

Recurrent Neural Networks

Zachary Lipton & Henry Chai

10701 — November 13th

Recurrent Neural Networks

Zachary Lipton & Henry Chai

10701 — November 13th



Bridging from CNNs + Images to RNNs + Text

Recap: Multi-Modal Pretraining

What exactly goes on inside this “text encoder”?

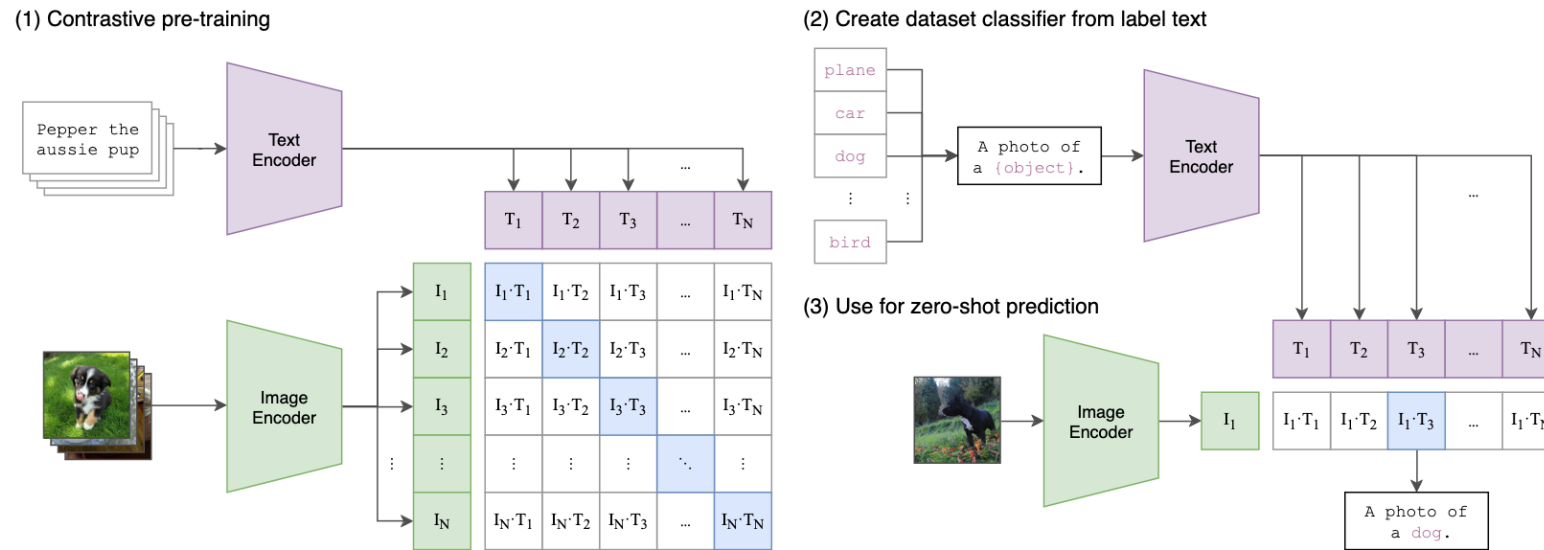
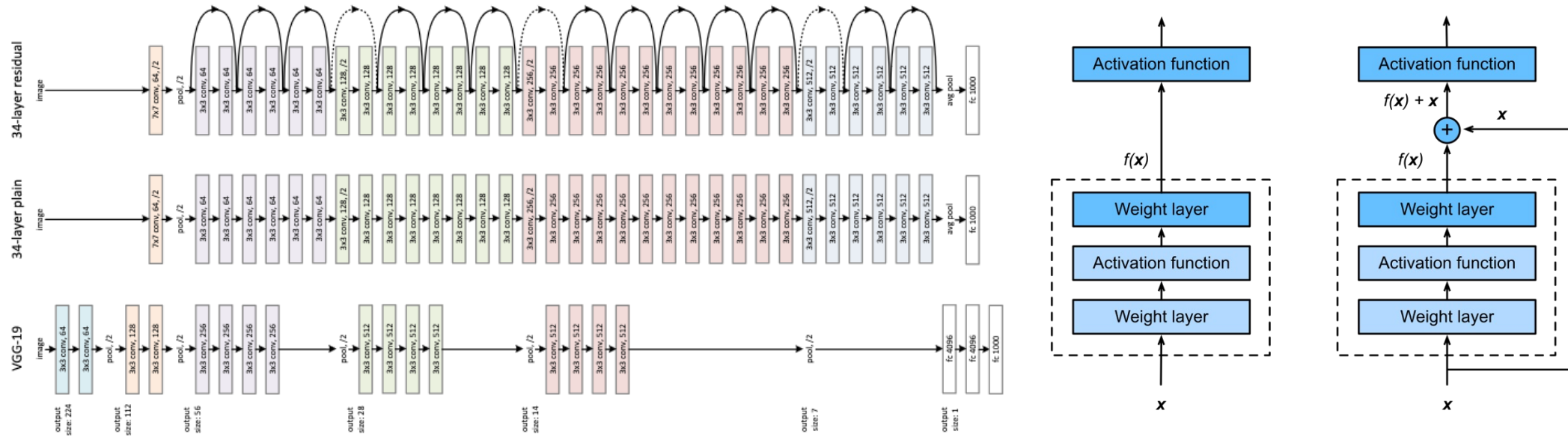


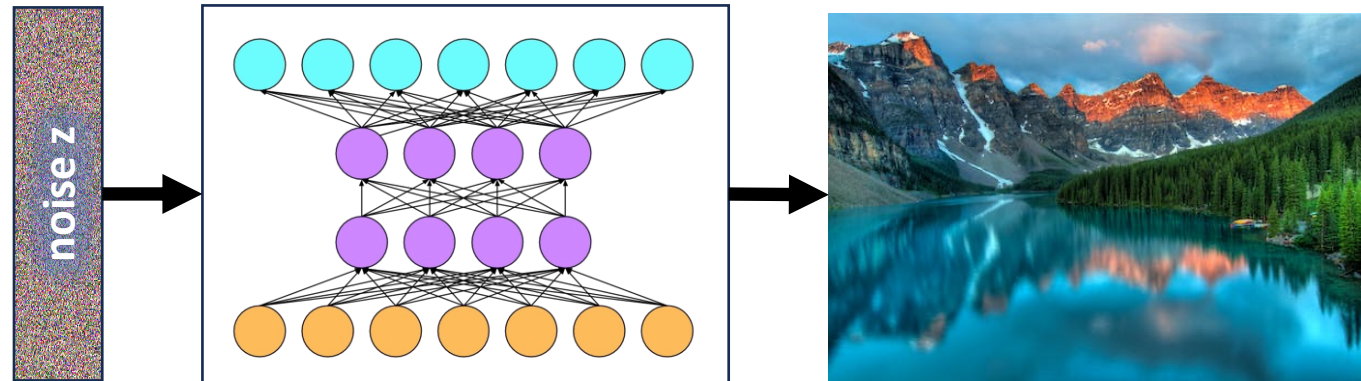
Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset’s classes.

ResNet as Gradient “Superhighway”



["Deep Residual Learning for Image Recognition" He et al. 2015](#)

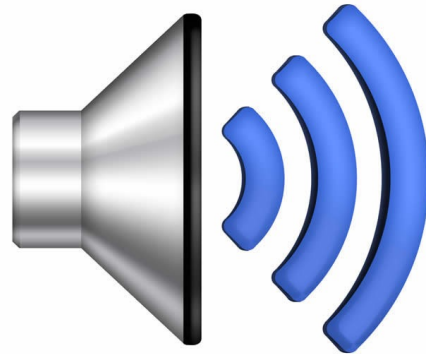
Question: How to Learn Distributions on Discrete Domains?



p ("Born and raised in North Carolina, Coltrane moved to Philadelphia after graduating from high school, where he studied music. Working in the bebop and hard bop idioms early in his career, Coltrane helped pioneer the use of modes and was one of the players at the forefront of free jazz. He led at least fifty recording sessions and appeared on many albums by other musicians, including trumpeter Miles Davis and pianist Thelonious Monk. Over the course of his career, Coltrane's music took on an increasingly spiritual dimension, as exemplified on his most acclaimed album A Love Supreme (1965) and others.")

Sequentially Structured Data

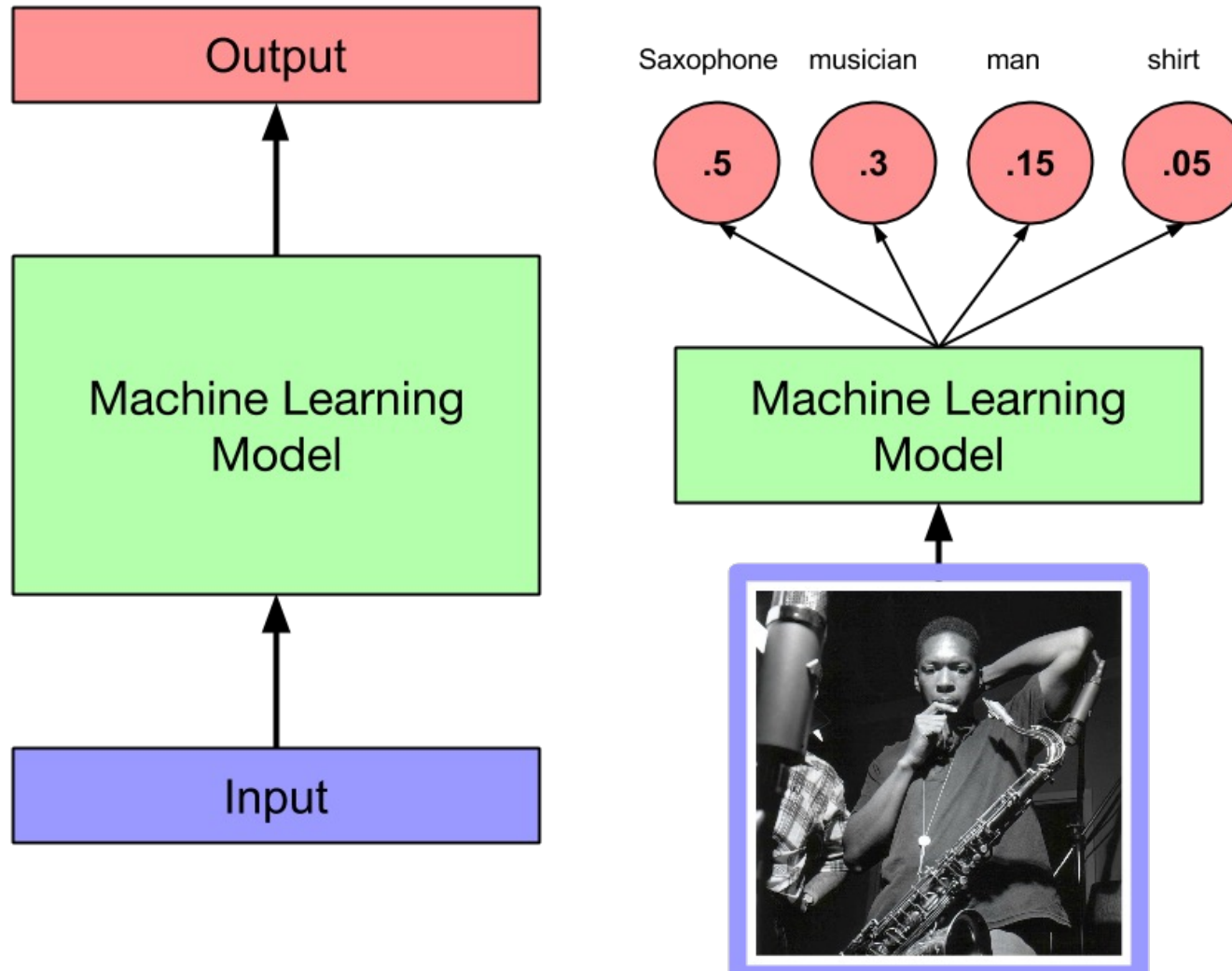
- Have been focusing on data with fixed-length inputs
- But how should we deal with video, text, audio, time series data?



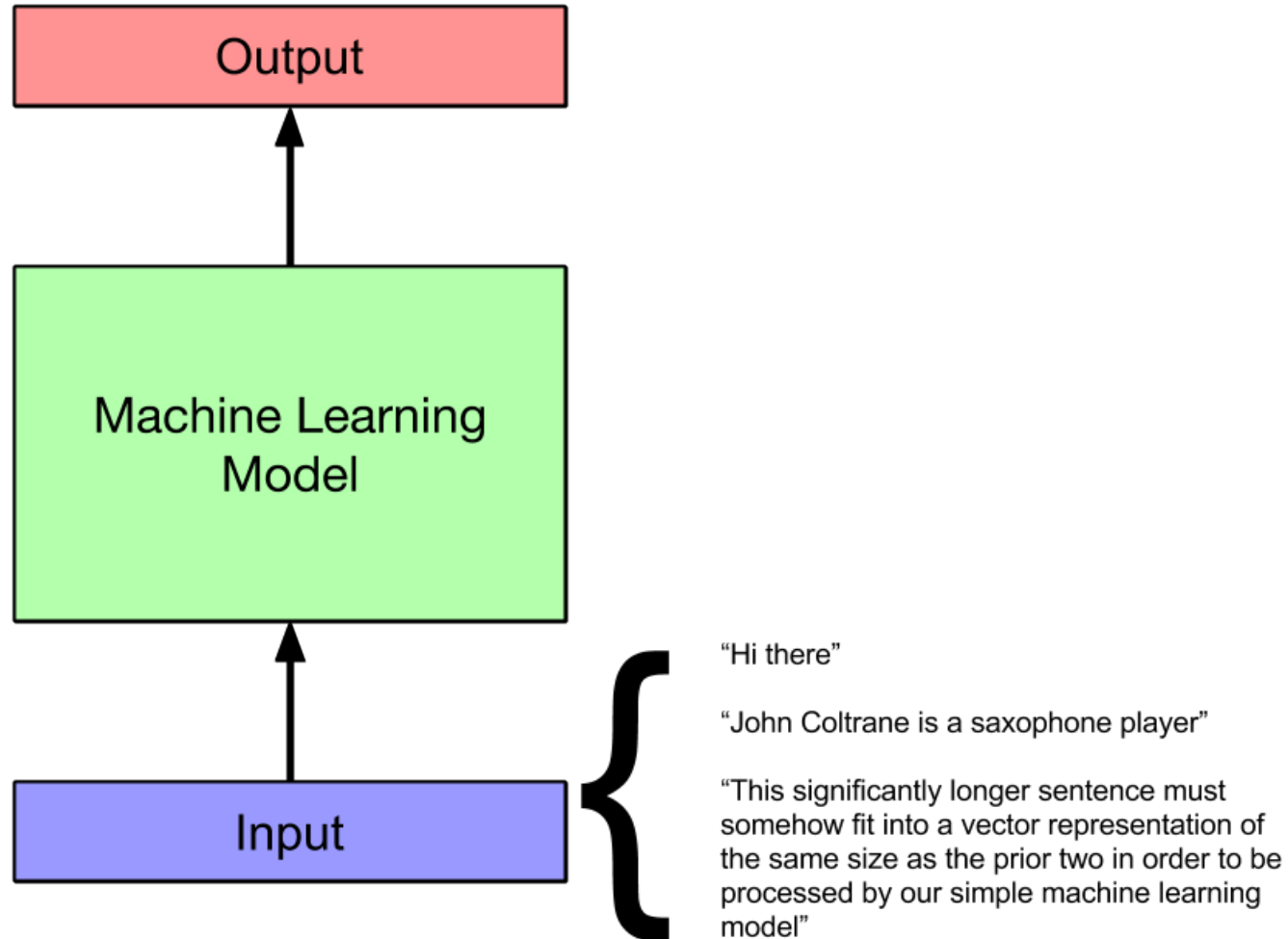
Enter Hamlet.
Pol. I heare him comming, with-draw my Lord.
Ham. To be, or not to be, that is the question,
Whether tis nobler in the minde to suffer
The slings and arrowes of outrageous fortune,
Or to take Armes against a sea of troubles,
And by opposing, end them, to die to sleepe
No more, and by a sleepe, to say we end
The hart-ake, and the thousand naturall shocks
That flesh is heire to; tis a confumation
Deuoutly to be wisht to die to sleepe,
To sleepe, perchance to dreame, I there's the rub,
For in that sleepe of death what dreames may come
When we haue shuffled off this mortall coyle
Must giue vs pause, there's the respect
That makes calamitie of so long life:



Feedforward Nets work for Fixed-Size Data



How to Handle Raw Text as Input (or output)?



Traditional Approach for Text Classification

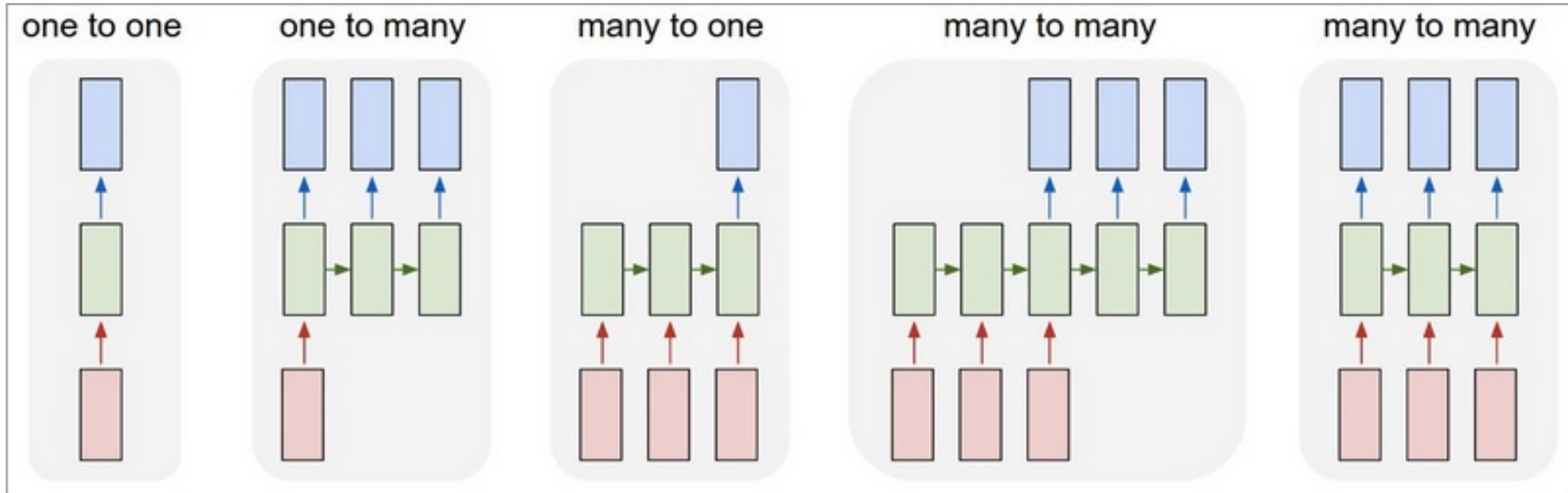
Bag-of-words representations

- Represent each document as vector of word counts $x \in \mathbb{R}^{|V|}$
- Dimension $|V|$ = vocab size
- Loses all information about word order.
- Reasonably effective for spam detection but cannot differentiate:
 - i. “Scientist exterminated by raging virus!”
 - ii. “Virus exterminated by raging scientist!”
- Limited ability to tackle complex language tasks (e.g., translation)
- No clear output mechanism

Complex Object → Sequence of Simple Ones

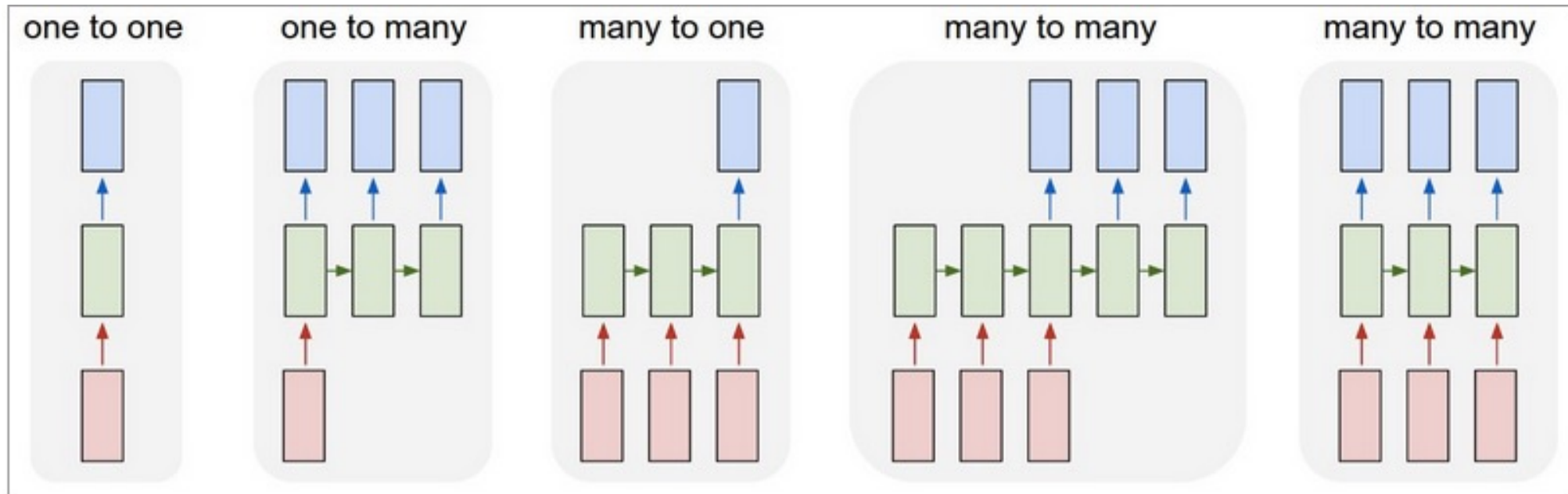
- Hard to represent sentence as fixed-length vector
- Easier to represent characters (or words) as fixed-length vectors
- Breaking up the inputs: ingest a word at each step
- Break up the outputs: predict next token at each step

Archetypal Sequence Learning Tasks



[slide credit: Andrej Karpathy](#)

Archetypal Sequence Learning Tasks



Simple prediction tasks: fixed-length input & output

Image captioning

Sentiment analysis

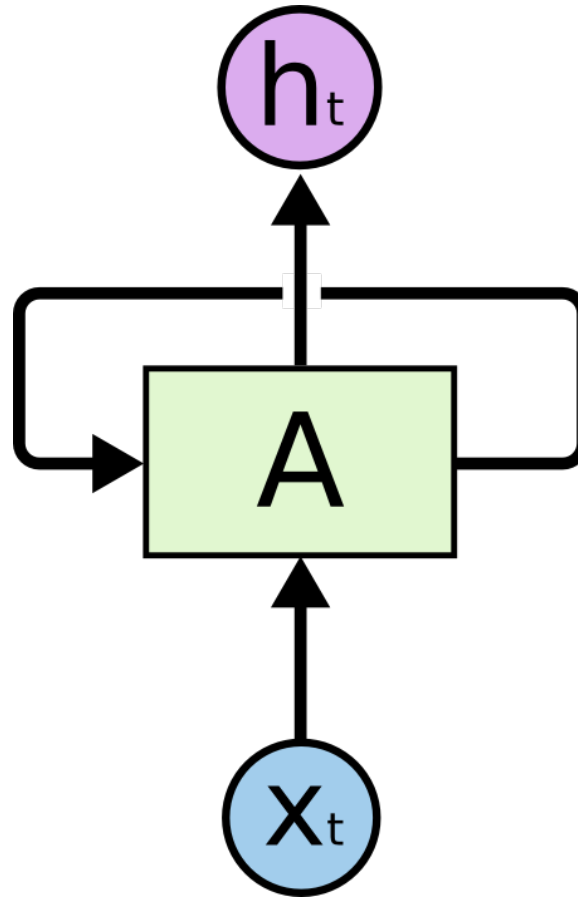
Video captioning,
Natural language translation

Text tagging tasks:
Part-of-speech detection,
Named entity recognition,
Semantic Role Labeling

Language Modeling

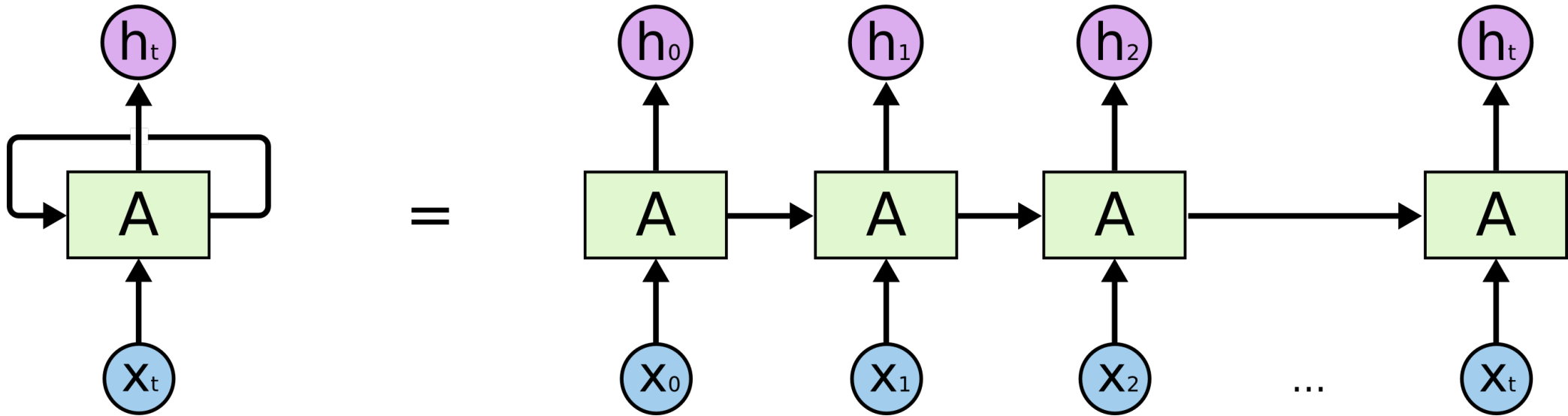
[slide credit: Andrej Karpathy](#)

Architecture Strategy: Recurrent Layers



[Img credit: Chris Olah blog post \(Understanding LSTMs\)](#)

Two Views (Self-connection vs Unrolled)



[Img credit: Chris Olah blog post \(Understanding LSTMs\)](#)

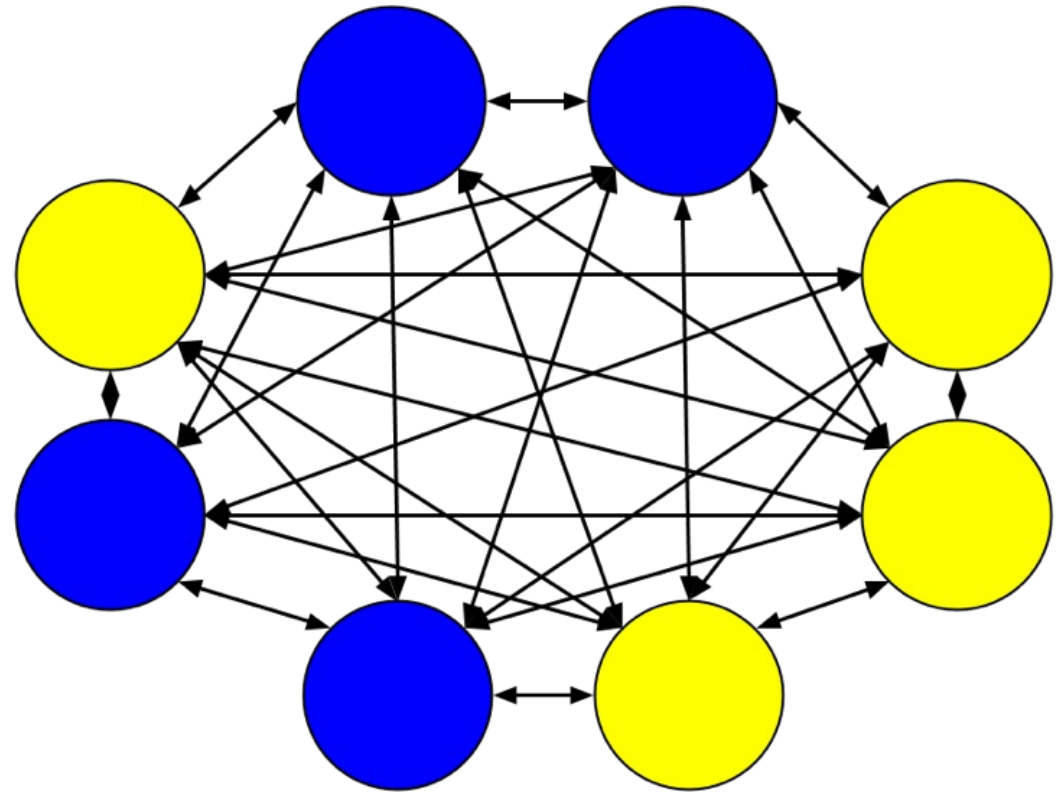
A grand, classical museum gallery with high ceilings and large columns. In the center, a statue of a man in a draped garment holds a glowing, complex neural network structure composed of many interconnected nodes and lines. The floor is polished and reflects the light. The text "A Brief Tour of Classical RNNs" is overlaid in white.

A Brief Tour of Classical RNNs

Room: Noxial Nervos

Hopfield Nets (Hopfield, 1982)

- Parametrization
 - Each node pair i, j connected by edge with weight w_{ij}
 - No self connections ($w_{ii} = 0$)
- Updates via Hebbian rule
 - Fire together \rightarrow wire together



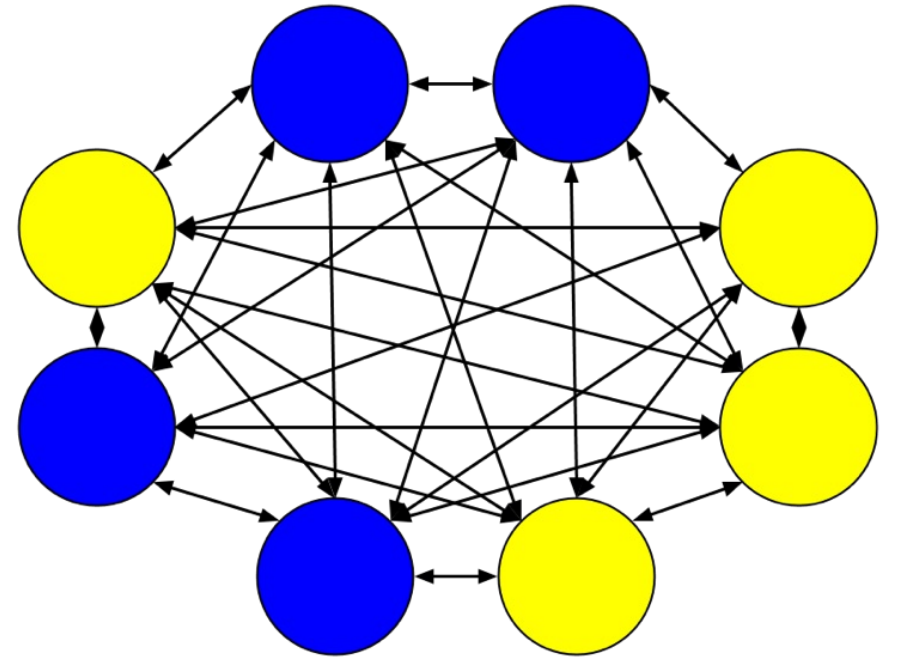
Hopfield nets (cont'd)

- Updates:
$$s_i \leftarrow \begin{cases} +1, & \text{for } \sum_j w_{ij} s_j > \theta_i \\ -1, & \text{otherwise} \end{cases}$$
- Two ways to “run” the net:
 - Asynchronously — update one at a time
 - Synchronously — simultaneously update all

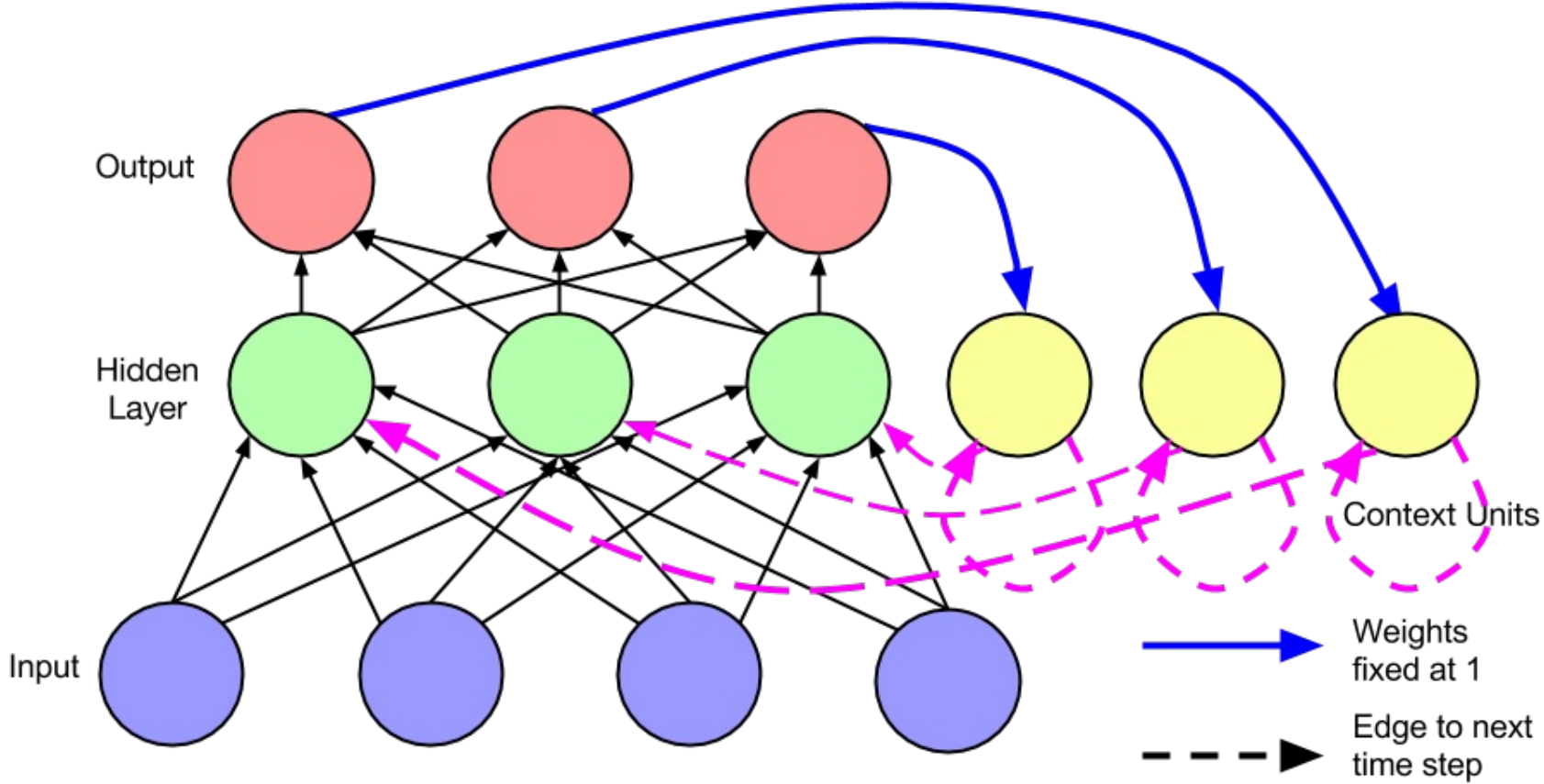
- Learning:

$$w_{ij} = 1/n \sum_{u=1}^n \epsilon_i^u \epsilon_j^u$$

- Idea: Minimize the “energy”



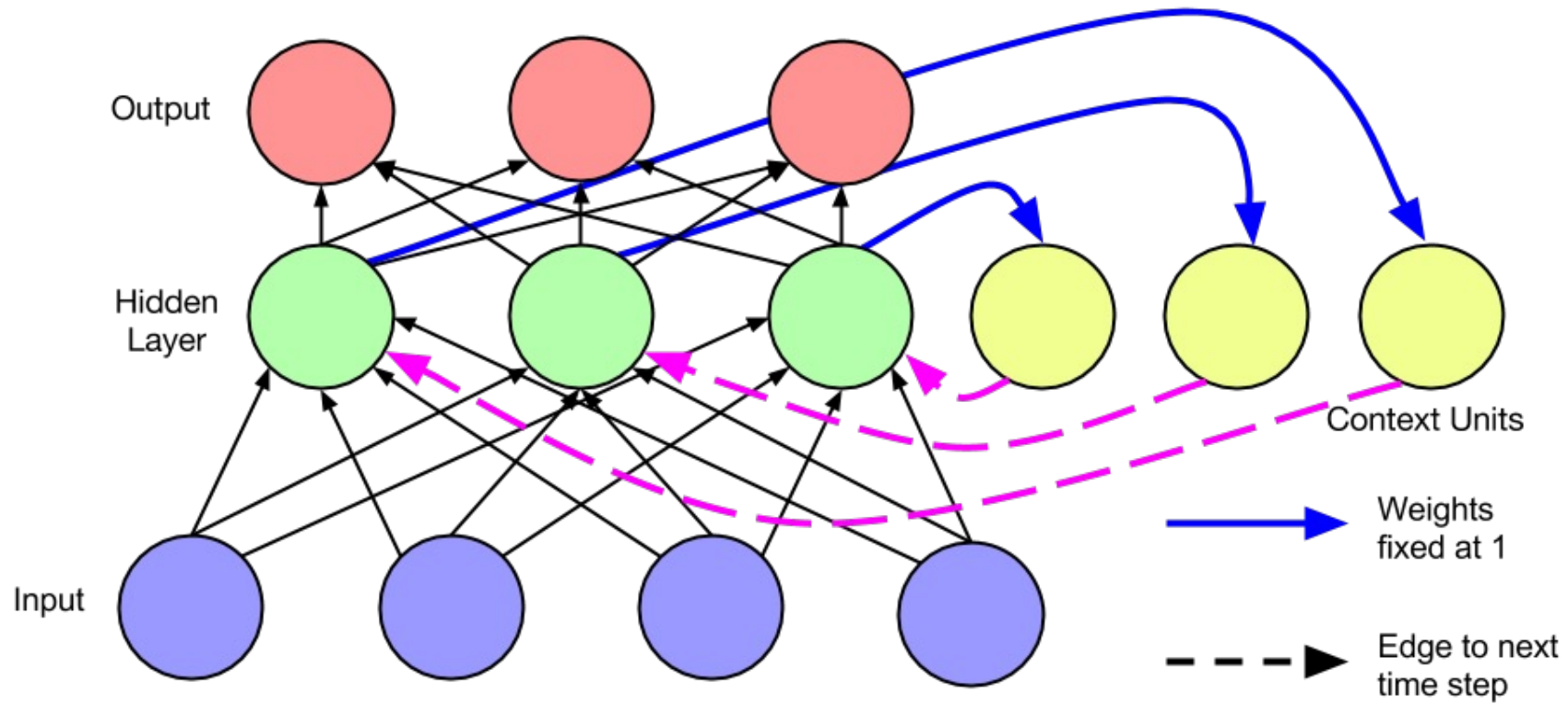
Jordan Nets (1986)



Jordan Nets (Michael I. Jordan, 1986)

- Pass output from previous step into hidden nodes at current step
- Used for planning
- Didn't use backpropagation owing to biological implausibility

Elman Nets (1990)



Elman's Experiments

(Finding Structure in Time, 1990)

- Trained net to perform temporal version of XOR
 - Input: 0, 0, 0, 1, 0, 1, 0, 1, 1
 - Every third character is the XOR of the previous two
 - Goal: predict next bit (first two are random)
 - System correctly guessed each third character with 85% accuracy
- Trained net to predict next character in text
 - Accurately predicted next char
 - 21 years before recent resurgence in auto-regressive text modeling

Sequence Model

- Dependent random variables

$$(x_1, \dots, x_T) \sim p(x)$$

- Conditional probability expansion

$$p(x) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots \cdot p(x_T|x_1, \dots, x_{T-1})$$

- **So why bother?**

$$p(x) = p(x_T) \cdot p(x_{T-1}|x_T) \cdot p(x_{T-2}|x_{T-1}, x_T) \cdot \dots \cdot p(x_1|x_2, \dots, x_T)$$

Sequence Model

$$p(x) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots p(x_T|x_1, \dots x_{T-1})$$



$$p(x) = p(x_T) \cdot p(x_{T-1}|x_T) \cdot p(x_{T-2}|x_{T-1}, x_T) \cdot \dots p(x_1|x_2, \dots x_T)$$



- Causality (physics) prevents the reverse direction
- 'wrong' direction often much more complex to model
- Typically more useful to predict text in left-to-right fashion

Sequence Model

$$p(x) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots p(x_T|x_1, \dots x_{T-1})$$

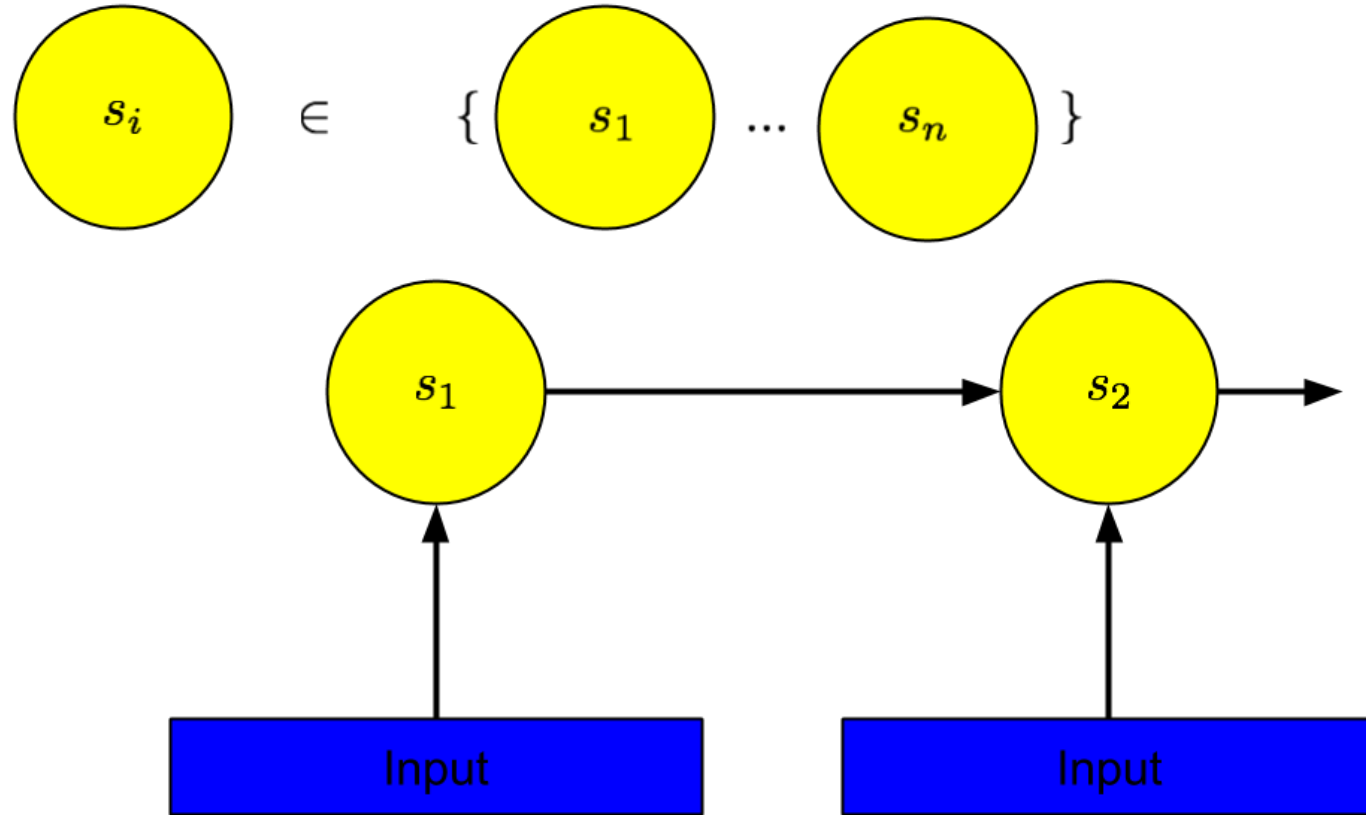


- Autoregressive model

$$p(x_t|x_1, \dots x_{t-1}) = p(x_t|f(x_1, \dots x_{t-1}))$$

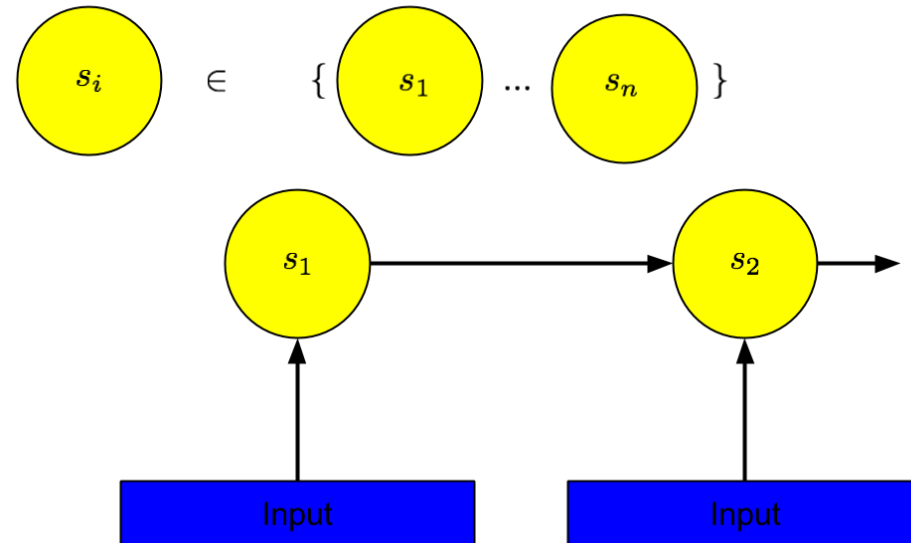
Some function of previously seen data

Why Not (Hidden) Markov Models?



Why Not (Hidden) Markov Models?

- For Markov models, inference is quadratic in number of states
- RNNs representational expressivity grows exponentially in latent dimension, computation scales linearly



Classic n-gram Language Models

- Make (incorrect) Markov assumption

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1}) \approx \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

- Estimate probabilities by counting exact occurrences of n-grams

$$p(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \quad p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

- Smoothing: give some mass to unseen n-grams so loss doesn't $\rightarrow \infty$

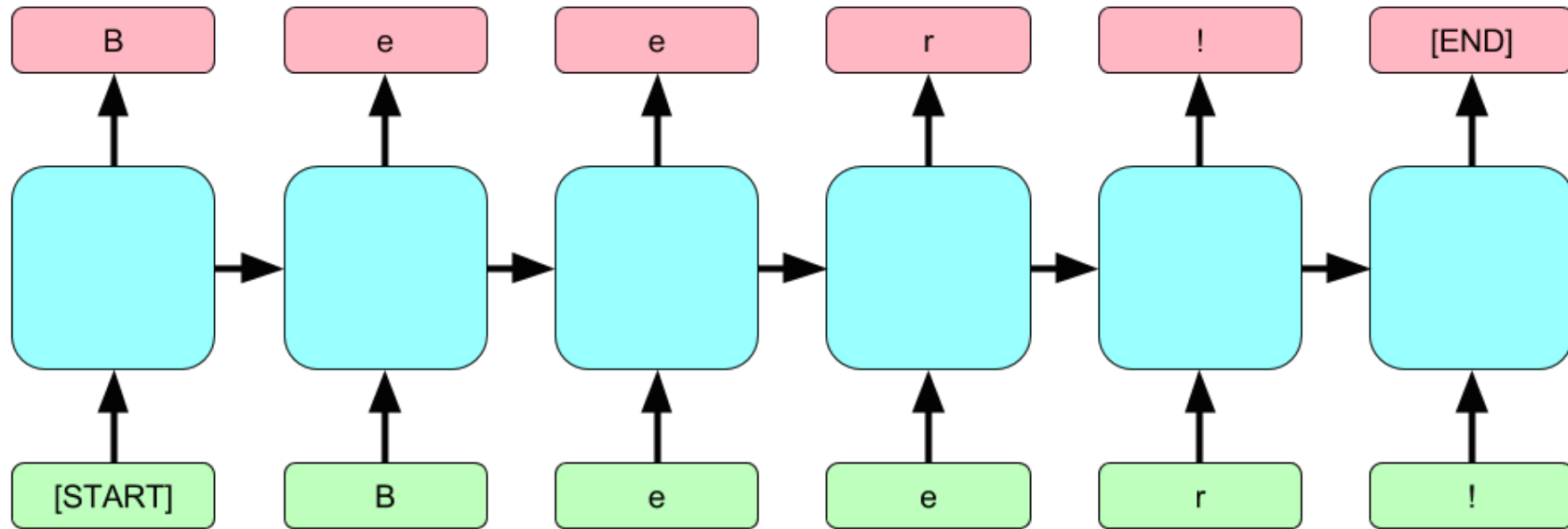
The equation for bigram probabilities is as follows:

$$p_{KN}(w_i|w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - \delta, 0)}{\sum_{w'} c(w_{i-1}, w')} + \lambda_{w_{i-1}} p_{KN}(w_i) \quad [4]$$

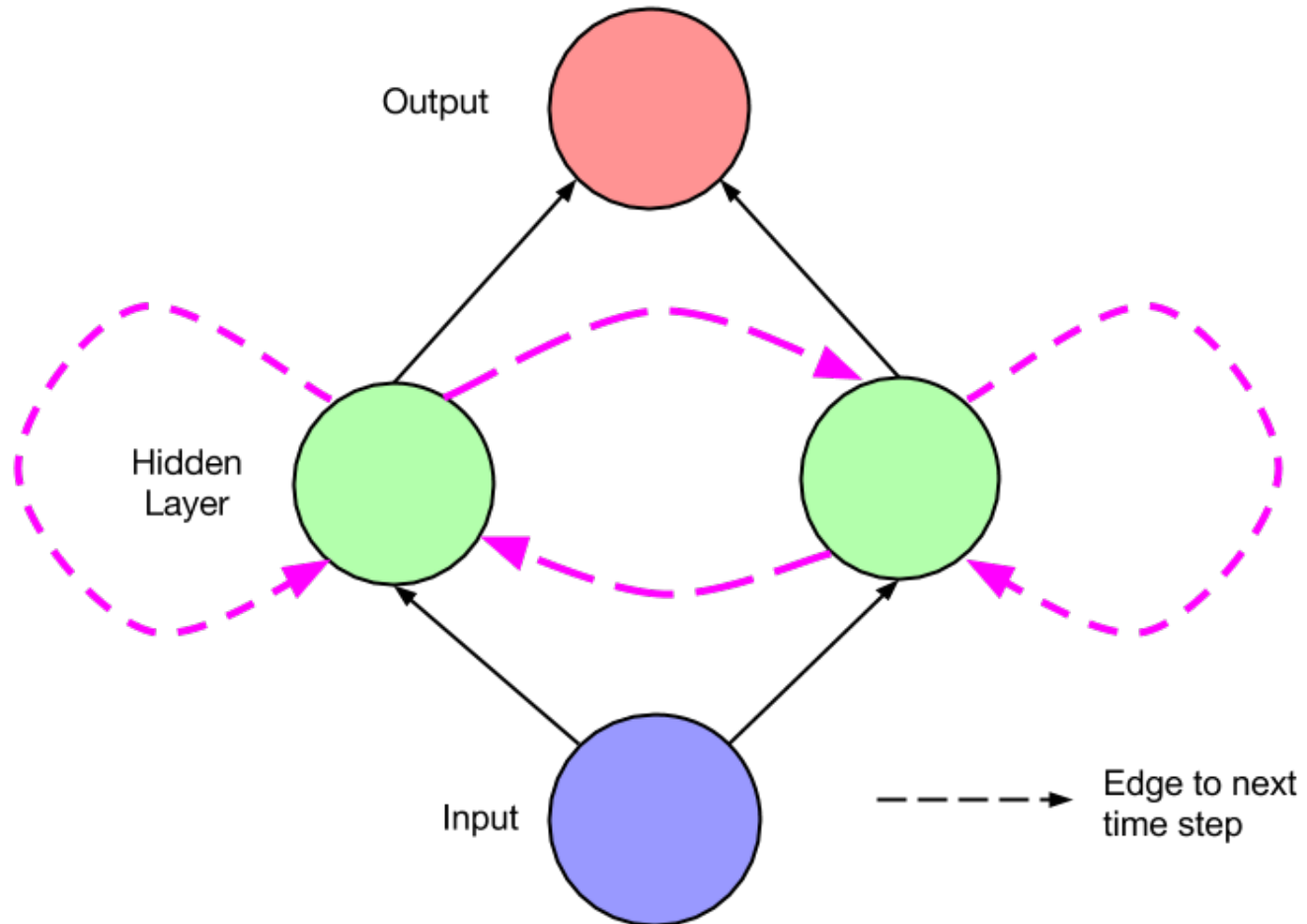
Where the unigram probability $p_{KN}(w_i)$ depends on how likely it is to see the word w_i in an unfamiliar context, which is estimated as the number of times it appears after any other word divided by the number of distinct pairs of consecutive words in the corpus:

[Kneser-Ney wiki page](#)

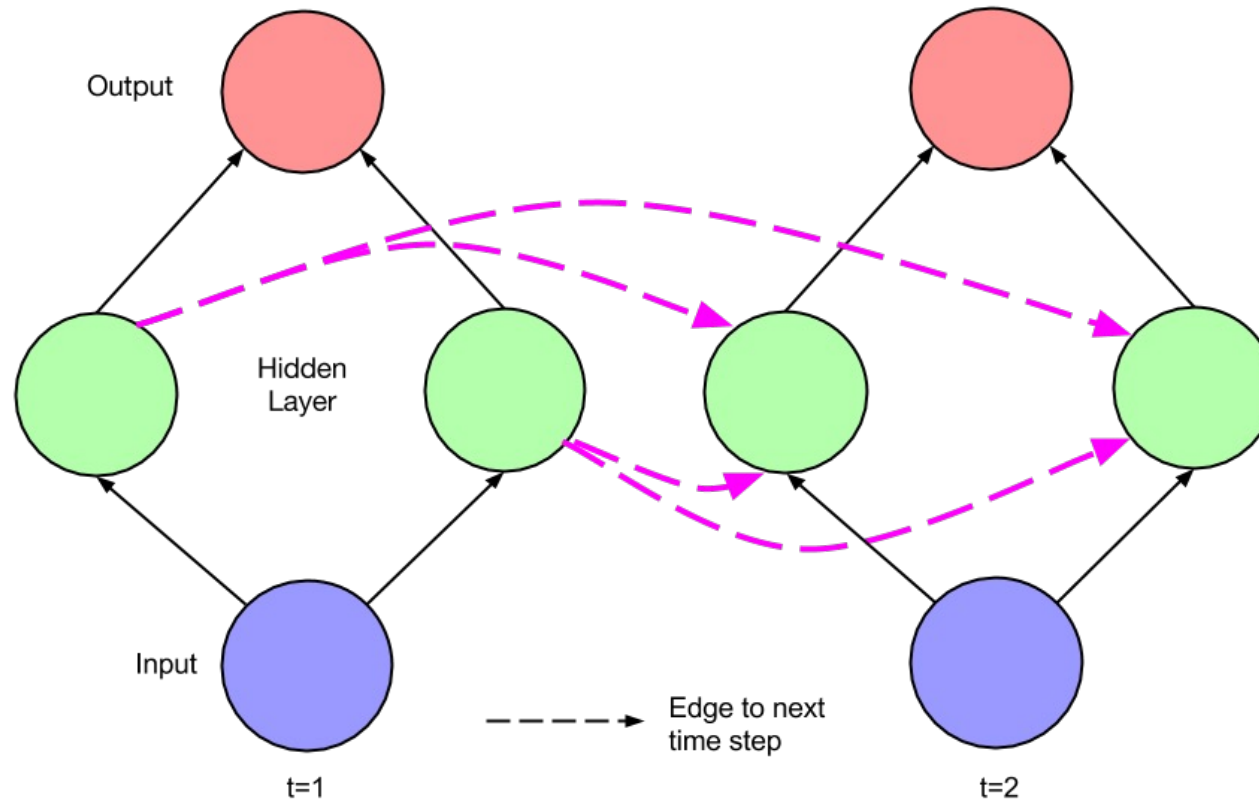
RNN-based Language Modeling



Recurrent Neural Networks



Recurrent Neural Network (Unfolded)



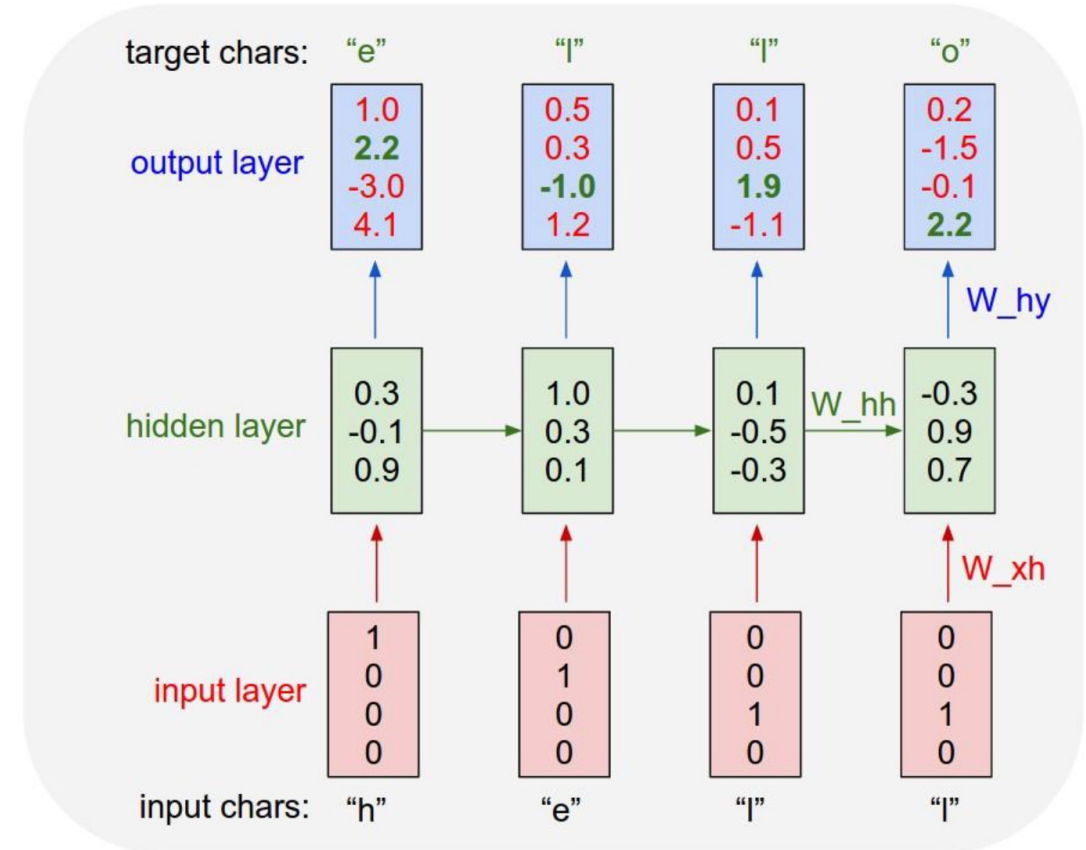
$$\mathbf{h}^{(t)} = \sigma(W_{hx}\mathbf{x}^{(t)} + W_{hh}\mathbf{h}^{(t-1)} + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W_{yh}\mathbf{h}^{(t)} + \mathbf{b}_y)$$

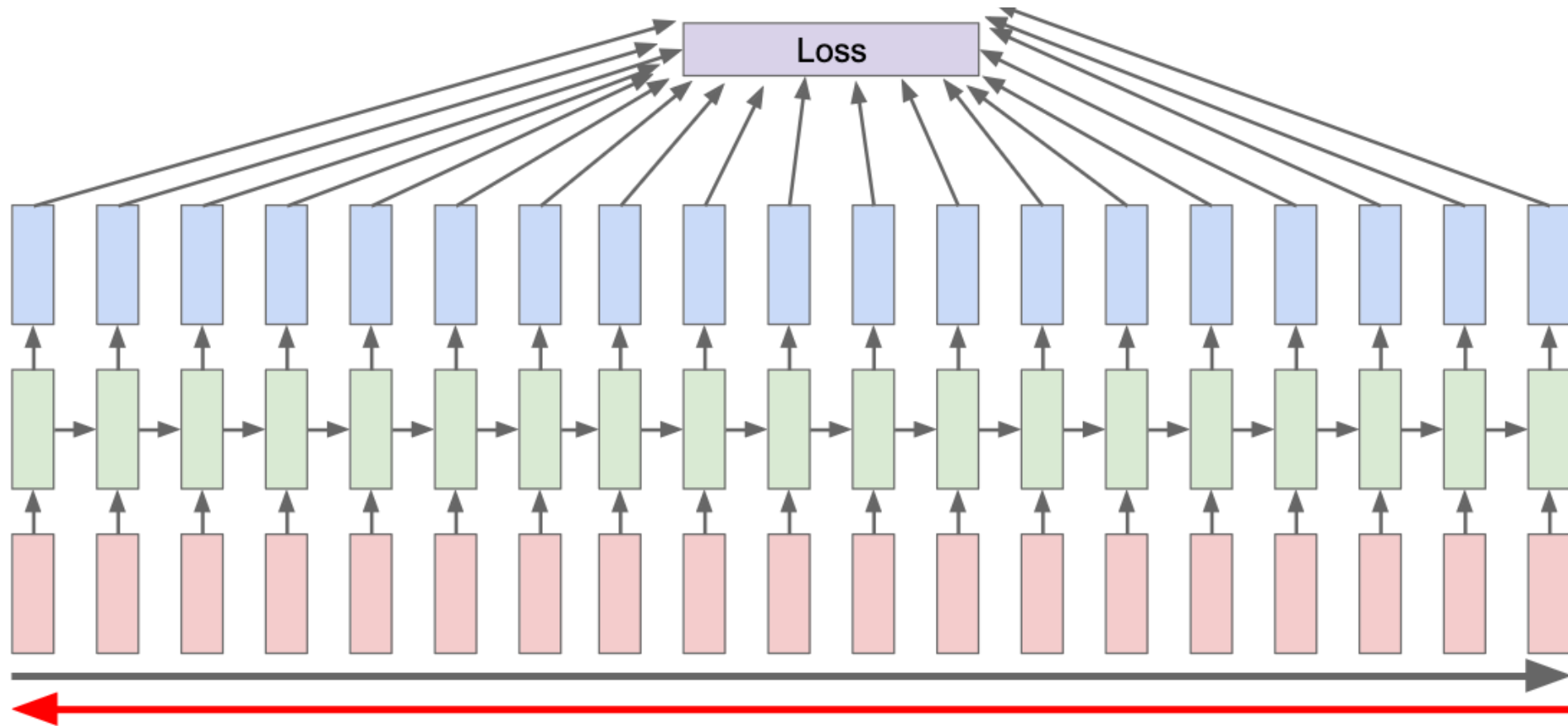
RNN-based LMs at Train Time

At Each sequence step t :

- Input is current token x_t
- Target is subsequent token x_{t+1}
- Sample sequences up to some maximum length
- Backprop across compute graph

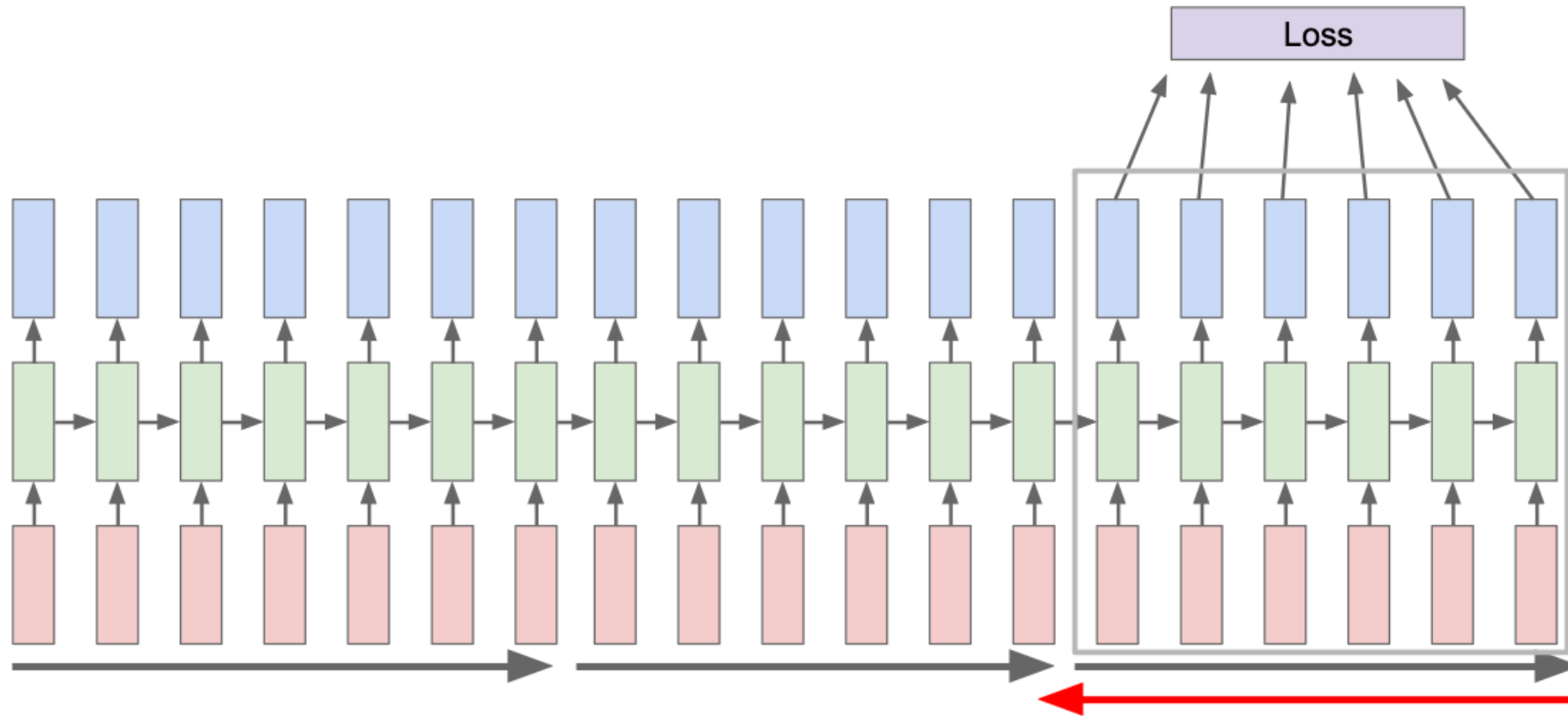


Backpropagation Through Time



[slide credit: Andrej Karpathy](#)

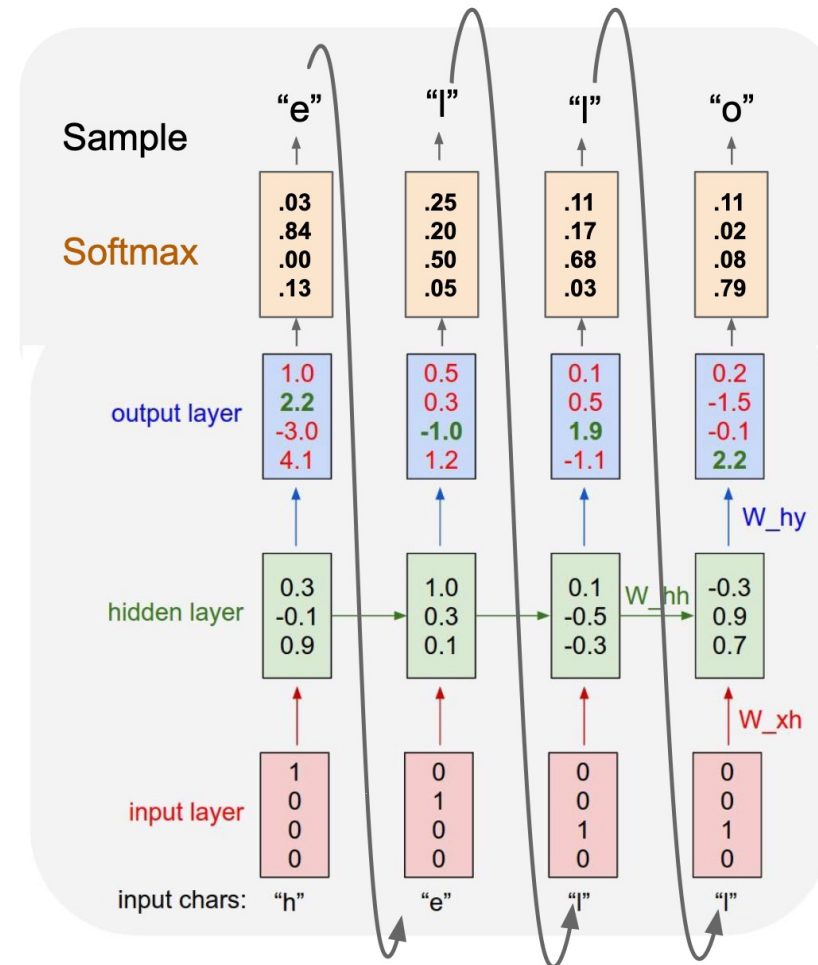
Truncated Backpropagation Through Time



[slide credit: Andrej Karpathy](#)

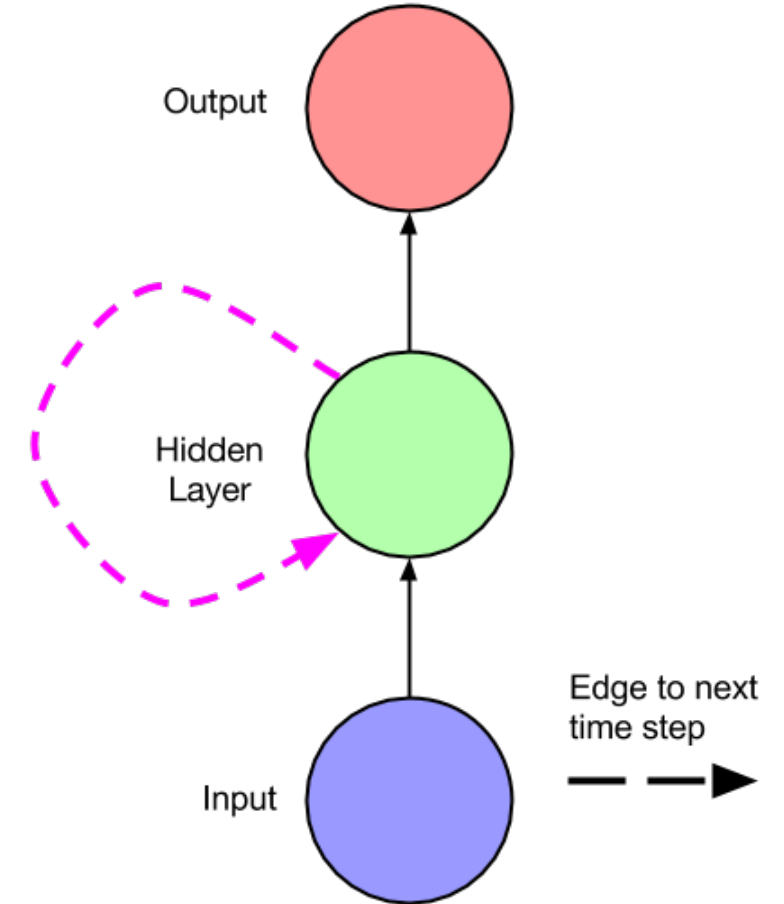
RNN-based LMs at Prediction Time

- Feed in some input prefix $x_{1:t}$
- At each step t , predict next token
- Sample from soft-max distribution
- Feed that token back to model as the subsequent input



Vanishing / Exploding Gradients

- Recurrent neural networks are very deep along the sequence axis
- Depending on properties of recurrent weights, can lead to gradients exploding or vanishing
- Many key innovations and training tricks in RNNs motivated by this problem

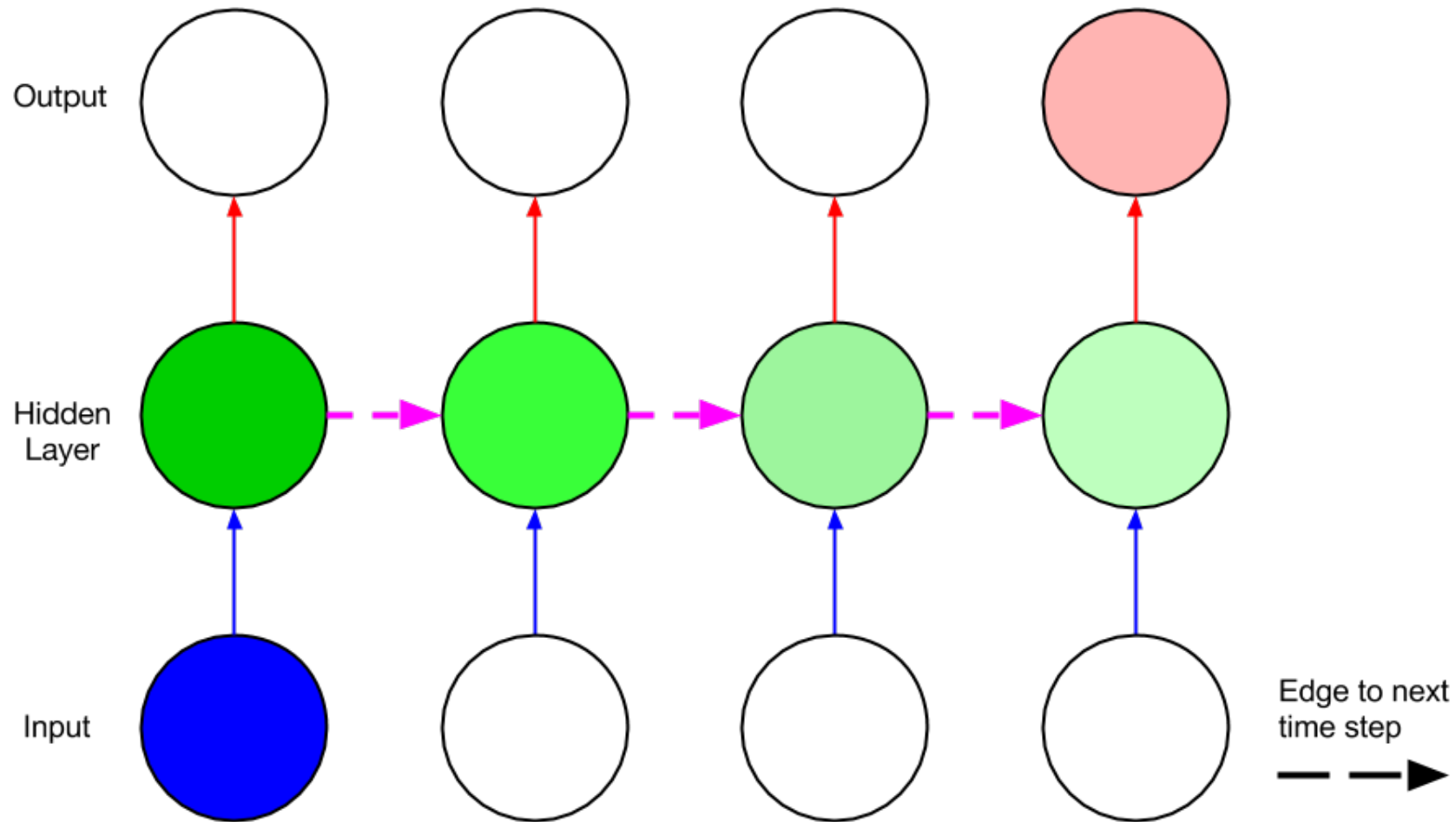


Gradient Clipping

- One strategy to mitigate the exploding gradient problem:
- Set an upper limit on the norm of the gradient
- Whenever gradient exceeds this norm, scale it down
- Apply SGD update on the rescaled gradient

$$\nabla_w^L \leftarrow \begin{cases} \nabla_w^L & \text{if } \|\nabla_w^L\|_1 \leq \max \\ \max \cdot \nabla_w^L / \|\nabla_w^L\|_1 & \text{if } \|\nabla_w^L\|_1 > \max \end{cases}$$

Vanishing Gradients



RNN Training Details

- Initial value for hidden state h_0 typically set to 0
- Traditional RNNs used tanh activation, random init
- Some modern RNNs initialize weights as Identity matrix & use ReLU

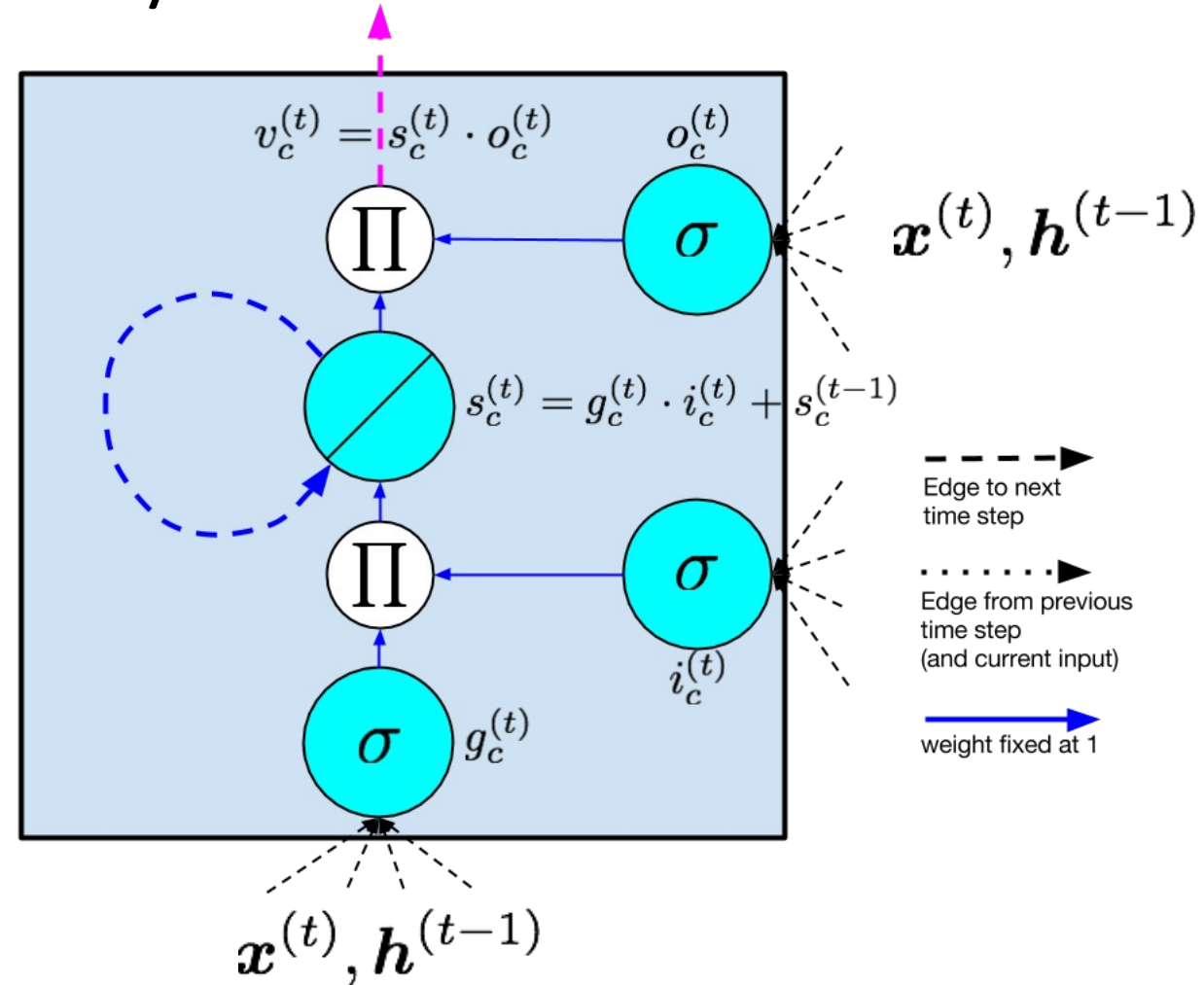
[Le, Jaitly, Hinton \(2015\)](#)

Long Short-Term Memory Cell

- Each **hidden node** replaced by a **memory cell**
- Designed to deal with learn long range time dependencies
- Internal state s gets residual self-connection (fixed weight of 1)
(similar idea to residual connections in ResNets)
- Input gate i determines when to let activation into the cell
- Output Gate o determines when to let activation out of the cell

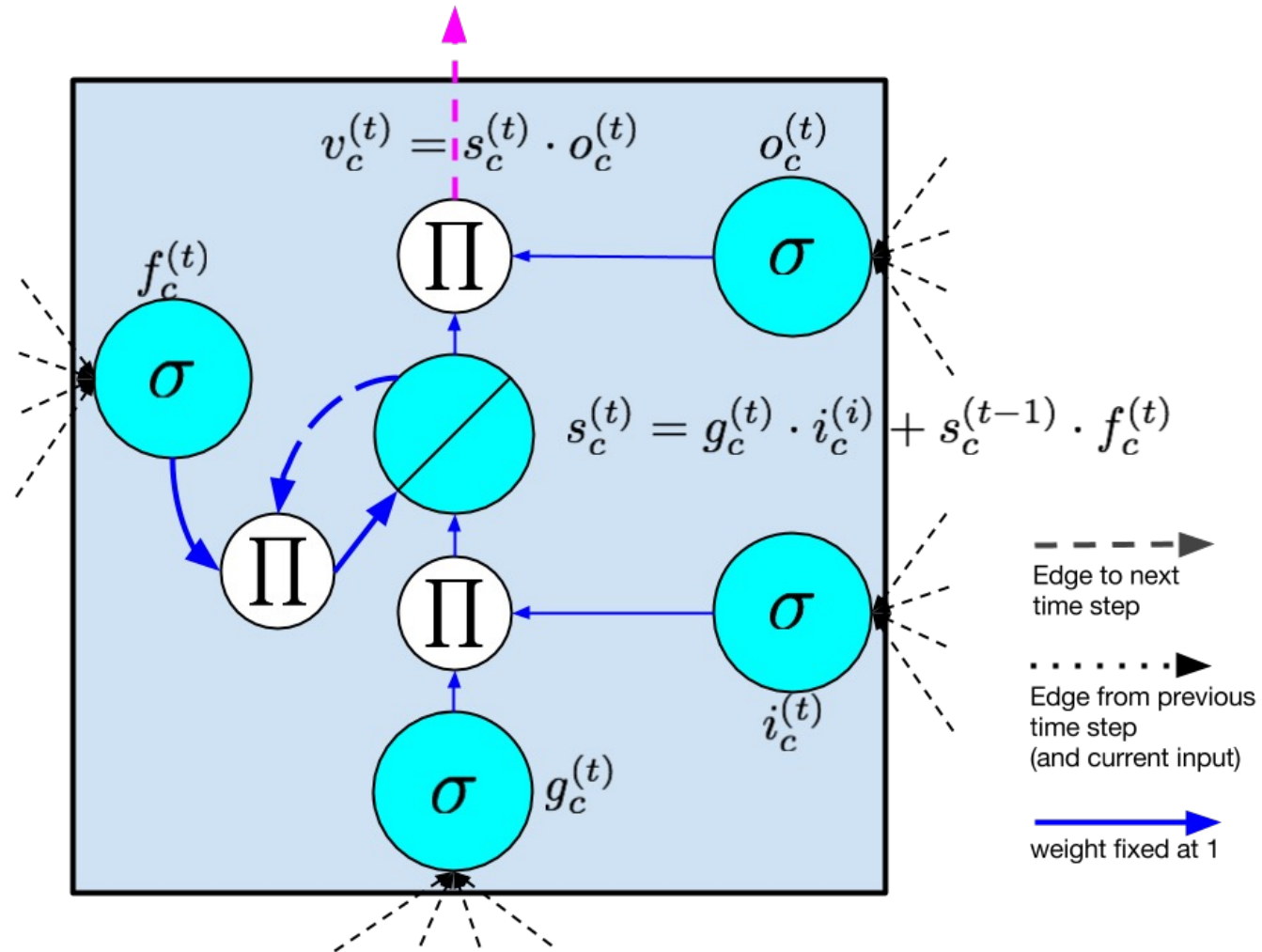
[Hochreiter & Schmidhuber 1997](#)

LSTM Memory Cell



[Hochreiter & Schmidhuber 1997](#)

Memory Cell with a Forget Gate



Gers et al., Neural Computation 2000

LSTM Forward Pass

$$\mathbf{g}^{(t)} = \phi(W_{gx}\mathbf{x}^{(t)} + W_{gh}\mathbf{h}^{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{i}^{(t)} = \sigma(W_{ix}\mathbf{x}^{(t)} + W_{ih}\mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}^{(t)} = \sigma(W_{fx}\mathbf{x}^{(t)} + W_{fh}\mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

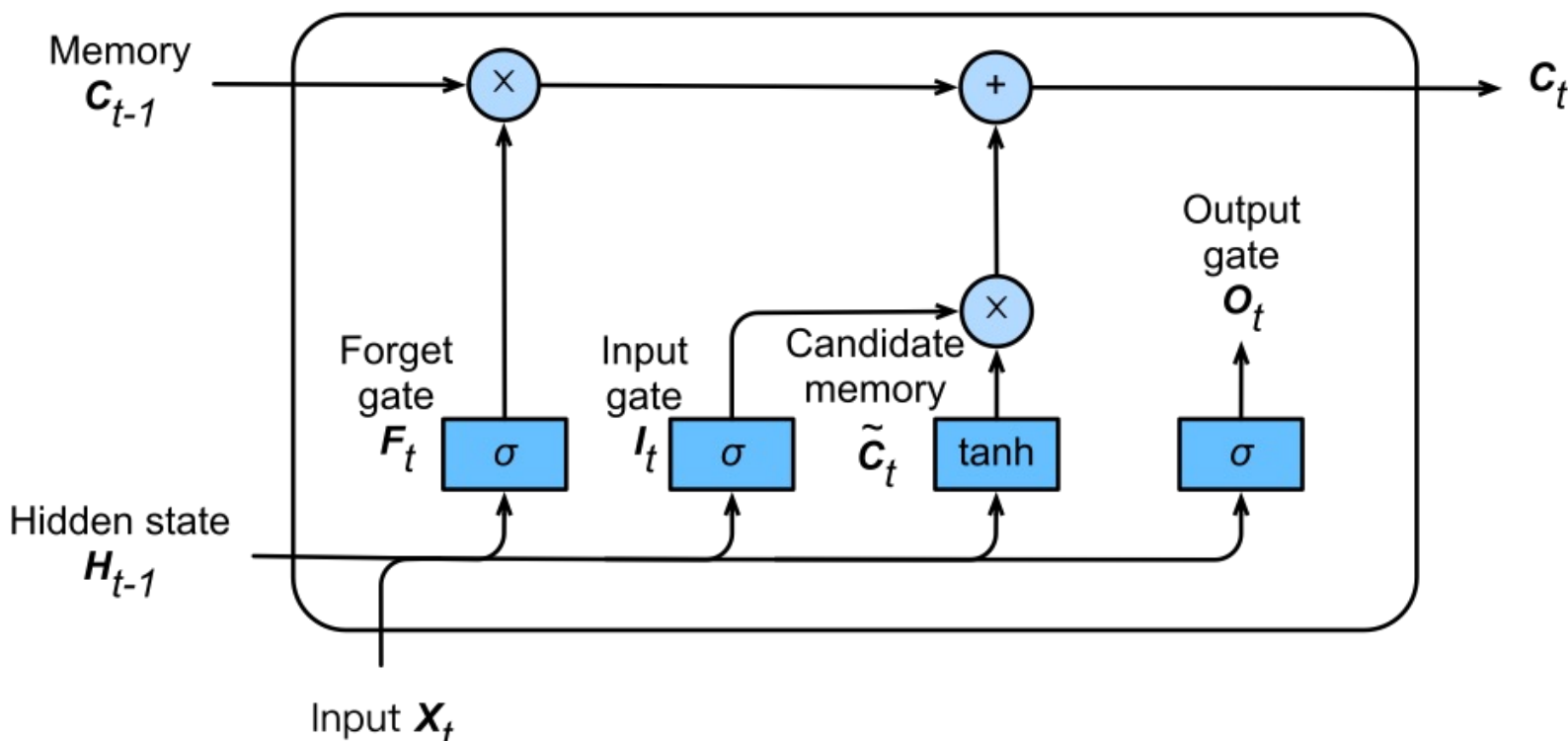
$$\mathbf{o}^{(t)} = \sigma(W_{ox}\mathbf{x}^{(t)} + W_{oh}\mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{s}^{(t)} = \mathbf{g}^{(t)} \odot \mathbf{i}^{(t)} + \mathbf{s}^{(t-1)} \odot \mathbf{f}^{(t)}$$

$$\mathbf{h}^{(t)} = \mathbf{s}^{(t)} \odot \mathbf{o}^{(t)}$$

Memory Cell Update

$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$



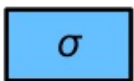
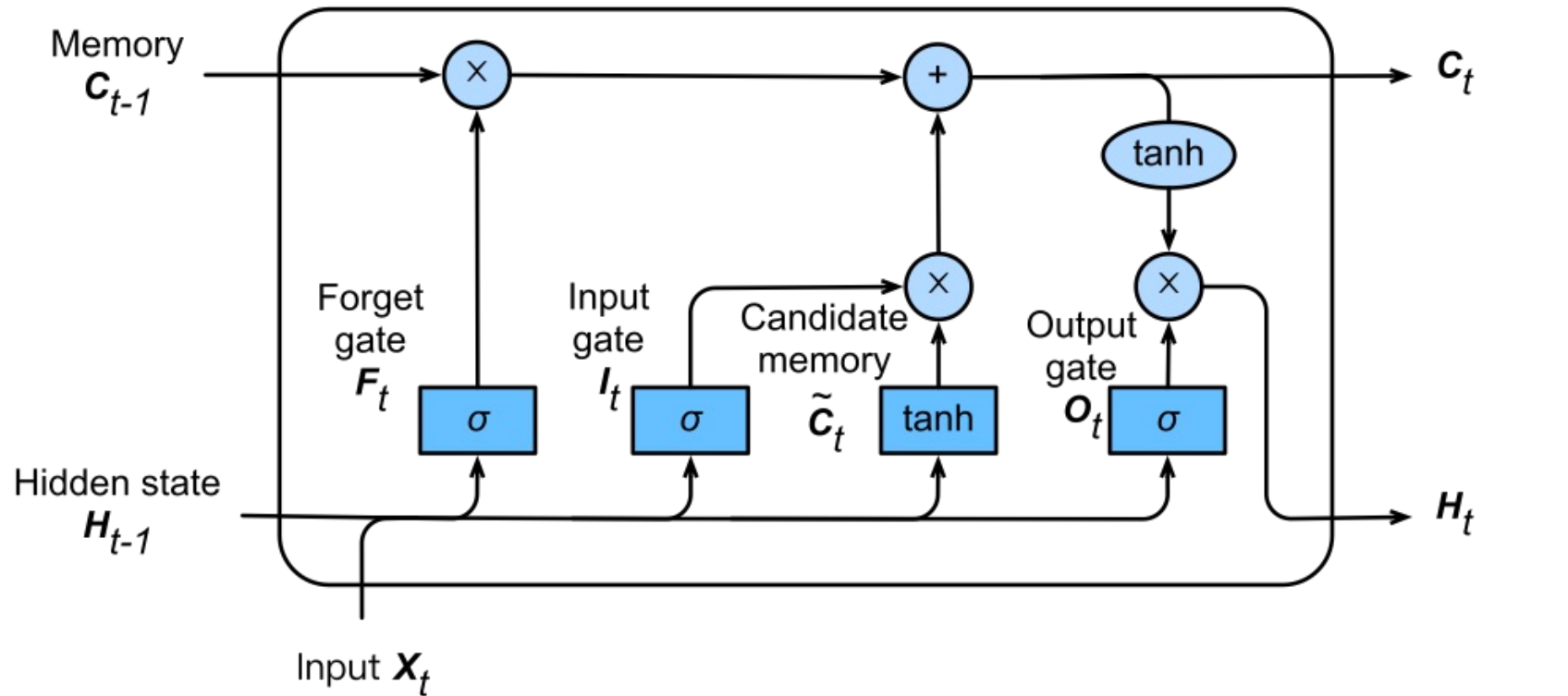
- σ

FC layer with activation function
- \times

Element-wise Operator
- Copy
- Concatenate

Generating Memory Cell Output

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t)$$



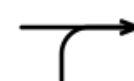
FC layer with activation function



Element-wise Operator

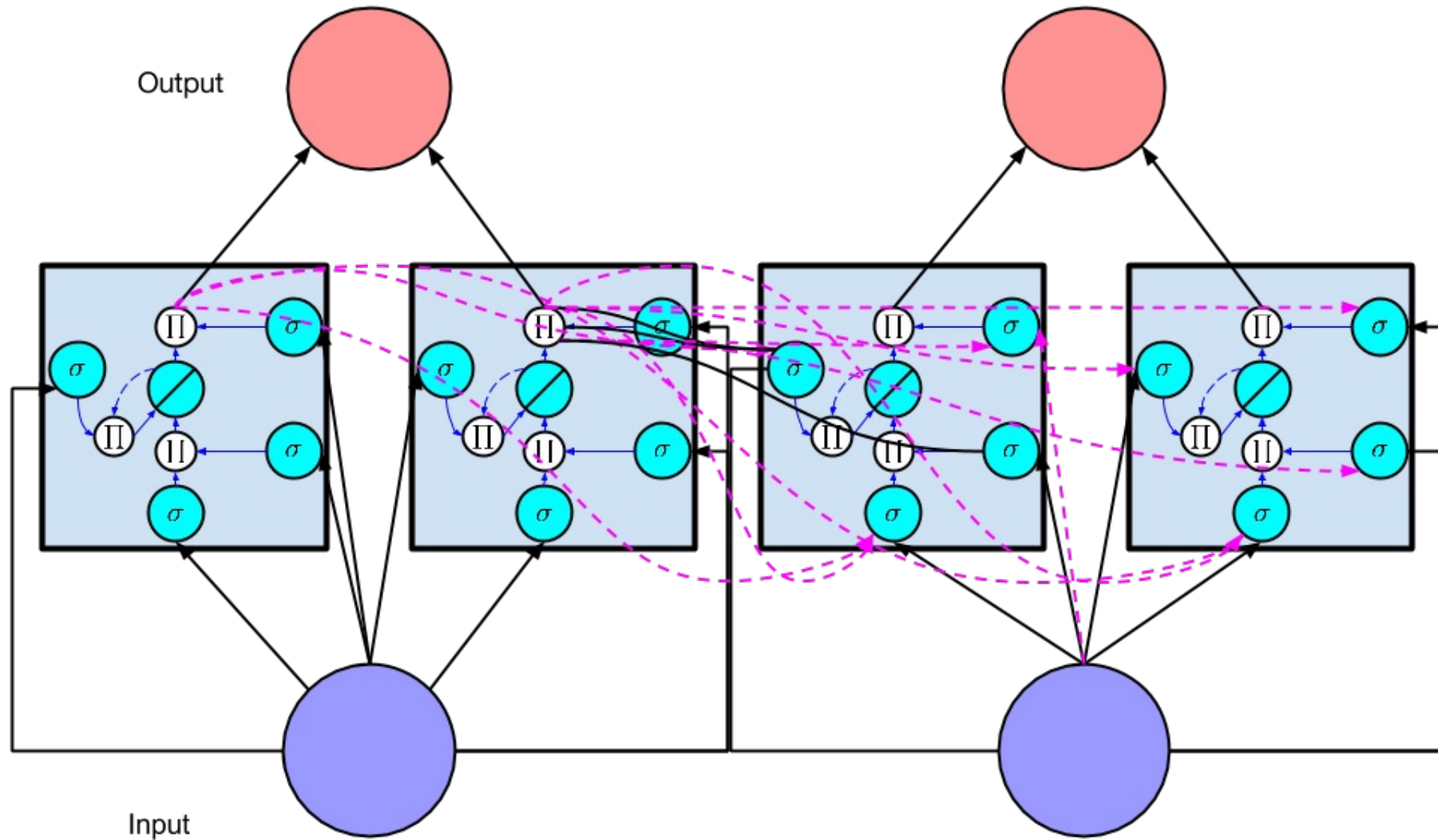


Copy

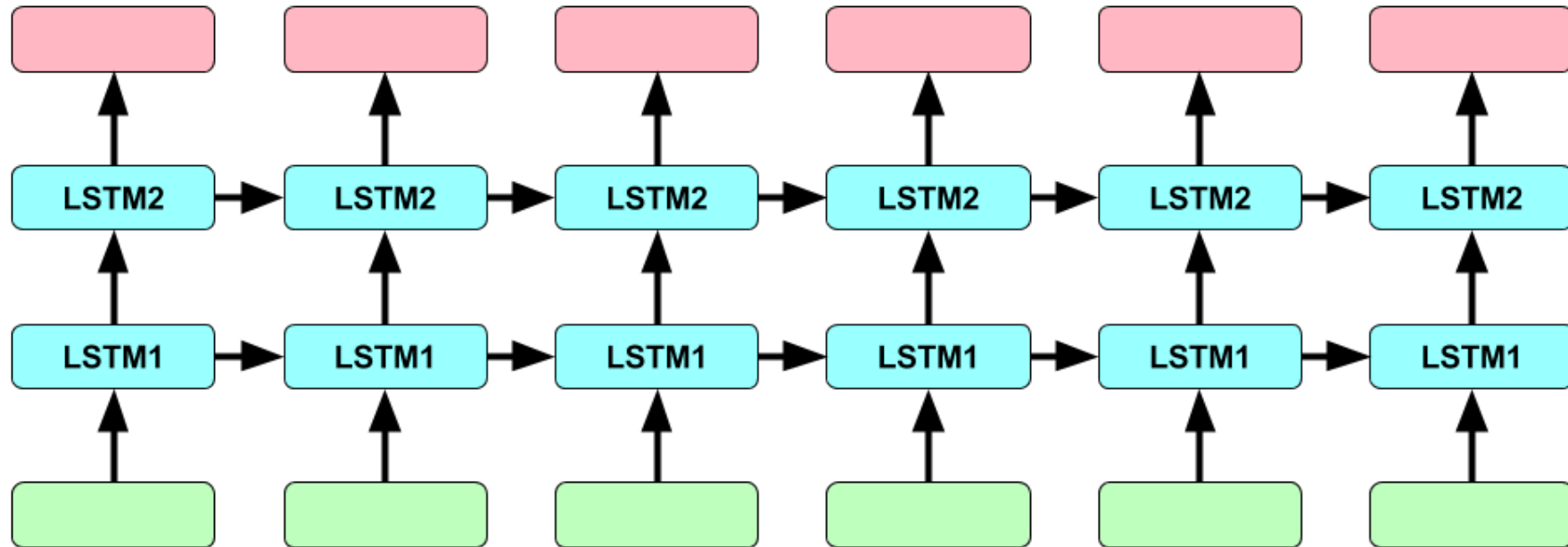


Concatenate

Wiring together the LSTM cells



Deep LSTMs (stack layers vertically)



Gated Recurrent Unit (GRUs)

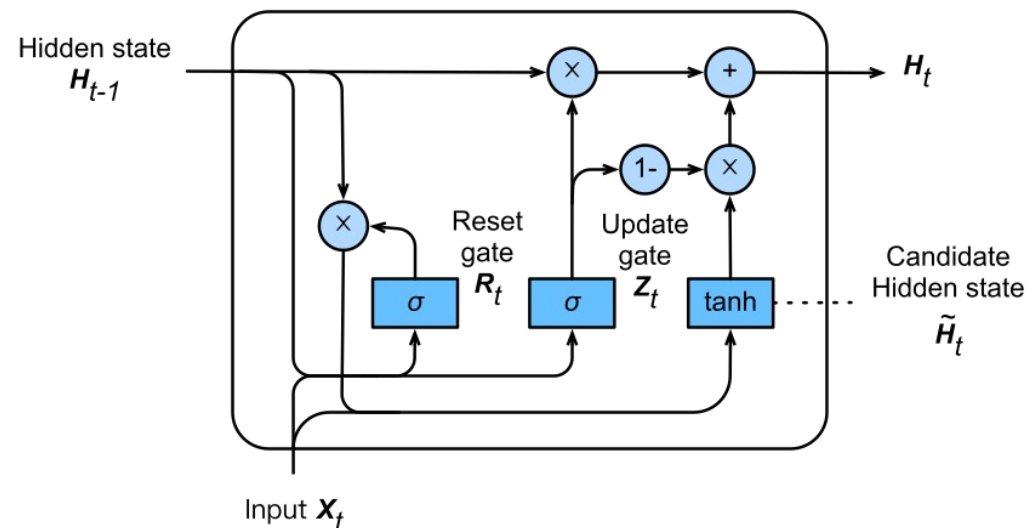
- Simplified version of LSTM cell
- Only two gates (not 3)
- Less parameters than LSTM
- Similar performance
- Keeps core ideas

$$\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r),$$

$$\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z)$$

$$\tilde{\mathbf{H}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xh} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{hh} + \mathbf{b}_h)$$

$$\mathbf{H}_t = \mathbf{Z}_t \odot \mathbf{H}_{t-1} + (1 - \mathbf{Z}_t) \odot \tilde{\mathbf{H}}_t$$



Training & Generating Shakespeare

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

Training Data:
All of Shakespeare

Size: 4.4MB

[Unreasonable Effectiveness of RNNs \(Karpathy 2015\)](#)

Further RNN Architecture Exploration

- [Stacked LSTM \(2013\)](#)
- [Grid LSTM \(2015\)](#)
- [Recurrent Highway Nets \(Zilly et al 2016\)](#)
- [NAS-based memory cell \(Zoph et al. 2016\)](#)
- [Pointer Sentinel LSTMs \(Merity et al. 2016\)](#)

ON THE STATE OF THE ART OF EVALUATION IN NEURAL LANGUAGE MODELS

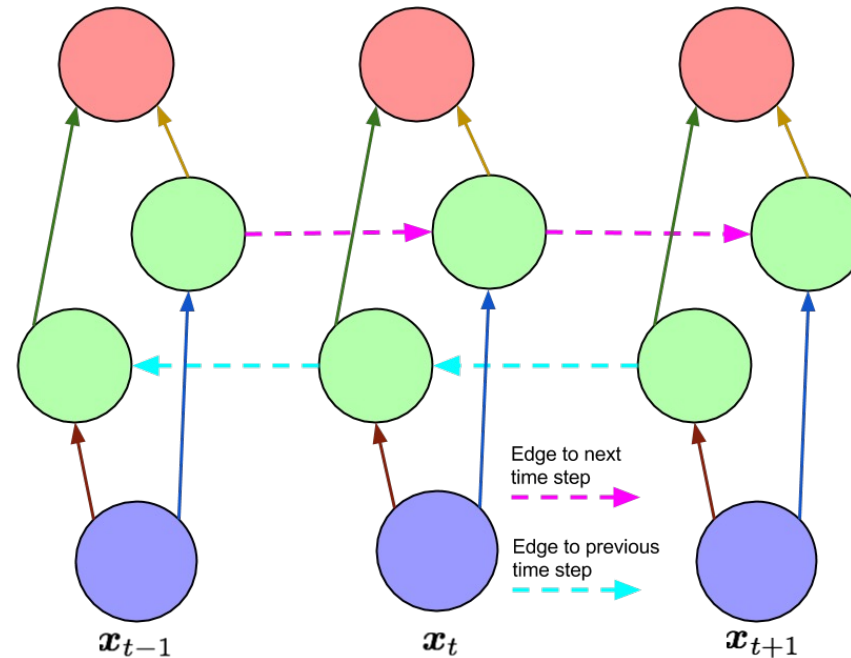
Gábor Melis[†], Chris Dyer[†], Phil Blunsom^{†‡}
{melisg1, cdyer, pblunsom}@google.com
[†]DeepMind
[‡]University of Oxford

ABSTRACT

Ongoing innovations in recurrent neural network architectures have provided a steady influx of apparently state-of-the-art results on language modelling benchmarks. However, these have been evaluated using differing codebases and limited computational resources, which represent uncontrolled sources of experimental variation. We reevaluate several popular architectures and regularisation methods with large-scale automatic black-box hyperparameter tuning and arrive at the somewhat surprising conclusion that standard LSTM architectures, when properly regularised, outperform more recent models. We establish a new state of the art on the Penn Treebank and Wikitext-2 corpora, as well as strong baselines on the Hutter Prize dataset.

BRNN

(Schuster & Paliwal, 1997)



$$\mathbf{h}_f^{(t)} = \sigma(W_{h_f x} \mathbf{x}^{(t)} + W_{h_f h_f} \mathbf{h}_f^{(t-1)} + \mathbf{b}_{h_f})$$

$$\mathbf{h}_b^{(t)} = \sigma(W_{h_b x} \mathbf{x}^{(t)} + W_{h_b h_b} \mathbf{h}_b^{(t+1)} + \mathbf{b}_{h_b}) \quad (1)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W_{y h_f} \mathbf{h}_f^{(t)} + W_{y h_b} \mathbf{h}_b^{(t)} + \mathbf{b}_y)$$

Applying Dropout to RNNs

- [Recurrent Dropout \(Gal & Ghahramani 2015\)](#):

- Apply same dropout mask at each time step.

- [Another variant Semeniuta et al. \(2016\)](#):

- Apply dropout to LSTM “update vector” g , not to hidden state.

$$\mathbf{c}_t = \mathbf{f}_t * \mathbf{c}_{t-1} + \mathbf{i}_t * d(\mathbf{g}_t)$$

- [Zoneout \(Krueger et al. 2017\)](#):

- Rather than zero-out units, replace with previous step’s activations.
- Gradient still flows through zone-out units.