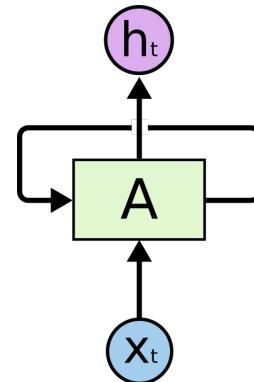


Recurrent Neural networks



Slides inspired by Abhishek Narwekar and Anusri Pampari at UIUC



Lecture Outline

1. Introduction
2. Learning Long Term Dependencies
3. Visualization for RNNs

Applications of RNNs

A person riding a motorcycle on a dirt road.



Image Captioning
[\[reference\]](#)
... and more!

VIOLA:

Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;

Write like Shakespeare
[\[reference\]](#)

In reply to Thomas Paine
 DeepDrumpf @DeepDrumpf · Mar 20
There will be no amnesty. It is going to pass because the people are going to be gone. I'm giving a mandate. #ComeyHearing @Thomas1774Paine

1 12 17

.. and Trump
[\[reference\]](#)

Applications of RNNs

Technically, an RNN models sequences

- Time series
- Natural Language, Speech

We can even convert non-sequences to sequences, eg: feed an image as a sequence of pixels!

Applications of RNNs

- RNN Generated TED Talks
 - [YouTube Link](#)
- RNN Generated Eminem rapper
 - [RNN Shady](#)
- RNN Generated Music
 - [Music Link](#)

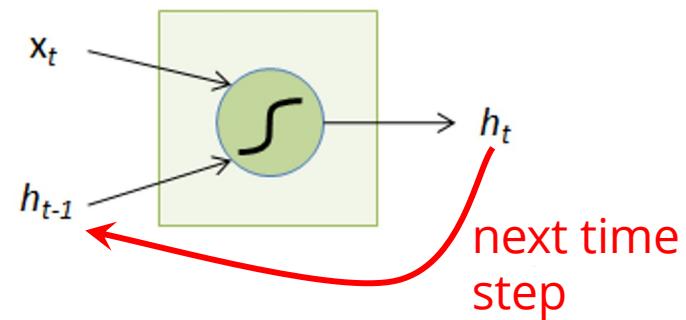
Why RNNs?

- Can model sequences having variable length
- **Efficient:** Weights shared across time-steps
- **They work!**
 - SOTA in several speech, NLP tasks
 - Not any more, now Transformers

(See for image stuffs, <https://openai.com/blog/dall-e/>)

The Recurrent Neuron

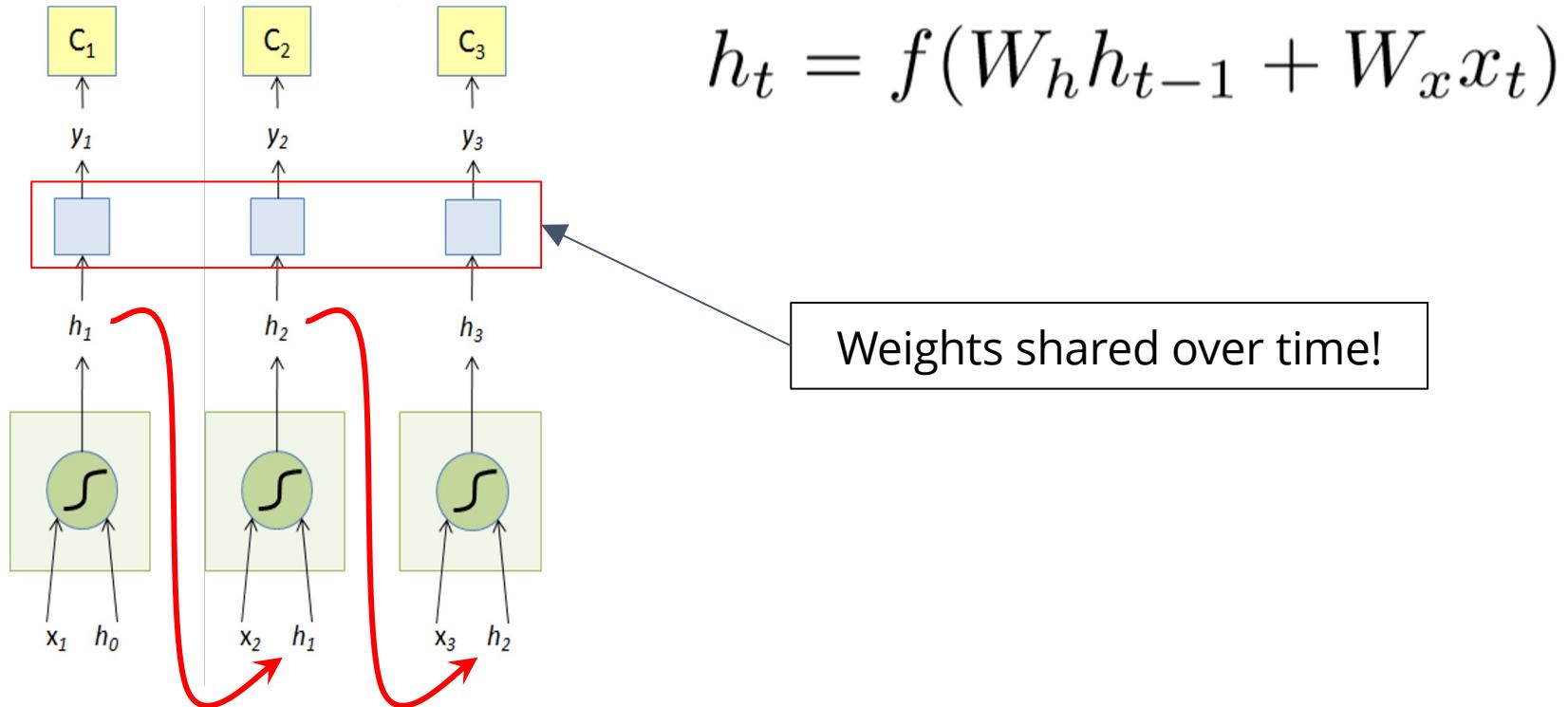
- x_t : Input at time t
- h_{t-1} : State at time t-1



$$h_t = f(W_h h_{t-1} + W_x x_t)$$

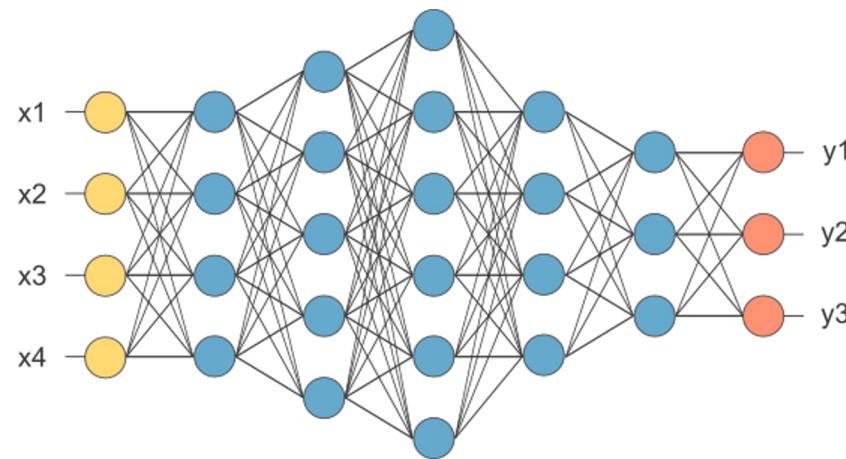
Source: [Slides by Arun](#)

Unfolding an RNN



Source: [Slides by Arun](#)

Making Feedforward Neural Networks Deep



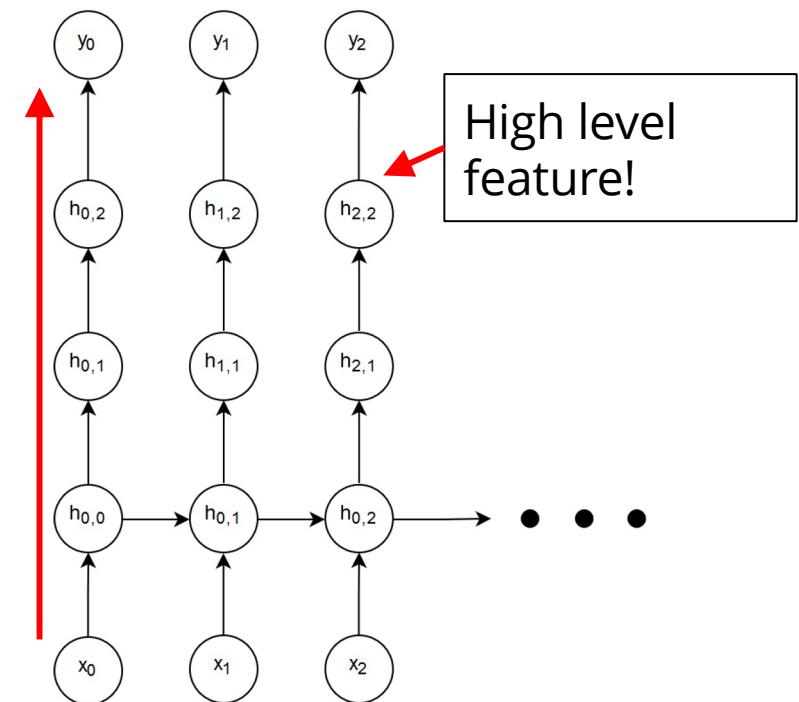
Source: http://www.opennn.net/images/deep_neural_network.png

Option 1: Feedforward Depth (d_f)

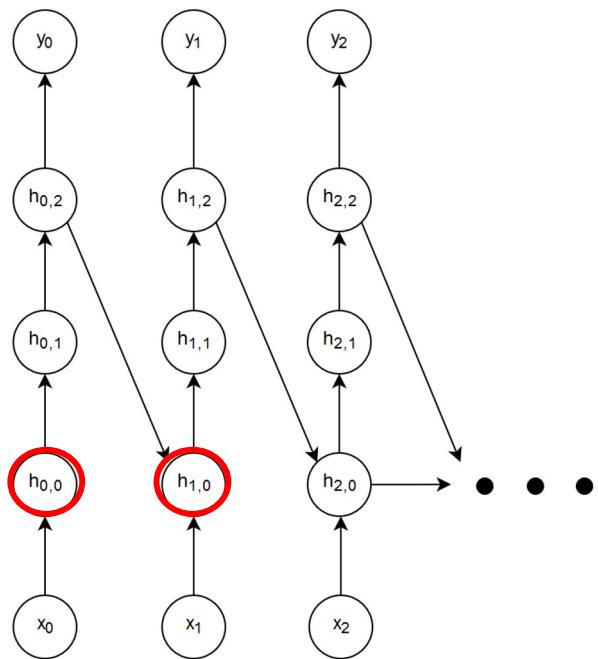
Notation: $h_{0,1} \Rightarrow$ time step 0, neuron #1

- **Feedforward depth:** longest path between an input and output at the same timestep

Feedforward depth = 4



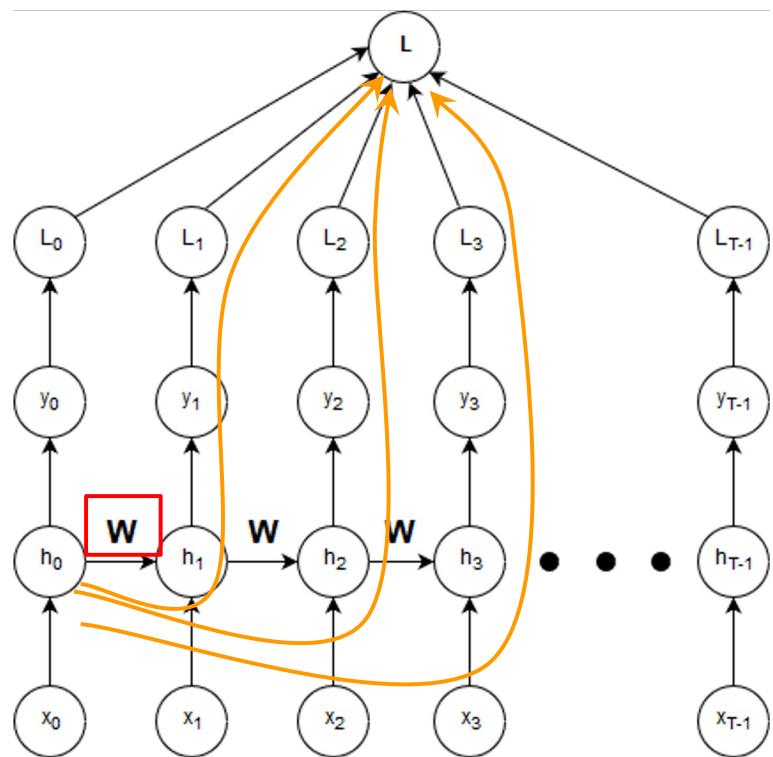
Option 2: Recurrent Depth (d_r)



- **Recurrent depth:** Longest path between **same hidden state** in **successive timesteps**

Recurrent depth = 3

Backpropagation Through Time (BPTT)



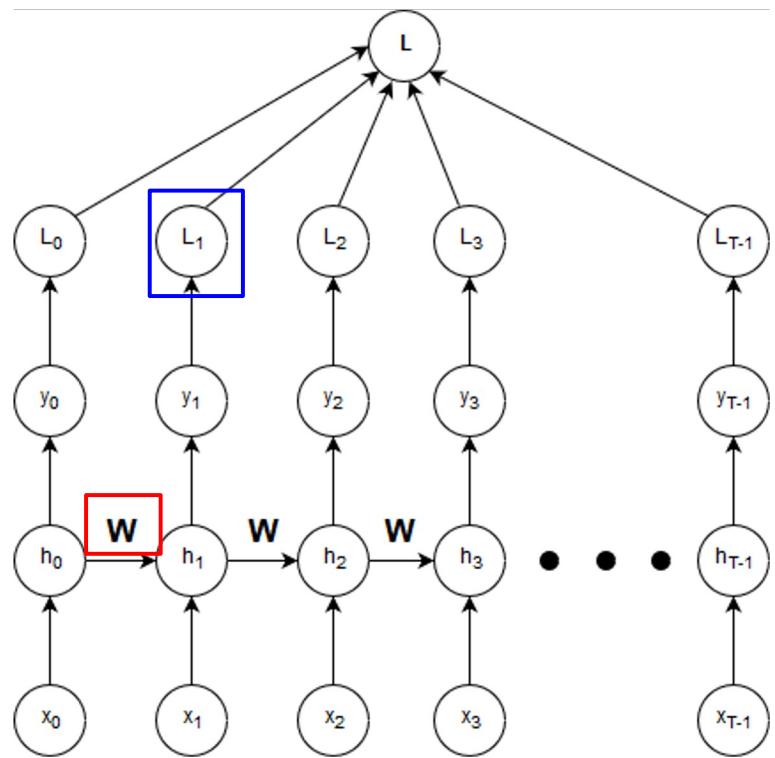
- Objective is to update the weight matrix:

$$\mathbf{W} \rightarrow \mathbf{W} - \alpha \frac{\partial L}{\partial \mathbf{W}}$$

- Issue: \mathbf{W} occurs each timestep
- **Every** path from \mathbf{W} to L is one dependency
- Find all paths from \mathbf{W} to L !

(note: dropping subscript h from \mathbf{W}_h for brevity)

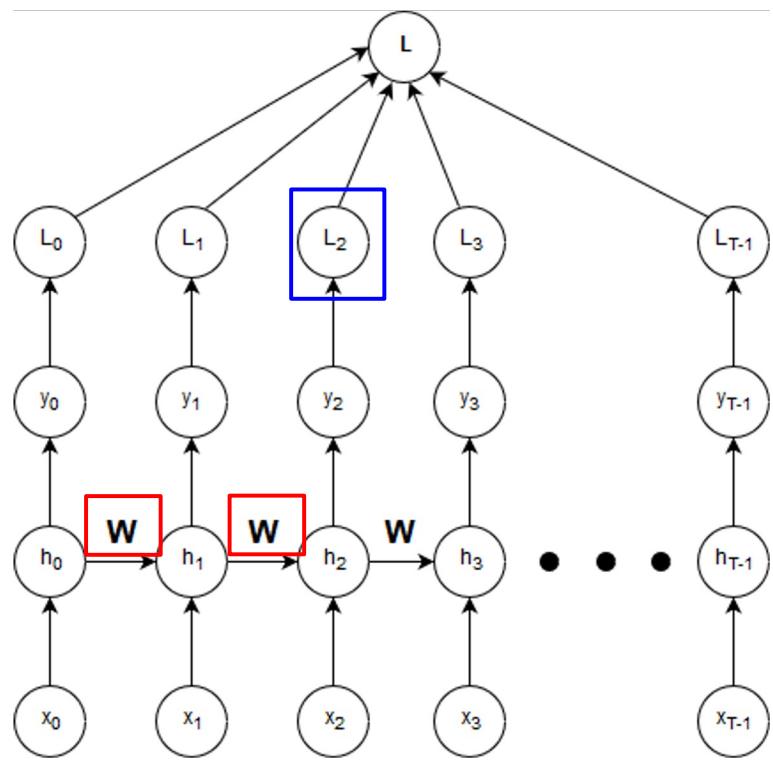
Systematically Finding All Paths



How many paths exist from W to L *through* L_1 ?

Just 1. Originating at h_0 .

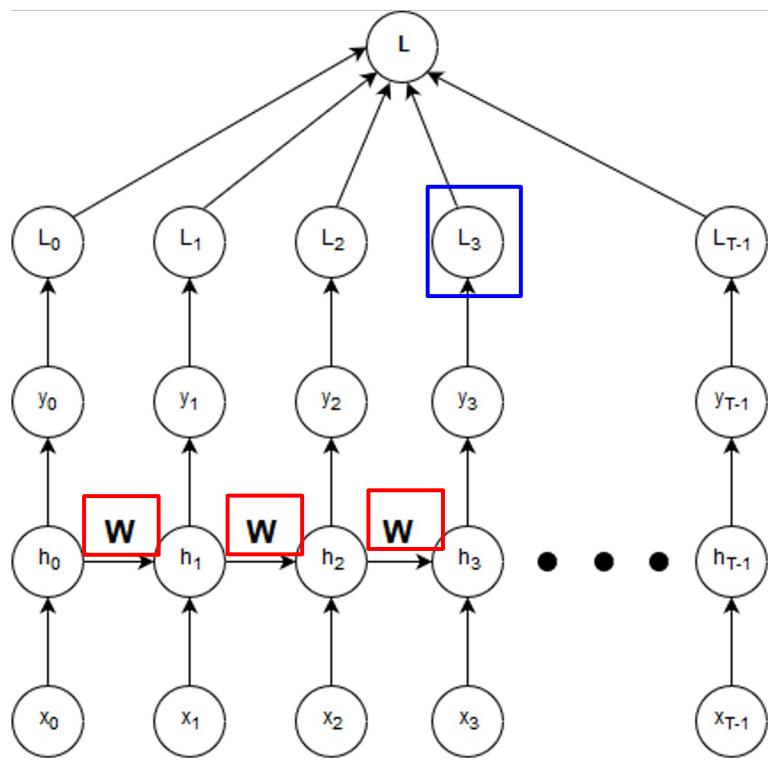
Systematically Finding All Paths



How many paths from **W** to **L** through L_2 ?

2. Originating at h_0 and h_1 .

Systematically Finding All Paths



And 3 in this case.

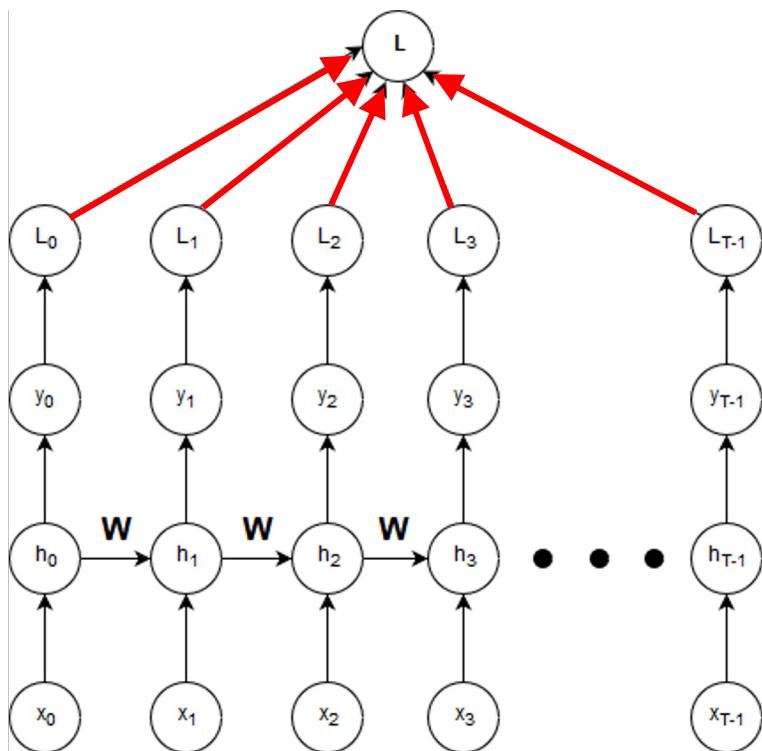
Origin of path = basis for Σ

$$\frac{\partial L}{\partial \mathbf{W}}$$

The gradient has two summations:

- 1: Over L_j
- 2: Over h_k

Backpropagation as two summations

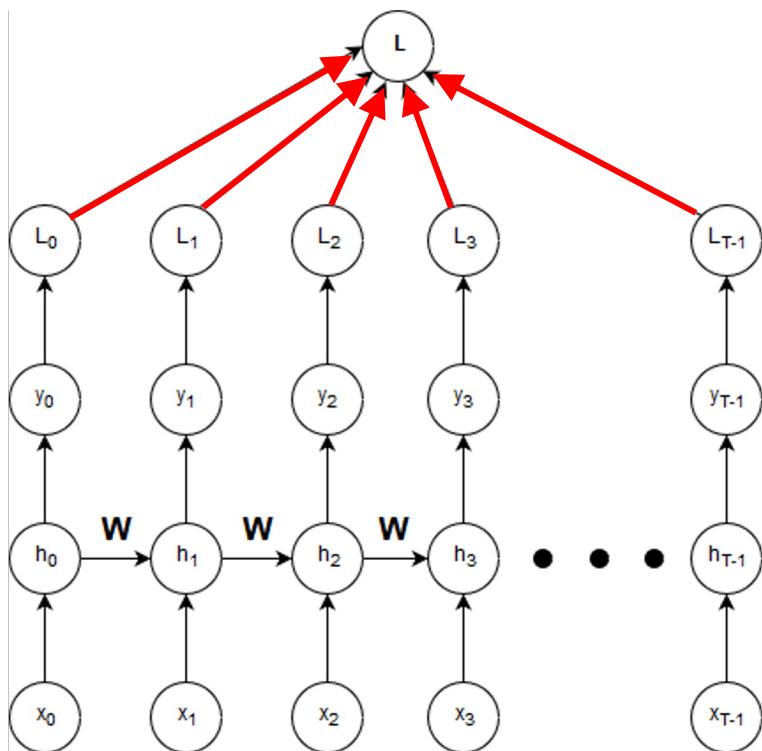


- First summation over L

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_{j=0}^{T-1} \frac{\partial L_j}{\partial \mathbf{W}} \frac{\partial L}{\partial L_j}$$

$$(L = L_1 + L_2 + \dots + L_{T-1})$$

Backpropagation as two summations

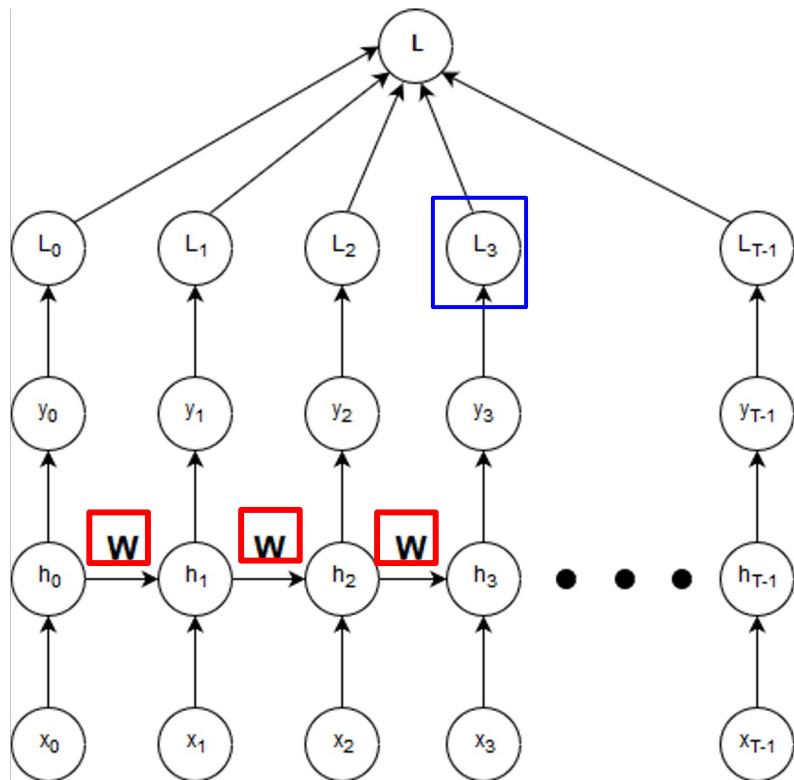


- First summation over L

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_{j=0}^{T-1} \frac{\partial L_j}{\partial \mathbf{W}} \cancel{\frac{\partial L}{\partial L_j}}$$

$$(L = L_1 + L_2 + \dots + L_{T-1})$$

Backpropagation as two summations

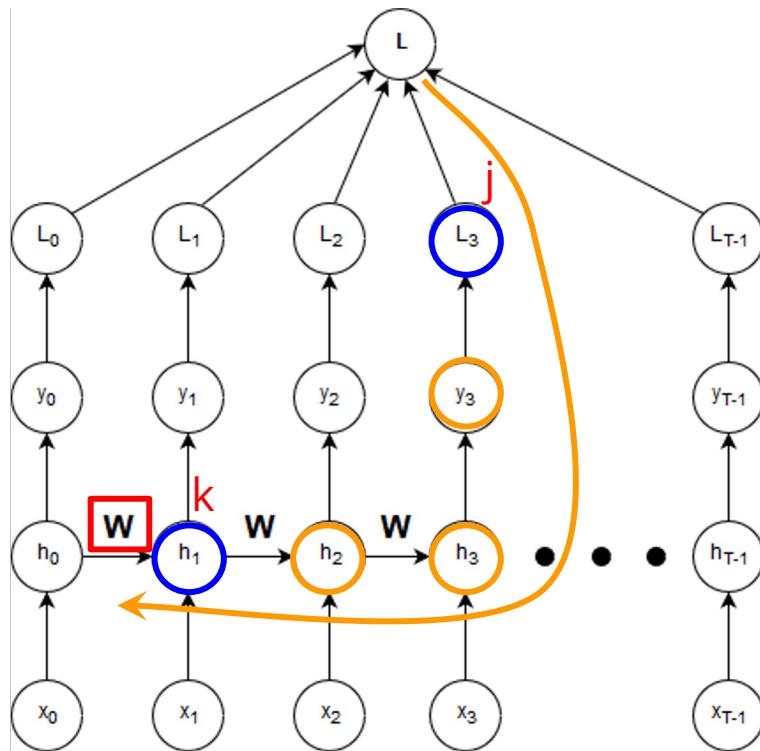


- **Second summation over h :**
Each L_j depends on the weight matrices *before it*

$$\frac{\partial L_j}{\partial \mathbf{W}} = \sum_{k=1}^j \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial \mathbf{W}}$$

L_j depends on all h_k before it.

Backpropagation as two summations

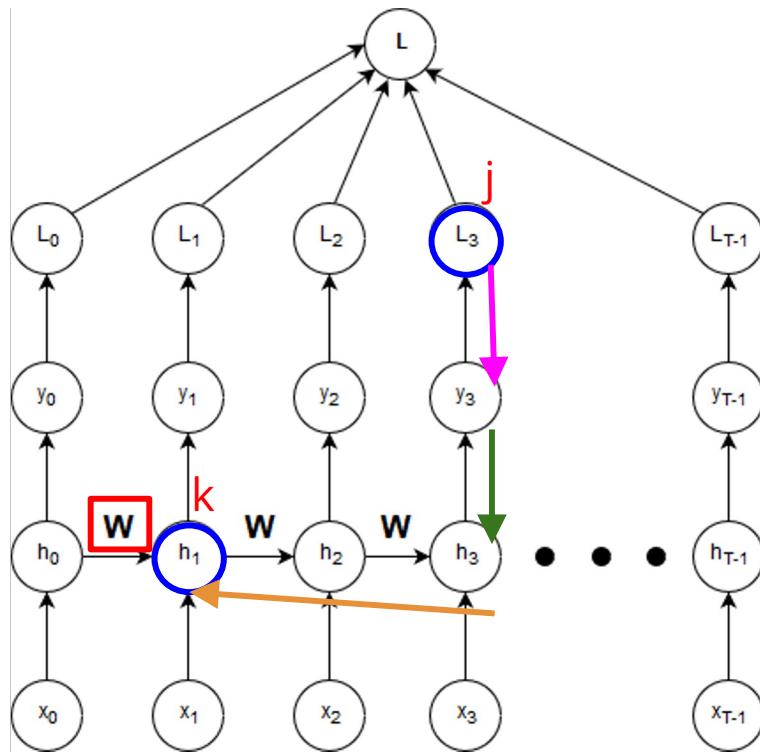


$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \boxed{\frac{\partial L_j}{\partial h_k}} \frac{\partial h_k}{\partial W}$$

- Use chain rule to fill missing steps

$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \boxed{\frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \frac{\partial h_j}{\partial h_k}} \frac{\partial h_k}{\partial W}$$

Backpropagation as two summations

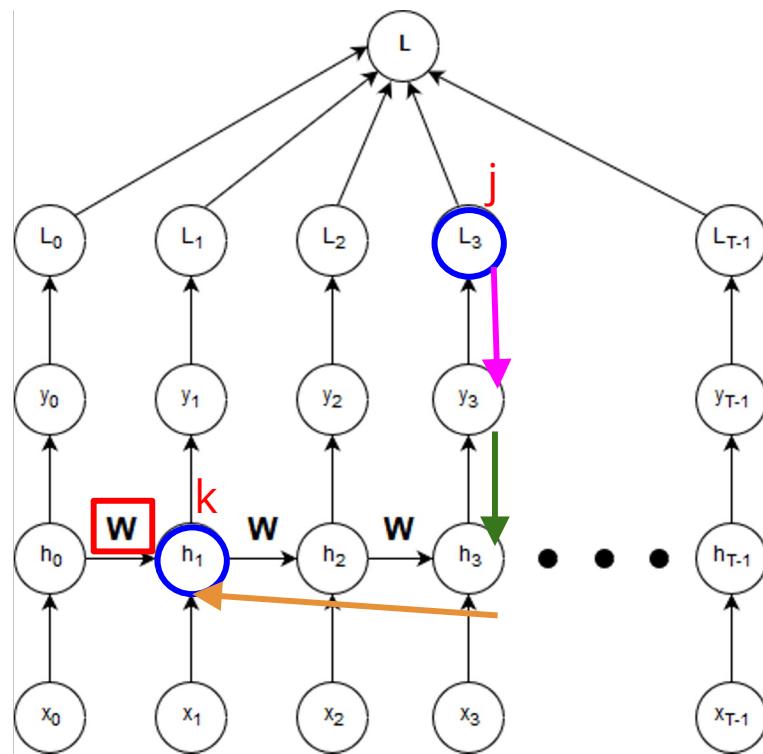


$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \boxed{\frac{\partial L_j}{\partial h_k}} \frac{\partial h_k}{\partial \mathbf{W}}$$

- Use chain rule to fill missing steps

$$\frac{\partial L_j}{\partial \mathbf{W}} = \sum_{k=1}^j \boxed{\frac{\partial L_j}{\partial y_j}} \boxed{\frac{\partial y_j}{\partial h_j}} \boxed{\frac{\partial h_j}{\partial h_k}} \frac{\partial h_k}{\partial \mathbf{W}}$$

The Jacobian



$$\frac{\partial L_j}{\partial \mathbf{W}} = \sum_{k=1}^j \left[\frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \frac{\partial h_j}{\partial h_k} \right] \frac{\partial h_k}{\partial \mathbf{W}}$$

Indirect dependency. One final use of the chain rule gives:

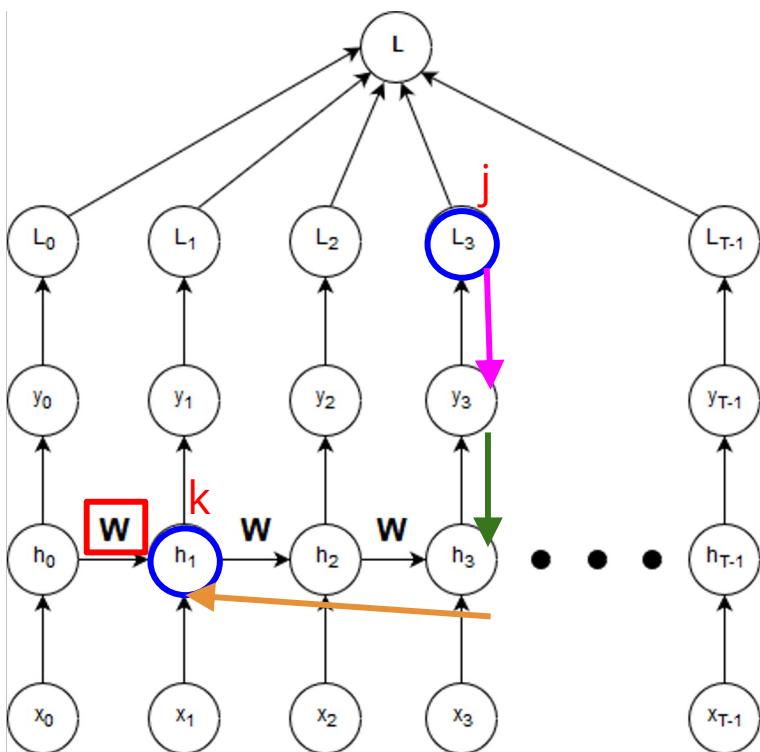
$$\frac{\partial h_j}{\partial h_k} = \prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}}$$

"The Jacobian"

The Final Backpropagation Equation

$$\frac{\partial L}{\partial \mathbf{W_h}} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left(\prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W_h}}$$

Backpropagation as two summations



$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{j=0}^{T-1} \sum_{k=1}^j \left[\frac{\partial L_j}{\partial y_j} \right] \left[\frac{\partial y_j}{\partial h_j} \right] \left(\prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \left[\frac{\partial h_k}{\partial \mathbf{W}_h} \right]$$

- Often, to reduce memory requirement, we truncate the network
- Inner summation runs from $j-p$ to j for some $p \Rightarrow$ truncated BP through time

Expanding the Jacobian

$$\frac{\partial L}{\partial W} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left(\prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}}$$

$$h_m = f(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m)$$

$$\frac{\partial h_m}{\partial h_{m-1}} = \mathbf{W}_h^T \text{diag}(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m))$$

The Issue with the Jacobian

$$\frac{\partial h_j}{\partial h_k} = \prod_{m=k+1}^j \boxed{\mathbf{W}_h^T} \boxed{\text{diag}(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m))}$$

Weight Matrix Derivative of activation function

Repeated matrix multiplications leads to **vanishing and exploding gradients**.
How? Let's take a slight detour.

Eigenvalues and Stability

- Consider identity activation function
- If Recurrent Matrix \mathbf{W}_h is a diagonalizable:

$$\mathbf{W}_h = Q * \Lambda * Q^{-1}$$

Computing powers of \mathbf{W}_h is simple:

$$\mathbf{W}_h^n = Q * \Lambda^n * Q^{-1}$$

Q matrix composed of eigenvectors of \mathbf{W}_h

Λ is a diagonal matrix with eigenvalues placed on the diagonals

Bengio et al, "On the difficulty of training recurrent neural networks." (2012)

Eigenvalues and stability

$$\Lambda = \begin{bmatrix} -0.6180 & 0 \\ 0 & 1.6180 \end{bmatrix} \rightarrow \Lambda^{10} = \begin{bmatrix} 0.0081 & 0 \\ 0 & 122.9919 \end{bmatrix}$$

Vanishing gradients

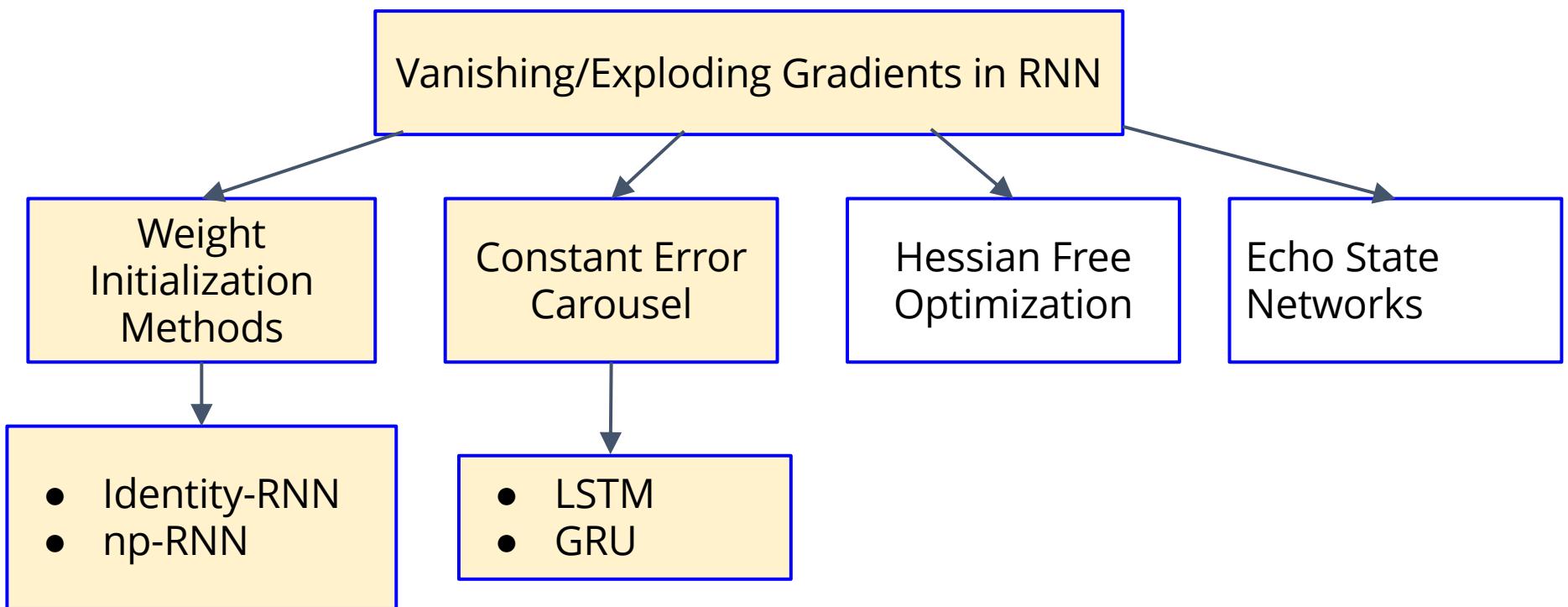
Exploding gradients

The diagram illustrates the effect of applying the matrix Λ repeatedly. The initial matrix Λ has eigenvalues -0.6180 and 1.6180 . After applying Λ 10 times, the resulting matrix Λ^{10} has elements 0.0081 and 122.9919 . The element 0.0081 is highlighted with a red box and labeled 'Vanishing gradients', indicating that the gradient flow is becoming increasingly slow. The element 122.9919 is highlighted with a red box and labeled 'Exploding gradients', indicating that the gradient flow is becoming increasingly unstable and large.

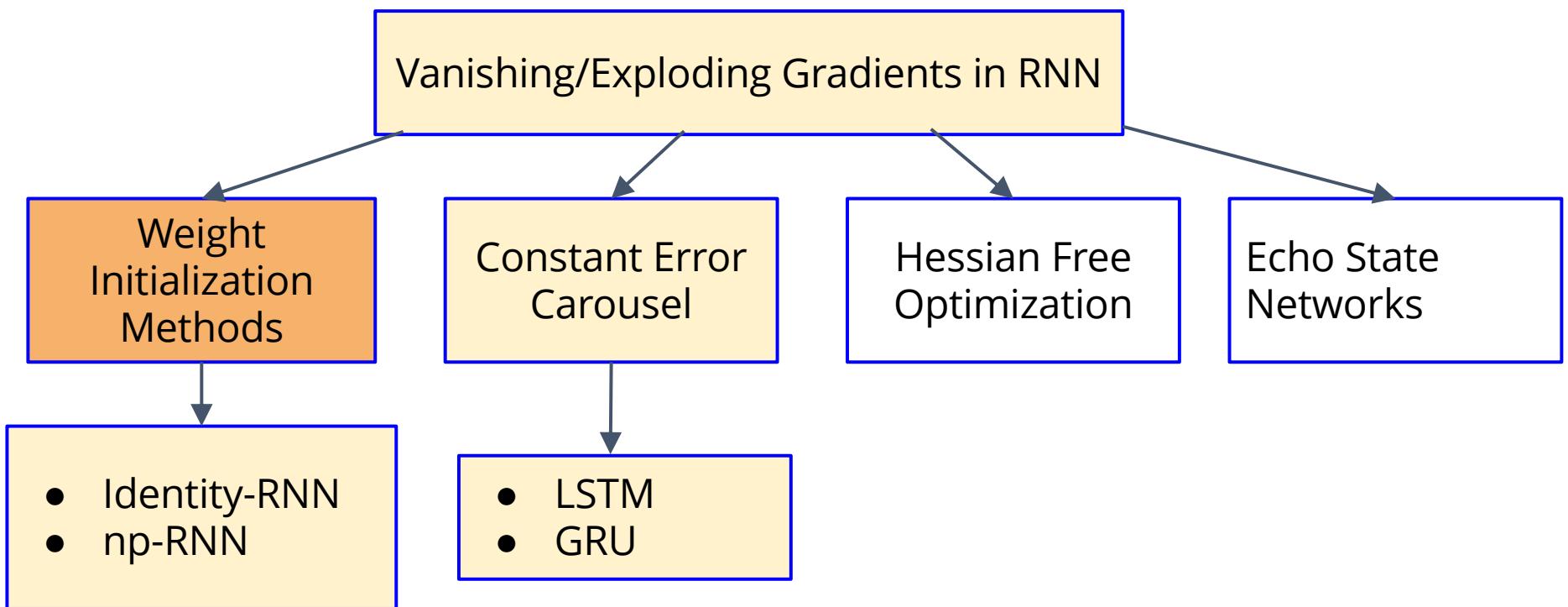
$$W_h^n = Q * \Lambda^n * Q^{-1}$$

2. Learning Long Term Dependencies

Outline



Outline



Weight Initialization Methods

$$\frac{\partial h_j}{\partial h_k} = \prod_{m=k+1}^j \mathbf{W}_h^T \text{diag}(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m))$$

- Activation function : ReLU

$$(\mathbf{W}_h^T)^n = Q^{-1} * \Lambda^n * Q^{-1}$$

Bengio et al., "On the difficulty of training recurrent neural networks." (2012)

Weight Initialization Trick #1: IRNN

- Careful initialization of W_h with suitable eigenvalues
 - ⇒ allows the RNN to learn in the initial epochs
 - ⇒ hence can generalize well for further iterations
- W_h initialized to Identity
- Activation function: ReLU

Quoc Le et al., "A Simple Way to Initialize Recurrent Networks of Rectified Linear Units" arXiv 2015

Weight Initialization Trick #2: np-RNN

- W_h positive definite (+ve real eigenvalues)
- At least one eigenvalue is 1, others all less than equal to one
- Activation function: ReLU

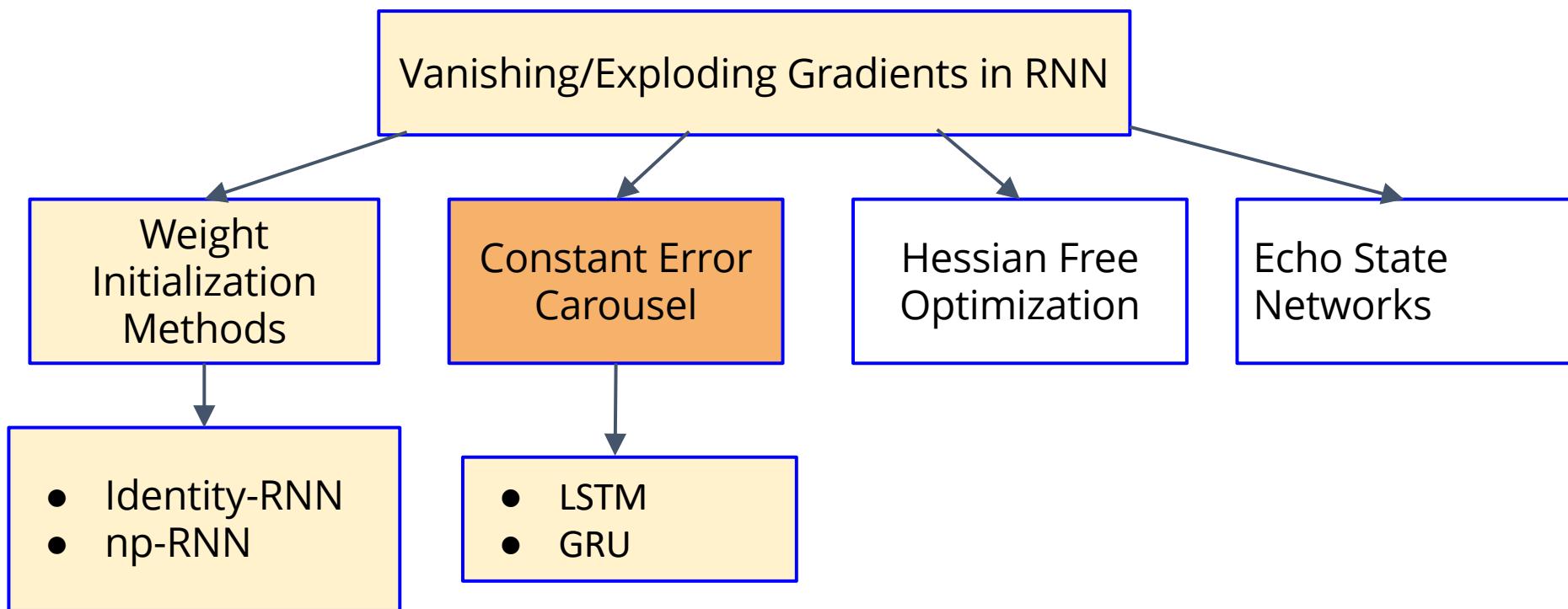
Talathi "Improving Performance of Recurrent Neural Network with ReLU nonlinearity" ICLR 2016

np-RNN vs IRNN

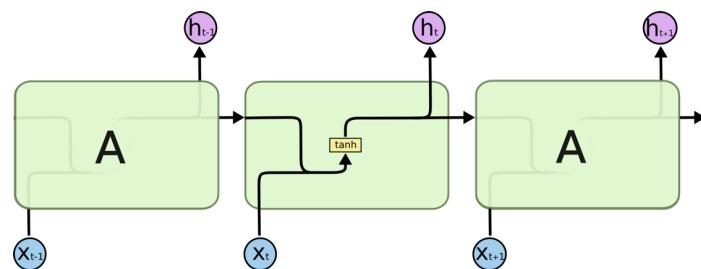
- Sequence Classification Task

RNN Type	Accuracy Test	Parameter Complexity Compared to RNN	Sensitivity to parameters
IRNN	67 %	x1	high
np-RNN	75.2 %	x1	low
LSTM (???)	78.5 %	x4	low

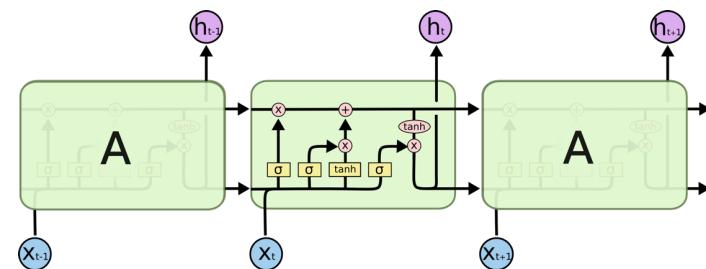
Outline



RNN vs. LSTM

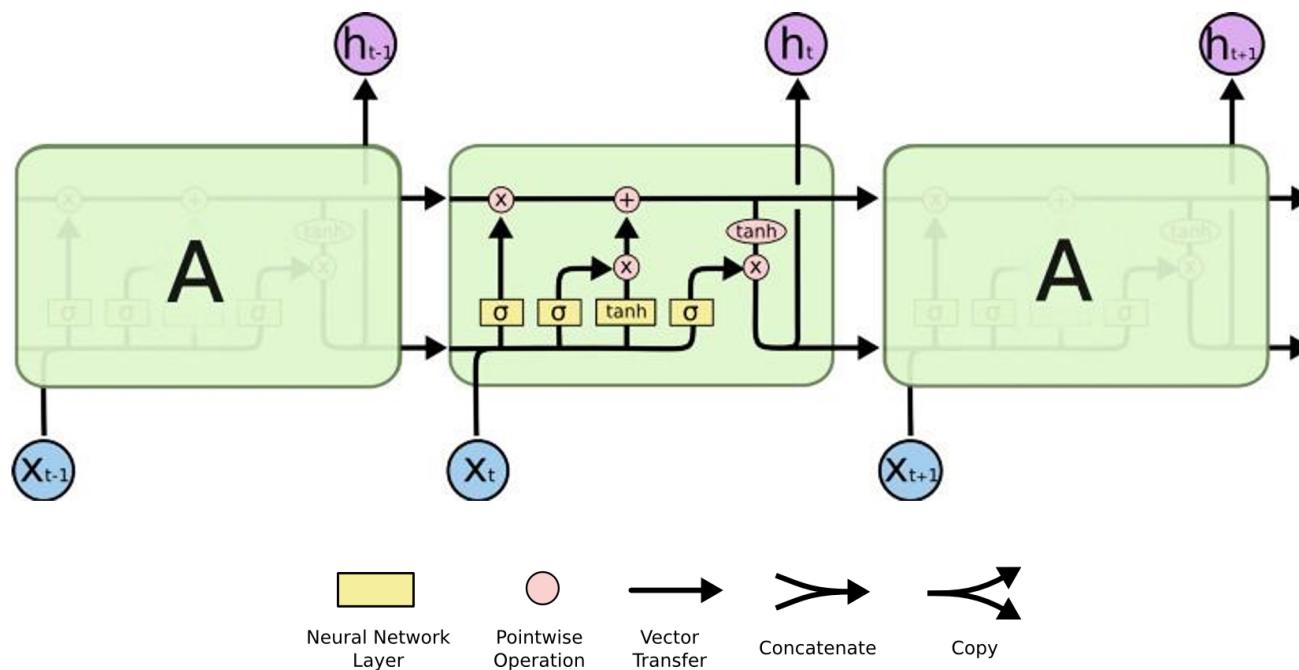


Vs.



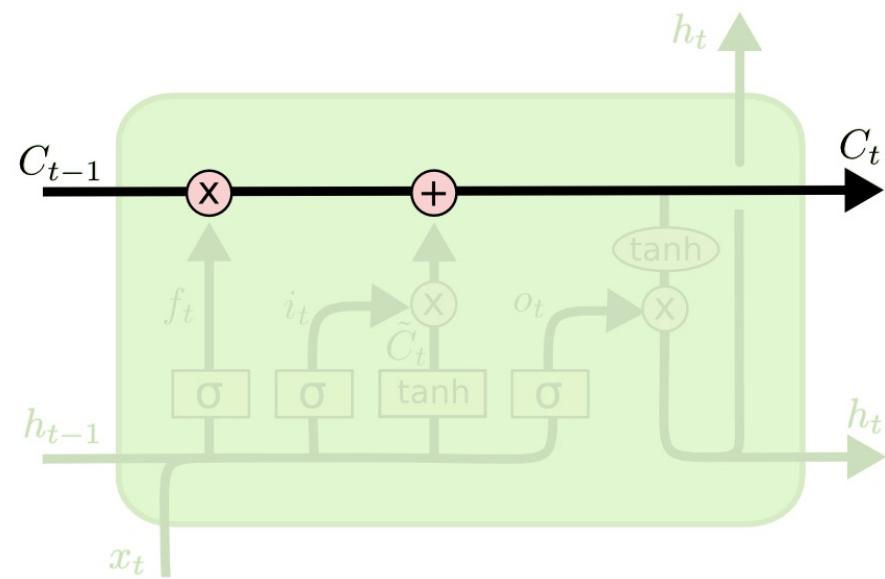
Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

The LSTM Network (Long Short-Term Memory)

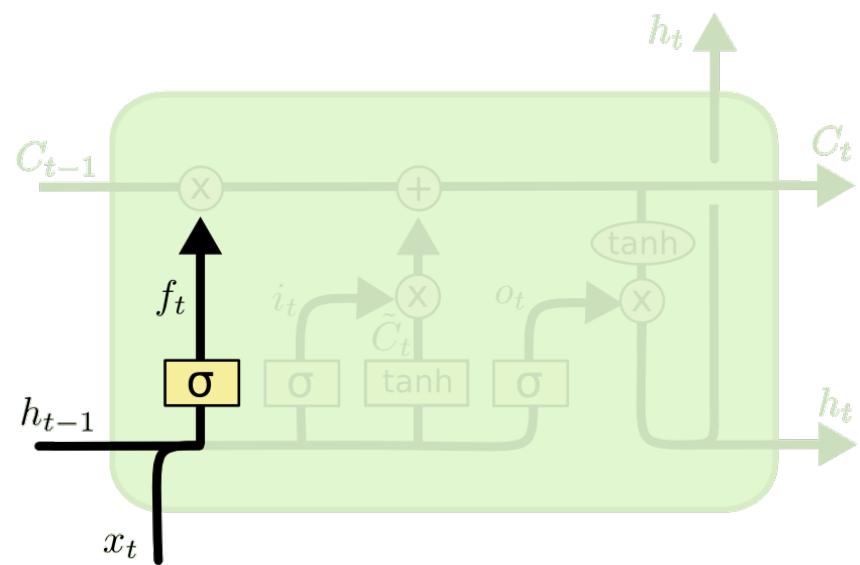
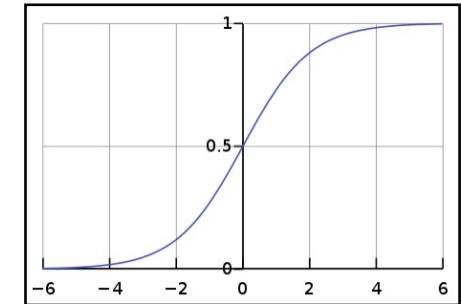


Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Avoid vanishing gradients

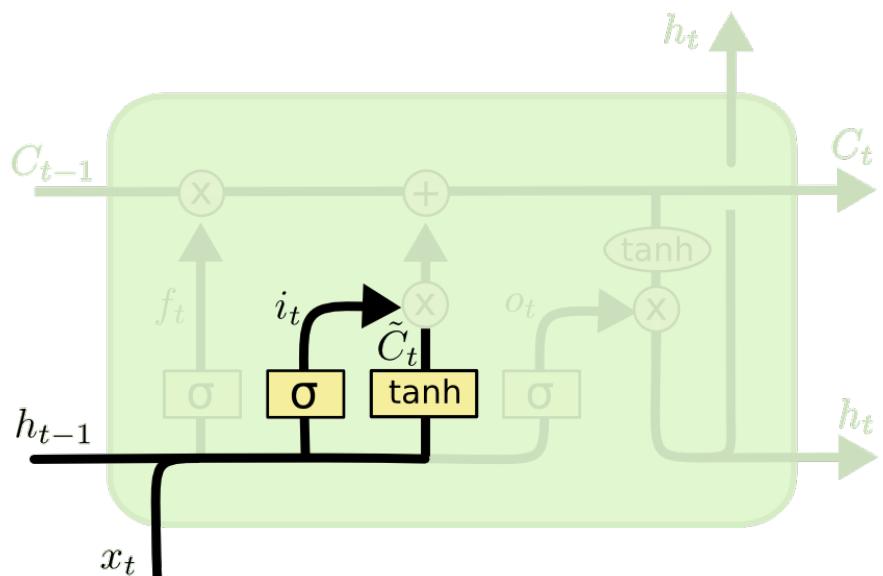
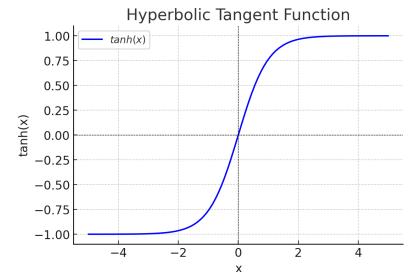


Decide “what to forget”



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

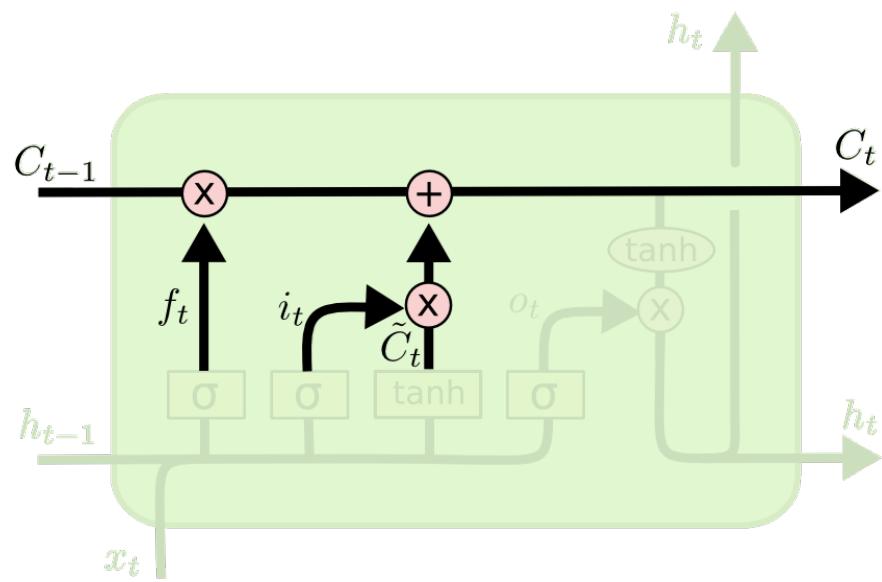
Decide “what to add”



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

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

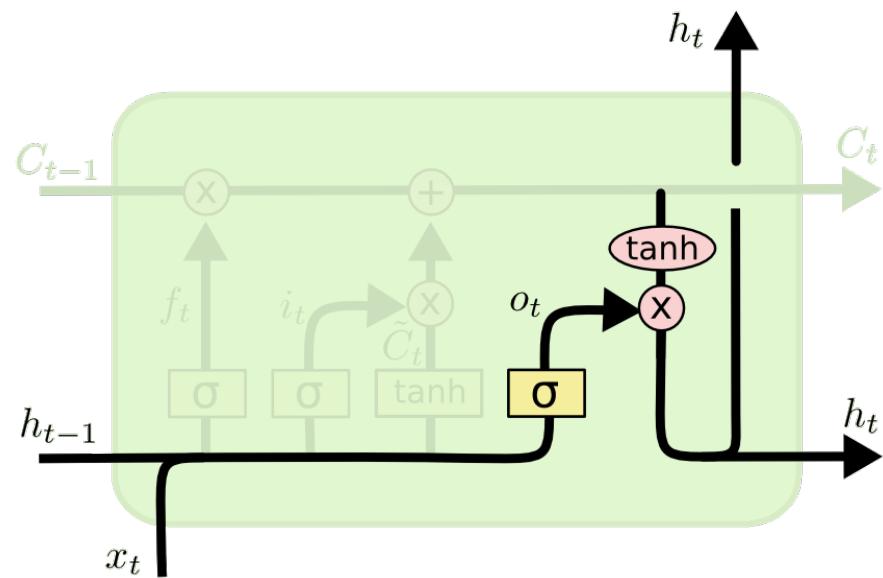
Update the state



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

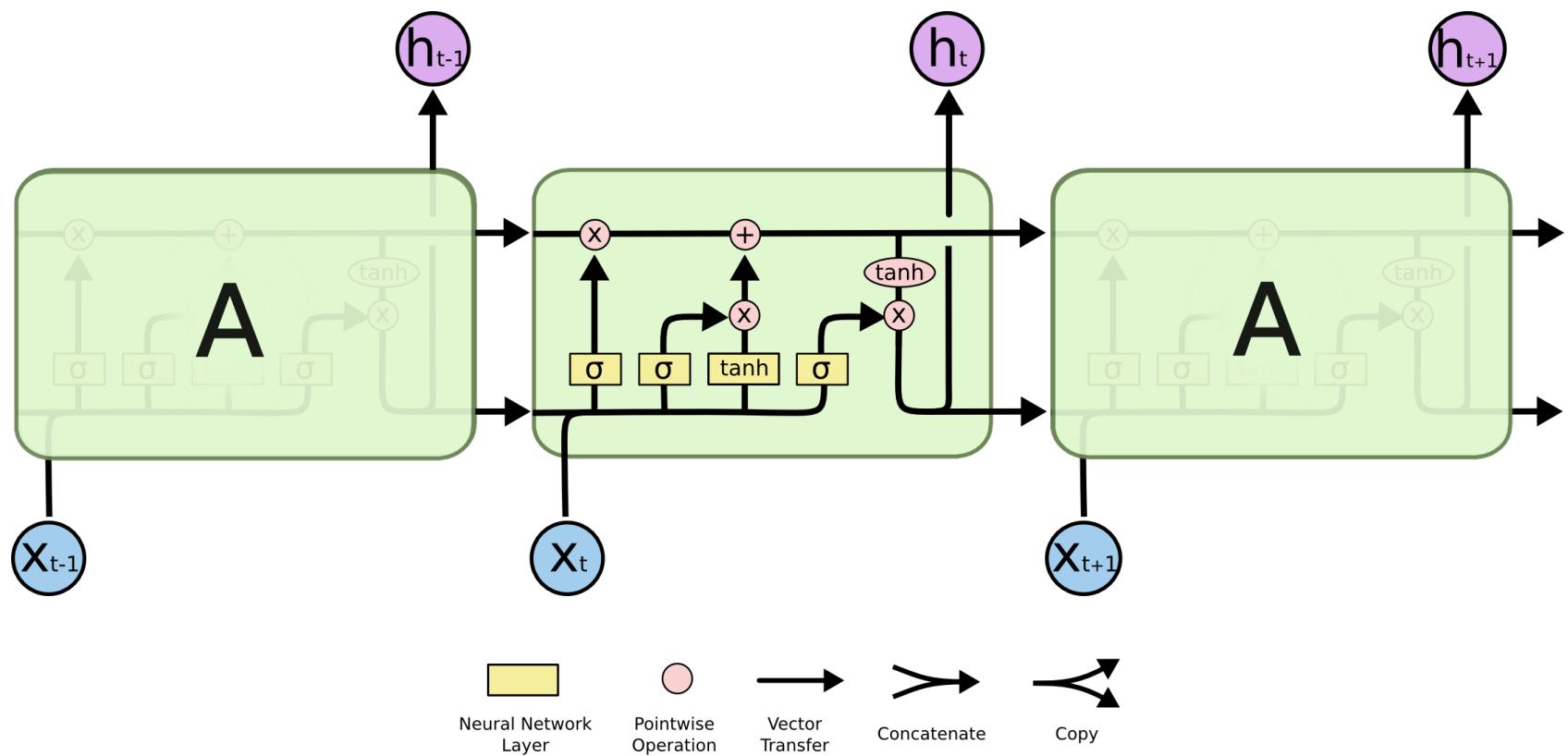
Prepare output

- Tanh re-normalizes output, and gate determines what to let through



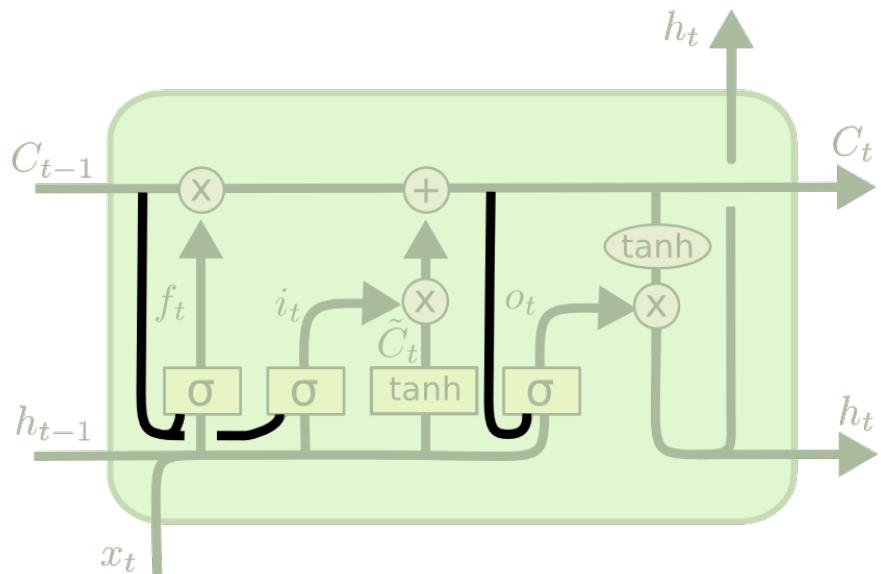
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh (C_t)$$

LSTM recap



LSTM variants (“peephole connections”)

- Let the gate layers “peep” at the cell state



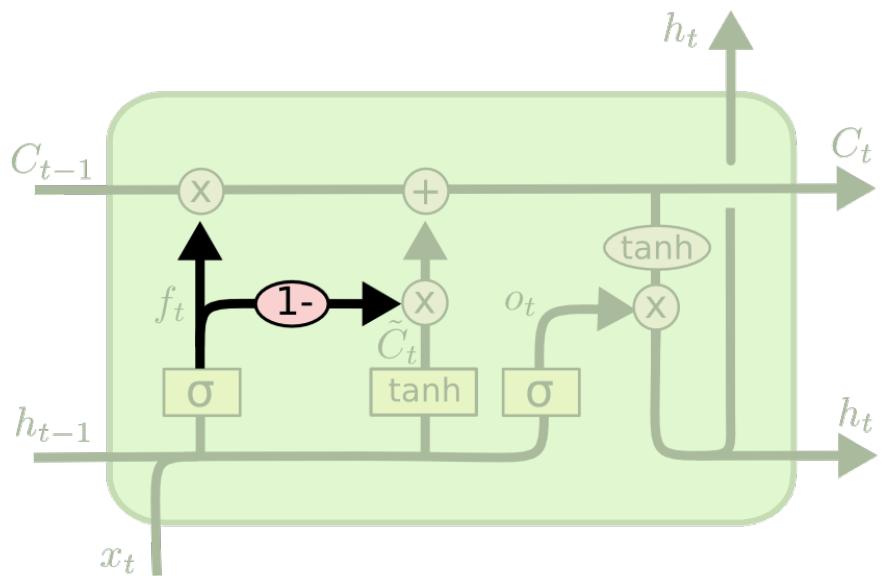
$$f_t = \sigma (W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

LSTM variants (coupled gate)

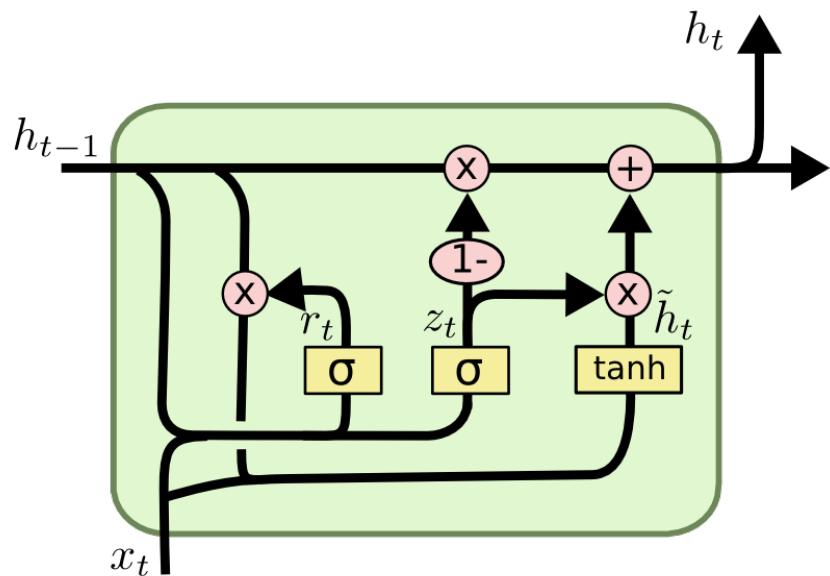
- Joint decision as to what delete/let through
(forget when we're going to input something in its place)



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

GRU (Gated Recurrent Units)

- Combines the forget and input gates into a single “update” gate



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

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

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

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

Comparing GRU and LSTM

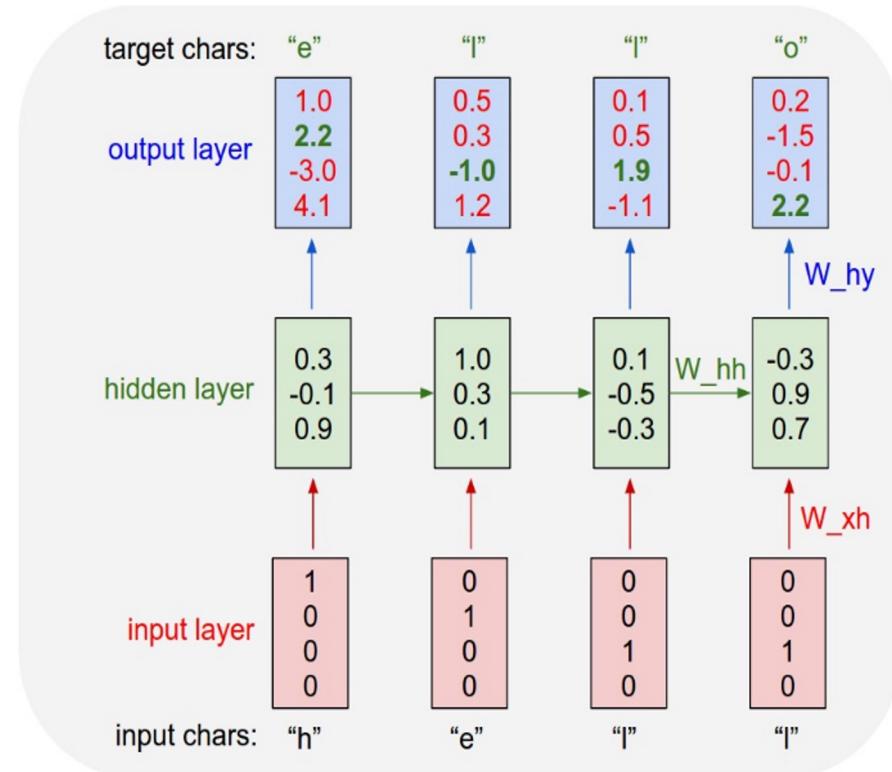
- Both **GRU and LSTM better than RNN with tanh**
(on music and speech modeling)
- GRU performs comparably to LSTM
- No clear consensus between GRU and LSTM

Source: Empirical evaluation of GRUs on sequence modeling, 2014

Section 3: Visualizing and Understanding Recurrent Networks

Character Level Language Modelling

Task: Predicting the next character given the current character



Andrej Karpathy, Blog on "Unreasonable Effectiveness of Recurrent Neural Networks"

Generated Text:

- Remembers to close a bracket
- Capitalize nouns
- 404 Page Not Found! :P The LSTM hallucinates it.

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servitious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS)[<http://www.humah.yahoo.com/guardian.cfm/7754800786d1/551963589.htm>] official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripad of aid exile.]]

Andrej Karpathy, Blog on "Unreasonable Effectiveness of Recurrent Neural Networks"

100 th iteration

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

300 th iteration

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwyl fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

700 th iteration

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

2000 th iteration

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the nacks and drove up his father-in-law women.

Andrej Karpathy, Blog on "Unreasonable Effectiveness of Recurrent Neural Networks"

Features RNN Captures in Common Language ?

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.



Cell Sensitive to Position in Line

- Can be interpreted as tracking the line length

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell That Turns On Inside Quotes

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

Andrej Karpathy, Blog on "Unreasonable Effectiveness of Recurrent Neural Networks"

Features RNN Captures in C Language?

```
for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
        pipe = (in_use & UMXTHREAD_UNCCA) +
            ((count & 0x00000000fffffff8) & 0x000000f) << 8;
    if (count == 0)
        sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
}
/* Free our user pages pointer to place camera if all dash */
```

Andrej Karpathy, Blog on “Unreasonable Effectiveness of Recurrent Neural Networks”

Cell That Activates Inside IF Statements

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
    siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

Andrej Karpathy, Blog on “Unreasonable Effectiveness of Recurrent Neural Networks”

Cell That Is Sensitive To Indentation

- Can be interpreted as tracking indentation of code.
- Increasing strength as indentation increases

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

Andrej Karpathy, Blog on “Unreasonable Effectiveness of Recurrent Neural Networks”

Non-Interpretable Cells

- Only 5% of the cells show such interesting properties
- Large portion of the cells are **not interpretable** by themselves

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
    */
```

Andrej Karpathy, Blog on "Unreasonable Effectiveness of Recurrent Neural Networks"

Key Takeaways

- **Recurrent neural network for sequential data processing**
 - Feedforward depth
 - Recurrent depth
- **Long term dependencies** are a major problem in RNNs. Solutions:
 - Intelligent weight initialization
 - LSTMs / GRUs
- **Visualization** helps
 - Analyze finer details of features produced by RNNs

References

- Survey Papers
 - Lipton, Zachary C., John Berkowitz, and Charles Elkan. [A critical review of recurrent neural networks for sequence learning](#). arXiv preprint arXiv:1506.00019 (2015).
 - Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. [Chapter 10: Sequence Modeling: Recurrent and Recursive Nets](#). MIT Press, 2016.
- Training
 - Semeniuta, Stanislau, Aliaksei Severyn, and Erhardt Barth. [Recurrent dropout without memory loss](#). arXiv preprint arXiv:1603.05118 (2016).
 - Arjovsky, Martin, Amar Shah, and Yoshua Bengio. [Unitary evolution recurrent neural networks](#). arXiv preprint arXiv:1511.06464 (2015).
 - Le, Quoc V., Navdeep Jaitly, and Geoffrey E. Hinton. [A simple way to initialize recurrent networks of rectified linear units](#). arXiv preprint arXiv:1504.00941 (2015).
 - Cooijmans, Tim, et al. [Recurrent batch normalization](#). arXiv preprint arXiv:1603.09025 (2016).

References (contd)

- Architectural Complexity Measures
 - Zhang, Saizheng, et al. [Architectural Complexity Measures of Recurrent Neural Networks](#). Advances in Neural Information Processing Systems. 2016.
 - Pascanu, Razvan, et al. [How to construct deep recurrent neural networks](#). arXiv preprint arXiv:1312.6026 (2013).
- RNN Variants
 - Zilly, Julian Georg, et al. [Recurrent highway networks](#). arXiv preprint arXiv:1607.03474 (2016)
 - Chung, Junyoung, Sungjin Ahn, and Yoshua Bengio. [Hierarchical multiscale recurrent neural networks](#), arXiv preprint arXiv:1609.01704 (2016).
- Visualization
 - Karpathy, Andrej, Justin Johnson, and Li Fei-Fei. [Visualizing and understanding recurrent networks](#). arXiv preprint arXiv:1506.02078 (2015).
 - Hendrik Strobelt, Sebastian Gehrmann, Bernd Huber, Hanspeter Pfister, Alexander M. Rush. [LSTMVis: Visual Analysis for RNN](#), arXiv preprint arXiv:1606.07461 (2016).