

Lecture 7: Recurrent Neural Networks

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 1

April 23, 2023

Training “Feedforward” Neural Networks

- 1. One time setup:** activation functions, preprocessing, weight initialization, regularization, gradient checking
- 2. Training dynamics:** babysitting the learning process, parameter updates, hyperparameter optimization
- 3. Evaluation:** model ensembles, test-time augmentation, transfer learning

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 2

April 23, 2023

More complexity...



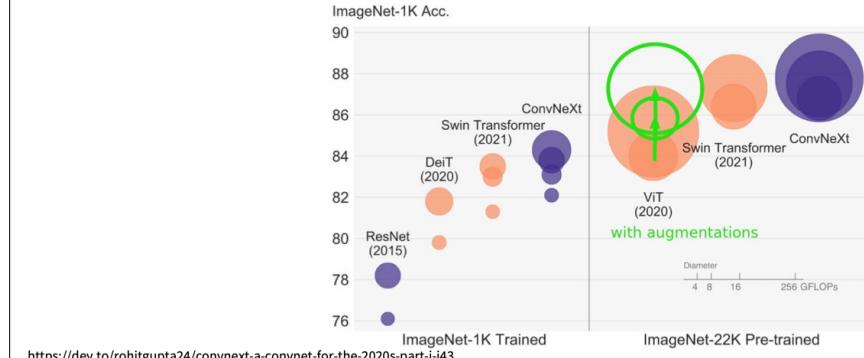
Liu, Zhuang, et al. "A convnet for the 2020s." CVPR 2022.

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 10

April 23, 2023

More complexity...



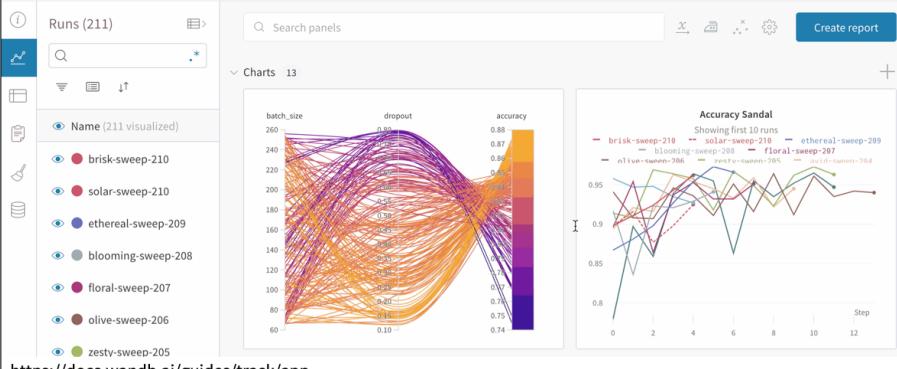
<https://dev.to/rohitgupta24/convnext-a-convnet-for-the-2020s-part-i-43>

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 11

April 23, 2023

Evaluate models and tune hyperparameters



<https://docs.wandb.ai/guides/track/app>

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 12 April 23, 2023

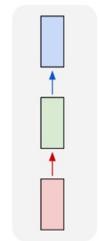
Today: Recurrent Neural Networks

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 13 April 23, 2023

“Vanilla” Neural Network

one to one



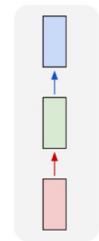
Vanilla Neural Networks

Fei-Fei Li, Ehsan Adeli

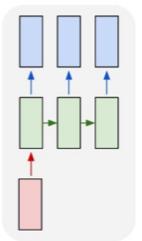
Lecture 8 - 14 April 23, 2023

Recurrent Neural Networks: Process Sequences

one to one



one to many

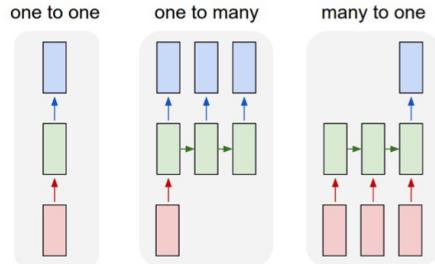


e.g. Image Captioning
image -> sequence of words

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 15 April 23, 2023

Recurrent Neural Networks: Process Sequences



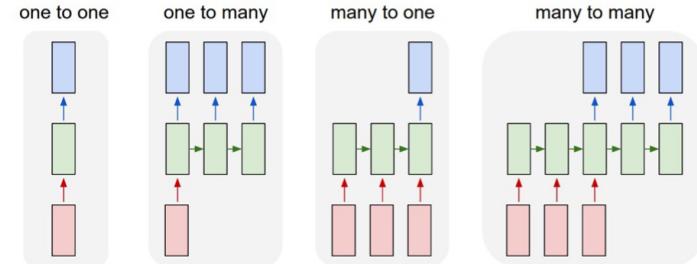
e.g. action prediction
sequence of video frames -> action class

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 16

April 23, 2023

Recurrent Neural Networks: Process Sequences



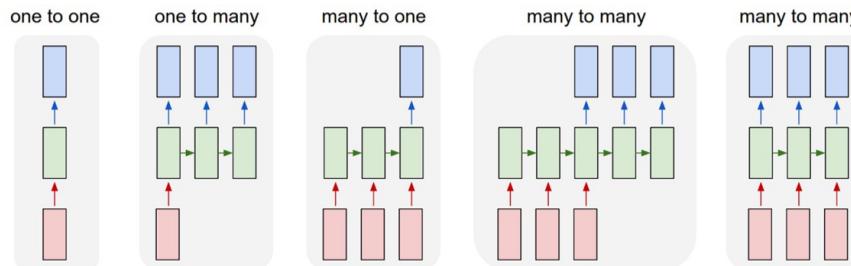
E.g. Video Captioning
Sequence of video frames -> caption

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 17

April 23, 2023

Recurrent Neural Networks: Process Sequences



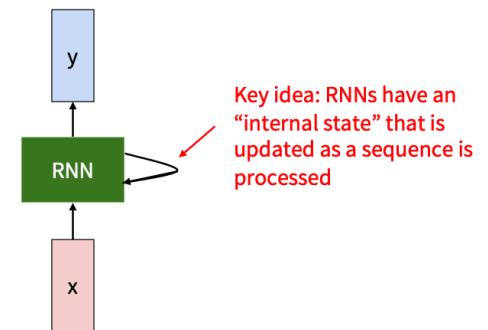
e.g. Video classification on frame level

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 18

April 23, 2023

Recurrent Neural Network

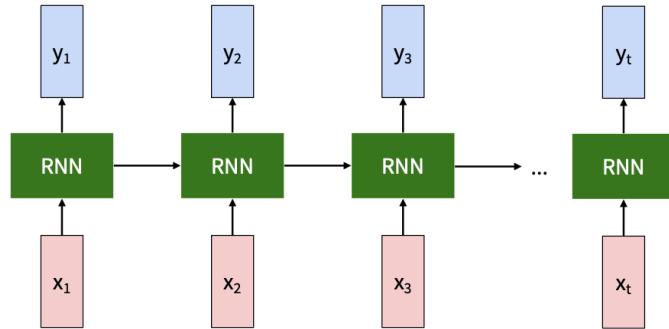


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 22

April 23, 2023

Unrolled RNN



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 23

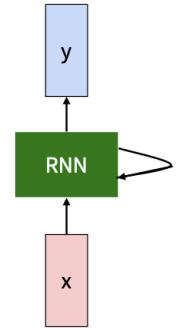
April 23, 2023

RNN hidden state update

We can process a sequence of vectors x by applying a recurrence formula at every time step:

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

new state old state input vector at some time step
 some function with parameters W



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 24

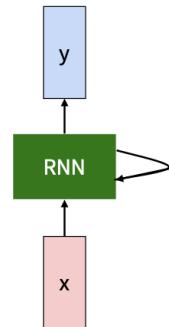
April 23, 2023

RNN output generation

We can process a sequence of vectors x by applying a recurrence formula at every time step:

$$y_t = f_{W_{hy}}(h_t)$$

output new state
 another function with parameters W_{hy}

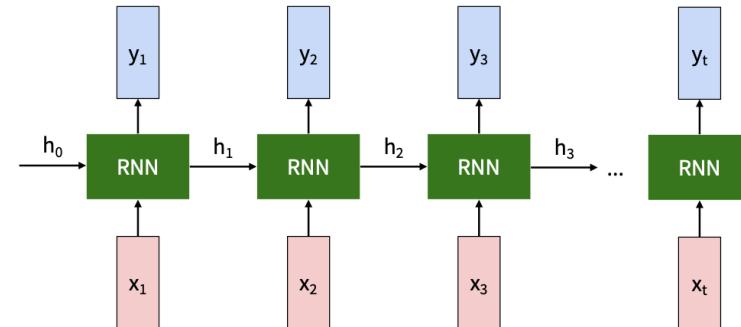


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 25

April 23, 2023

Recurrent Neural Network



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 26

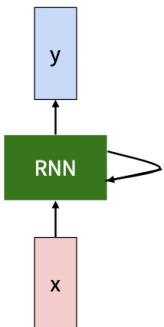
April 23, 2023

Recurrent Neural Network

We can process a sequence of vectors x by applying a recurrence formula at every time step:

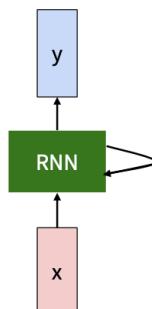
$$h_t = f_W(h_{t-1}, x_t)$$

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



(Vanilla) Recurrent Neural Network

The state consists of a single “hidden” vector h :



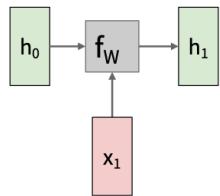
$$h_t = f_W(h_{t-1}, x_t)$$

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

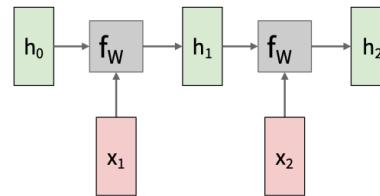
$$y_t = W_{hy}h_t$$

Sometimes called a “Vanilla RNN” or an “Elman RNN” after Prof. Jeffrey Elman

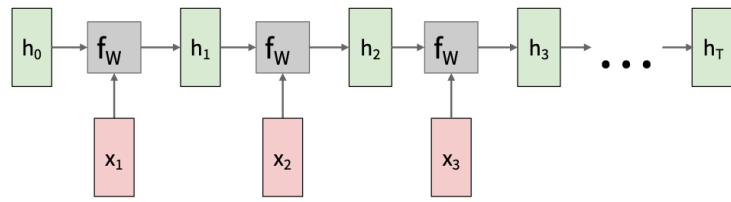
RNN: Computational Graph



RNN: Computational Graph



RNN: Computational Graph

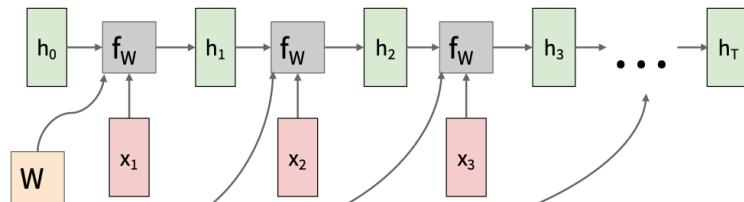


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 31 April 23, 2023

RNN: Computational Graph

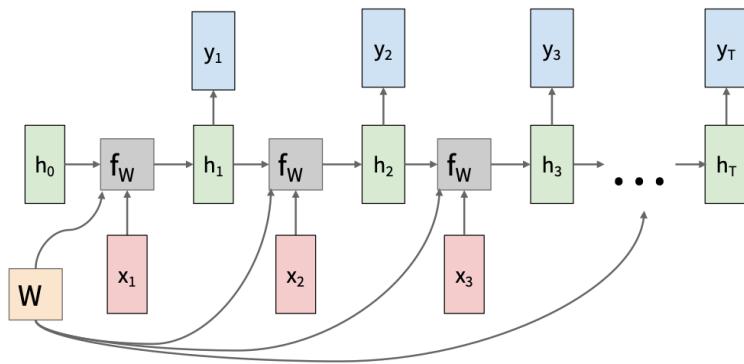
Re-use the same weight matrix at every time-step



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 32 April 23, 2023

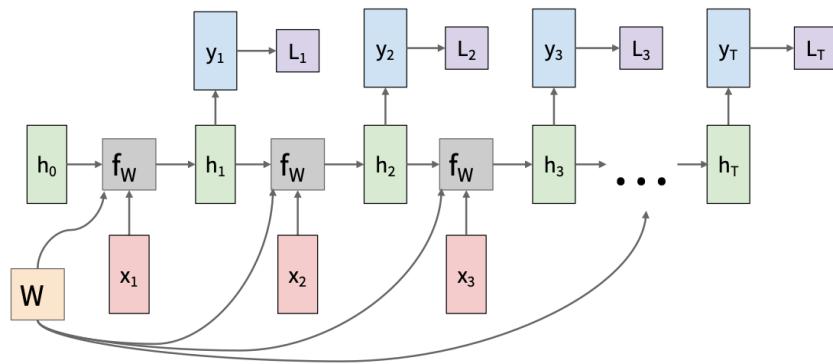
RNN: Computational Graph: Many to Many



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 33 April 23, 2023

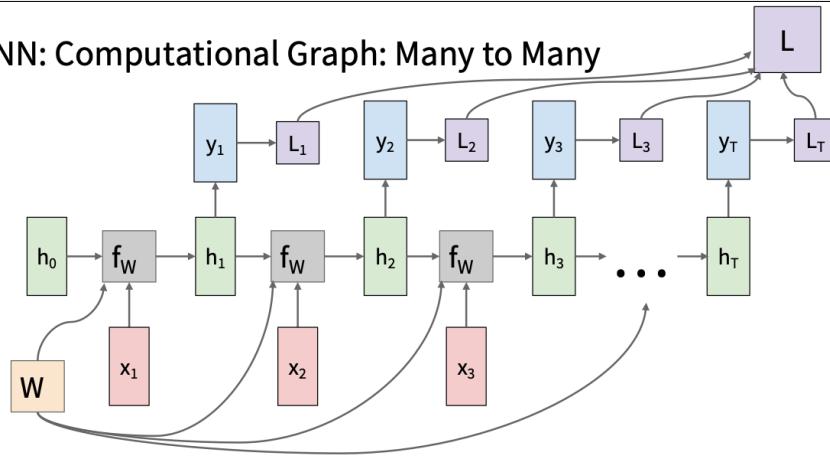
RNN: Computational Graph: Many to Many



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 34 April 23, 2023

RNN: Computational Graph: Many to Many

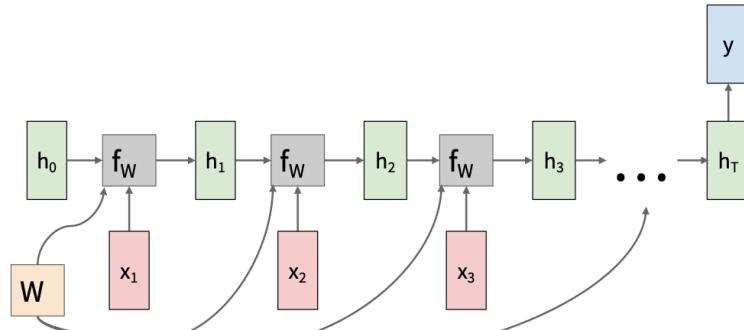


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 35

April 23, 2023

RNN: Computational Graph: Many to One

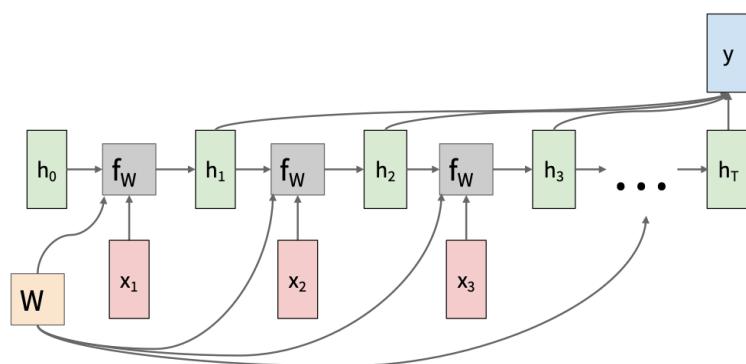


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 36

April 23, 2023

RNN: Computational Graph: Many to One

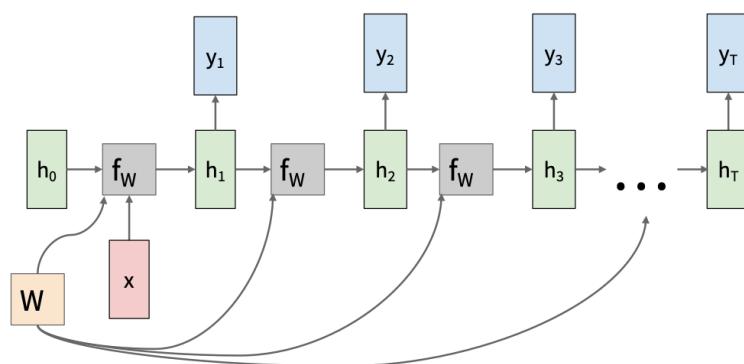


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 37

April 23, 2023

RNN: Computational Graph: One to Many

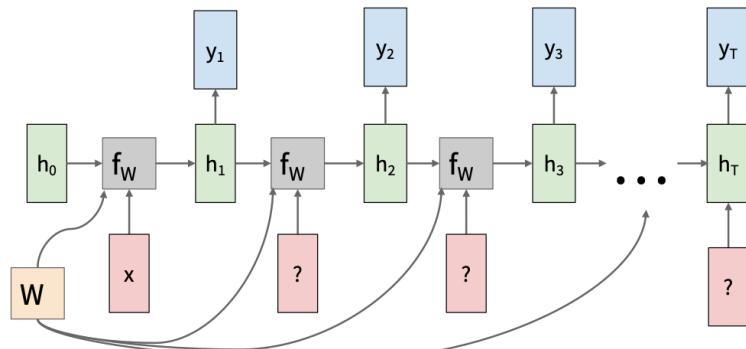


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 38

April 23, 2023

RNN: Computational Graph: One to Many

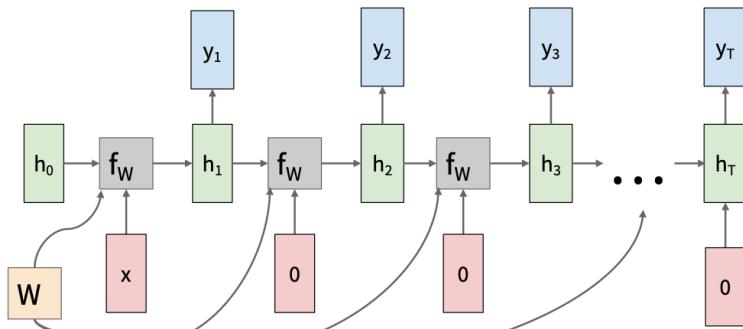


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 39

April 23, 2023

RNN: Computational Graph: One to Many

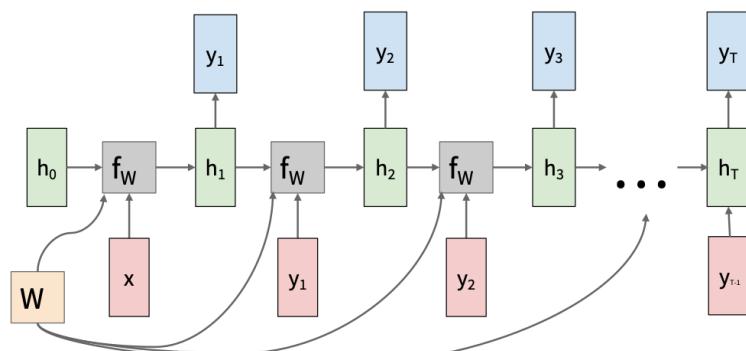


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 40

April 23, 2023

RNN: Computational Graph: One to Many



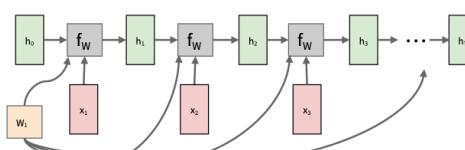
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 41

April 23, 2023

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector



Sutskever et al., "Sequence to Sequence Learning with Neural Networks", NIPS 2014

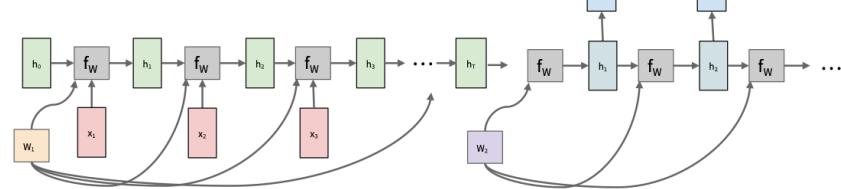
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 42

April 23, 2023

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector



Sutskever et al., "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Fei-Fei Li, Ehsan Adeli

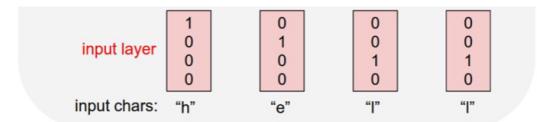
Lecture 8 - 43

April 23, 2023

Example:
Character-level
Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"



input layer 1 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0

input chars: "h" "e" "l" "l"

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 44

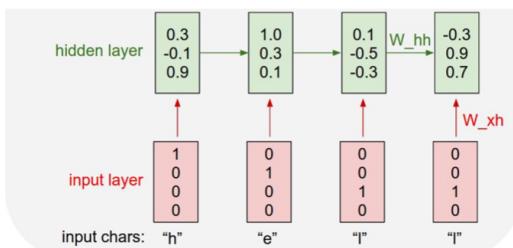
April 23, 2023

Example:
Character-level
Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"

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



Fei-Fei Li, Ehsan Adeli

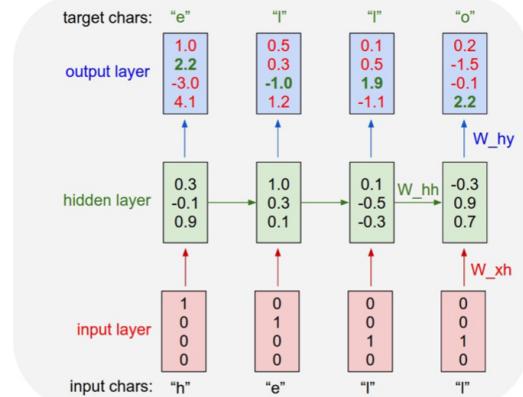
Lecture 8 - 45

April 23, 2023

Example:
Character-level
Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"



target chars: "e" "l" "l" "o"
output layer 1.0 2.2 -3.0 4.1 0.5 0.3 -1.0 1.2 0.1 0.5 1.9 -1.1 0.2 -1.5 -0.1 2.2

hidden layer 0.3 -0.1 0.9 1.0 0.3 0.1 0.1 -0.5 -0.3 -0.3 0.9 0.7

input layer 1 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0

Fei-Fei Li, Ehsan Adeli

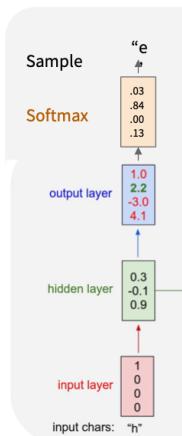
Lecture 8 - 46

April 23, 2023

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample characters one at a time, feed back to model



Fei-Fei Li, Ehsan Adeli

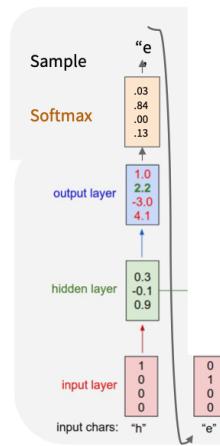
Lecture 8 - 47

April 23, 2023

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample characters one at a time, feed back to model



Fei-Fei Li, Ehsan Adeli

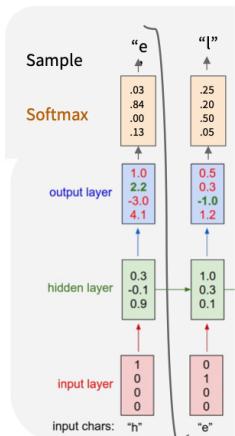
Lecture 8 - 48

April 23, 2023

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample characters one at a time, feed back to model



Fei-Fei Li, Ehsan Adeli

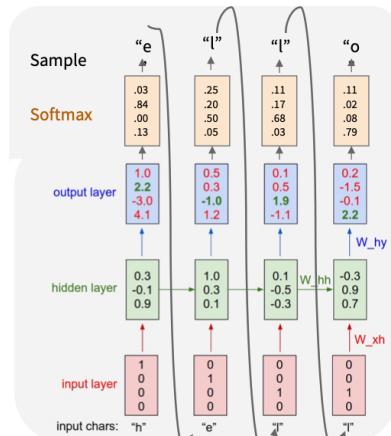
Lecture 8 - 49

April 23, 2023

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample characters one at a time, feed back to model



Fei-Fei Li, Ehsan Adeli

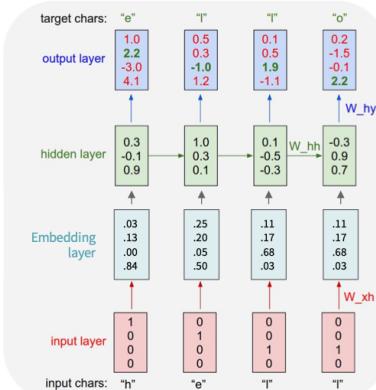
Lecture 8 - 50

April 23, 2023

Example: Character-level Language Model Sampling

$$\begin{aligned} [w_{11} w_{12} w_{13} w_{14}] [1] &= [w_{11}] \\ [w_{21} w_{22} w_{23} w_{14}] [0] &= [w_{21}] \\ [w_{31} w_{32} w_{33} w_{14}] [0] &= [w_{31}] \\ [w_{41} w_{42} w_{43} w_{44}] [0] &= [w_{41}] \end{aligned}$$

Matrix multiplication with a one-hot vector just extracts a column from the weight matrix. We often put a separate embedding layer between the input and hidden layers.



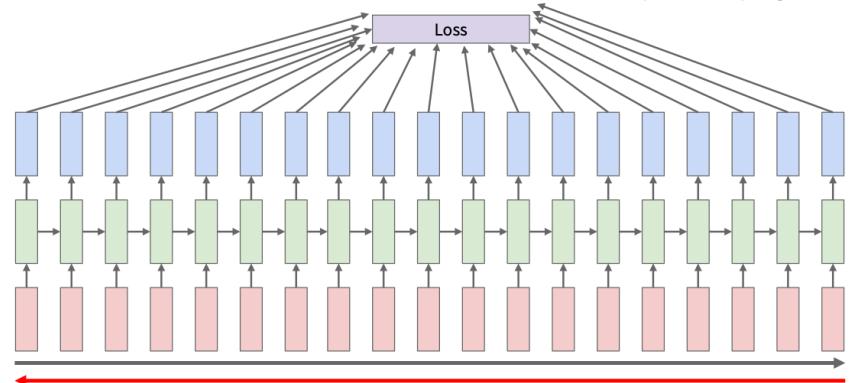
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 51

April 23, 2023

Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

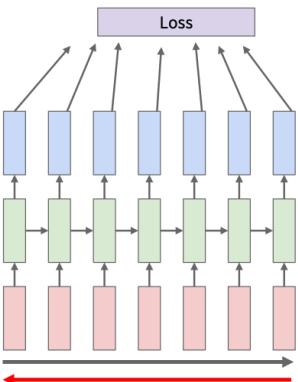


Fei-Fei Li, Ehsan Adeli

Lecture 8 - 52

April 23, 2023

Truncated Backpropagation through time



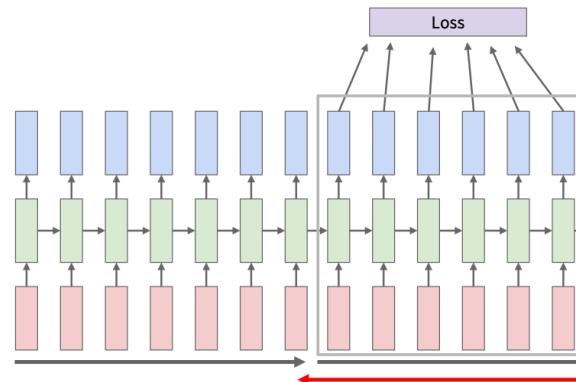
Run forward and backward through chunks of the sequence instead of whole sequence

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 53

April 23, 2023

Truncated Backpropagation through time



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 54

April 23, 2023

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

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

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

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

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

VIOLA:
I'll drink it.

Fei-Fei Li, Ehsan Adeli

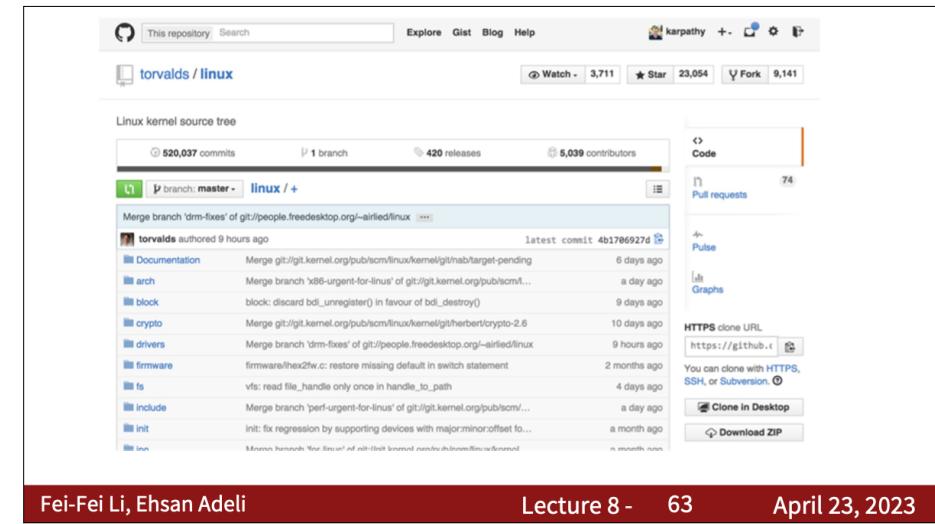
Lecture 8 - 59 April 23, 2023

Generated C code

```
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffff8) & 0x000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &offset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}
```

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 64 April 23, 2023



Fei-Fei Li, Ehsan Adeli

Lecture 8 - 63 April 23, 2023

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteuv.h>
#include <asm/pgproto.h>
```

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 65 April 23, 2023

```

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteov.h>
#include <asm/pgproto.h>

#define REG_PG    vesa_slot_addr_pack
#define PFM_NOCOMP AFCSR(0, load)
#define STACK_DDR(type)   (func)

#define SNAP_ALLOCATE(nr)      (e)
#define emulate_sigs() arch_get_unaligned_child()
#define access_rw(TST) asm volatile("movd %esp, %0, %3 : : \"r\" (0); \
if (_type & DO_READ)
static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
pC>[1]);
static void
os_prefix(unsigned long sys)
{
#endif CONFIG_PREEMPT
PUT_PARAM_RAID(2, sel) = get_state_state();
set_pid_sum((unsigned long)state, current_state_str(),
(unsigned long)-1->lr_full, low;
}

```

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 66

April 23, 2023

RNN tradeoffs

RNN Advantages:

- Can process any length of the input
- Computation for step t can (in theory) use information from many steps back
- Model size does not increase for longer input
- The same weights are applied on every timestep, so there is symmetry in how inputs are processed.

RNN Disadvantages:

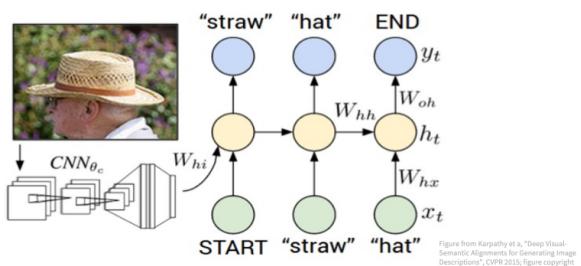
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 75

April 23, 2023

Image Captioning



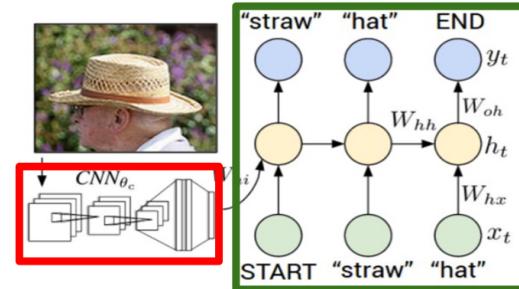
Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 76

April 23, 2023

Recurrent Neural Network



Convolutional Neural Network

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 77

April 23, 2023



test image

This image is CC0 public domain.

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 78

April 23, 2023

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax



test image

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 79

April 23, 2023



test image

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

softmax

~~FC-1000~~
~~softmax~~

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 80

April 23, 2023

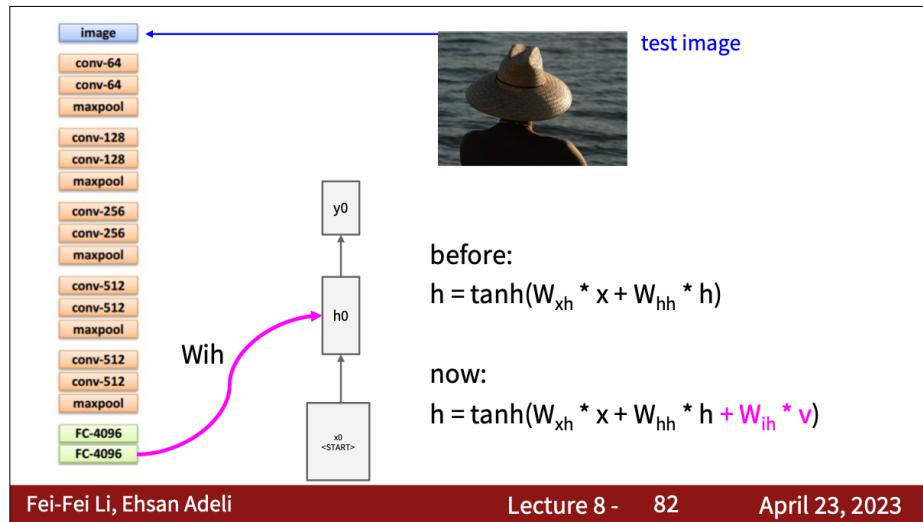


test image

Fei-Fei Li, Ehsan Adeli

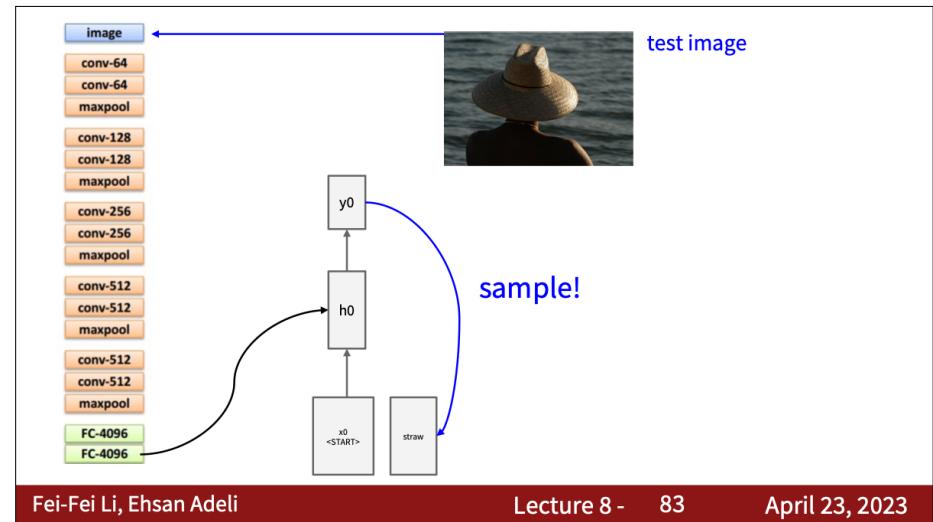
Lecture 8 - 81

April 23, 2023



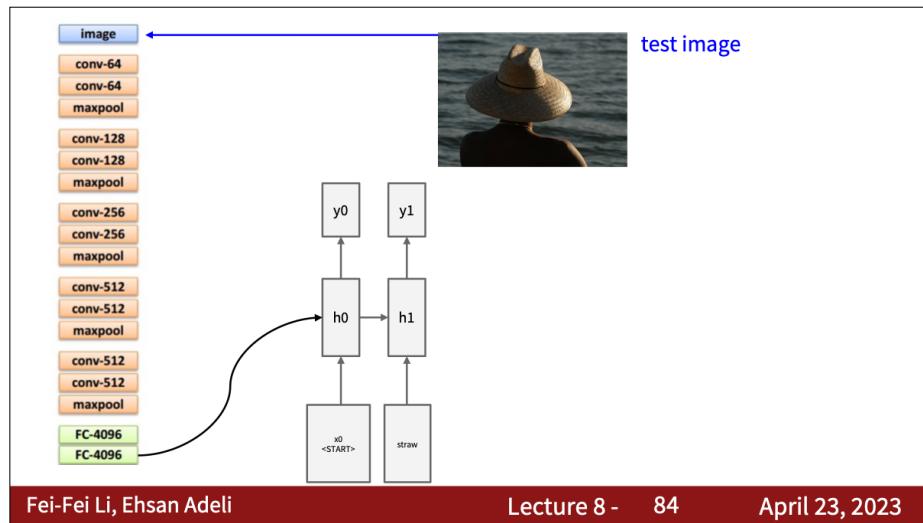
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 82 April 23, 2023



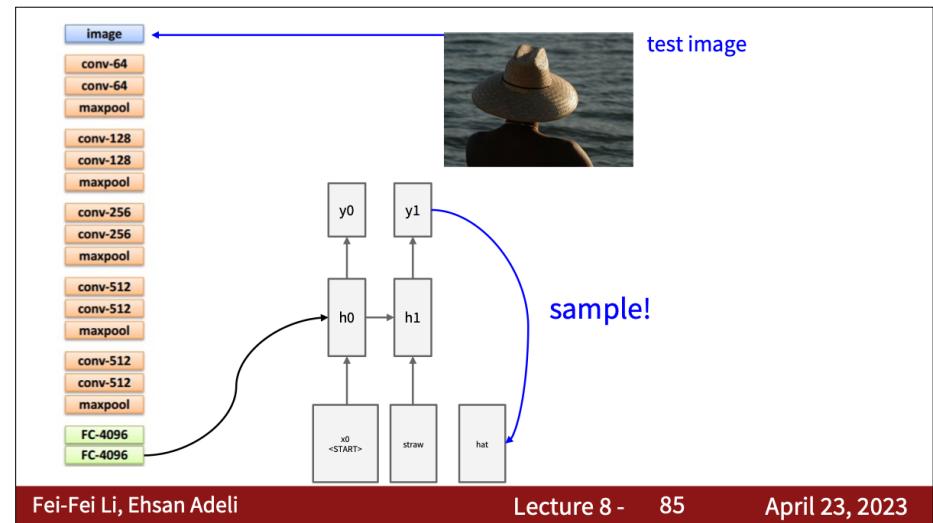
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 83 April 23, 2023



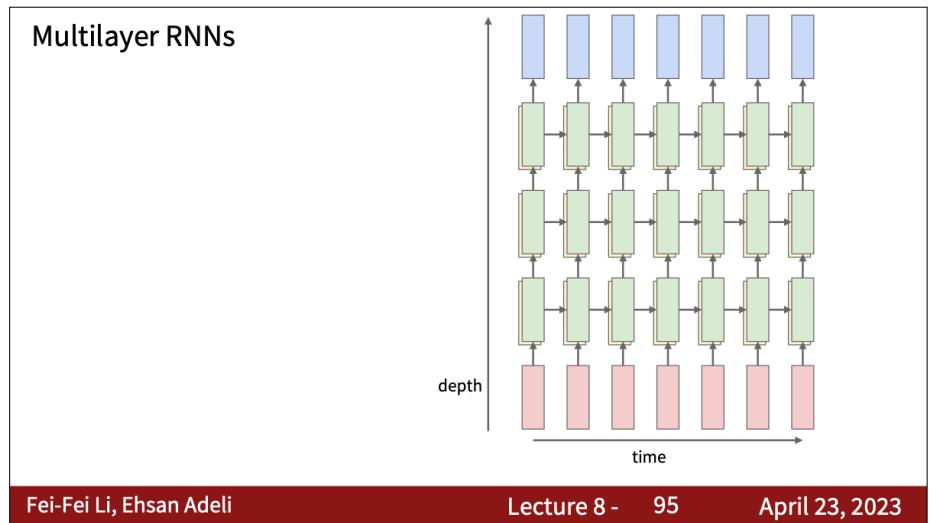
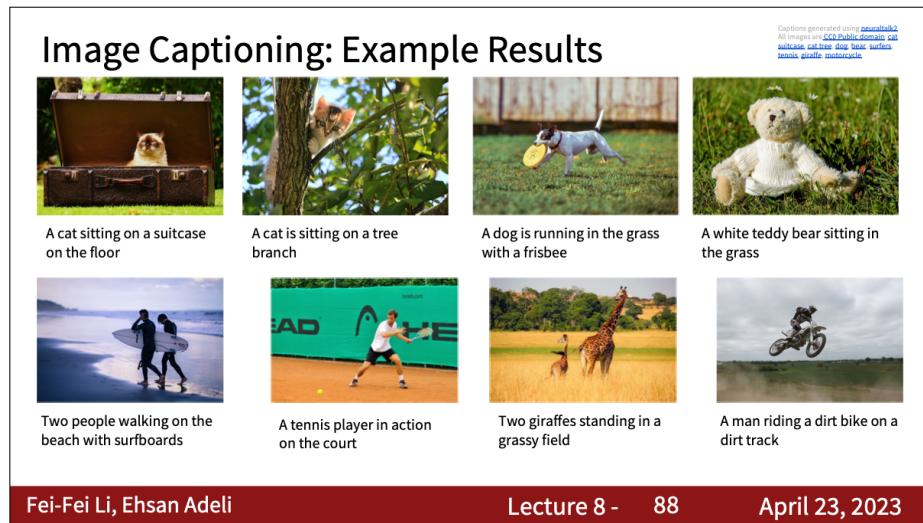
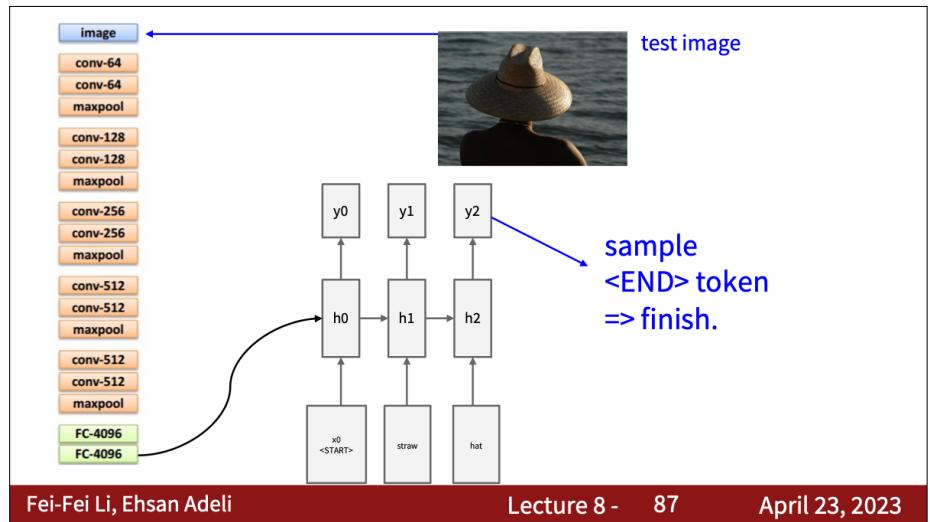
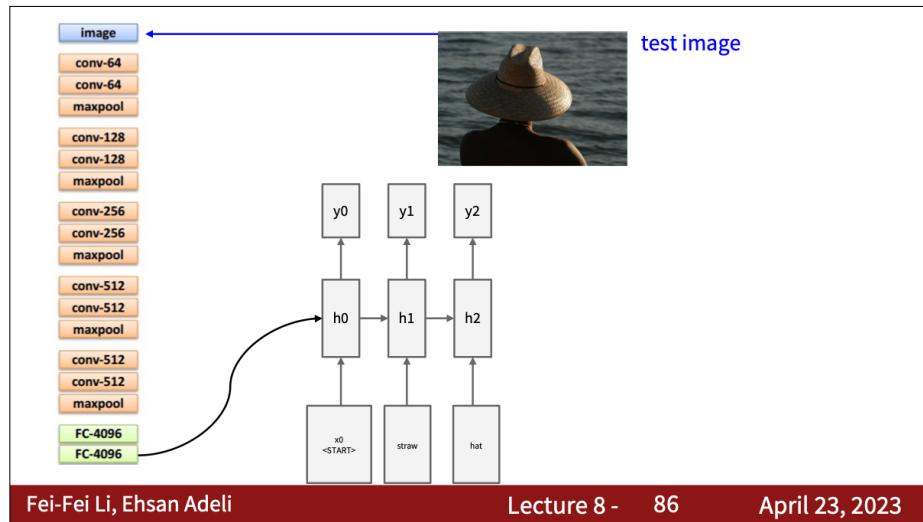
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 84 April 23, 2023



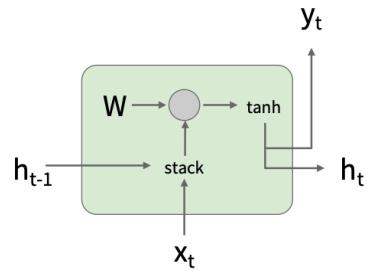
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 85 April 23, 2023



Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

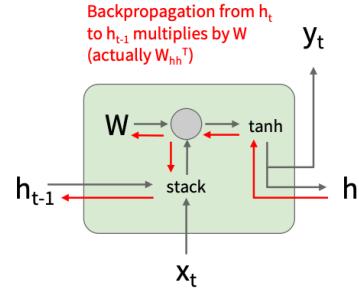
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 97

April 23, 2023

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

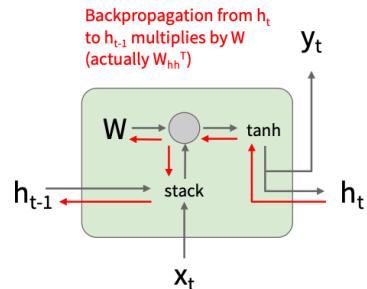
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 98

April 23, 2023

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}$$

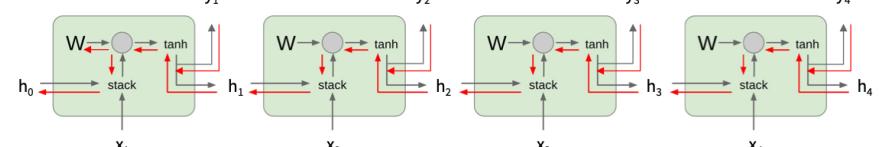
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 99

April 23, 2023

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

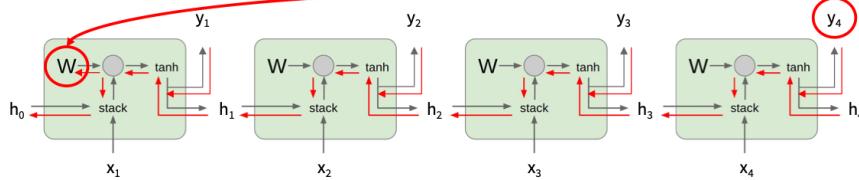
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 100

April 23, 2023

Vanilla RNN Gradient Flow

Gradients over multiple time steps:



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

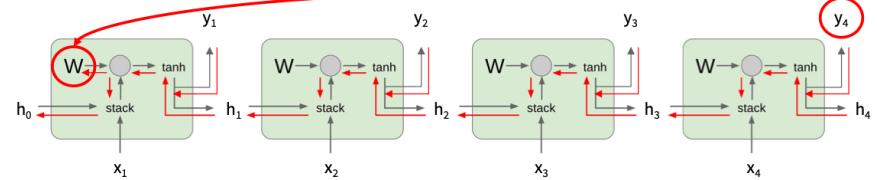
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W}$$

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 101 April 23, 2023

Vanilla RNN Gradient Flow

Gradients over multiple time steps:



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

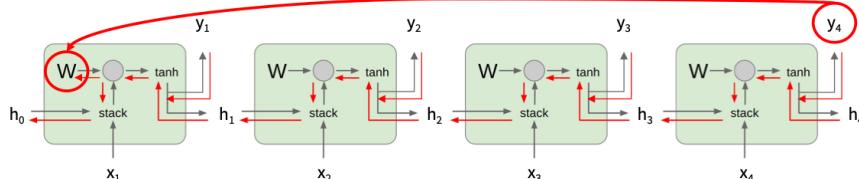
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 102 April 23, 2023

Vanilla RNN Gradient Flow

Gradients over multiple time steps:



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W} \quad \frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$

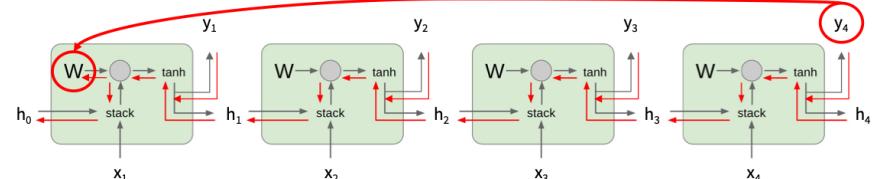
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 103 April 23, 2023

Vanilla RNN Gradient Flow

Gradients over multiple time steps:



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W} \quad \text{Almost always } < 1 \\ \text{Vanishing gradients}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) \right) W_{hh}^{T-1} \frac{\partial h_1}{\partial W}$$

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 104 April 23, 2023

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

Four gates

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

Cell state: $c_t = f \odot c_{t-1} + i \odot g$

Hidden state: $h_t = o \odot \tanh(c_t)$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

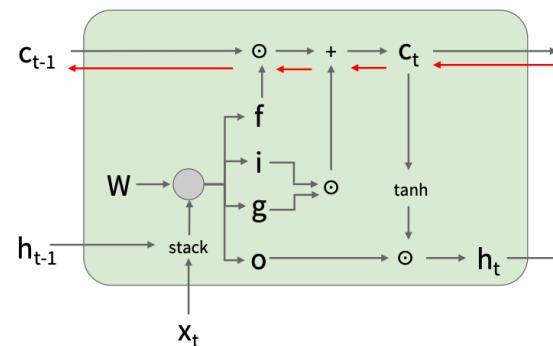
Fei-Fei Li, Ehsan Adeli

LSTM

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Backpropagation from c_t to c_{t-1} only elementwise multiplication by f , no matrix multiply by W



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

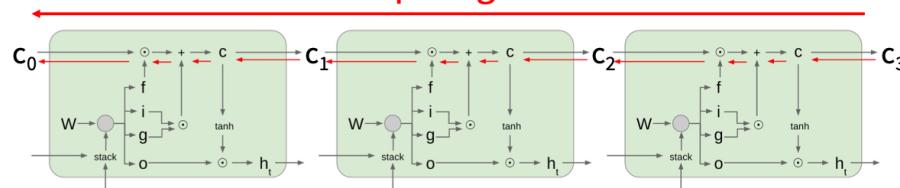
Fei-Fei Li, Ehsan Adeli

Lecture 8 - 117 April 23, 2023

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Fei-Fei Li, Ehsan Adeli

Lecture 8 - April 23, 2023

Do LSTMs solve the vanishing gradient problem?

The LSTM architecture makes it easier for the RNN to preserve information over many timesteps

- e.g. if the $f = 1$ and the $i = 0$, then the information of that cell is preserved indefinitely.
- By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix W_h that preserves info in hidden state

LSTM doesn't guarantee that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

Fei-Fei Li, Ehsan Adeli

Lecture 8 - 119 April 23, 2023

Summary

- RNNs allow a lot of flexibility in architecture design
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better/simpler architectures are a hot topic of current research, as well as new paradigms for reasoning over sequences
- Better understanding (both theoretical and empirical) is needed.