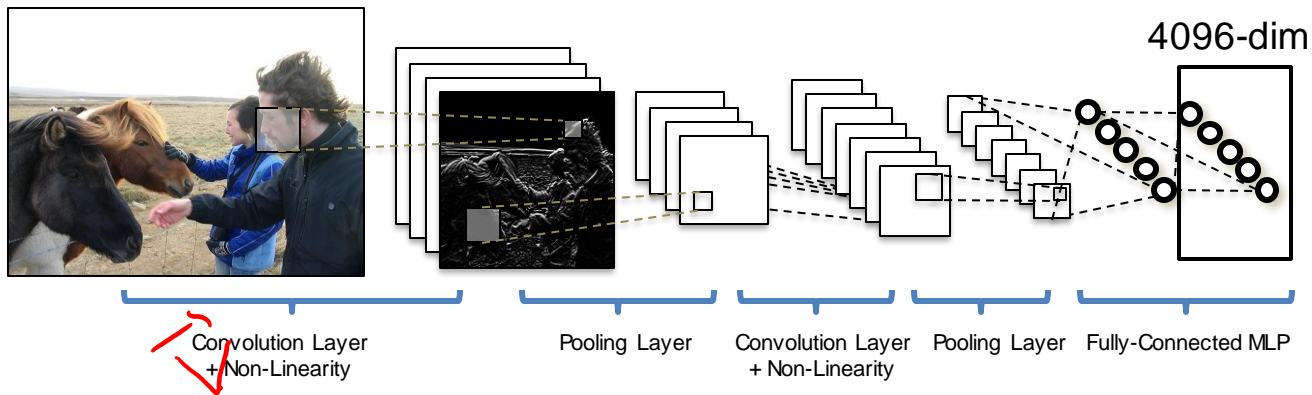
- LSTMs (and other RNN variants)

Neural Image Captioning

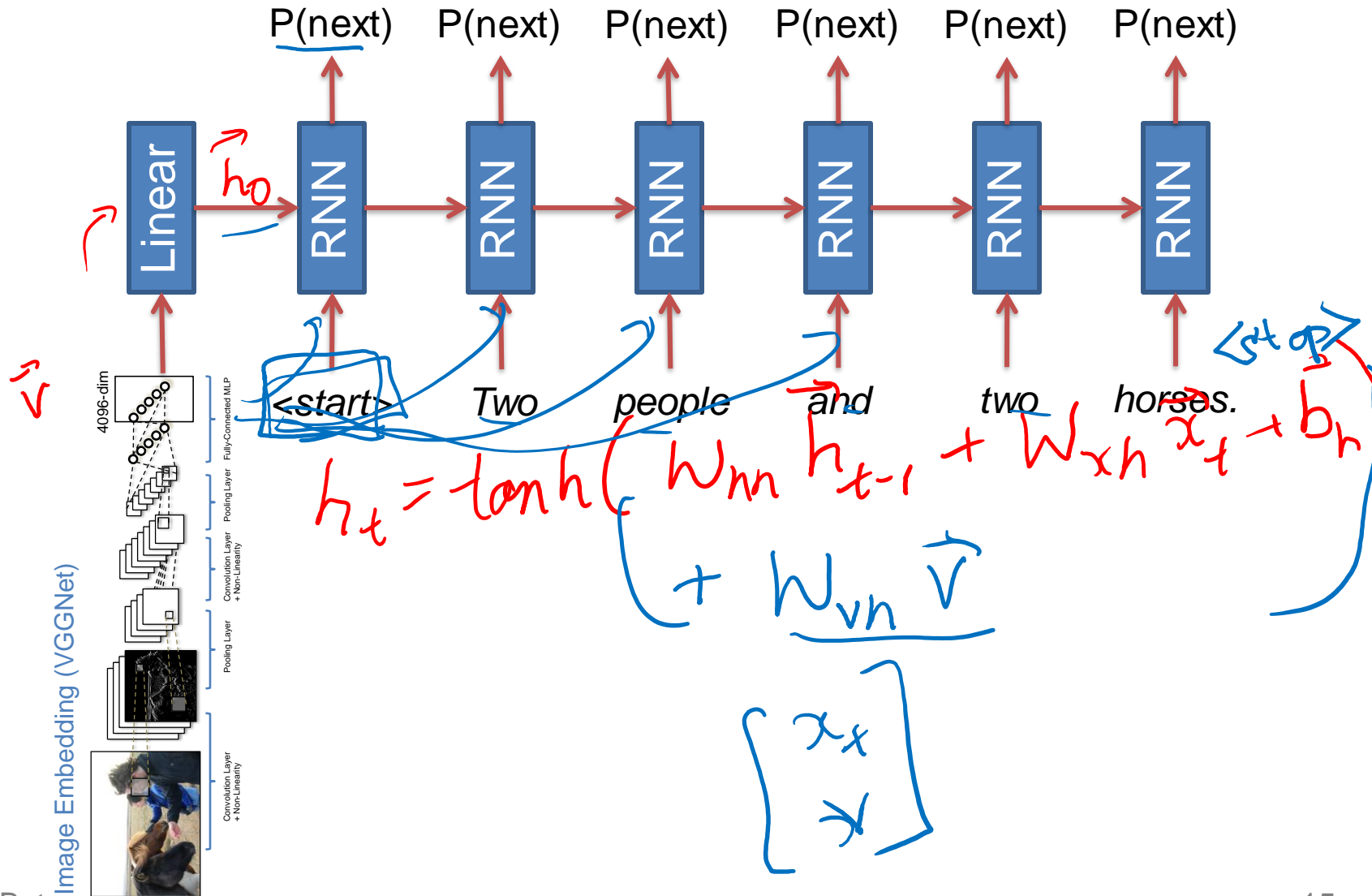


Neural Image Captioning

Image Embedding (VGGNet)



Neural Image Captioning



Beam Search Demo

- <http://dbs.cloudcv.org/captioning&mode=interactive>

Image Captioning: Example Results

Captions generated using
[neuraltalk2](#)
All images are [CC0 Public domain](#)
[cat](#) [suitcase](#) [cat tree](#) [dog](#) [bear](#)
[surfers](#) [tennis](#) [giraffe](#) [motorcycle](#)



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



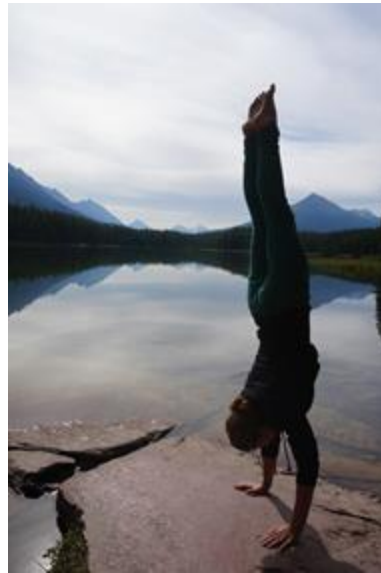
A man riding a dirt bike on a dirt track

Image Captioning: Failure Cases

Captions generated using [neuraltalk2](#).
All images are [CC0 Public domain](#): [fur coat](#), [handstand](#), [spider web](#), [baseball](#)



A woman is holding a cat in her hand



A woman standing on a beach holding a surfboard



A bird is perched on a tree branch



A person holding a computer mouse on a desk



A man in a baseball uniform throwing a ball

Plan for Today

- Recurrent Neural Networks (RNNs)
 - Inference: Beam Search
 - Example: Image Captioning
 - Multilayer RNNs
 - Problems with gradients in “vanilla” RNNs
 - LSTMs (and other RNN variants)

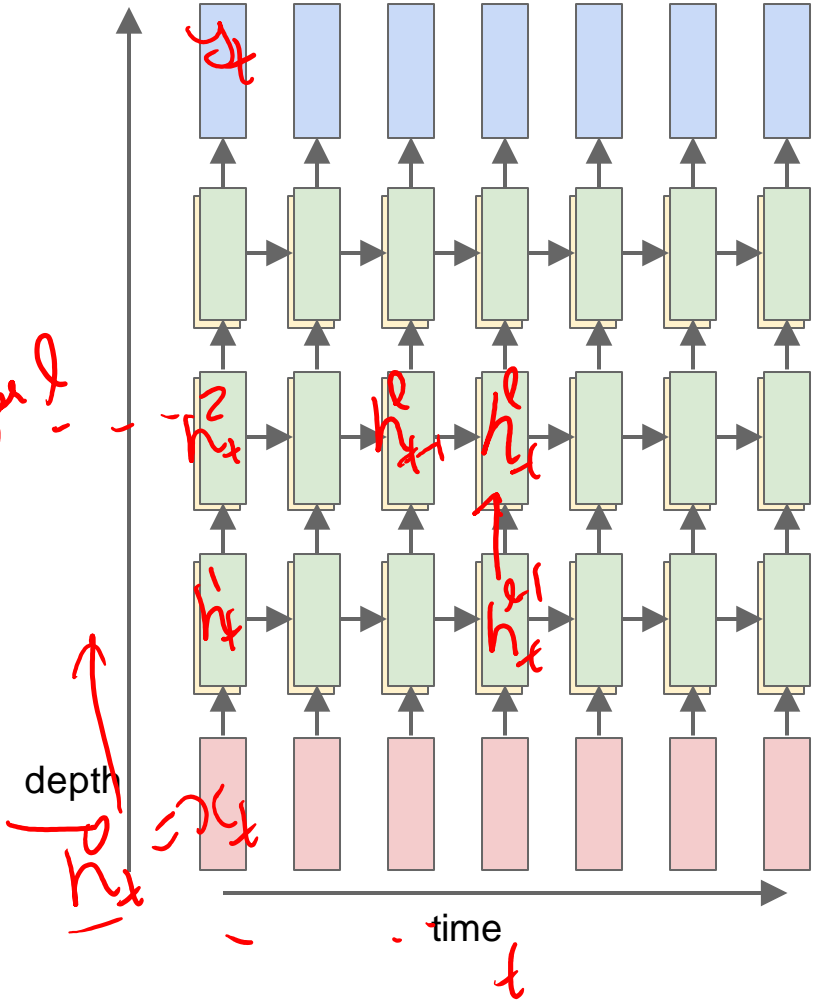
Multilayer RNNs

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix} + b$$

$h \in \mathbb{R}^n$ $W^l [n \times 2n]$

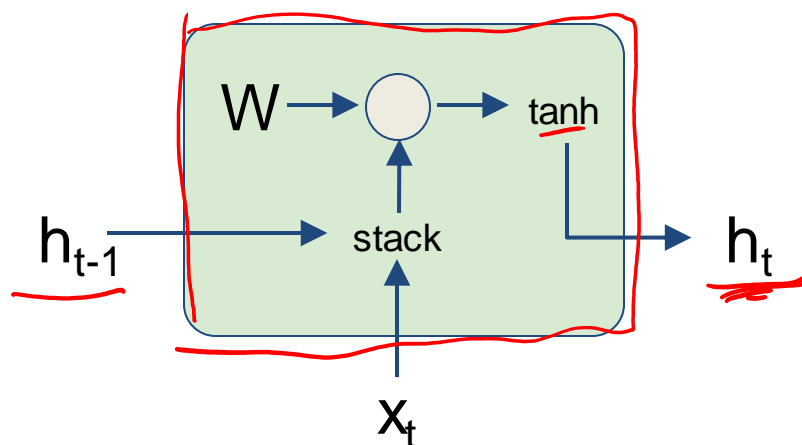
h_t^l

layer 1



Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

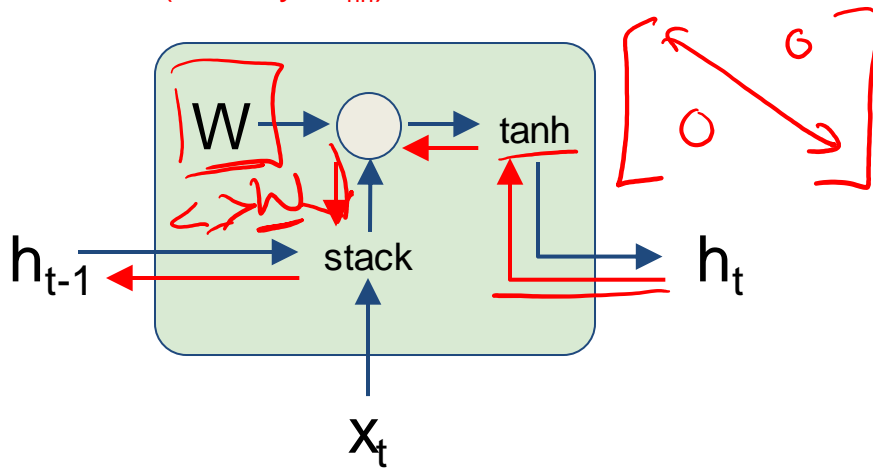


$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(\underline{W} \begin{pmatrix} h_{t-1} \\ \underline{x_t} \end{pmatrix}\right) \end{aligned}$$

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

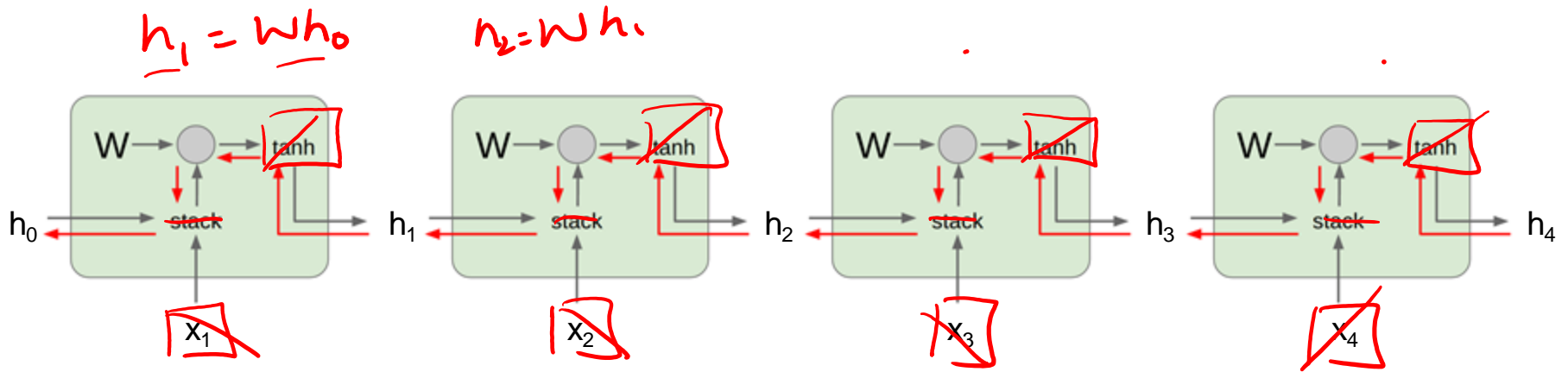
Backpropagation from h_t
to h_{t-1} multiplies by W
(actually W_{hh})



$$\begin{aligned}h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)\end{aligned}$$

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h_0 involves many factors of W (and repeated \tanh)

$$\frac{\partial h_T}{\partial h_0} = \left[\frac{\partial h_T}{\partial h_{T-1}} \right] \dots \frac{\partial h_1}{\partial h_0}$$

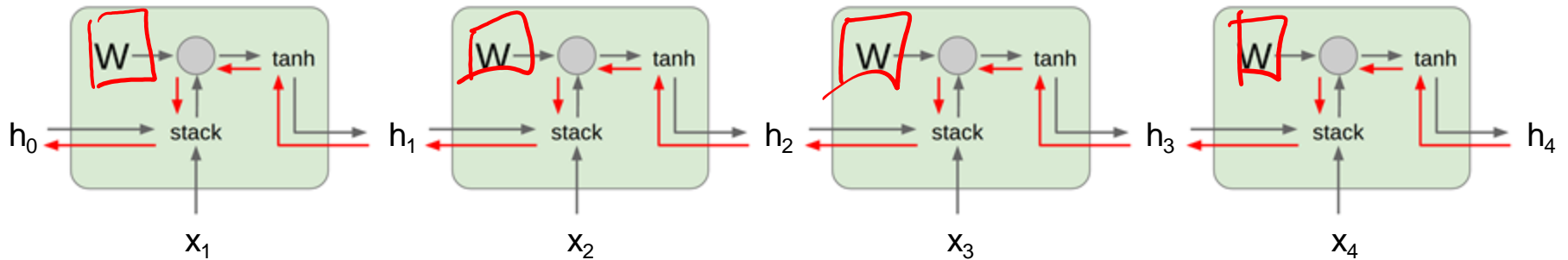
$$= W \cdot W \cdot \dots \cdot W$$

$$= W^{(T+1)}$$

If $h_t \in \mathbb{R}^1$
 $W \in \mathbb{R}^1$
 $(W)^{T+1}$

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



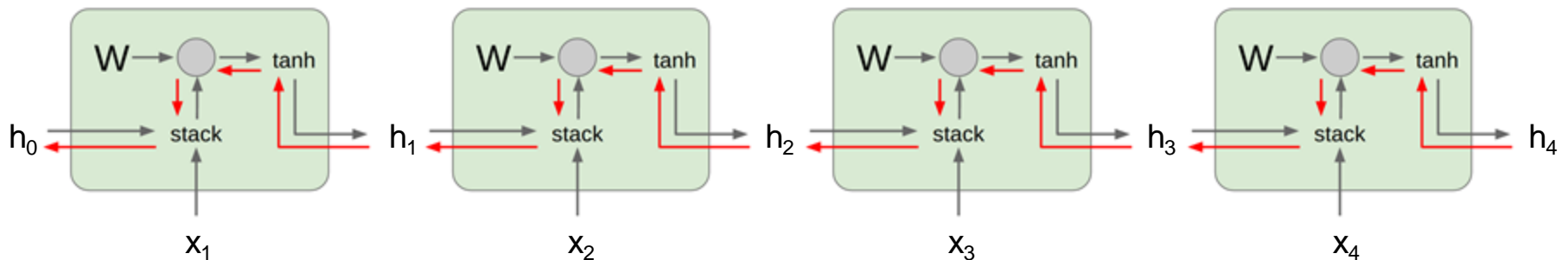
Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

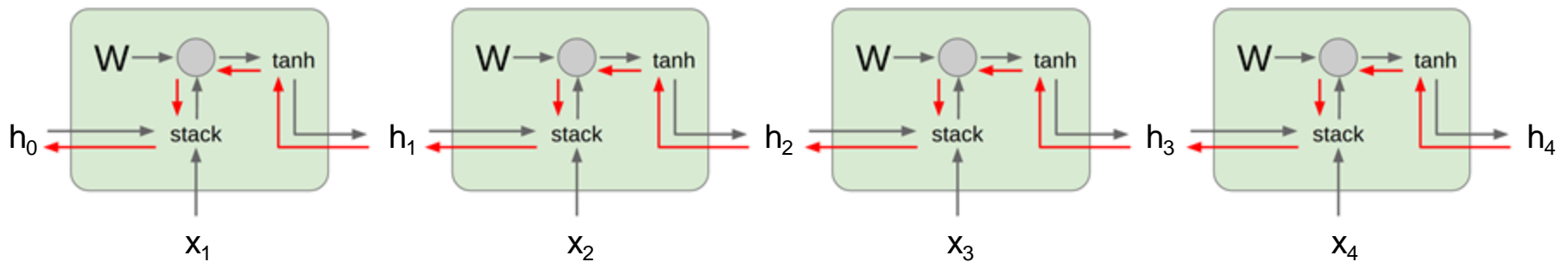
Largest singular value < 1 :
Vanishing gradients

Gradient clipping: Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

→ Change RNN architecture

Long Short Term Memory (LSTM)

Vanilla RNN

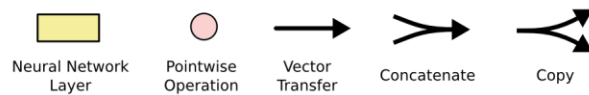
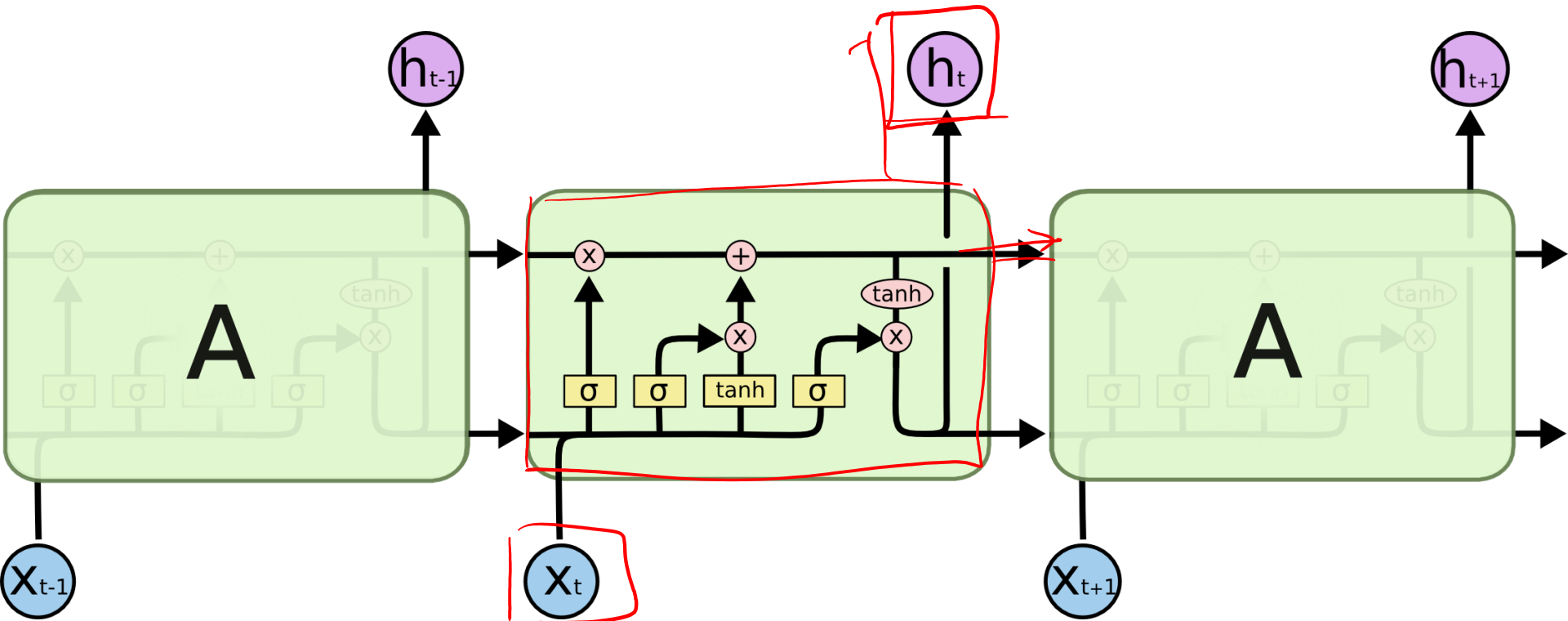
$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

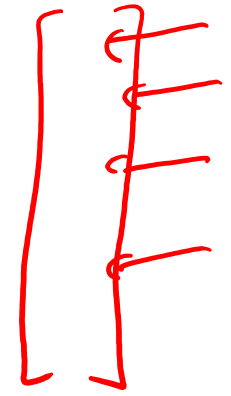
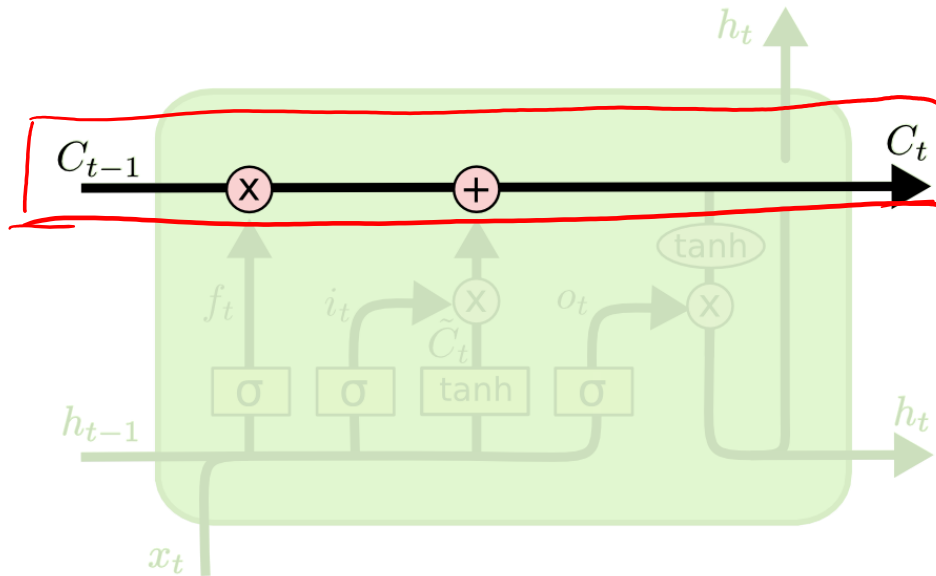
Meet LSTMs



LSTMs Intuition: Memory

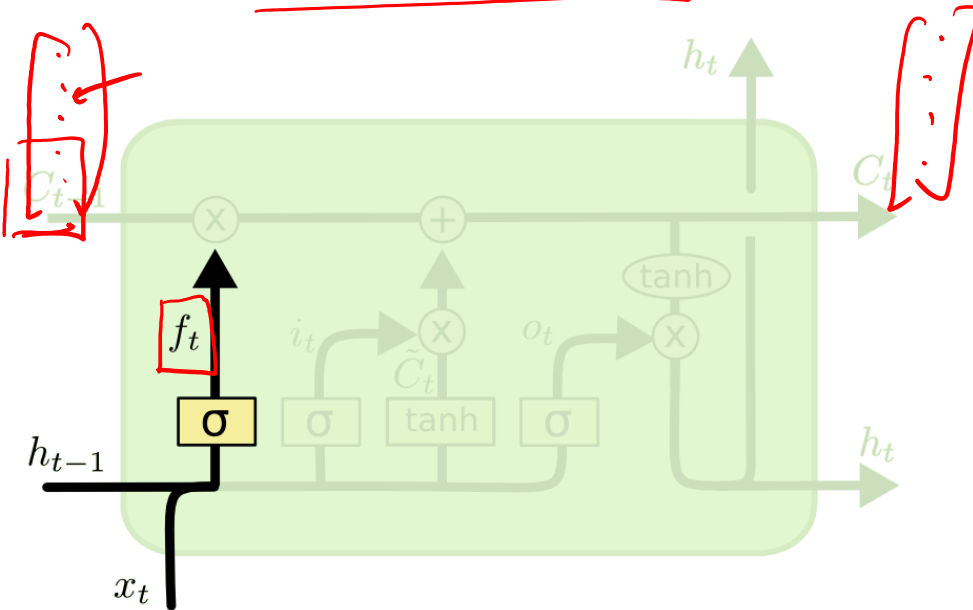
- Cell State / Memory

C_t : h_t hidden state
 $C_t \in \mathbb{R}^d$
 $h_t \in \mathbb{R}^d$



LSTMs Intuition: Forget Gate

- Should we continue to remember this “bit” of information or not?



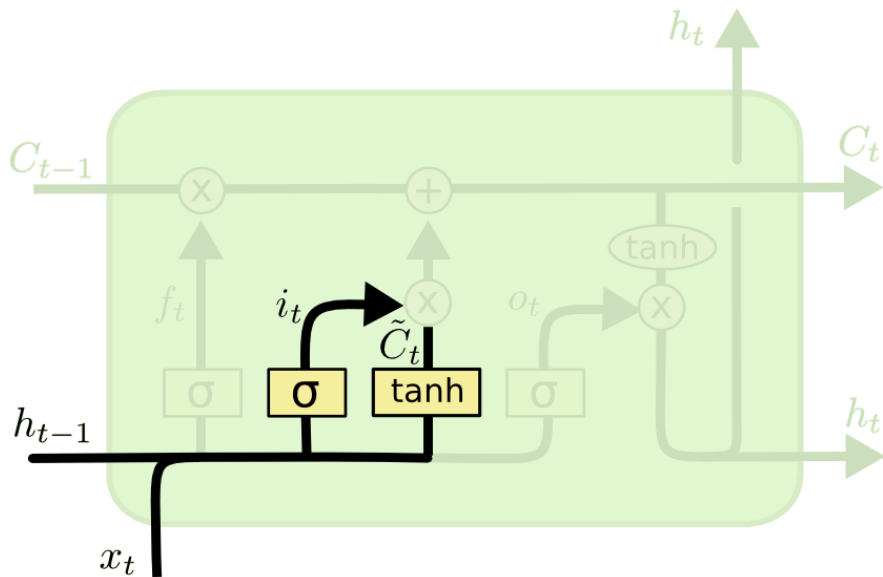
$$f_t \in \mathbb{R}^d \quad \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$\underline{f}_t = \underline{\sigma}(\underline{W}_f \cdot \underline{[h_{t-1}, x_t]} + \underline{b}_f)$$

LSTMs Intuition: Input Gate

- Should we update this “bit” of information or not?
 - If so, with what?

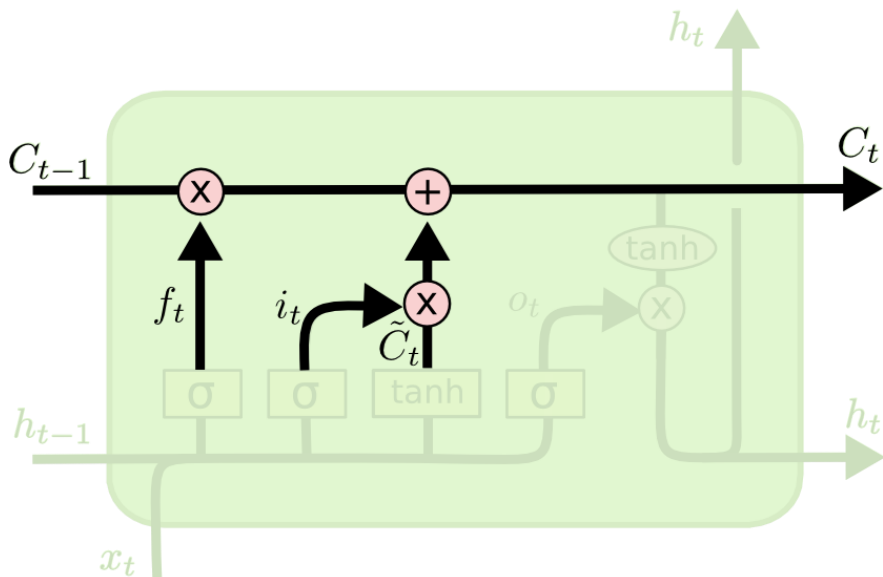
$$i_t \in \mathbb{R}^d \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$



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

LSTMs Intuition: Memory Update

- Forget that + memorize this

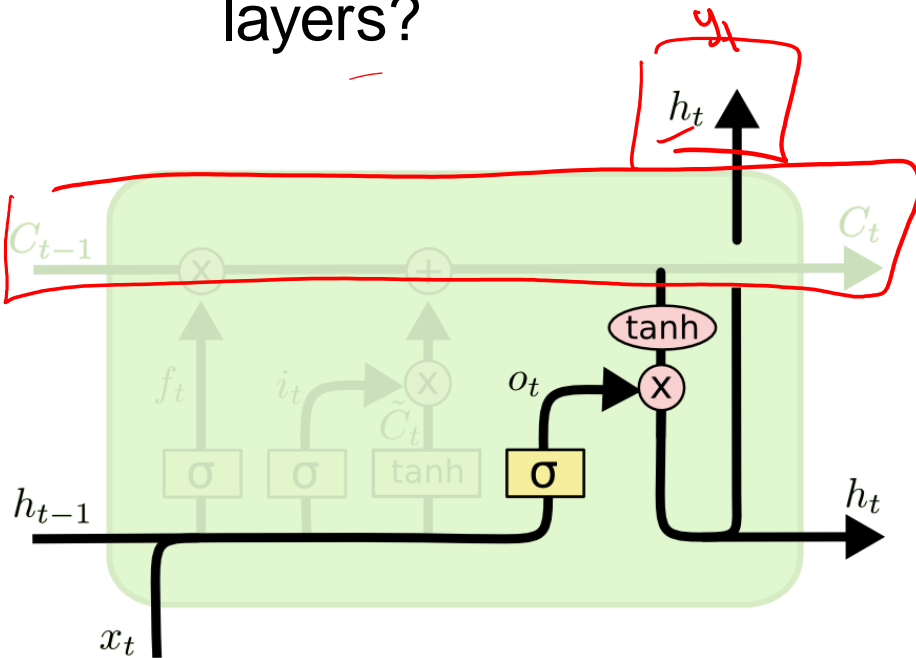


$$C_t = \boxed{f_t} * \underline{C_{t-1}} + \boxed{i_t} * \tilde{C}_t$$

Handwritten red annotations showing matrix representations of the forget and input gates. The forget gate is represented as a column vector $\begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix}$ and the input gate as a column vector $\begin{bmatrix} i_1 \\ i_2 \\ \vdots \\ i_n \end{bmatrix}$. The word "forget" is written below the forget gate vector.

LSTMs Intuition: Output Gate

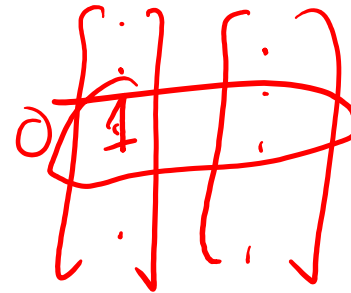
- Should we output this “bit” of information to “deeper” layers?



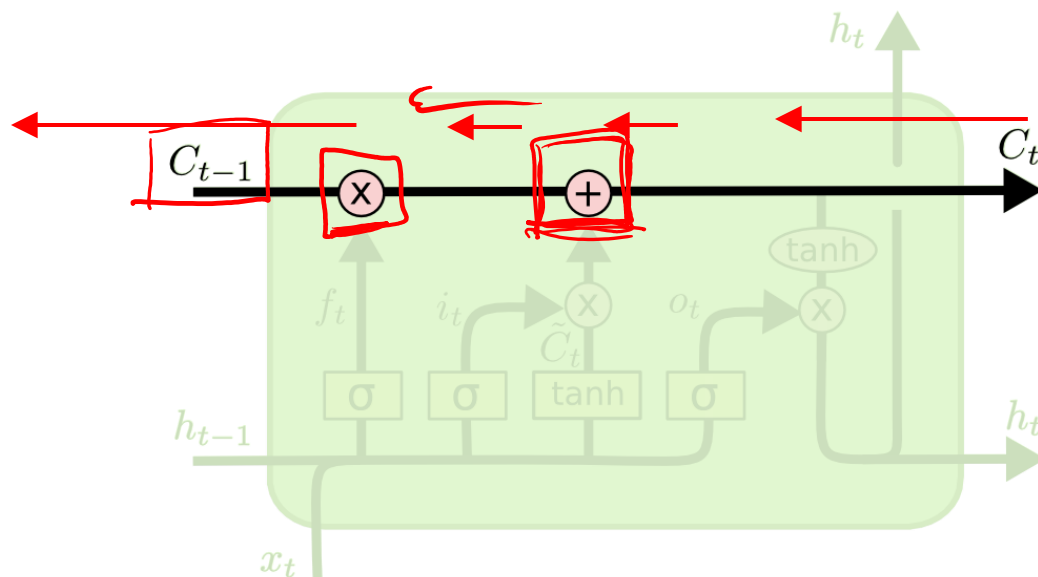
$$o_t \in \mathbb{R}^d$$

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

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

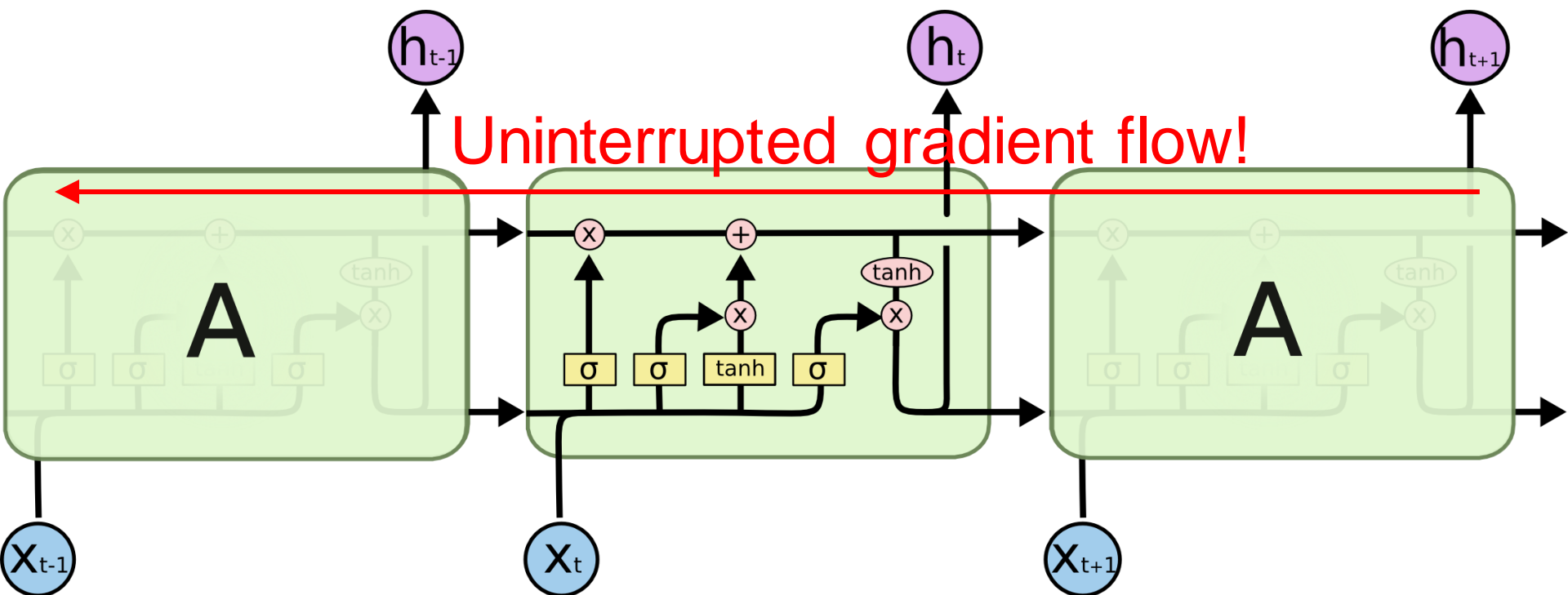


LSTMs Intuition: Additive Updates



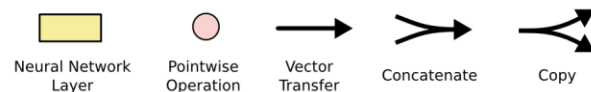
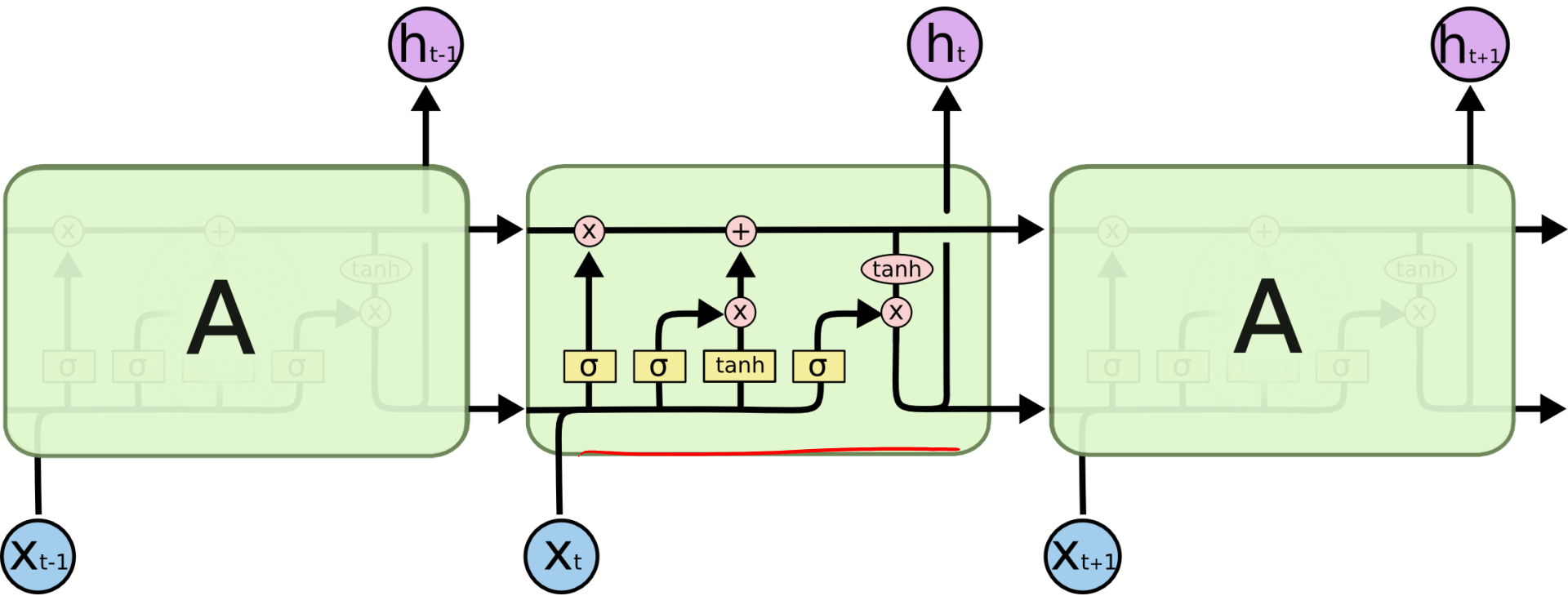
Backpropagation from c_t to c_{t-1} only
elementwise
multiplication by f , no
matrix multiply by W

LSTMs Intuition: Additive Updates



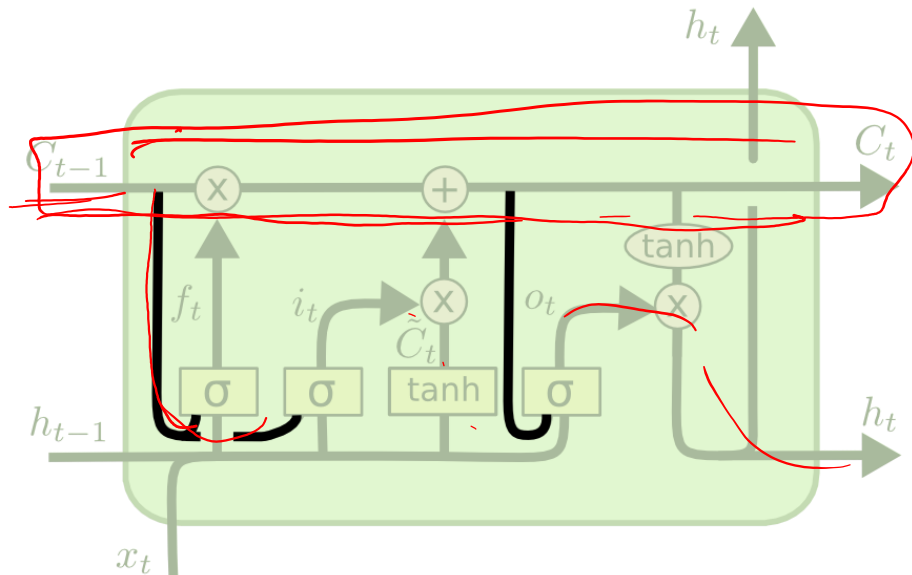
LSTMs

- A pretty sophisticated cell



LSTM Variants #1: Peephole Connections

- Let gates see the cell state / memory



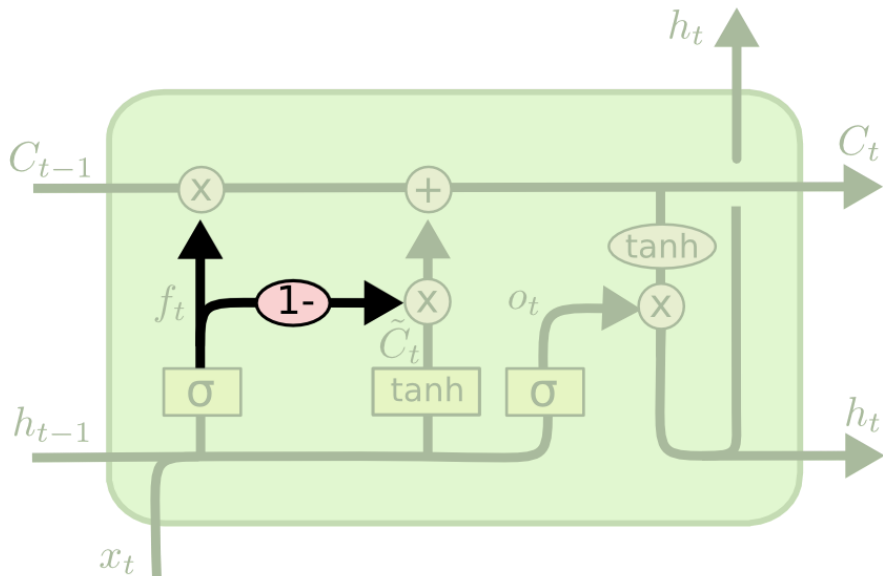
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

LSTM Variants #2: Coupled Gates

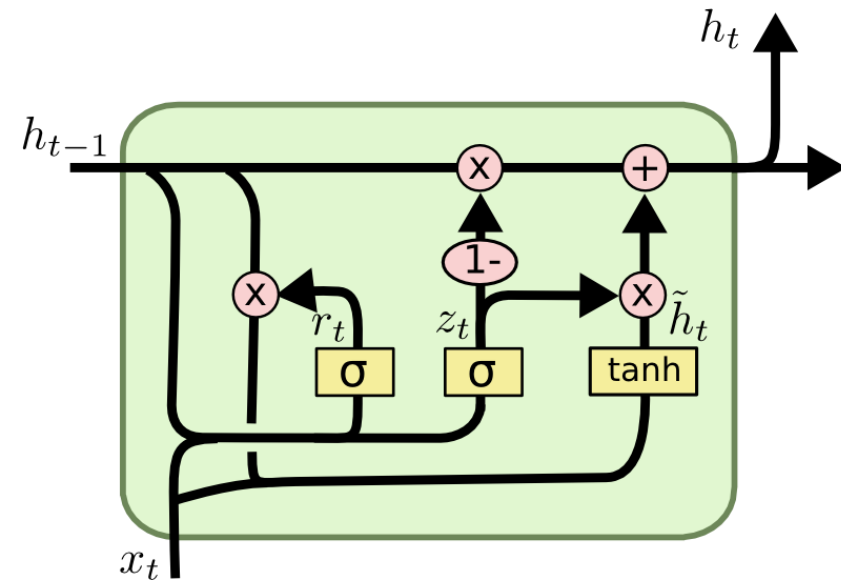
- Only memorize new if forgetting old



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

LSTM Variants #3: Gated Recurrent Units

- Changes:
 - No explicit memory; memory = hidden output
 - Z = memorize new and forget old



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

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

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

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

Other RNN Variants

[*An Empirical Exploration of Recurrent Network Architectures*,
Jozefowicz et al., 2015]

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz} \tanh(h_t) + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

Summary

- RNNs allow a lot of flexibility in architecture design
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better/simpler architectures are a hot topic of current research
- Better understanding (both theoretical and empirical) is needed.

Administrativa

- Guest Lecture: Nirbhay Modhe
 - Reinforcement Learning



<https://arjunmajum.github.io/>