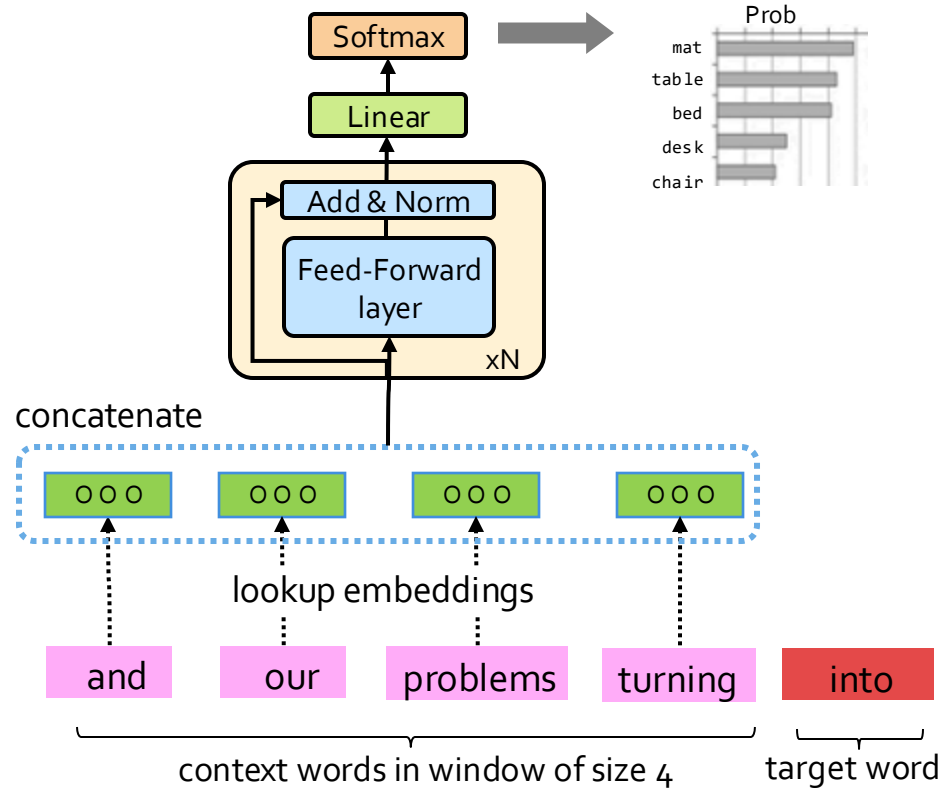


Recurrent Neural Language Models

CSCI 601-471/671 (NLP: Self-Supervised Models)

Recap

- **Neural Language Models:** neural networks trained with LM objective.
- **Fixed-window Neural LM:** first of many neural LMs we will see in this class.



What Changed from N-Gram LMs to Neural LMs?

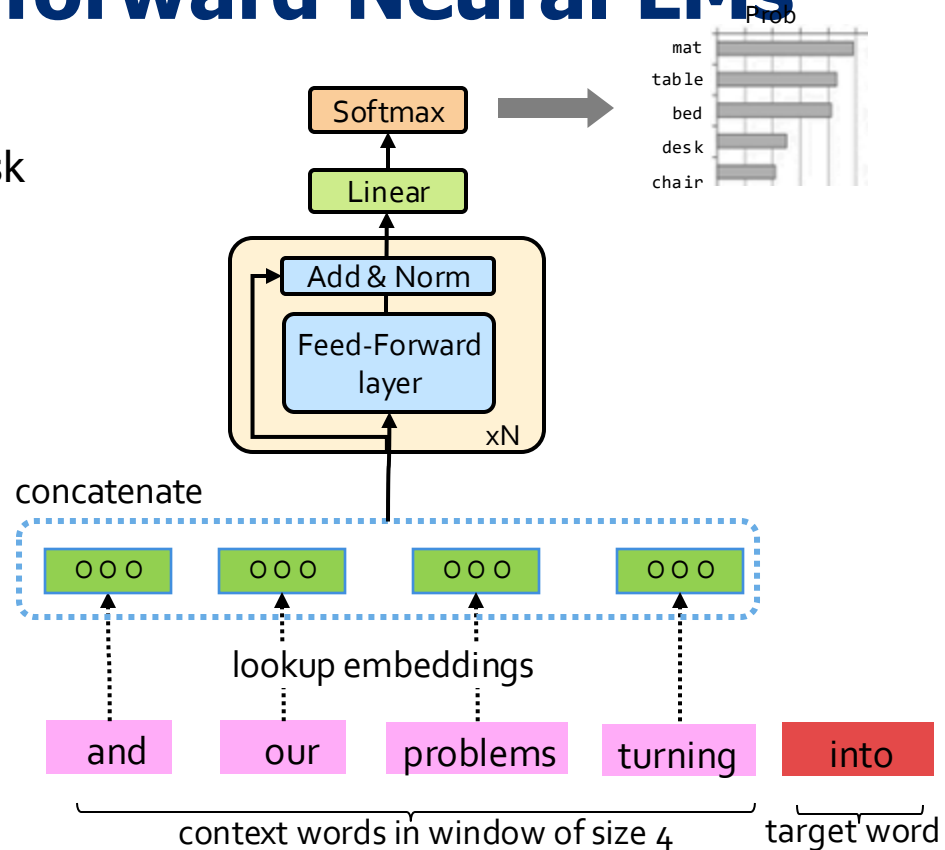
- What is the source of Neural LM's **strength**?
- Why **sparsity** is less of an issue for Neural LMs?
- **Answer:** In n-grams, we treat all prefixes independently of each other! (even those that are semantically similar)

students opened their ____
pupils opened their ____
scholars opened their ____
undergraduates opened their ____
students turned the pages of their ____
students attentively perused their ____
...

Neural LMs are able to **share information across these semantically-similar prefixes** and overcome the sparsity issue.

Moving Beyond Feedforward Neural LMs

- Are competitive at language modeling task
- However, they
 - have difficulty in remembering long range dependencies
 - have a fixed window size
- **Key question:** how to better capture long-range dependencies?
- Alternative here: a new family of neural networks: recurrent nets



Chapter Goals

1. Introducing Recurrent Neural Language Models
2. RNNs: Pros and Cons
3. Algorithms for sampling from LMs
4. Bonus: Pre-trained RNN language models

Chapter goals — Getting comfortable with RNNs for language modeling and the use of LMs for solving down-stream tasks.

Recurrent Neural Nets

Infinite Use of Finite Model

- Main question: how can a **finite** model a **long** (infinite) context?
- Solution: recursion! (recursive use of a model)
- RNNs are a family of neural networks introduced to **learn sequential data** via **recursive** dynamics.
- Inspired by the temporality of human thoughts.

Recurrent Neural Networks (RNNs)

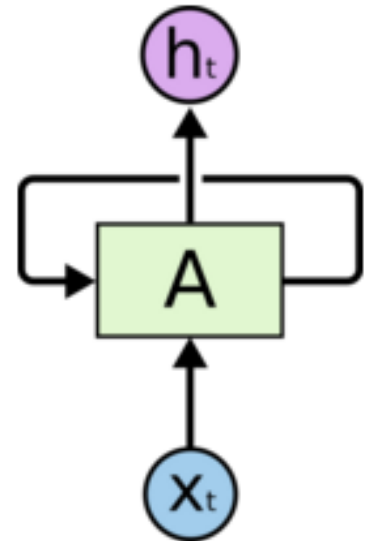
new state

old state

Input vector at t

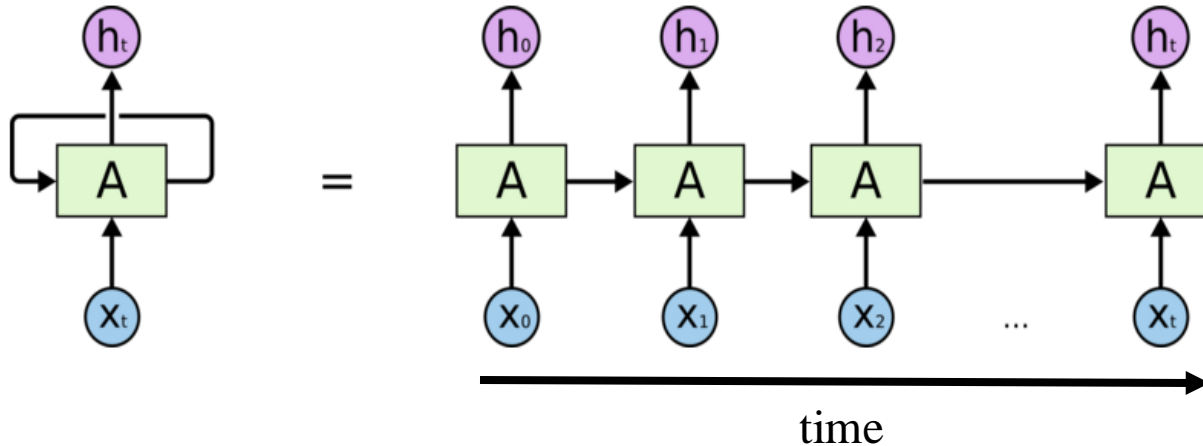
$$h_t = f(h_{t-1}, x_t)$$

- In the diagram, $f(\cdot)$ looks at some **input** x_t and its **previous hidden state** h_{t-1} and outputs a **revised state** h_t .
- A loop allows information to be passed from one step of the network to the next.



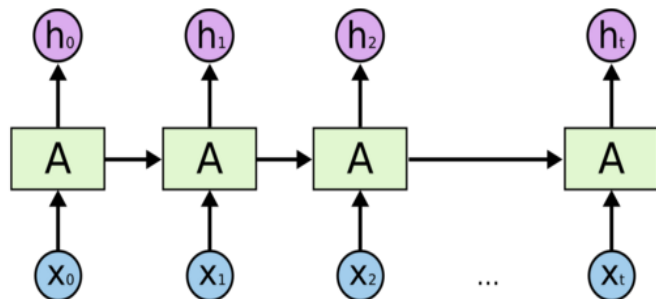
Unrolling RNN

- The diagram above shows what happens if we **unroll the loop**.



- A recurrent neural network can be thought of as **multiple copies of the same network**, each **passing a message to a successor**.

LMs w/ Recurrent Neural Nets

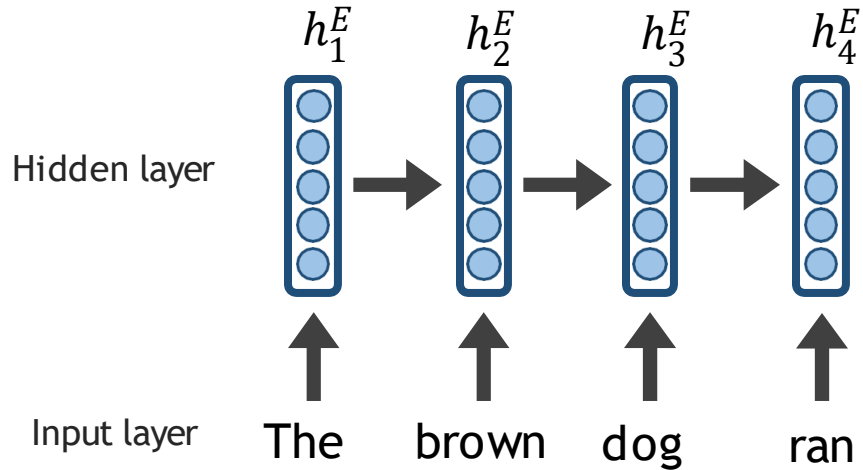


$$P(X_t | X_1, \dots, X_{t-1})$$

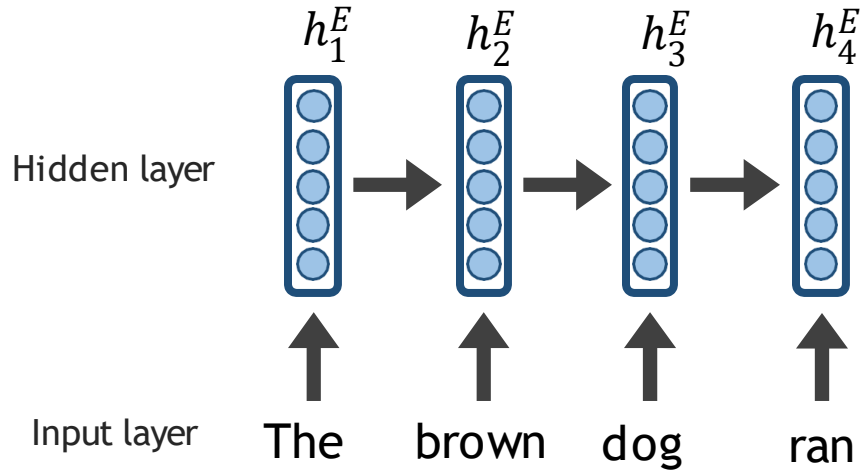
next word context

- We feed the **words one at a time** to the RNN.
- A **predictive head** uses the latest embedding vector to produce a **probability over the vocabulary**.

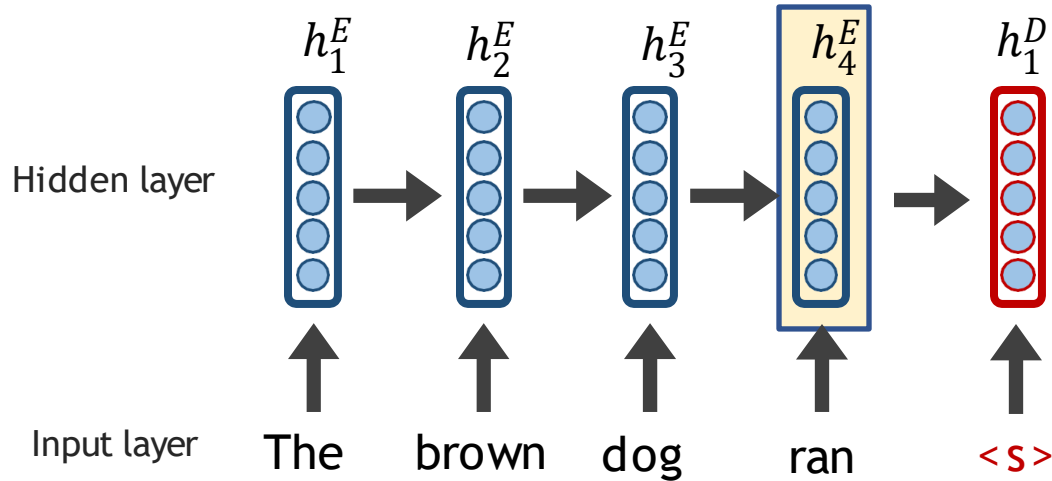
Example Inference (Generation) with RNN



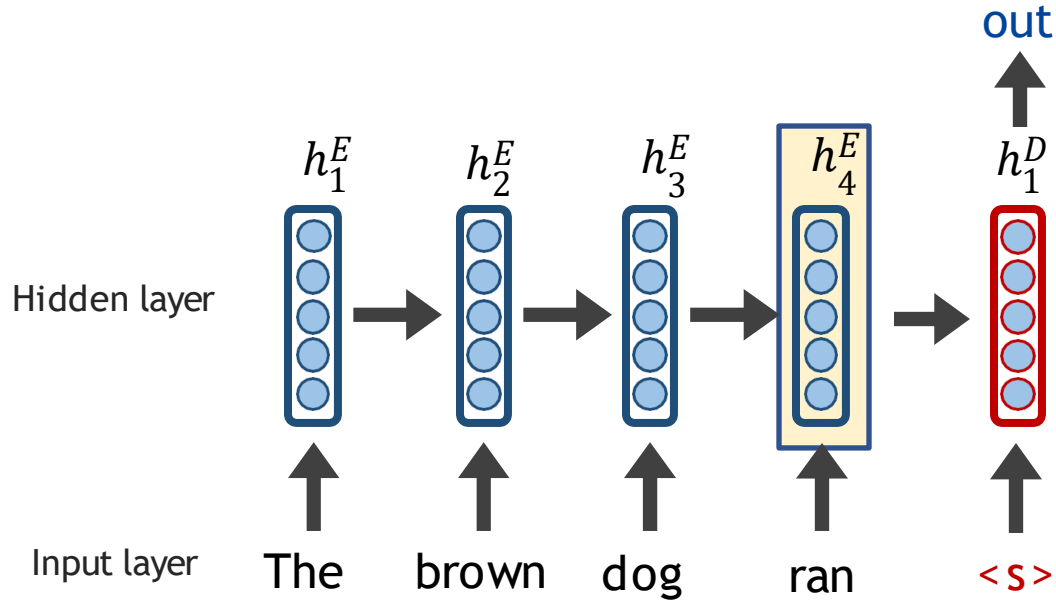
Example Inference (Generation) with RNN



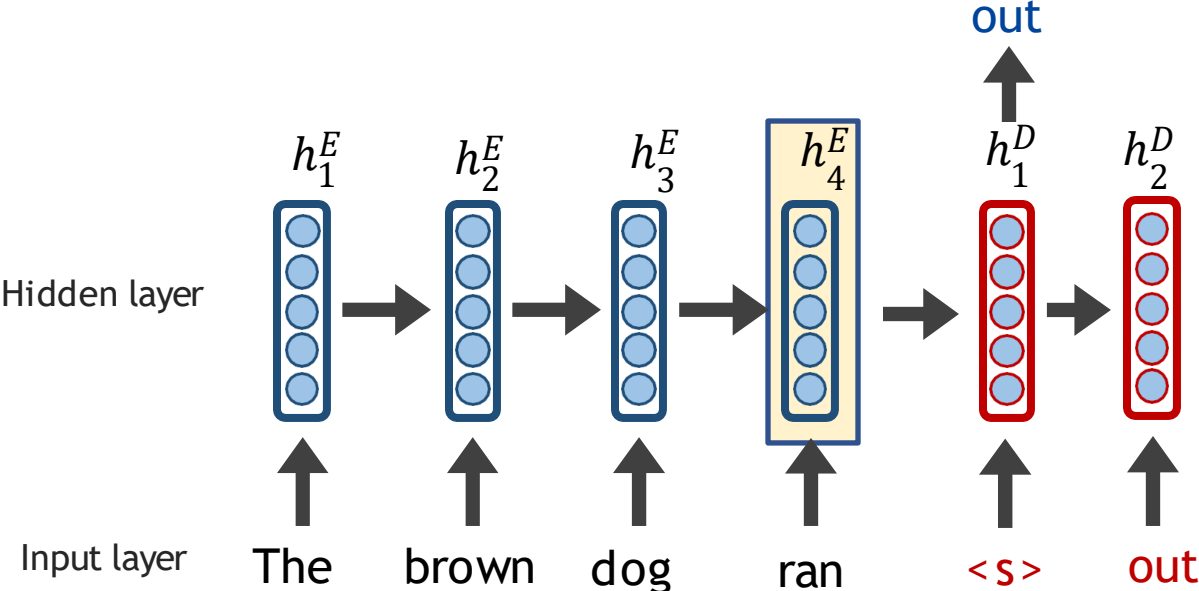
Example Inference (Generation) with RNN



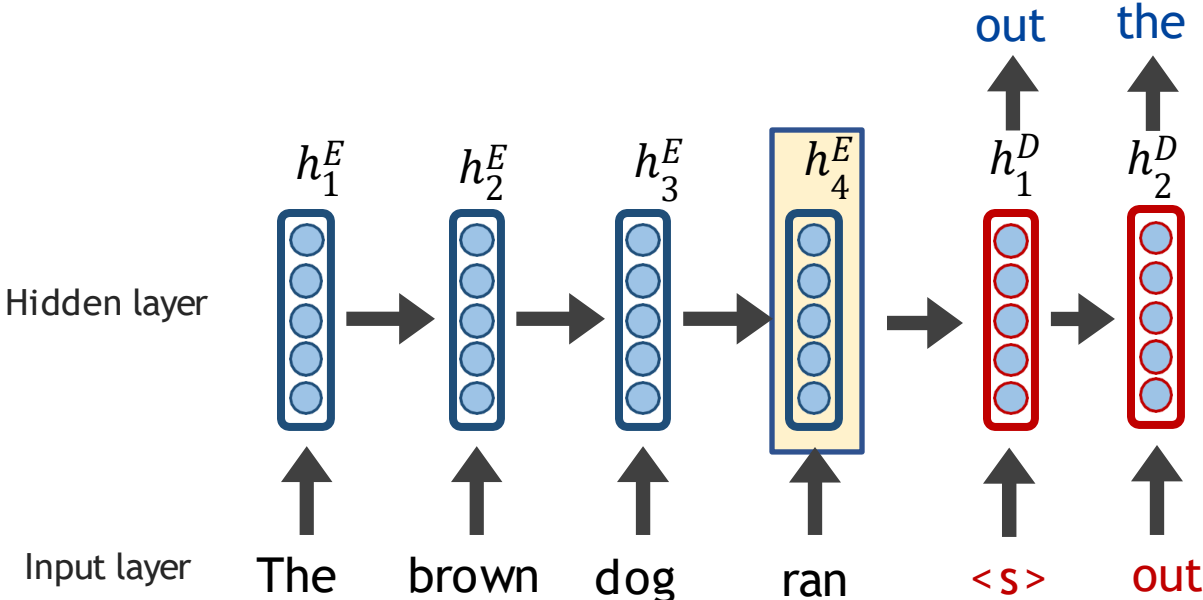
Example Inference (Generation) with RNN



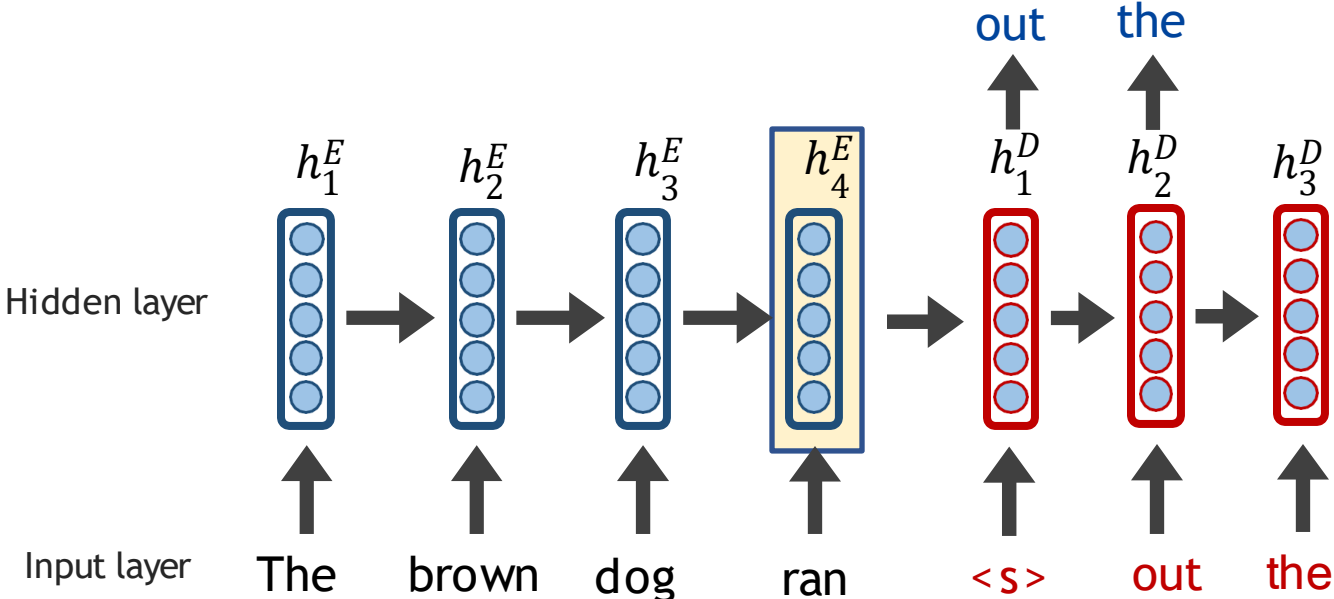
Example Inference (Generation) with RNN



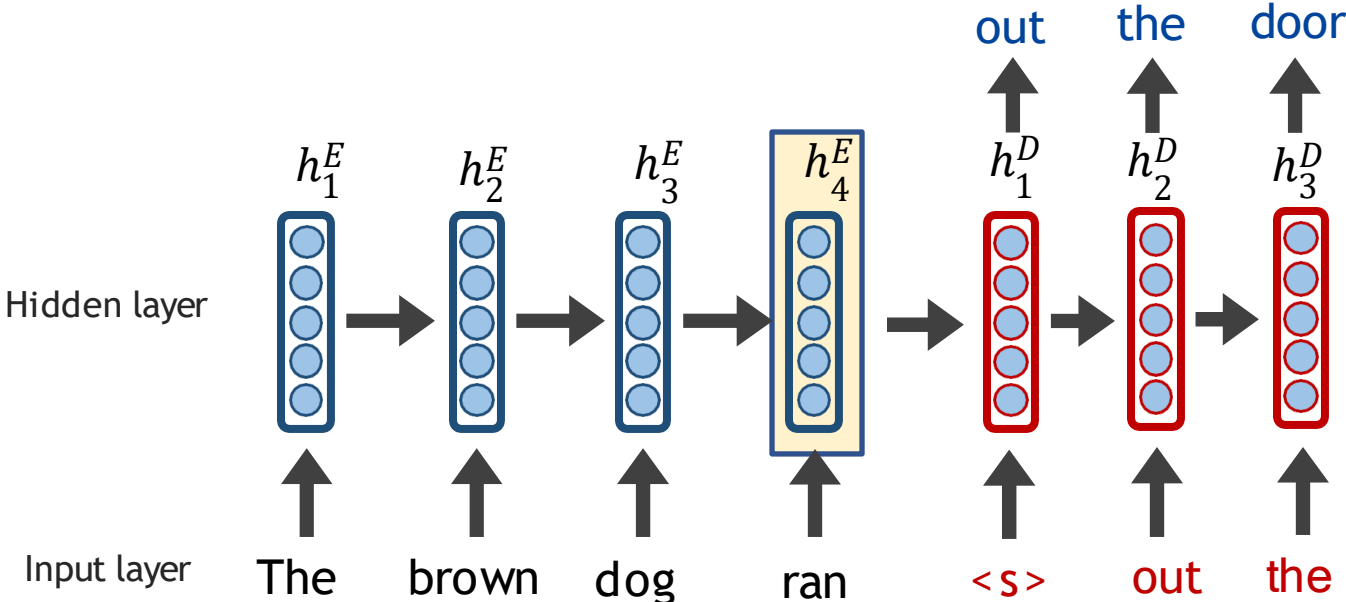
Example Inference (Generation) with RNN



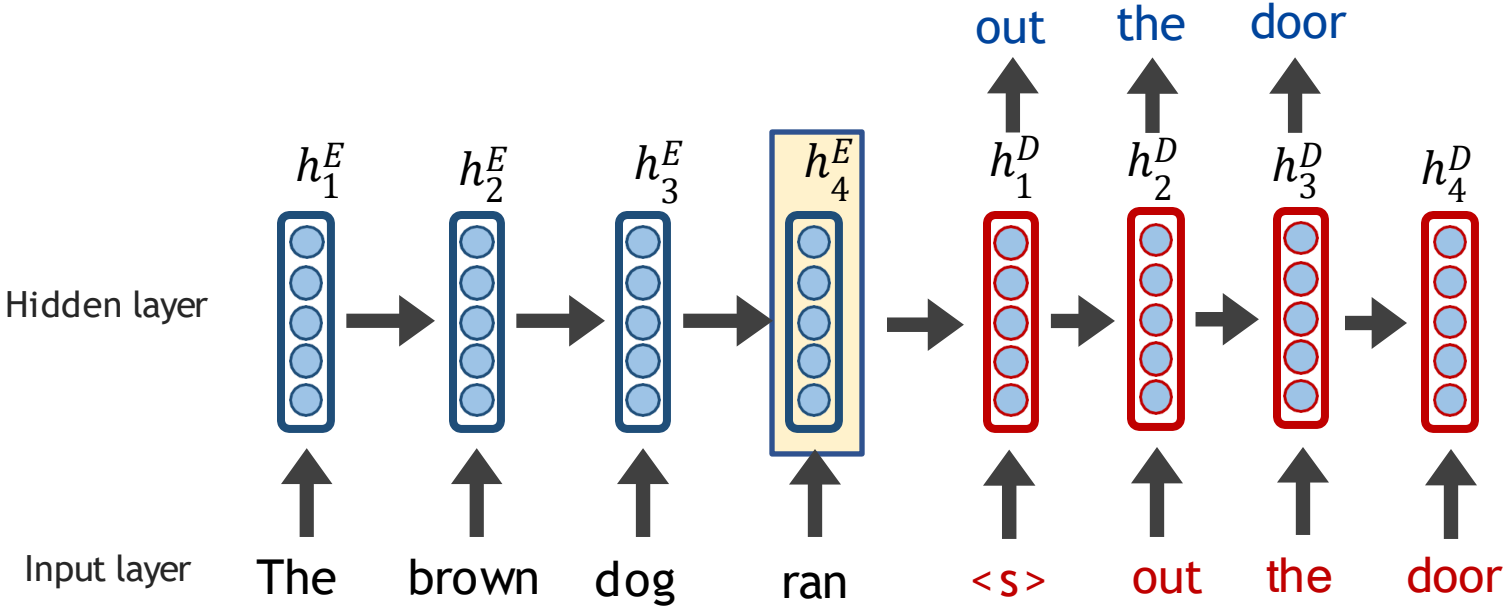
Example Inference (Generation) with RNN



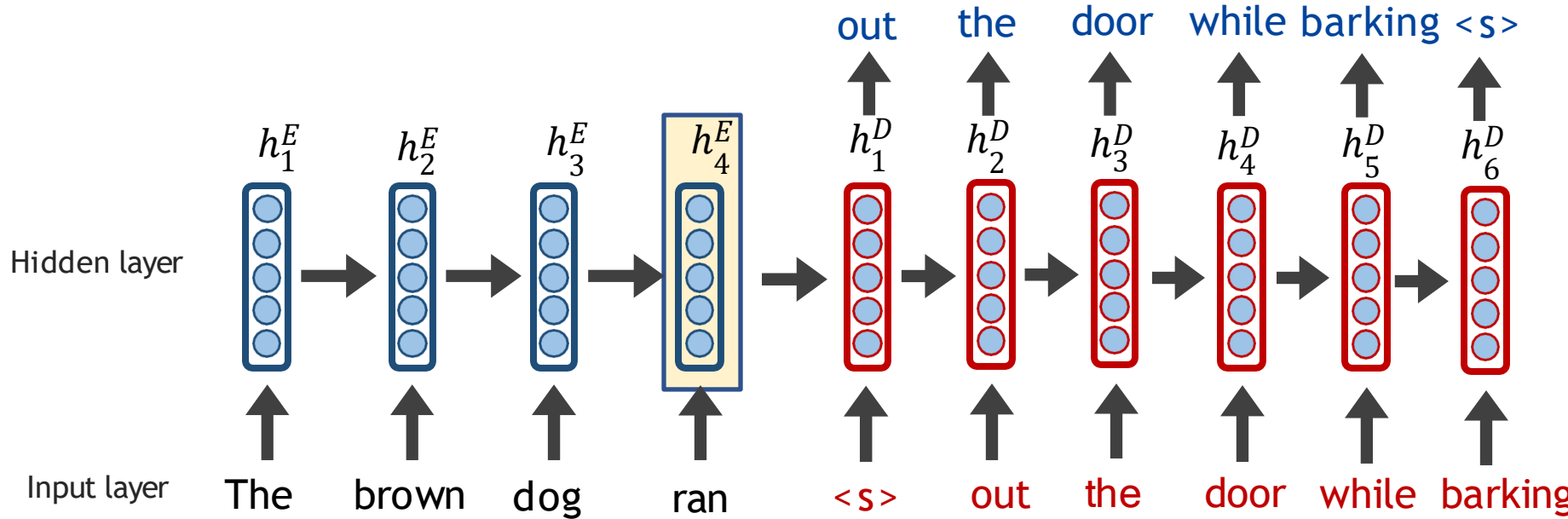
Example Inference (Generation) with RNN



Example Inference (Generation) with RNN



Example Inference (Generation) with RNN



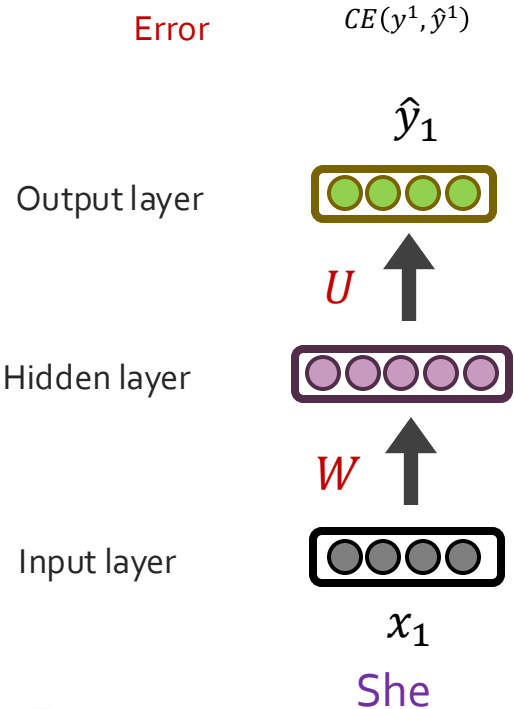


Training RNNs



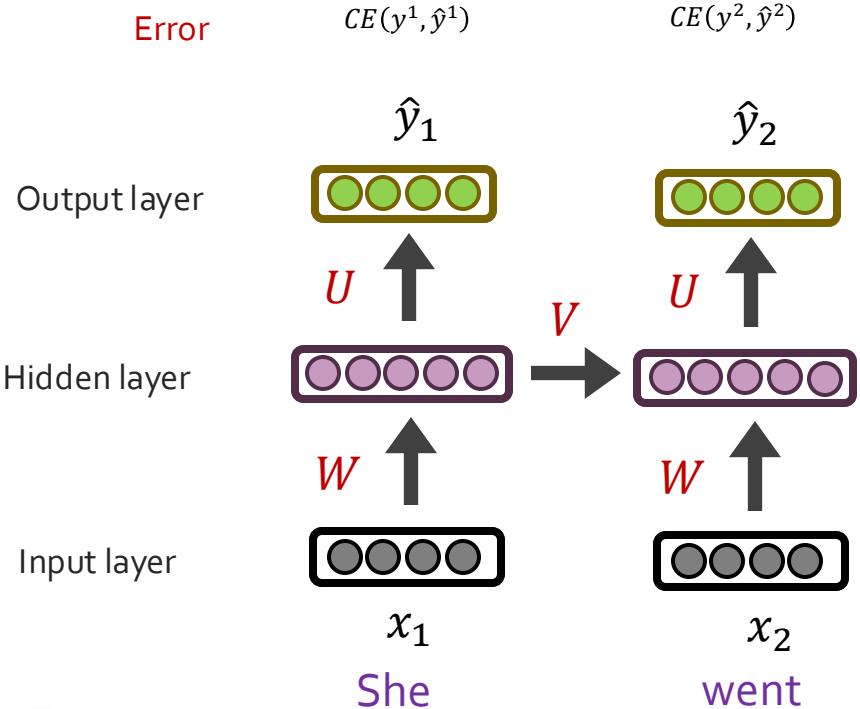
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



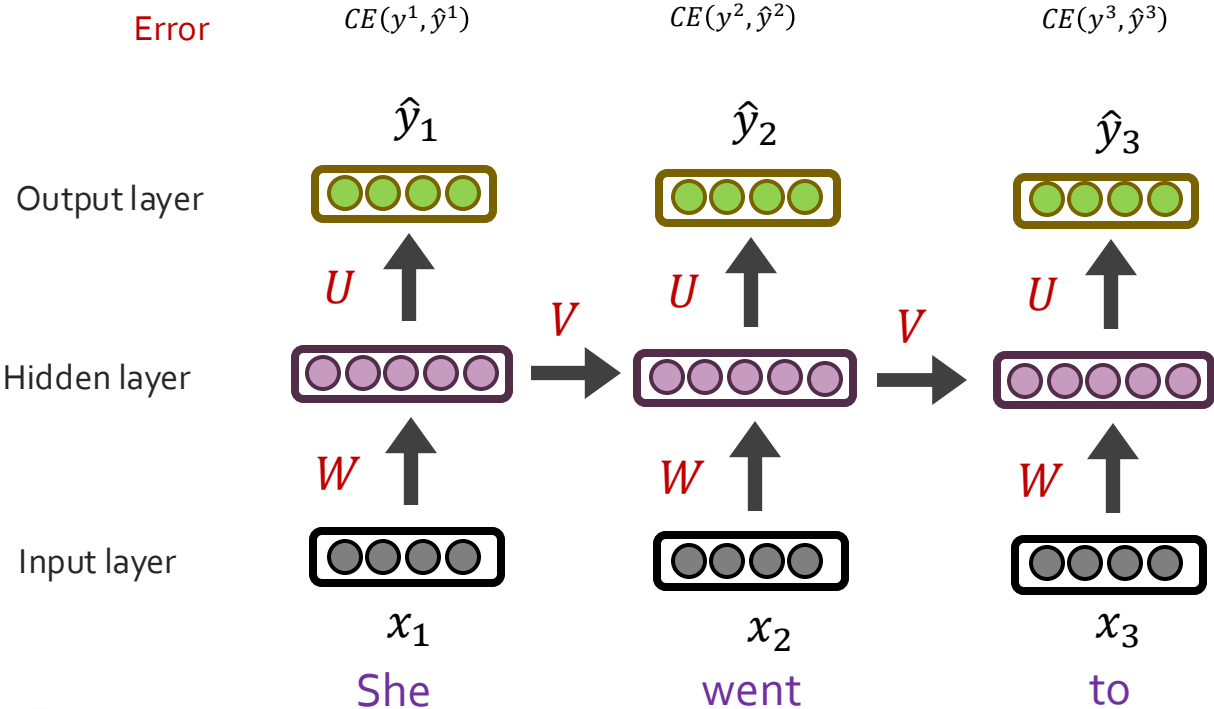
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



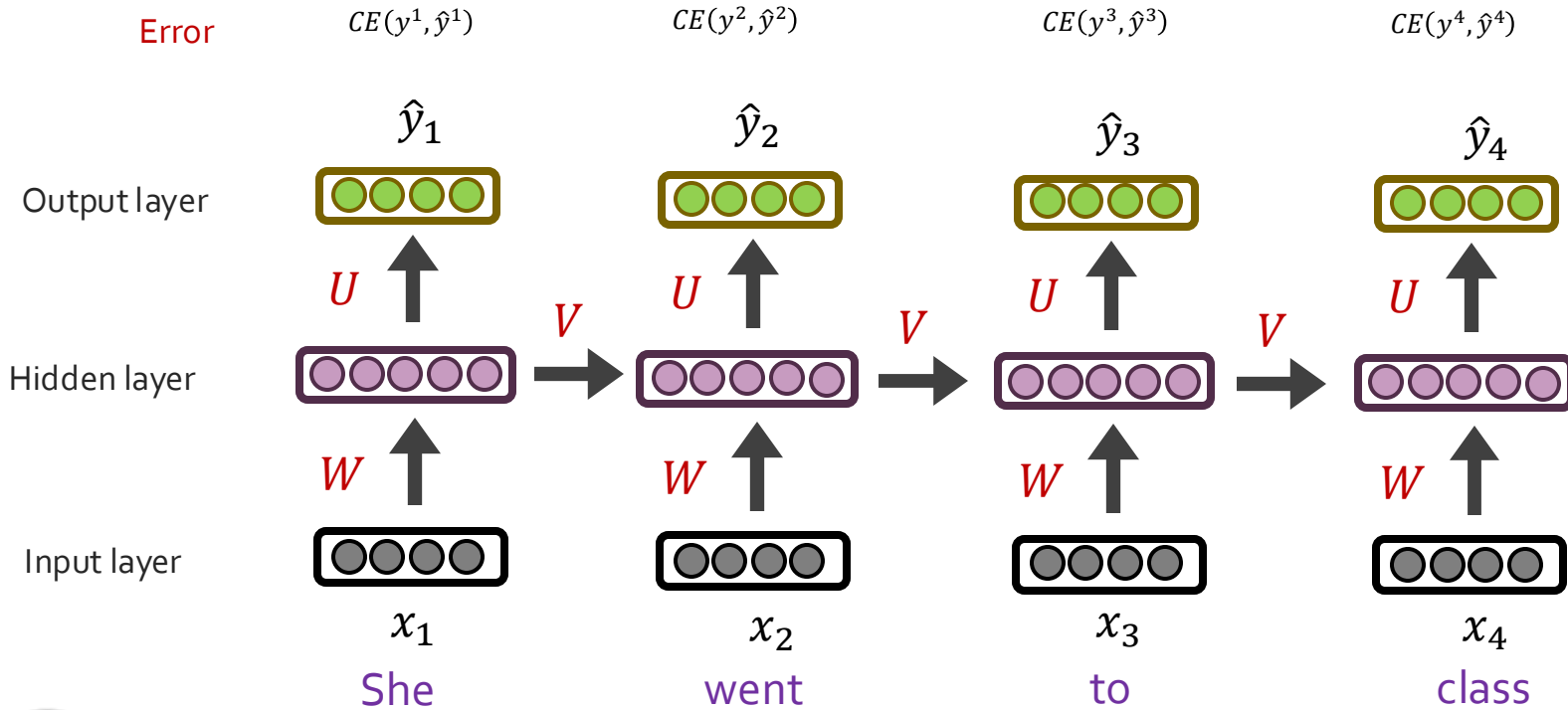
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



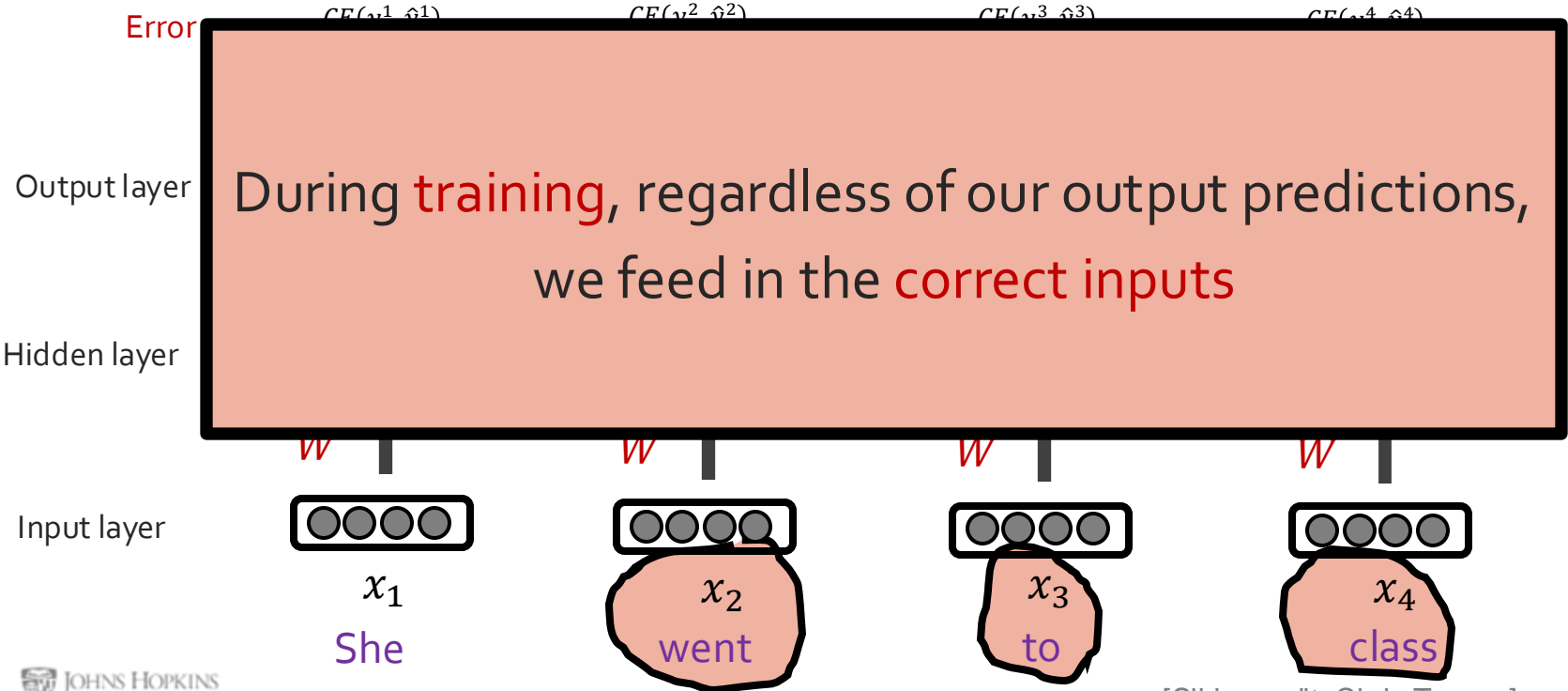
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



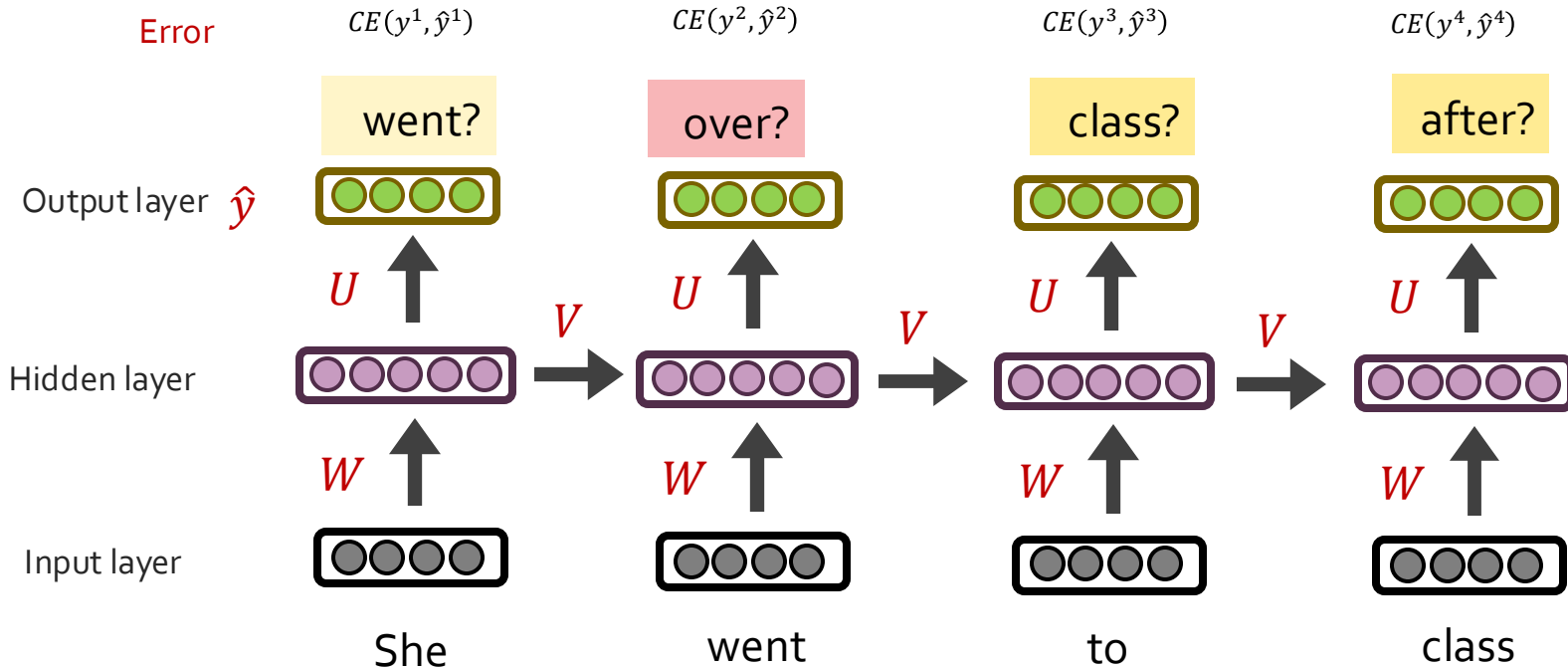
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



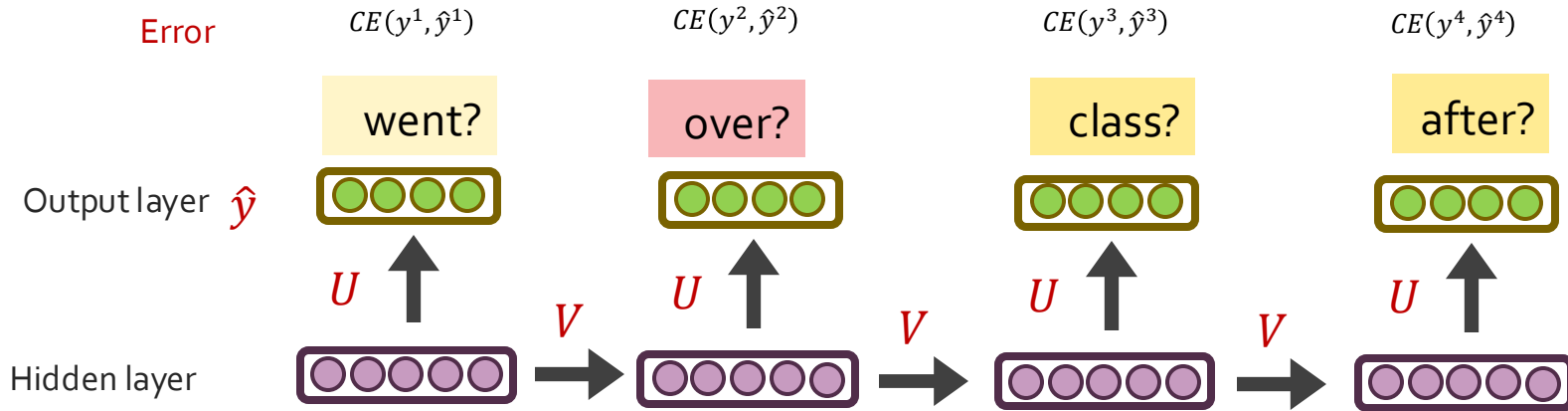
RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



RNN: Forward Propagation

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$



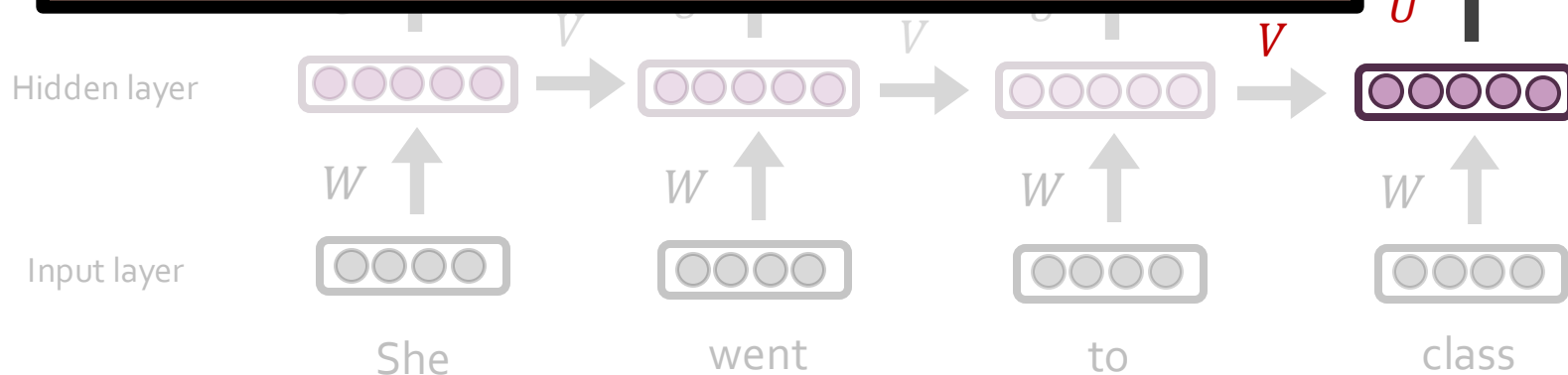
Our **total loss** is simply the **average loss** across all T time steps

Backward Step

$$CE(y^i, \hat{y}^i) = - \sum_{w \in V} y_w^i \log(\hat{y}_w^i)$$

To update our weights (e.g. Θ), we calculate the gradient of our loss w.r.t. the repeated weight matrix (e.g., $\frac{\partial L}{\partial \Theta}$).

Using the chain rule, we trace the derivative all the way **back to the beginning**, while summing the results.



Backward Step

To update our weights (e.g. Θ), we calculate the gradient of our loss w.r.t. the repeated weight matrix (e.g., $\frac{\partial L}{\partial \Theta}$).

Using the chain rule, we trace the derivative all the way **back to the beginning**, while summing the results.

Hidden layer



Input layer



She

went

to

class

$$\frac{\partial L}{\partial V}$$

$$CE(y^4, \hat{y}^4)$$

U

V^3

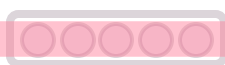
6

Backward Step

To update our weights (e.g. Θ), we calculate the gradient of our loss w.r.t. the repeated weight matrix (e.g., $\frac{\partial L}{\partial \Theta}$).

Using the chain rule, we trace the derivative all the way **back to the beginning**, while summing the results.

Hidden layer



Input layer



She

went

to

class

$$\frac{\partial L}{\partial V}$$

$$CE(y^4, \hat{y}^4)$$

U

V^2

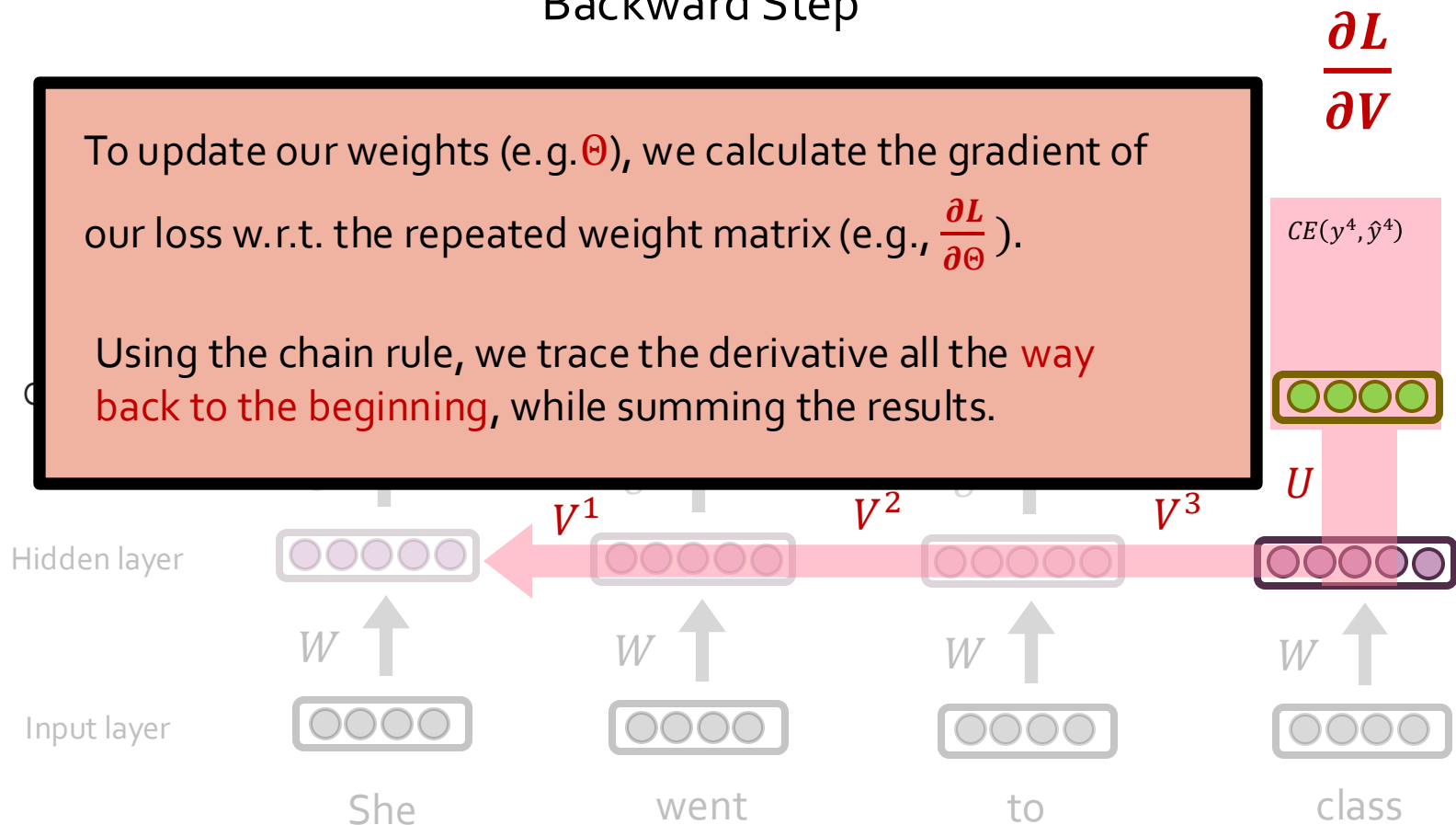
V^3

○

Backward Step

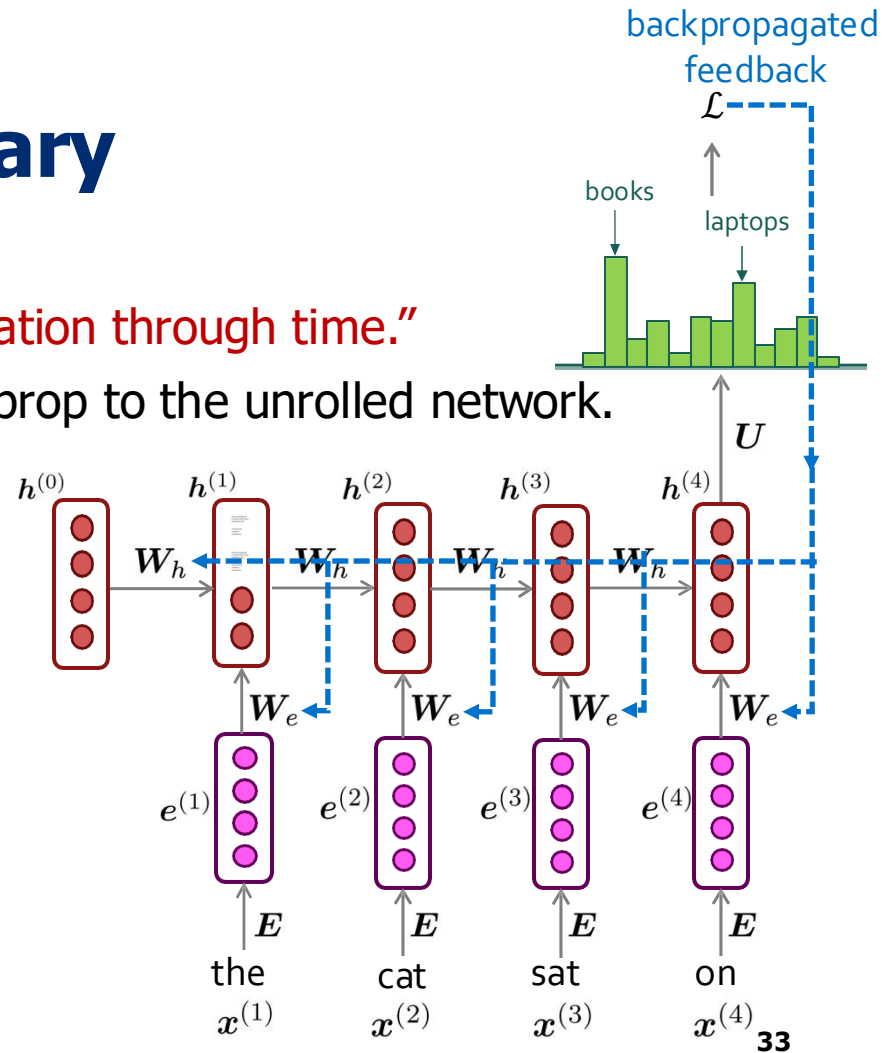
To update our weights (e.g. Θ), we calculate the gradient of our loss w.r.t. the repeated weight matrix (e.g., $\frac{\partial L}{\partial \Theta}$).

Using the chain rule, we trace the derivative all the way **back to the beginning**, while summing the results.



Training RNNs: Summary

- RNNs can be trained using “backpropagation through time.”
- Can be viewed as applying normal backprop to the unrolled network.
- Model’s learnable parameters Θ
 1. Compute $\mathcal{L}(\Theta)$ for a batch of sentences
 2. Compute gradients $\nabla_{\Theta} \mathcal{L}(\Theta)$
 3. Update the weights and then repeat





Examples



RNN: Generation



- When trained on Harry Potter text, it generates:

“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

RNNs: Generation



- RNN-LM trained on Obama speeches:

The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

RNNs in Practice



- RNN-LM trained on **food recipes**:

Title: CHOCOLATE RANCH BARBECUE

Categories: Game, Casseroles, Cookies, Cookies

Yield: 6 Servings

2 tb Parmesan cheese -- chopped

1 c Coconut milk

3 Eggs, beaten

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese. Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.

Evaluation LMs with Perplexity (2016)

n-gram model →

Increasingly
complex RNNs

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

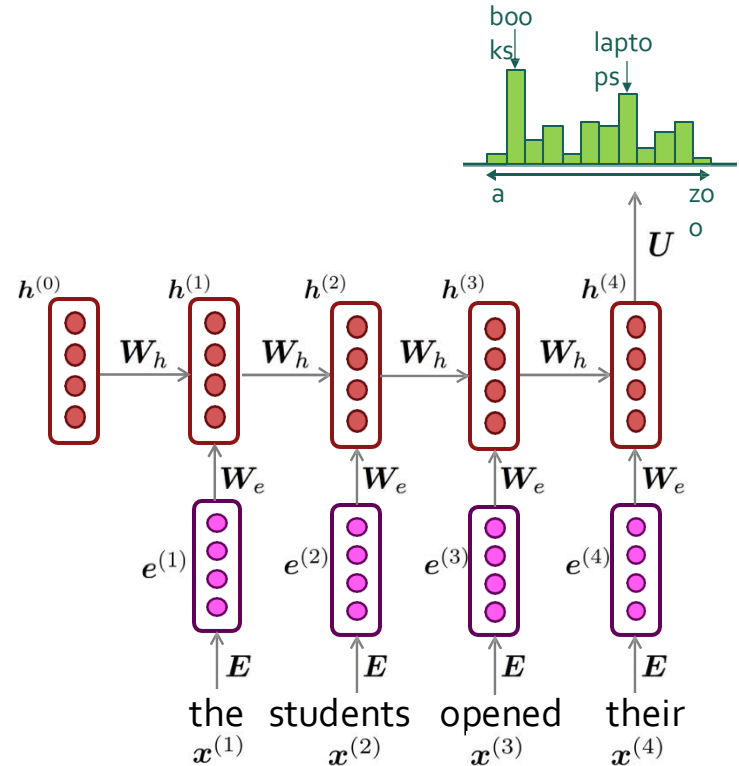
Summary

- RNNs: Repeated use of finite structure.
- A natural fit for language modeling.
- Next: let's summarize the pros and cons.

RNN-LMs: Pros and Cons

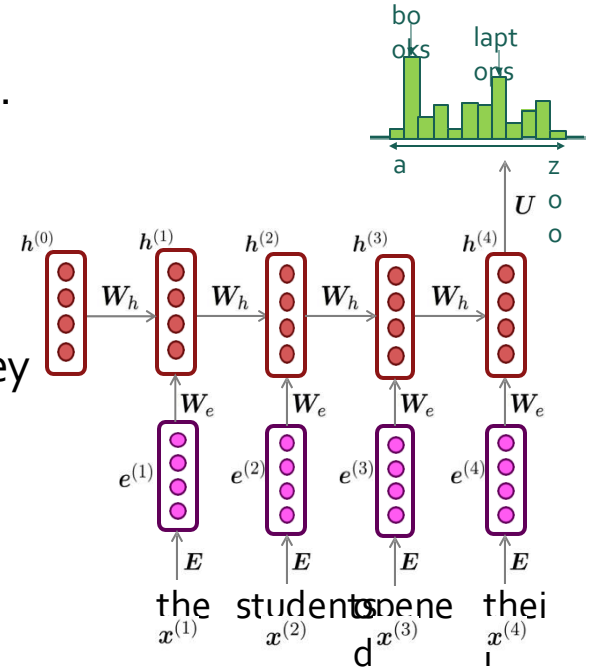
RNNs: Advantages

- Model size doesn't increase for longer inputs — reusing a compact set of model parameters.
- Computation for step t can (in theory) use information from many steps back



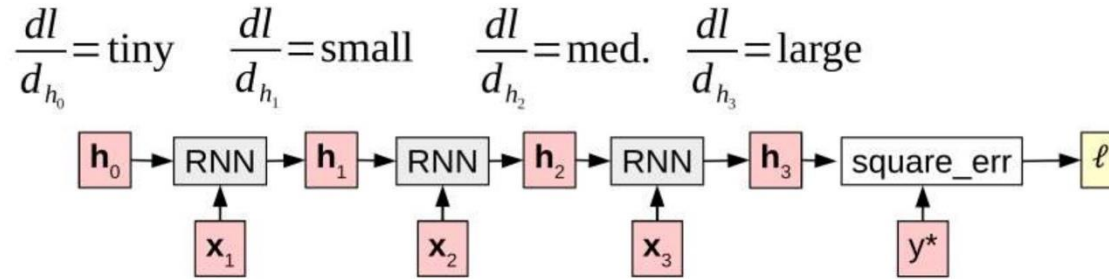
RNNs: Weaknesses

- Recurrent computation is **slow** and difficult to parallelize.
 - Next week: self-attention mechanism, better at representing long sequences and also parallelizable.
- While RNNs in theory can represent long sequences, they quickly **forget** portions of the input.
- Vanishing/exploding gradients.



Vanishing/Exploding Gradient Problem: Intuition

- Backpropagated errors multiply at each layer, resulting in

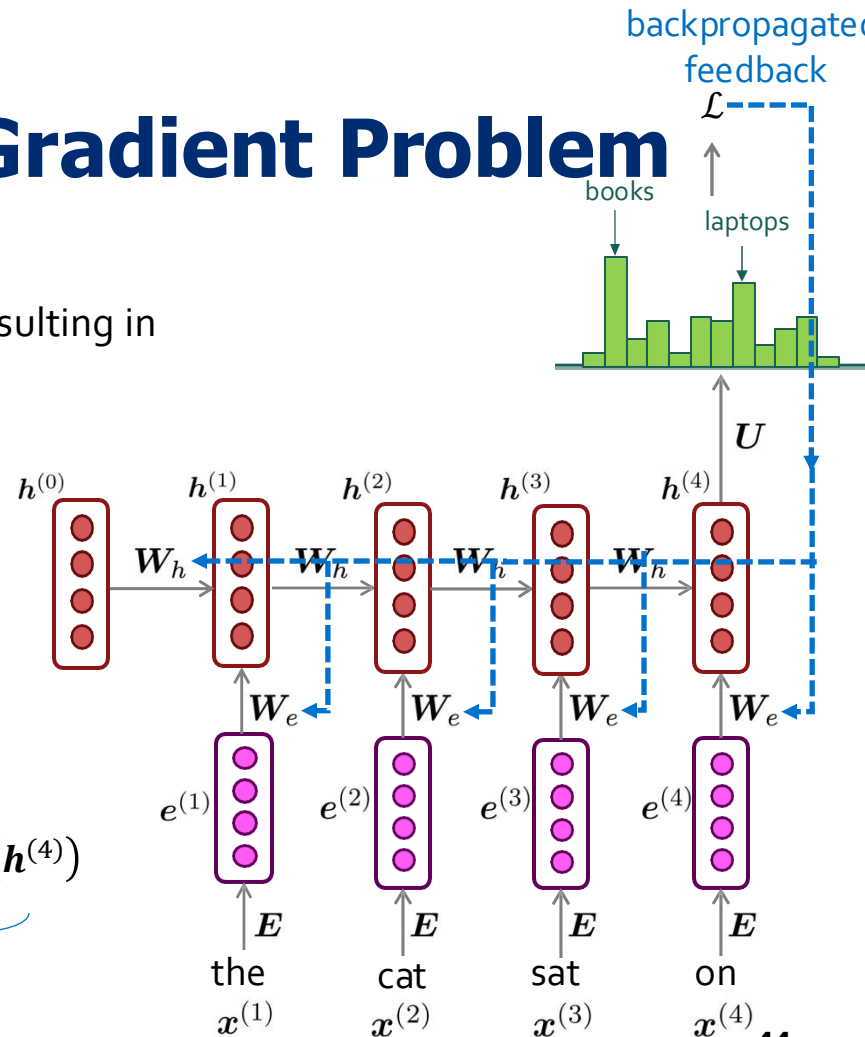


Vanishing/Exploding Gradient Problem

- Backpropagated errors multiply at each layer, resulting in exponential decay (if derivative is small) or growth (if derivative is large).
- Makes it very difficult to train deep networks, or simple recurrent networks over many time steps.

$$\nabla_{\mathbf{W}_h} \mathcal{L} = (\mathbf{J}_{\mathcal{L}}(\mathbf{W}_{L-1}))^T = \sum_{t=0} \left(\mathbf{J}_{\mathcal{L}}(\mathbf{h}^{(t)}) \mathbf{J}_{\mathbf{h}^{(t)}}(\mathbf{W}_h) \right)^T$$

$$\mathbf{J}_{\mathcal{L}}(\mathbf{h}^{(0)}) = \underbrace{\mathbf{J}_{\mathbf{h}^{(1)}}(\mathbf{h}^{(0)}) \mathbf{J}_{\mathbf{h}^{(2)}}(\mathbf{h}^{(1)}) \times \dots \times \mathbf{J}_{\mathbf{h}^{(4)}}(\mathbf{h}^{(3)}) \mathbf{J}_{\mathcal{L}}(\mathbf{h}^{(4)})}_{\text{chain rule}}$$



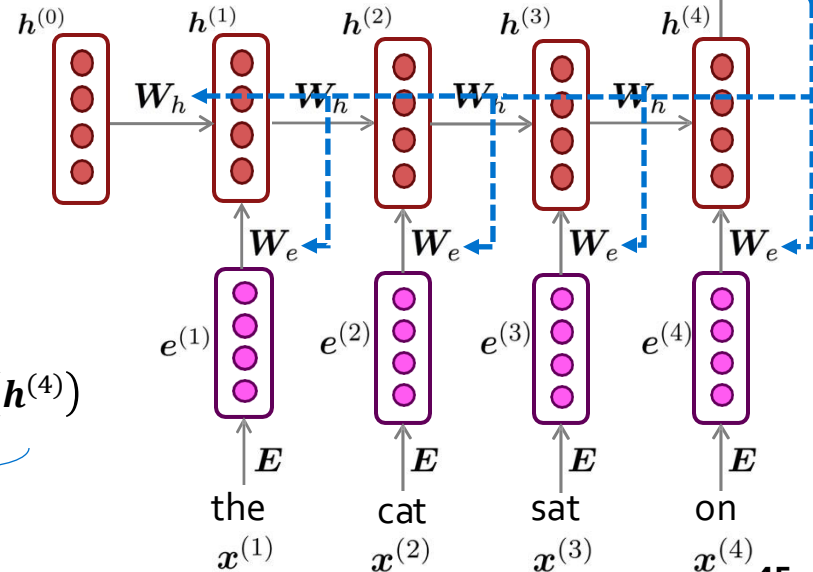
Vanishing/Exploding Gradient Problem

- Note:** instability of matrix powers can be determined from their eigenvalues.

Gradient signal from far away is lost.
So, model weights are updated only
with respect to near effects, not long-
term effects.

$$J_{\mathcal{L}}(\mathbf{h}^{(0)}) = J_{\mathbf{h}^{(1)}}(\mathbf{h}^{(0)}) J_{\mathbf{h}^{(2)}}(\mathbf{h}^{(1)}) \times \dots \times J_{\mathbf{h}^{(4)}}(\mathbf{h}^{(3)}) J_{\mathcal{L}}(\mathbf{h}^{(4)})$$

chain rule



Coping with Exploding Gradients

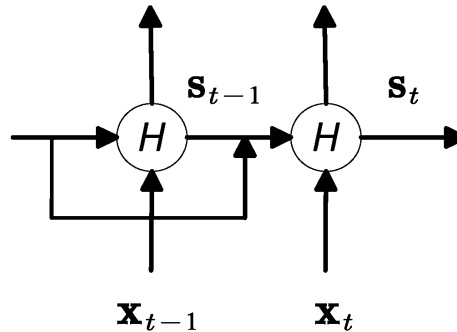
- **Gradient clipping:**
 - If the norm of the gradient is greater than some threshold, scale it down before applying SGD update.
- **Intuition:** take a step in the same direction, but a smaller step

Algorithm 1 Pseudo-code for norm clipping

```
 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$   
if  $\|\hat{\mathbf{g}}\| \geq \textit{threshold}$  then  
     $\hat{\mathbf{g}} \leftarrow \frac{\textit{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$   
end if
```

Coping with Vanishing Gradients

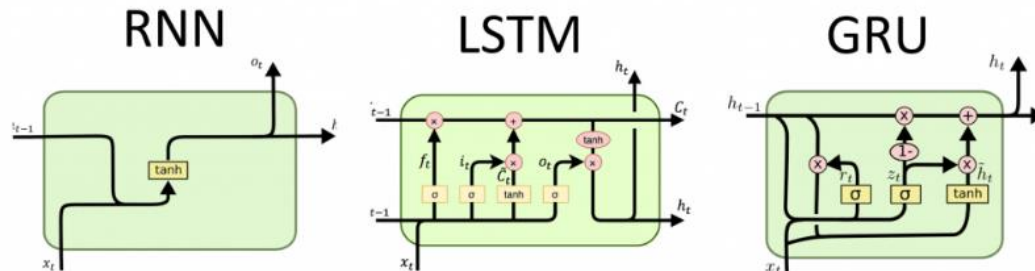
- **Using residual connections:**
 - lots of new deep architectures (RNN or otherwise) add direct connections, thus allowing the gradient to flow)



Res-RNN

RNNs: Difficulty in Learning Long-Range Dependencies (3)

- While RNNs in theory can represent long sequences, in practice teaching them about long-range dependencies is **non-trivial**.
- Changes to the architecture makes it easier for the RNN to preserve information over many timesteps
 - Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber 1997, Gers+ 2000]
 - Gated Recurrent Units (GRU) [Cho+ 2014]



RNNs: Difficulty in Learning Long-Range Dependencies (3)

- While RNNs in theory can represent long sequences, in practice teaching them about long-range dependencies is **non-trivial**.
- Changes to the architecture makes it easier for the RNN to preserve information over many timesteps
 - Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber 1997, Gers+ 2000]
 - Gated Recurrent Units (GRU) [Cho+ 2014]
- Many of these variants were the dominant architecture of In 2013–2015.
- We will not cover these alternative architecture in favor or spending more time on more modern developments.

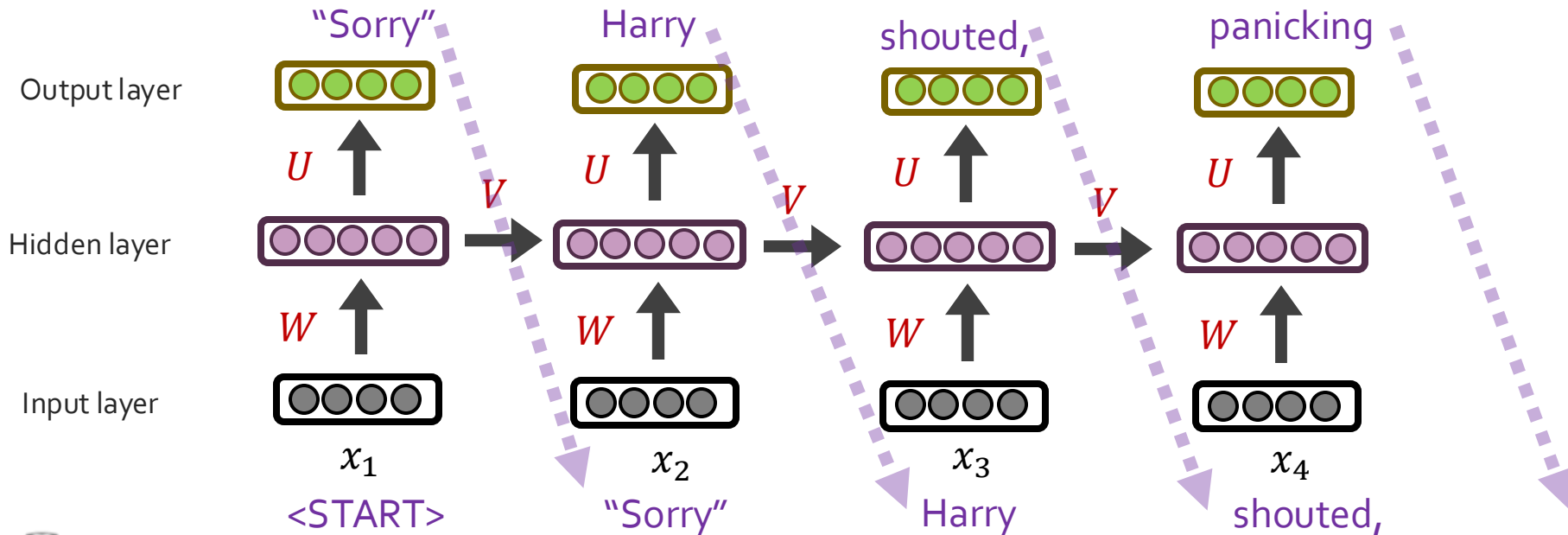
Summary

- RNNs provide a compact model, regardless of sequence size. In theory, this is great!
- In practice, however:
 - They still struggle with remembering long-range dependencies.
 - Training them is not difficult because of vanishing/exploding gradients.
- Despite these limits, RNNs provided improvements at the time that they were introduced and laid the foundation for the future progress.

Bonus: Pre-training RNN Language Models

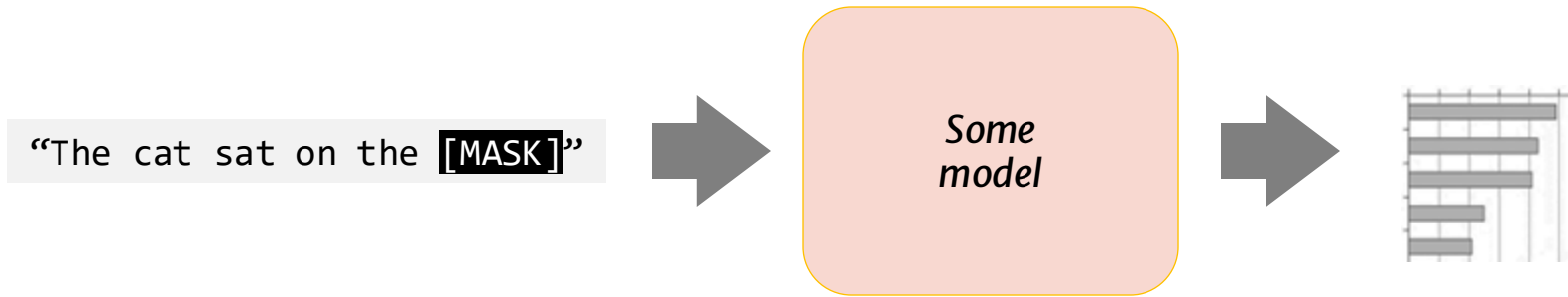
Recap: Recurrent Neural Networks

- Repeated use of a **finite** model

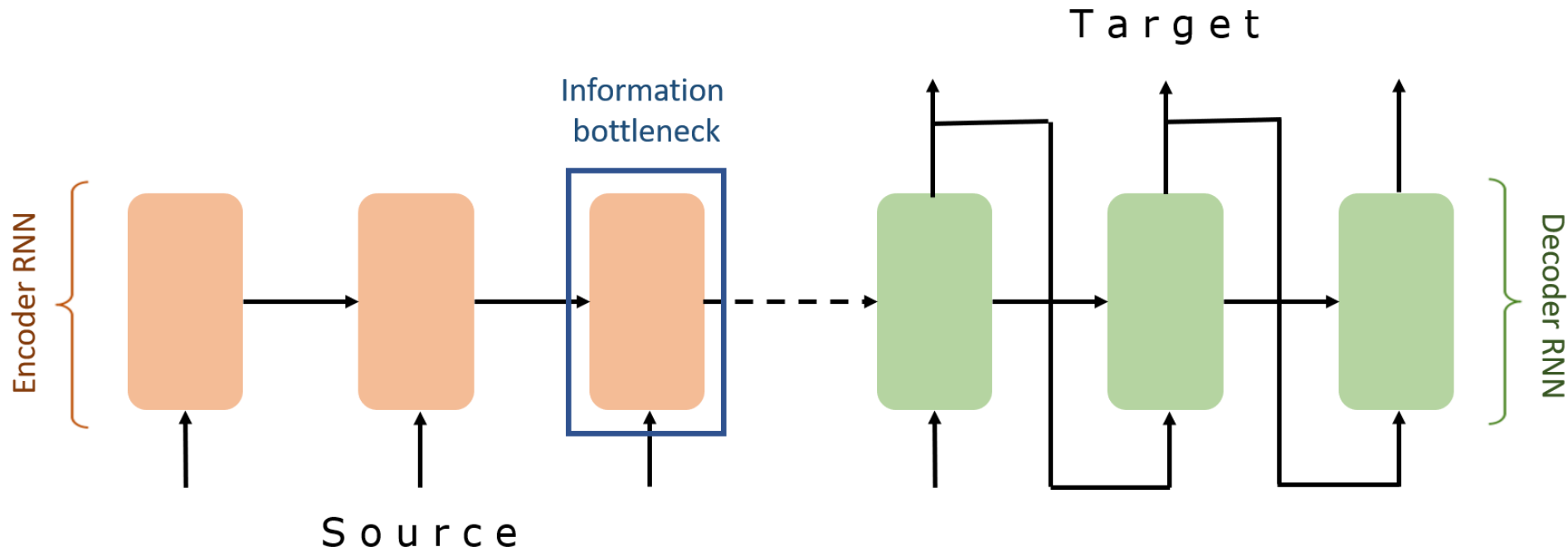


Recap: Encoder-Decoder Architectures

- It is useful to think of generative models as two sub-models.



Recap: Encoder-Decoder Architectures



Contextual Meaning of Words



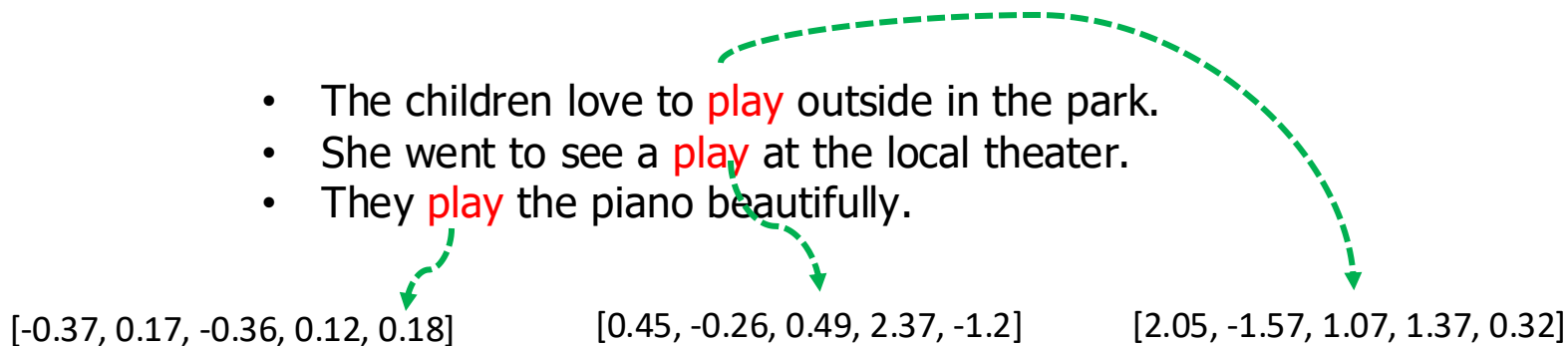
- Earlier word embedding methods (e.g., Word2Vec, GloVe) learn a single “static” vector for each word.
 - Static embeddings are **not** flexible and expressive enough.
 - The children love to **play** outside in the park.
 - She went to see a **play** at the local theater.
 - They **play** the piano beautifully.

Information from context is necessary to capture the correct meaning of the word.

ELMo: First Major Self-Supervised LM



- **Goal:** get highly rich, contextualized embeddings (word tokens) that depend on the entire sentence in which a word is used.

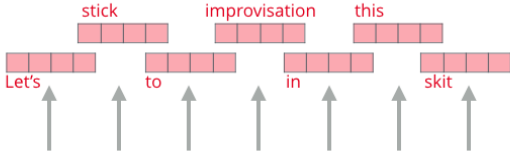




ELMo: First Major Self-Supervised LM

- **Goal:** get highly rich, contextualized embeddings (word tokens) that depend on the entire sentence in which a word is used.

ELMo
Embeddings

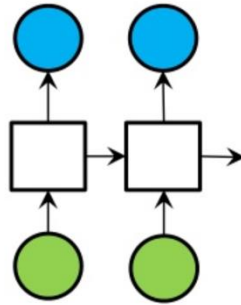


Words to embed

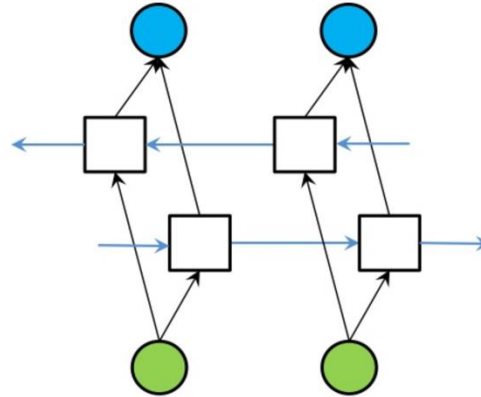


Extending RNNs to Both Directions

- An RNN limitation: Hidden variables capture **only one side of the context**.
- Solution: Bi-Directional RNNs



RNN

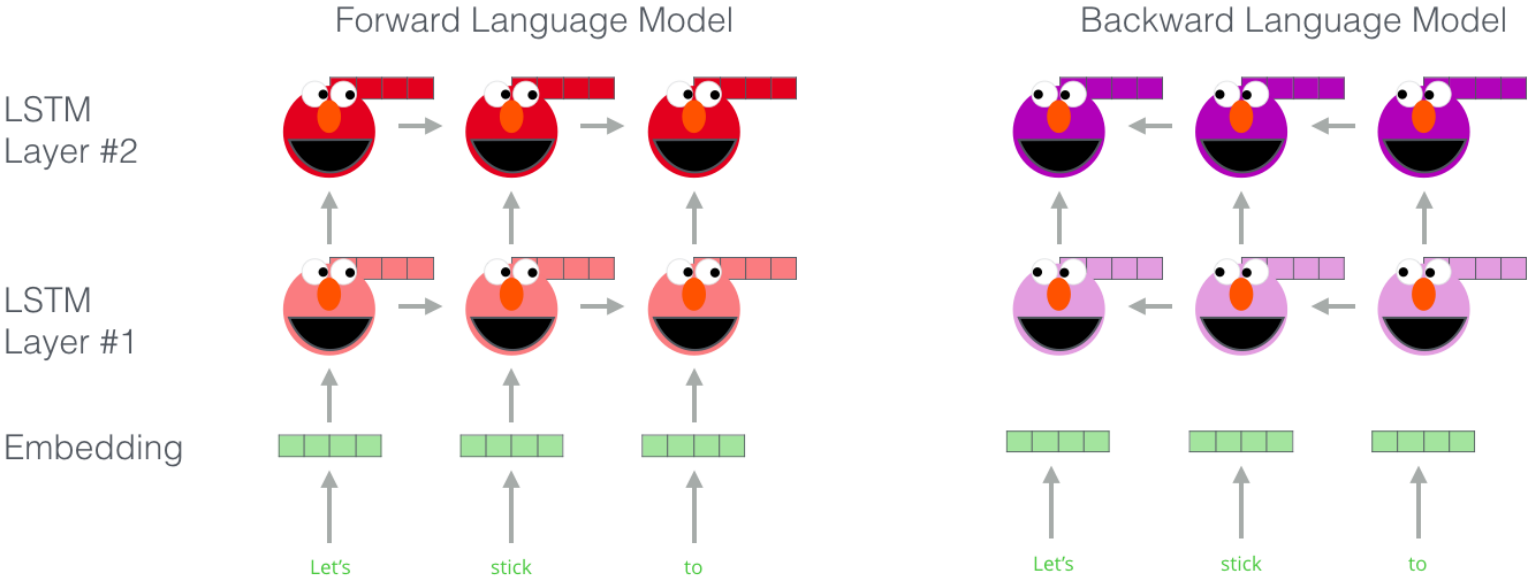


Bi-directional RNN



ELMo: First Major Self-Supervised LM

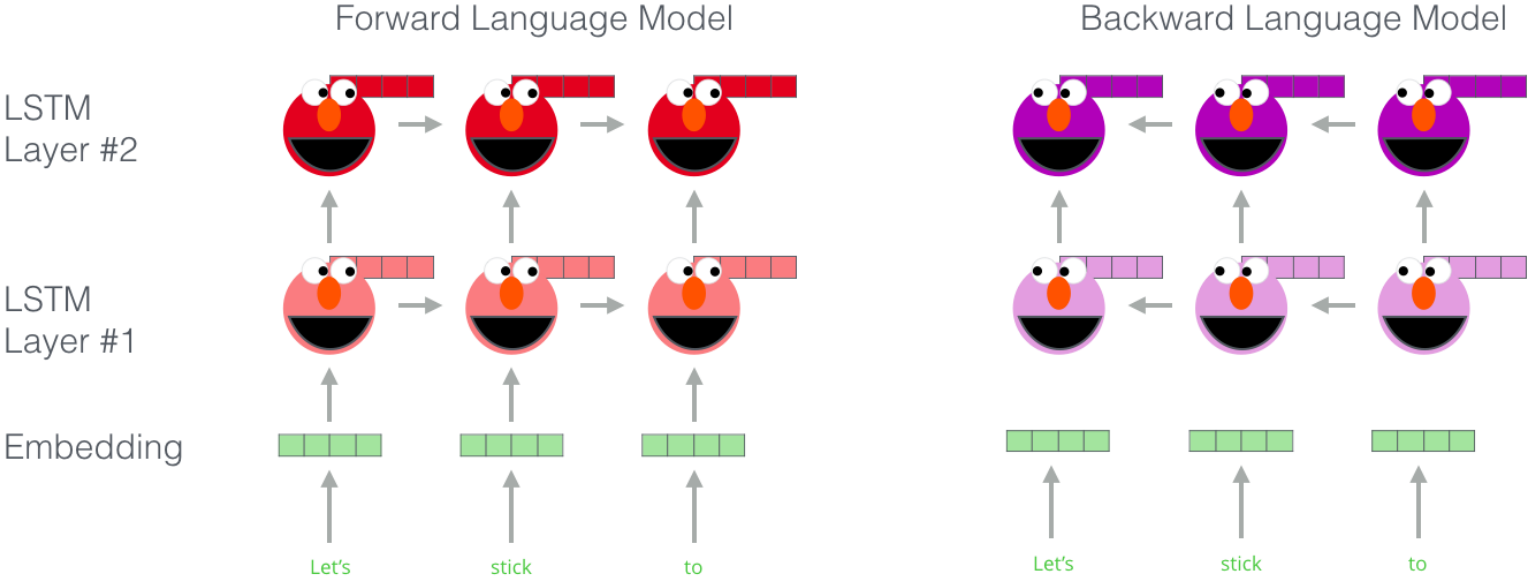
- Use both directions of context (bi-directional), with increasing abstractions (stacked)
 - Two LSTMs in different directions — capture both directions





ELMo: First Major Self-Supervised LM

- Linearly combine all abstract representations (hidden layers) and optimize w.r.t. a particular task (e.g., sentiment classification)



ELMo: Some Details

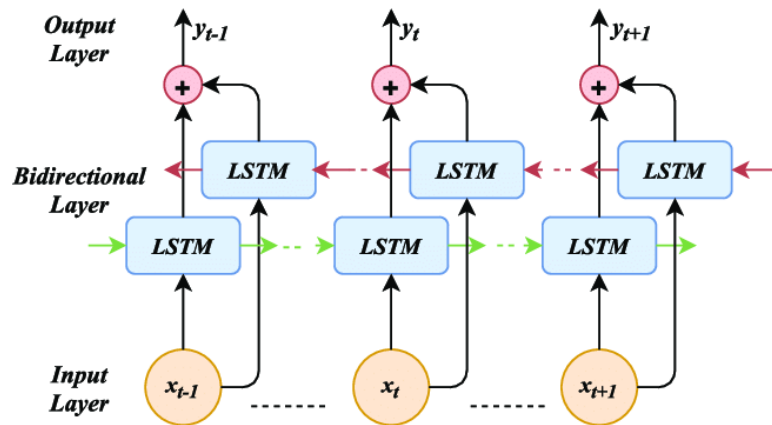
- Train a forward language model by modeling prob of each word, given its left context.

$$p(t_1, \dots, t_k) = \prod_{k=1}^N p(t_k | t_1, \dots, t_{k-1})$$

- Similarly, train a backward language model, conditioned on the right context.

- Some training details:

- Use 4096 dim hidden states
- Residual connections from the first to second layer
- Trained 10 epochs on 1B Word Benchmark
- Results in perplexity of ~ 39

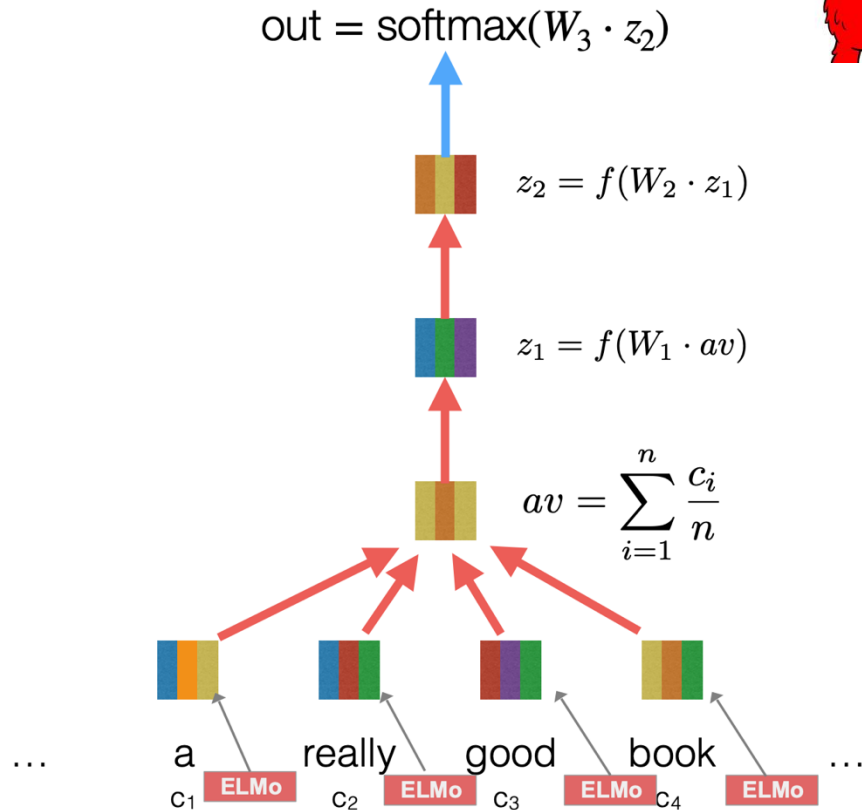


$$p(t_1, \dots, t_k) = \prod_{k=1}^N p(t_k | t_{k+1}, \dots, t_N)$$

Adapting ELMo Representations for Tasks



- Fine-tune classifiers using contextualized word representations extracted from ELMo.



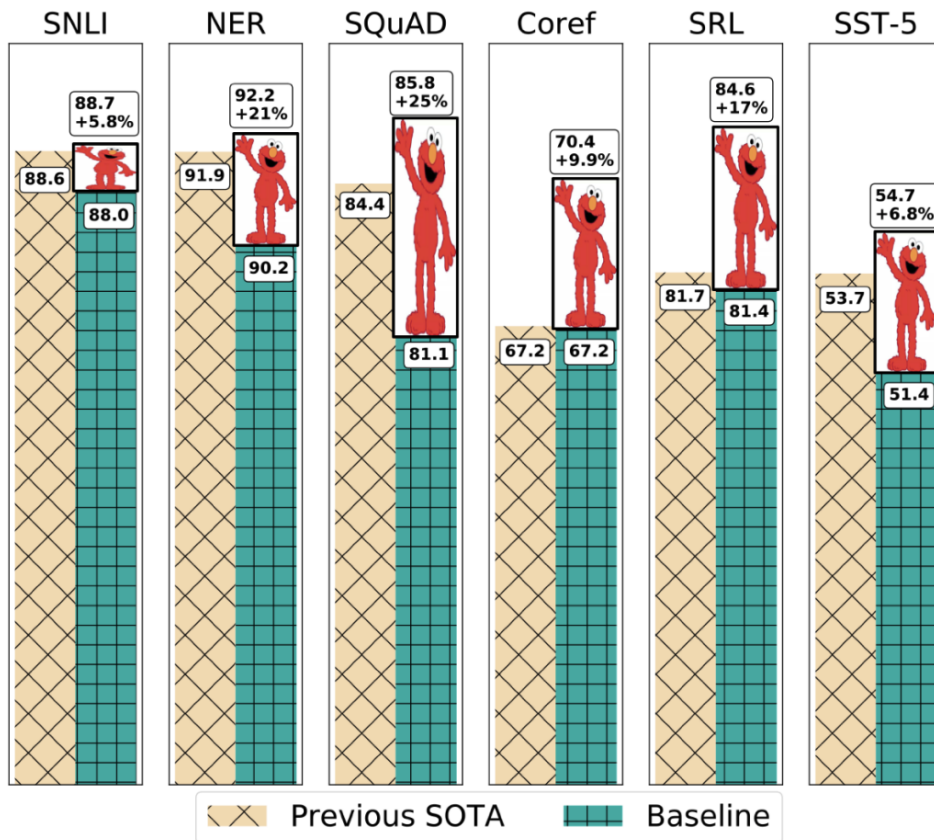
ELMo: Evaluation

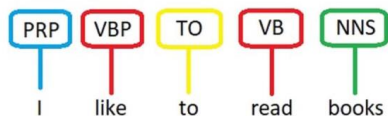
- SQuAD: question answering
- SNLI: textual entailment
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis



Barack Obama nominated **Hillary Rodham Clinton** as **his secretary of state** on Monday. **He** chose **her** because **she** had foreign affairs experience as a former **First Lady**.

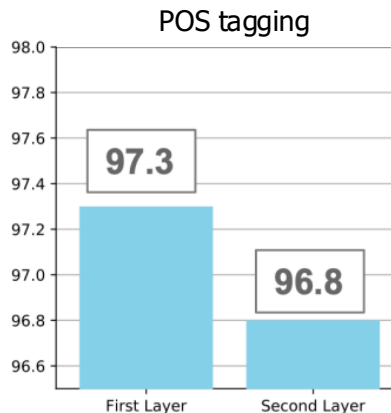
Experimental Results



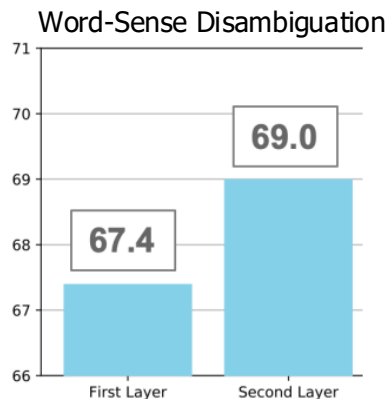


The **bank** can guarantee **deposits** will eventually cover future tuition costs because it invests in adjustable-rate **mortgage** securities.

bank ¹	Gloss:	a financial institution that accepts deposits and channels the money into lending activities
	Examples:	"he cashed a check at the bank", "that bank holds the mortgage on my home"
bank ²	Gloss:	sloping land (especially the slope beside a body of water)
	Examples:	"they pulled the canoe up on the bank", "he sat on the bank of the river and watched the currents"



First Layer > Second Layer

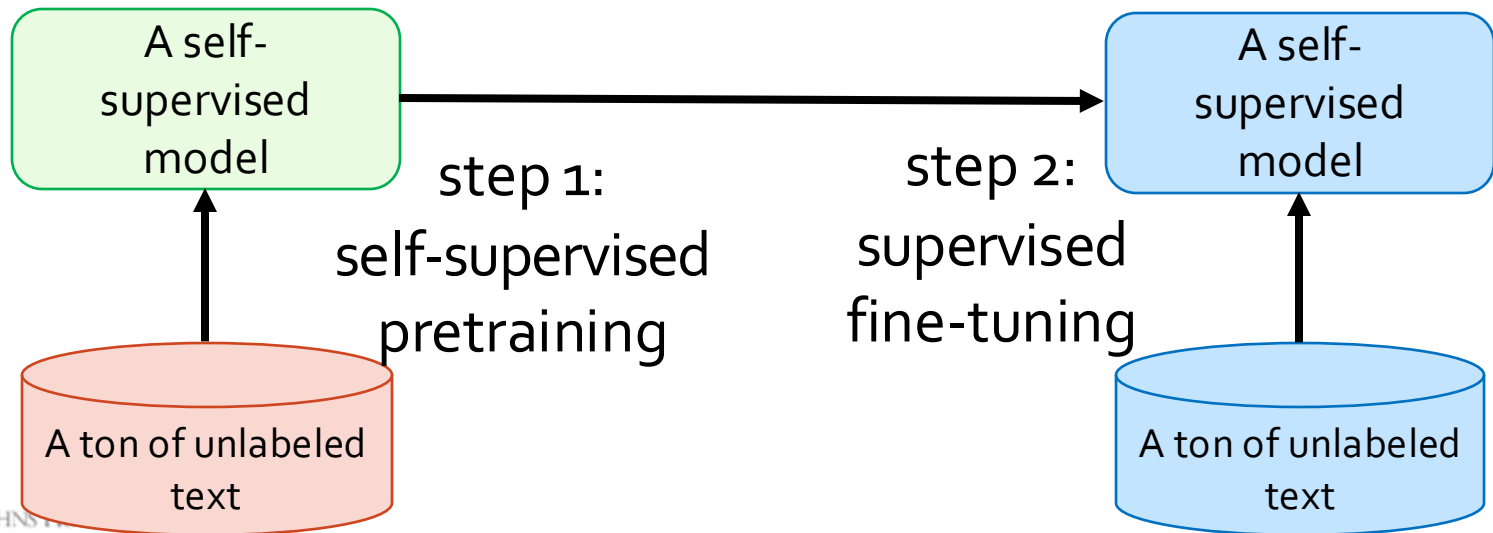


Second Layer > First Layer

Syntactic information is better represented at **lower layers** while **semantic** information is captured a **higher layers**

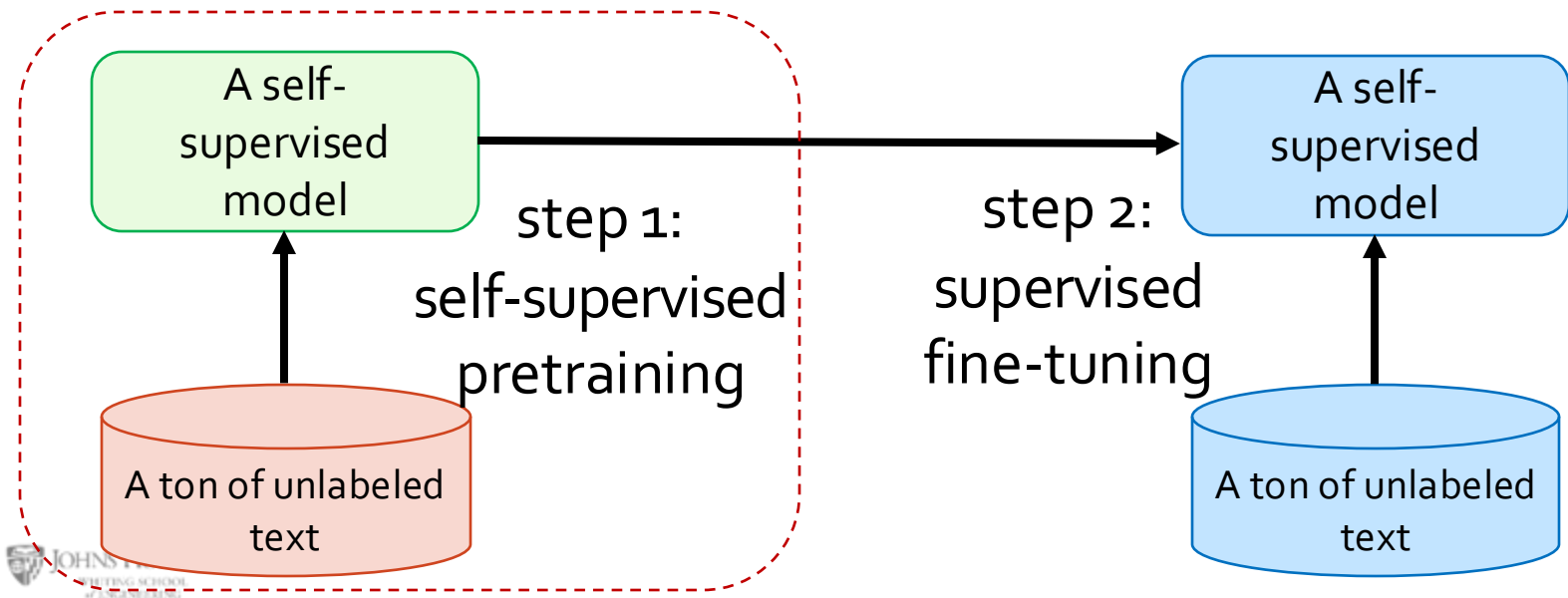
A broader lesson: Pre-training + Fine-tuning

- Let's say I want to train a model for sentiment analysis
- In the past, I would simply train a supervised model on Word2Vec representation of review sentences (e.g., HW2).



A broader lesson: Pre-training + Fine-tuning

Contextual representation of LMs is much **stronger** than word-level ones.
Now that we have Fixed-Window LM, we can use it to build better classifiers!



Summary



- ELMo: Stacked Bi-directional LSTMs
- ELMo yielded incredibly good **contextualized** embeddings, which yielded **SOTA** results when applied to **many** NLP tasks.
- **Main ELMo takeaway:** given **enough [unlabeled] training data**, having **tons of parameters** model is useful — the system can determine **how to best use context**.

Summary

- Recurrent Neural Networks

- A family of neural networks that allow architecture for inputs of variable length

- RNN-LM: LM based on RNNs

- A notable example: **ELMo**

- Cons:

- Sequential processing
- While in theory it maintain infinite history, in practice it suffers from long-range dependencies.

