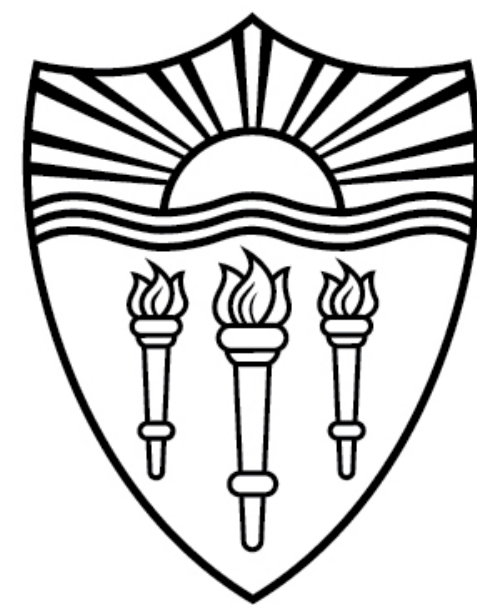


CSCI 544: Applied Natural Language Processing

# **Recurrent Neural Networks for Sequence Labeling**

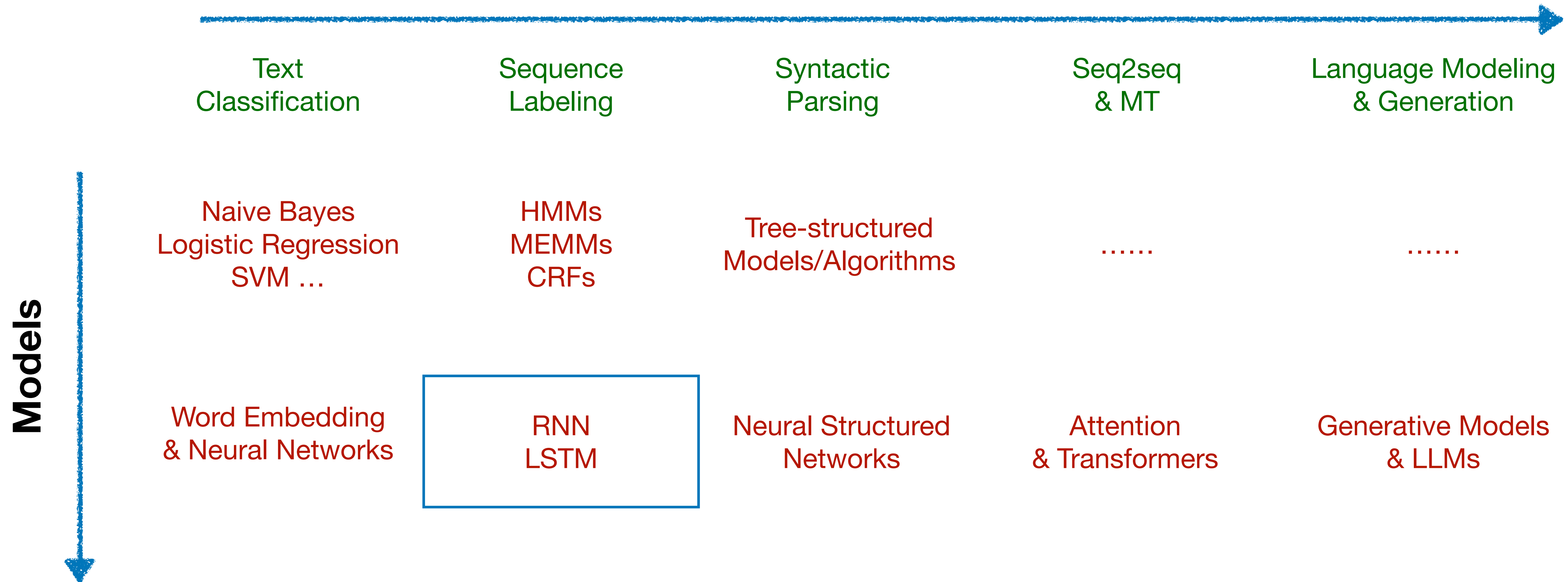
Xuezhe Ma (Max)



**USC** University of  
Southern California

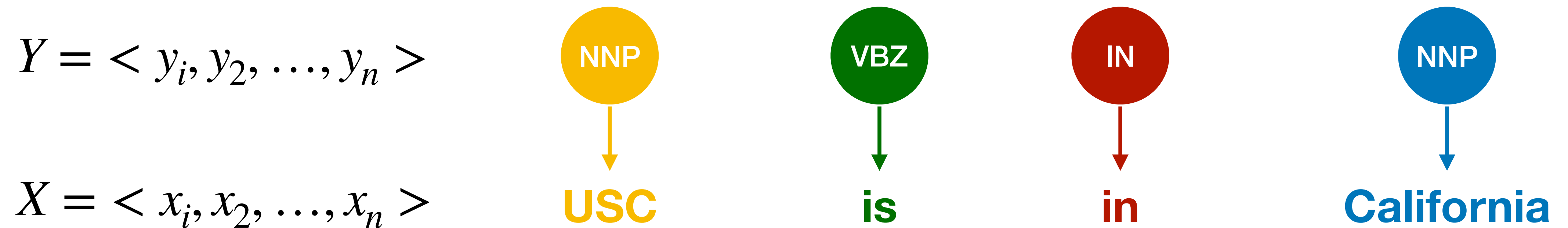
# Course Organization

## NLP Tasks



# Recap: What is Sequence Labeling?

**A type of structured prediction tasks**



Assigning each token of  $X$ , e.g.  $x_i$  a corresponding label  $y_i$

# Recap: Classical Models for Sequence Labeling

- **Structured Models**
  - HMMs
  - MEMMs
  - CRFs
- **Decoding Algorithms**
  - Greedy decoding
  - Viterbi decoding

**Motivation:** modeling dependencies between multiple labels/tags

# An Essential Question

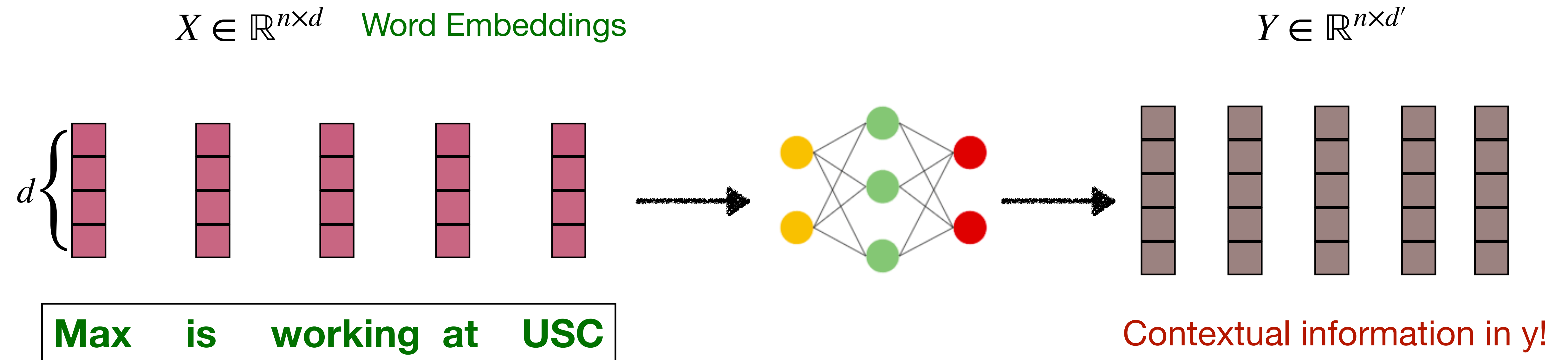
- Do we need structured models if the feature representations of the input sentence is perfect



$$P(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{j=1}^n P(y_j | x_1, \dots, x_n) \quad ?$$

# Our Goal: Sentence Representations

- One feature vector for each word



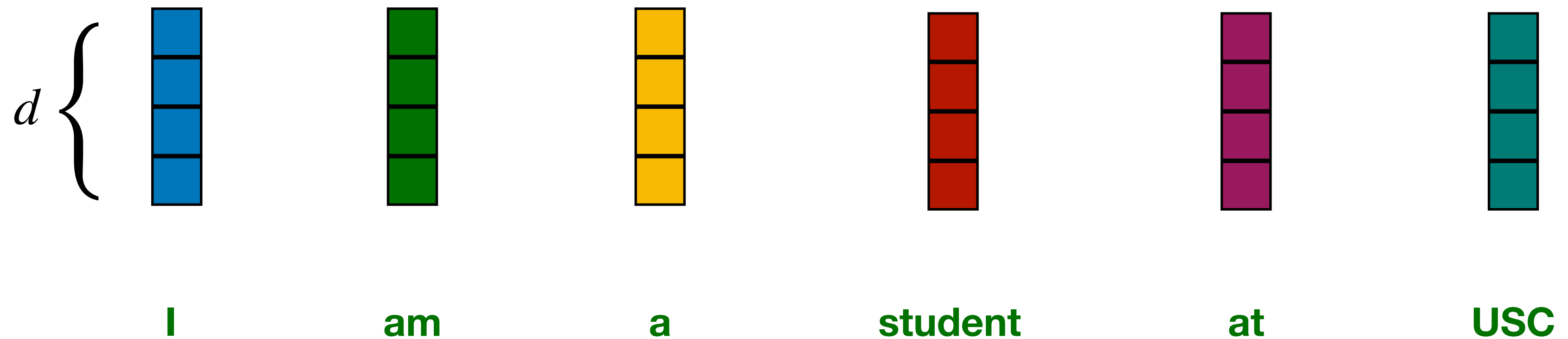
$$P(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{j=1}^n P(y_j | x_1, \dots, x_n)$$

# Recap: One Vector to Represent a Document

- Classification

- We need **a single feature vector** to feed into ML classifiers

$X \in \mathbb{R}^{n \times d}$  Word Embeddings

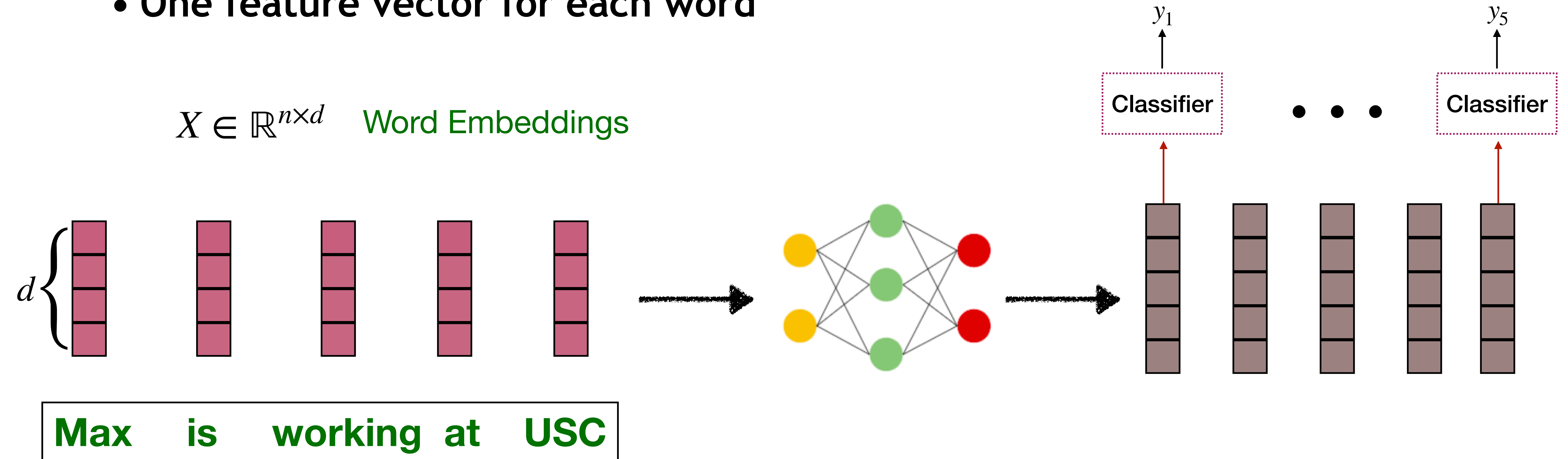


$$y = f(x_1, \dots, x_n) \in \mathbb{R}^{d'}$$



# Our Goal: Sentence Representations

- One feature vector for each word

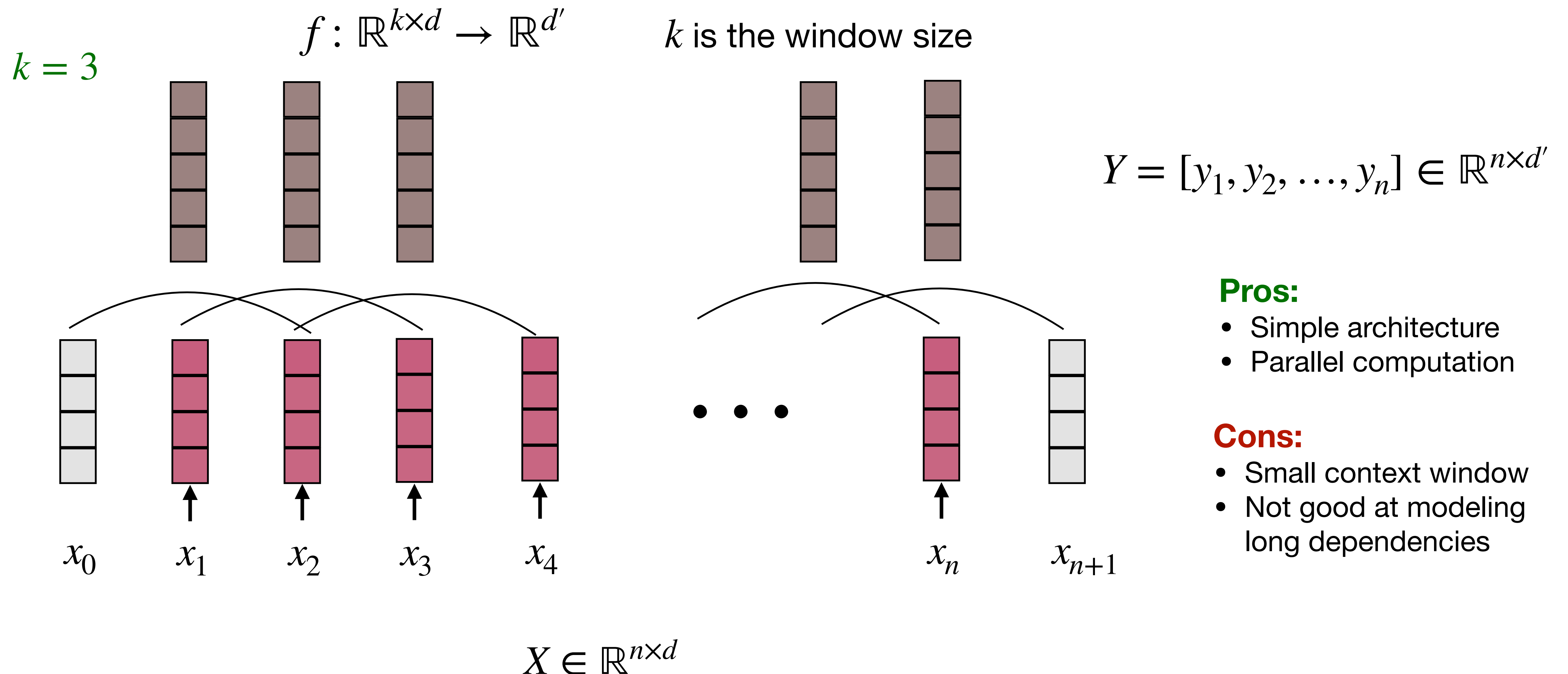


$$P(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{j=1}^n P(y_j | x_1, \dots, x_n)$$



# Convolutional Neural Networks (CNNs)

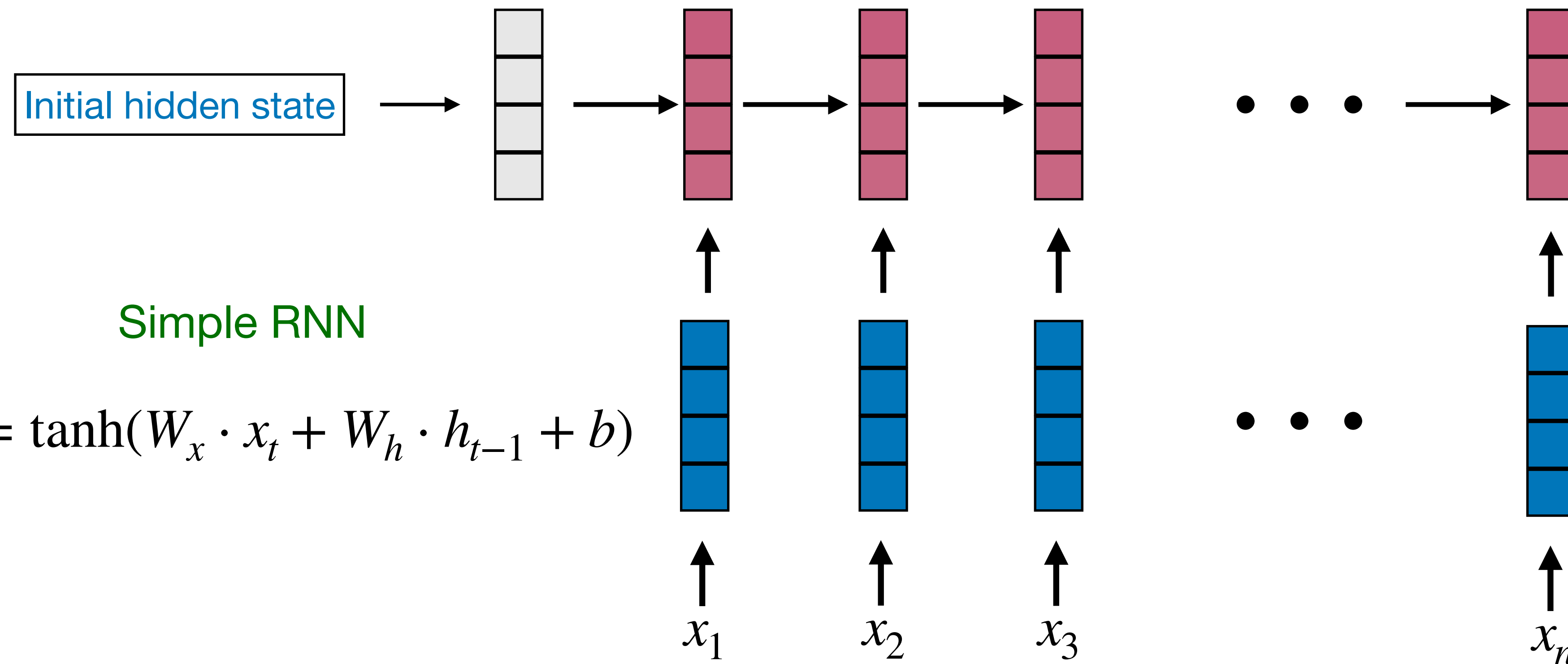
- **Basic Idea:** only model a segment of input with **fixed window-size**



# Recurrent Neural Networks

# Neural Networks for Sentence Representations

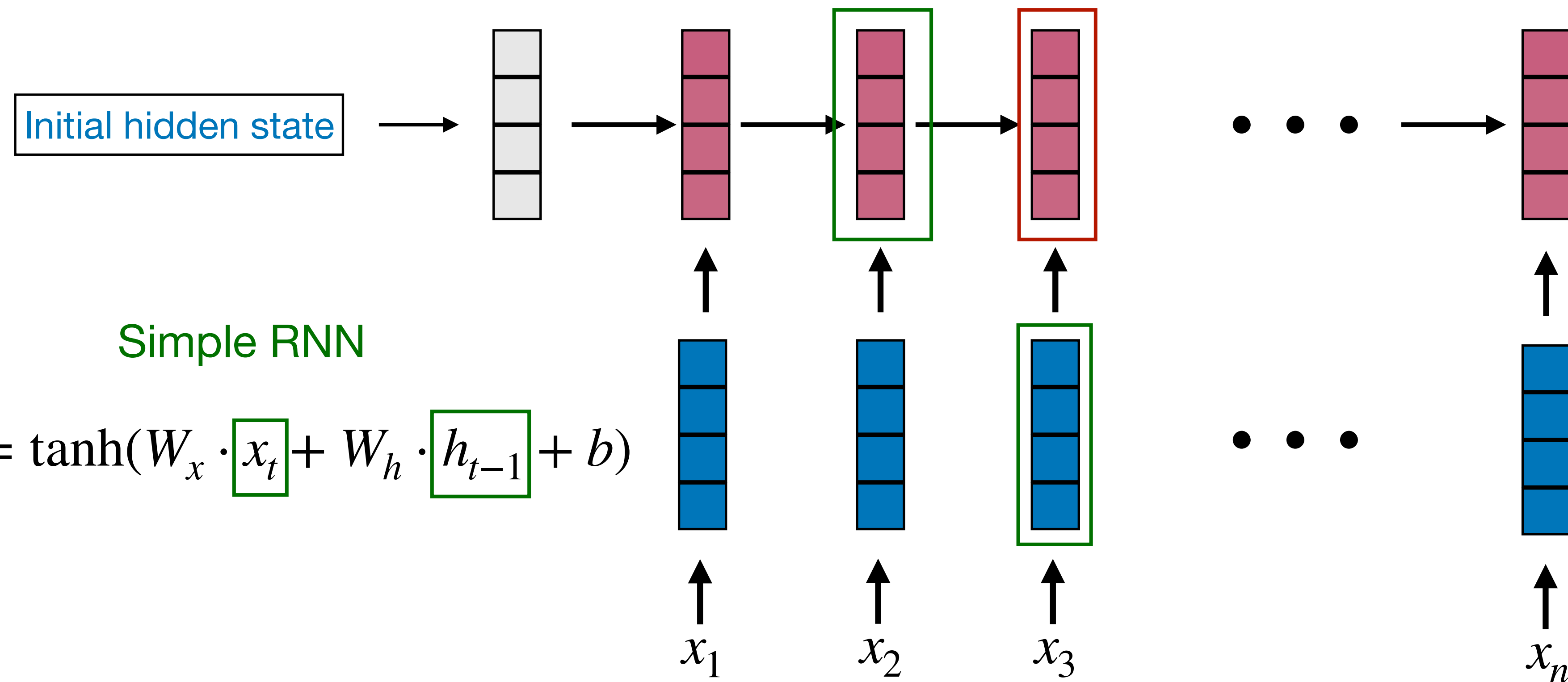
- Recurrent Neural Networks (RNNs)
  - Re-using one feed forward network in a recurrent way



# Neural Networks for Sentence Representations

- Recurrent Neural Networks (RNNs)
  - Re-using one feed forward network in a recurrent way

## Inherent Markov Property



# Problems of Simple RNN

- No future contexts

- Hard to train

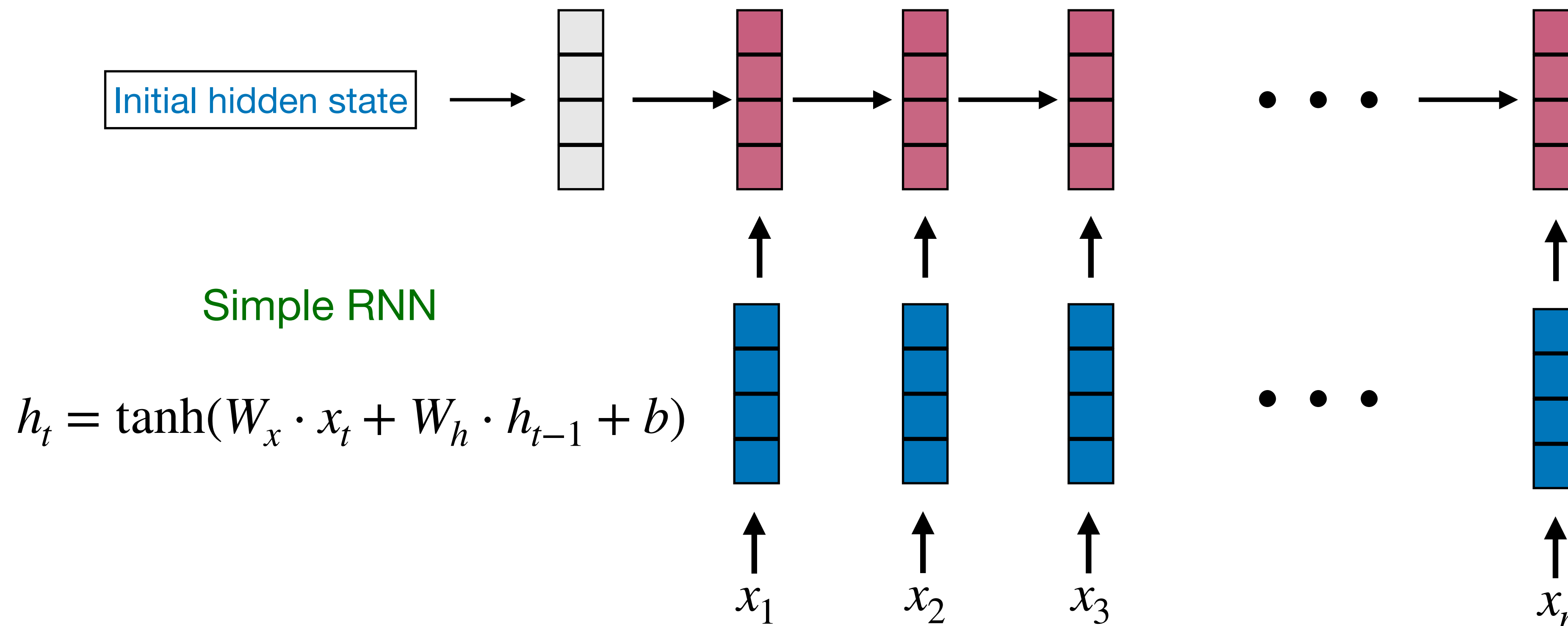
- Gradient vanishing/exploding

- Inefficient

- Sequential computation

- Limited size of hidden states

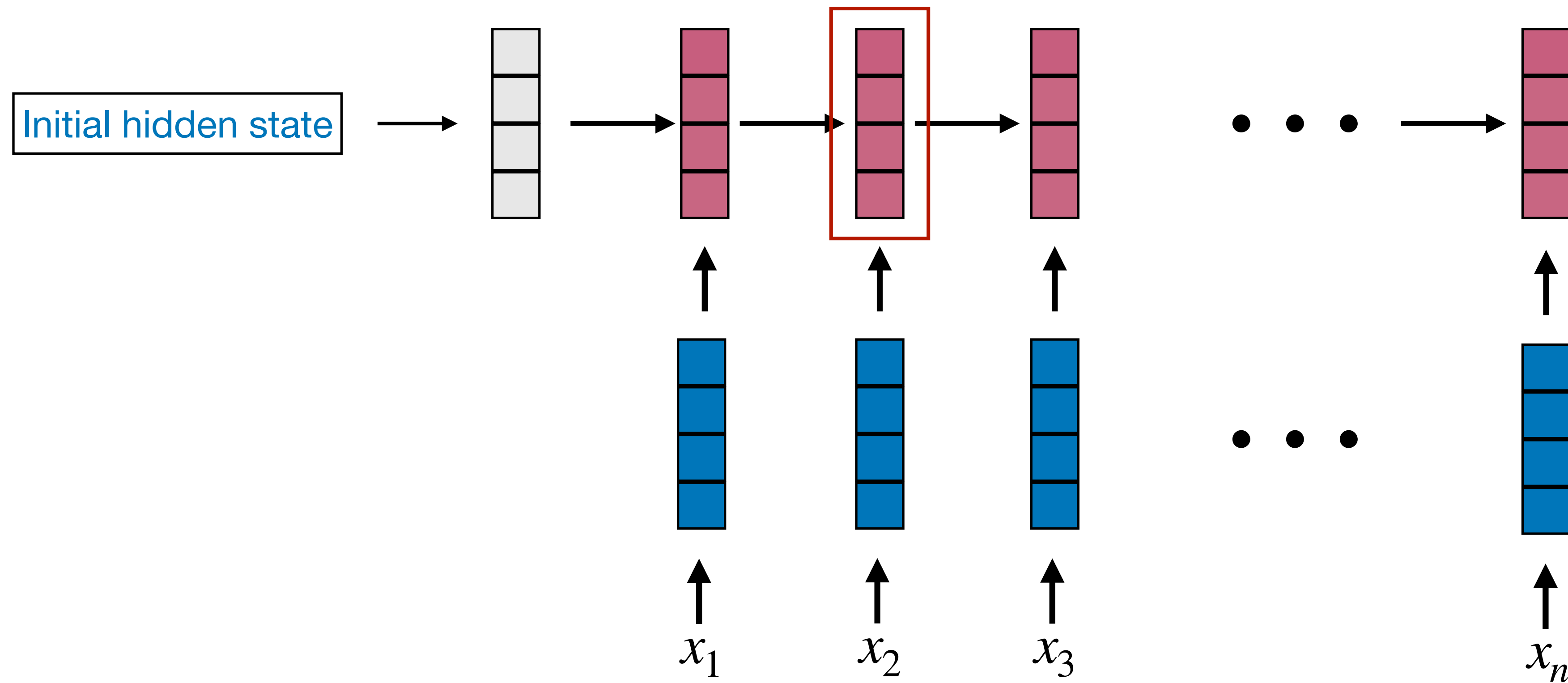
- Memory cost



# Bidirectional RNNs

- RNN **cannot** model **future** information

$h_2$  cannot access the information in  $x_3, \dots, x_n$



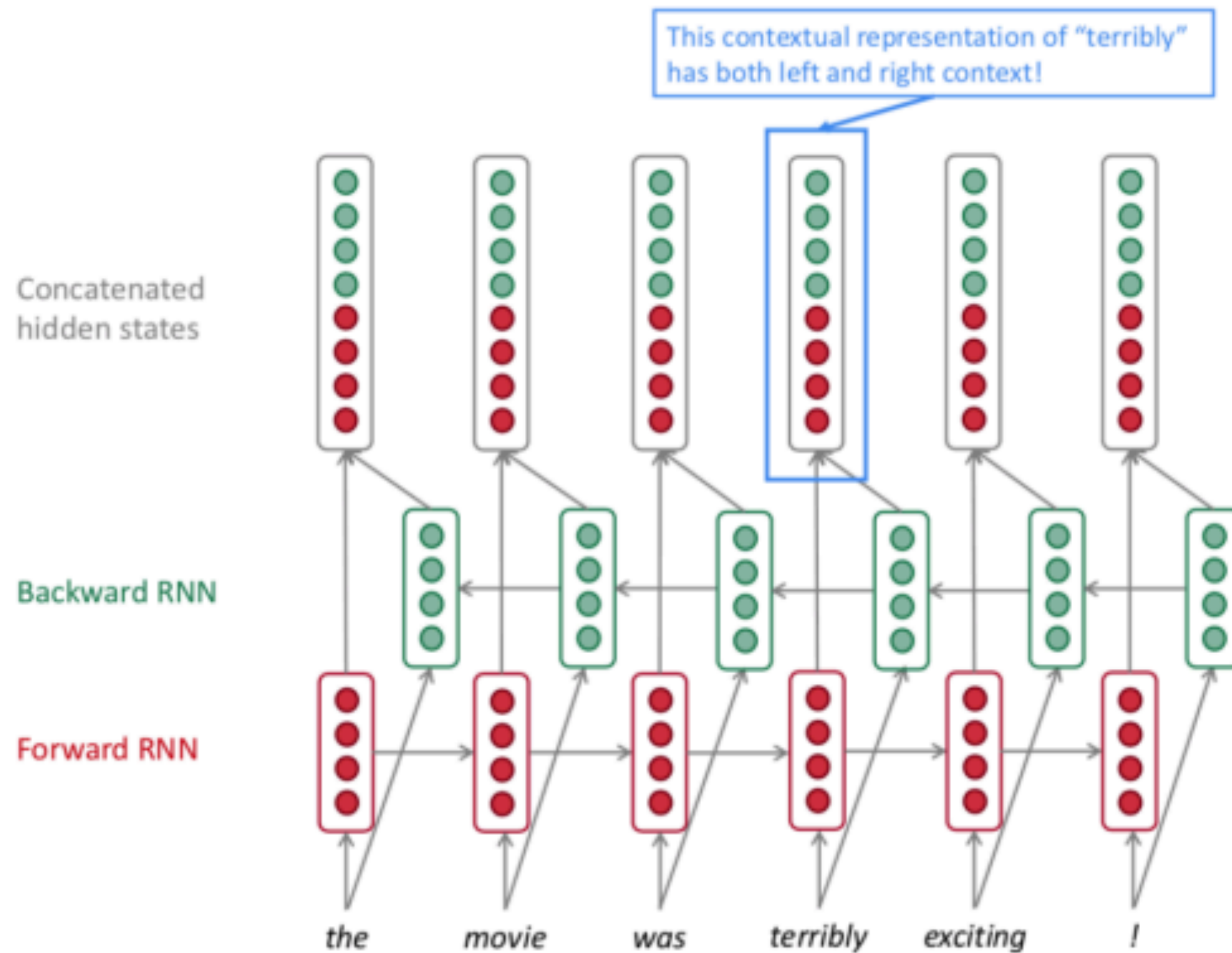
# Bidirectional RNNs

- Future information is important for sequence labeling tasks!

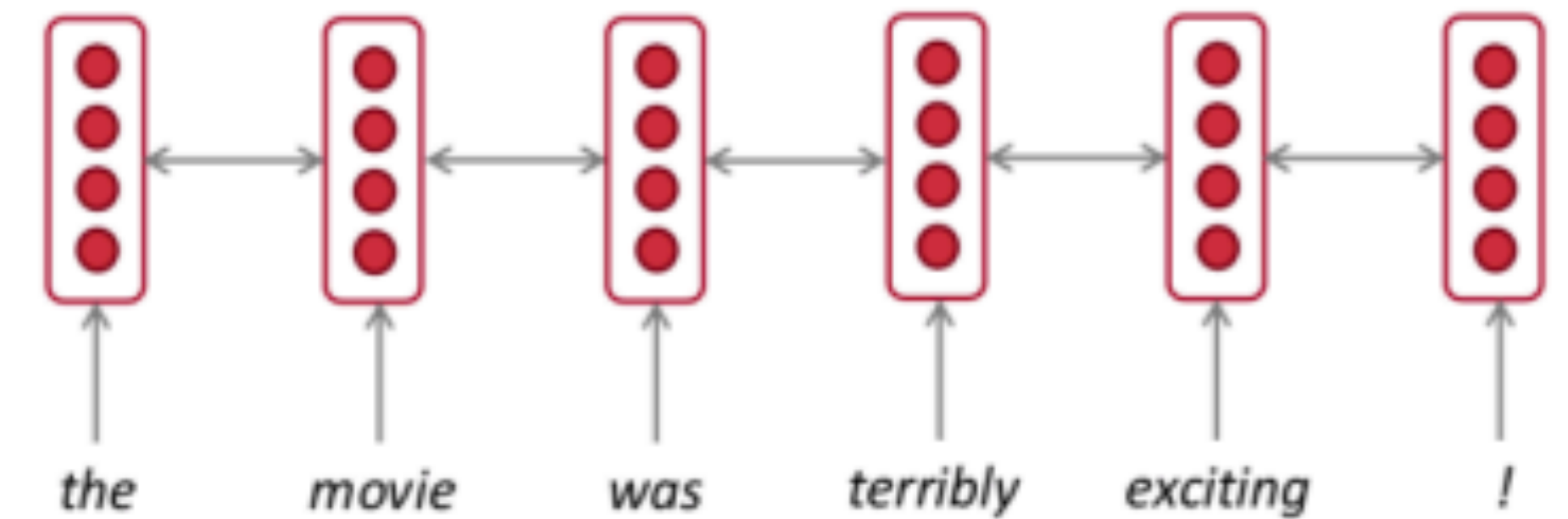


$$P(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{j=1}^n P(y_j | x_1, \dots, x_n)$$

# Bidirectional RNNs



=



$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t) \in \mathbb{R}^h$$

$$\vec{\mathbf{h}}_t = f_1(\vec{\mathbf{h}}_{t-1}, \mathbf{x}_t), t = 1, 2, \dots, n$$

$$\overleftarrow{\mathbf{h}}_t = f_2(\overleftarrow{\mathbf{h}}_{t+1}, \mathbf{x}_t), t = n, n-1, \dots, 1$$

$$\mathbf{h}_t = [\overleftarrow{\mathbf{h}}_t, \vec{\mathbf{h}}_t] \in \mathbb{R}^{2h}$$



# Bidirectional RNNs

- Bidirectional RNNs are only applicable if we have access to the **entire input sequence**
- If we do have entire input sequence, bidirectionality is powerful (and should be the **default choice**)
- A very common choice for sentence/document encoding: multi-layer bidirectional RNNs

# Problems of Simple RNN

- No future contexts

- Hard to train

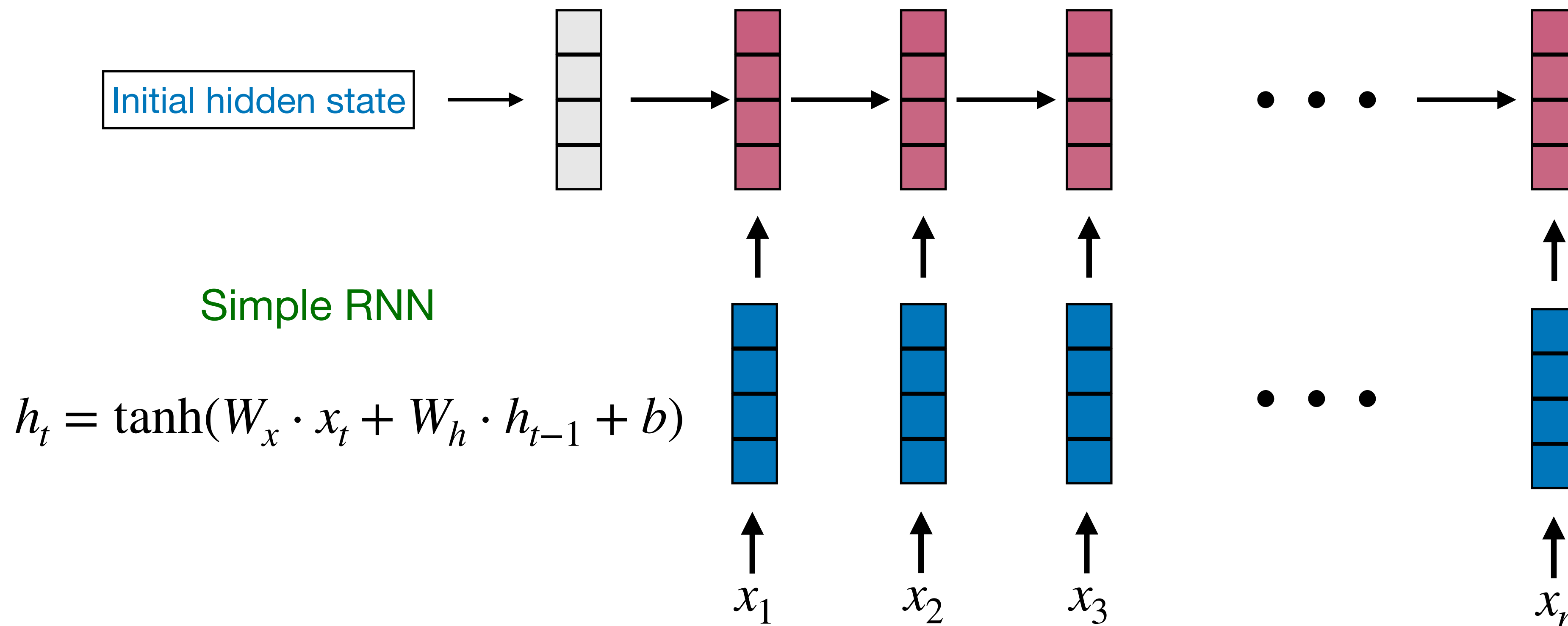
- Gradient vanishing/exploding

- Inefficient

- Sequential computation

- Limited size of hidden states

- Memory cost



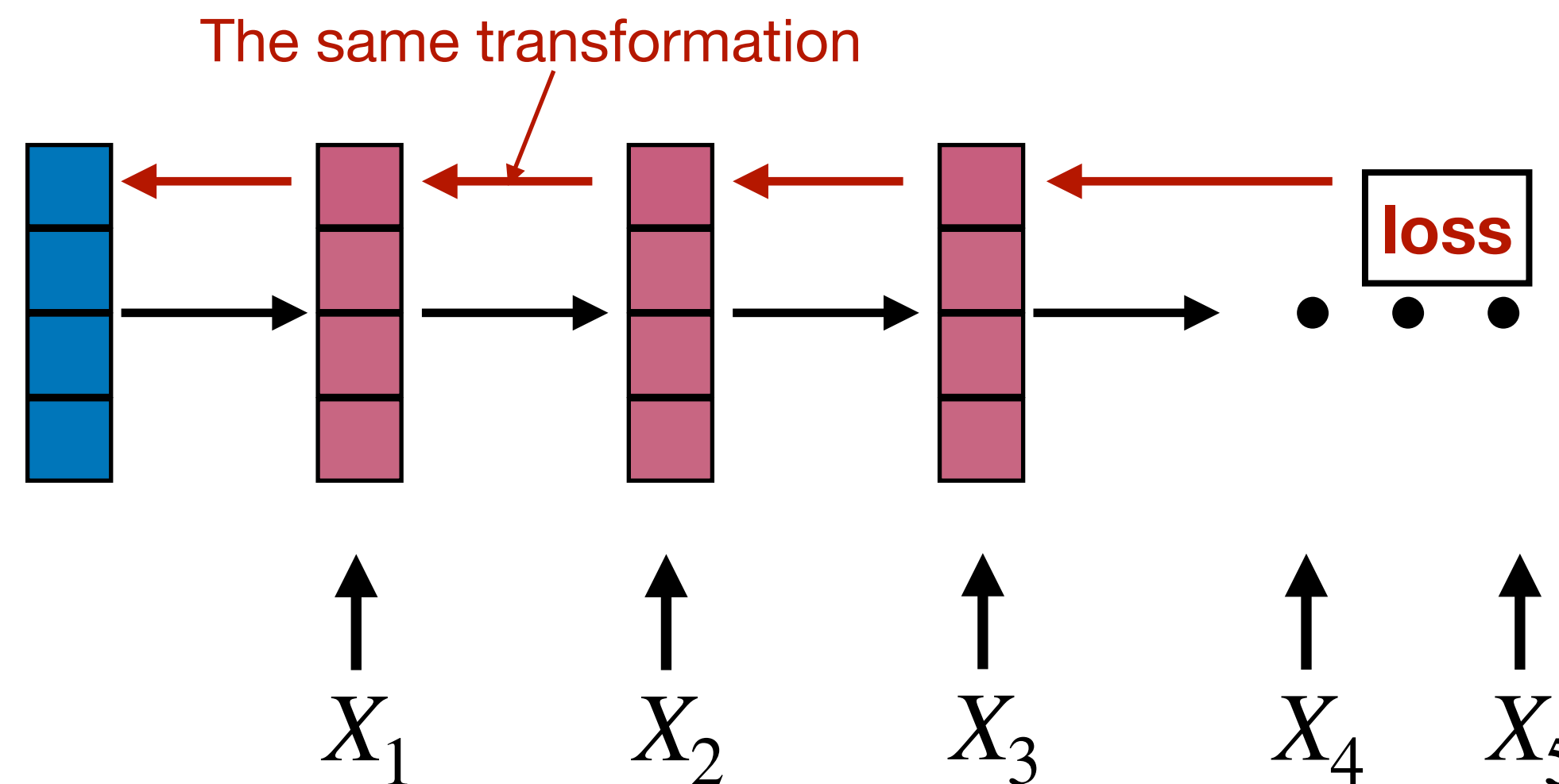
# Simple RNN is hard to train

- Hard to capture long-distance information: **vanishing/exploding gradients**

$$h_t = \tanh(W_x \cdot x_t + W_h \cdot h_{t-1} + b)$$

$$\left\| \frac{\partial h_t}{\partial h_{t-1}} \right\| \sim \|W_h\|$$

$$\left\| \frac{\partial h_t}{\partial h_{t-m}} \right\| \sim \|W_h\|^m$$



Why is this not a serious problem for multi-layer FFN/CNN?

# Advanced RNN Variants

- Long-short Term Memory (LSTMs)
- Gated Recurrent Units (GRUs)

LSTMs

$$\mathbf{i}_t = \sigma(\mathbf{W}^i \mathbf{h}_{t-1} + \mathbf{U}^i \mathbf{x}_t + \mathbf{b}^i)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}^f \mathbf{h}_{t-1} + \mathbf{U}^f \mathbf{x}_t + \mathbf{b}^f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}^o \mathbf{h}_{t-1} + \mathbf{U}^o \mathbf{x}_t + \mathbf{b}^o)$$

$$\mathbf{g}_t = \tanh(\mathbf{W}^g \mathbf{h}_{t-1} + \mathbf{U}^g \mathbf{x}_t + \mathbf{b}^g)$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t + \mathbf{g}_t \odot \mathbf{i}_t$$

$$\mathbf{h}_t = \tanh(\mathbf{c}_t) \odot \mathbf{o}_t$$

GRUs

$$\mathbf{r}_t = \sigma(\mathbf{W}^r \mathbf{h}_{t-1} + \mathbf{U}^r \mathbf{x}_t + \mathbf{b}^r)$$

$$\mathbf{z}_t = \sigma(\mathbf{W}^z \mathbf{h}_{t-1} + \mathbf{U}^z \mathbf{x}_t + \mathbf{b}^z)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{U} \mathbf{x}_t + \mathbf{b})$$

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

LSTM is more an art than a science

# LSTM: Long Short-Term Memory

$$\mathbf{i}_t = \sigma(\mathbf{W}^i \mathbf{h}_{t-1} + \mathbf{U}^i \mathbf{x}_t + \mathbf{b}^i)$$

Input gate:  $i_t \in (0,1)$

$$\mathbf{f}_t = \sigma(\mathbf{W}^f \mathbf{h}_{t-1} + \mathbf{U}^f \mathbf{x}_t + \mathbf{b}^f)$$

Forget gate:  $f_t \in (0,1)$

$$\mathbf{o}_t = \sigma(\mathbf{W}^o \mathbf{h}_{t-1} + \mathbf{U}^o \mathbf{x}_t + \mathbf{b}^o)$$

Output gate:  $o_t \in (0,1)$

$$\mathbf{g}_t = \tanh(\mathbf{W}^g \mathbf{h}_{t-1} + \mathbf{U}^g \mathbf{x}_t + \mathbf{b}^g)$$

Simple RNN

$$\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t + \mathbf{g}_t \odot \mathbf{i}_t$$

Cell state vector (internal memory)

$$\mathbf{h}_t = \tanh(\mathbf{c}_t) \odot \mathbf{o}_t$$

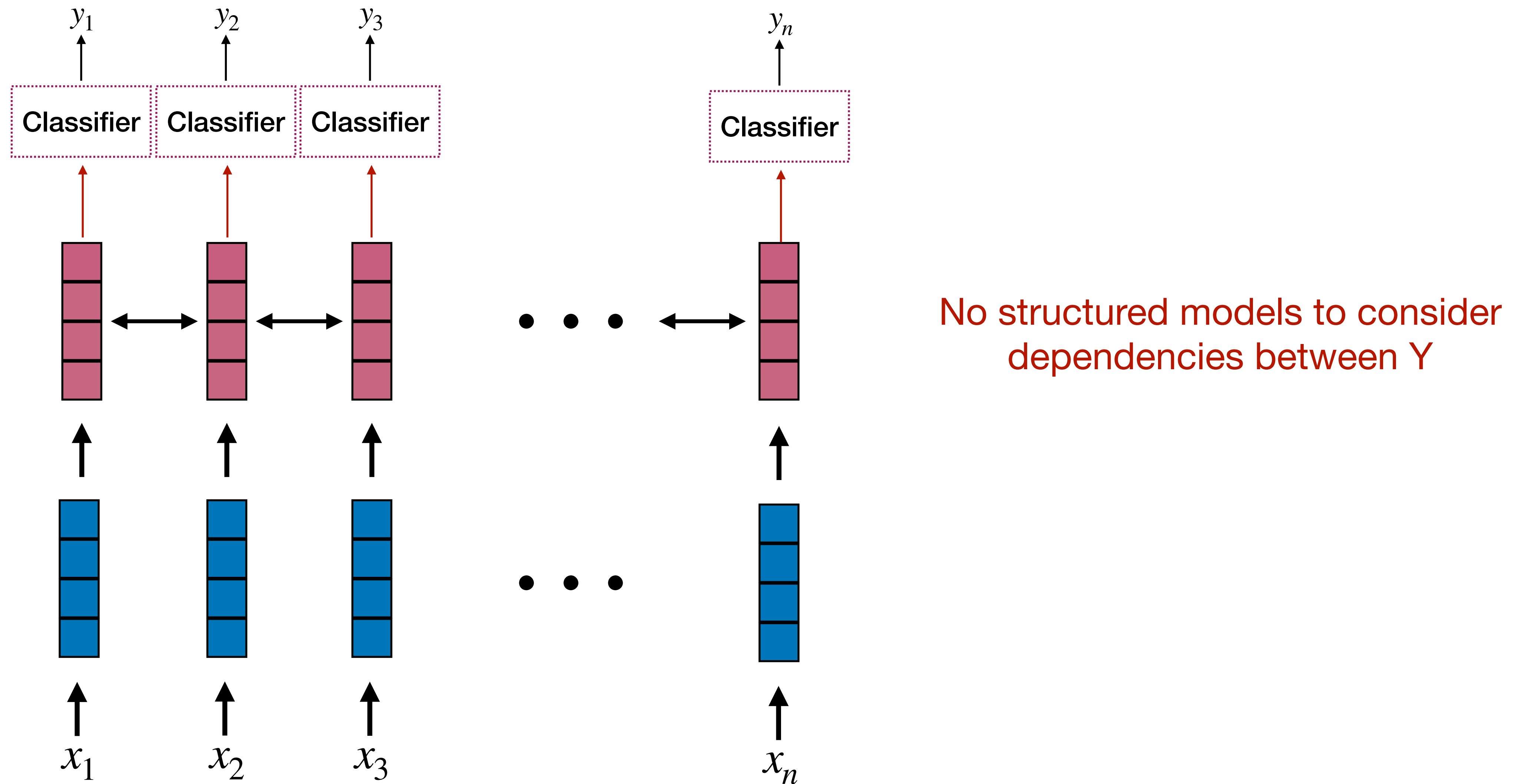
Final hidden vector

All the components are directly learned during training

# LSTMs for Sequence Labeling

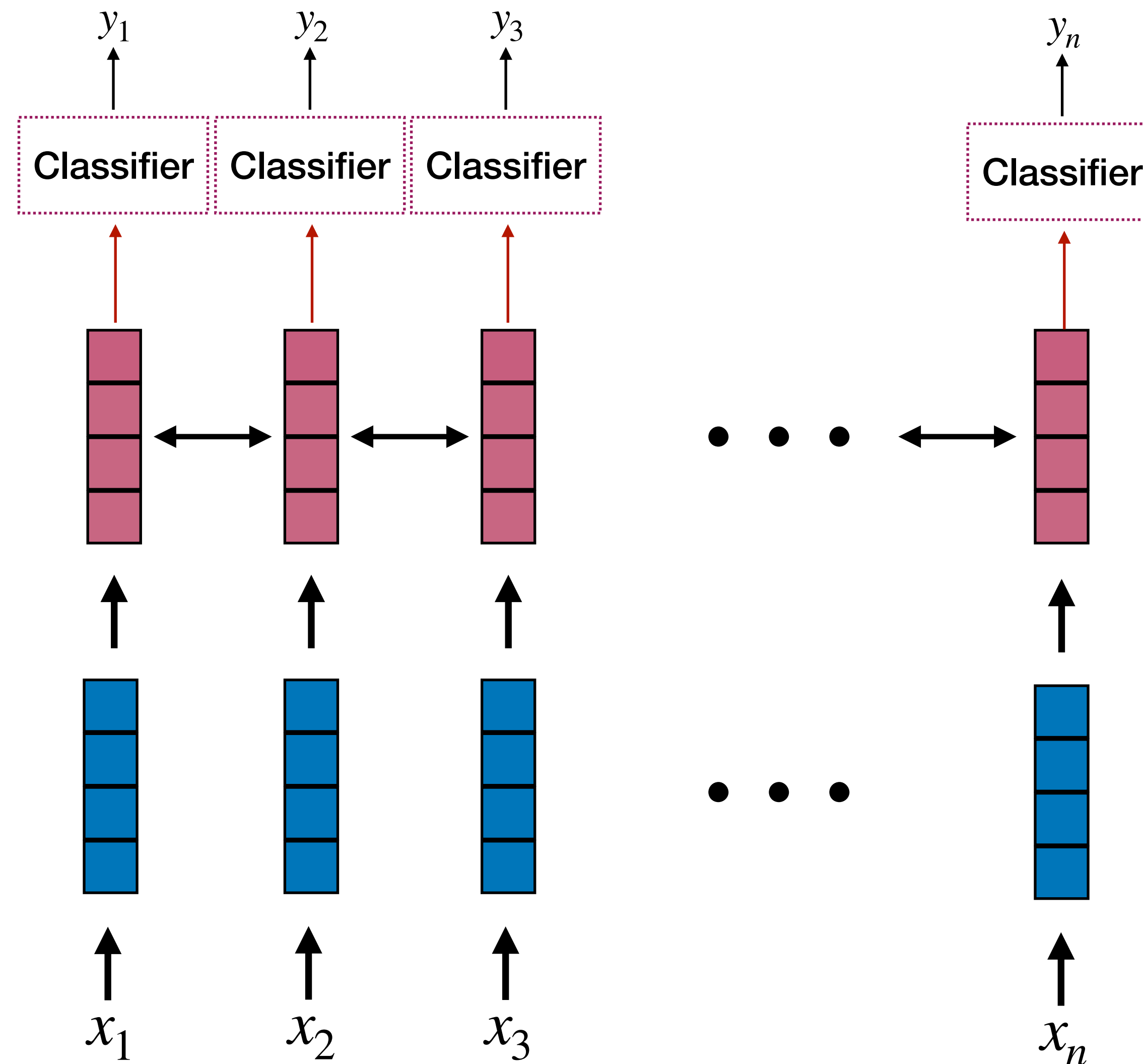
# LSTM for Sequence Labeling

- A simple bidirectional LSTM model



# Sequence Labeling

- How to initialize parameters?



**LSTM:**  $\text{Unif}(-\sqrt{1/d}, \sqrt{1/d})$

**Embedding:** Word2Vec or GloVe



# Sequence Labeling

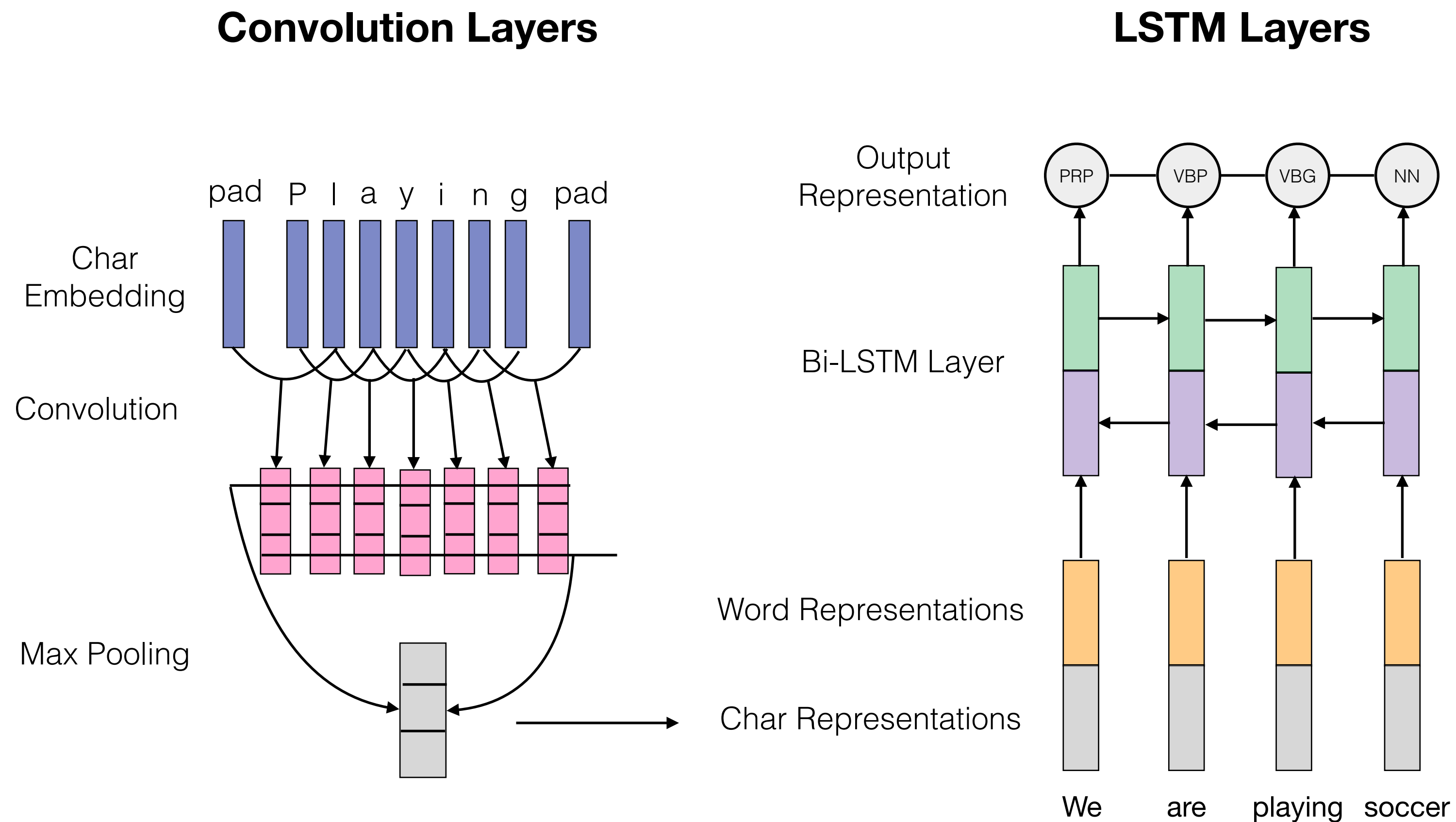
- A simple bidirectional LSTM model

	POS Tagging	NER
<b>HMM</b>	96.4%	75.3%
<b>CRF</b>	97.0%	88.7%
<b>CRF+external resources</b>	97.3%	91.2%
<b>BLSTM-Random</b>	96.1%	80.7%
<b>BLSTM-GloVe</b>	96.9%	87.0%

BLSTM is not good enough!

# Bidirectional LSTMs + CNNs

- Bidirectional LSTM only encodes word-level information
- Spelling features are important
  - Using CNN to model character-level information



# Sequence Labeling

- **BLSTM-CNN**

	POS Tagging	NER
<b>CRF</b>	97.0%	88.7%
<b>CRF+external resources</b>	97.3%	91.2%
<b>BLSTM</b>	96.9%	87.0%
<b>BLSTM-CNN</b>	97.3%	89.4%

**BLSTM-CNN is better than CRF!**

# An Essential Question

- Do we need structured models if the feature representations of the input sentence is perfect



$$P(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{j=1}^n P(y_j | x_1, \dots, x_n) \quad ?$$

BiLSTM + Char-level CNN is not strong enough...

Structured models are still useful?

Combining structured models with LSTM?

# Recap: Log-Linear Models

- MEMMs

$$p(y_j | y_{j-1}, x_1, \dots, x_m) = \frac{\exp(v \cdot f(x_1, \dots, x_m, j, y_{j-1}, y_j))}{\sum_{y'_j \in \mathbb{Y}} \exp(v \cdot f(x_1, \dots, x_m, j, y_{j-1}, y'_j))}$$

- CRFs

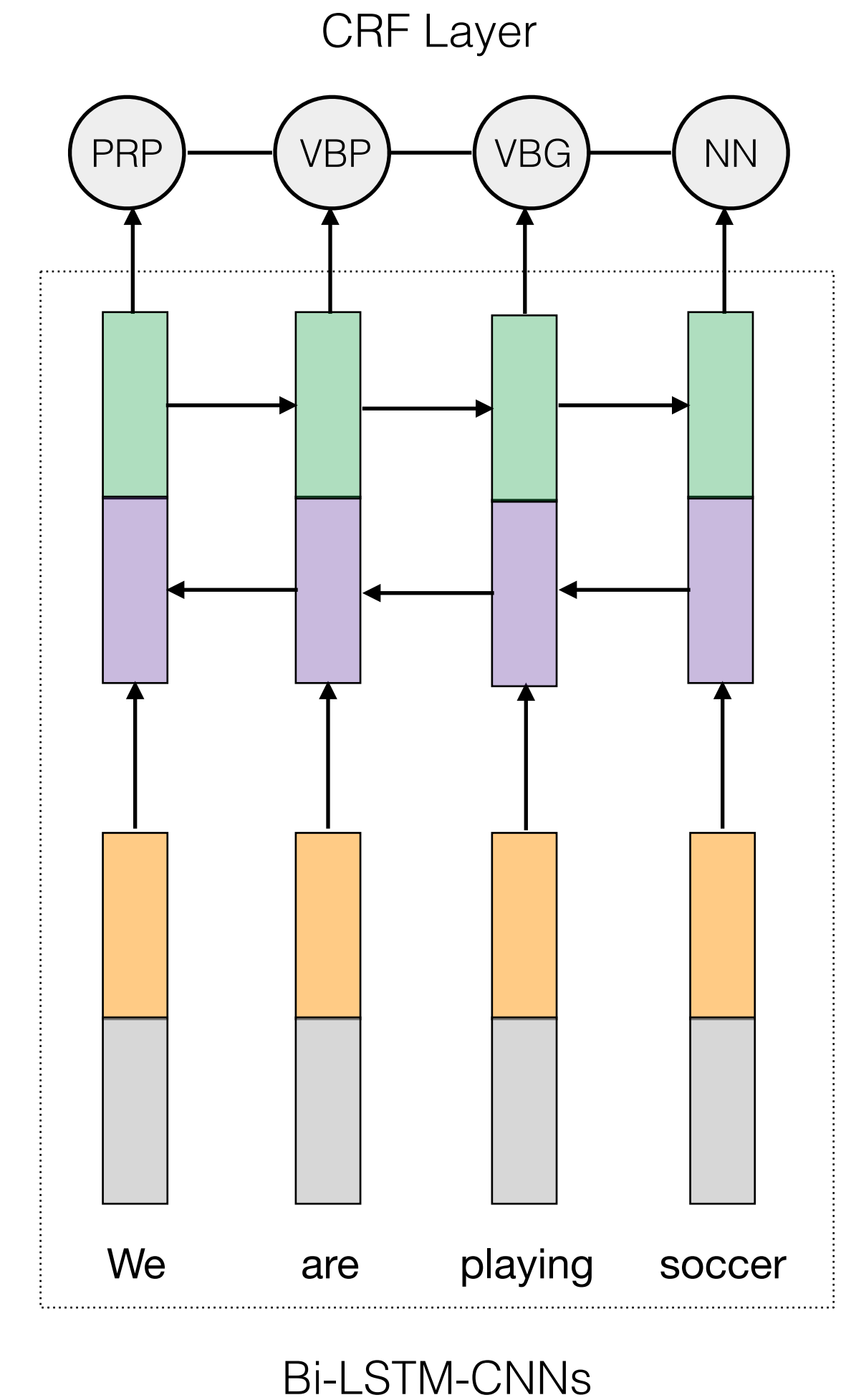
$$p(y_1, \dots, y_m | x_1, \dots, x_m) = \frac{\prod_{j=1}^m \exp(v \cdot f(x_1, \dots, x_m, j, y_{j-1}, y_j))}{\sum_{y'_1, \dots, y'_m \in \mathbb{Y}} \prod_{j=1}^m \exp(v \cdot f(x_1, \dots, x_m, j, y'_{j-1}, y'_j))}$$

# Sequence Labeling: BLSTM-CNNs-CRF

$$v \cdot f(x_1, \dots, x_m, j, y_{j-1}, y_j) = \mathbf{W}_{y_{j-1}, y_j}^T \mathbf{h}_j$$

Parameters w.r.t  $y_{j-1}, y_j$ 
LSTM hidden state at j

$$p(\mathbf{y} | \mathbf{x}; \mathbf{W}) = \frac{\prod_{j=1}^n \exp(\mathbf{W}_{y_{j-1}, y_j} \mathbf{h}_j)}{\sum_{\mathbf{y}' \in \mathcal{Y}} \prod_{j=1}^n \exp(\mathbf{W}_{y'_{j-1}, y'_j} \mathbf{h}_j)}$$



# Sequence Labeling

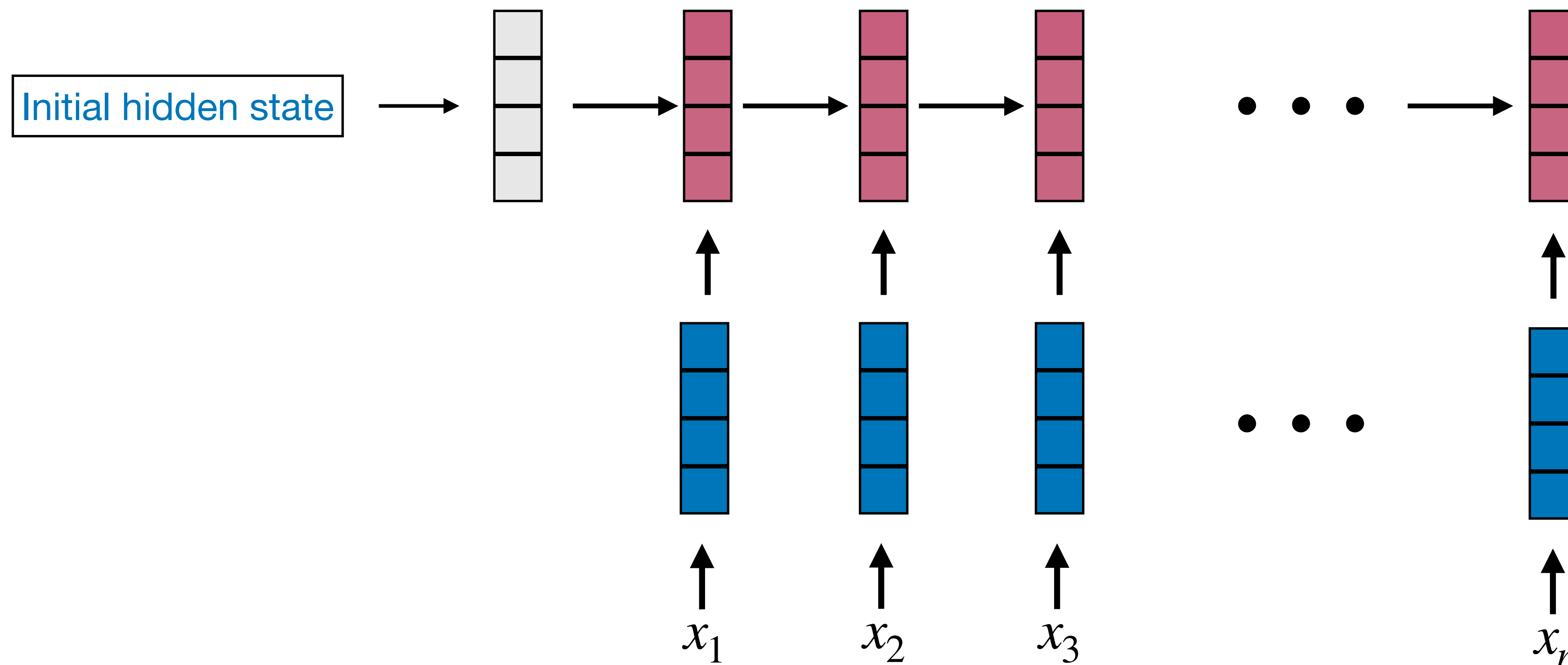
- **BLSTM-CNN-CRF**

	POS Tagging	NER
<b>CRF</b>	97.0%	88.7%
<b>CRF+external resources</b>	97.3%	91.2%
<b>BLSTM</b>	96.9%	87.0%
<b>BLSTM-CNN</b>	97.3%	89.4%
<b>BLSTM-CNN-CRF</b>	97.6%	91.2%

Structured models are still useful!

# Problems of Simple RNN

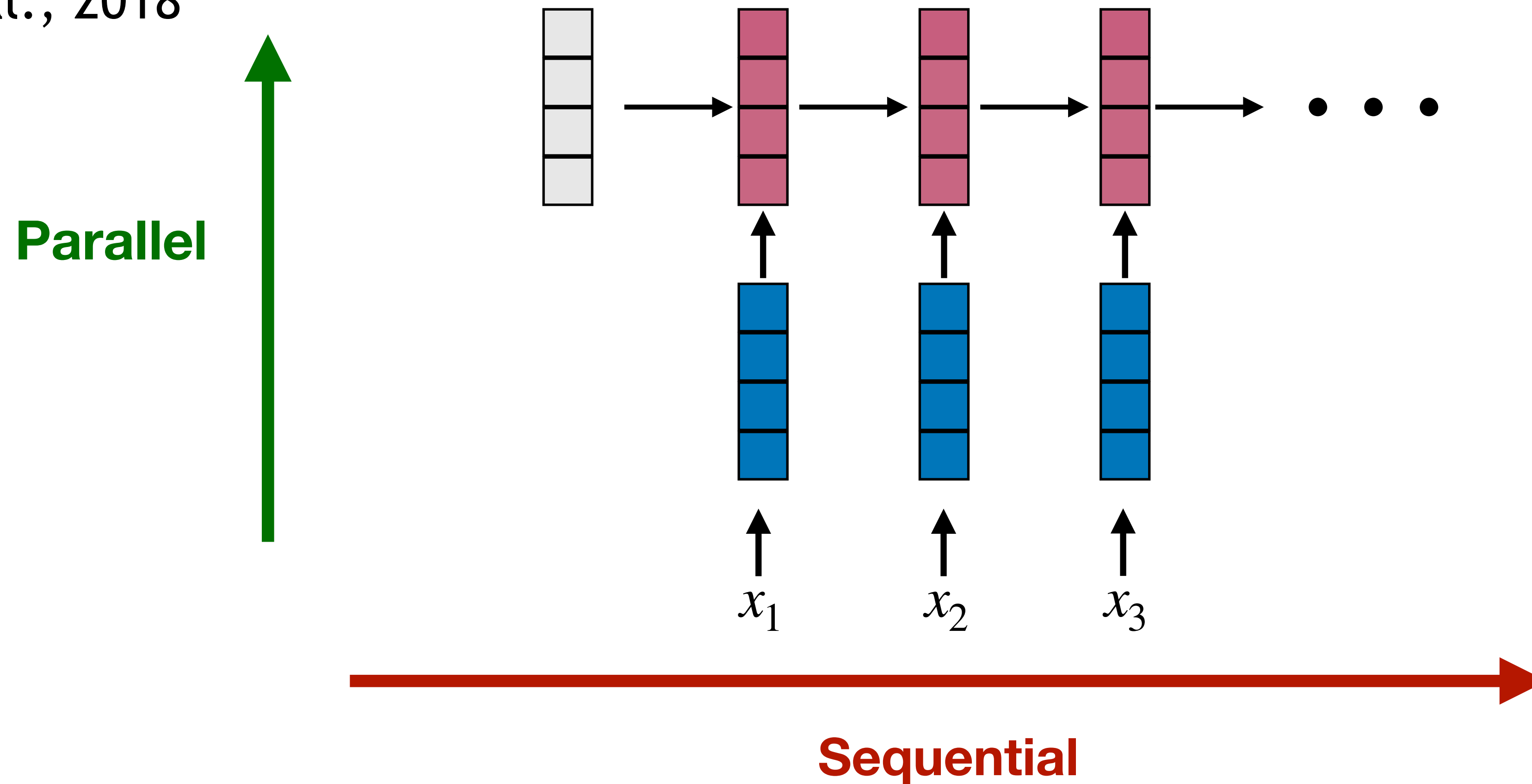
- No future contexts
- Hard to train
  - Gradient vanishing/exploding
- Inefficient
  - Sequential computation
- Limited size of hidden states
  - Memory cost





# Improving RNN efficiency

- Simple Recurrent Unit
  - Lei et al., 2018



$$h_t = \tanh(W_x \cdot x_t + \boxed{W_h \cdot h_{t-1}} + b)$$



$$V_h \odot h_{t-1}$$

# Problems of Simple RNN

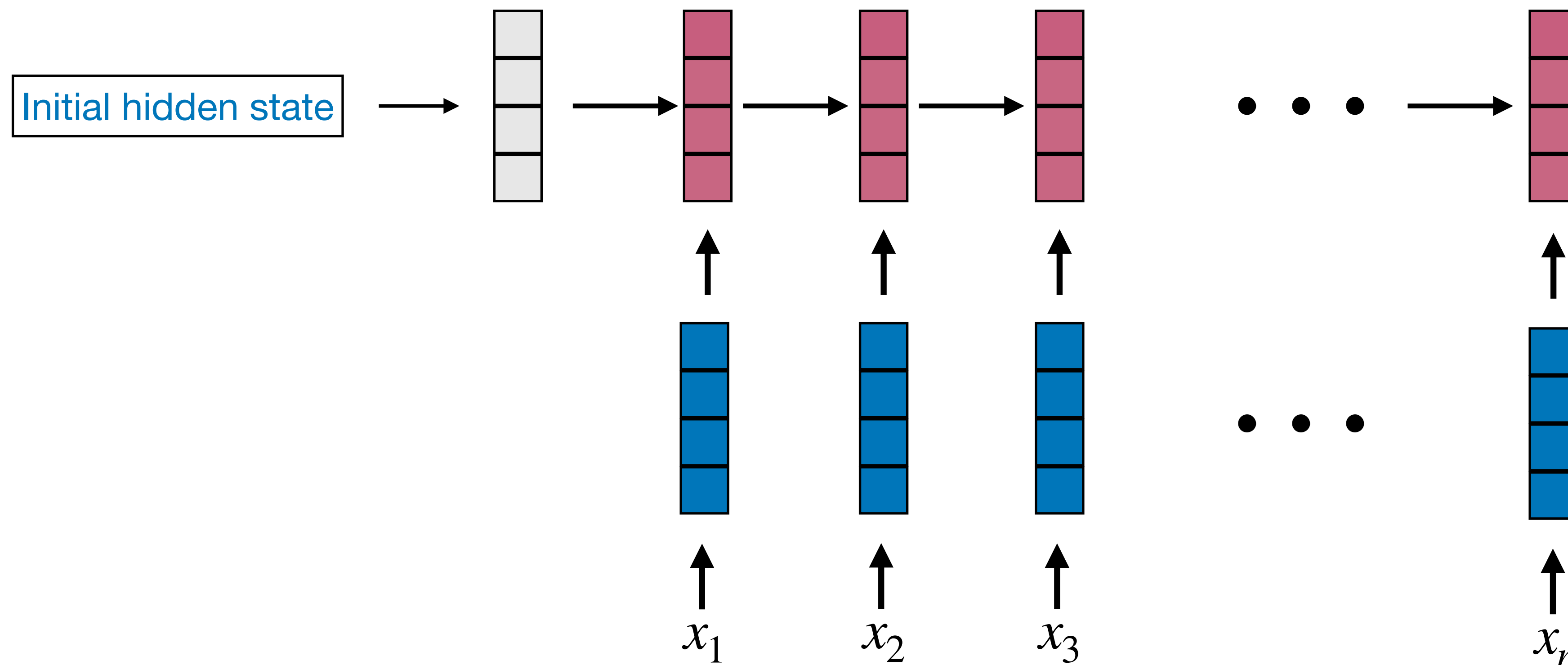
- No future contexts
- Hard to train
  - Gradient vanishing/exploding

- Inefficient

- Sequential computation

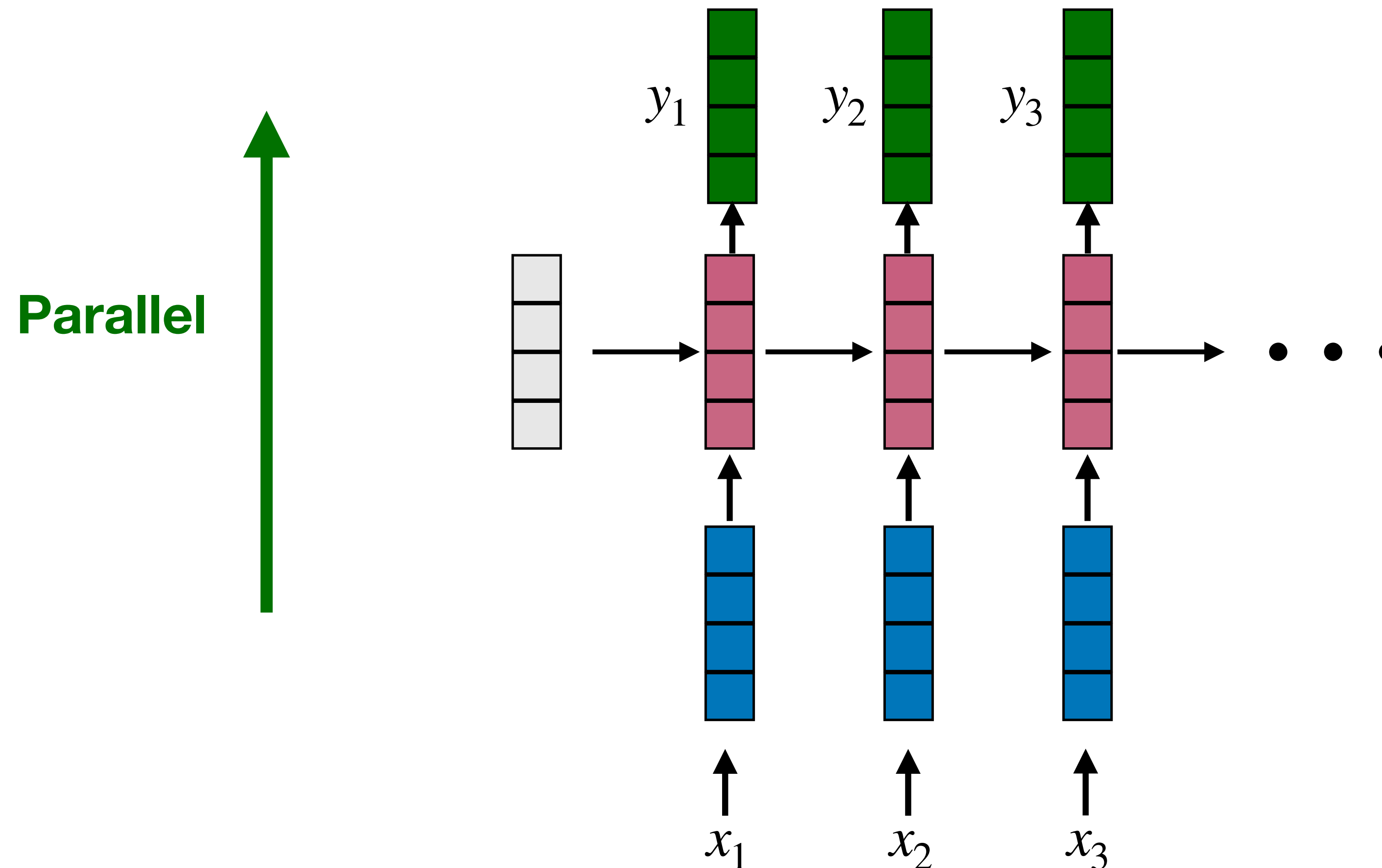
- Limited size of hidden states

- Memory cost



# Enlarging Hidden Size

- Linear Recurrent Neural Networks
  - Directly computing output  $y_t$  by skipping the explicit computation of  $h_t$



# Linear RNN

- Taking one dimension of  $x$  as an example, e.g.  $x_t \in \mathbb{R}$

- Expanding  $x_t$  to  $n$  dimensions

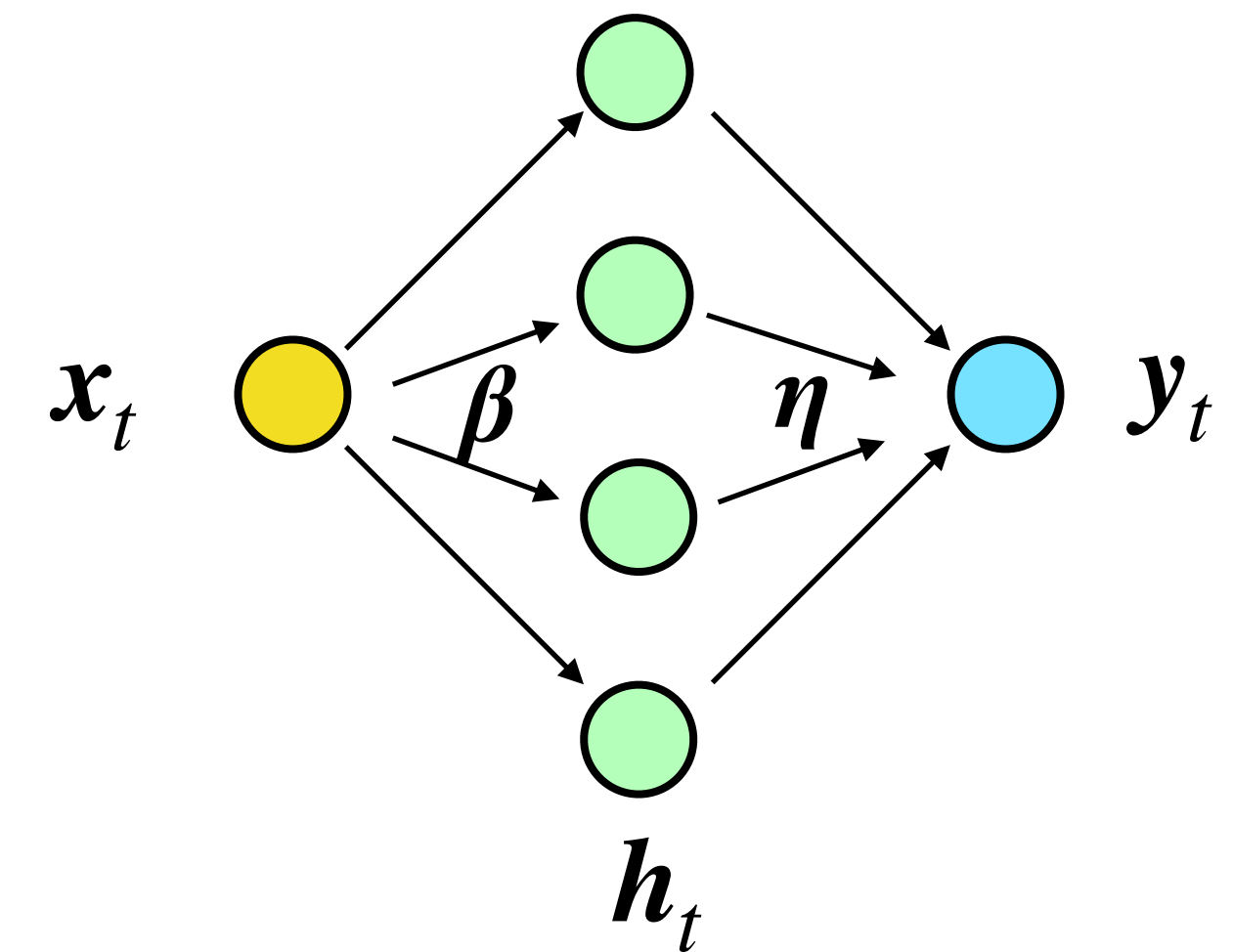
$$\mathbf{u}_t = \beta x_t \in \mathbb{R}^n$$

- Applying linear operation on hidden state

$$\mathbf{h}_t = A\mathbf{h}_{t-1} + \mathbf{u}_t = A\mathbf{h}_{t-1} + \beta x_t \in \mathbb{R}^n$$

- Mapping the  $h$ -dimensional vector back to 1 dimension

$$y_t = \eta^T \mathbf{h}_t \in \mathbb{R}$$



# Linear RNN

- Unrolling the computation of  $y_t$

$$\mathbf{h}_t = A\mathbf{h}_{t-1} + \mathbf{u}_t = A\mathbf{h}_{t-1} + \beta x_t \in \mathbb{R}^n$$

$$\mathbf{y}_t = \boldsymbol{\eta}^T \mathbf{h}_t \in \mathbb{R}$$

$$\mathbf{h}_1 = A\mathbf{h}_0 + \beta x_1$$

$$\mathbf{y}_1 = \boldsymbol{\eta}^T \mathbf{h}_1 = (\boldsymbol{\eta}^T A)\mathbf{h}_0 + (\boldsymbol{\eta}^T \beta)x_1$$

$$\mathbf{h}_2 = A\mathbf{h}_1 + \beta x_2 = A^2\mathbf{h}_0 + A\beta x_1 + \beta x_2$$

$$\begin{aligned}\mathbf{y}_2 &= \boldsymbol{\eta}^T \mathbf{h}_2 \\ &= (\boldsymbol{\eta}^T A^2)\mathbf{h}_0 + (\boldsymbol{\eta}^T A\beta)x_1 + (\boldsymbol{\eta}^T \beta)x_2\end{aligned}$$

•  
•  
•

# Reading Materials

- **Relavant Papers**
  - BLSTM-CNNs-CRF