

A Sequence Modeling Problem: Predict the Next Word

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

given these words predict the next word



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Idea #1: Use a Fixed Window

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

given these two words predict the next word



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Idea #1: Use a Fixed Window

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

given these two words predict the next word

One-hot feature encoding: tells us what each word is

[1 0 0 0 0 0 1 0 0 0]
for a
prediction



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Problem #1: Can't Model Long-Term Dependencies

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



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



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Idea #2: Use Entire Sequence as Set of Counts

"This morning I took my cat for a"

↓
"bag of words"
[0 1 0 0 1 0 0 ... 0 0 1 1 0 0 0 1]
↓
prediction

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Problem #2: Counts Don't Preserve Order

The food was good, not bad at all.
vs.
The food was bad, not good at all.

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Idea #3: Use a Really Big Fixed Window

"This morning I took my cat for a walk"

given these words predict the next word
[1 0 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 1 0 ...]
morning I took this cat
↓
prediction

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Problem #3: No Parameter Sharing

[1 0 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 1 0 ...]
this morning took the cat

Each of these inputs has a separate parameter:

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018. 1/27/20

Problem #3: No Parameter Sharing

[1 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 1 0 ...]
this morning took the cat

Each of these inputs has a **separate parameter**:

[0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 1 ...]
this morning



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018 1/27/20

Problem #3: No Parameter Sharing

[1 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 1 0 ...]
this morning took the cat

Each of these inputs has a **separate parameter**:

[0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 1 ...]
this morning

Things we learn about the sequence **won't transfer** if they appear **elsewhere** in the sequence.



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018 1/27/20

Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018 1/27/20

Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences
2. Track **long-term** dependencies



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H. Suresh, 6.S191 2018 1/27/20

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**



6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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

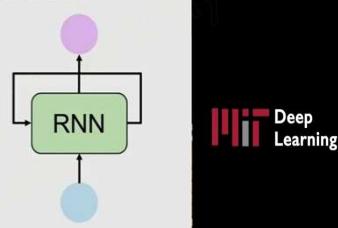


6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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

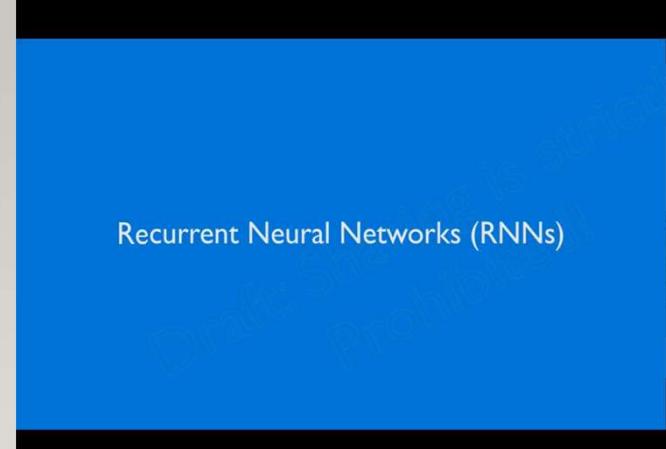


Today: **Recurrent Neural Networks (RNNs)** as an approach to sequence modeling problems

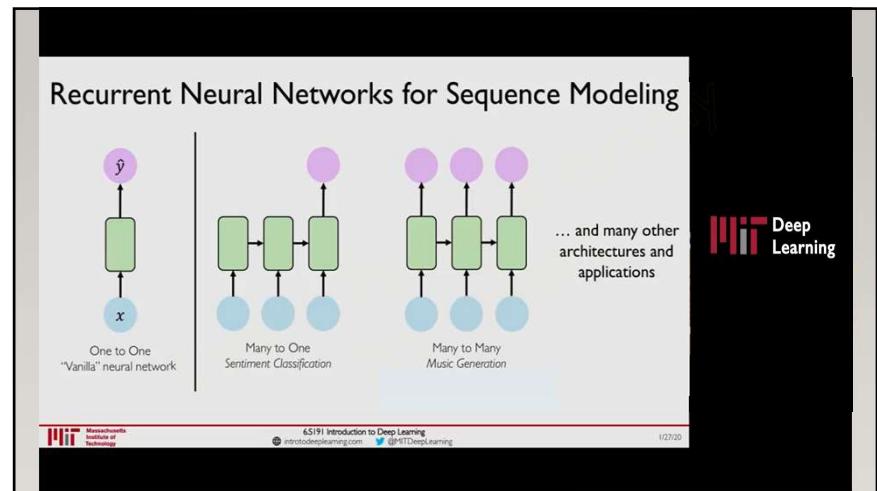
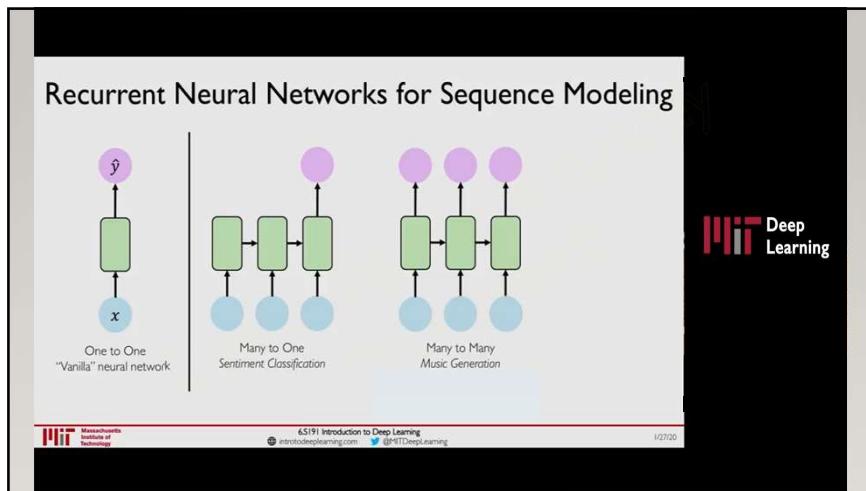
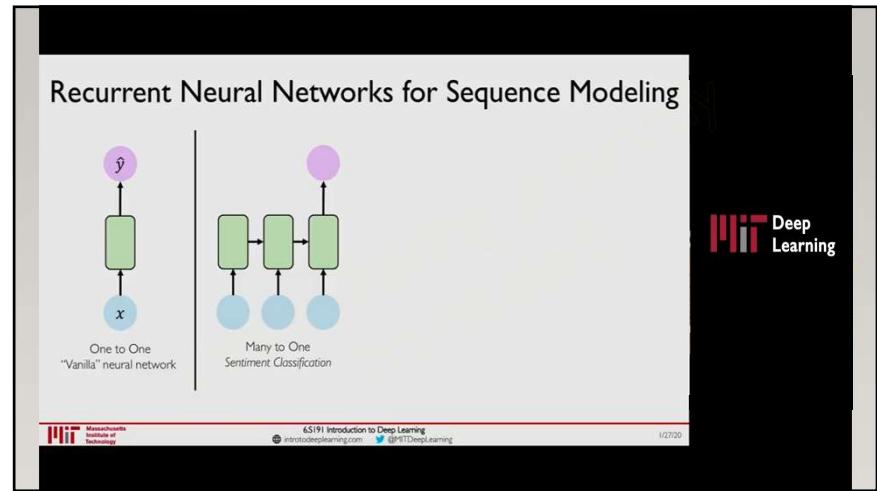
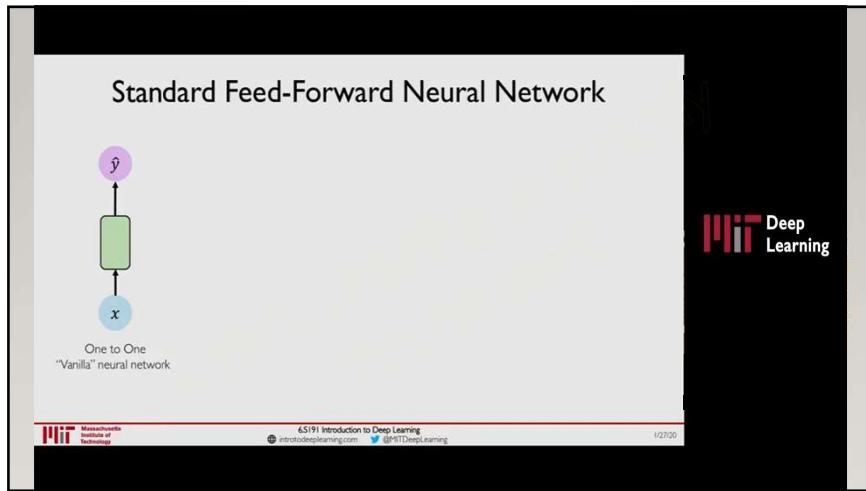


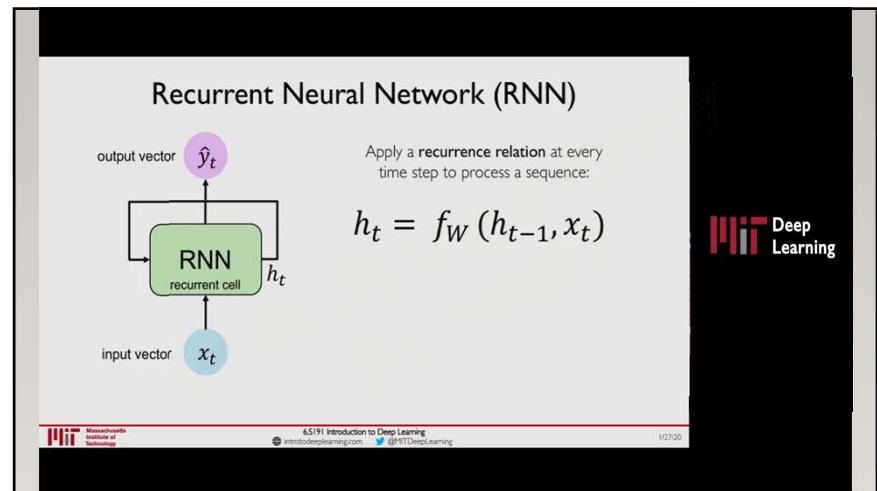
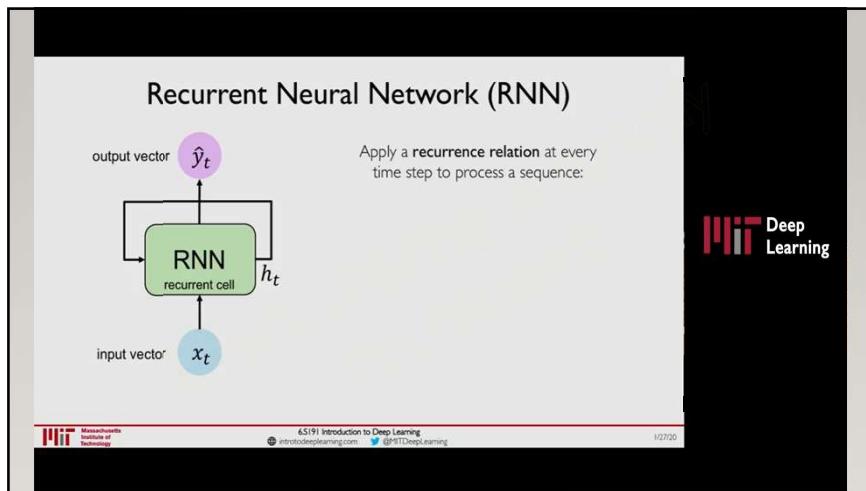
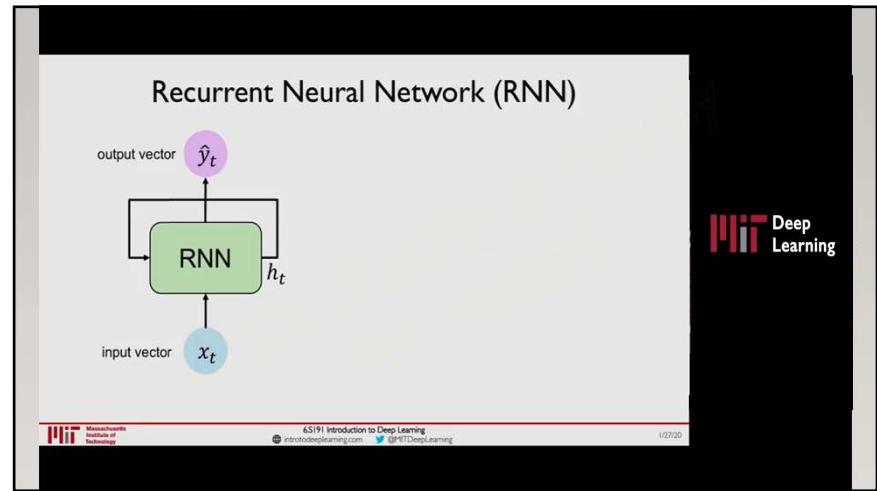
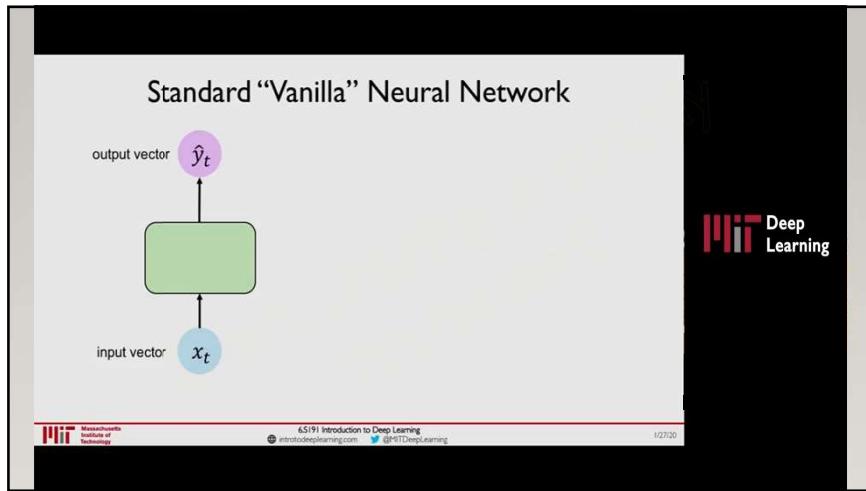
6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

Recurrent Neural Networks (RNNs)




6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20





Recurrent Neural Network (RNN)

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

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

cell state function parameterized by W

output vector \hat{y}_t
RNN recurrent cell
input vector x_t
 h_t

MIT Deep Learning

6.S191 Introduction to Deep Learning introdeeplearning.com @MITDeepLearning 1/27/20

Recurrent Neural Network (RNN)

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

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

cell state function parameterized by W

old state

output vector \hat{y}_t
RNN recurrent cell
input vector x_t
 h_t

MIT Deep Learning

6.S191 Introduction to Deep Learning introdeeplearning.com @MITDeepLearning 1/27/20

Recurrent Neural Network (RNN)

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

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

cell state function parameterized by W

old state

input vector at time step t

output vector \hat{y}_t
RNN recurrent cell
input vector x_t
 h_t

MIT Deep Learning

6.S191 Introduction to Deep Learning introdeeplearning.com @MITDeepLearning 1/27/20

Recurrent Neural Network (RNN)

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

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

cell state function parameterized by W

old state

input vector at time step t

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

output vector \hat{y}_t
RNN recurrent cell
input vector x_t
 h_t

MIT Deep Learning

6.S191 Introduction to Deep Learning introdeeplearning.com @MITDeepLearning 1/27/20

RNN Intuition

```

my_rnn = RNN()
hidden_state = [0, 0, 0, 0]

sentence = ["I", "love", "recurrent", "neural"]

for word in sentence:
    prediction, hidden_state = my_rnn(word, hidden_state)

next_word_prediction = prediction
# >>> "networks!"

```

MIT Deep Learning

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

RNN Intuition

```

my_rnn = RNN()
hidden_state = [0, 0, 0, 0]

sentence = ["I", "love", "recurrent", "neural"]

for word in sentence:
    prediction, hidden_state = my_rnn(word, hidden_state)

next_word_prediction = prediction
# >>> "networks!"

```

MIT Deep Learning

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

RNN Intuition

```

my_rnn = RNN()
hidden_state = [0, 0, 0, 0]

sentence = ["I", "love", "recurrent", "neural"]

for word in sentence:
    prediction, hidden_state = my_rnn(word, hidden_state)

next_word_prediction = prediction
# >>> "networks!"

```

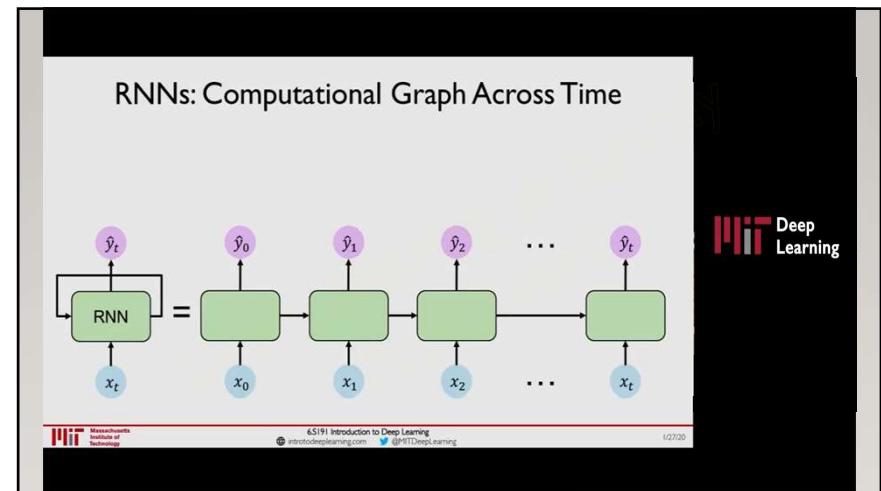
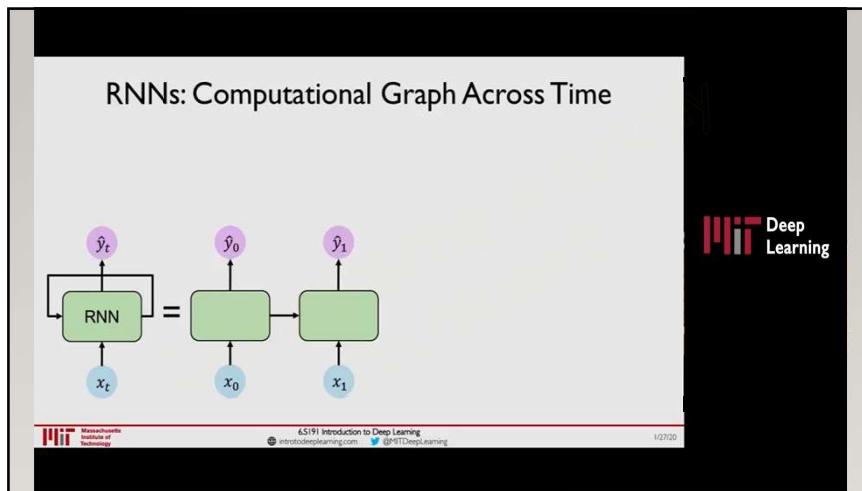
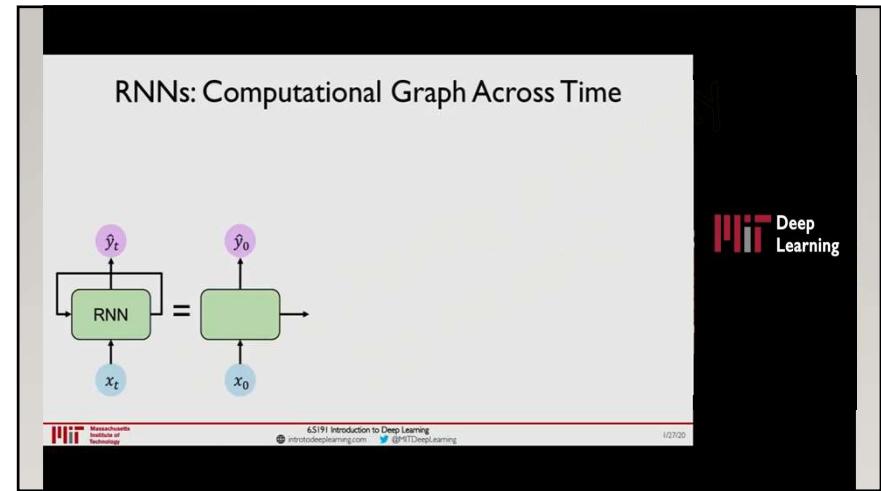
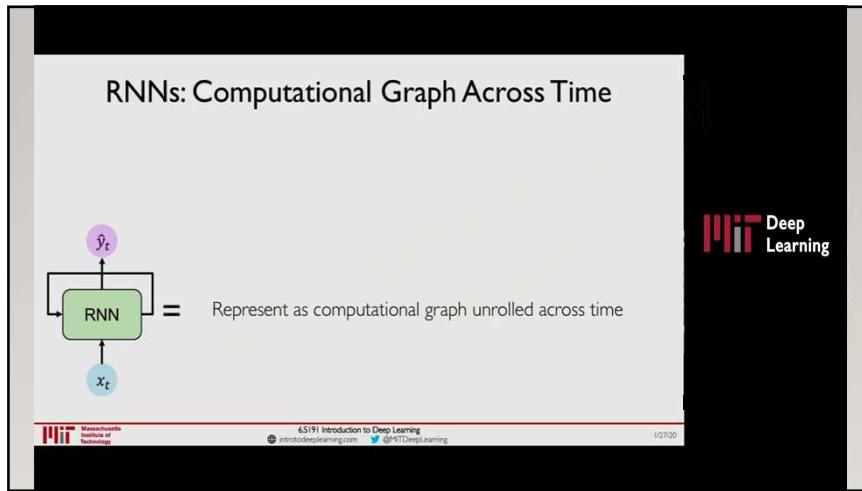
MIT Deep Learning

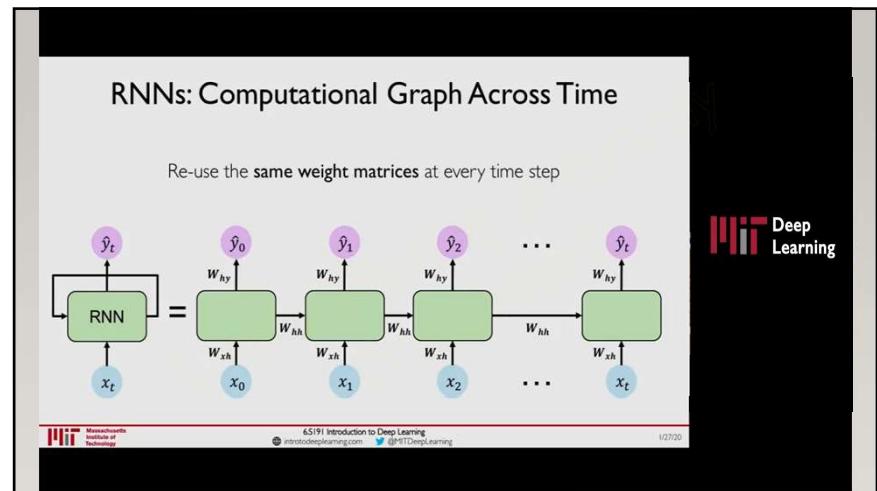
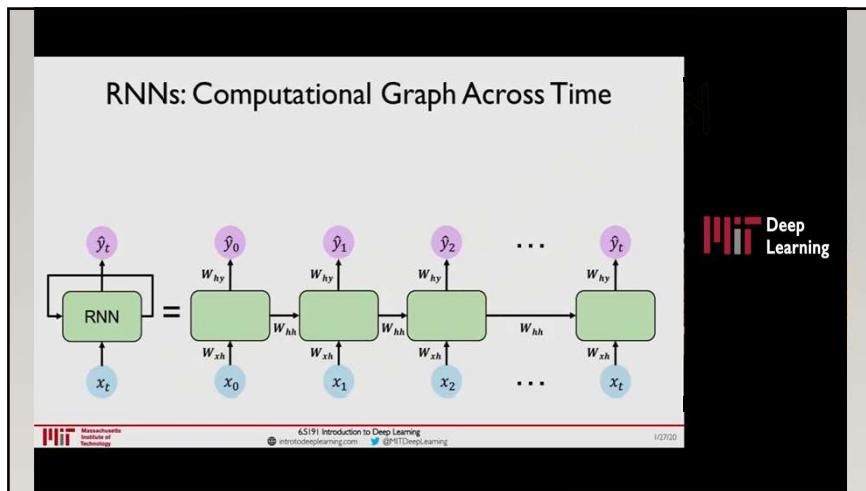
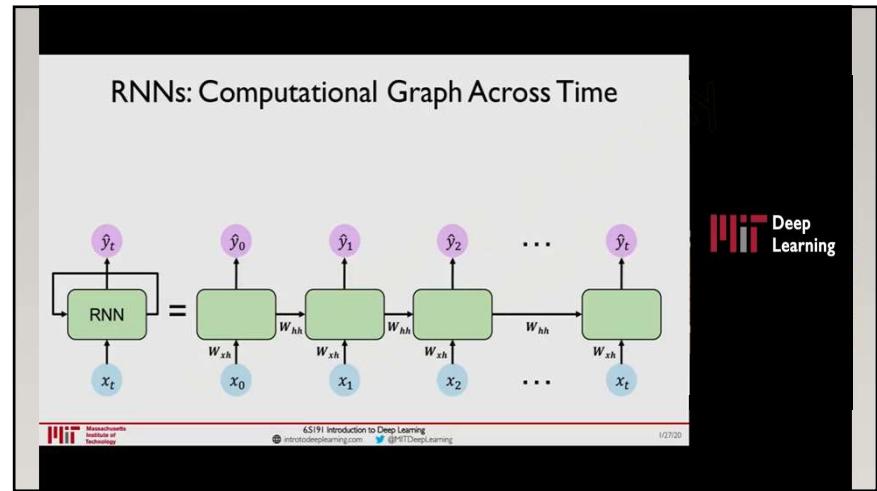
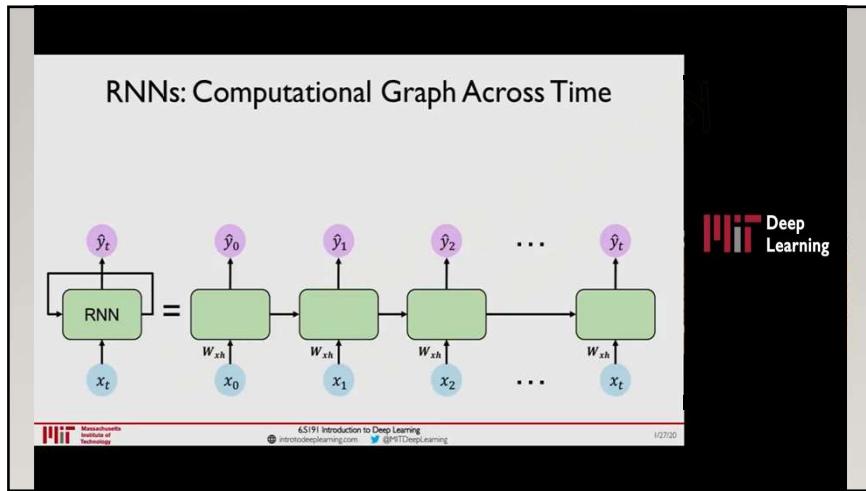
6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

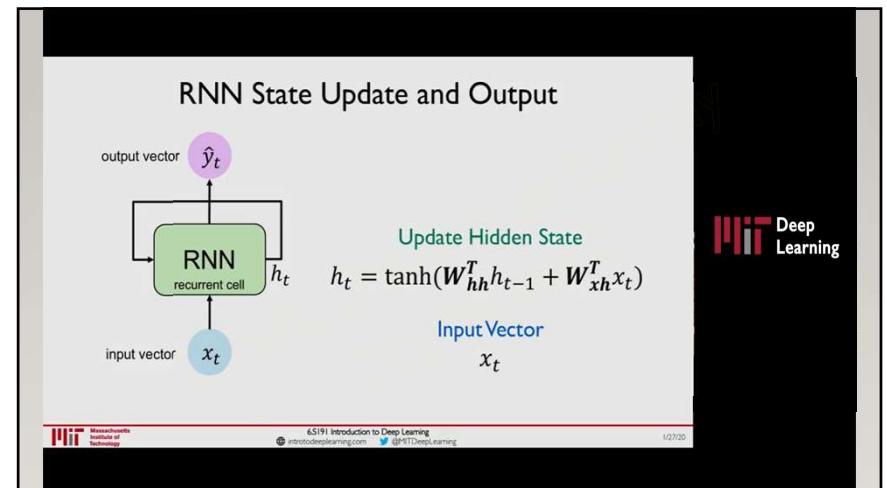
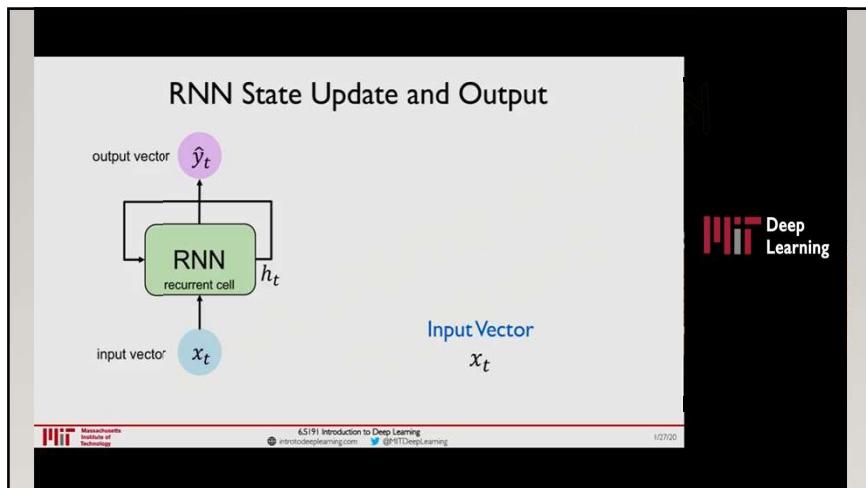
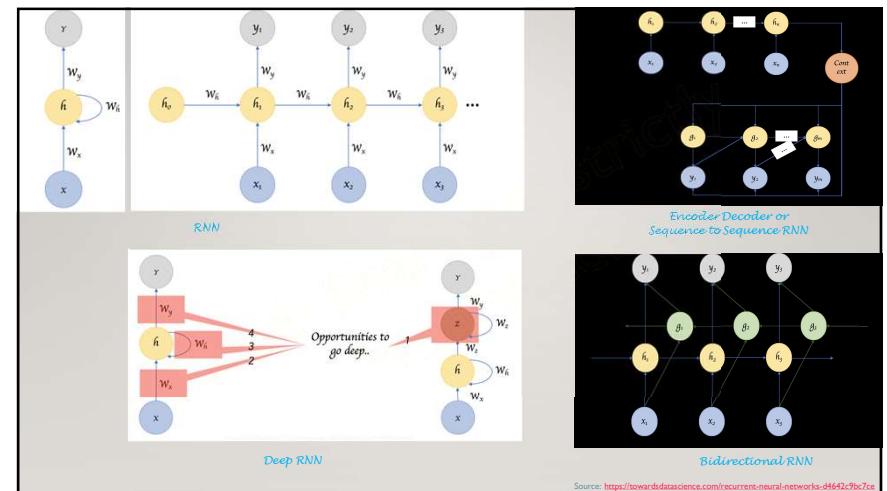
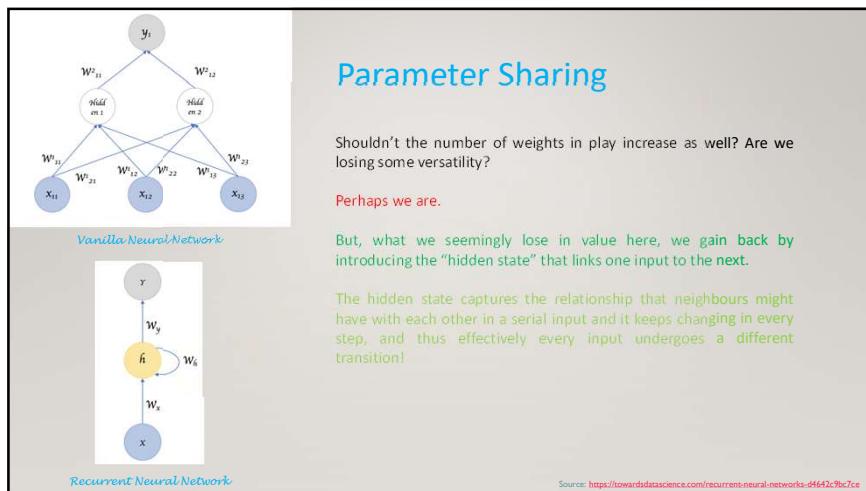
RNNs: Computational Graph Across Time

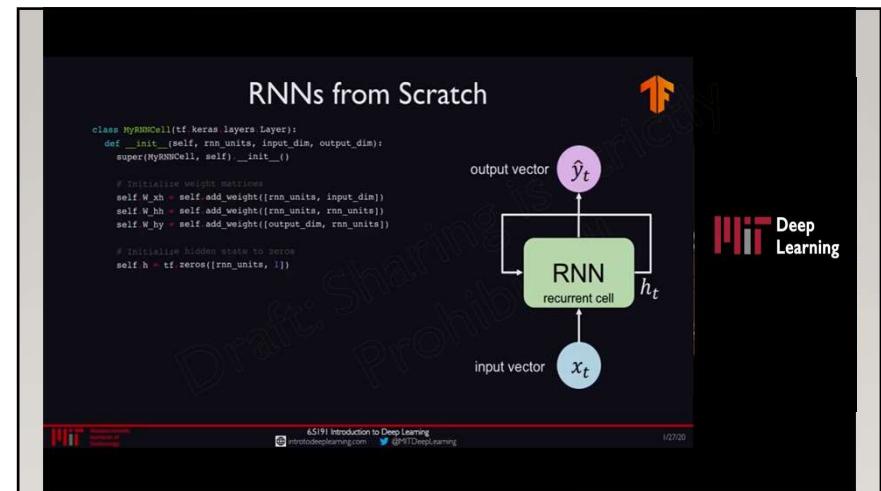
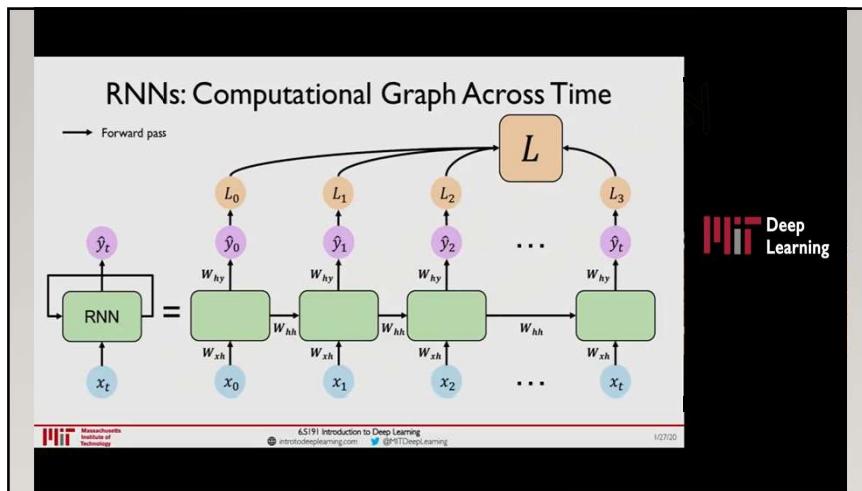
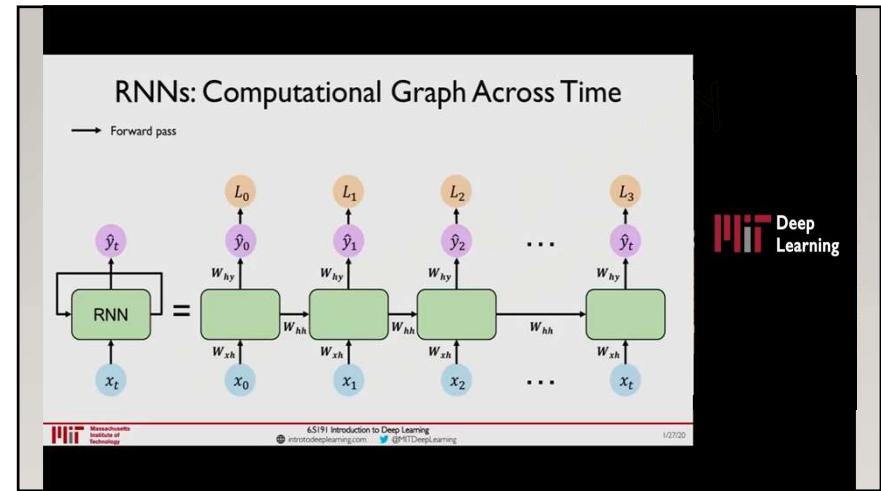
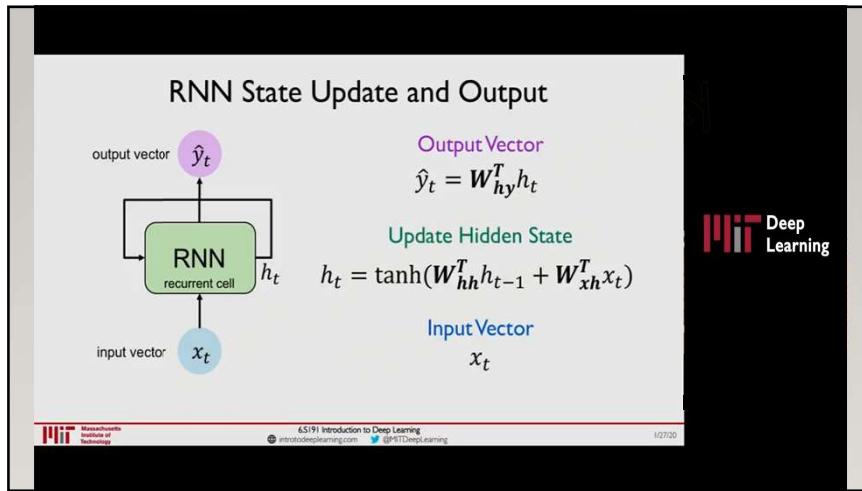
MIT Deep Learning

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20









RNNs from Scratch

```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_by = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh(self.W_hh * self.h + self.W_xh * x)

        # Compute the output
        output = self.W_by * self.h

        # Return the current output and hidden state
        return output, self.h
```

6.519 Introduction to Deep Learning
introdeeplearning.com @MITDeepLearning

1/27/20

RNNs from Scratch

```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_by = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh(self.W_hh * self.h + self.W_xh * x)

        # Compute the output
        output = self.W_by * self.h

        # Return the current output and hidden state
        return output, self.h
```

6.519 Introduction to Deep Learning
introdeeplearning.com @MITDeepLearning

1/27/20

RNNs from Scratch

```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_by = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh(self.W_hh * self.h + self.W_xh * x)

        # Compute the output
        output = self.W_by * self.h

        # Return the current output and hidden state
        return output, self.h
```

6.519 Introduction to Deep Learning
introdeeplearning.com @MITDeepLearning

1/27/20

RNNs from Scratch

```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_by = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

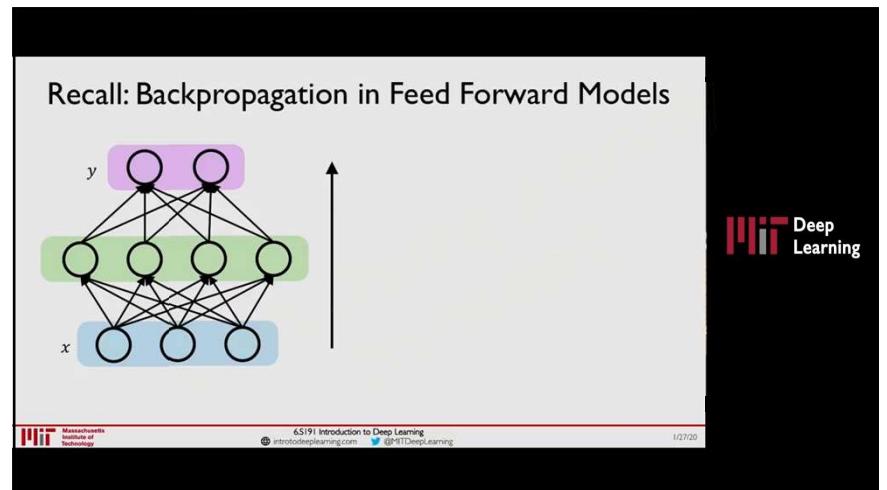
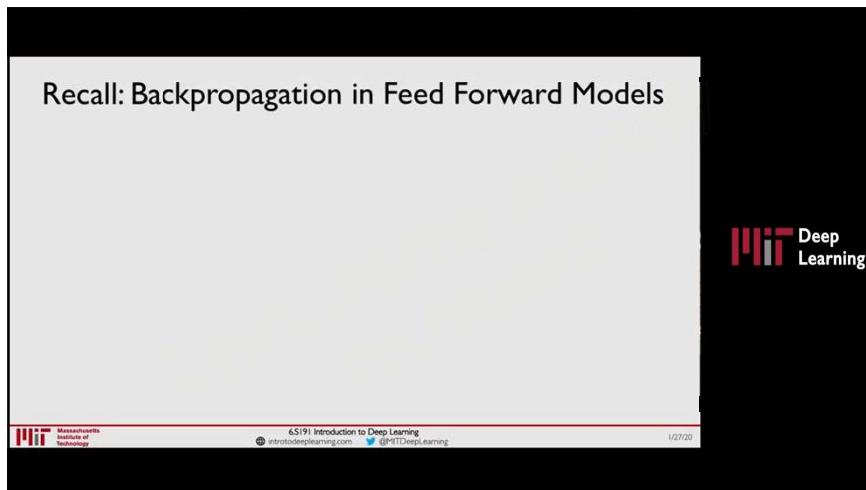
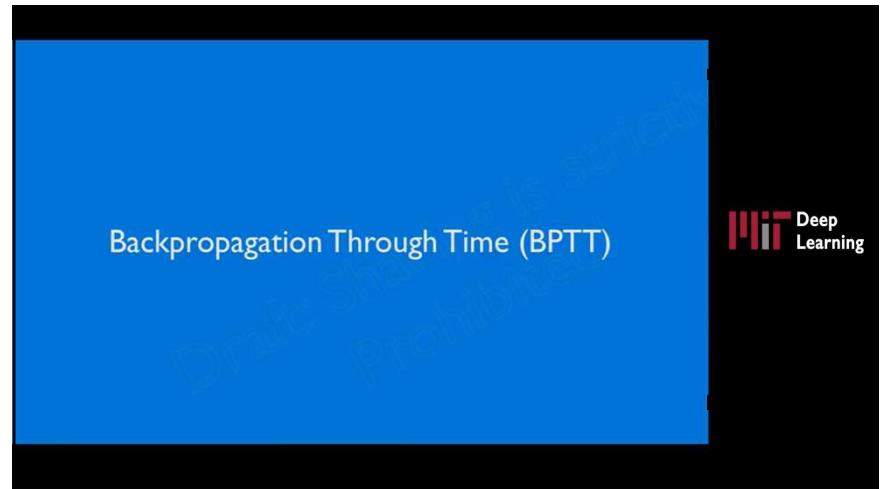
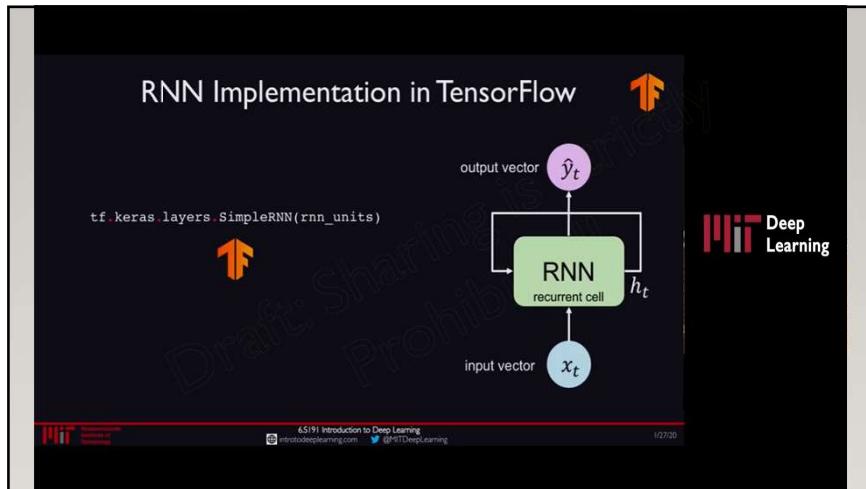
    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh(self.W_hh * self.h + self.W_xh * x)

        # Compute the output
        output = self.W_by * self.h

        # Return the current output and hidden state
        return output, self.h
```

6.519 Introduction to Deep Learning
introdeeplearning.com @MITDeepLearning

1/27/20



Recall: Backpropagation in Feed Forward Models

MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

Recall: Backpropagation in Feed Forward Models

Backpropagation algorithm:

MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

Recall: Backpropagation in Feed Forward Models

Backpropagation algorithm:

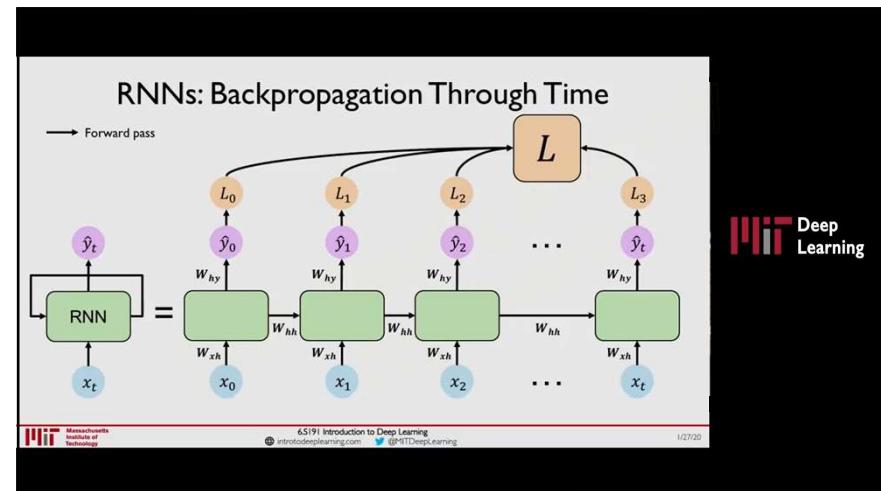
- Take the derivative (gradient) of the loss with respect to each parameter
- Shift parameters in order to minimize loss

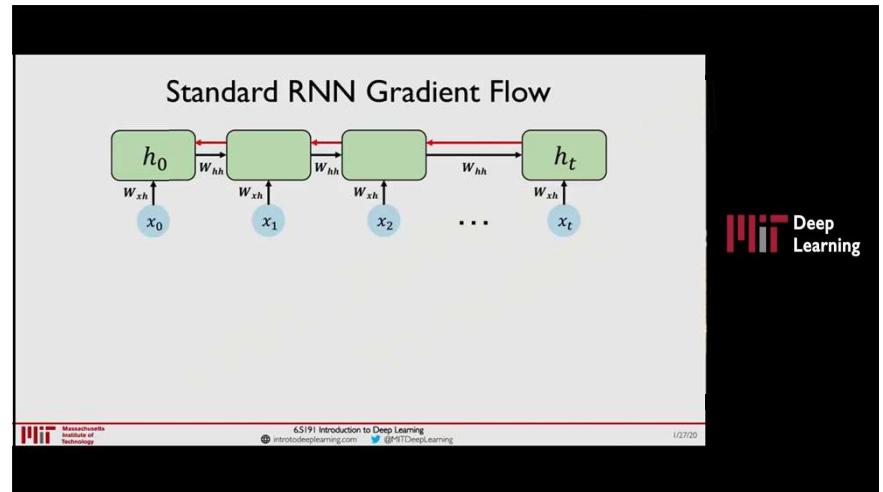
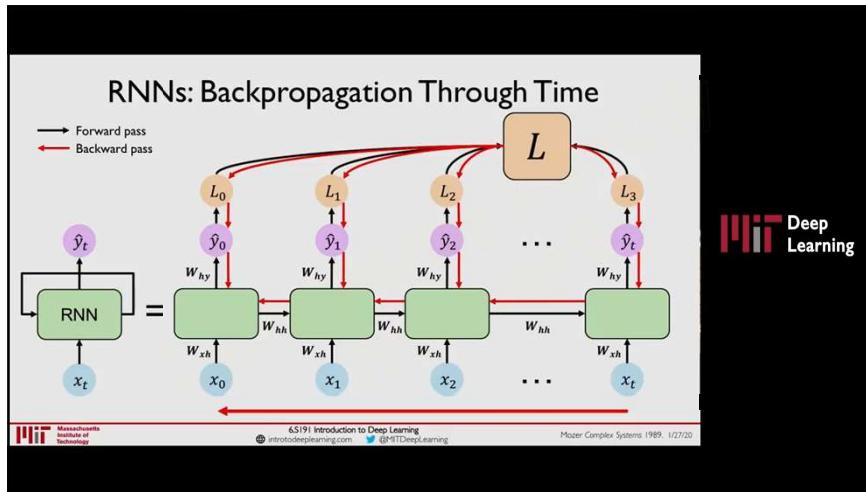
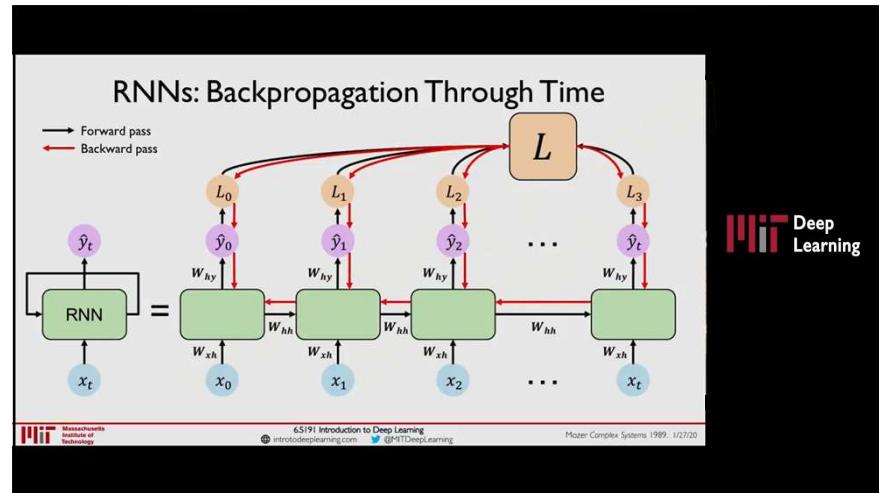
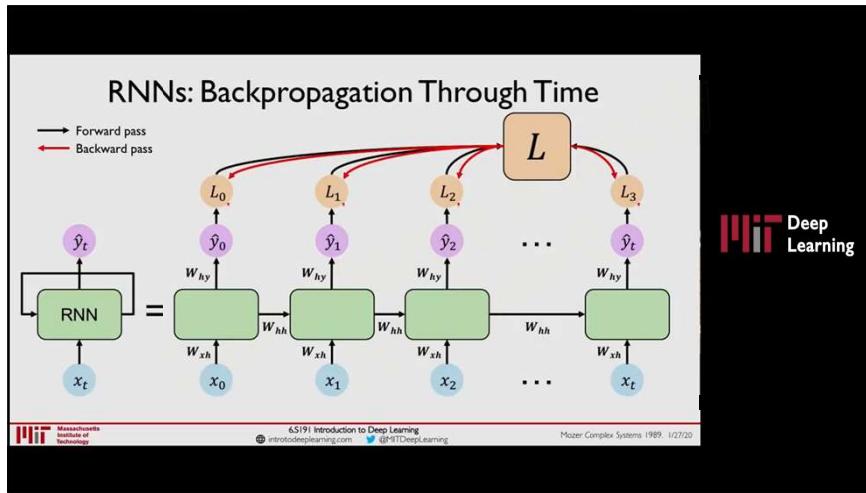
MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20





Standard RNN Gradient Flow

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning 1/27/20

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning 1/27/20

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 prevent overflow

Many values < 1: vanishing gradients

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning 1/27/20

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 prevent overflow

Many values < 1: vanishing gradients

- Activation function
- Weight initialization
- Network architecture

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning 1/27/20

The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

MIT Deep Learning

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

Multiply many small numbers together

MIT Deep Learning

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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

MIT Deep Learning

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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

MIT Deep Learning

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

The Problem of Long-Term Dependencies

"The clouds are in the ____"

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning
1/27/20

The Problem of Long-Term Dependencies

"The clouds are in the ____"

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning
1/27/20

The Problem of Long-Term Dependencies

"The clouds are in the ____"

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning
1/27/20

The Problem of Long-Term Dependencies

"The clouds are in the ____"

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

MIT Deep Learning

Massachusetts Institute of Technology
6.S191 Introduction to Deep Learning
introtodeeplearning.com @MITDeepLearning
1/27/20

Trick #1: Activation Functions

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

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H.Suresh, 6.S191 2018. 1/27/20

Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

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

Initialize **biases** to zero

This helps prevent the weights from shrinking to zero.

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H.Suresh, 6.S191 2018. 1/27/20

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.

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H.Suresh, 6.S191 2018. 1/27/20

Solution #3: Gated Cells

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

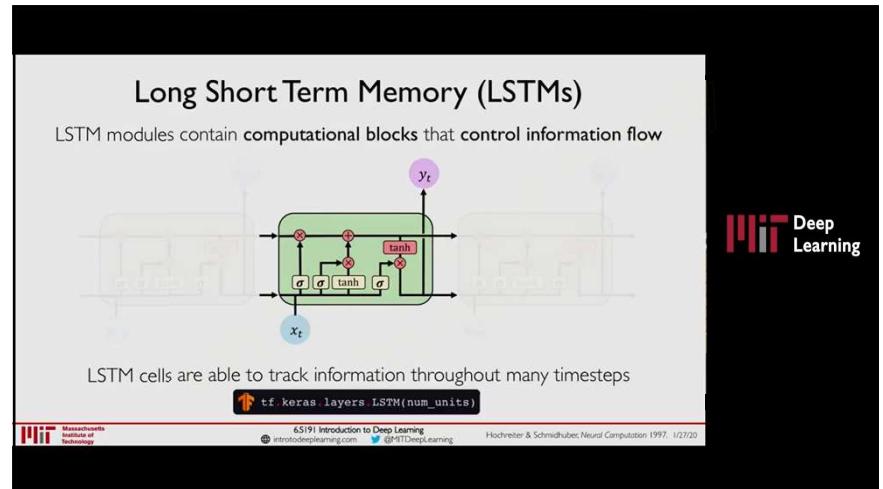
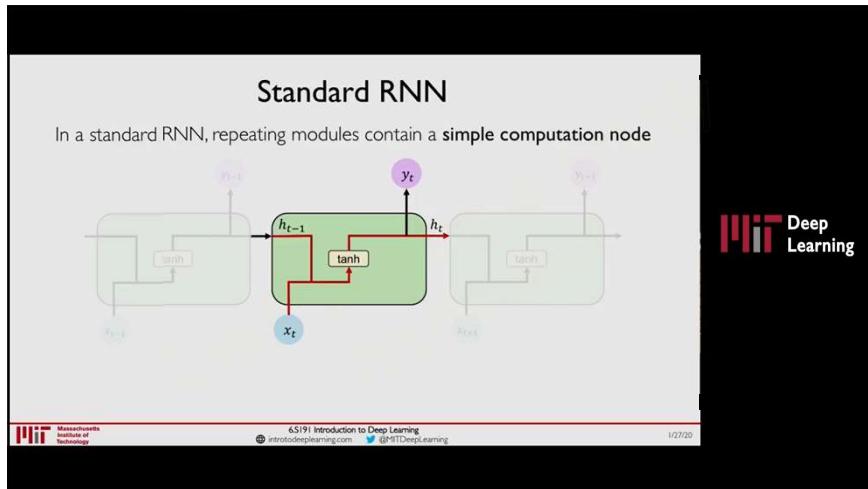
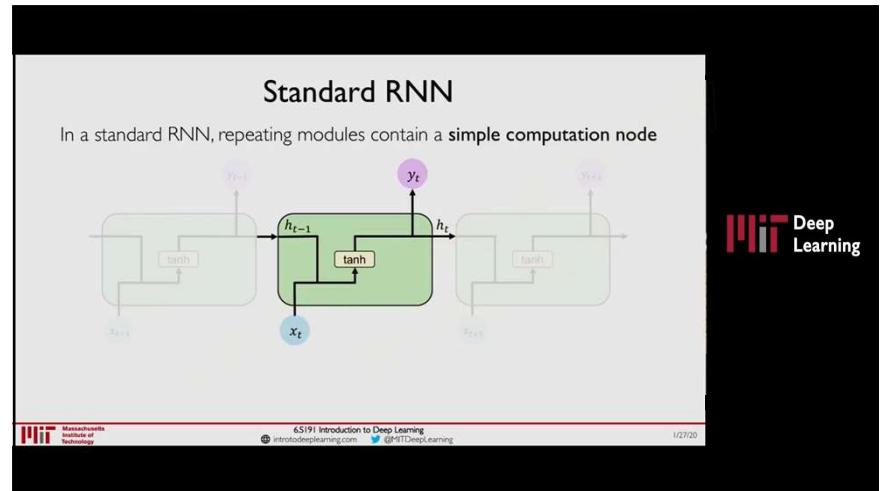
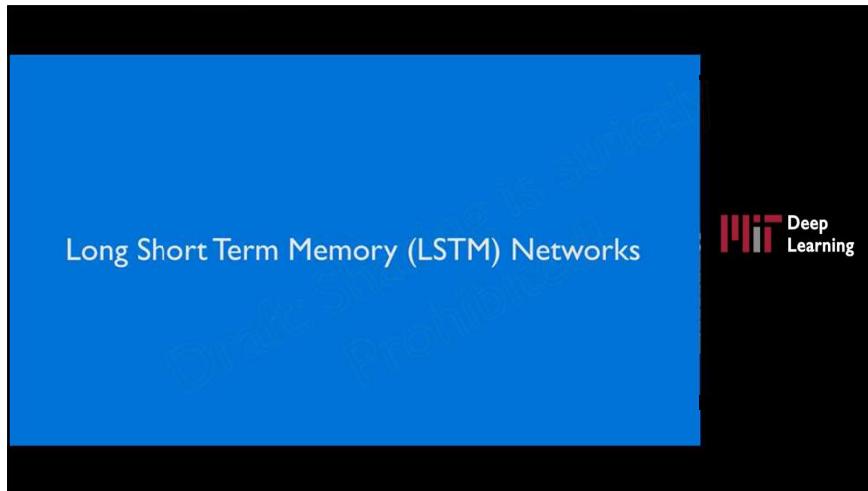
gated cell
LSTM, GRU, etc.

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

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning H.Suresh, 6.S191 2018. 1/27/20



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

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

Long Short Term Memory (LSTMs)

How do LSTMs work?

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

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning Hochreiter & Schmidhuber: Neural Computation 1997, 1/21/20

Long Short Term Memory (LSTMs)

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

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

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning Olah, "Understanding LSTMs", 1/27/20

Long Short Term Memory (LSTMs)

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

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

MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning Olah, "Understanding LSTMs", 1/27/20

Long Short Term Memory (LSTMs)

1) Forget 2) Store 3) Update 4) Output
LSTMs selectively update cell state values

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning Olah, "Understanding LSTMs", 1/27/20

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

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning Olah, "Understanding LSTMs", 1/27/20

Long Short Term Memory (LSTMs)

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

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

LSTM Gradient Flow

Uninterrupted gradient flow!

MIT Deep Learning

Massachusetts Institute of Technology

6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

LSTMs: Key Concepts



MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

LSTMs: Key Concepts

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



MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

LSTMs: Key Concepts

1. Maintain a **separate cell state** from what is outputted
2. Use **gates** to control the **flow of information**



MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

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



MIT Deep Learning

Massachusetts Institute of Technology

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

1/27/20

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



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

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



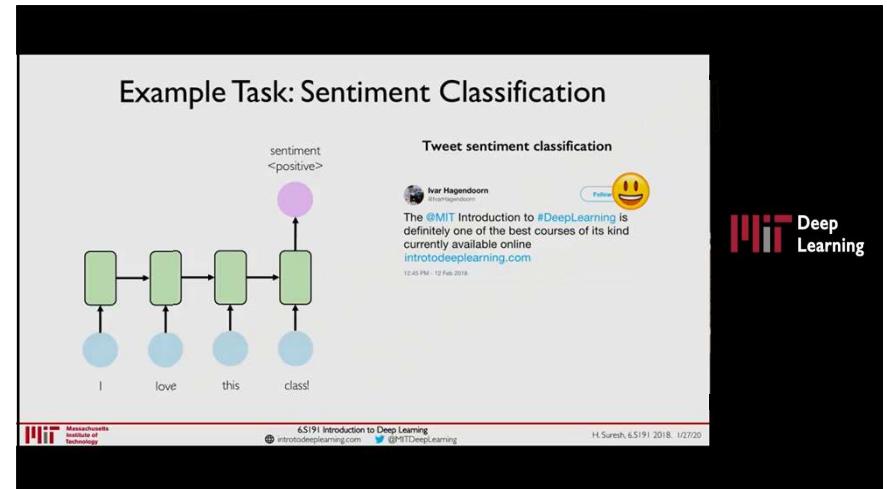
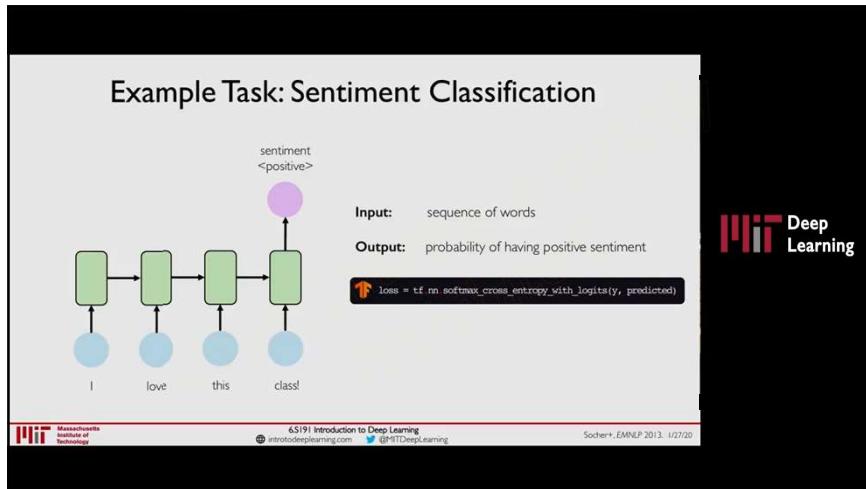
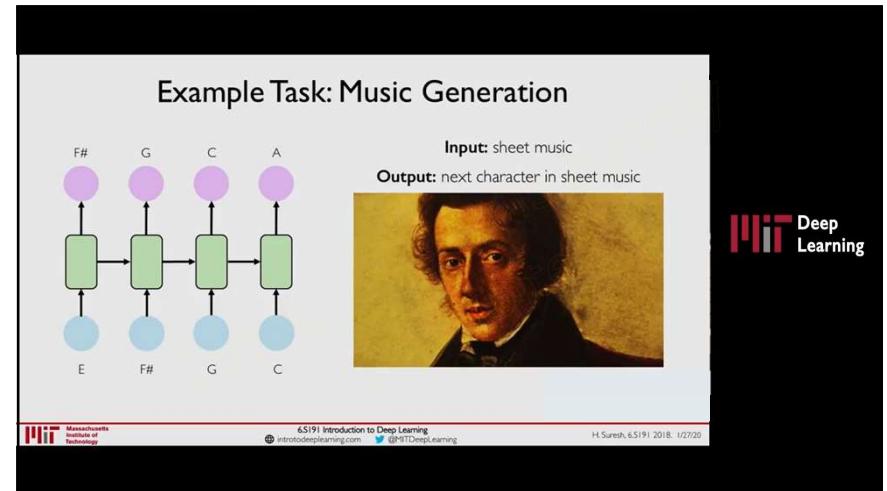
MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

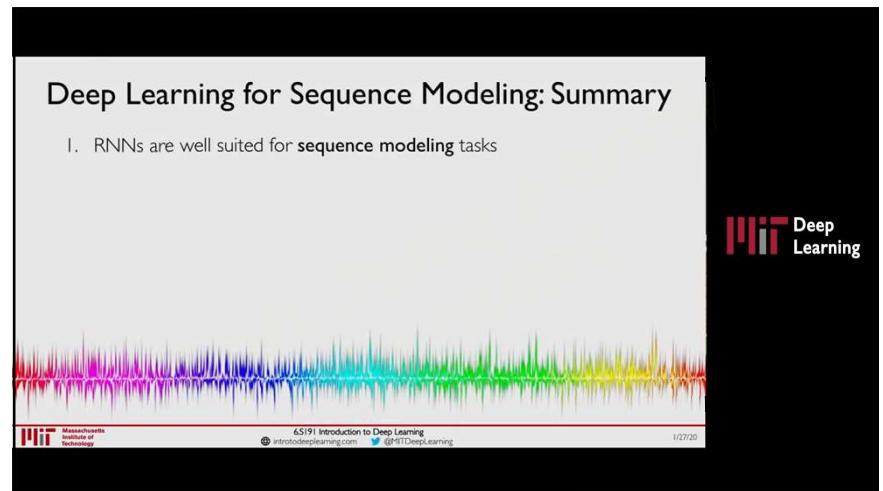
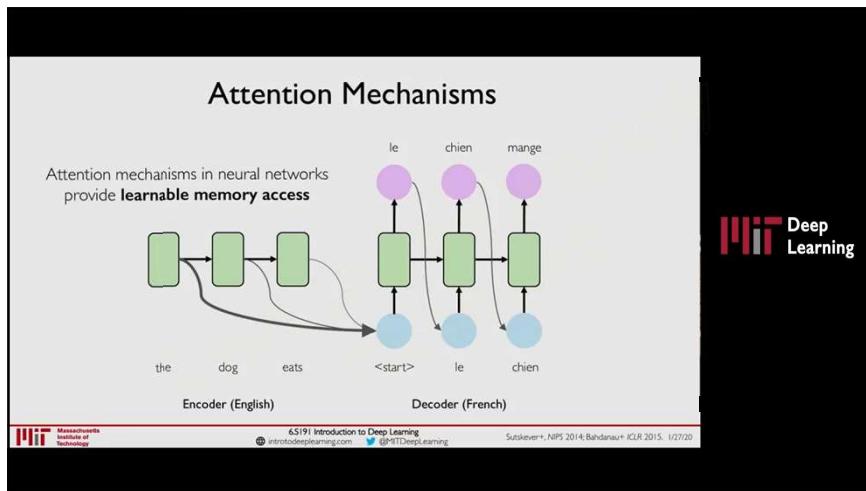
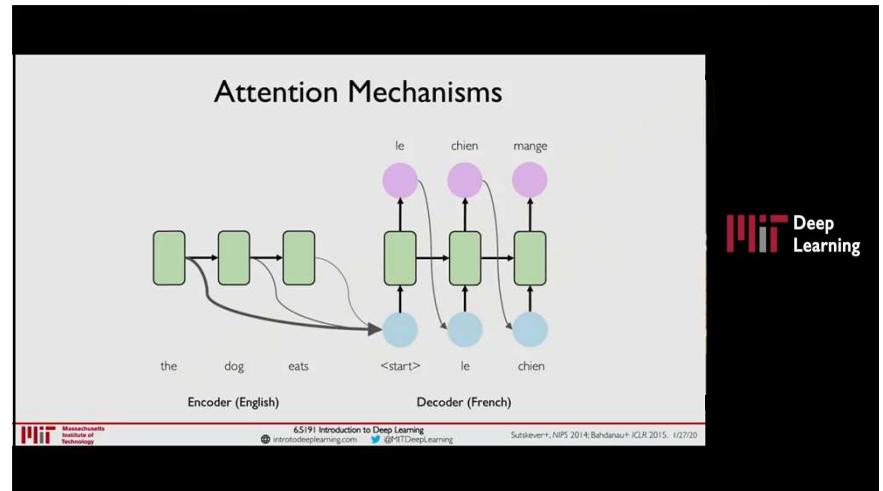
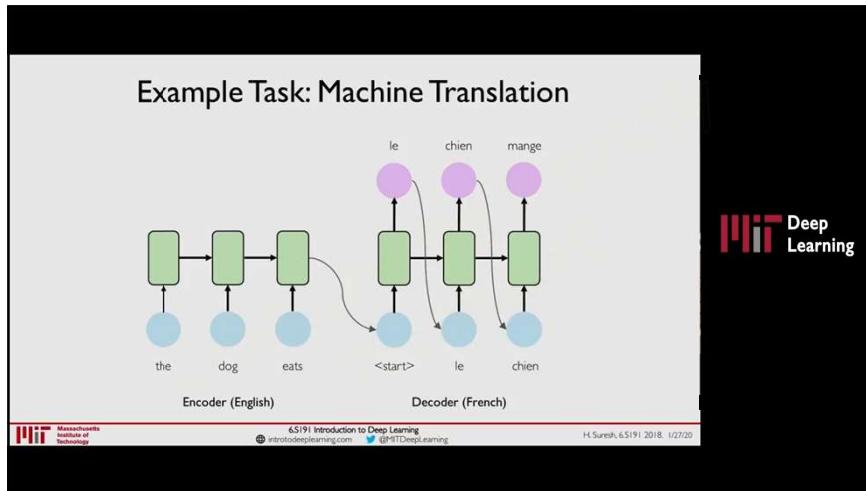
LSTMs: Key Concepts

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**



MIT Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20





Deep Learning for Sequence Modeling: Summary

- 1. RNNs are well suited for **sequence modeling** tasks
- 2. Model sequences via a **recurrence relation**



Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

Deep Learning for Sequence Modeling: Summary

- 1. RNNs are well suited for **sequence modeling** tasks
- 2. Model sequences via a **recurrence relation**
- 3. Training RNNs with **backpropagation through time**



Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

Deep Learning for Sequence Modeling: Summary

- 1. RNNs are well suited for **sequence modeling** tasks
- 2. Model sequences via a **recurrence relation**
- 3. Training RNNs with **backpropagation through time**
- 4. Gated cells like **LSTMs** let us model **long-term dependencies**



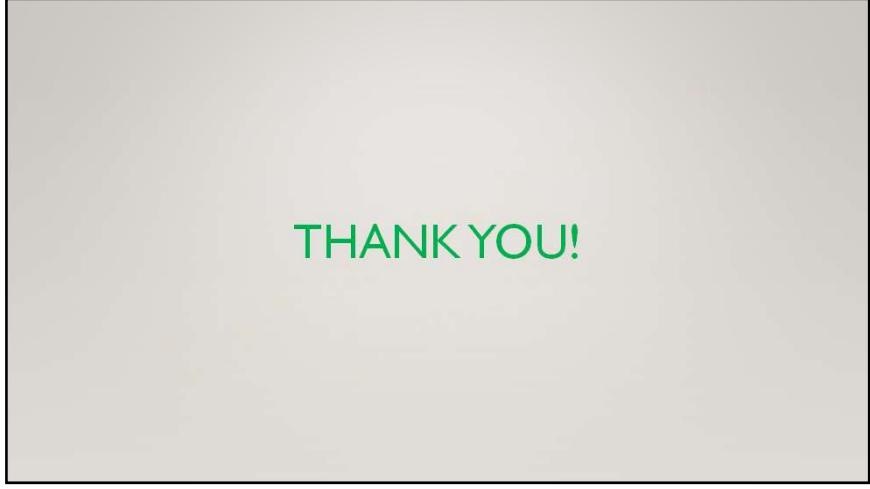
Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20

Deep Learning for Sequence Modeling: Summary

- 1. RNNs are well suited for **sequence modeling** tasks
- 2. Model sequences via a **recurrence relation**
- 3. Training RNNs with **backpropagation through time**
- 4. Gated cells like **LSTMs** let us model **long-term dependencies**
- 5. Models for **music generation**, classification, machine translation, and more



Massachusetts Institute of Technology 6.S191 Introduction to Deep Learning introtodeeplearning.com @MITDeepLearning 1/27/20



THANK YOU!