# Day 6 - From Recurrent Nets to Transformers

Advanced Text as Data: Natural Language Processing
Essex Summer School in Social Science Data Analysis

Burt L. Monroe (Instructor) & Sam Bestvater (TA)
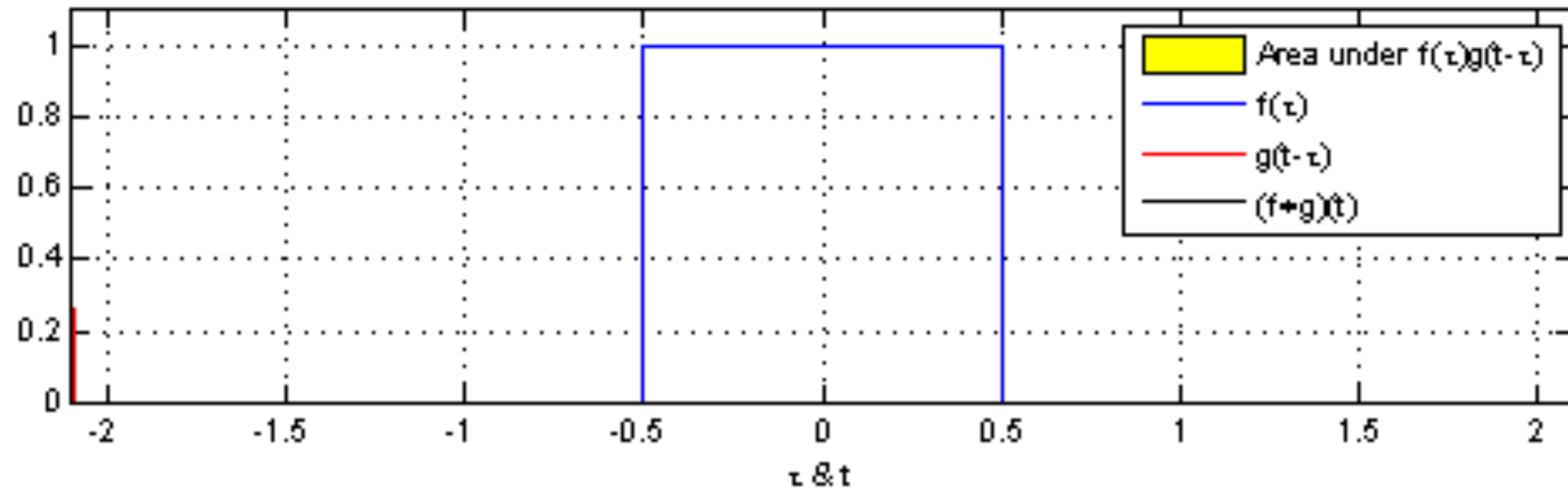Pennsylvania State University

August 3, 2021

# Today

- Convolutional Neural Nets (CNNs) - Convolution, filters/kernels, higher-level features

- Recurrent Neural Nets (RNNs) - Recurrence / sequence, encoder-decoder seq2seq

- Gating in recurrent networks (LSTMs / bi-LSTMs)

- Attention mechanism in seq2seq models

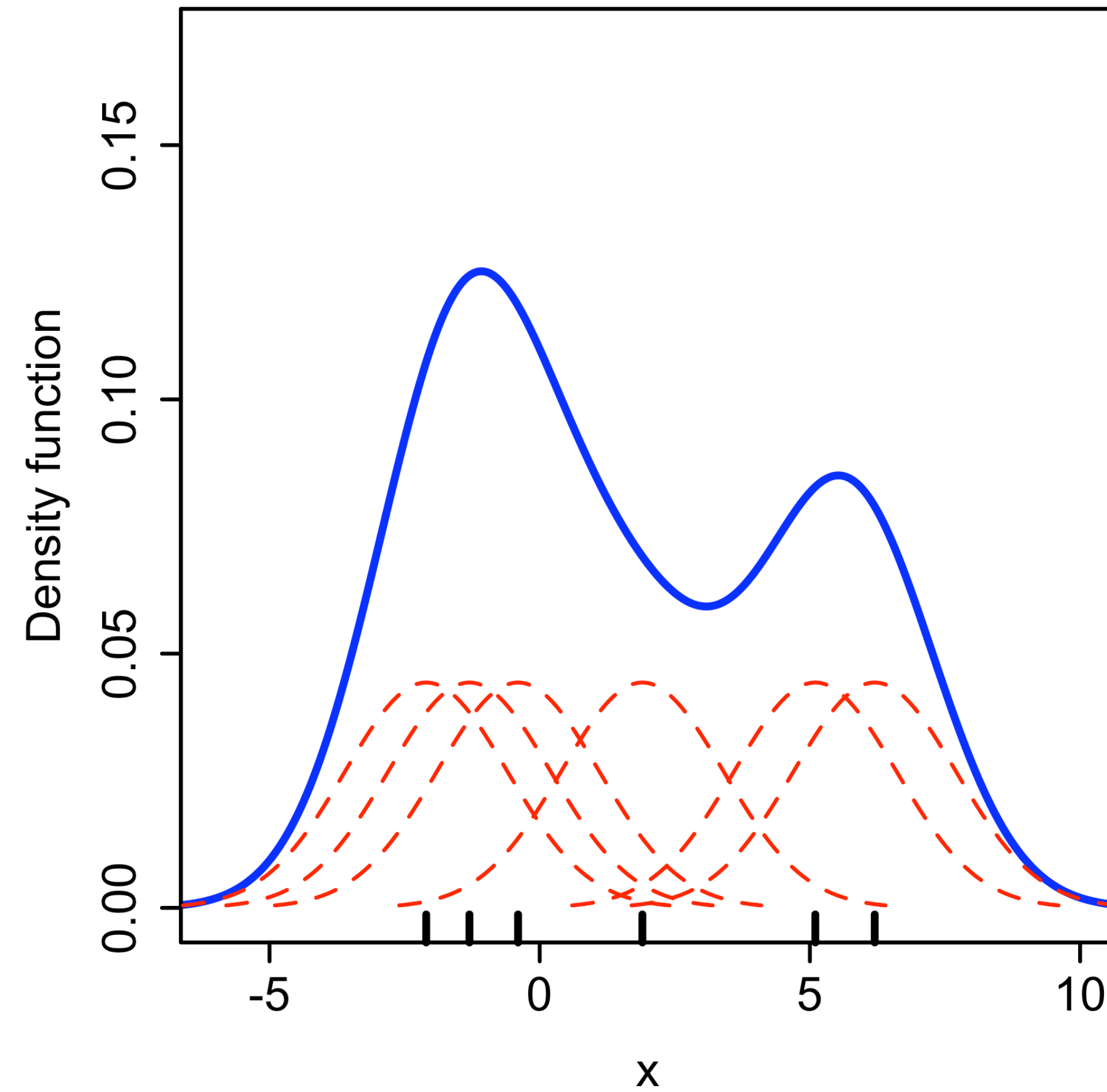- Self-attention & positional encodings (transformer)
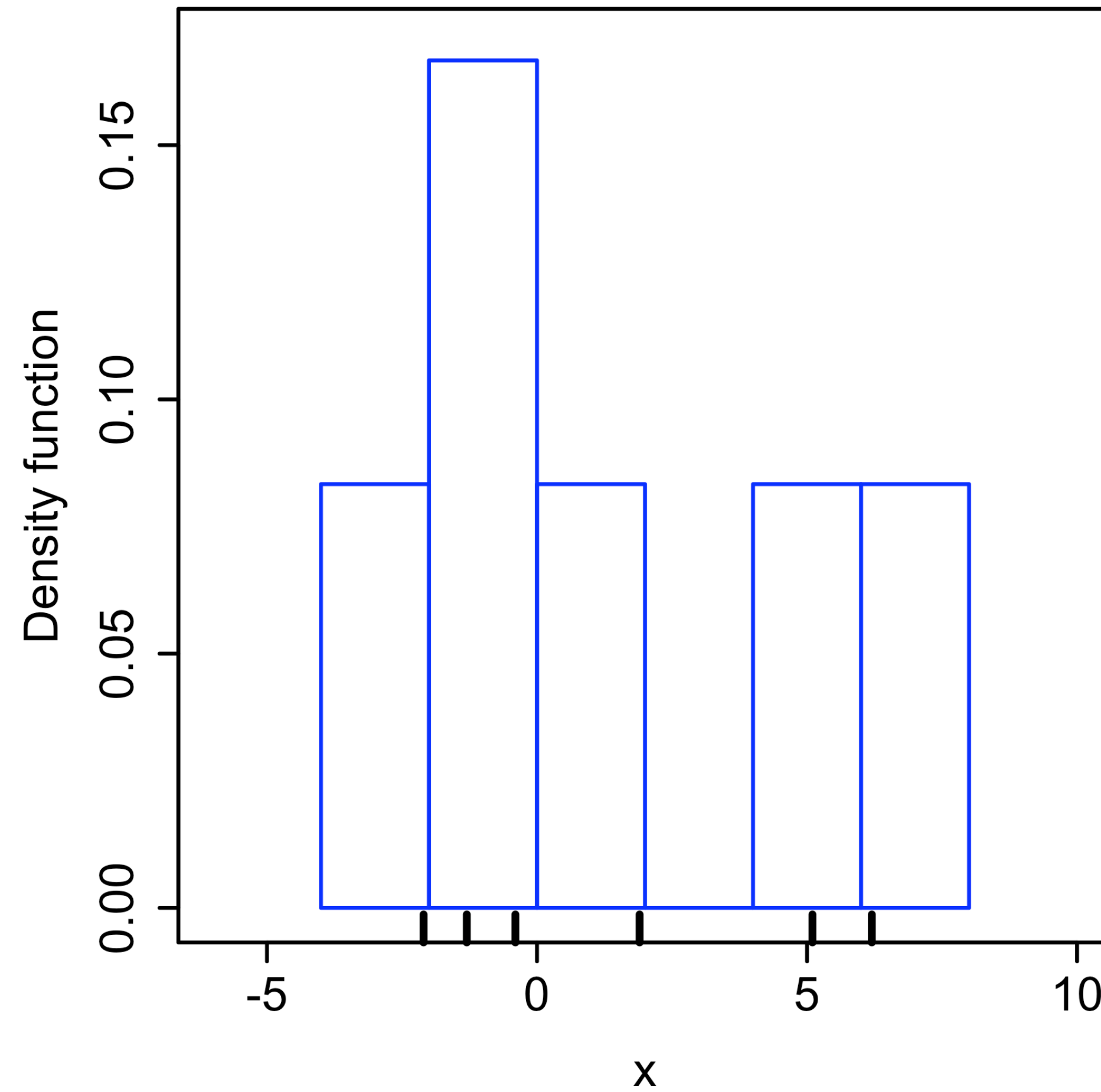
# Today (unlikely)

- Convolutional Neural Nets (CNNs) - Convolution, filters/kernels, higher-level features

- Recurrent Neural Nets (RNNs) - Recurrence / sequence, encoder-decoder seq2seq

- Gating in recurrent networks (LSTMs / bi-LSTMs)

- Attention mechanism in seq2seq models

- Self-attention & positional encodings (transformer)

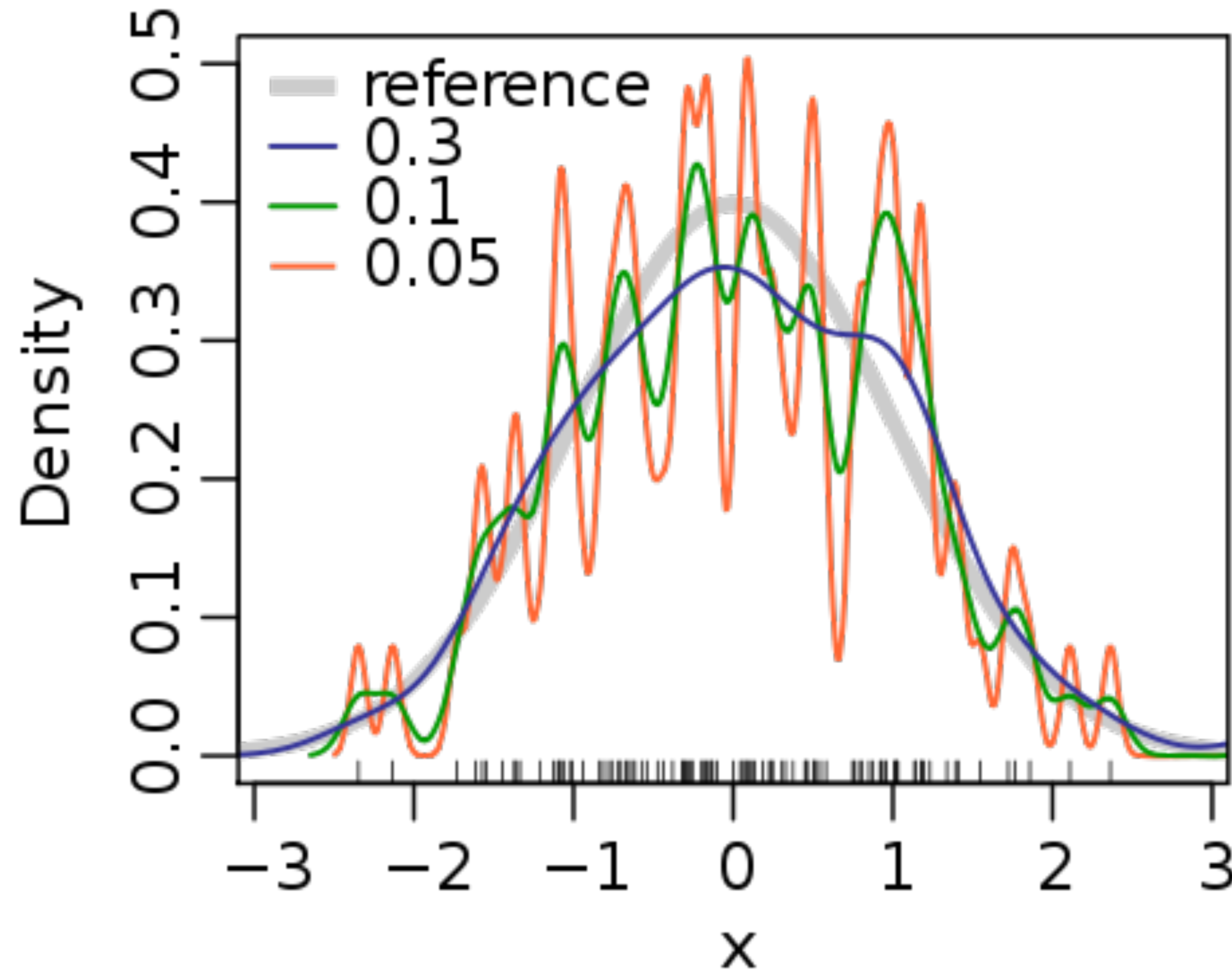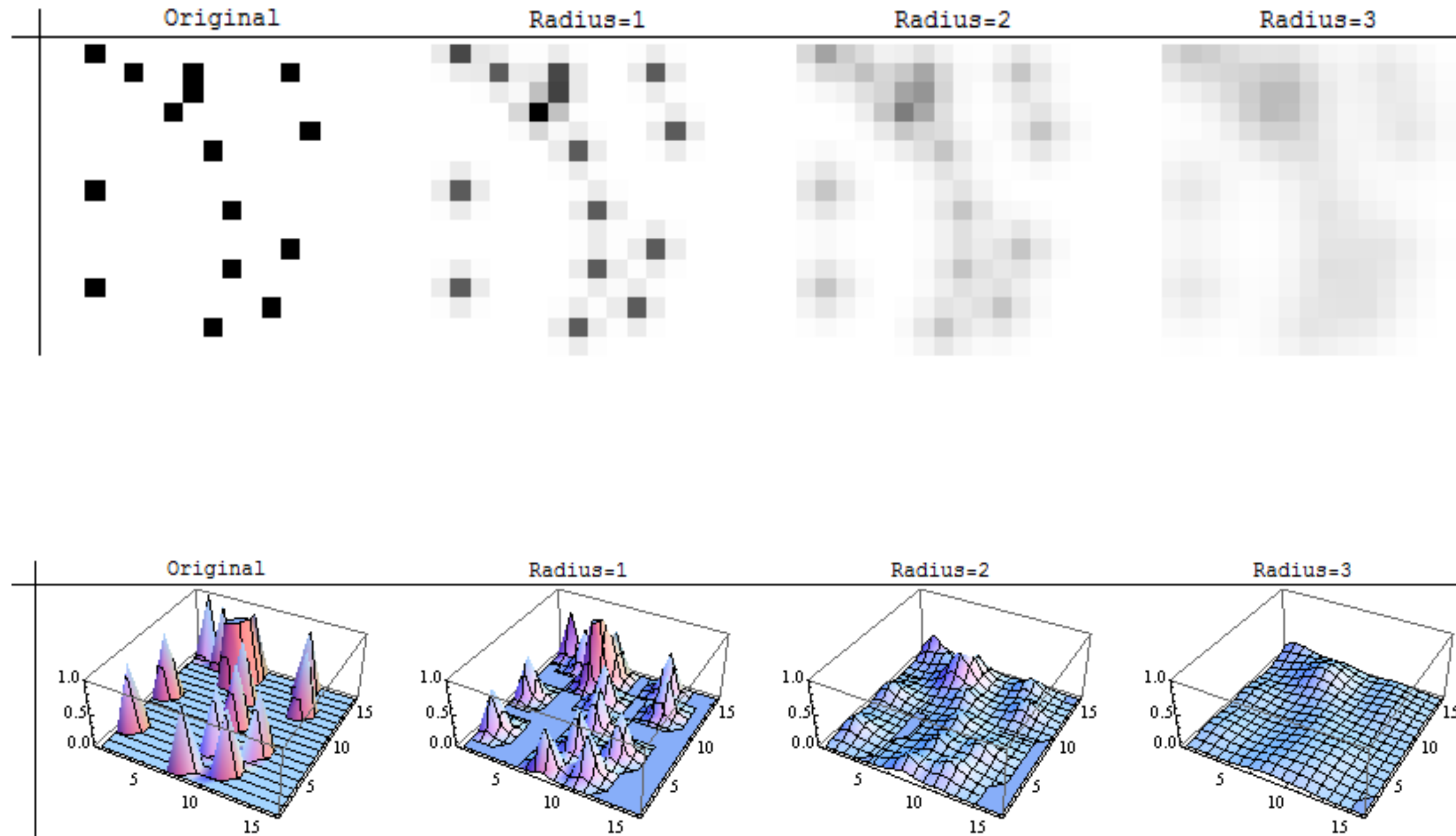# Convolution, Convolutional Neural Nets, and CNNs in NLP

# Convolution

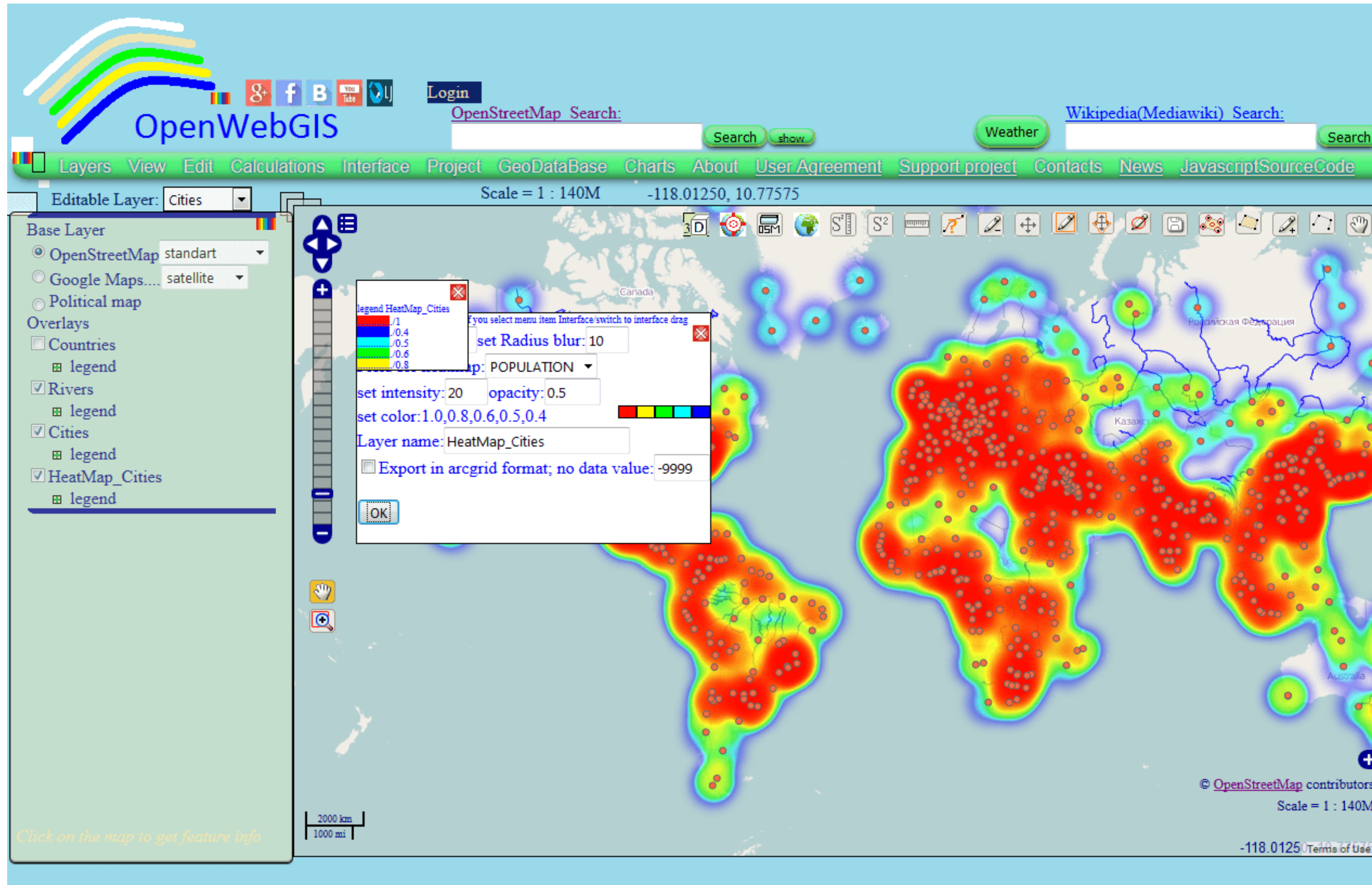# Kernel density — smooth histogram by convolving a Gaussian over observations

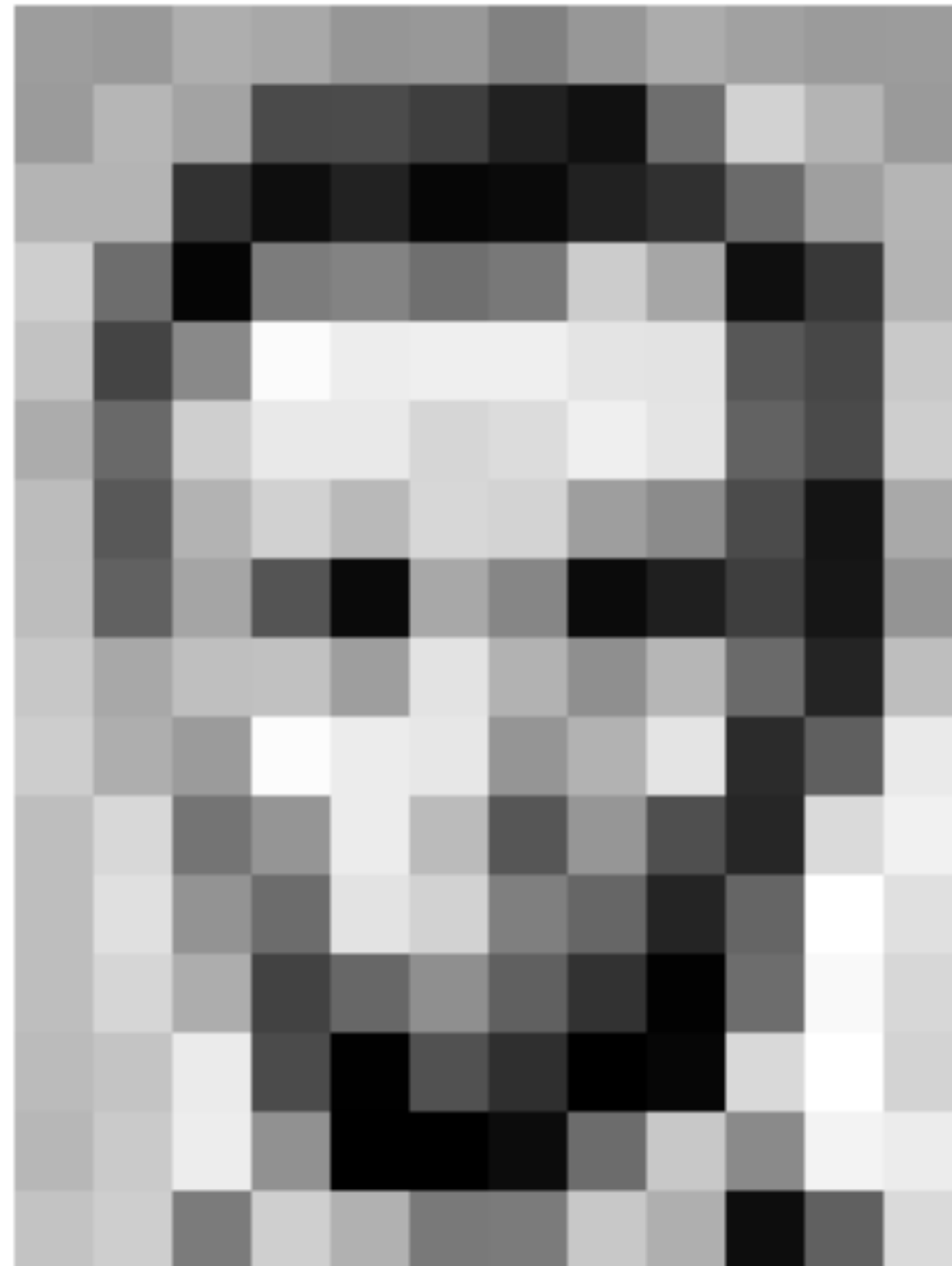# Kernel density — how smooth depends on variance / "width" of the Gaussian

# Kernel density — smooth in two dimensions

# Kernel density — this is familiar in "heatmaps"

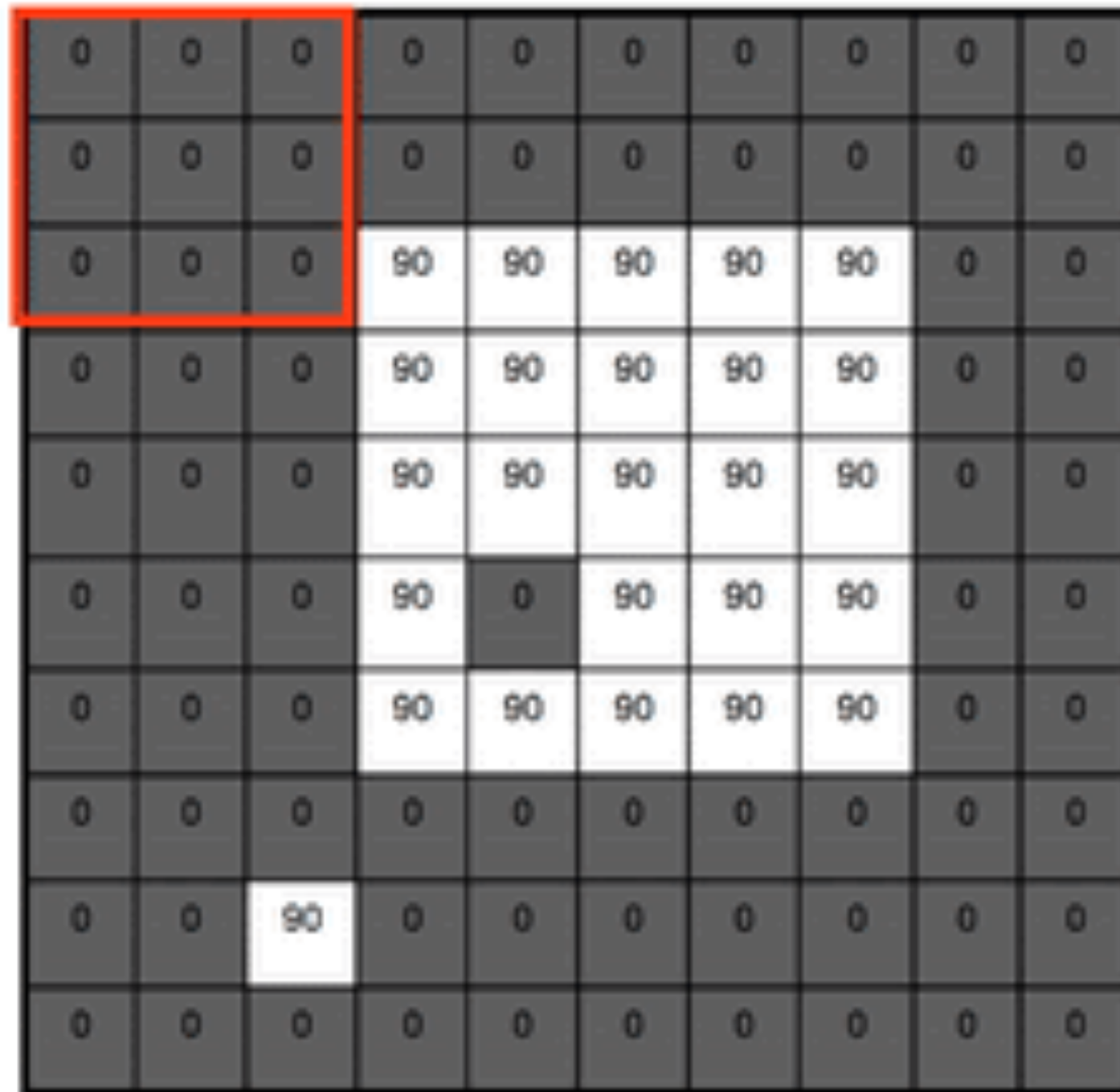# Now imagine an image as two-dimensional data — a grid of pixel intensities

An image *filter* is a kernel - a small window we *convolve* over an image.
The filter illustrated here averages the nine pixels in the window.



$F[x, y]$

$G[x, y]$

# A ~Gaussian kernel (high in the middle, lower away from the middle) acts as a smoothing or "blur filter"

| 0.0625 | 0.125 | 0.0625 |
|--------|-------|--------|
| 0.125  | 0.25  | 0.125  |
| 0.0625 | 0.125 | 0.0625 |



| -1 | -1 | -1 |
|----|----|----|
| -1 | 8  | -1 |
| -1 | -1 | -1 |

(a) Blur kernel.          (b) Blur kernel applied.

# The kernel on the right acts as an "edge filter"



| 0.0625 | 0.125 | 0.0625 |
|--------|-------|--------|
| 0.125  | 0.25  | 0.125  |
| 0.0625 | 0.125 | 0.0625 |

(a) Blur kernel.

(b) Blur kernel applied.

| -1 | -1 | -1 |
|----|----|----|
| -1 | 8  | -1 |
| -1 | -1 | -1 |

(c) Edge kernel.

(d) Edge kernel applied.

Figure 8: Effect of convolutional image kernels.

Input

Feature maps

f.maps

f.maps

Output

Convolutions

Subsampling

Convolutions

Subsampling

Fully connected

*Source:* Wikipedia, "Convolutional Neural Nets"

# CNN layers learn filters to detect and combine higher level "features"

# CNN layers learn filters to detect and combine higher level "features"



Layer 1

Layer 2

Layer 3

Convolutional Neural Network Visualization (Images)

http://scs.ryerson.ca/~aharley/vis/

# Typical CNN architecture for NLP



Figure 1: Illustration of a CNN architecture for sentence classification. We depict three filter region sizes: 2, 3 and 4, each of which has 2 filters. Filters perform convolutions on the sentence matrix and generate (variable-length) feature maps; 1-max pooling is performed over each map, i.e., the largest number from each feature map is recorded. Thus a univariate feature vector is generated from all six maps, and these 6 features are concatenated to form a feature vector for the penultimate layer. The final softmax layer then receives this feature vector as input and uses it to classify the sentence; here we assume binary classification and hence depict two possible output states.

*Source:* Zhang and Wallace (2015)

# 7-gram features detected by CNN



*Source:* Kalchbrenner, et al. (2014)

# Typical CNN architecture for NLP



activation function

convolution

1-max pooling

softmax function regularization in this layer

Sentence matrix 7 × 5

3 region sizes: (2,3,4) 2 filters for each region size totally 6 filters

2 feature maps for each region size

6 univariate vectors concatenated together to form a single feature vector

2 classes

d=5

I like this movie very much !

Figure 1: Illustration of a CNN architecture for sentence classification. We depict three filter region sizes: 2, 3 and 4, each of which has 2 filters. Filters perform convolutions on the sentence matrix and generate (variable-length) feature maps; 1-max pooling is performed over each map, i.e., the largest number from each feature map is recorded. Thus a univariate feature vector is generated from all six maps, and these 6 features are concatenated to form a feature vector for the penultimate layer. The final softmax layer then receives this feature vector as input and uses it to classify the sentence; here we assume binary classification and hence depict two possible output states.

*Source:* Zhang and Wallace (2015)

# 7-gram features detected by CNN



**POSITIVE**

| | | | | | | |
|---|---|---|---|---|---|---|
| lovely | comedic | moments | and | several | fine | performances |
| good | script | , | good | dialogue | , | funny |
| sustains | throughout | is | daring | , | inventive | and |
| well | written | , | nicely | acted | and | beautifully |
| remarkably | solid | and | subtly | satirical | tour | de |

**NEGATIVE**

| | | | | | | |
|---|---|---|---|---|---|---|
| , | nonexistent | plot | and | pretentious | visual | style |
| it | fails | the | most | basic | test | as |
| so | stupid | , | so | ill | conceived | , |
| , | too | dull | and | pretentious | to | be |
| hood | rats | butt | their | ugly | heads | in |

**'NOT'**

| | | | | | | |
|---|---|---|---|---|---|---|
| n't | have | any | huge | laughs | in | its |
| no | movement | , | no | , | not | much |
| n't | stop | me | from | enjoying | much | of |
| not | that | kung | pow | is | n't | funny |
| not | a | moment | that | is | not | false |

**'TOO'**

| | | | | | | |
|---|---|---|---|---|---|---|
| , | too | dull | and | pretentious | to | be |
| either | too | serious | or | too | lighthearted | , |
| too | slow | , | too | long | and | too |
| feels | too | formulaic | and | too | familiar | to |
| is | too | predictable | and | too | self | conscious |

*Source:* Kalchbrenner, et al. (2014)

```python
# Build the model
inputs = keras.Input(shape=(None,), dtype="int32")
x = layers.Embedding(max_features, 16)(inputs) #
x = layers.Conv1D(filters=128, kernel_size=5, strides=1, padding='same', activation='relu')(x) #
x = layers.GlobalMaxPooling1D()(x)
x = layers.Dense(16, activation = 'relu')(x)
outputs = layers.Dense(1, activation="sigmoid")(x) #
model = keras.Model(inputs, outputs) #
model.summary()
```
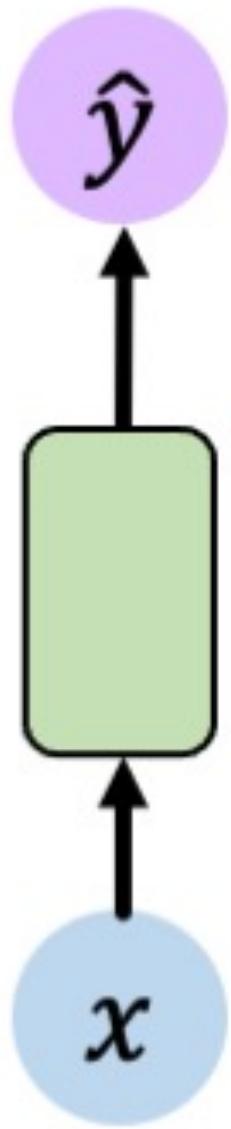
```
Model: "model_1"
_____
Layer (type)                 Output Shape              Param #
=================================================================
input_4 (InputLayer)         [(None, None)]            0

embedding_3 (Embedding)      (None, None, 16)          80000

conv1d_3 (Conv1D)            (None, None, 128)         10368

global_max_pooling1d_2 (Glob (None, 128)               0

dense_4 (Dense)              (None, 16)                2064

dense_5 (Dense)              (None, 1)                 17
=================================================================
Total params: 92,449
Trainable params: 92,449
Non-trainable params: 0
_____
```
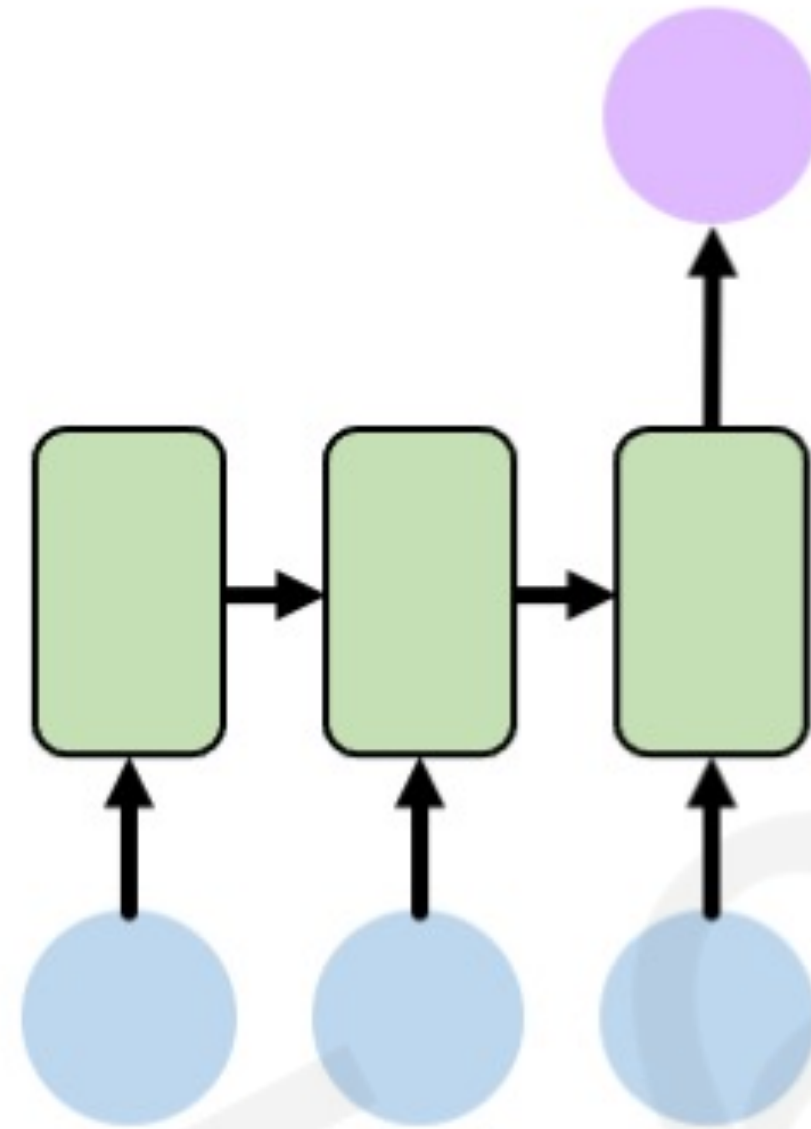
# Modeling sequence with recurrence
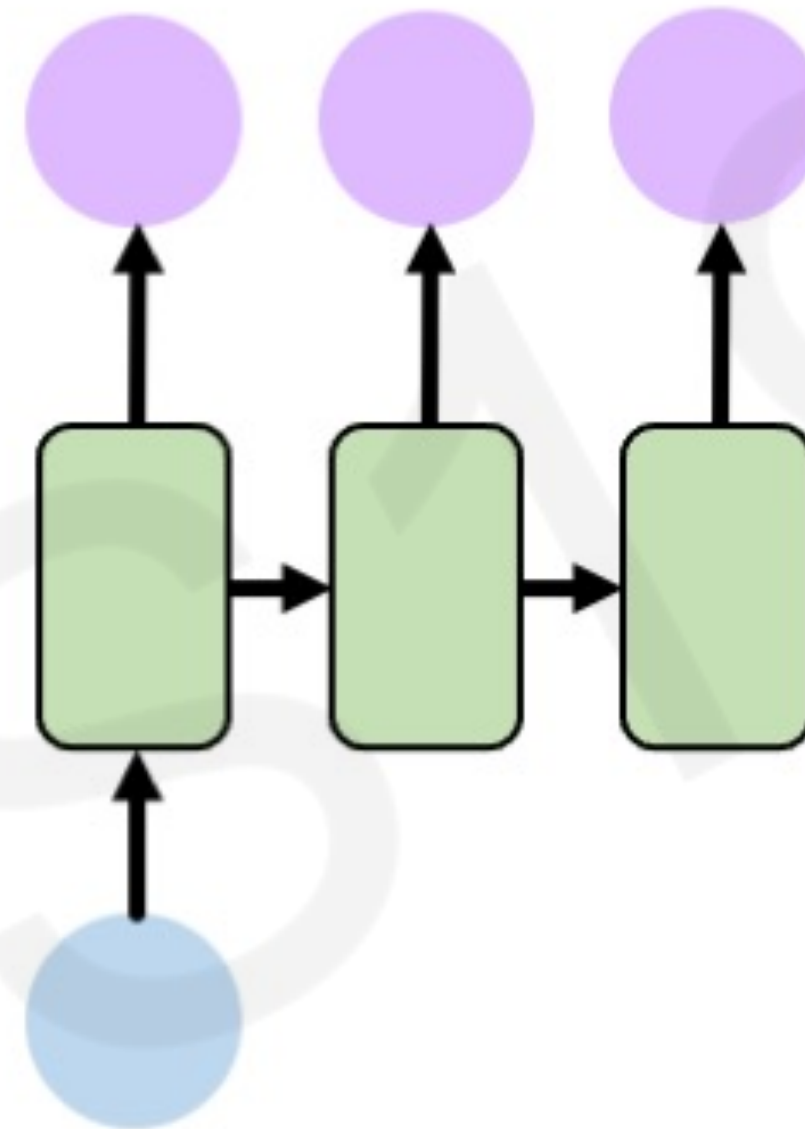
# Sequence Modeling Applications



One to One
**Binary Classification**

"Will I pass this class?"
Student → Pass?
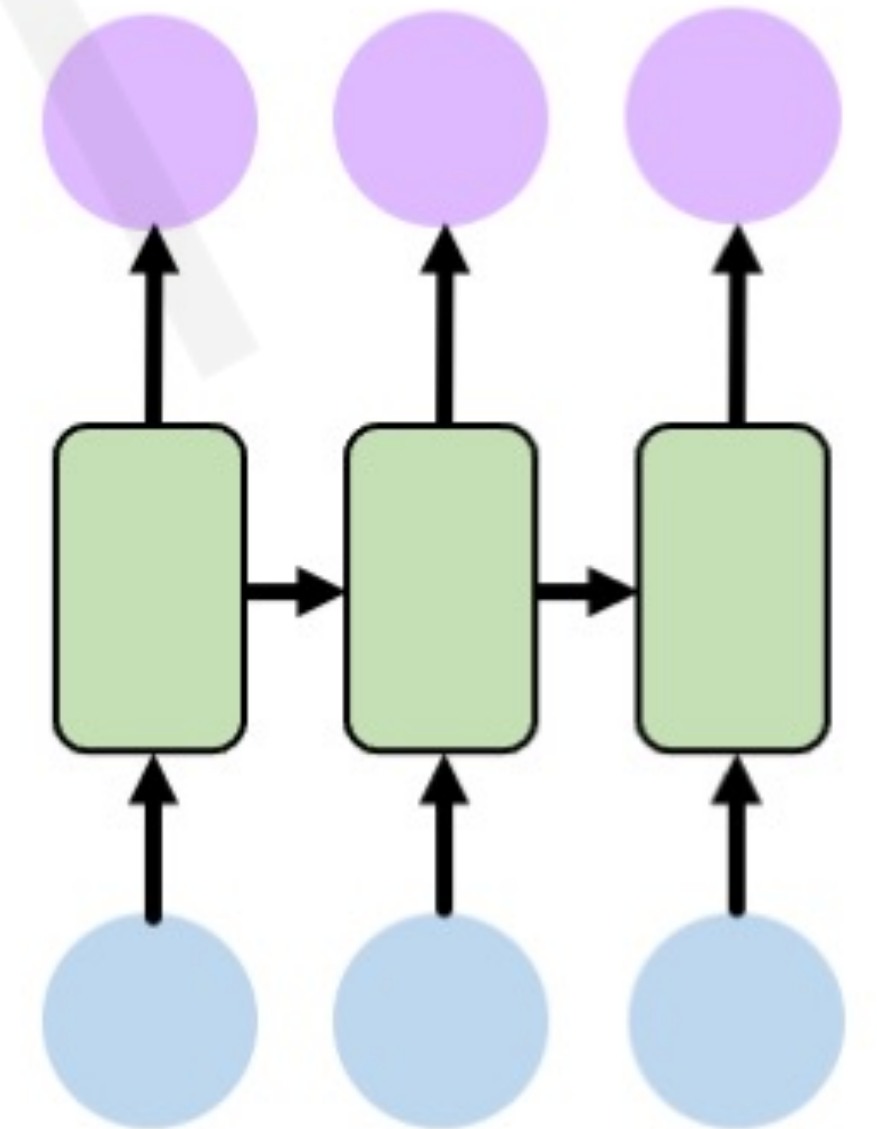
Many to One
**Sentiment Classification**

Ivar Hagendoorn
@IvarHagendoorn

Follow

The @MIT Introduction to #DeepLearning is
definitely one of the best courses of its kind
currently available online
introtodeeplearning.com

12:45 PM - 12 Feb 2018

😃

One to Many
**Image Captioning**

"A baseball player throws a ball."

Many to Many
**Machine Translation**

# The Perceptron Revisited

Massachusetts Institute of Technology

# Feed-Forward Networks Revisited



$$\boldsymbol{x} \in \mathbb{R}^m \qquad\qquad \hat{\boldsymbol{y}} \in \mathbb{R}^n$$

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/19/21

# Feed-Forward Networks Revisited



$$\boldsymbol{x_t} \in \mathbb{R}^m \qquad\qquad \boldsymbol{\hat{y}_t} \in \mathbb{R}^n$$

# Handling Individual Time Steps



output vector

$\hat{y}_t$     $\hat{y}_0$     $\hat{y}_1$     $\hat{y}_2$

input vector

$x_t$     $x_0$     $x_1$     $x_2$

$$\hat{y}_t = f(x_t)$$

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com  @MITDeepLearning

1/19/21

# Neurons with Recurrence



output vector $\hat{y}_t$

$\hat{y}_0$       $\hat{y}_1$       $\hat{y}_2$

$h_0$       $h_1$

input vector $x_t$

$x_0$       $x_1$       $x_2$

$$\hat{y}_t = f(x_t, h_{t-1})$$

output       input   past memory

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/19/21

# Neurons with Recurrence



$$\hat{y}_t = f(x_t, h_{t-1})$$

output     input   past memory

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/19/21

# Recurrent Neural Networks (RNNs)



output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

Apply a **recurrence relation** at every time step to process a sequence:

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

cell state

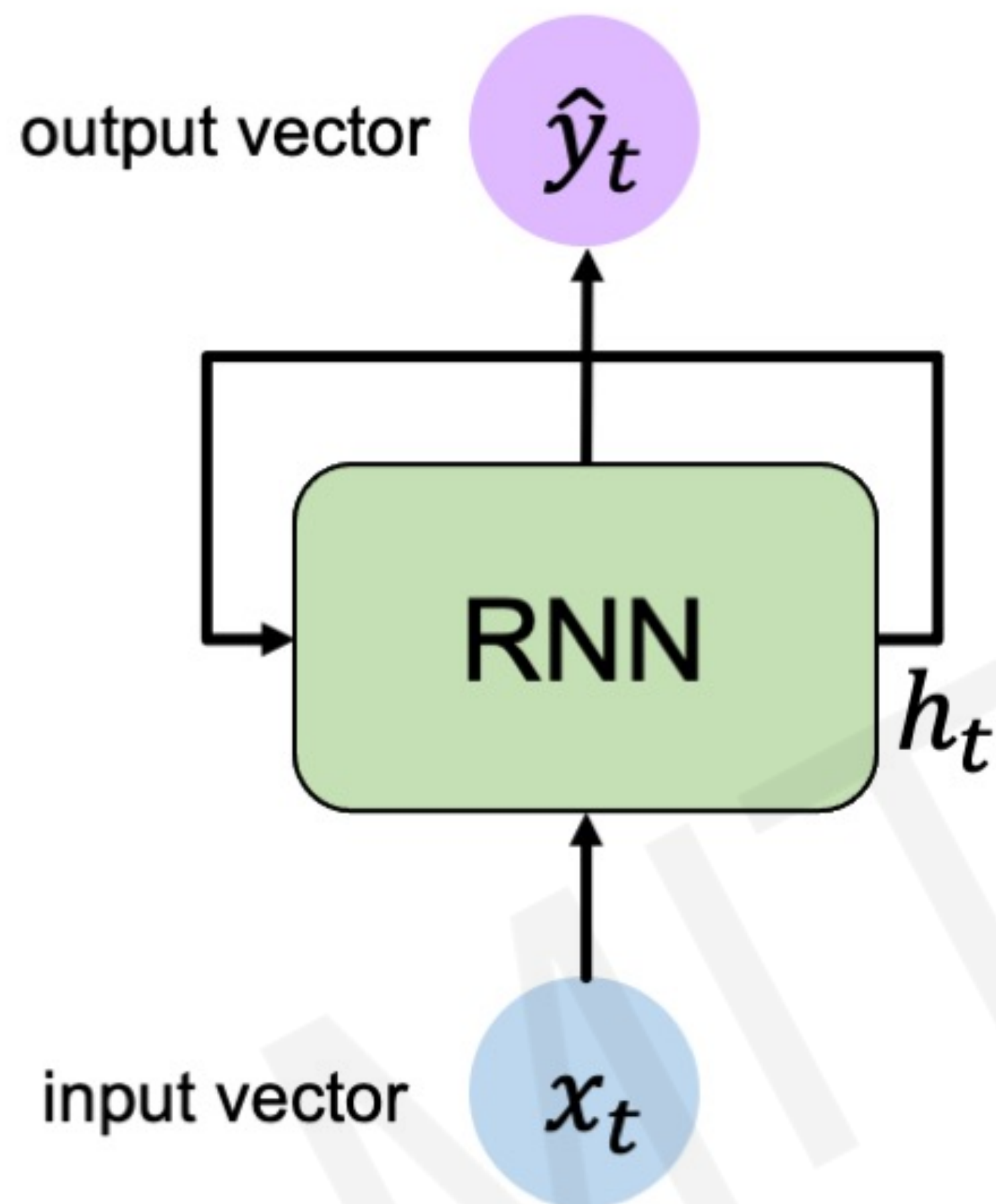function with weights W

input

old state

**Note:** the same function and set of parameters are used at every time step

RNNs have a **state**, $h_t$, that is updated **at each time step** as a sequence is processed
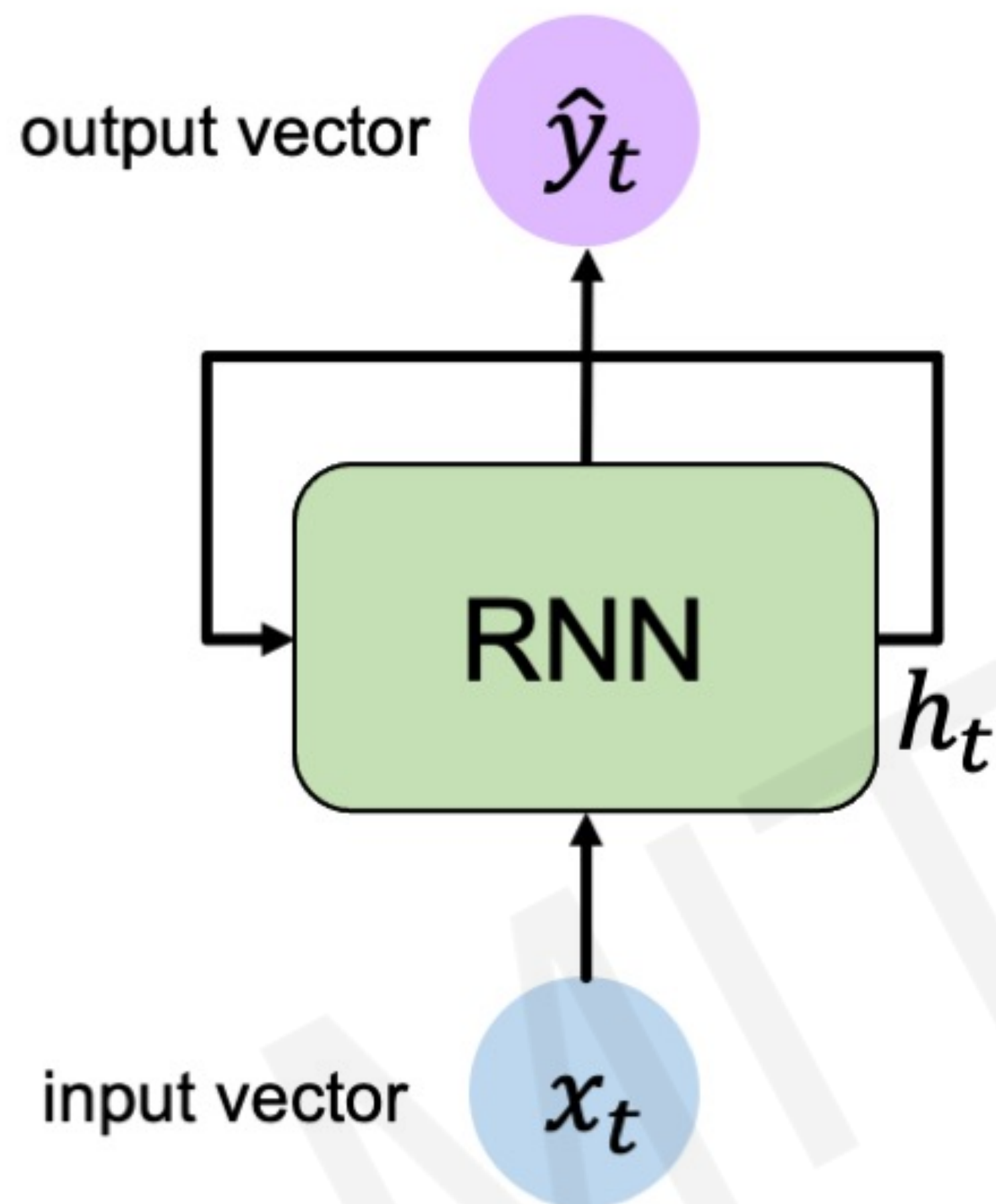
# RNN State Update and Output



output vector  $\hat{y}_t$

RNN  $h_t$

input vector  $x_t$

# RNN State Update and Output

output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

**Input Vector**

$x_t$

# RNN State Update and Output



output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

**Update Hidden State**

$$h_t = \tanh(\boldsymbol{W_{hh}^T} h_{t-1} + \boldsymbol{W_{xh}^T} x_t)$$

**Input Vector**

$x_t$

# RNN State Update and Output



output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

**Output Vector**

$$\hat{y}_t = \boldsymbol{W}_{\boldsymbol{hy}}^{\boldsymbol{T}} h_t$$

**Update Hidden State**

$$h_t = \tanh(\boldsymbol{W}_{\boldsymbol{hh}}^{\boldsymbol{T}} h_{t-1} + \boldsymbol{W}_{\boldsymbol{xh}}^{\boldsymbol{T}} x_t)$$

**Input Vector**

$$x_t$$

# RNNs: Computational Graph Across Time



$\hat{y}_t$

RNN = Represent as computational graph unrolled across time

$x_t$

# RNNs: Computational Graph Across Time

→ Forward pass

Re-use the **same weight matrices** at every time step

Massachusetts
Institute of
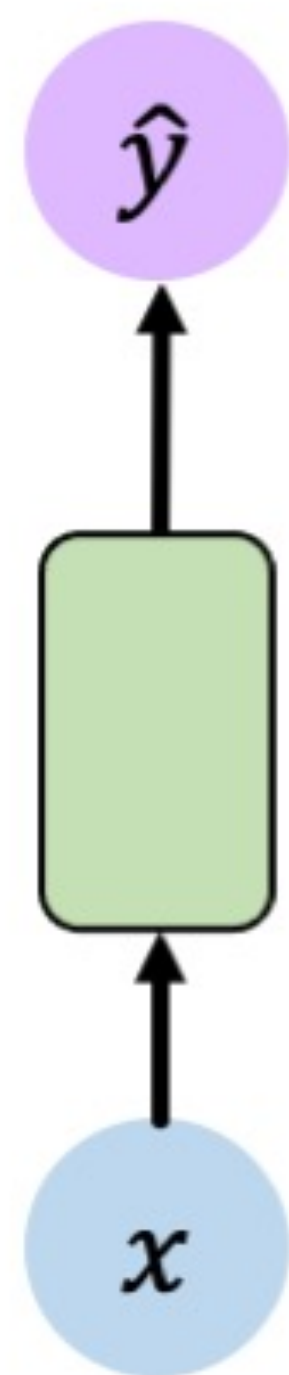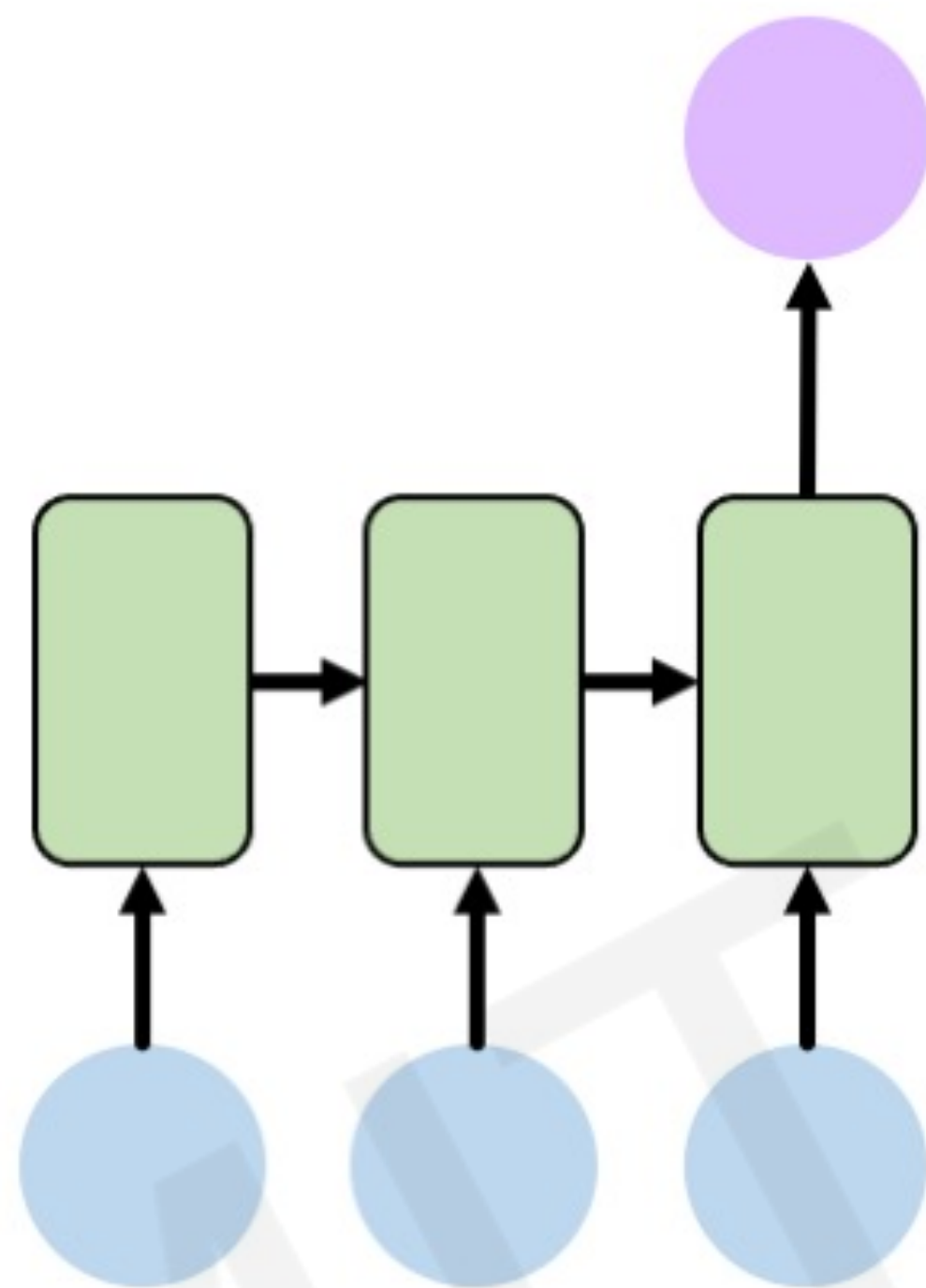Technology

# RNN Implementation in TensorFlow

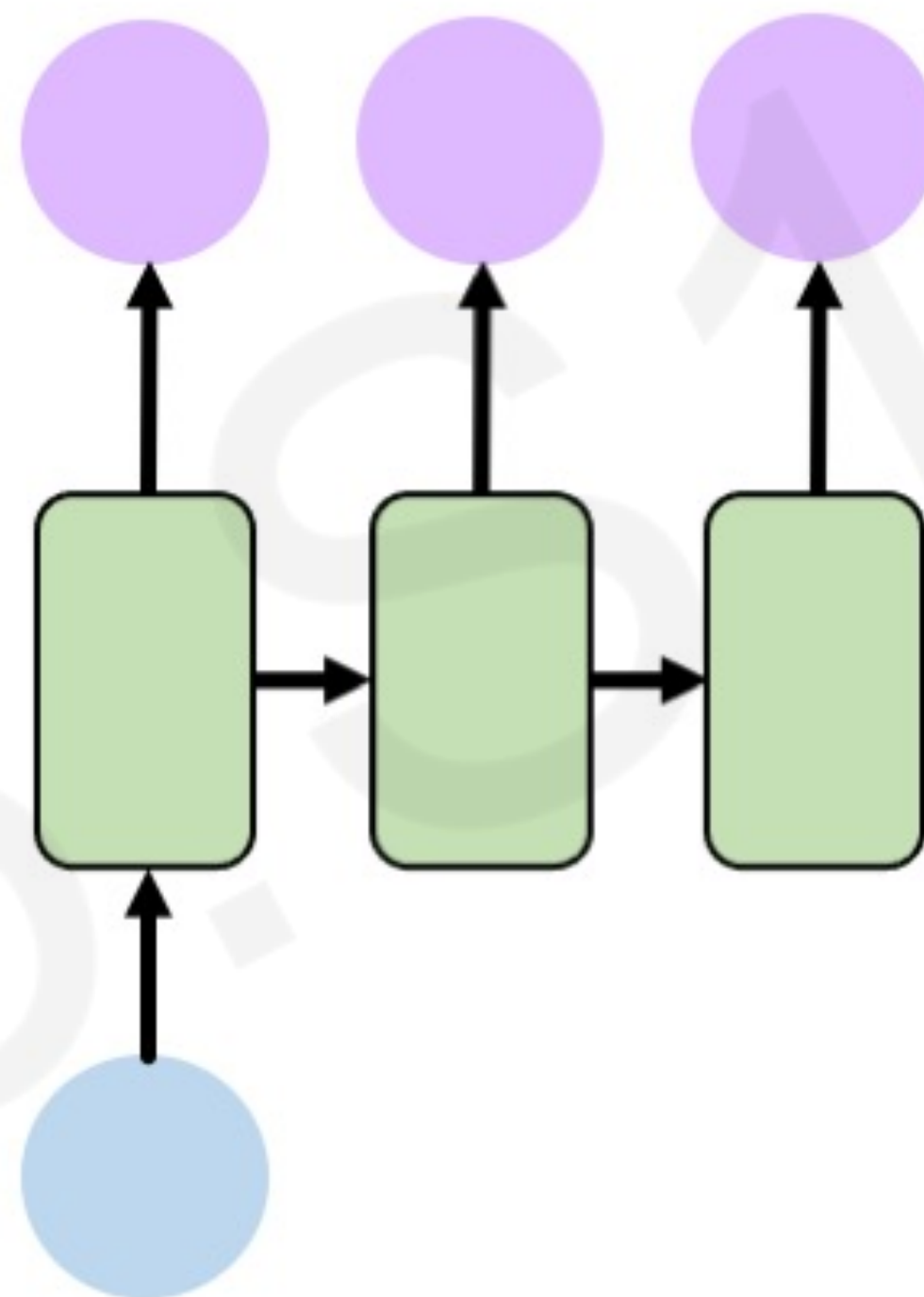`tf.keras.layers.SimpleRNN(rnn_units)`

# RNNs for Sequence Modeling



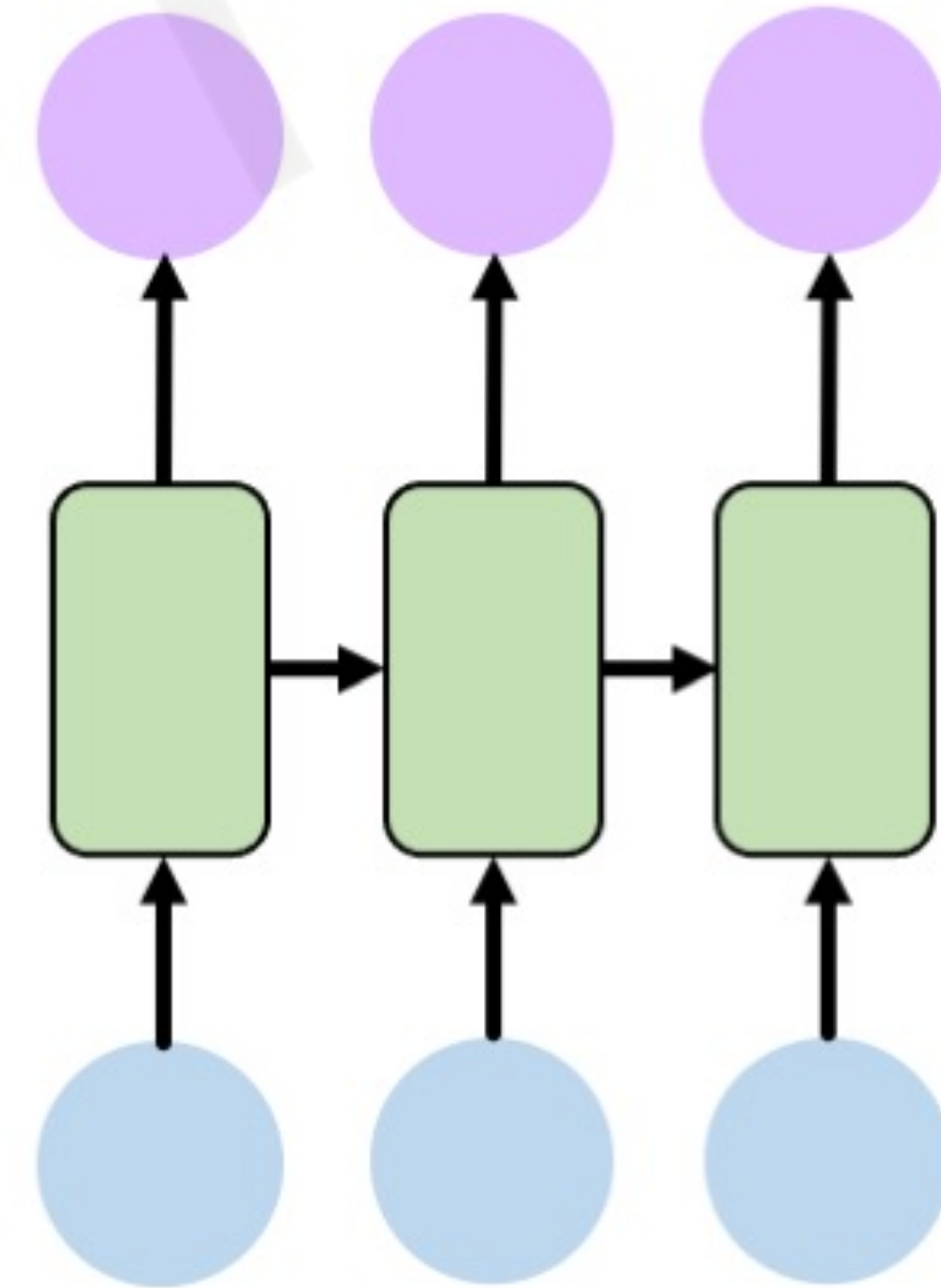One to One
"Vanilla" NN
*Binary classification*

Many to One
*Sentiment Classification*

One to Many
*Text Generation*
*Image Captioning*

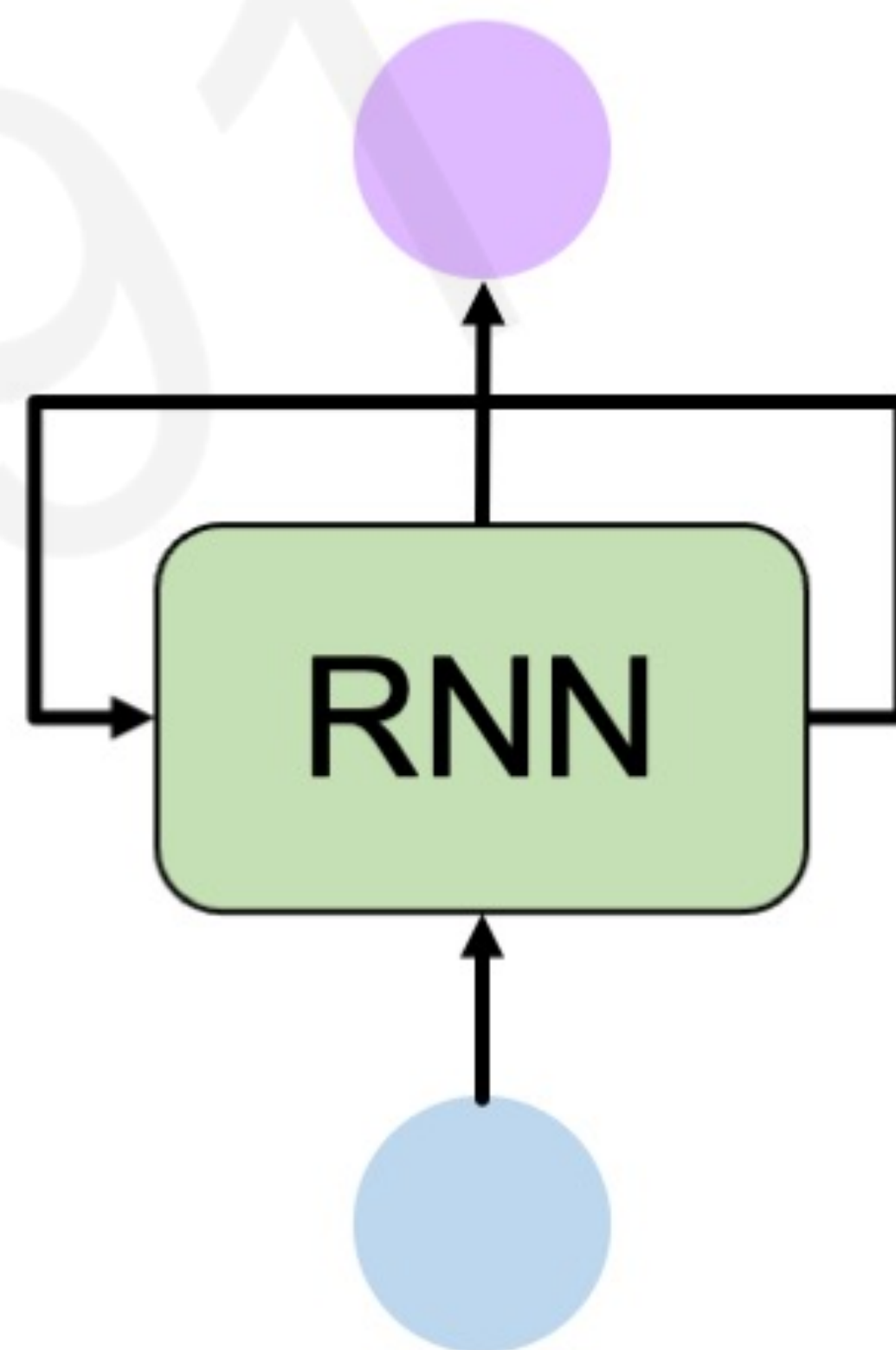Many to Many
*Translation & Forecasting*
*Music Generation*

**… and many other architectures and applications**   ⭐ **6.S191 Lab!**

# Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences

2. Track **long-term** dependencies

3. Maintain information about **order**

4. **Share parameters** across the sequence



**Recurrent Neural Networks (RNNs)** meet
these sequence modeling design criteria

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com  @MITDeepLearning

1/19/21

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

**given these words**

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

given these words                    predict the
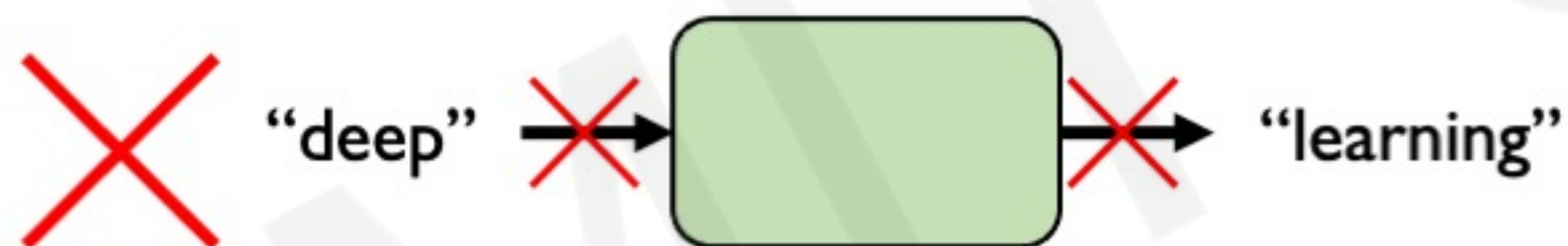                                     next word

# A Sequence Modeling Problem: Predict the Next Word
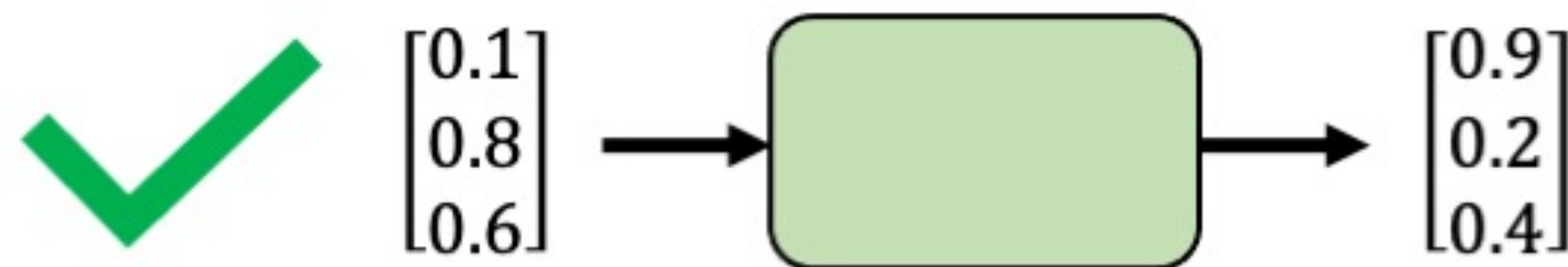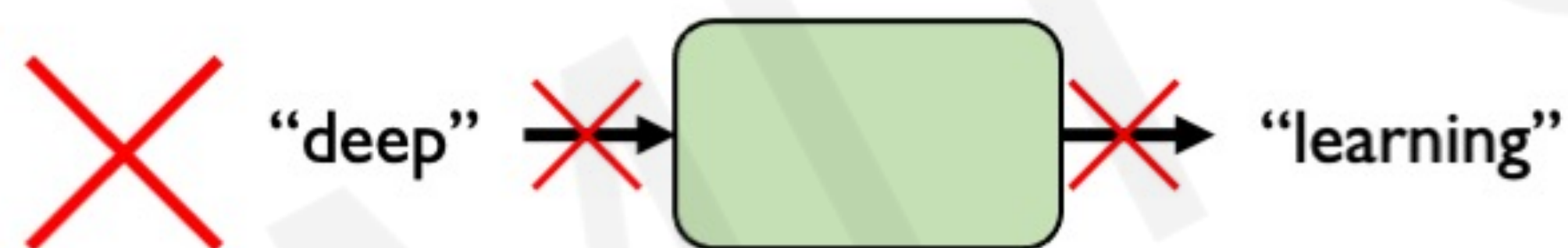
"This morning I took my cat for a walk."

given these words    predict the
next word

## Representing Language to a Neural Network



"deep" → → "learning"

*Neural networks cannot interpret words*

$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix}$ → → $\begin{bmatrix} 0.9 \\ 0.2 \\ 0.4 \end{bmatrix}$

*Neural networks require numerical inputs*

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

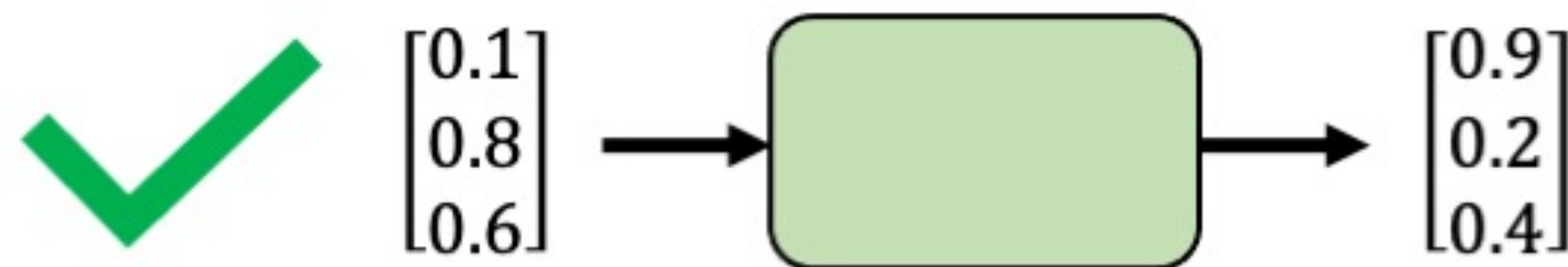given these words          predict the
                           next word
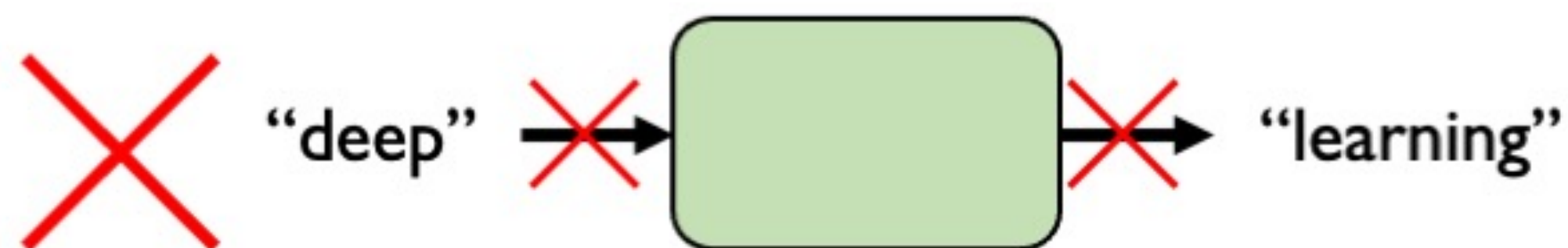
## Representing Language to a Neural Network



"deep" ✗→ ✗ "learning"

*Neural networks cannot interpret words*

$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix}$ → → $\begin{bmatrix} 0.9 \\ 0.2 \\ 0.4 \end{bmatrix}$

*Neural networks require numerical inputs*

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com     @MITDeepLearning

1/19/21

# Encoding Language for a Neural Network

❌ "deep" ❌ → [green box] ❌ → "learning"

*Neural networks cannot interpret words*

✅ $\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix}$ → [green box] → $\begin{bmatrix} 0.9 \\ 0.2 \\ 0.4 \end{bmatrix}$

*Neural networks require numerical inputs*

## Embedding: transform indexes into a vector of fixed size.

| | |
|---|---|
| this    cat    for <br> my   took <br>    I    walk <br> a     morning | a → 1 <br> cat → 2 <br> …    … <br> walk → N |

**One-hot embedding**

"cat" = [ 0, 1, 0, 0, 0, 0 ]

$i$-th index

**Learned embedding**

run walk     dog cat <br> day    sun     happy sad

**1. Vocabulary:** Corpus of words

**2. Indexing:** Word to index

**3. Embedding:** Index to fixed-sized vector

# Handle Variable Sequence Lengths

The food was great

vs.

We visited a restaurant for lunch

vs.

We were hungry but cleaned the house before eating

# Model Long-Term Dependencies

"**France** is where I grew up, but I now live in Boston. I speak fluent ___."

*J'aime 6.S191!*

We need information from **the distant past** to accurately predict the correct word.

# Capture Differences in Sequence Order
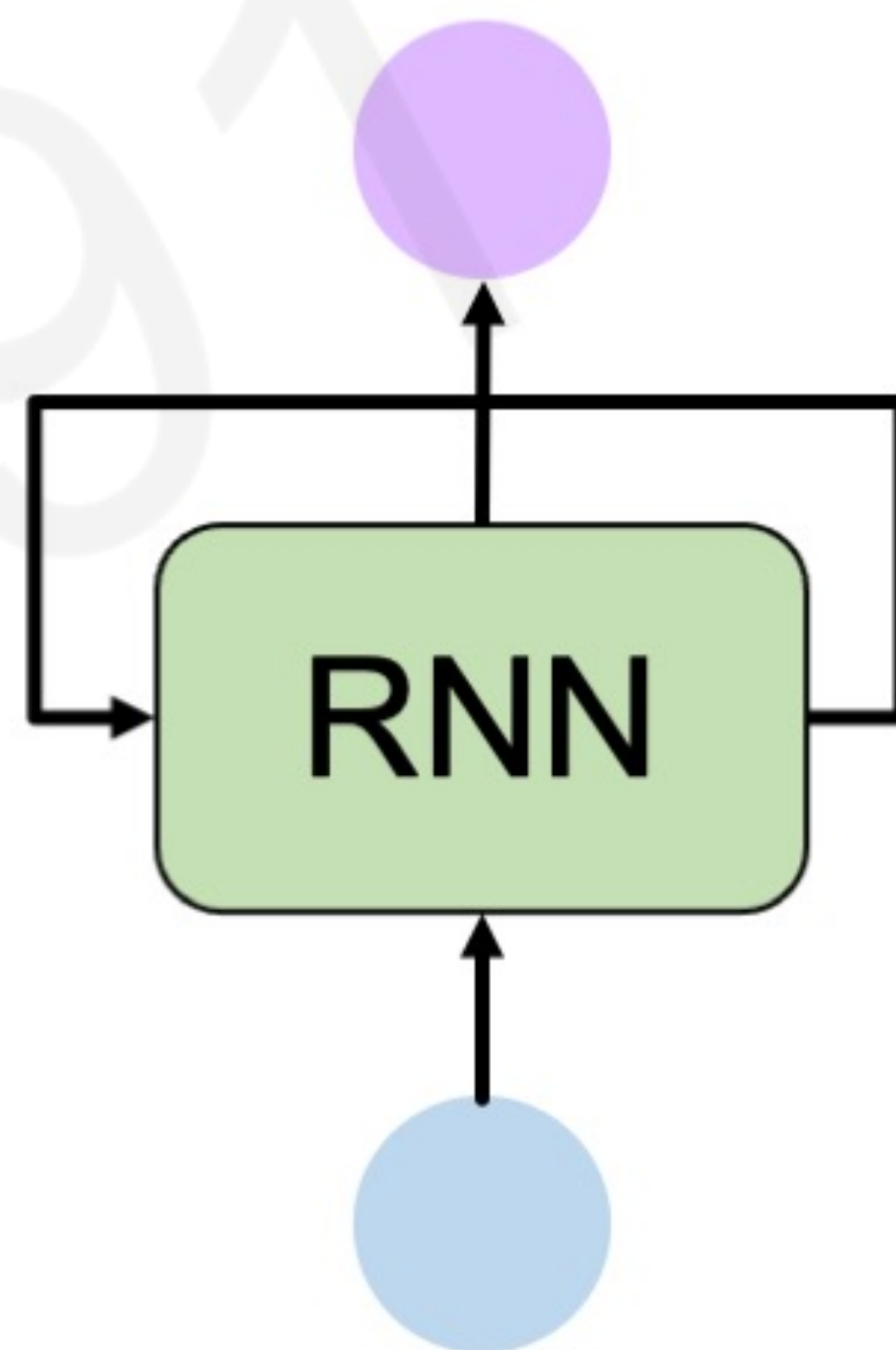


The food was good, not bad at all.

vs.

The food was bad, not good at all.

# Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences

2. Track **long-term** dependencies

3. Maintain information about **order**

4. **Share parameters** across the sequence

**Recurrent Neural Networks (RNNs)** meet
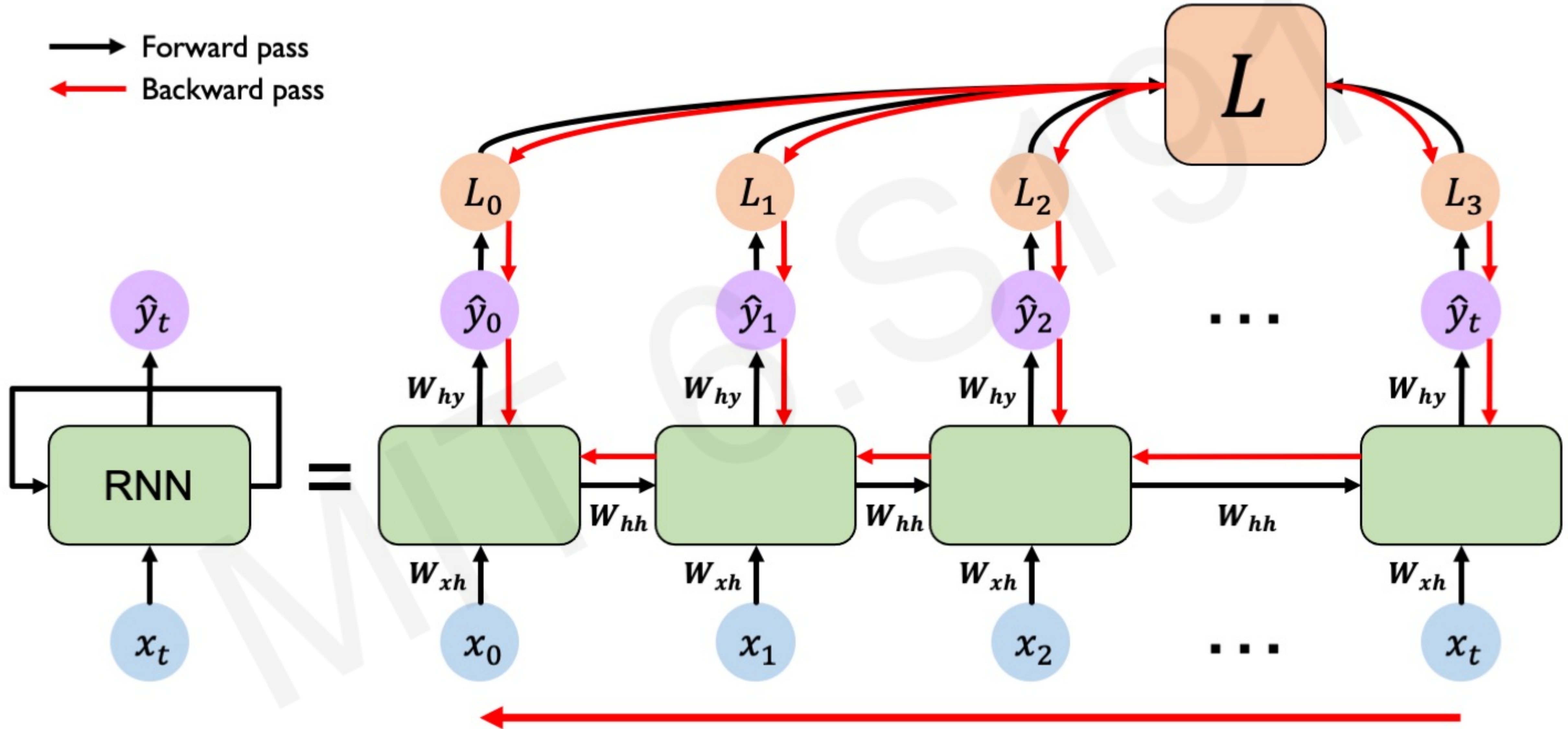these sequence modeling design criteria

# Recall: Backpropagation in Feed Forward Models



**Backpropagation algorithm:**

1. Take the derivative (gradient) of the loss with respect to each parameter
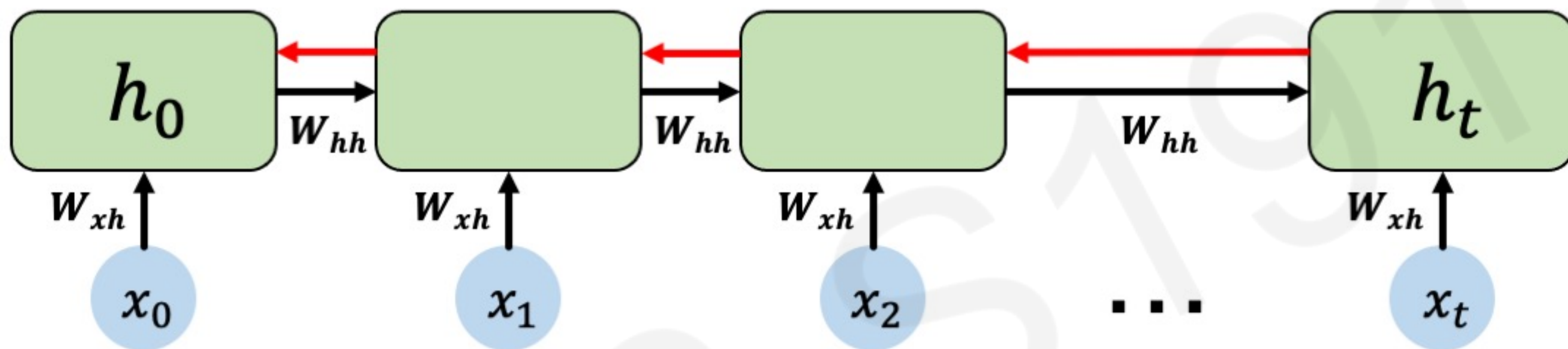
2. Shift parameters in order to minimize loss

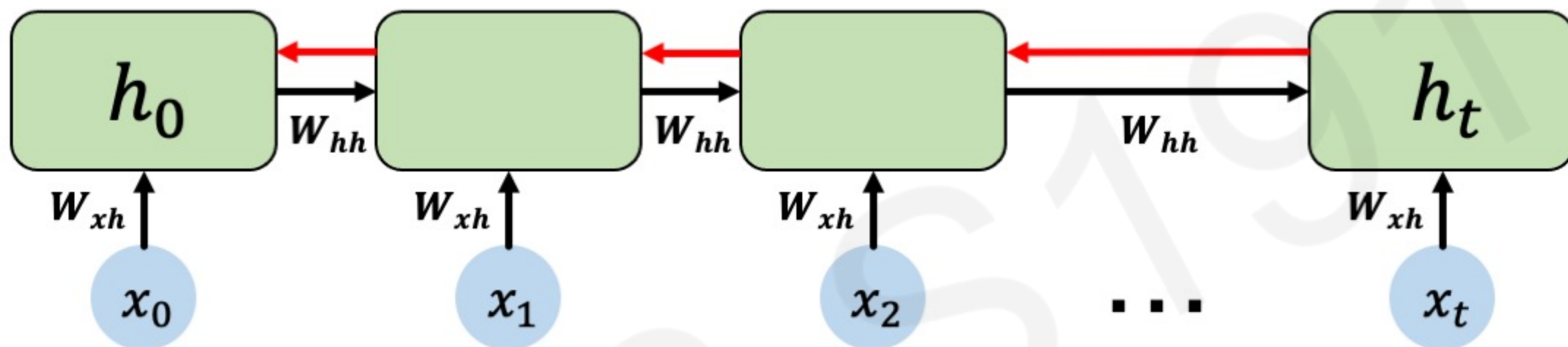# RNNs: Backpropagation Through Time



Forward pass

$L_0$  $L_1$  $L_2$  $L_3$

$\hat{y}_t$  $\hat{y}_0$  $\hat{y}_1$  $\hat{y}_2$  ...  $\hat{y}_t$

$W_{hy}$  $W_{hy}$  $W_{hy}$  $W_{hy}$

RNN  =  $W_{hh}$  $W_{hh}$  $W_{hh}$

$W_{xh}$  $W_{xh}$  $W_{xh}$  $W_{xh}$

$x_t$  $x_0$  $x_1$  $x_2$  ...  $x_t$

# RNNs: Backpropagation Through Time

# Standard RNN Gradient Flow

# Standard RNN Gradient Flow



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + **repeated gradient computation!**
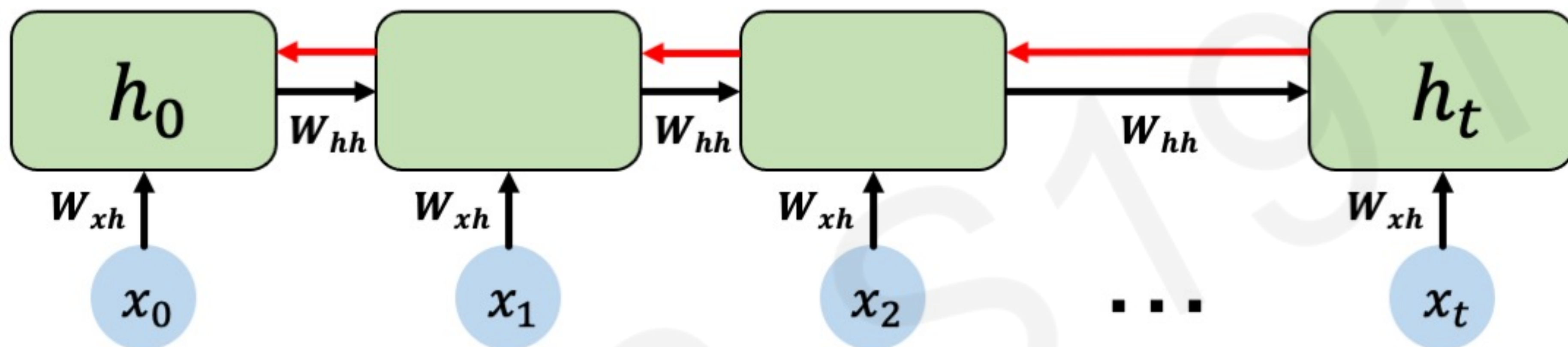
# Standard RNN Gradient Flow: Exploding Gradients



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + **repeated gradient computation!**

Many values > 1:
**exploding gradients**

**Gradient clipping** to scale big gradients

# Standard RNN Gradient Flow: Vanishing Gradients



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + **repeated gradient computation!**

Many values > 1:
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1:
**vanishing gradients**

1. Activation function
2. Weight initialization
3. Network architecture

# The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**
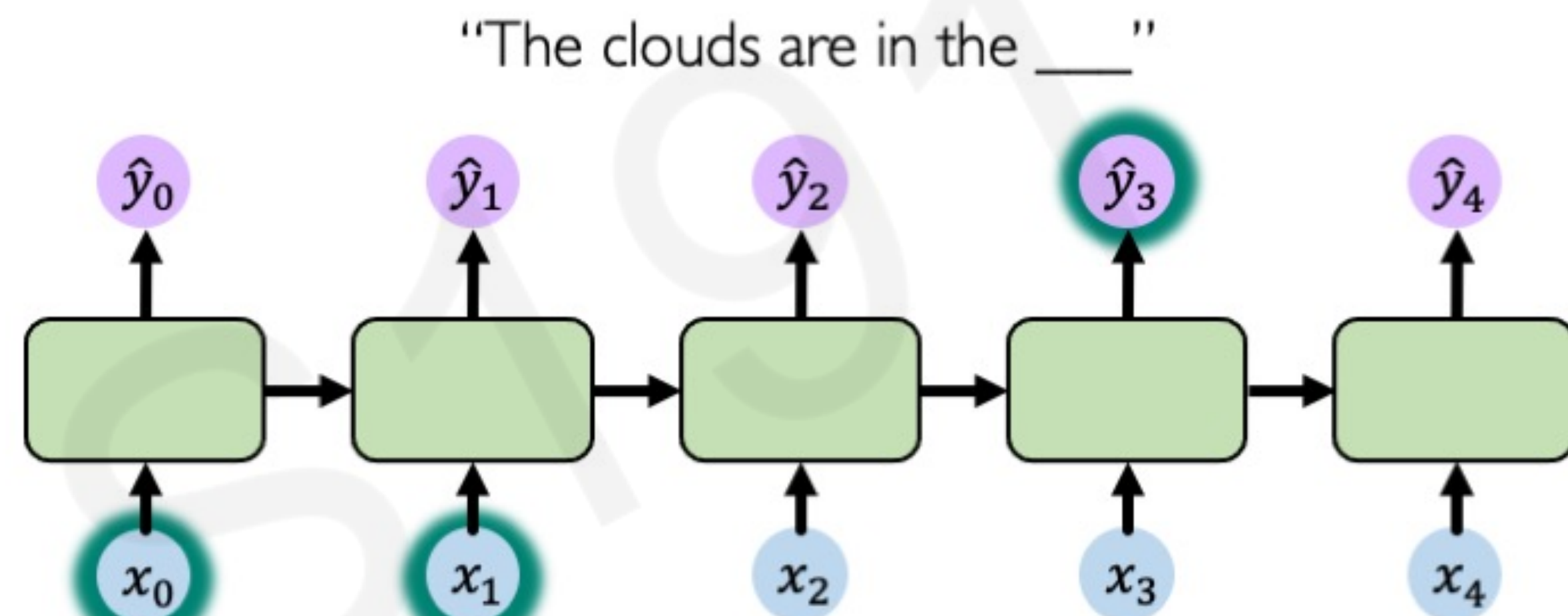
Multiply many **small numbers** together

↓

Errors due to further back time steps
have smaller and smaller gradients

# The Problem of Long-Term Dependencies

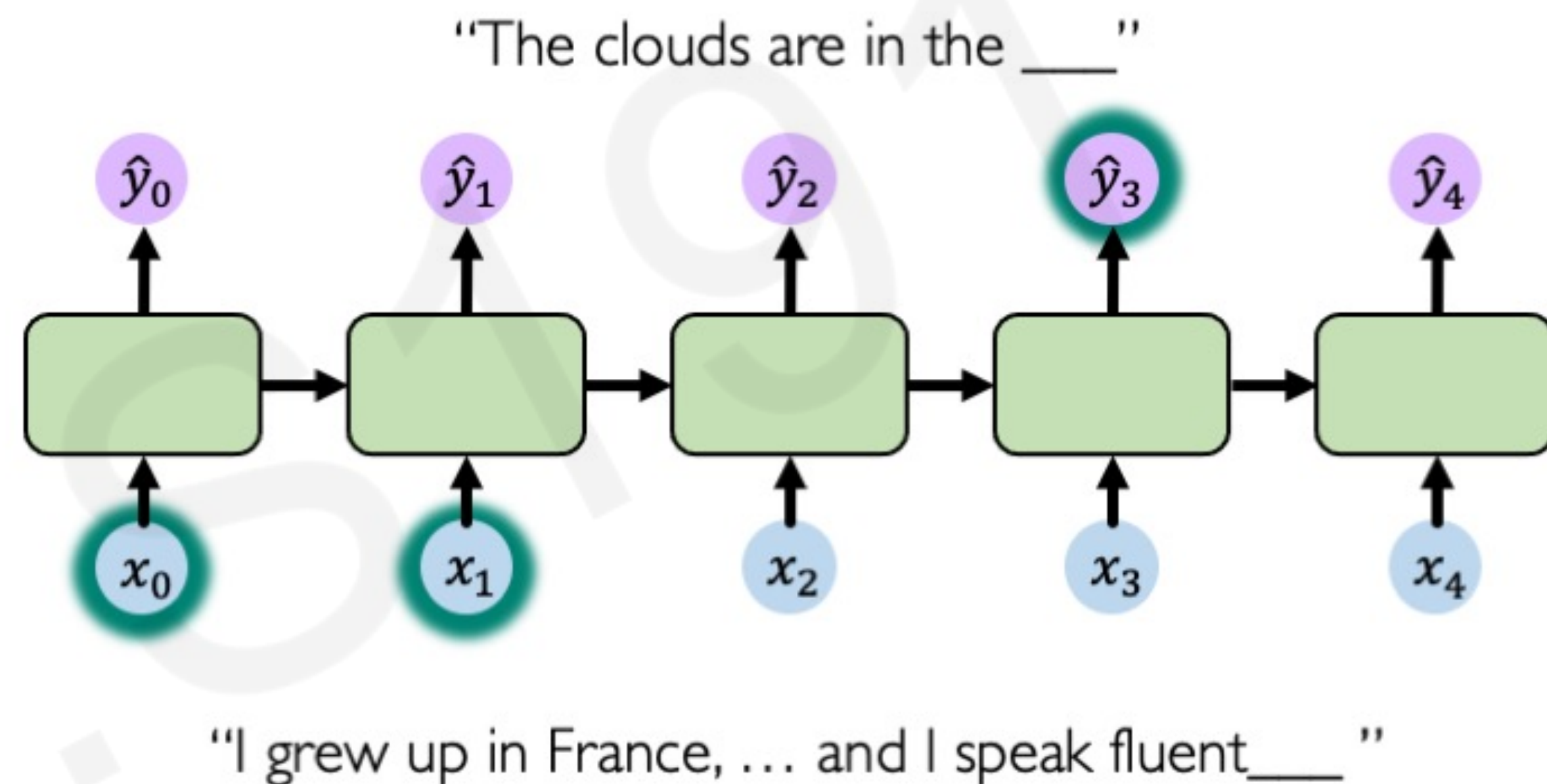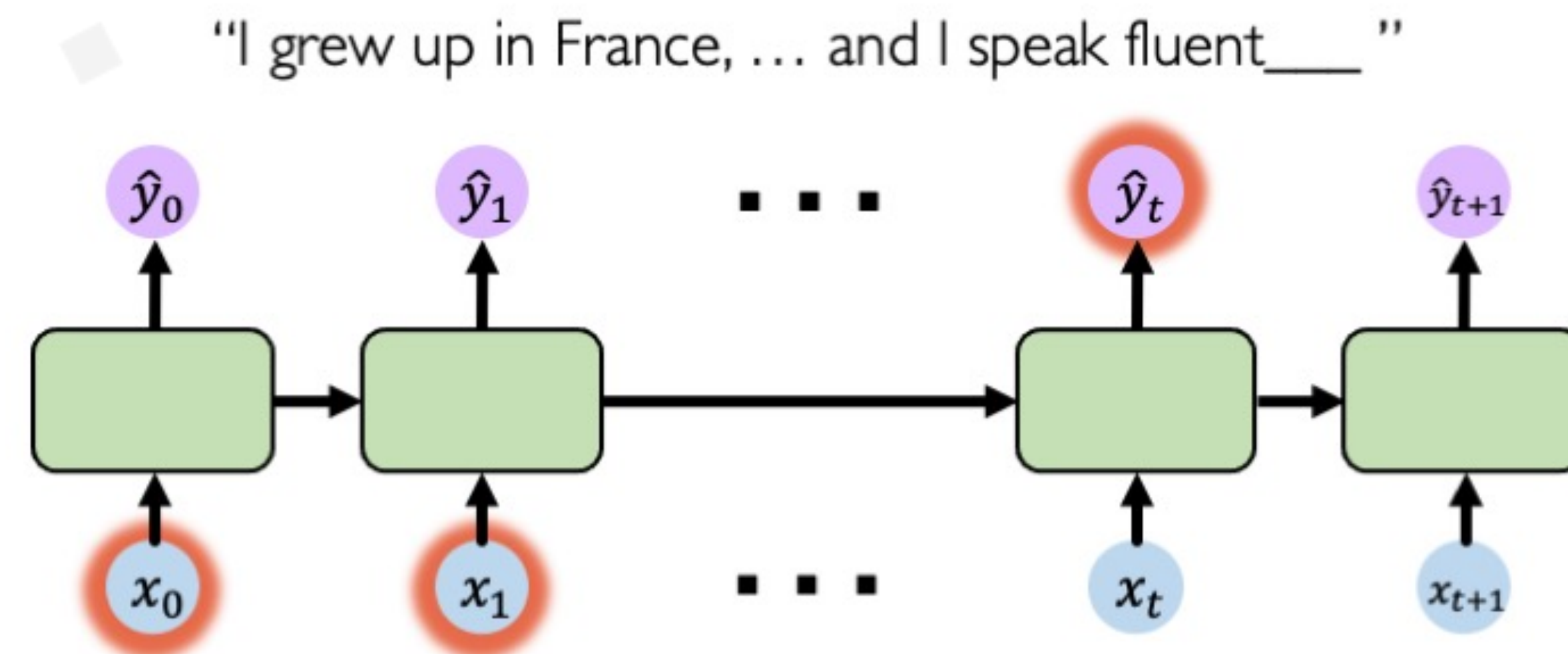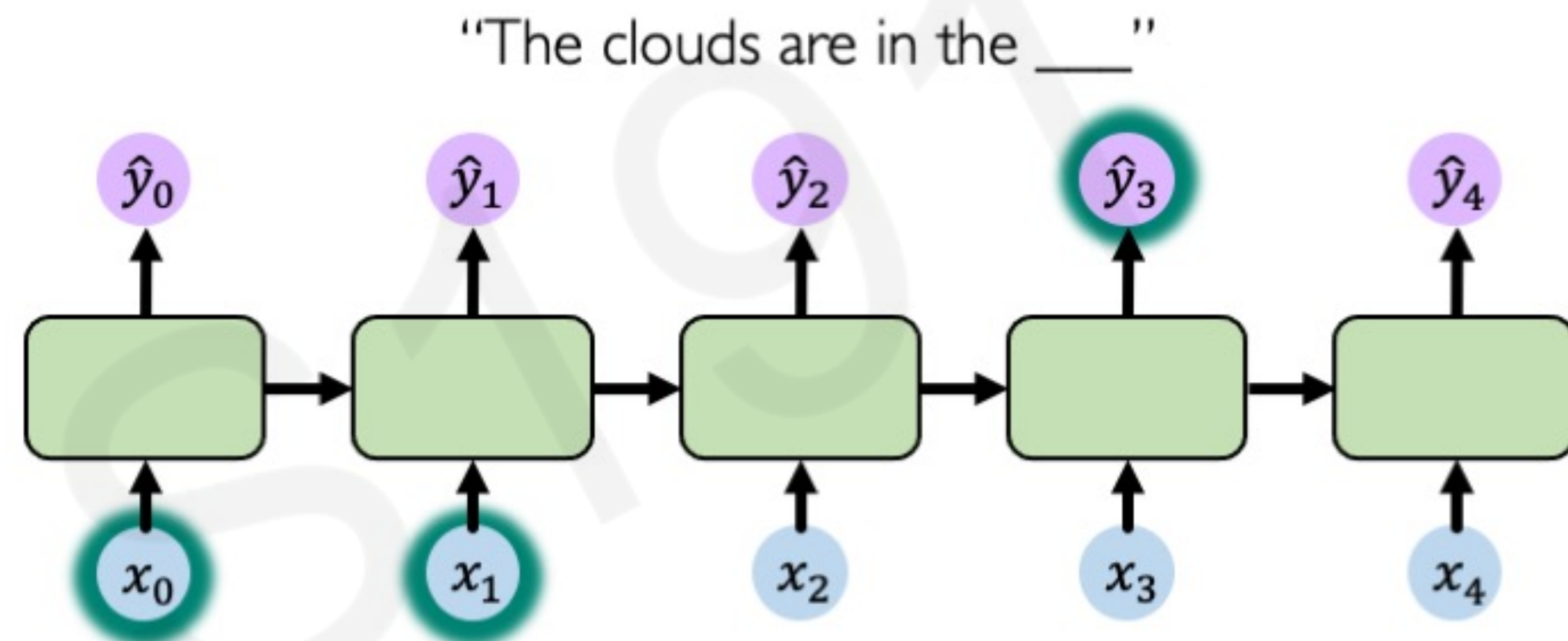**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

↓

Errors due to further back time steps
have smaller and smaller gradients

↓

Bias parameters to capture short-term
dependencies

Massachusetts
Institute of
Technology

# The Problem of Long-Term Dependencies

"The clouds are in the ___"

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

$\downarrow$

Errors due to further back time steps
have smaller and smaller gradients

$\downarrow$

Bias parameters to capture short-term
dependencies

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**
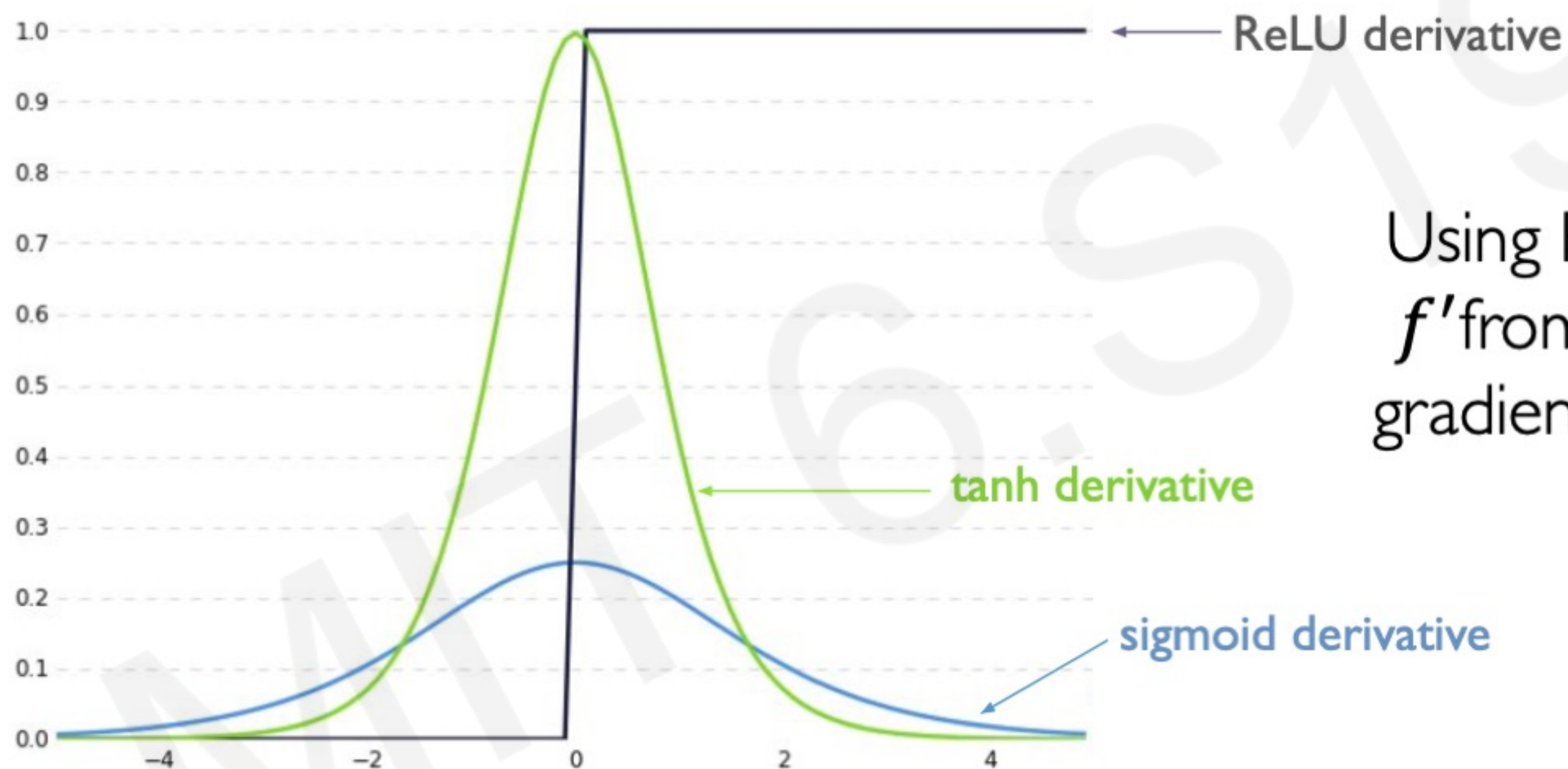
Multiply many **small numbers** together

↓

Errors due to further back time steps have smaller and smaller gradients

↓

Bias parameters to capture short-term dependencies

"The clouds are in the ___"

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

⬇

Errors due to further back time steps
have smaller and smaller gradients

⬇

Bias parameters to capture short-term
dependencies

"The clouds are in the ___"



"I grew up in France, … and I speak fluent___"

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

$\downarrow$

Errors due to further back time steps
have smaller and smaller gradients

$\downarrow$

Bias parameters to capture short-term
dependencies

"The clouds are in the ___"



"I grew up in France, … and I speak fluent___"

# Trick #1: Activation Functions



ReLU derivative

Using ReLU prevents $f'$ from shrinking the gradients when $x > 0$

tanh derivative

sigmoid derivative

# Trick #2: Parameter Initialization

SKIP

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

Massachusetts
Institute of
Technology

# Solution #3: Gated Cells

Idea: use a more **complex recurrent unit with gates** to control what information is passed through



gated cell

LSTM, GRU, etc.

**Long Short Term Memory (LSTMs)** networks rely on a gated cell to track information throughout many time steps.

# Standard RNN

In a standard RNN, repeating modules contain a **simple computation node**

# Long Short Term Memory (LSTMs)

LSTM modules contain **computational blocks** that **control information flow**



LSTM cells are able to track information throughout many timesteps

`tf.keras.layers.LSTM(num_units)`

# Long Short Term Memory (LSTMs)

Information is **added** or **removed** through structures called **gates**



Gates optionally let information through, for example via a sigmoid neural net layer and pointwise multiplication

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/19/21

# Long Short Term Memory (LSTMs)

## How do LSTMs work?

## 1) Forget   2) Store   3) Update   4) Output

# Long Short Term Memory (LSTMs)

**1) Forget**   2) Store   3) Update   4) Output

LSTMs **forget irrelevant** parts of the previous state

# Long Short Term Memory (LSTMs)

1) Forget   **2) Store**   3) Update   4) Output

LSTMs **store relevant** new information into the cell state

# Long Short Term Memory (LSTMs)

1) Forget   2) Store   **3) Update**   4) Output

LSTMs **selectively update** cell state values

# Long Short Term Memory (LSTMs)

1) Forget   2) Store   3) Update   **4) Output**

The **output gate** controls what information is sent to the next time step

# Long Short Term Memory (LSTMs)

## 1) Forget   2) Store   3) Update   4) Output

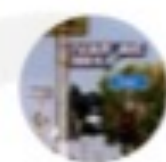# LSTM Gradient Flow

## Uninterrupted gradient flow!

# LSTMs: Key Concepts

1. Maintain a **separate cell state** from what is outputted

2. Use **gates** to control the **flow of information**

   - **Forget** gate gets rid of irrelevant information

   - **Store** relevant information from current input

   - Selectively **update** cell state

   - **Output** gate returns a filtered version of the cell state

3. Backpropagation through time with **uninterrupted gradient flow**

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/19/21

# Example Task: Sentiment Classification

sentiment
<positive>



**Tweet sentiment classification**

**Ivar Hagendoorn**
@IvarHagendoorn

Follow

The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online introtodeeplearning.com

12:45 PM - 12 Feb 2018

**Angels-Cave**
@AngelsCave

Follow

Replying to @Kazuki2048

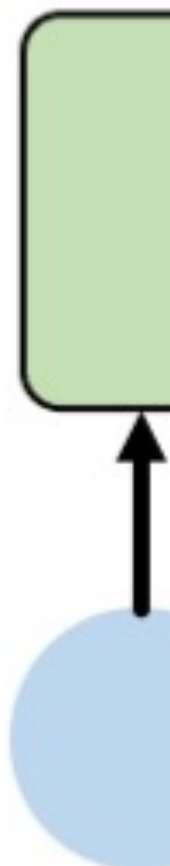I wouldn't mind a bit of snow right now. We haven't had any in my bit of the Midlands this winter! :(

2:19 AM - 25 Jan 2019

I    love    this    class!

# Example Task: Sentiment Classification

```
[ ] inputs = keras.Input(shape=(None,), dtype="int32")
    x = layers.Embedding(max_features, 16)(inputs) # define a 16-dimensional Embedding layer that acts on "inputs
    x = layers.LSTM(16)(x) # Add a 16-node LSTM that acts on the output of the Embedding layer
    outputs = layers.Dense(1, activation="sigmoid")(x) # define a 1-node sigmoid / classifier layer that acts on
    model = keras.Model(inputs, outputs) # define the model as inputs -> outputs
    model.summary()
```

```
Model: "model"
_____
Layer (type)                    Output Shape             Param #
===============================================================
input_3 (InputLayer)            [(None, None)]           0
_____
embedding_1 (Embedding)         (None, None, 16)         80000
_____
lstm_1 (LSTM)                   (None, 16)               2112
_____
dense_1 (Dense)                 (None, 1)                17
===============================================================
Total params: 82,129
Trainable params: 82,129
Non-trainable params: 0
_____
```

I    love    this    class!

winter! :(
2:19 AM - 25 Jan 2019

# Example Task: Sentiment Classification

```
[ ]  inputs = keras.Input(shape=(None,), dtype="int32")
     x = layers.Embedding(max_features, 16)(inputs) # define a 16-dimensional Embedding layer that acts on "input
     x = layers.Bidirectional(layers.LSTM(16))(x) # Add a 16-node bi-LSTM that acts on the output of the Embeddin
     outputs = layers.Dense(1, activation="sigmoid")(x) # define a 1-node sigmoid / classifier layer that acts on
     model = keras.Model(inputs, outputs) # define the model as inputs -> outputs
     model.summary()
```

```
Model: "model_1"

_____
Layer (type)                 Output Shape              Param #
=================================================================
input_4 (InputLayer)         [(None, None)]            0

_____
embedding_2 (Embedding)      (None, None, 16)          80000

_____
bidirectional (Bidirectional (None, 32)                4224

_____
dense_2 (Dense)              (None, 1)                 33
=================================================================
Total params: 84,257
Trainable params: 84,257
Non-trainable params: 0
_____
```

2:19 AM - 25 Jan 2019

# One big bi-LSTM success was ELMo — contextual embeddings

Embedding of "stick" in "Let's stick to" - Step #1

# Embedding of "stick" in "Let's stick to" - Step #2

## 1- Concatenate hidden layers

Forward Language Model

Backward Language Model

## 2- Multiply each vector by a weight based on the task

x $s_2$

x $s_1$

x $s_0$

Let's          stick          to

Let's          stick          to

## 3- Sum the (now weighted) vectors

ELMo embedding of "stick" for this task in this context

# Example Task: Machine Translation



Encoder (English)

Decoder (French)

# Example Task: Machine Translation

Sequence-to-sequence, or seq2seq



Encoder (English)

Decoder (French)

Massachusetts
Institute of
Technology

# Example Task: Machine Translation

## Potential Issues



Encoding bottleneck

le     chien     mange

the     dog     eats            le     chien

**Encoder (English)**            **Decoder (French)**

# Example Task: Machine Translation

## Potential Issues

Encoding bottleneck

Slow, no parallelization

le      chien      mange

the      dog      eats            le      chien

**Encoder (English)**            **Decoder (French)**

# Example Task: Machine Translation

## Potential Issues

🔽 Encoding bottleneck

🕐 Slow, no parallelization

🧠 Not long memory



Encoder (English)

Decoder (French)

# Example Task: Machine Translation

Attention mechanisms in neural networks
provide **learnable memory access**



Encoder (English)

Decoder (French)

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com   @MITDeepLearning

Sutskever+, *NeurIPS* 2014; Bahdanau+ *ICLR* 2015;
Vaswani+, *NeurIPS* 2017.   1/19/21

# Attention

Following images/videos from Jay Alammar, "The Illustrated Transformer" and "Visualizing a Neural Machine Translation Model (Mechanics of Seq2seq Models with Attention")

## SEQUENCE TO SEQUENCE MODEL



## Neural Machine Translation

### SEQUENCE TO SEQUENCE MODEL

The context is a vector of floats. Later in this post we will visualize vectors in color by assigning brighter colors to the cells with higher values.

We need to turn the input words into vectors before processing them. That transformation is done using a word embedding algorithm. We can use pre-trained embeddings or train our own embedding on our dataset. Embedding vectors of size 200 or 300 are typical, we're showing a vector of size four for simplicity.

# Recurrent Neural Network

**Time step #1:**

An RNN takes two input vectors:



hidden
state #0

input vector #1

hidden
state #0

input #1

# Time step:

## Neural Machine Translation
### SEQUENCE TO SEQUENCE MODEL

Je   suis   étudiant → **ENCODER** → **DECODER** →

## Neural Machine Translation
### SEQUENCE TO SEQUENCE MODEL

**Encoding Stage**

Encoder RNN →

**Decoding Stage**

Decoder RNN

Je          suis          étudiant

# Neural Machine Translation

**SEQUENCE TO SEQUENCE MODEL WITH ATTENTION**



Encoding Stage

Encoder RNN

Decoding Stage

Attention Decoder RNN

Je          suis          étudiant

**Attention at time step 4**

1. The attention decoder RNN takes in the embedding of the <END> token, and an initial decoder hidden state.
2. The RNN processes its inputs, producing an output and a new hidden state vector (h4). The output is discarded.
3. Attention Step: We use the encoder hidden states and the h4 vector to calculate a context vector (C4) for this time step.
4. We concatenate h4 and C4 into one vector.
5. We pass this vector through a feedforward neural network (one trained jointly with the model).
6. The output of the feedforward neural networks indicates the output word of this time step.
7. Repeat for the next time steps

# Neural Machine Translation
## SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

**Encoding Stage**
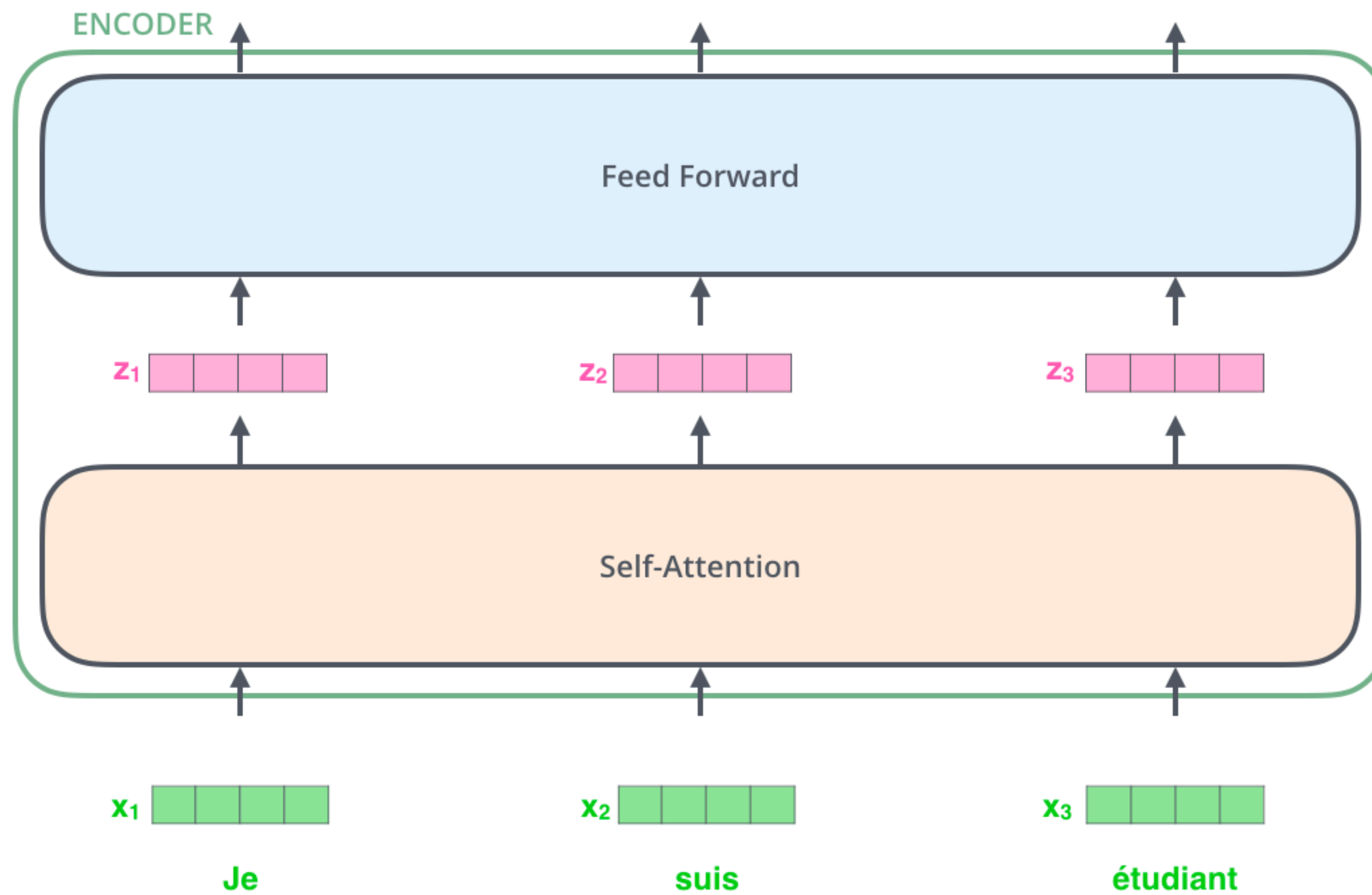
**Attention Decoding Stage**

$h_1\,h_2\,h_3$

$h_{init}$

<END>

4

The transformer ("Attention is All You Need")

INPUT

Je suis étudiant

THE
TRANSFORMER

OUTPUT

I am a student

OUTPUT

I am a student

ENCODERS → DECODERS

INPUT

Je suis étudiant

ENCODER #2

ENCODER #1

$r_1$

$r_2$

Feed Forward
Neural Network

Feed Forward
Neural Network

$z_1$

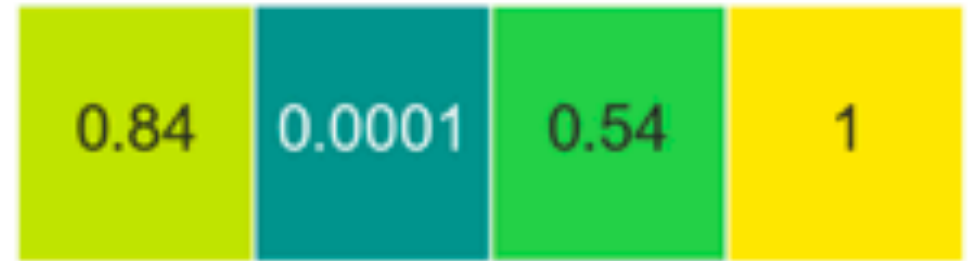$z_2$

Self-Attention

$x_1$

$x_2$

**Thinking**

**Machines**

Matrix math of self-attention - yada yada yada

Multi-headed attention - blah blah blah