# ELEC/COMP 576:
# Understanding and Visualizing Convnets & Introduction to Recurrent Neural Networks
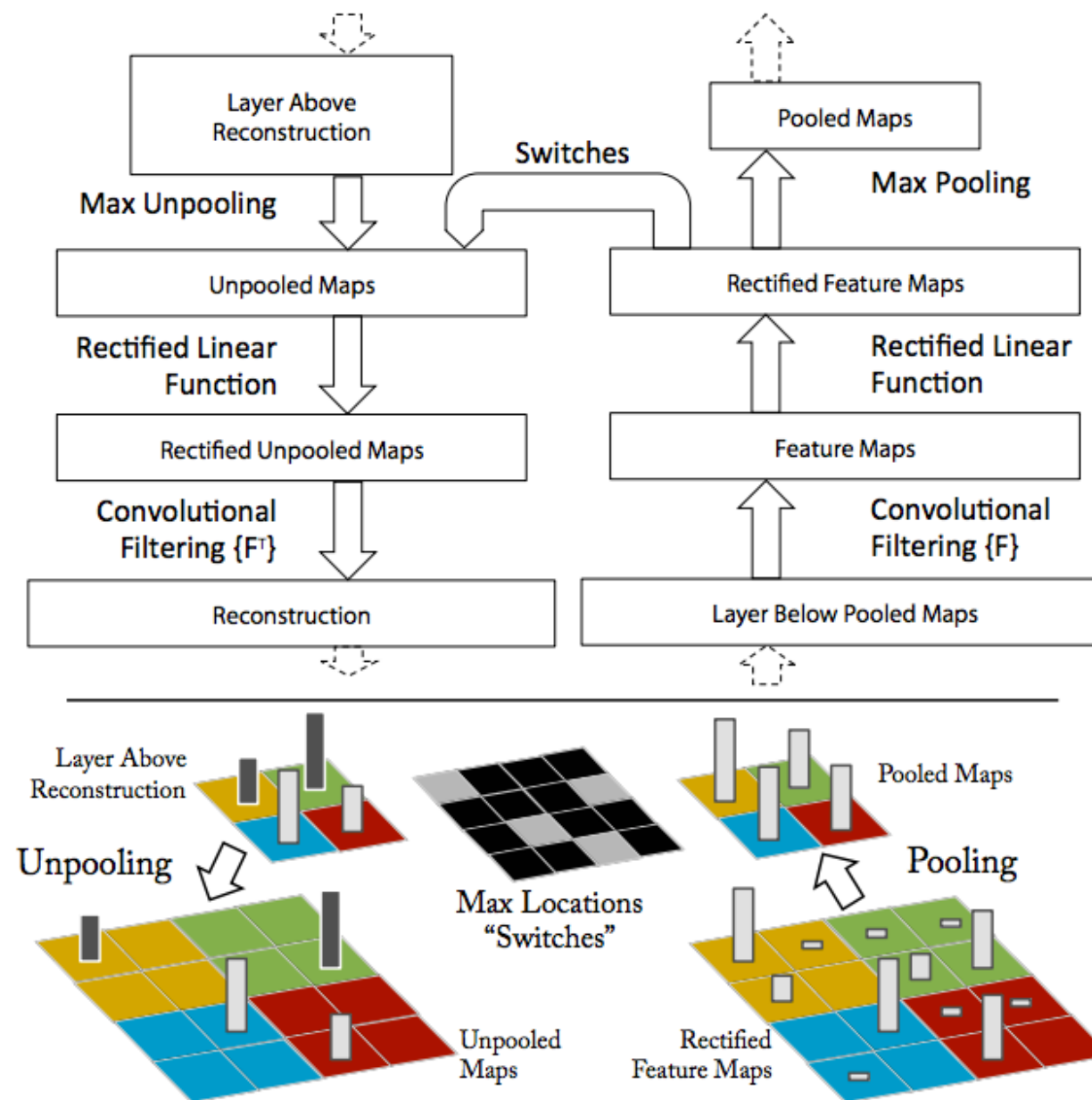
**Ankit B. Patel**

*Baylor College of Medicine (Neuroscience Dept.)*
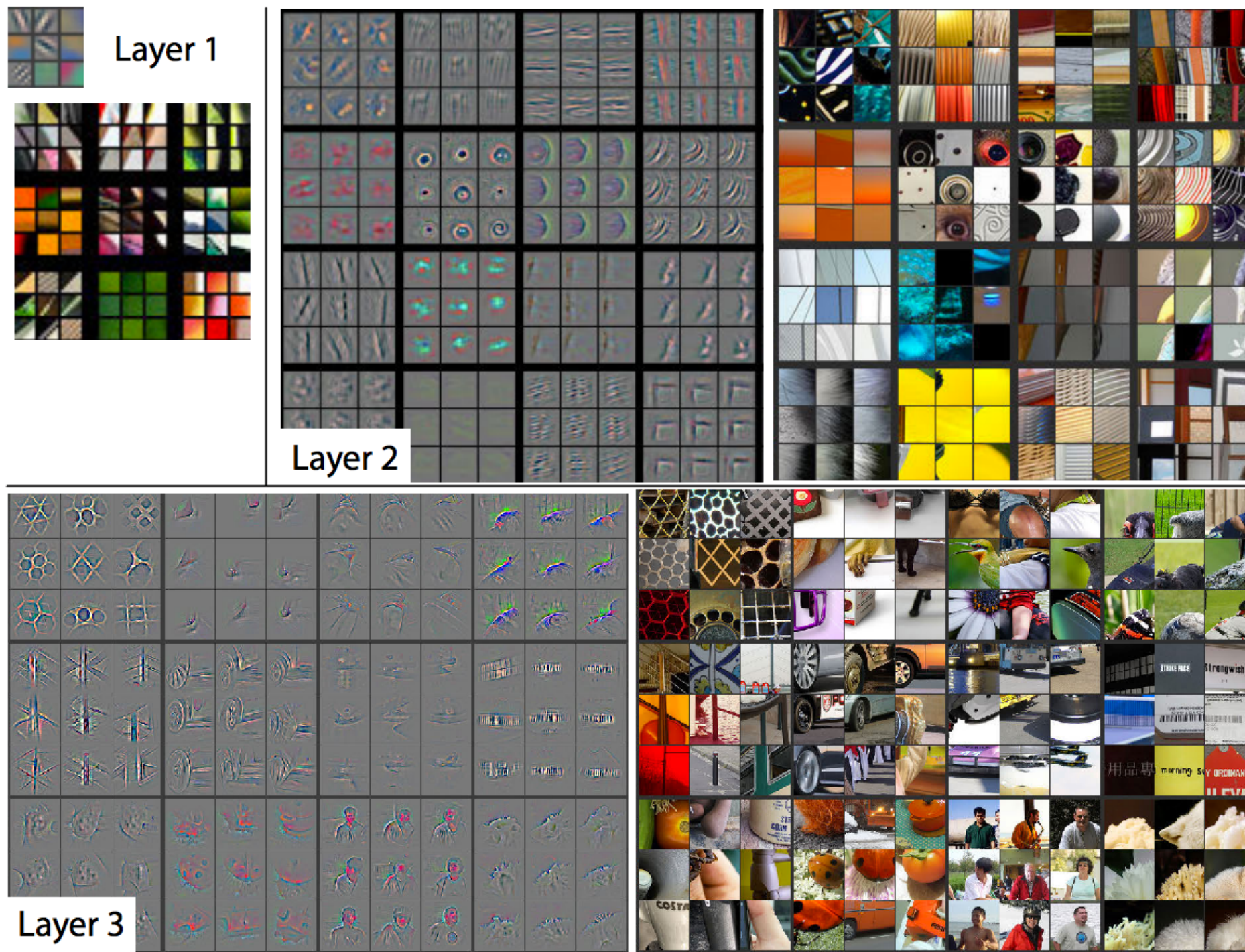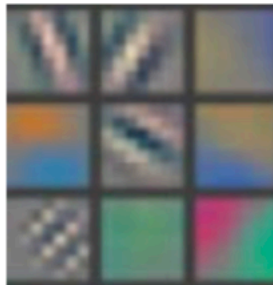*Rice University (ECE Dept.)*

# Understand & Visualizing Convnets
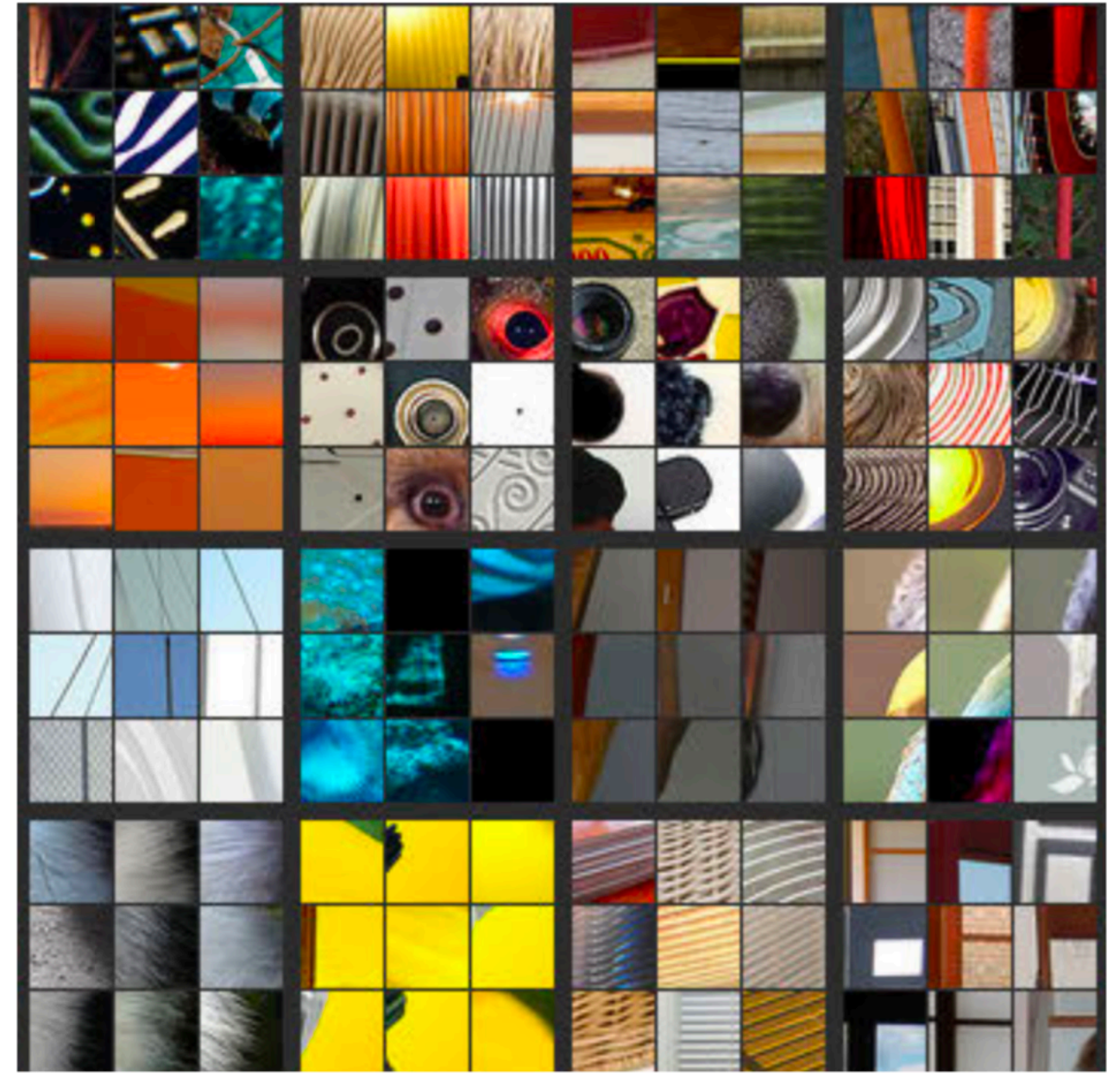
# Deconvolutional Net
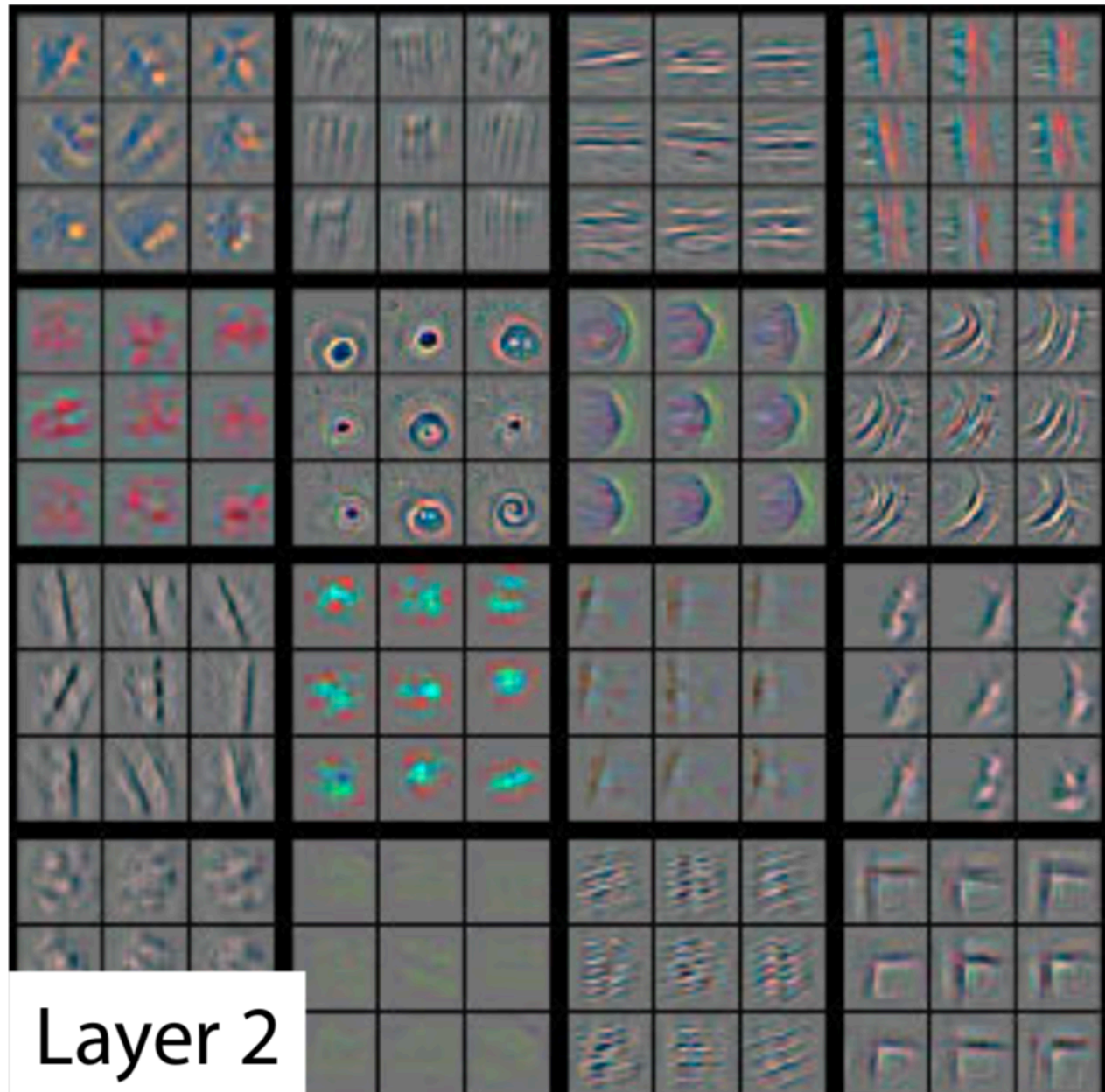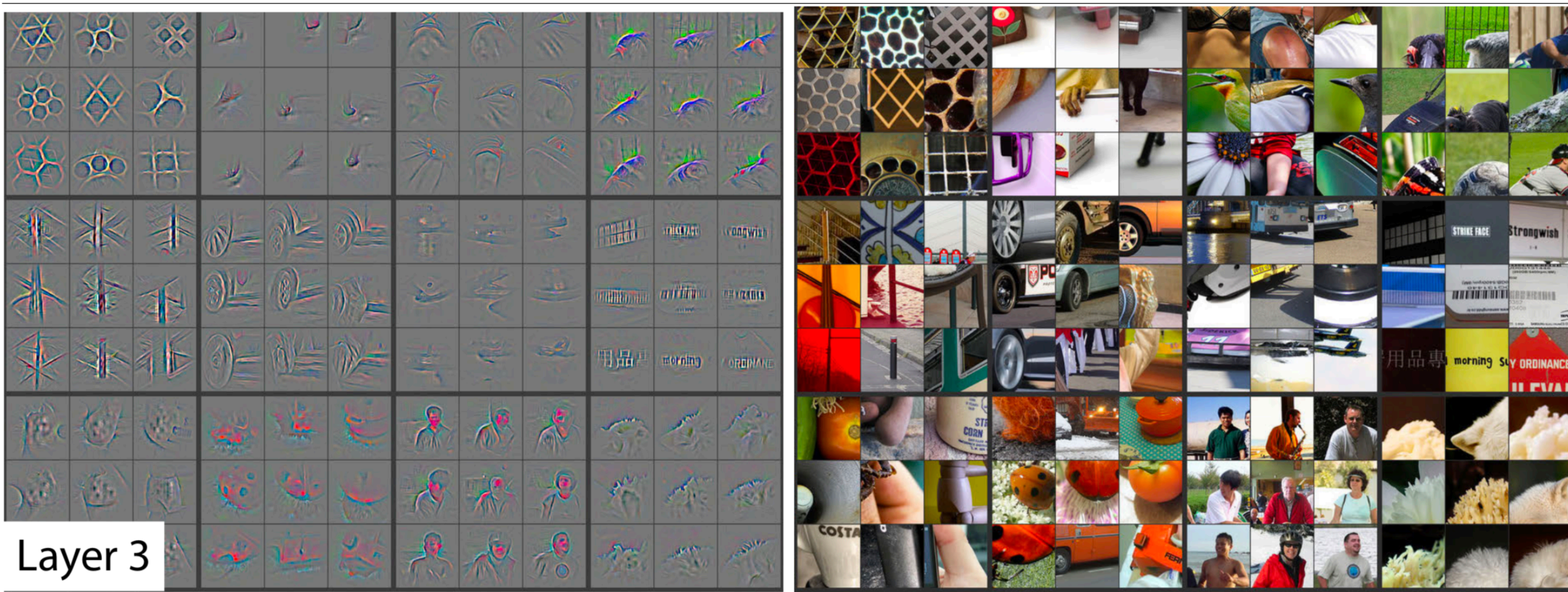
# Feature Visualization



Layer 1

Layer 2

Layer 3

# Feature Visualization



Layer 1

# Feature Visualization
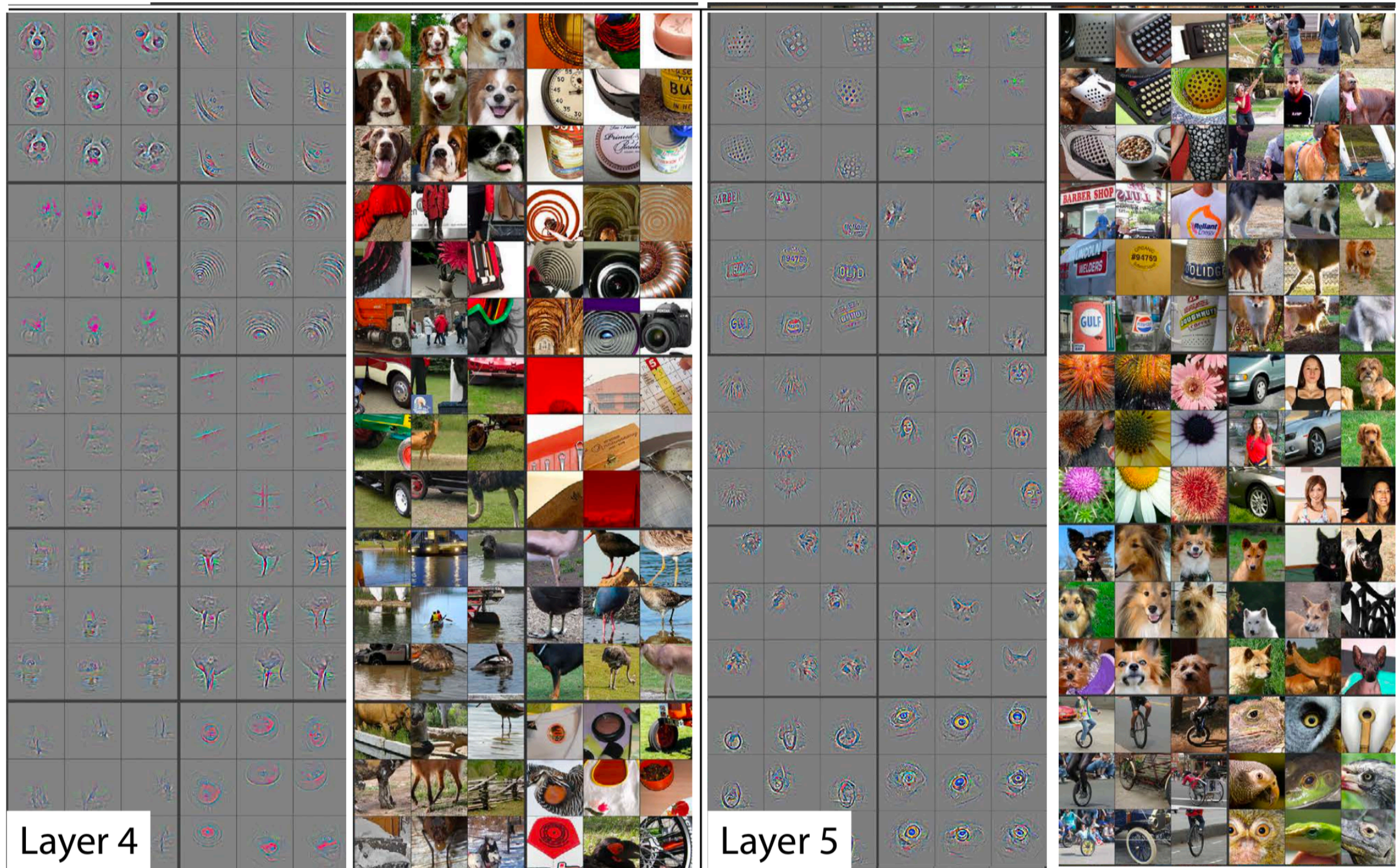


Layer 2

[Zeiler and Fergus]

# Feature Visualization



Layer 3

# Feature Visualization
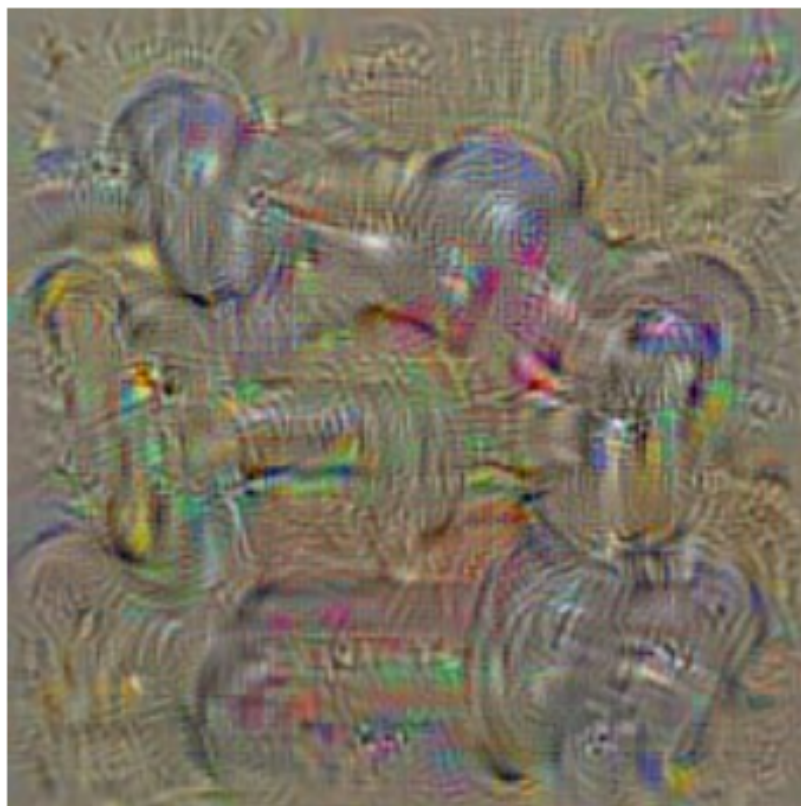


Layer 4

Layer 5

[Zeiler and Fergus]

# Activity Maximization (aka Saliency Maps)

$$\arg\max_I S_c(I) - \lambda\|I\|_2^2, \qquad S_c(I) \approx w^T I + b,$$



dumbbell

cup

dalmatian

[Simonyan et al.]

# Deep Dream Visualization

# Deep Dream Visualization

- To produce human viewable images, need to

  - Activity maximization (gradient ascent)

  - L2 regularization

  - Gaussian blur

  - Clipping

  - Multiple scales (octaves)

  - **Code:** https://github.com/google/deepdream/blob/master/dream.ipynb

# Example Image

# Dumbbell Deep Dream



AlexNet



VGGNet



GoogleNet

# Deep Dream Video



Class: goldfish, Carassius auratus

# Infinite Zoom-In
# on Deep Dream

https://www.youtube.com/watch?v=SCE-QeDfXtA

# Texture Synthesis



$$E_L = \sum \left( \hat{G}^L - G^L \right)^2$$

$$\hat{G}^L_{ij} = \sum_k \hat{F}^L_{ik} \hat{F}^L_{jk}$$

$$\mathcal{L}(\vec{x}, \hat{\vec{x}}) = \sum_{l=0}^{L} w_l E_l$$

$$\hat{\vec{x}} := \hat{\vec{x}} - \alpha \frac{\partial \mathcal{L}}{\partial \hat{\vec{x}}}$$

**[Leon Gatys, Alexander Ecker, Matthias Bethge]**

# Generated Textures



[Leon Gatys, Alexander Ecker, Matthias Bethge]

# DeepStyle Examples



**[Leon Gatys, Alexander Ecker, Matthias Bethge]**

# DeepStyle: Combining Style + Content from Distinct Images

$$\mathcal{L}_{content}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} \left(F_{ij}^l - P_{ij}^l\right)^2 \, .$$

The derivative of this loss with respect to the activations in layer $l$ equals

$$\frac{\partial \mathcal{L}_{content}}{\partial F_{ij}^l} = \begin{cases} \left(F^l - P^l\right)_{ij} & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0 \, . \end{cases}$$

$$E_l = \frac{1}{4 N_l^2 M_l^2} \sum_{i,j} \left(G_{ij}^l - A_{ij}^l\right)^2$$

and the total loss is

$$\mathcal{L}_{style}(\vec{a}, \vec{x}) = \sum_{l=0}^{L} w_l E_l$$

$$\frac{\partial E_l}{\partial F_{ij}^l} = \begin{cases} \frac{1}{N_l^2 M_l^2} \left((F^l)^{\mathrm{T}} \left(G^l - A^l\right)\right)_{ji} & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0 \, . \end{cases}$$

**[Leon Gatys, Alexander Ecker, Matthias Bethge]**

# Introduction to
# Recurrent Neural Networks
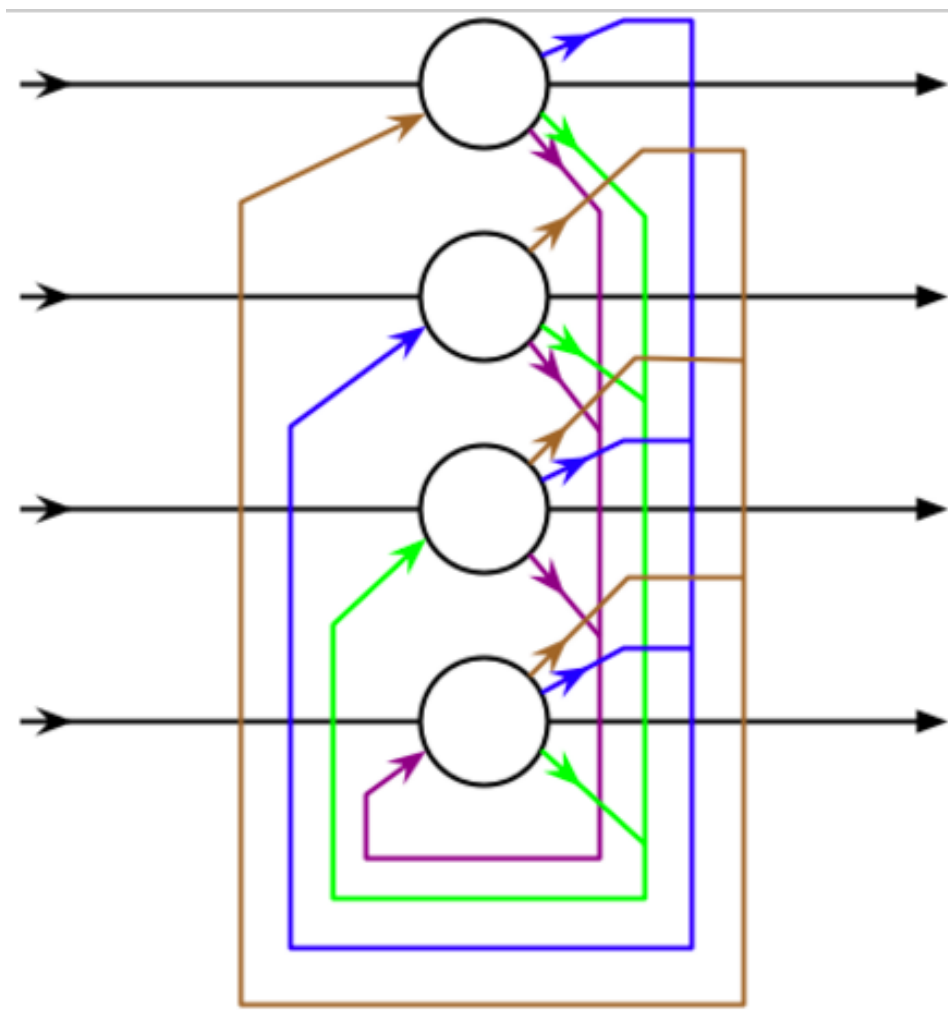
# What Are Recurrent Neural Networks?

- Recurrent Neural Networks (RNNs) are networks that have feedback

  - Output is feed back to the input

  - Sequence processing

- Ideal for time-series data or sequential data

# History of RNNs

# Important RNN Architectures

- Hopfield Network

- Jordan and Elman Networks

- Echo State Networks

- Long Short Term Memory (LSTM)

- Bi-Directional RNN

- Gated Recurrent Unit (GRU)

- Neural Turing Machine
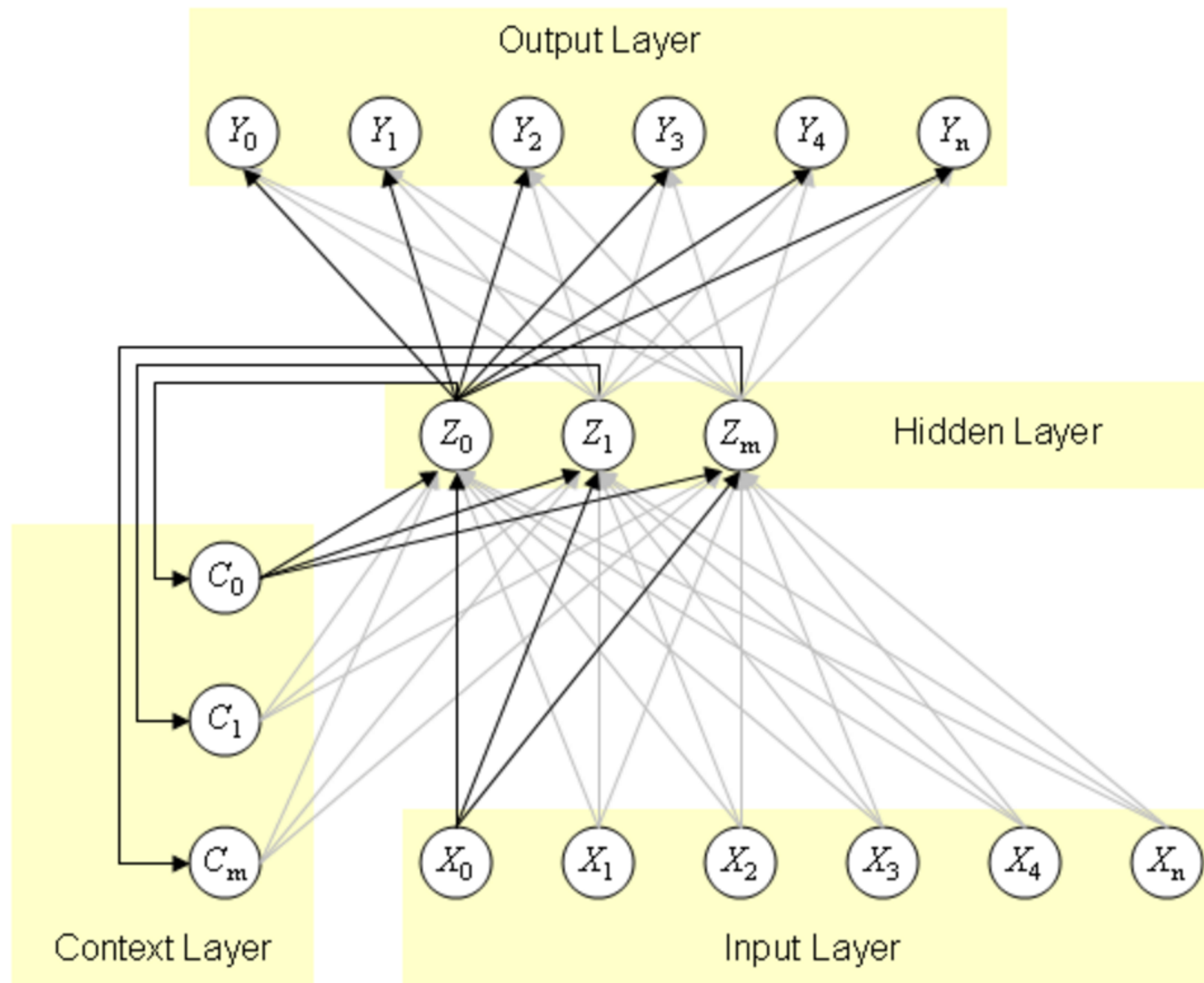
# Hopfield Network



$$s_i \leftarrow \begin{cases} +1 & \text{if } \sum_j w_{ij} s_j \geq \theta_i, \\ -1 & \text{otherwise.} \end{cases}$$
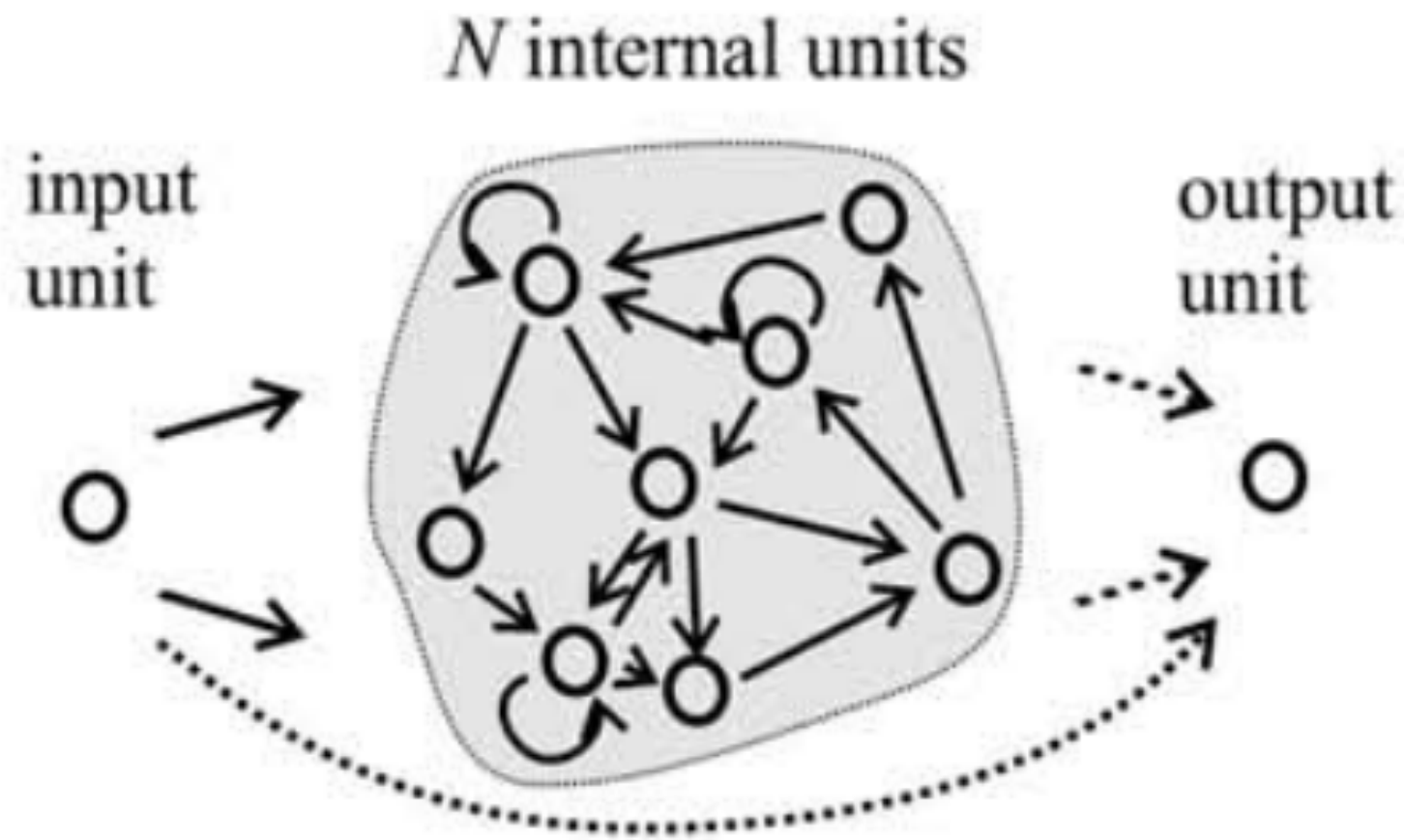
where:

- $w_{ij}$ is the strength of the connection weight from unit j to unit i (the weight of the connection).
- $s_j$ is the state of unit j.
- $\theta_i$ is the threshold of unit i.

# Elman Networks



[John McCullock]

# Echo State Networks



[Herbert Jaeger]

# Definition of RNNs

# RNN Formulation
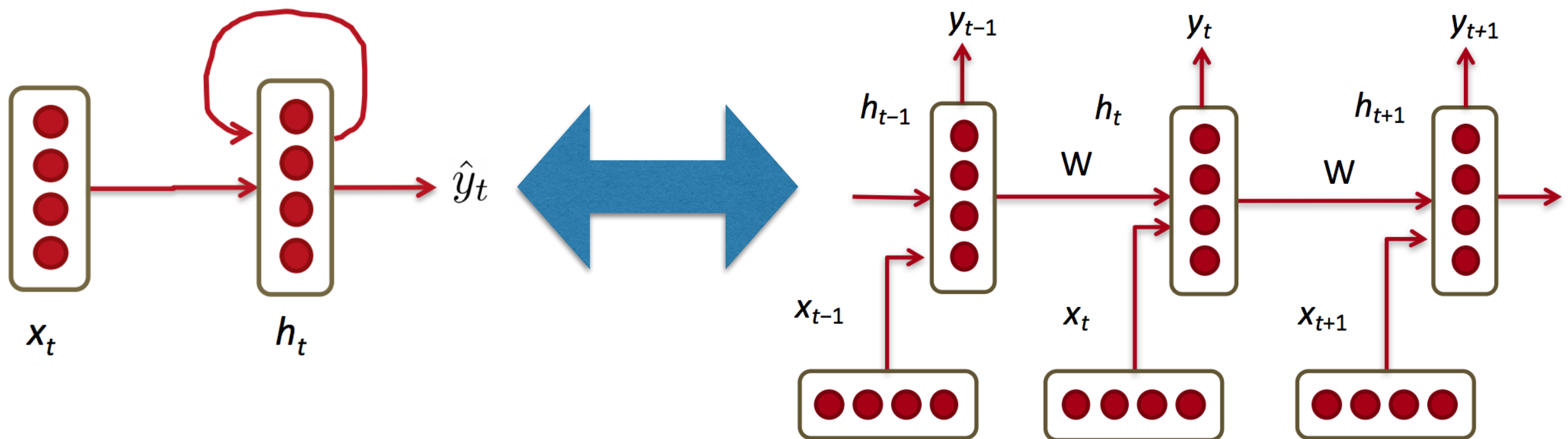
$$x_1, \ldots, x_{t-1}, x_t, x_{t+1}, \ldots, x_T$$

$$h_t = \sigma\left(W^{(hh)}h_{t-1} + W^{(hx)}x_{[t]}\right)$$

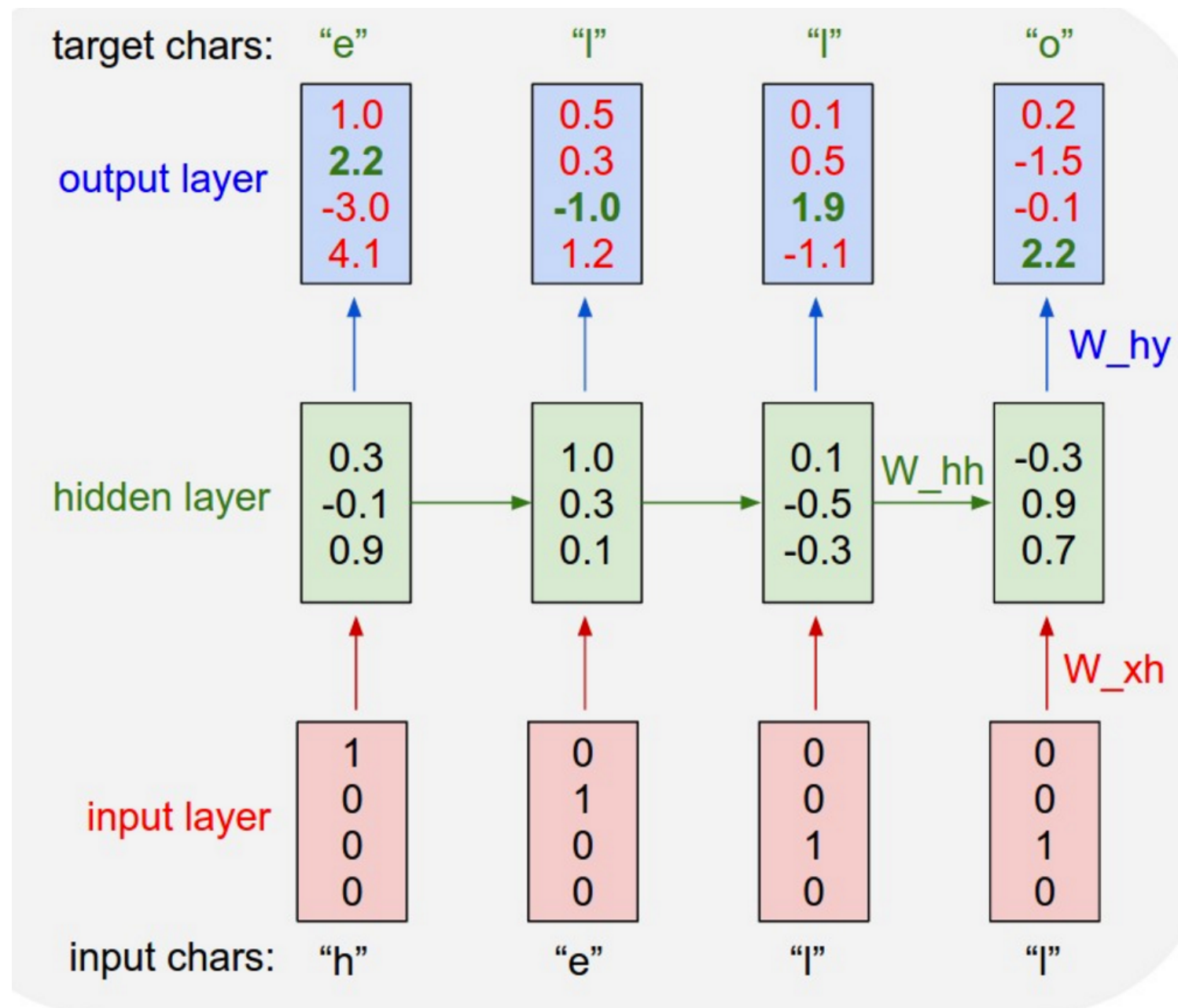$$\hat{y}_t = \operatorname{softmax}\left(W^{(S)}h_t\right)$$

**[Richard Socher]**
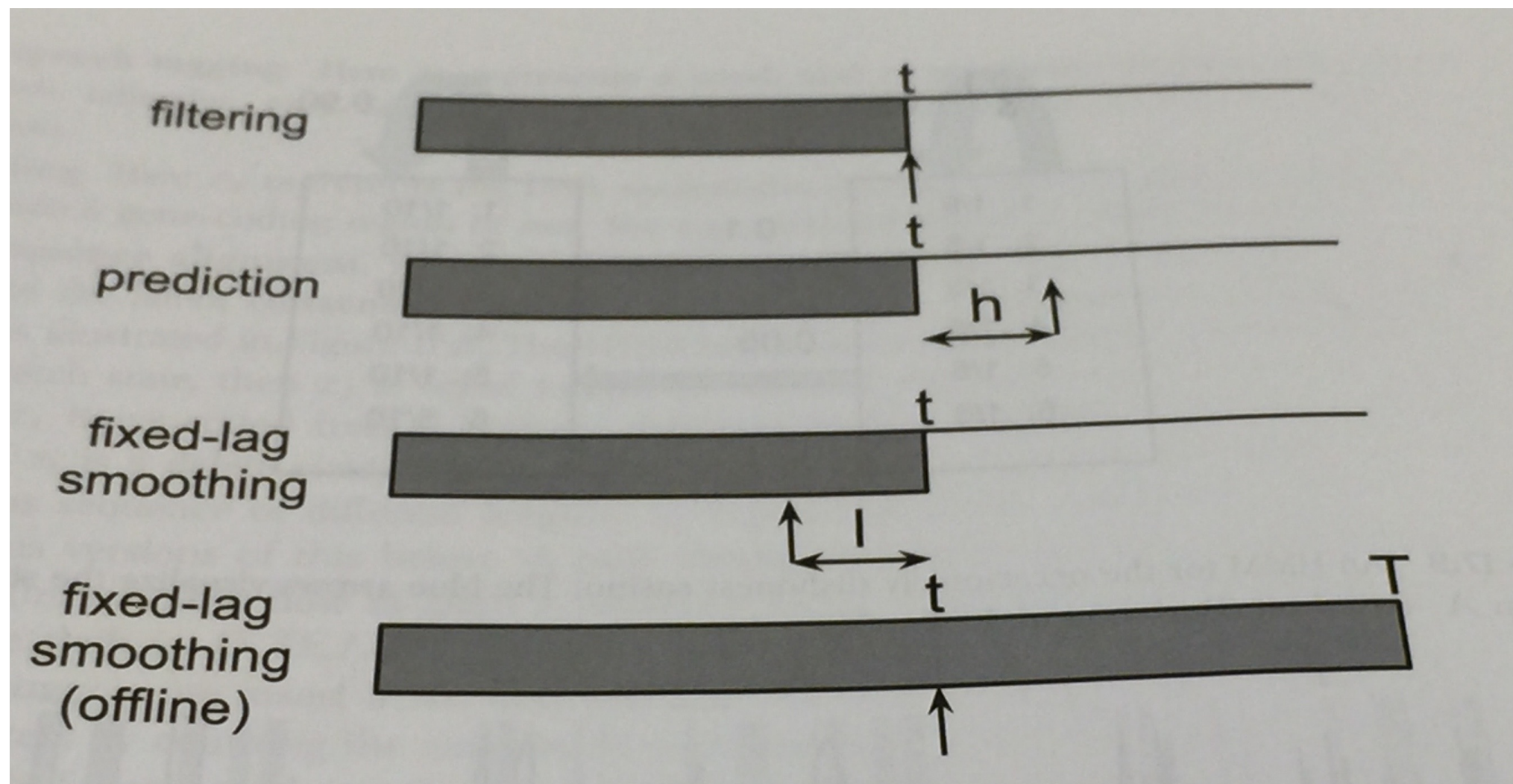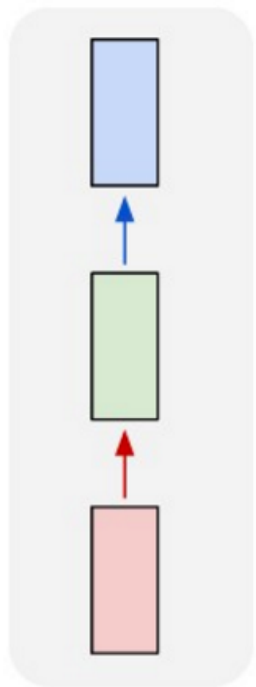
# RNN Diagram

Unrolled into FF NN

# RNN Example



[Andrej Karpathy]

# Different Inference Tasks —> Different RNN Architectures



[Kevin Murphy]

# Different Structures for Filtering/Prediction Tasks
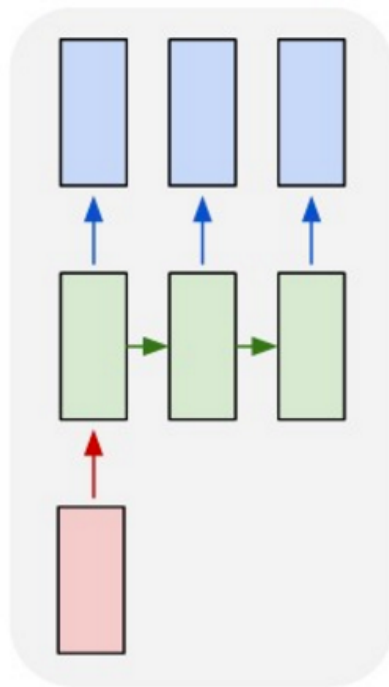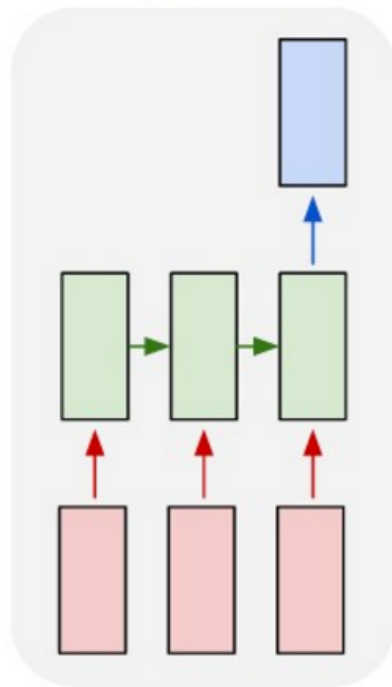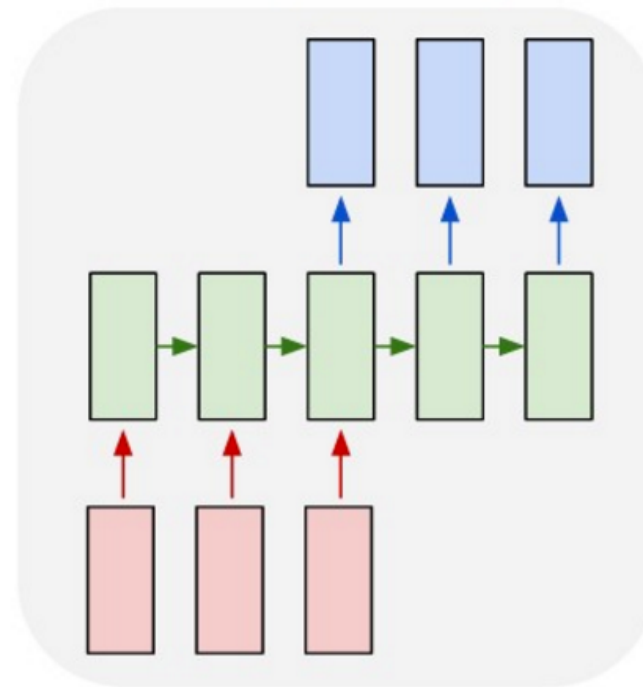


| one to one | one to many | many to one | many to many | many to many |
|:---:|:---:|:---:|:---:|:---:|
| Object Recognition | Image Captioning | Action Recognition | Machine Translation | Object Tracking |

[Andrej Karpathy]

# Universal Expressive Power Results

The *Universal Approximation Theorem* tells us that:

Any non-linear dynamical system can be approximated to any accuracy by a recurrent neural network, with no restrictions on the compactness of the state space, provided that the network has enough sigmoidal hidden units.

This underlies the computational power of recurrent neural networks.

**[John Bullinaria]**

# Training RNNs

# Training an RNN

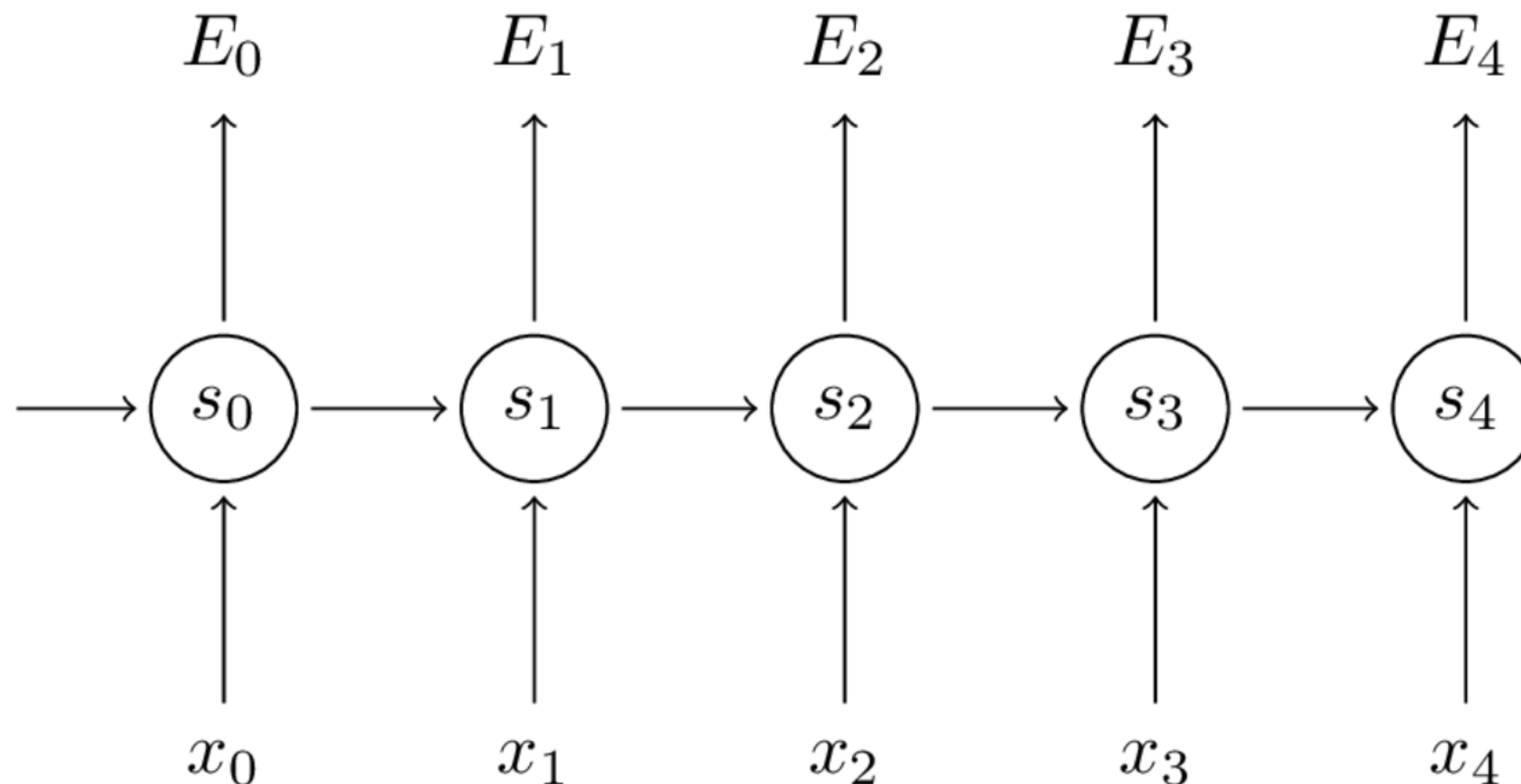- Use back propagation through time (BPTT)

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

$$= -\sum_t y_t \log \hat{y}_t$$



**[Denny Britz]**

# Back Propagation through Time

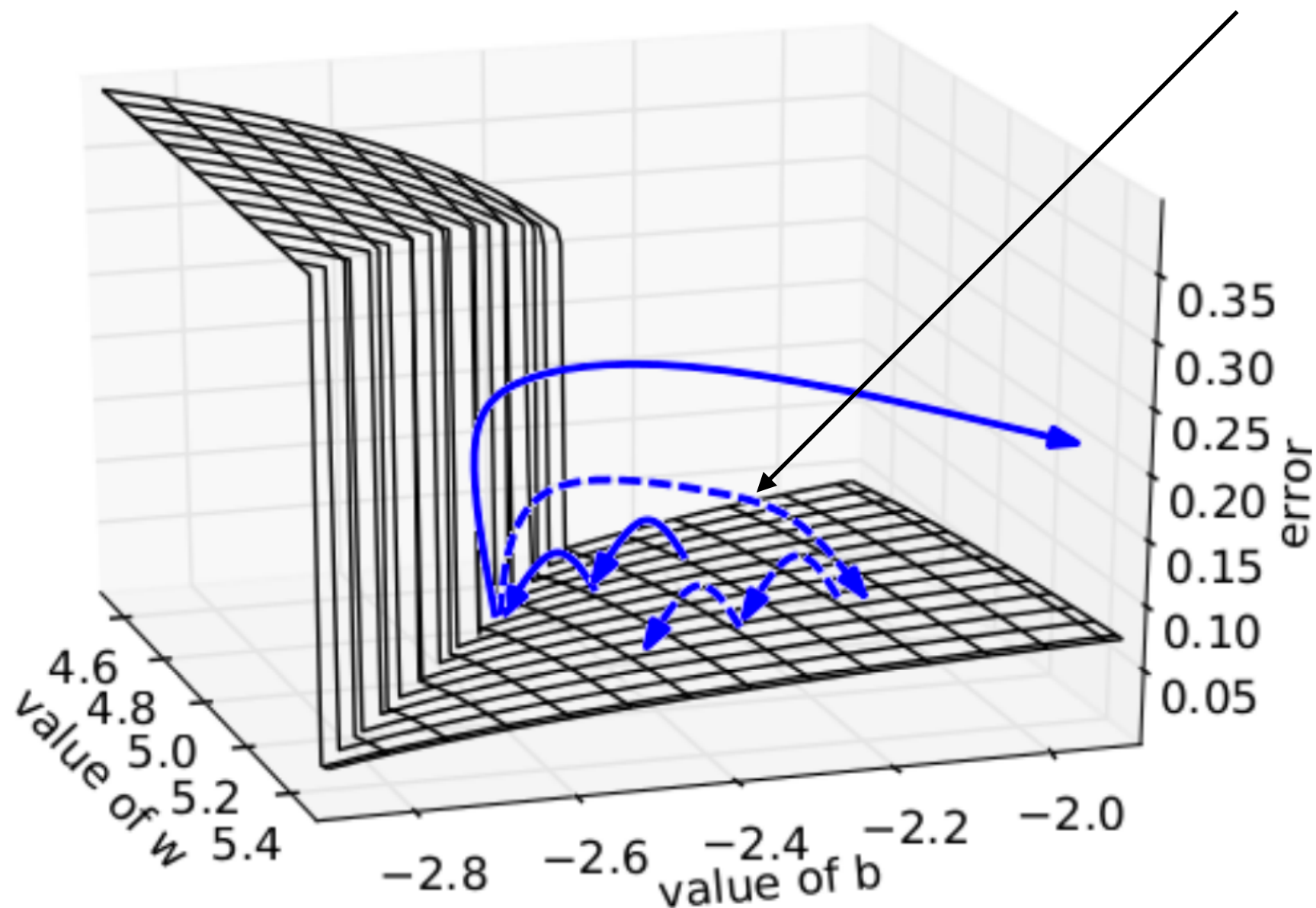$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W}.$$

# RNN Training Issues

- Exploding/Vanishing gradients
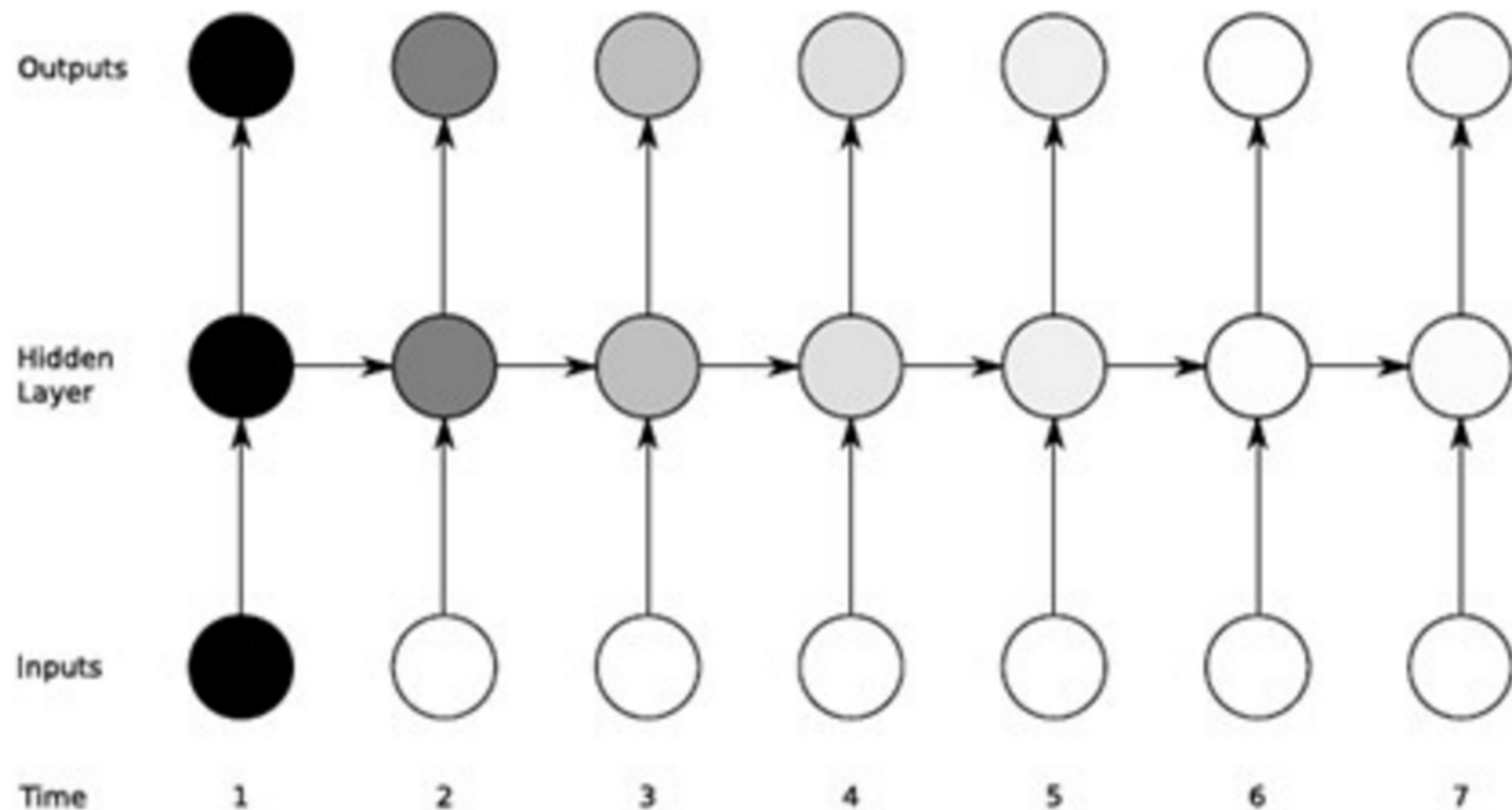
- Exploding/Vanishing activations

# Exploding Gradients

Solution: Gradient Clipping



**[Richard Socher]**

# Vanishing Gradients/ Activations



[Hochreiter, Schmidhuber]

# Why Training is Unstable

$$x^{(l)} = W^{(l-1)} y^{(l-1)} + b^{(l-1)}$$

$$y^{(l)} = f(x^{(l)})$$

Let the activation function $f(x) = \alpha x + \beta$,

$$Var\left(y^{(l)}\right) = \alpha^2 n_{l-1} \sigma_{l-1}^2 \left(Var\left(y^{(l-1)}\right) + \beta^2 I_{n_l}\right).$$

$$Var\left(\frac{\partial cost}{\partial y^{(l-1)}}\right) = \alpha^2 n_l \sigma_{l-1}^2 Var\left(\frac{\partial cost}{\partial y^{(l)}}\right).$$

**Variance of activations/gradients grows multiplicatively**

[Xu, Huang, Li]

# Interesting Question

- Are there modifications to an RNN such that it can combat these activations/gradient problems?
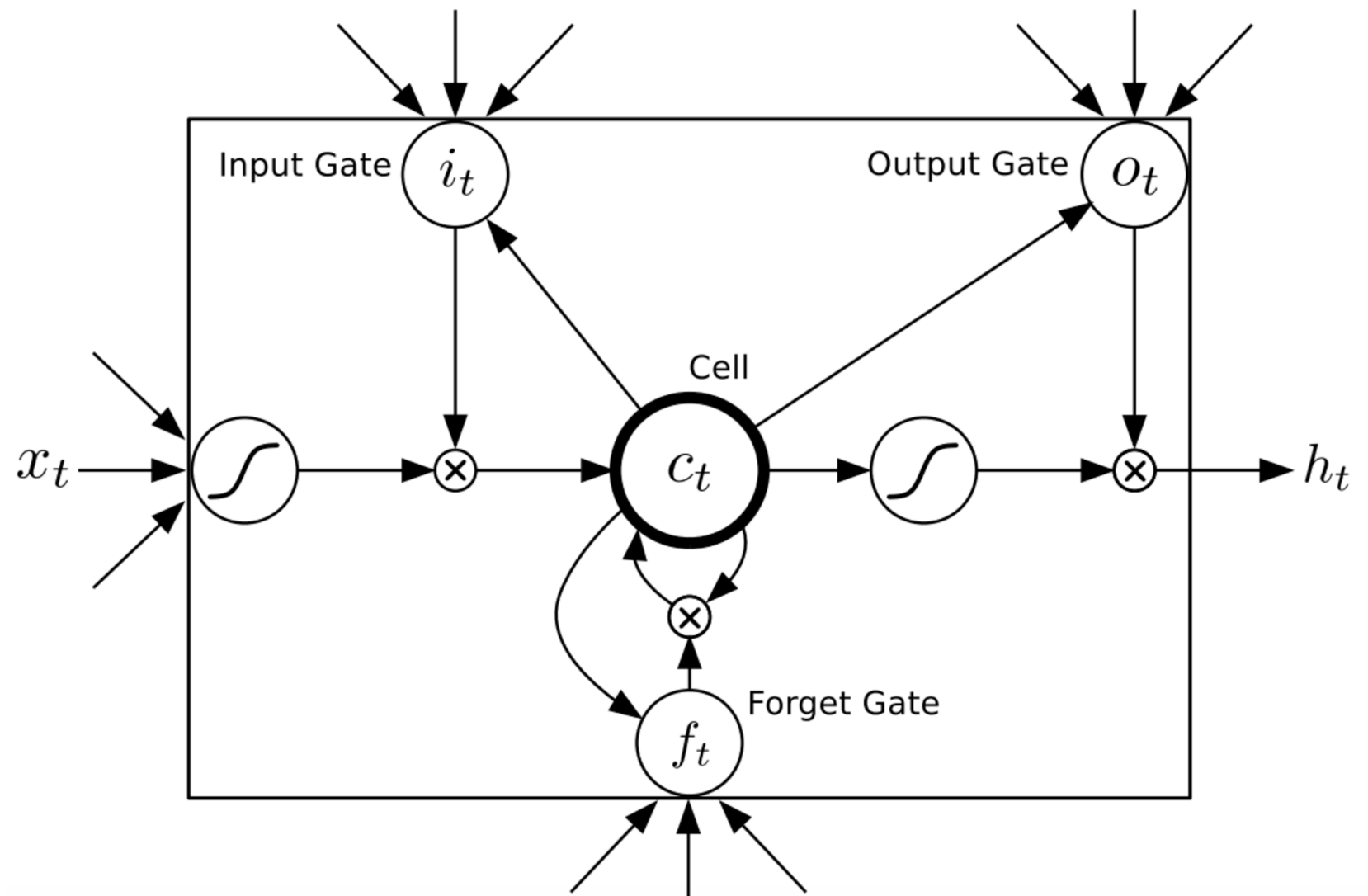
# RNNs with Longer Term Memory

# Motivation

- The need to remember certain events for arbitrarily long periods of time (Non-Markovian)

- The need to forget certain events

# Long Short Term Memory

- 3 gates

  - Input

  - Forget

  - Output



**[Zygmunt Z.]**

# LSTM Formulation

$$i_t = \sigma \left( W_{xi} x_t + W_{hi} h_{t-1} + W_{ci} c_{t-1} + b_i \right)$$

$$f_t = \sigma \left( W_{xf} x_t + W_{hf} h_{t-1} + W_{cf} c_{t-1} + b_f \right)$$

$$c_t = f_t c_{t-1} + i_t \tanh \left( W_{xc} x_t + W_{hc} h_{t-1} + b_c \right)$$

$$o_t = \sigma \left( W_{xo} x_t + W_{ho} h_{t-1} + W_{co} c_t + b_o \right)$$

$$h_t = o_t \tanh(c_t)$$

$$y_t = W_{ho} h_t + b_o$$

**[Alex Graves, Navdeep Jaitly]**
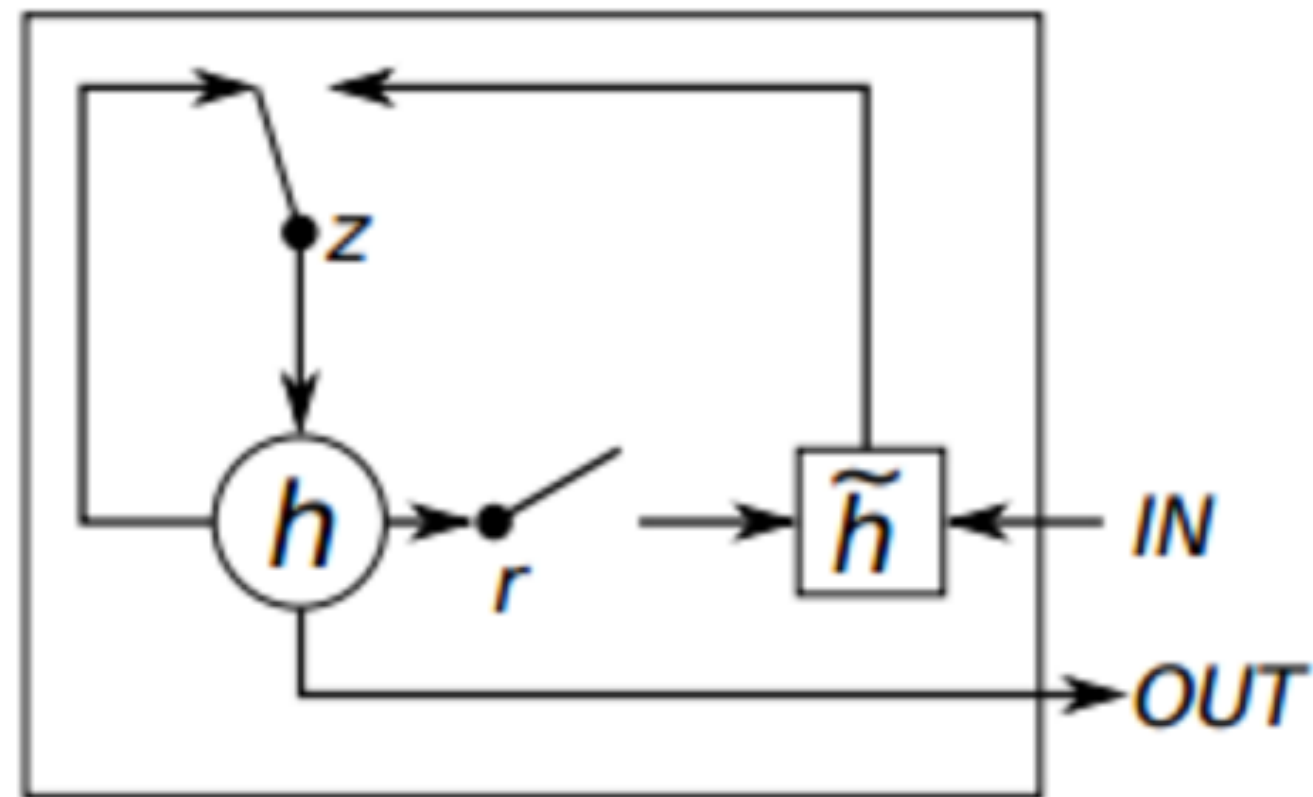
# Preserving Gradients

# Gated Recurrent Unit

- 2 gates

  - Reset

    - Combine new input with previous memory

  - Update

    - How long the previous memory should stay



**[Zygmunt Z.]**

# GRU Formulation

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$

$$r = \sigma(x_t U^r + s_{t-1} W^r)$$

$$h = tanh(x_t U^h + (s_{t-1} \circ r) W^h)$$

$$s_t = (1 - z) \circ h + z \circ s_{t-1}$$
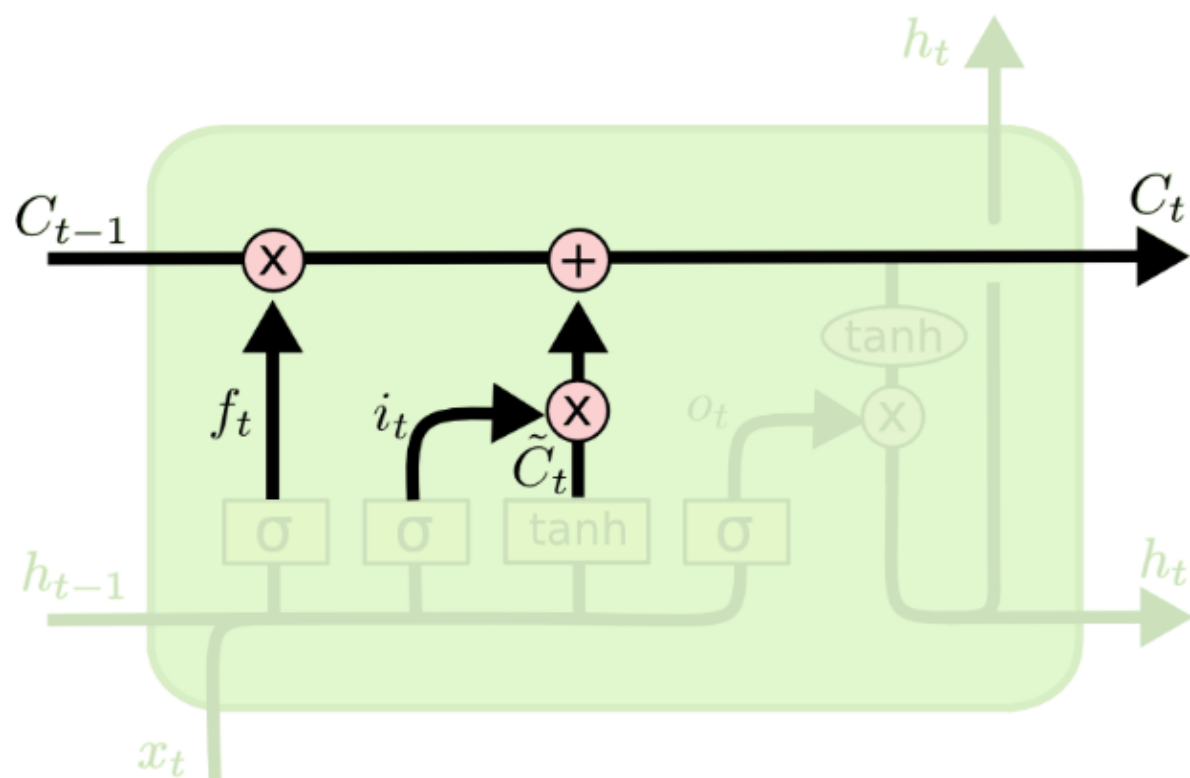
**[Danny Britz]**

# LSTM & GRU Benefits

- Remember for longer temporal durations

  - RNN has issues for remembering longer durations

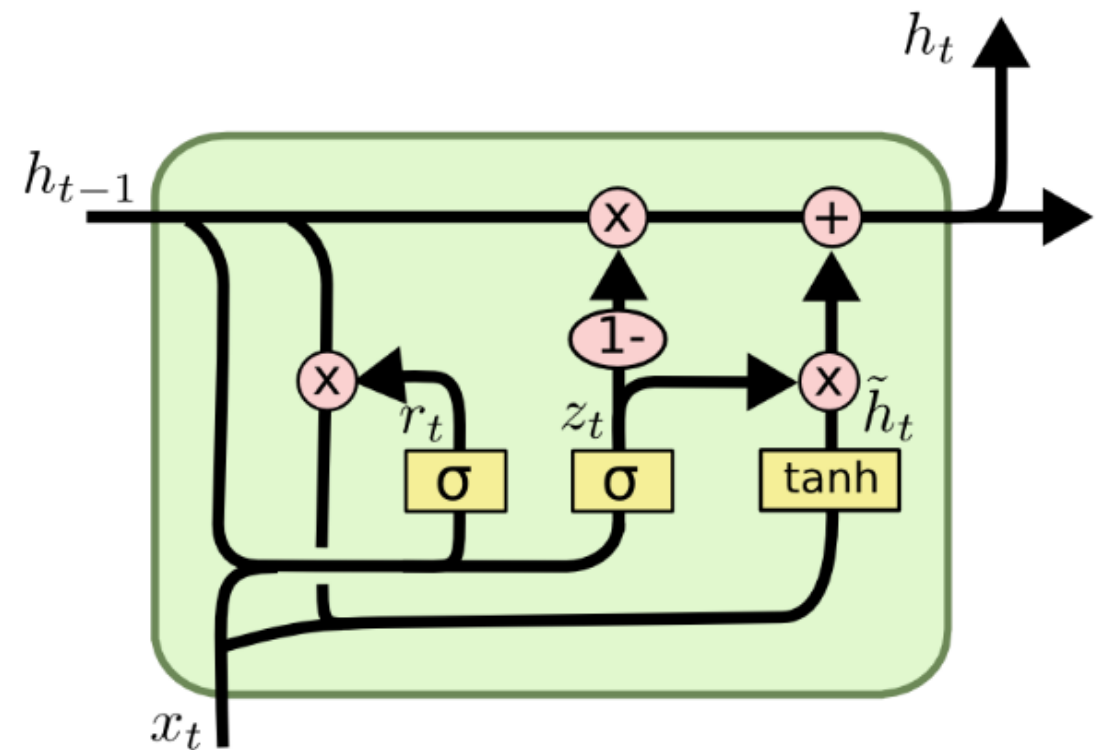- Able to have feedback flow at different strengths depending on inputs

# Differences between LSTM & GRU

- GRU has two gates, while LSTM has three gates

- GRU does not have internal memory

- GRU does not use a second nonlinearity for computing the output

# Visual Difference of LSTM & GRU



LSTM

GRU

# LSTM vs GRU Results

| | | | tanh | GRU | LSTM |
|---|---|---|---|---|---|
| **Music Datasets** | Nottingham | train | 3.22 | 2.79 | 3.08 |
| | | test | **3.13** | 3.23 | 3.20 |
| | JSB Chorales | train | 8.82 | 6.94 | 8.15 |
| | | test | 9.10 | **8.54** | 8.67 |
| | MuseData | train | 5.64 | 5.06 | 5.18 |
| | | test | 6.23 | **5.99** | 6.23 |
| | Piano-midi | train | 5.64 | 4.93 | 6.49 |
| | | test | 9.03 | **8.82** | 9.03 |
| **Ubisoft Datasets** | Ubisoft dataset A | train | 6.29 | 2.31 | 1.44 |
| | | test | 6.44 | 3.59 | **2.70** |
| | Ubisoft dataset B | train | 7.61 | 0.38 | 0.80 |
| | | test | 7.62 | **0.88** | 1.26 |

**[Chung, Gulcehre, Cho, Bengio]**

# Other Methods for Stabilizing RNN Training

# Why Training is Unstable

$$x^{(l)} = W^{(l-1)}y^{(l-1)} + b^{(l-1)}$$

$$y^{(l)} = f(x^{(l)})$$

*Let the activation function* $f(x) = \alpha x + \beta,$

$$Var\left(y^{(l)}\right) = \alpha^2 n_{l-1}\sigma_{l-1}^2 \left(Var\left(y^{(l-1)}\right) + \beta^2 I_{n_l}\right).$$

$$Var\left(\frac{\partial cost}{\partial y^{(l-1)}}\right) = \alpha^2 n_l \sigma_{l-1}^2 Var\left(\frac{\partial cost}{\partial y^{(l)}}\right).$$

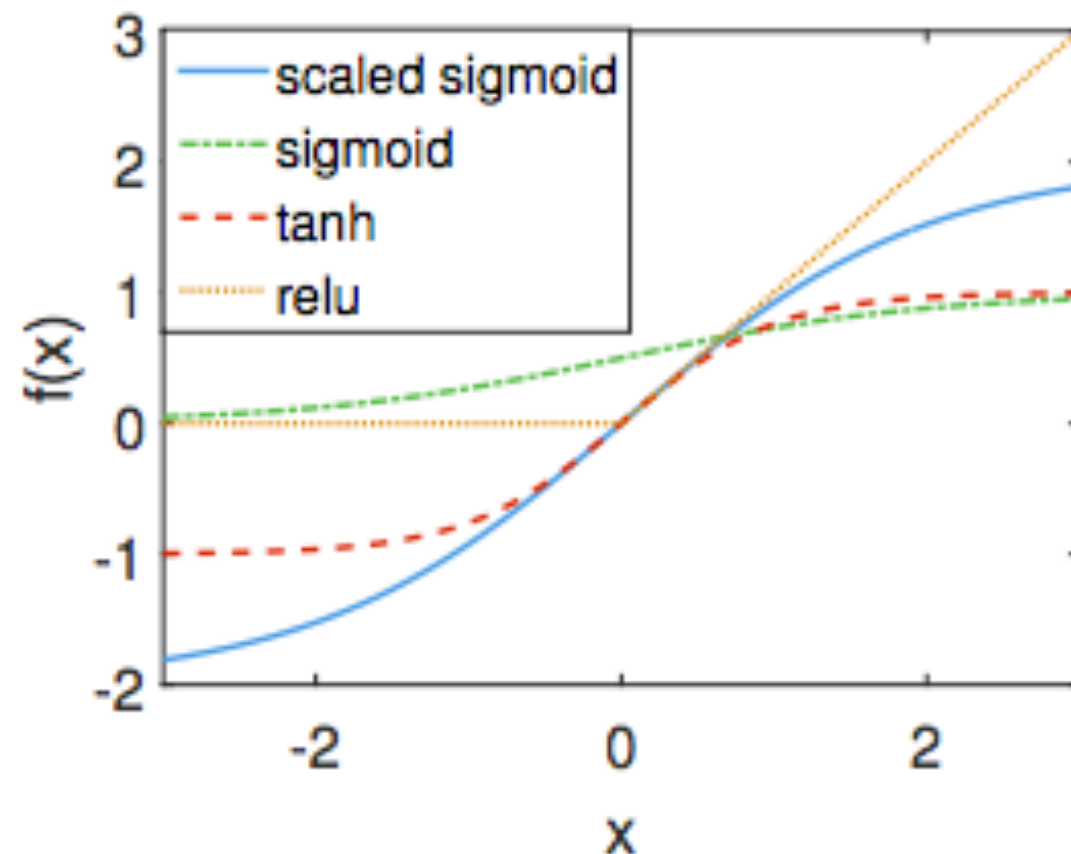**Variance of activations/gradients grows multiplicatively**

**[Xu, Huang, Li]**

# Stabilizing Activations & Gradients

$$\text{Var}\left(y^{(l)}\right) = \text{Var}\left(y^{(l-1)}\right) \quad \text{and} \quad \text{Var}\left(\frac{\partial \text{cost}}{\partial y^{(l)}}\right) = \text{Var}\left(\frac{\partial \text{cost}}{\partial y^{(l-1)}}\right);$$

$$n_l \sigma_{l-1}^2 \approx 1 \quad \text{and} \quad n_{l-1} \sigma_{l-1}^2 \approx 1;$$

We want $\alpha = 1$ and $\beta = 0$.

[Xu, Huang, Li]

# Taylor Expansions of Different Activation Functions



$$\text{sigmoid}(x) = \frac{1}{2} + \frac{x}{4} - \frac{x^3}{48} + O(x^5)$$

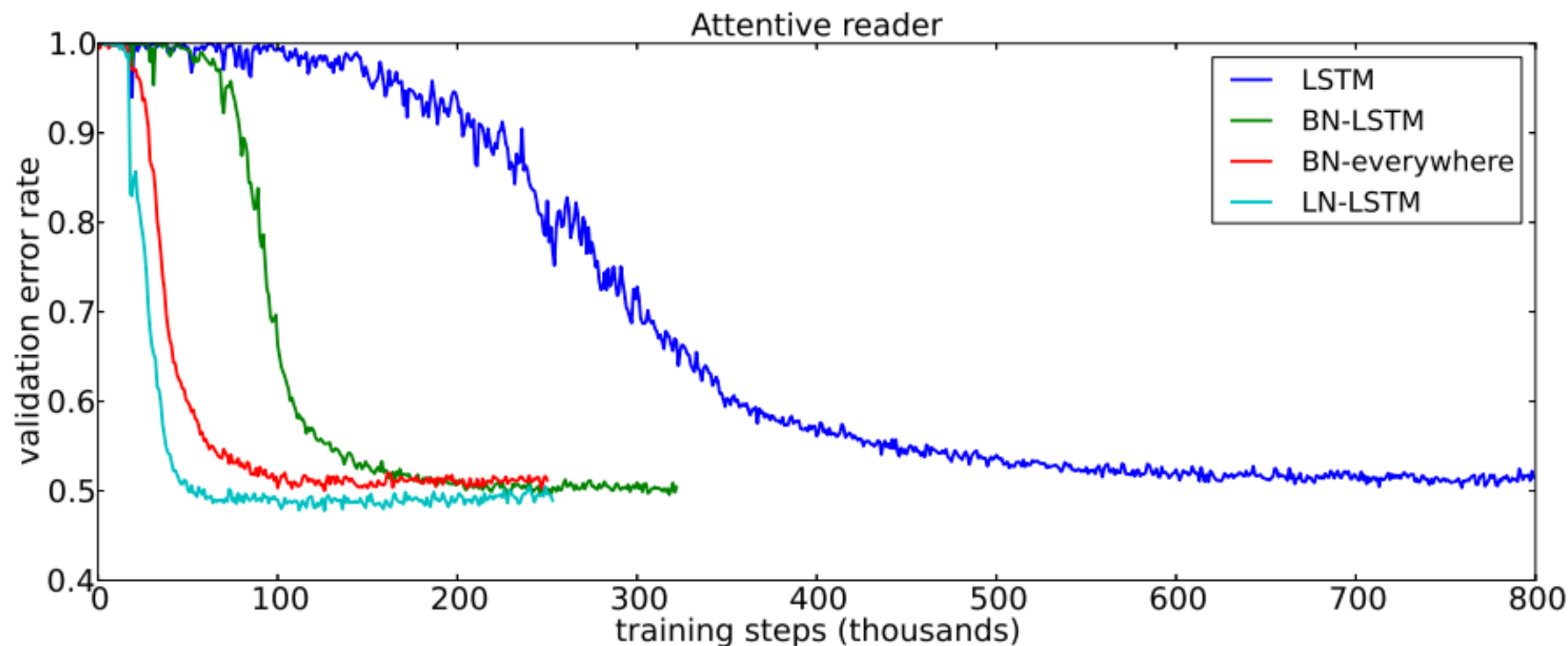$$\tanh(x) = 0 + x - \frac{x^3}{3} + O(x^5)$$

$$\text{relu}(x) = 0 + x \quad \text{for } x \geq 0.$$

[Xu, Huang, Li]

# Layer Normalization

- Similar to batch normalization

  - Apply it to RNNs to stabilize the hidden state dynamics

$$\mathbf{h}^t = f\left[\frac{\mathbf{g}}{\sigma^t} \odot (\mathbf{a}^t - \mu^t) + \mathbf{b}\right] \qquad \mu^t = \frac{1}{H}\sum_{i=1}^{H} a_i^t \qquad \sigma^t = \sqrt{\frac{1}{H}\sum_{i=1}^{H}(a_i^t - \mu^t)^2}$$
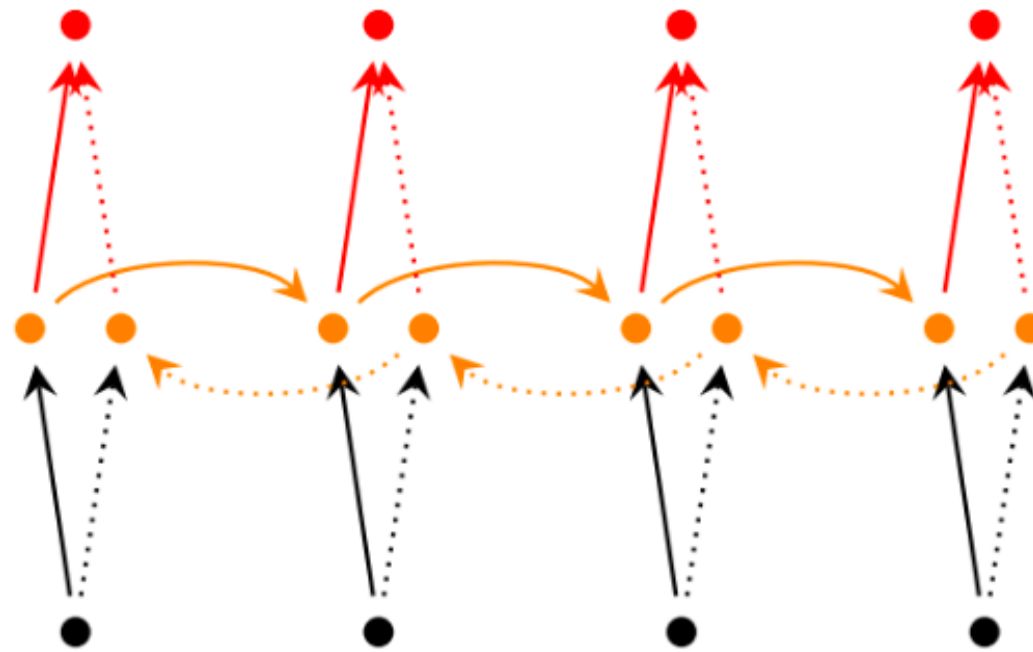
**[Ba, Kiros, Hinton]**

# Layer Normalization Results



[Ba, Kiros, Hinton]

# Variants of RNNs

# Bidirectional RNNs

- The output at time *t* does not depend on previous time steps but also the future

  - Two RNNs stacked on top of each other

# Deep RNNs

- Stack them on top of each other

  - The output of the previous RNN is the input to the next one

# The Power of RNNs: Understanding and Visualizing

# The Effectiveness of an RNN

```c
#define REG_PG      vesa_slot_addr_pack
#define PFM_NOCOMP  AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)       (e)
#define emulate_sigs()  arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0));   \
  if (__type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
          pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
  PUT_PARAM_RAID(2, sel) = get_state_state();
  set_pid_sum((unsigned long)state, current_state_str(),
          (unsigned long)-1->lr_full; low;
}
```

[Andrej Karpathy]

# The Effectiveness of an RNN

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 servicious, 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/7754800786d17551963s89.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]**

# The Effectiveness of an RNN



[Andrej Karpathy]

# The Effectiveness of an RNN

Trained on *War & Peace*

Iteration: 100

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

Iteration: 300

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

Iteration: 2000

```
"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.
```

# Visualize the Neurons of an RNN

# Visualize the Neurons of an RNN



Cell sensitive to position in line:

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."
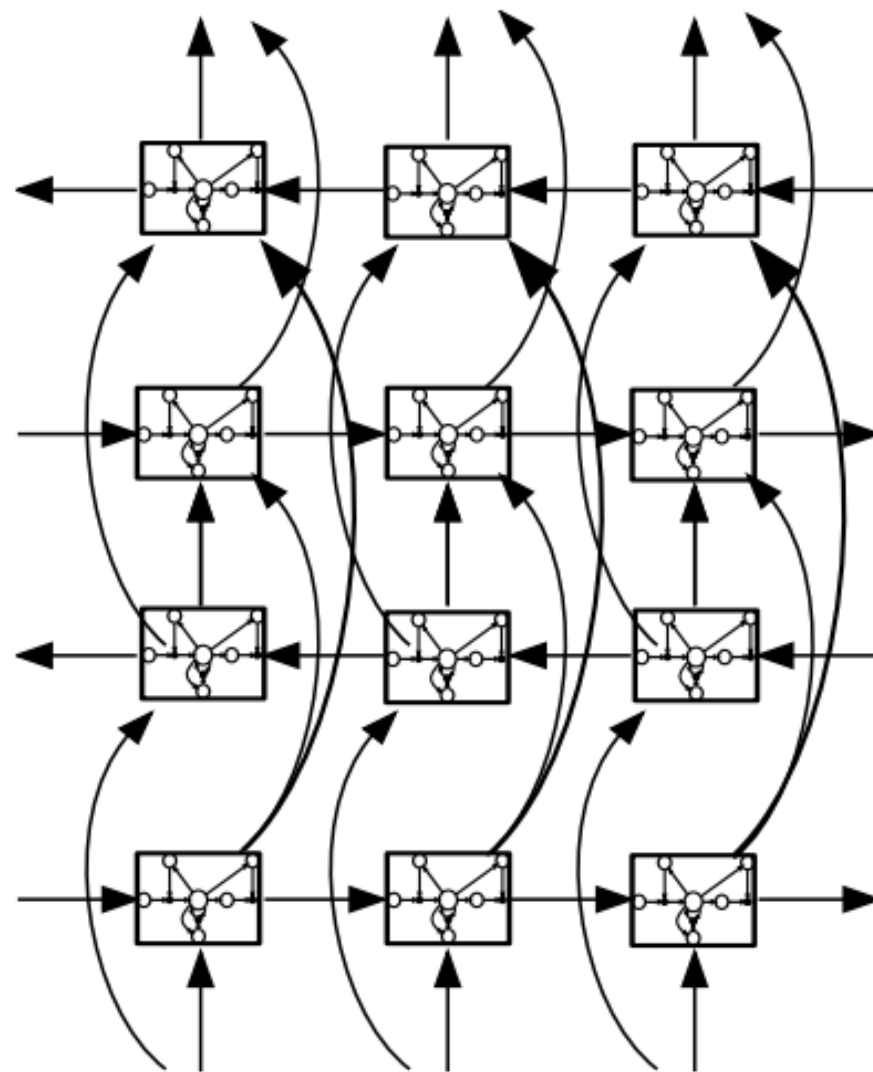
**[Andrej Karpathy]**

# Applications

# RNN Applications

- Speech Recognition

- Natural Language Processing

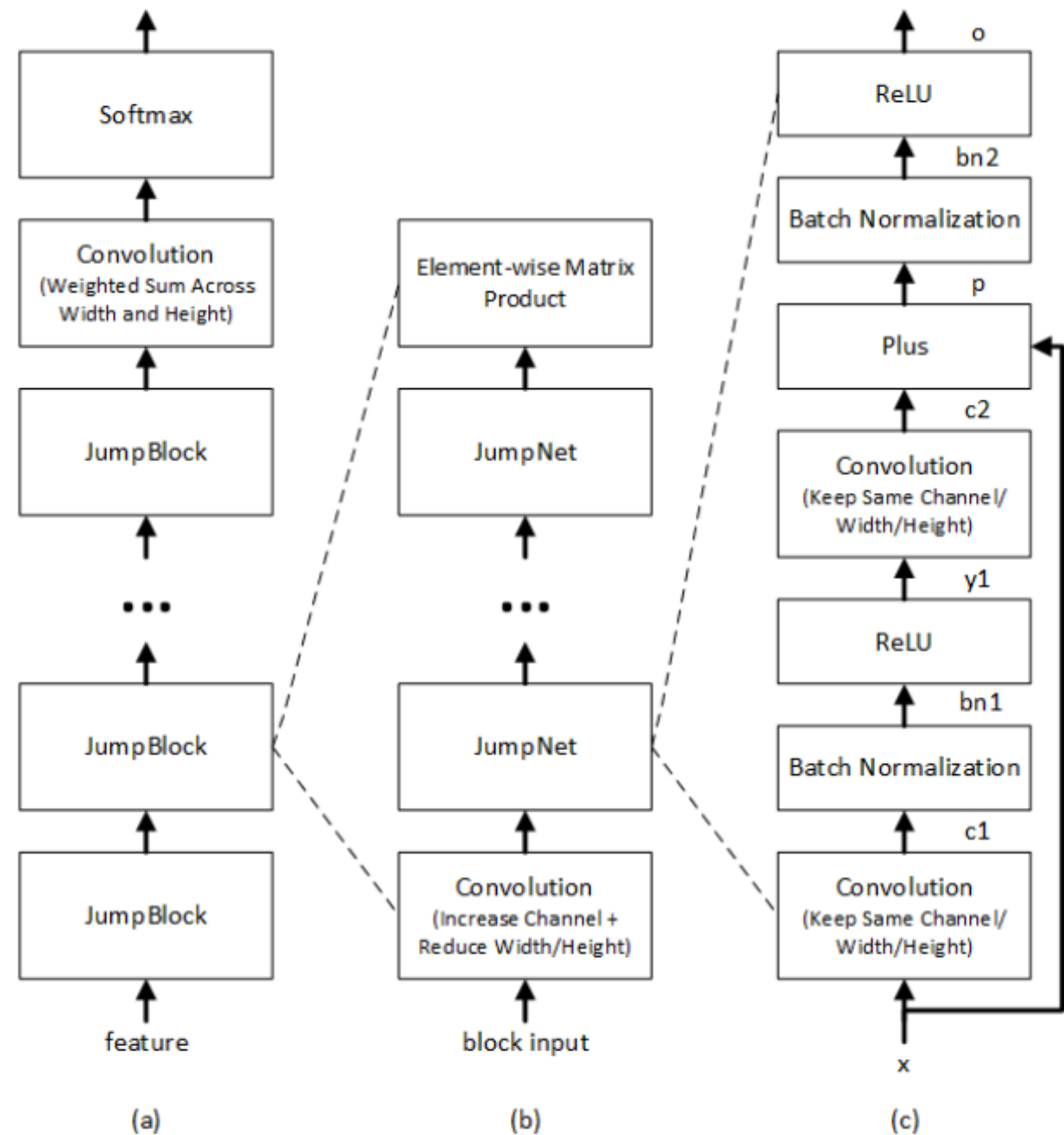- Action Recognition

- Machine Translation

- Many more to come

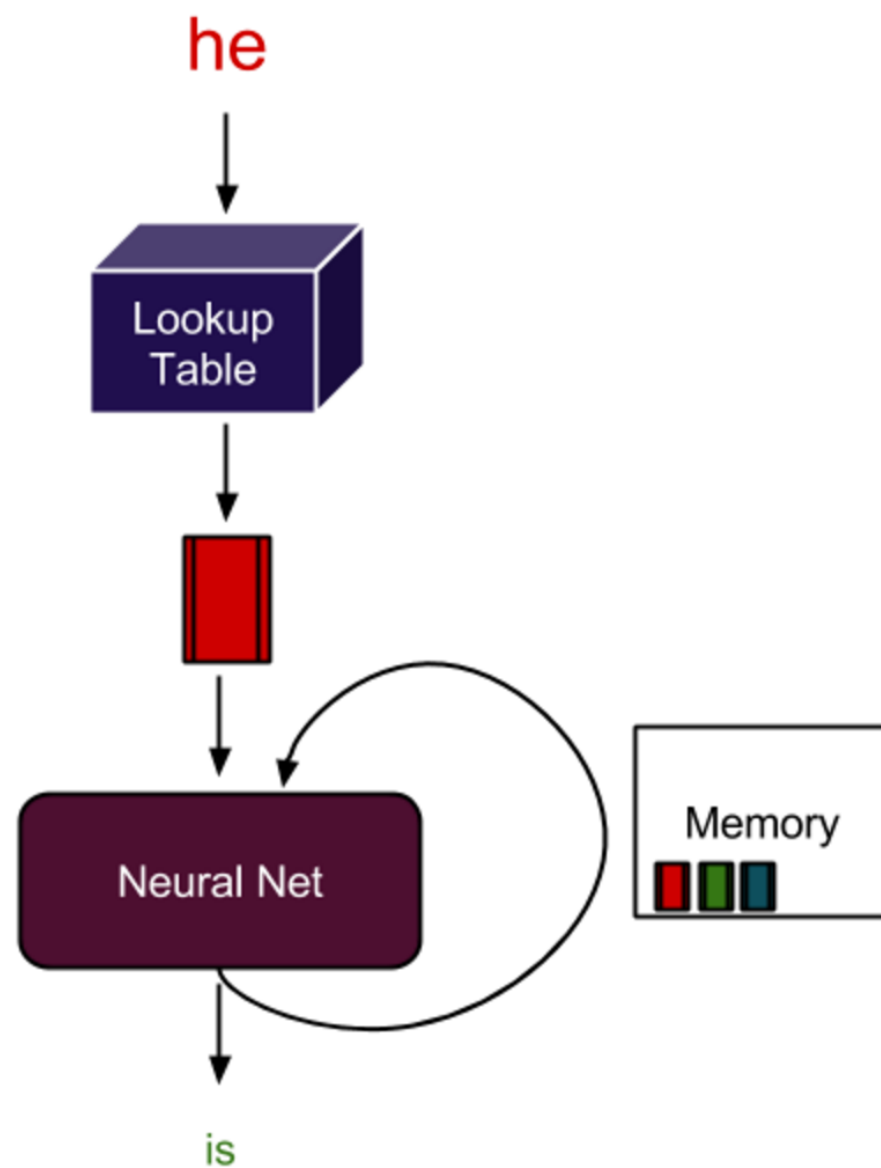# Speech Recognition

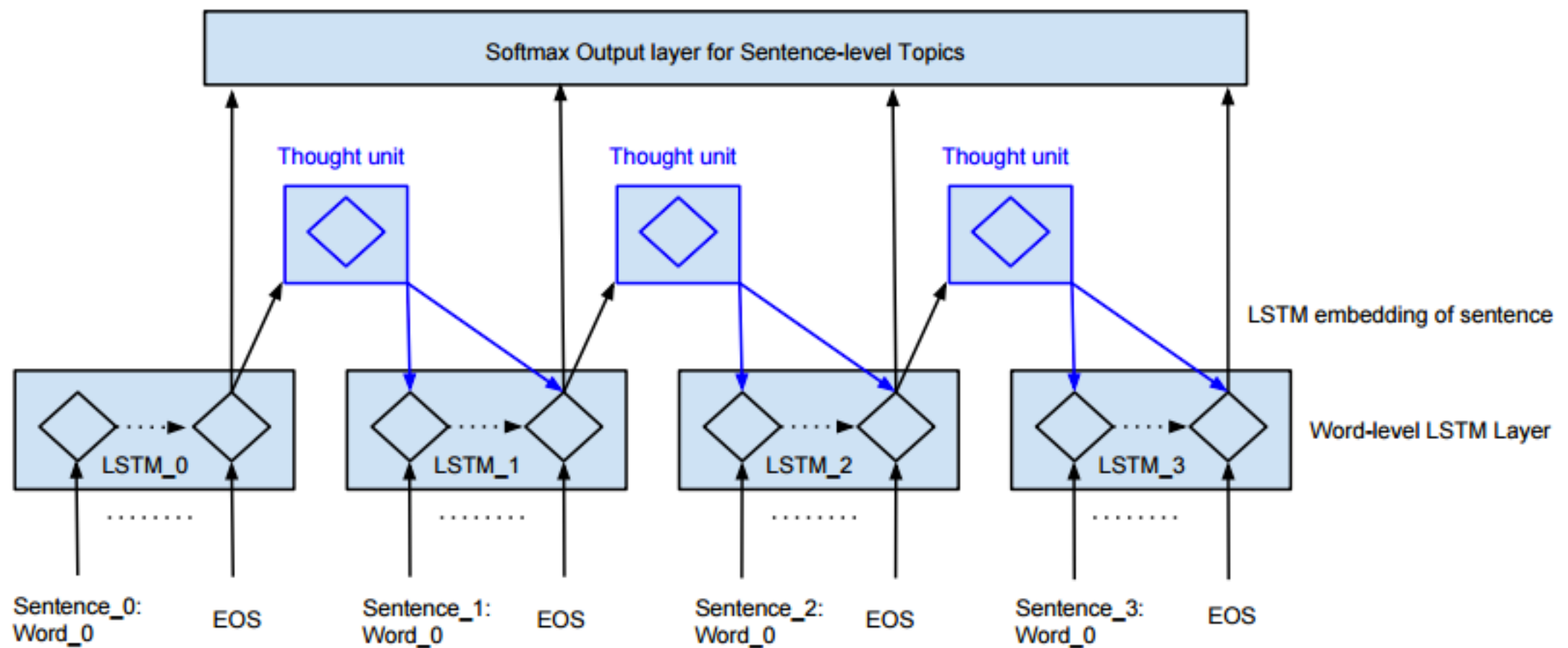- Deep Bidirectional LSTM

# Conversational Speech Recognition

- Achieving human parity



[Xiong et al.]
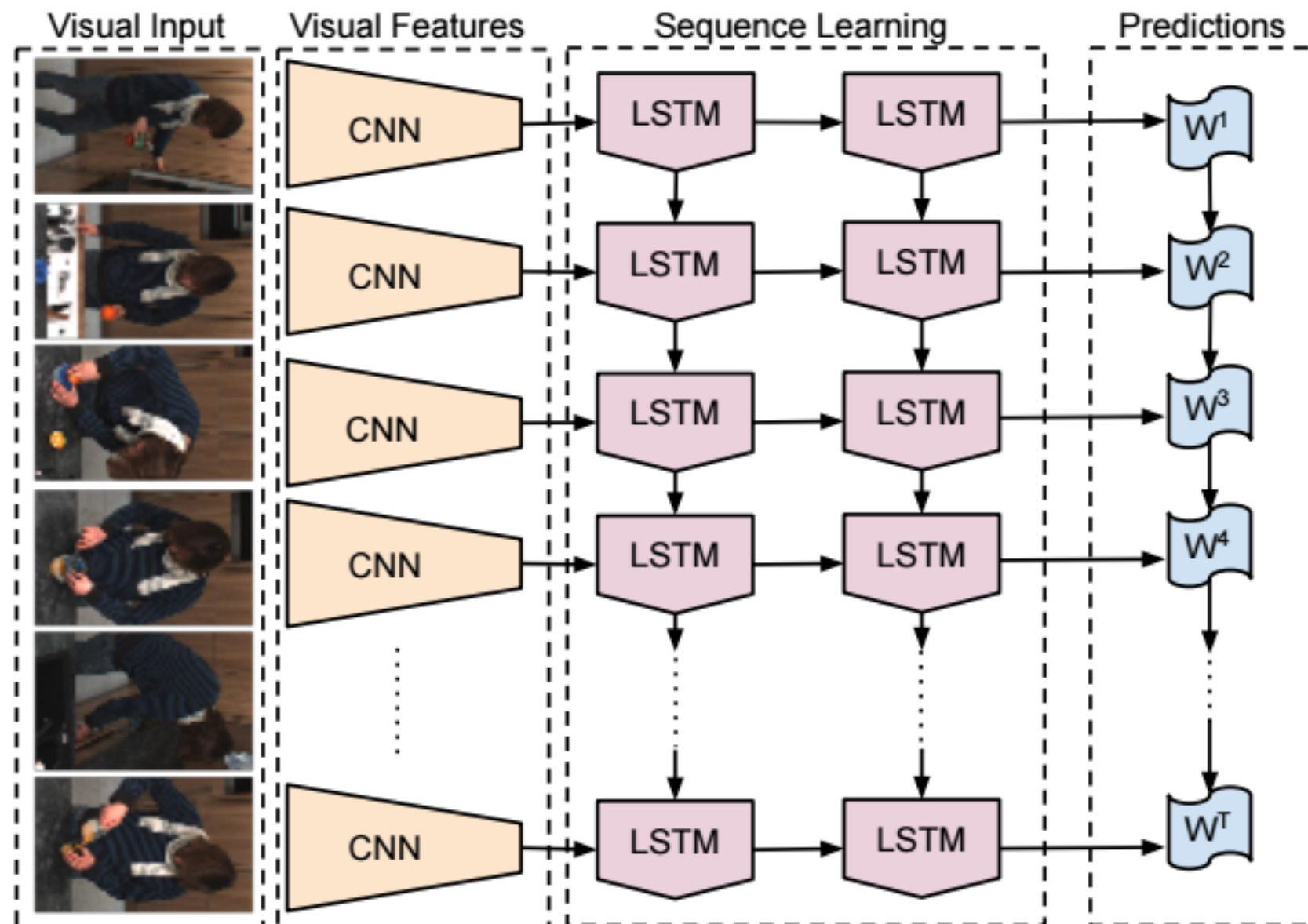
# Natural Language Processing

# Contextual LSTM for NLP Tasks
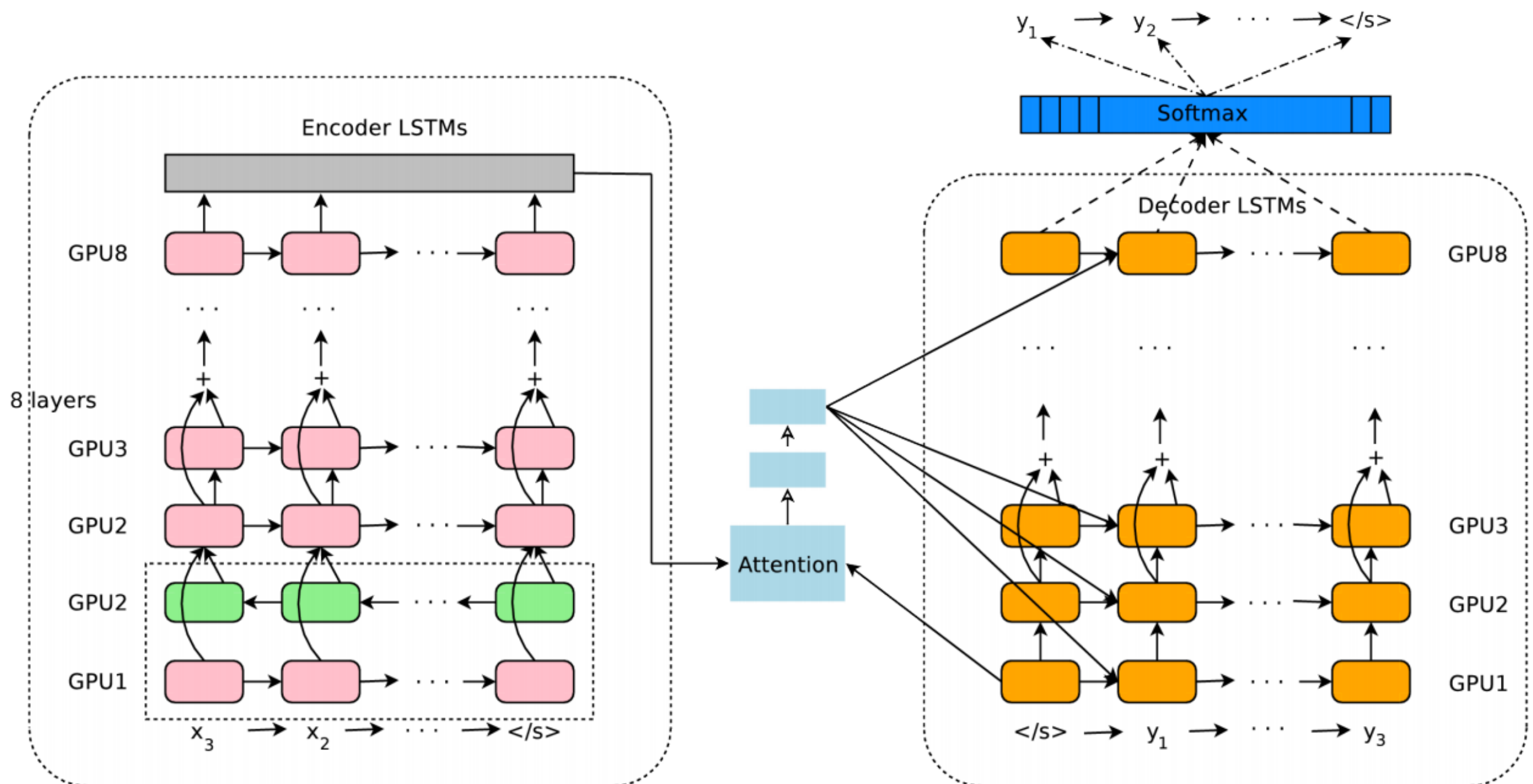


[Ghost et al.]

# Action Recognition

- Long-term Recurrent Convnet



[Donahue et al.]

# Google's Neural Machine Translation System
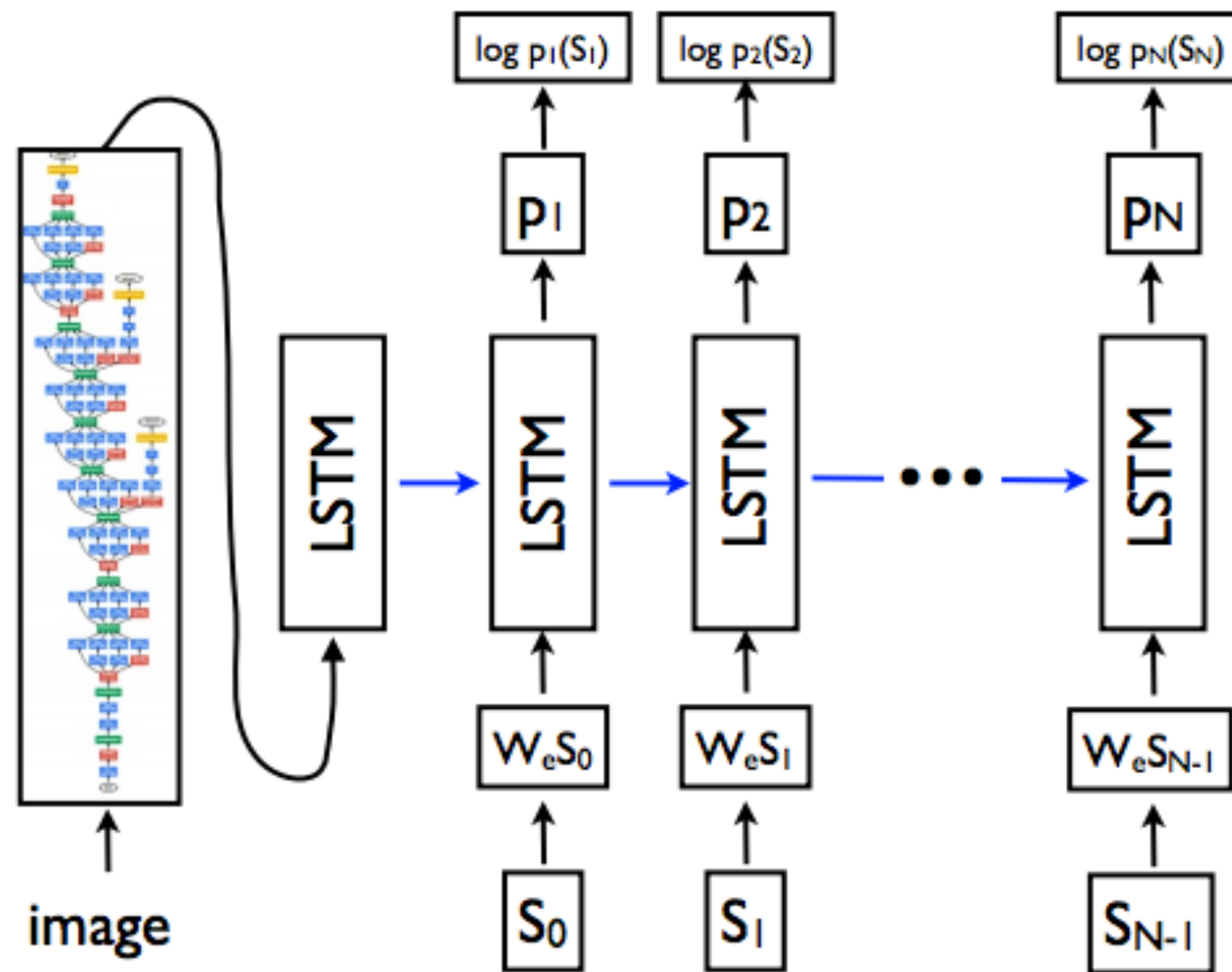


[Yonghui Wu et al.]

# Image Captioning



[Vinyals, Toshev, Bengio, Erhan]

# Image Captioning



[Vinyals, Toshev, Bengio, Erhan]

# Object Tracking



[Ning, Zhang, Huang, He, Wang]

# Neural Turing Machines



Memory is an array of vectors

Network A writes and reads from this memory each step

# WaveNet



[van den Oord et al.]

# DoomBot

- Doom Competition

  - Facebook won 1st place (F1)

  - https://www.youtube.com/watch?v=94EPSjQH38Y

# $\texttt{ODE2RNN}$: Parameter Estimation for Systems of Ordinary Differential Equations



**ode2rnn framework**

**A**

Dynamical Model (ODEs)

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t); \theta_M)$$
$$\mathbf{x}(t = 0) = \mathbf{x}_0$$
$$\mathbf{o}(t) = Observe(\mathbf{x}(t))$$

Numerical Simulation

$$\mathbf{x}_{j+1} = \mathbf{x}_j + \sum_{l=1}^{m} w^{(l)} f(\mathbf{x}_j^{(l)}, \theta_M)\Delta t + o(\Delta t^m)$$

RNN

$$\mathbf{x}_{j+1} = g(\mathbf{x}_j; \theta_{RNN})$$

**ODE-based RNN Recurrent Update**

Loss Function

$$\ell(\theta_M; \mathcal{D}) \equiv \sum_{k=1}^{N} \|\mathbf{o}_k - \mathbf{x}_k\|_2^2$$

Estimated Parameters

Optimization
- gradient-based
- L-BFGS,
- Nelder-Mead, ...

Train

$$\hat{\theta}_M$$

Parameter Gradients

Scientific Questions

Experimental Data

$$\mathbf{o}_1, \ldots, \mathbf{o}_N$$

**B**

RK4

$$\mathbf{x}_j \quad \boxed{f} \rightarrow \mathbf{x}_j^{(1)}$$
$$\boxed{f} \rightarrow \mathbf{x}_j^{(2)}$$
$$\boxed{f} \rightarrow \mathbf{x}_j^{(3)} \quad \mathbf{x}_{j+1}$$
$$\boxed{f} \rightarrow \mathbf{x}_j^{(4)}$$

$$\mathbf{h}_j \longrightarrow \sigma \longrightarrow \mathbf{h}_{j+1}$$