

lec10

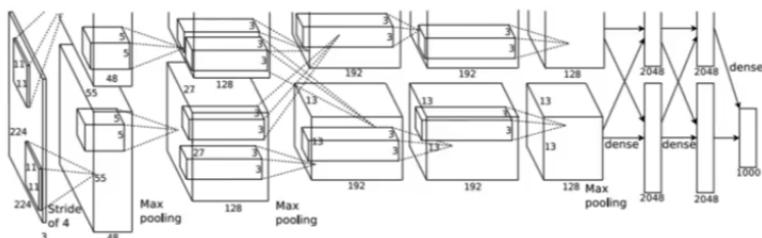
Creation Date: 22/01/2020 22:31

Last Modified Date: 23/01/2020 13:52

Lec 10 : Recurrent Neural Networks

Last Time: CNN Architectures

AlexNet



Revolution of Depth

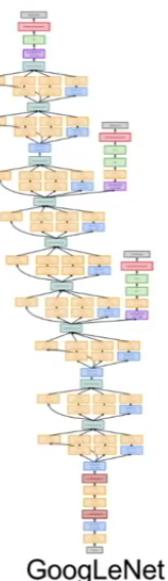
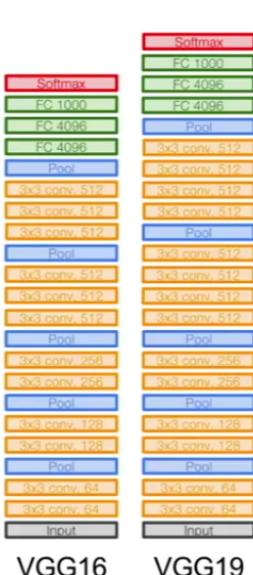


Figure copyright Kaiming He, 2016 Revolution it is permitted

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 4 University, May 4, 2017

Last Time: CNN Architectures



Revolution of Depth

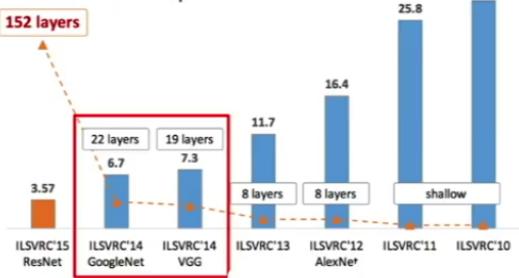
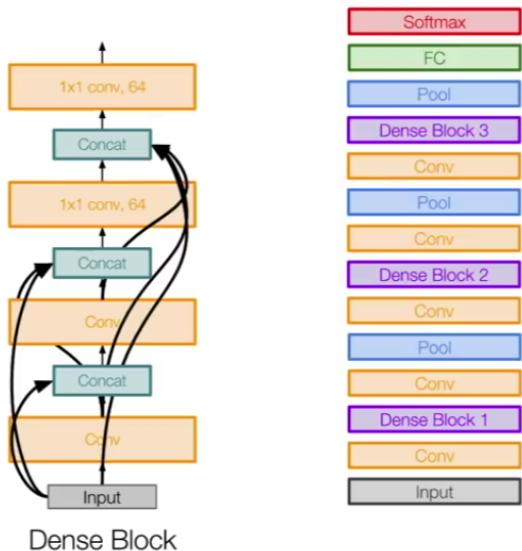


Figure copyright Kaiming He, 2016 Revolution it is permitted

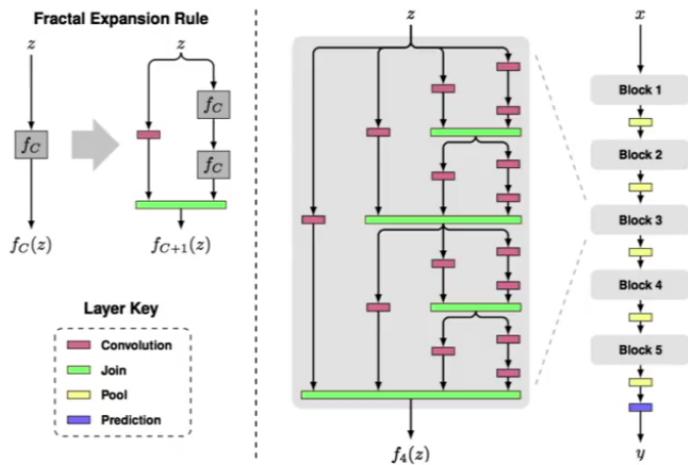
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 5 University, May 4, 2017

DenseNet



FractalNet

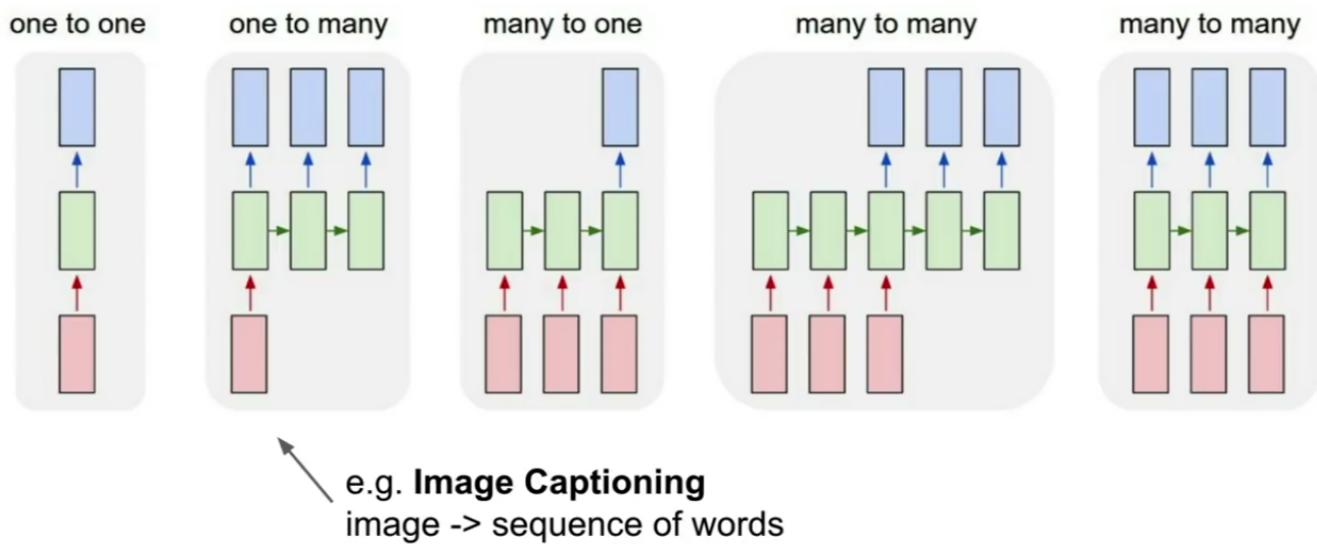


Figures copyright Larsson et al., 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 7 University, May 4, 2017

Recurrent Neural Networks: Process Sequences



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 12 University, May 4, 2017

- One-to-many: Image Captioning
- Many-to-one: Sentiment analysis given a string of words
- Many-to-many: Machine Translation(sequence of words to sequence of words)

Sequential Processing of Non-Sequence Data

Classify images by taking a series of “glimpses”

Ba, Mnih, and Kavukcuoglu, “Multiple Object Recognition with Visual Attention”, ICLR 2015.
Gregor et al, “DRAW: A Recurrent Neural Network For Image Generation”, ICML 2015
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 16 University May 4, 2017

Sequential Processing of Non-Sequence Data

Generate images one piece at a time!

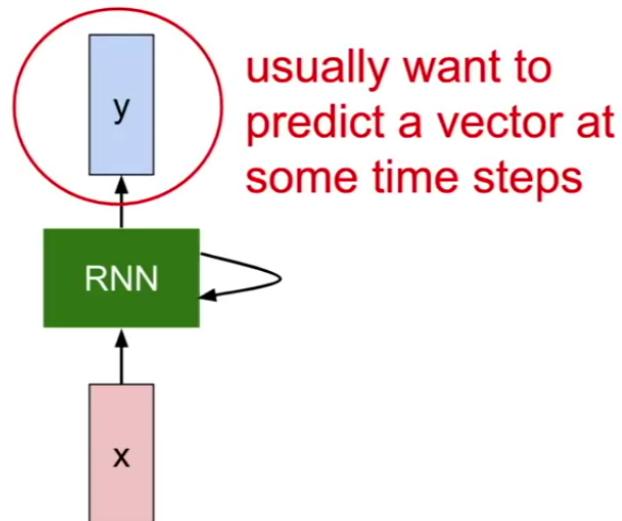
Gregor et al, “DRAW: A Recurrent Neural Network For Image Generation”, ICML 2015
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 17 University May 4, 2017

Recurrent Neural Network



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 19 University, May 4, 2017

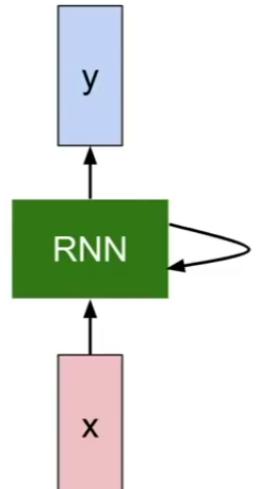
- The RNN block maintains an internal state which updates for every input.

Recurrent Neural Network

We can process a sequence of vectors \mathbf{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 & Justin Johnson & Serena Yeung

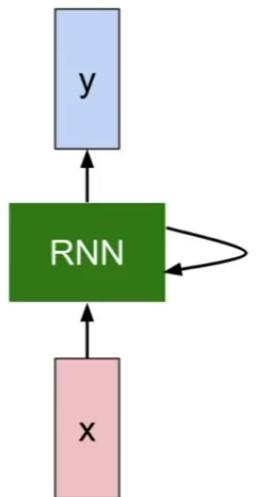
Lecture 10 - 20 University, May 4, 2017

- In the RNN block we are learning a recurrence relation.

Recurrent Neural Network

We can process a sequence of vectors \mathbf{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.

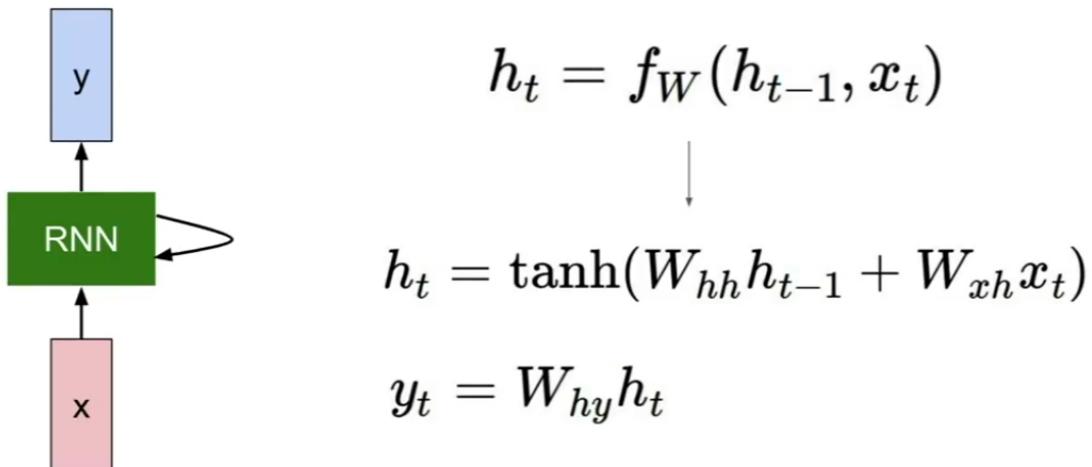
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 21 University of California, Berkeley May 4, 2017

NOTE: At every time step the same function and the same set of parameters are used.

(Vanilla) Recurrent Neural Network

The state consists of a single “hidden” vector \mathbf{h} :

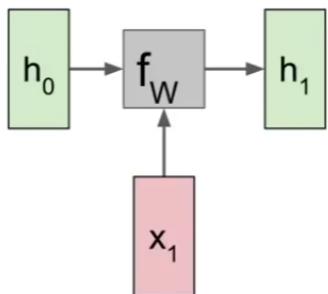


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 22 University of California, Berkeley May 4, 2017

- If we want to produce y_t at each time step then we use W_{hy} .

RNN: Computational Graph

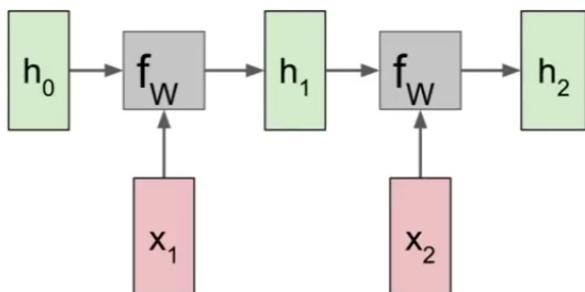


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 23 University of California, Berkeley May 4, 2017

- Usually the initial previous step h_0 is initialized to 0

RNN: Computational Graph

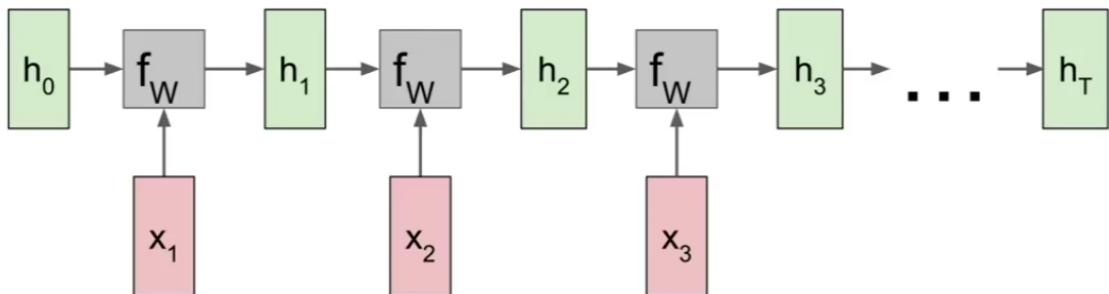


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 24 University of California, Berkeley May 4, 2017

- NOTE: The function f_W remains the same over time steps.

RNN: Computational Graph

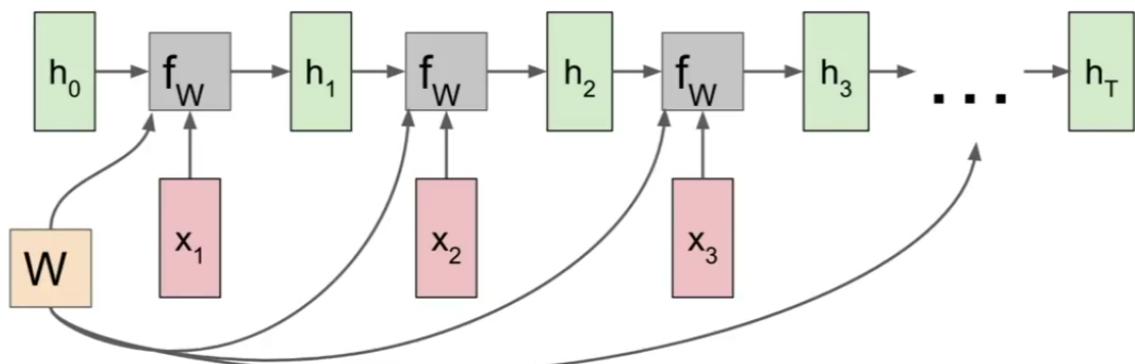


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 25 University of California, Berkeley, May 4, 2017

RNN: Computational Graph

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

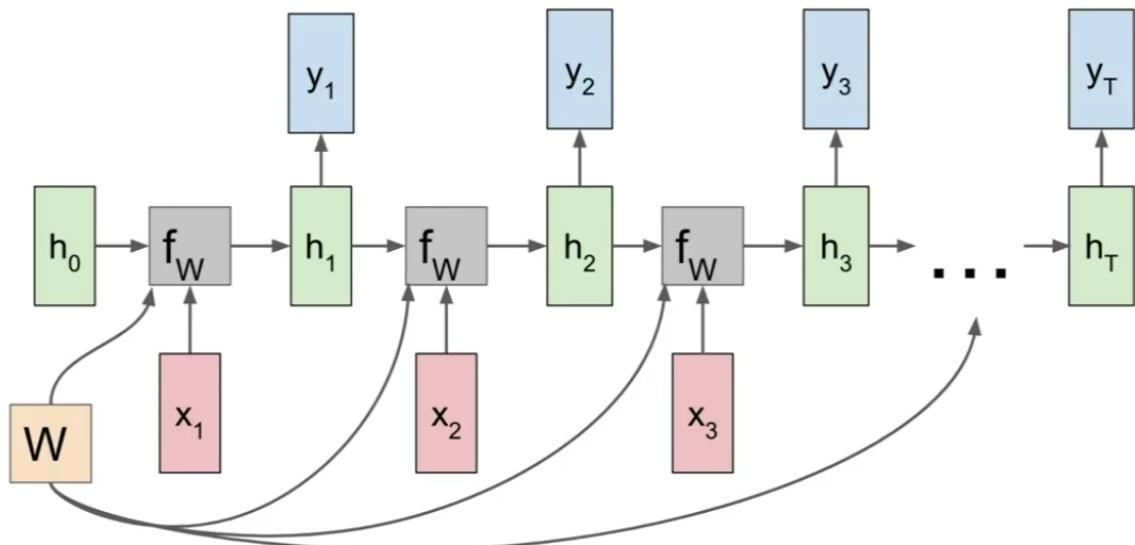


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 26 University of California, Berkeley, May 4, 2017

- WE RE-USE THE SAME WEIGHT MATRIX AT EVERY TIME-STEP

RNN: Computational Graph: Many to Many

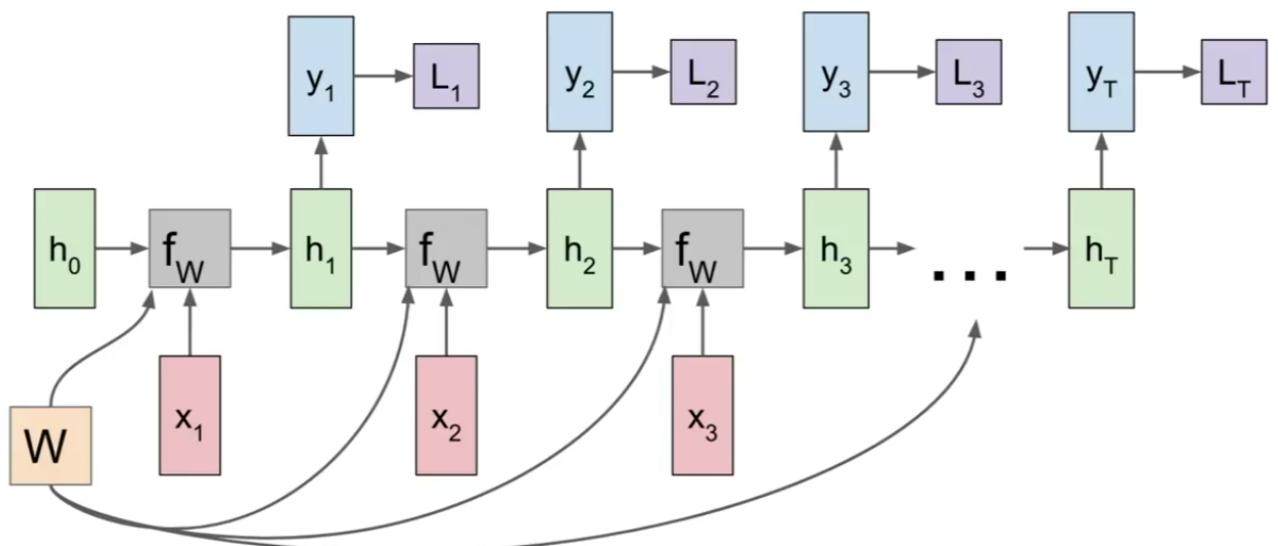


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 27 University May 4, 2017

- If we need y_t we can generate it using h_t at each time step.

RNN: Computational Graph: Many to Many

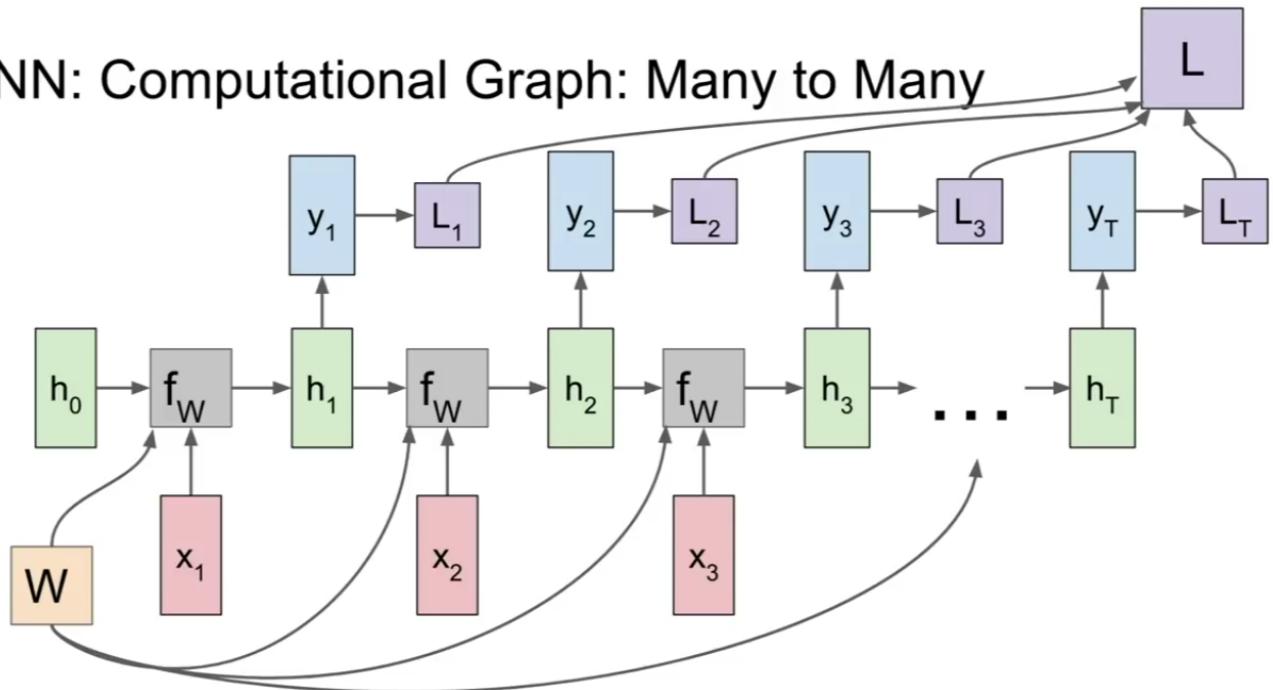


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 28 University May 4, 2017

- If we have ground truth label present at each time step then we can use it at corresponding time steps.

RNN: Computational Graph: Many to Many

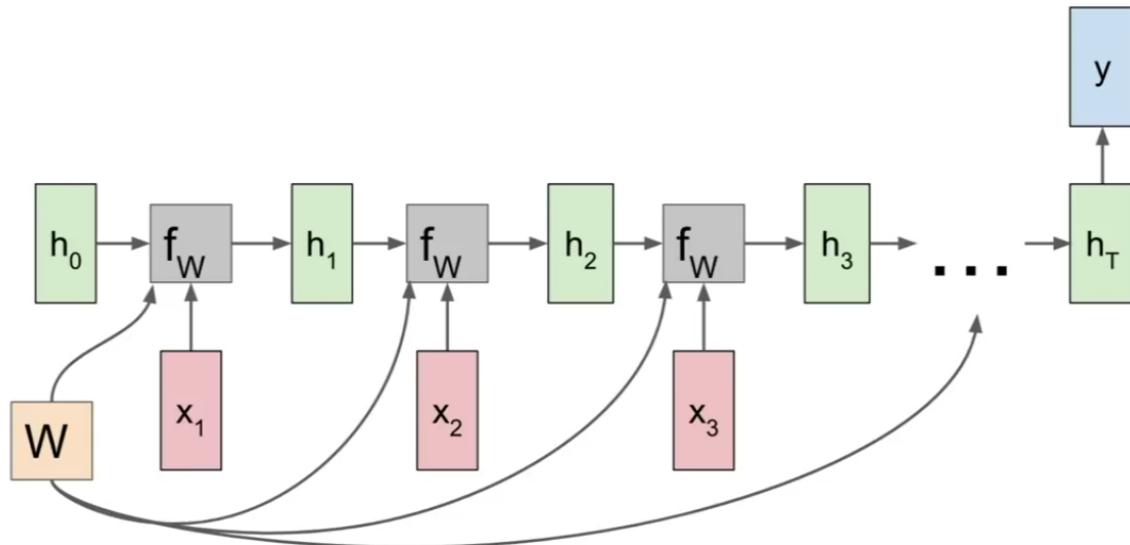


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 29 University May 4, 2017

- The overall loss will be the combined loss at each times step.
- During backpropagation the gradients will flow to each time step, and the corresponding time step will compute local gradients on the weight W

RNN: Computational Graph: Many to One

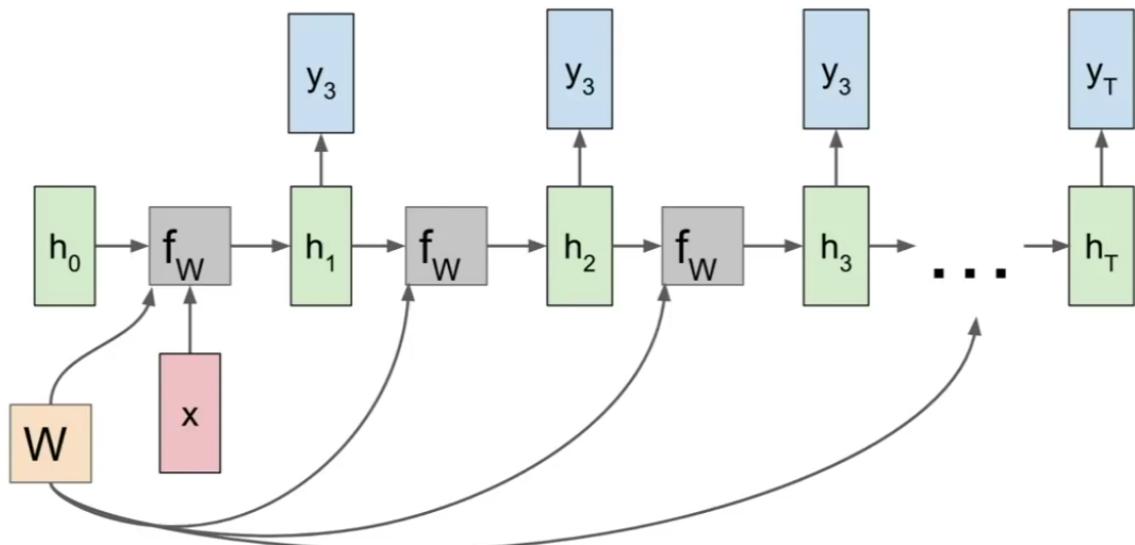


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 30 University May 4, 2017

- During Sentiment Analysis we can use the output of the last state since it is a many-to-one problem.

RNN: Computational Graph: One to Many



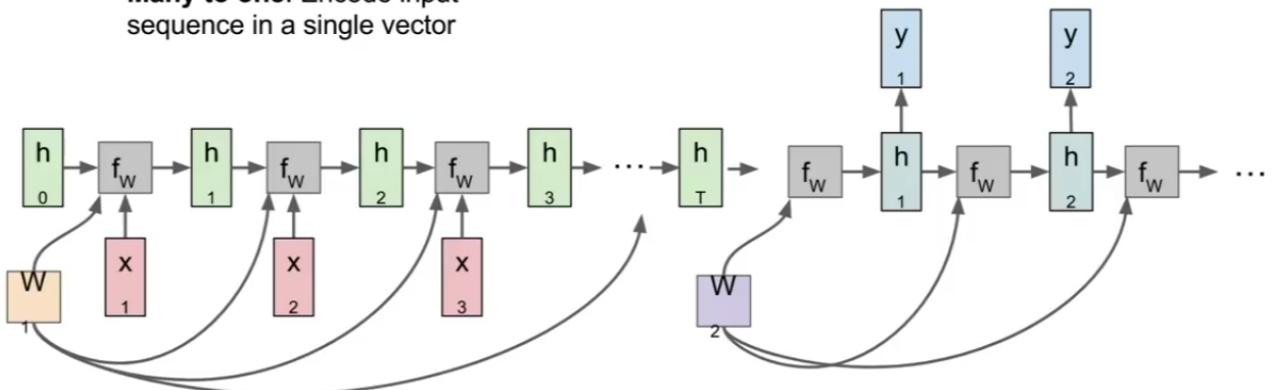
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 31 University May 4, 2017

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

Many to one: Encode input sequence in a single vector

One to many: Produce output sequence from single input vector



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 33 University May 4, 2017

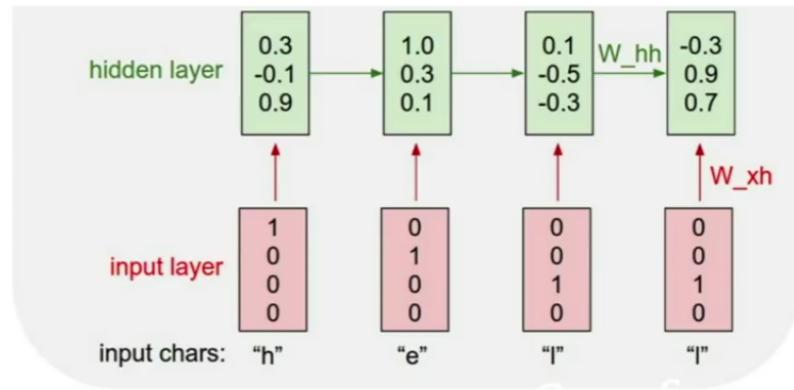
- This is like an encoder-decoder network for Sequence to Sequence translation where the encoder will summarize the input into a single vector which is fed to decoder
- RNNs are commonly used in language modeling where we can learn for either character level data or word level data and generate characters and words.

Example: Character-level Language Model

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

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

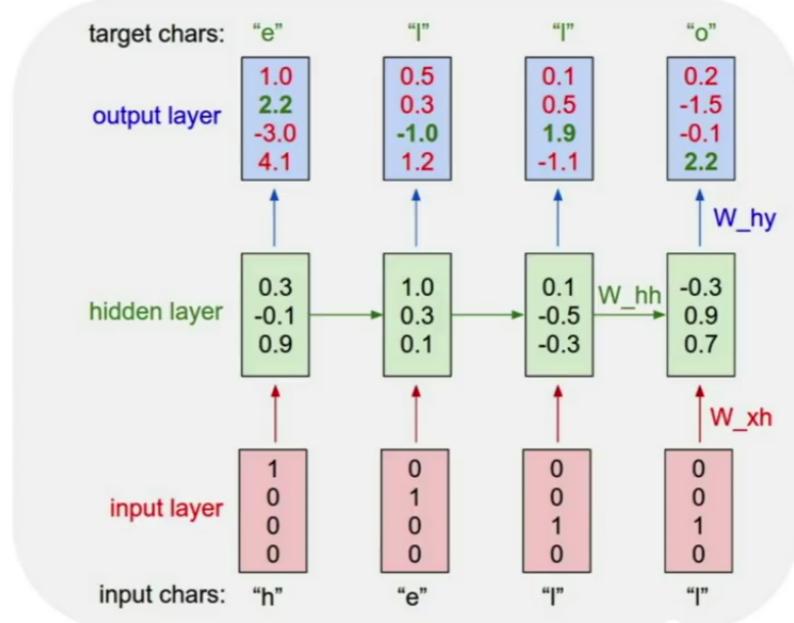


- Since the given vocabulary(all possible characters in the given language) has only 4 letters, each of the input vector will be of the length 4.

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

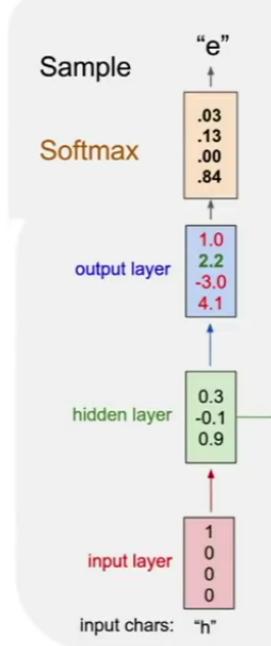


- During the forward pass at the first time step, it will receive the input letter h. It will go to the RNN cell and this cell will produce the output y_t which for the given problem is the next character from the given vocabulary.
- We calculate the loss based on the target and the predicted vector.

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

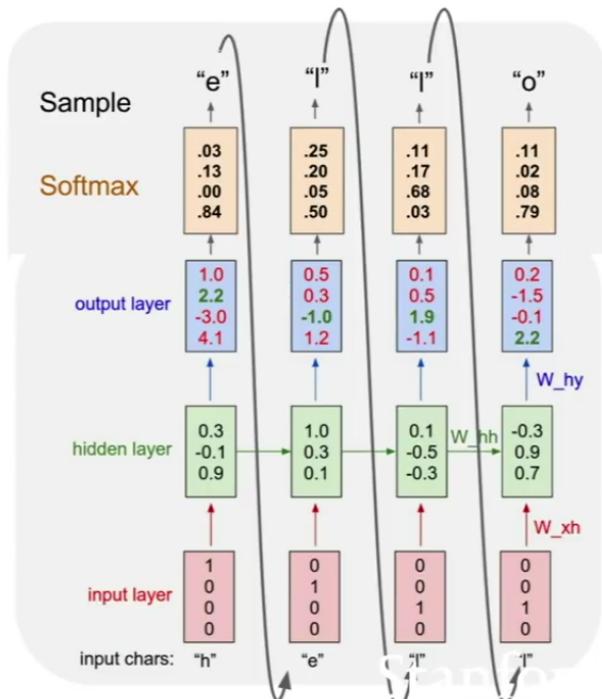
At test-time sample
characters one at a time,
feed back to model



Example: Character-level Language Model Sampling

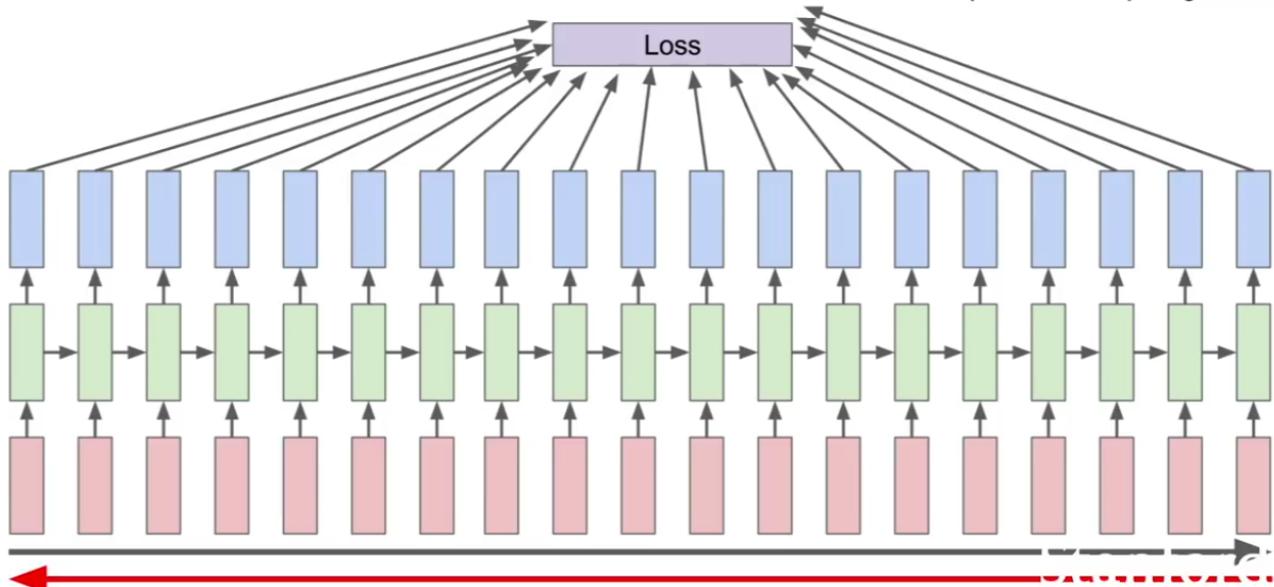
Vocabulary:
[h,e,l,o]

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



Backpropagation through time

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



Fei-Fei Li & Justin Johnson & Serena Yeung

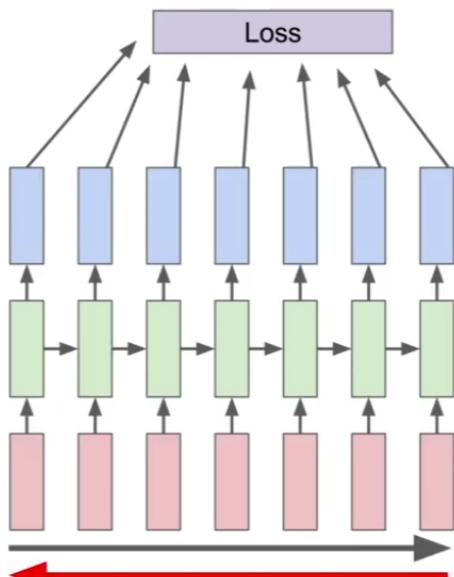
Lecture 10 - 41 University May 4, 2017

Backpropagation through time: During forward pass we use the entire sequence to compute loss, then during backpropagation we compute gradients at each time step.

- Problem: If we have a very large sequence (like entire Wikipedia) then both the forward pass (to compute total loss) and backward pass (to compute W) will take huge amount of time.
- Solution: Truncated Backpropagation through time.

Truncated Backpropagation through time

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

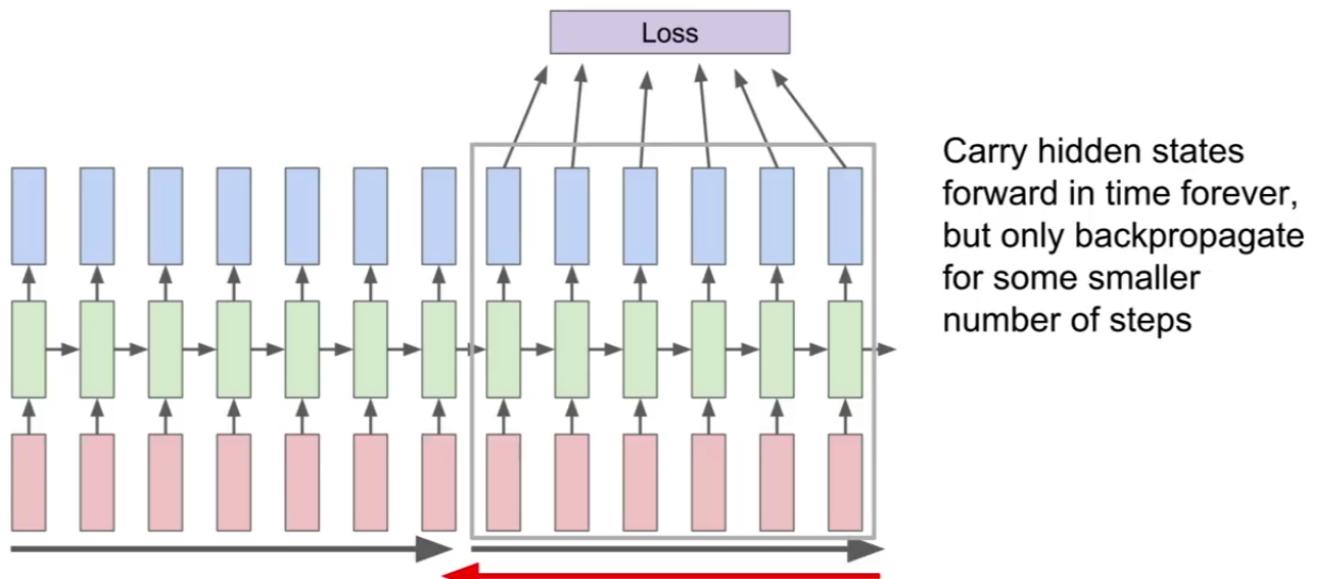


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 42 University May 4, 2017

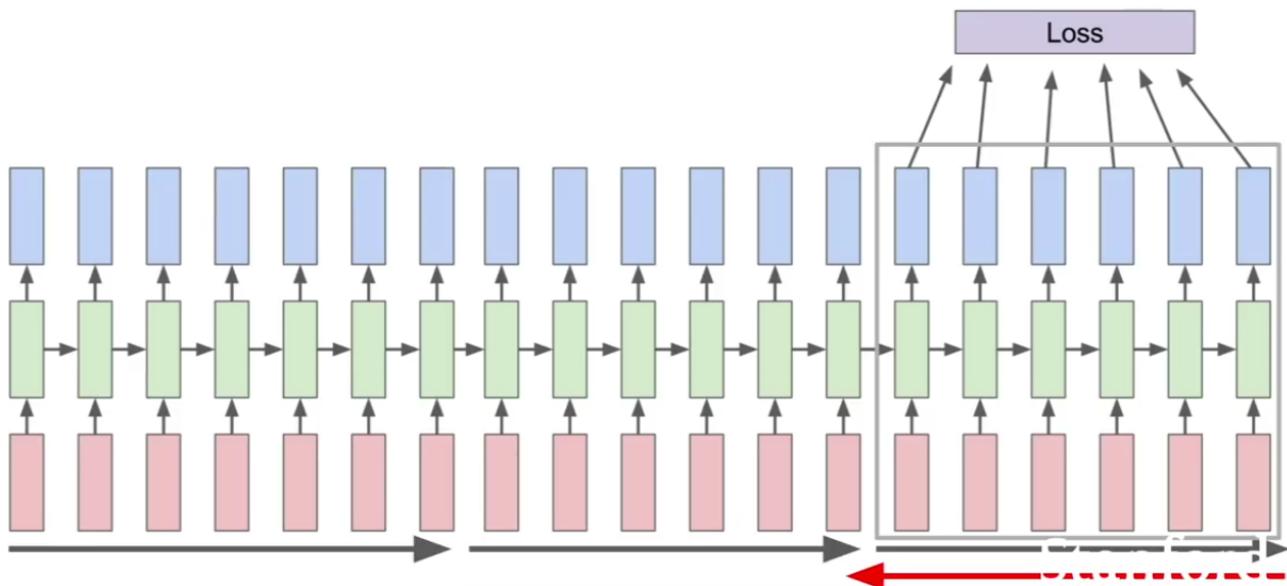
- In this approach we divide the sequence into chunks and perform forward and backward pass through these chunks.

Truncated Backpropagation through time



- We carry forward the hidden states to the next step of "chunk" processing but the backpropagation will happen only in the "chunk".

Truncated Backpropagation through time



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 44 University May 4, 2017

[min-char-rnn.py](#) gist: 112 lines of Python

```

***

Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)

BSD License

Copyright (c) 2014 Andrej Karpathy

import numpy as np

# data I/O
data = open("input.txt", 'r').read() # should be simple plain text file
chars = list(set(data))
data_size, vocab_size = len(data), len(chars)
print "data has %d characters, %d unique" % (data_size, vocab_size)

char_to_ix = {ch:i for i, ch in enumerate(chars)}
ix_to_char = {i:ch for i, ch in enumerate(chars)}

# hyperparameters
hidden_size = 200 # size of hidden layer of neurons
seq_length = 20 # number of steps to unroll the RNN for
learning_rate = 1e-1

# model parameters
wh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
bh = np.zeros(hidden_size) # bias for hidden state
w-hidden = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
bh-hidden = np.zeros(hidden_size) # bias for hidden state
by = np.zeros(vocab_size) # output bias

def lossFun(inputs, targets, hprev):
    """ Compute loss and gradient over all time steps. 

    inputs,targets are both list of integers.
    hprev is initial hidden state.
    returns the loss, gradients on model parameters, and last hidden state
    """
    loss = 0.0
    dhprev = np.zeros_like(hprev)
    for t in range(len(inputs)):
        x = inputs[t]
        y = targets[t]
        hprev = hprev if t == 0 else hprev
        hprev, wh, bh, by = forward(x, hprev, wh, bh, by)
        loss += np.square(y - hprev).sum()
        dhprev = backward(y, hprev, dhprev, wh, bh, by)
    return loss, dhprev

def forward(x, hprev, wh, bh, by):
    """ Compute h = tanh(wx + bh) and y = softmax(hw + by) """
    h = np.tanh(np.dot(x, wh) + bh)
    y = np.exp(hw + by) / np.sum(np.exp(hw))
    return h, y

def backward(y, hprev, dhprev, wh, bh, by):
    """ Compute gradients w.r.t. previous hidden state, input, and weights """
    dh = np.dot(y, wh.T) + dhprev
    dw = np.outer(dhprev, x)
    db = dhprev.sum(0)
    dy = dh * (1 - h**2)
    dx = dy * np.dot(dy, wh.T)
    dhprev = np.dot(dx, wh)
    return dh, dw, db, dx, dy

def sample(h, seed_ix, n):
    """ Sample a sequence of integers from the model
    h is memory state, seed_ix is seed letter for first time step
    """
    if n == 0:
        return []
    ixes = []
    for t in range(n):
        h, y = forward(seed_ix, h, wh, bh, by)
        p = np.exp(y) / np.sum(np.exp(y))
        ix = np.random.choice(range(vocab_size), p=prob.ravel())
        x = np.zeros(vocab_size)
        x[ix] = 1
        ixes.append(ix)
    return ixes

n, p = 0, 0.0
batch_size, numh = np.zeros_like(wh), np.zeros_like(bh), np.zeros_like(by)
mem_vars = [batch_size, numh, by] # memory variables for Adagrad
smooth_loss = -np.log(1.0/vocab_size)*seq_length # base at iteration 0
while True:
    if p > seq_length:
        # prepare inputs (we're reusing same inputs)
        inputs = np.zeros((batch_size, seq_length))
        hprev = np.zeros((batch_size, hidden_size))
        g = g + 0 # start of data
        inputs = [char_to_ix[ch] for ch in data[p-seq_length:p]]
        targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]

        # sample from the model now and then
        if t % 1000 == 0:
            hprev = sample(inputs[0], 200)
            txt = ''.join(ix_to_char[i] for i in hprev)
            print '-----%s-----' % (txt, )
        g = g + 0 # end of data

        # forward seq_length characters through the net and fetch gradient
        loss, dh, dw, db, dx, dy, hprev = lossFun(inputs, targets, hprev)
        smooth_loss = smooth_loss * 0.9 + loss * 0.001
        if t % 1000 == 0: print 'iter %d, loss: %f' % (t, smooth_loss)
        print progress

        # perform parameter update with Adagrad
        for param, dparam, mem in zip([wh, bh, by], [dw, db, dx], [batch_size, numh, by]):
            m = mem + dparam * dparam
            param -= learning_rate * dparam / np.sqrt(m + 1e-8) # adagrad update
            mem = m

        p += seq_length # move data pointer
    n += 1 # iteration counter

```

Fei-Fei Li & Justin Johnson & Serena Yeung

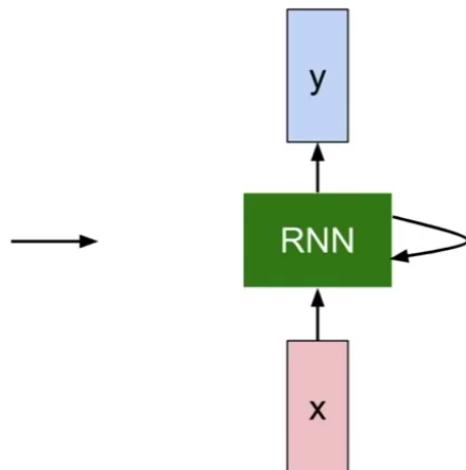
Lecture 10 - 45 University May 4, 2017

THE SONNETS

by William Shakespeare

From fairest creatures we desire increase,
That thereby beauty's rose might never die,
But as the riper should by time decease,
His tender heir might bear his name:
But thou, contraried to thine own bright eyes,
Feed'st thy light's flame with self-substantial fuel,
Making a famine where abundance lies,
Thyself thy foe, to thy sweet self too cruel:
Though that art now the world's fresh ornament,
And only hold to the gaudy spring,
Within thine own bed buried to content,
And tender churl mak'st waste in niggarding:
Pity the world, or else this glutton be,
To eat the world's due, by the grave and thee.

When forty winters shall besiege thy brow,
And dig deep trenches in thy beauty's field,
Thy youth's proud livery so gazed on now,
Will be a tatter'd weed of small worth held:
Then being asked, where all thy beauty lies,
Where all the treasure of thy lusty days;
To say, within thine own deep sunken eyes,
Were an all-eating shame, and thriftless praise.
How much more praise deserved thy beauty's use,
If thou couldst answer 'This fair child of mine
Shall sum my count, and make my old excuse,'
Proving his beauty by succession thine!
This were to be new made when thou art old,
And see thy blood warm when thou feel'st it cold.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 46 University May 4, 2017

at first:

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoanns lng

↓ train more

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwyl fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her heary, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 47 University May 4, 2017

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.

VIOLA:

Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

For $\bigoplus_{m=1,\dots,m} \mathcal{L}_m = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ???. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sh}(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\text{GL}_{S'}(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of X' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widehat{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)^{opp}_{fppf}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longrightarrow (U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example ???. It may replace S by $X_{\text{spaces},\text{etale}}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ???. Namely, by Lemma ?? we see that R is geometrically regular over S .

Lemma 0.1. Assume (3) and (3) by the construction in the description.

Suppose $X = \lim X_i$ (by the formal open covering X and a single map $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$ over U compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If T is surjective we may assume that T is connected with residue fields of S . Moreover there exists a closed subspace $Z \subset X$ of X where U in X' is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

Proof. This is form all sheaves of sheaves on X . But given a scheme U and a surjective étale morphism $U \rightarrow X$. Let $U \cap U_i = \coprod_{i=1,\dots,n} U_i$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \lim_i X_i$. \square

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{X,\dots,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S , $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq p$ is a subset of $\mathcal{J}_{n,0} \circ \mathcal{A}_2$ works.

Lemma 0.3. In Situation ???. Hence we may assume $q' = 0$.

Proof. We will use the property we see that p is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F -algebra where δ_{n+1} is a scheme over S . \square

```

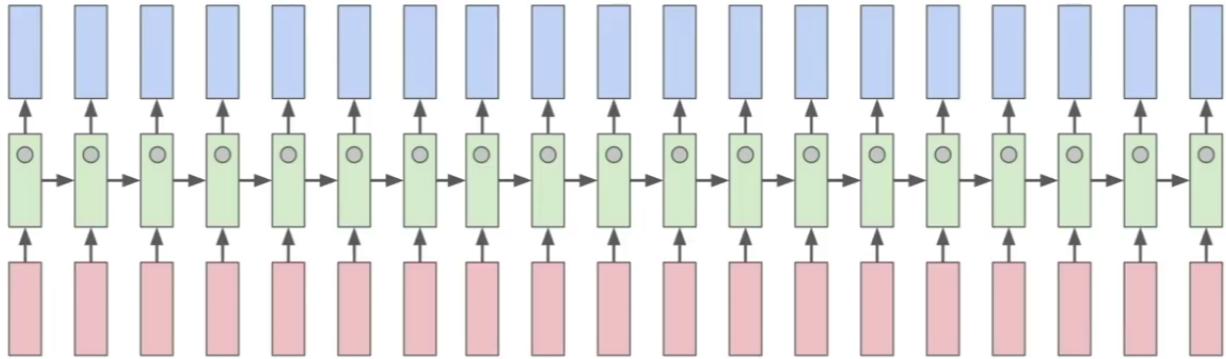
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << i))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffffff8) & 0x0000000f) << 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, &soffset);
    /* 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 ");
}

```

Generated C code

- To understand what the RNN is learning over the course of steps we can see how the RNN (a vector in the RNN cell) changes over time steps.

Searching for interpretable cells



Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 56 University, May 4, 2017

Searching for interpretable cells

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
    */
```

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 57 University, May 4, 2017

- Initially the vectors output some gibberish, but over the course of time we can notice that certain vector is responsible to learn some kind of a hidden feature.

Subtitle track: Disable Searching for interpretable cells

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

quote detection cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 58 University May 4, 2017

- Here we can notice that one of the vector is learning to interpret 'quotes'.
- We can notice that when the model detects a quotation(") it turns on (high activation) and once it finds the end of a quotation(") it turns off until it detects the next quotation character.

Searching for interpretable cells

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

line length tracking cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 59 University May 4, 2017

- Here the cell is tracking length of a line

Searching for interpretable cells

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
    siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

if statement cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 60 University May 4, 2017

- Here the cell is active when it detects an 'if' statement

Searching for interpretable cells

Cell that turns on inside comments and quotes:

```
/* Duplicate LSM field information. The lsm_rule is opaque, so
 * re-initialized. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
    struct audit_field *sf)
{
    int ret = 0;
    char *lsm_str;
    /* our own copy of lsm_str */
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);
    if (unlikely(!lsm_str))
        return -ENOMEM;
    df->lsm_str = lsm_str;
    /* our own (refreshed) copy of lsm_rule */
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
        (void *)adf->lsm_rule);
    /* Keep currently invalid fields around in case they
     * become valid after a policy reload. */
    if (ret == -EINVAL) {
        pr_warn("audit rule for LSM '\\%s\\' is invalid\n",
            df->lsm_str);
        ret = 0;
    }
    return ret;
}
```

quote/comment cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 61 University May 4, 2017

- Here the cell is activated when it detects comments

Searching for interpretable cells

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

code depth cell

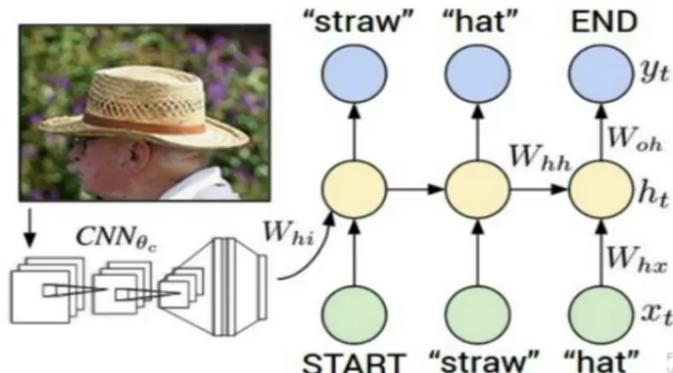
Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 62 University May 4, 2017

- Here the cell is activated for indentation in the code.

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.

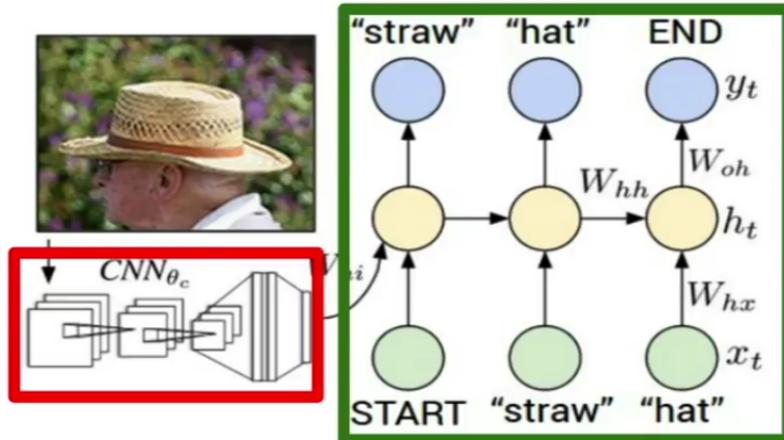
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 63 University May 4, 2017

- In CV for the problem of Image Captioning we need to use RNN to generate the caption given an Image (one_to_many RNN)

Recurrent Neural Network



Convolutional Neural Network

Fei-Fei Li & Justin Johnson & Serena Yeung

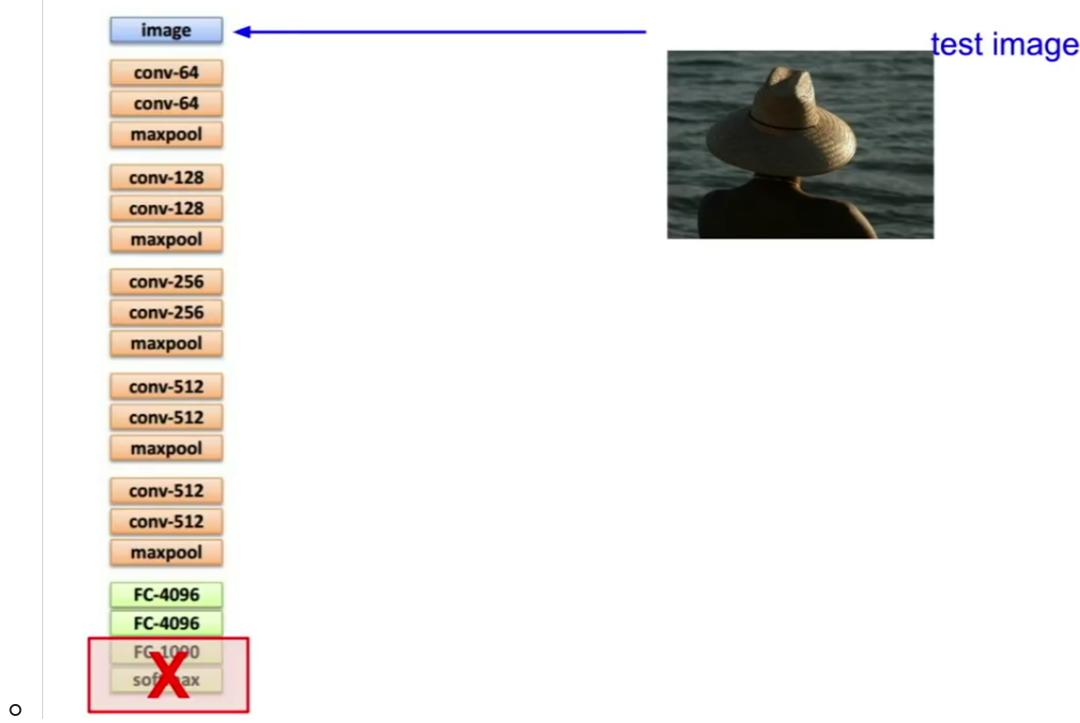
Lecture 10 - 64 University May 4, 2017

- The CNN will produce a summary(feature) vector which will be fed to the first step of the RNN language model which will produce words of the caption one at a time.
- How does Image Captioning work during Test Time?

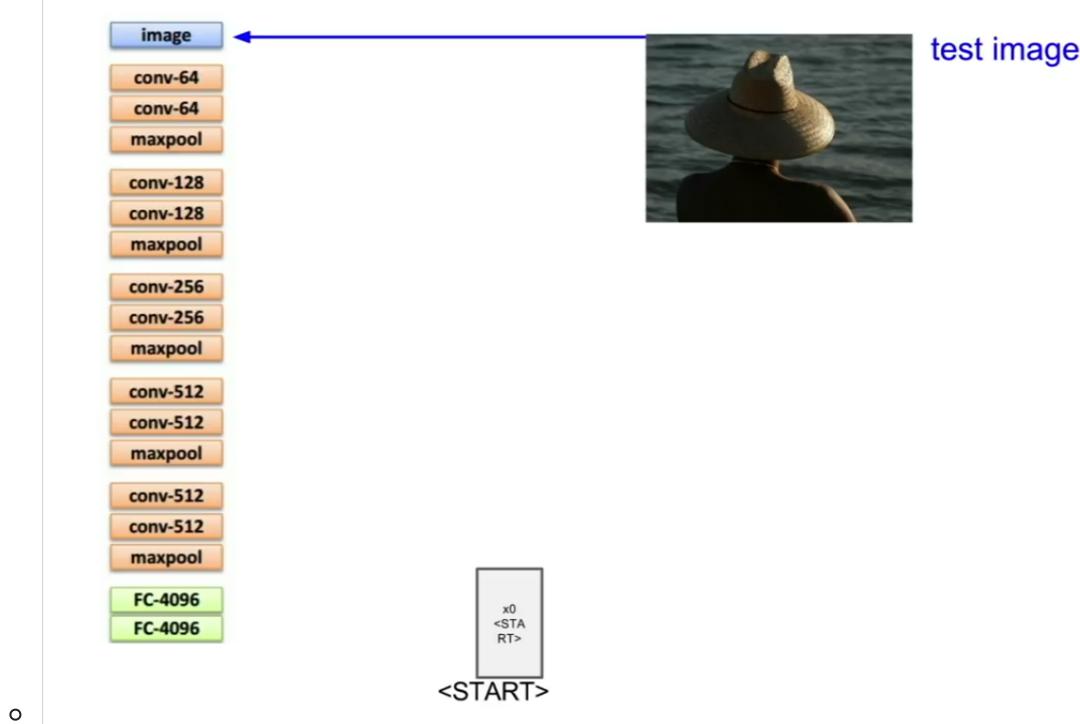


This image is CC0 public domain

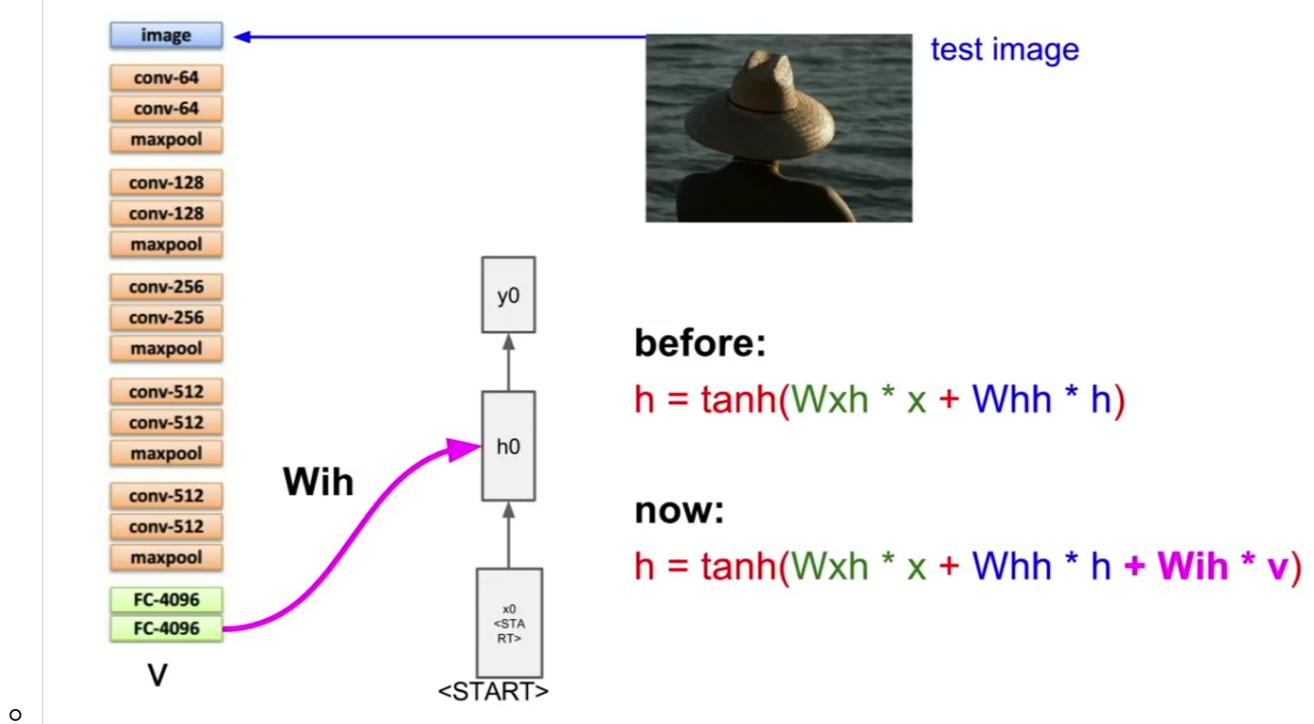
- Given a test image we feed it through the CNN model, but instead of passing it through the final FC layers which are commonly used for classification, we only use the final feature, here it is FC-4096 layer.



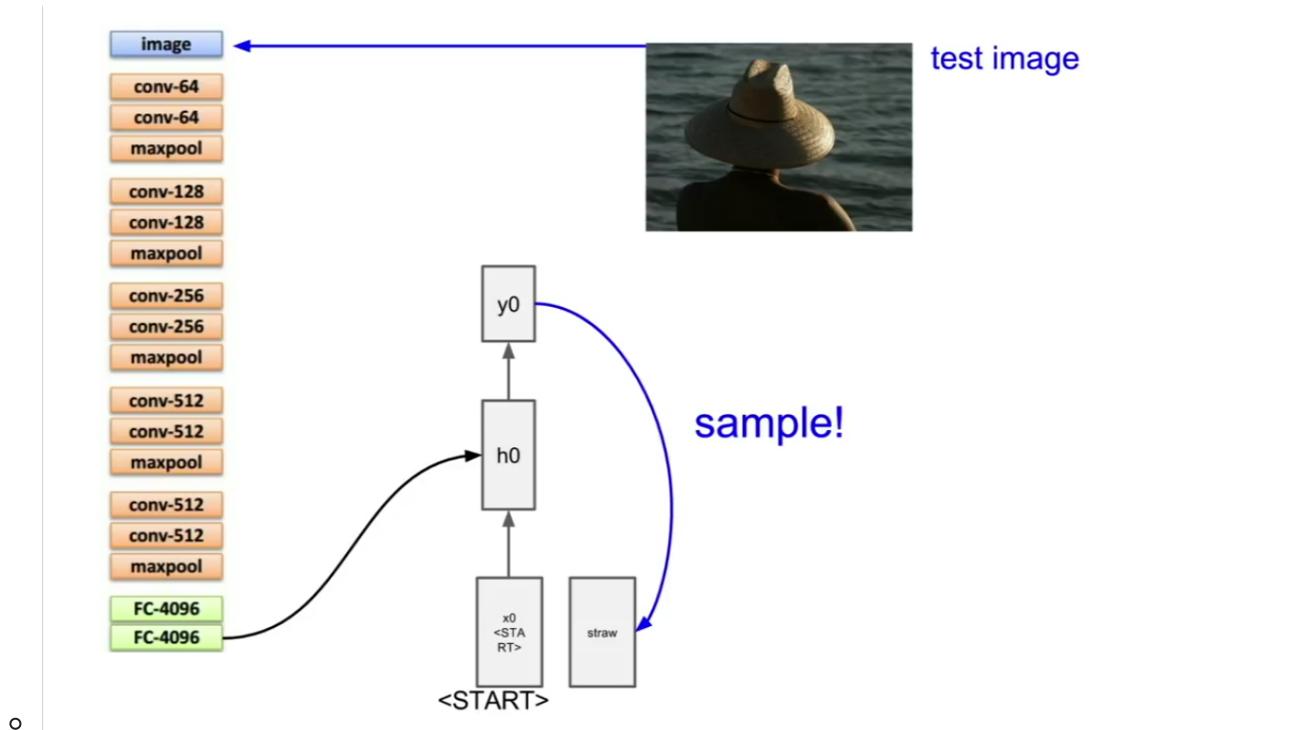
- We will use the FC-4096 vector to summarize the image.



- To start generating the caption we need to seed the language model. Here the token is the initial seed which is fed into the language model.

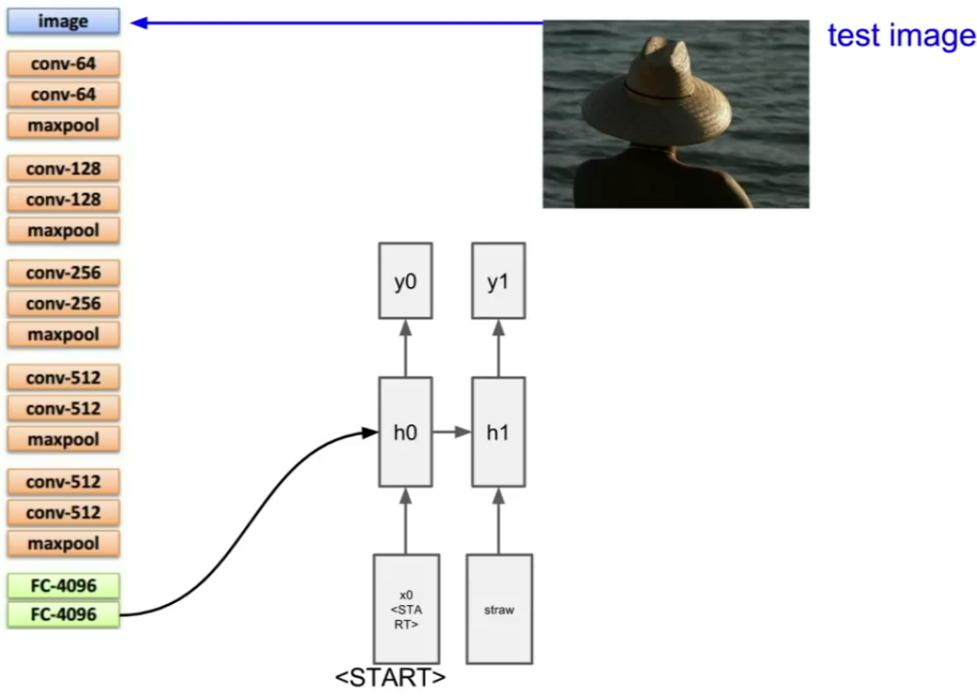


- Previously we were feeding the current input x which has its corresponding matrix W_{xh} and previous hidden state h_{t-1} which has its corresponding matrix W_{hh} .
- Now since we have the input image's feature vector v our function changes as shown above.

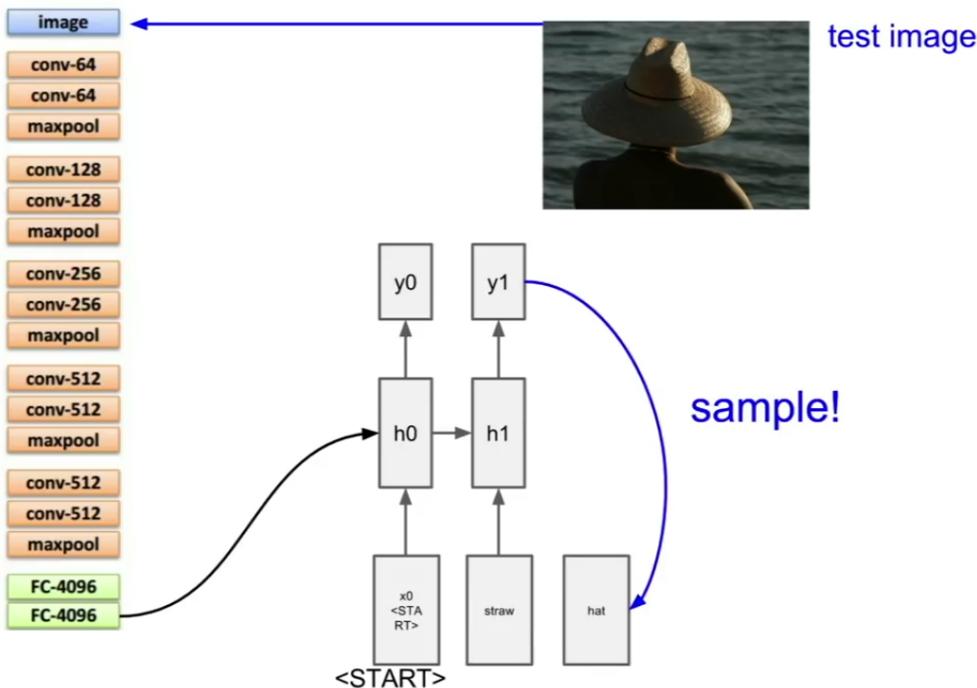


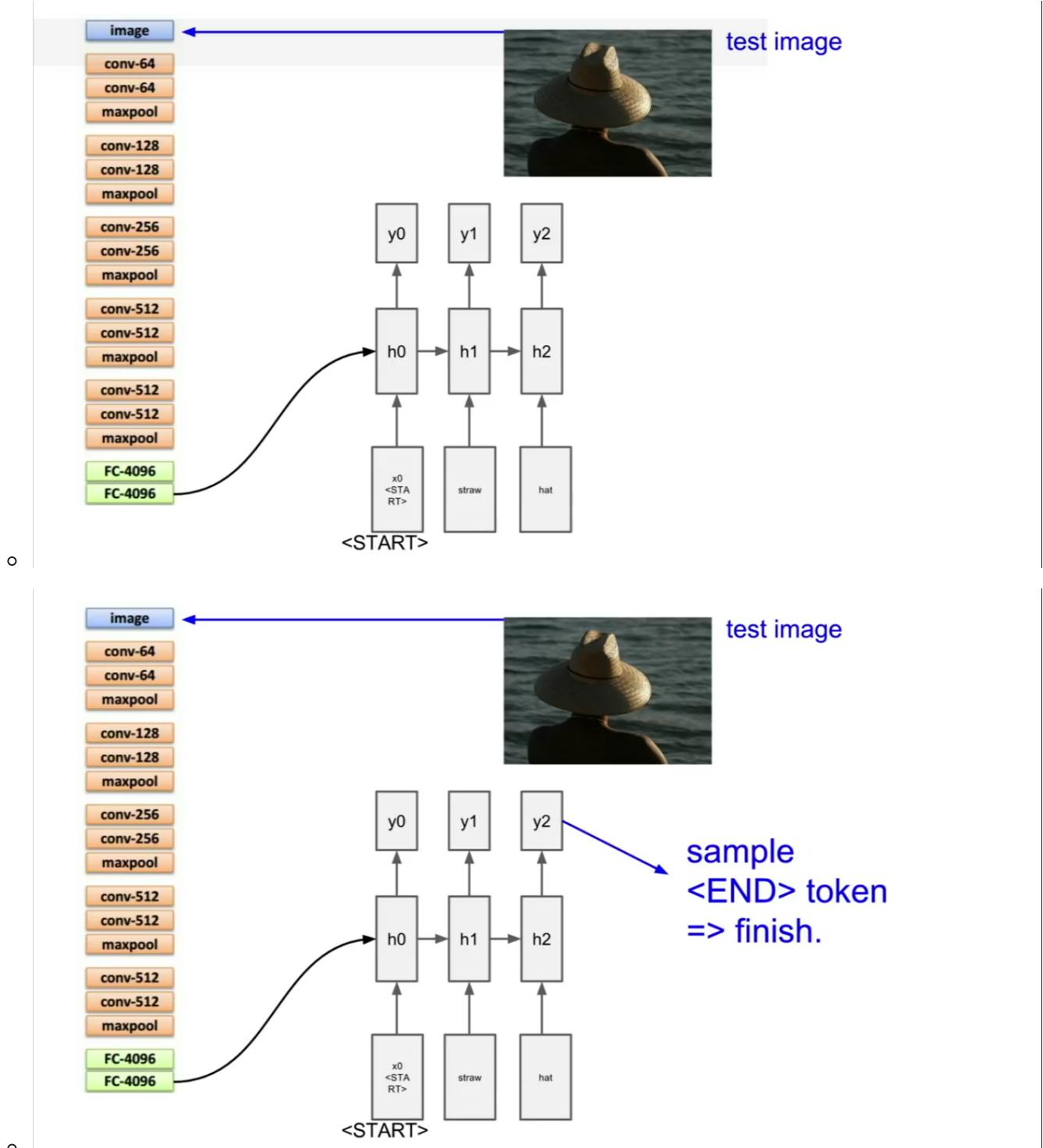
- The output(word) generated by the model at that time step will be the input to the next time step.

o



o





- We stop the generation once we generate the token, which is similar to the '.' to a sentence.
- **During Training we intentionally have the token to denote the end of caption. So the network learns that the token comes at the end of sequences.**
- Microsoft COCO is the widely used dataset which has image and captions.
- We pass the gradient even through the CNN to fine tune the entire CNN+RNN combination

Image Captioning: Example Results

Captions generated using neuraltalk2
All images are CC0 Public domain:
cat suitcase, cat tree, dog, bear,
surfers, tennis, giraffe, motorcycle



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 75 University May 4, 2017

Image Captioning: Failure Cases

Captions generated using neuraltalk2
All images are CC0 Public domain: fur coat, handstand, spider web, baseball



A woman is holding a cat in her hand



A woman standing on a beach holding a surfboard



A bird is perched on a tree branch



A person holding a computer mouse on a desk



A man in a baseball uniform throwing a ball

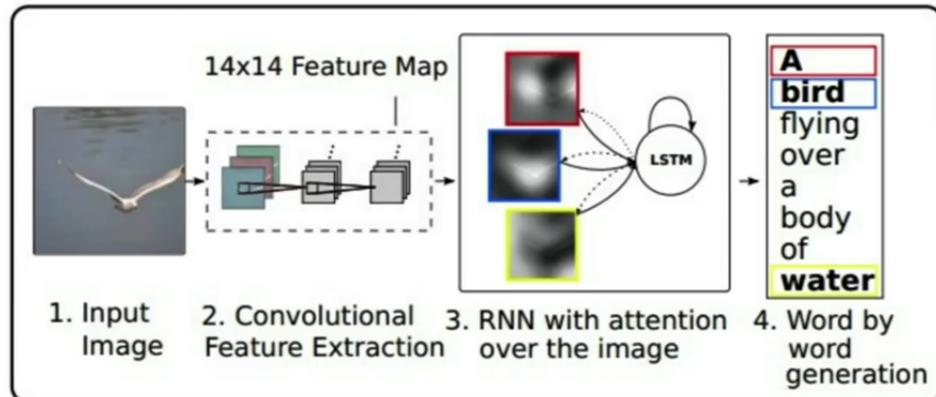
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 76 University May 4, 2017

- There are many failure cases because the model was not trained on all possible images. It was trained on many birds images so in many failure cases, the false results having 'bird' comes up often.

Image Captioning with Attention

RNN focuses its attention at a different spatial location when generating each word



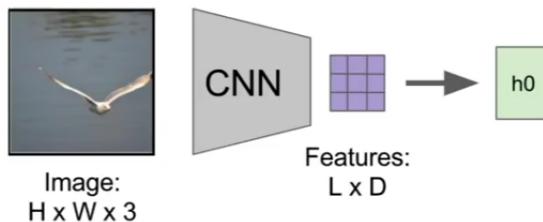
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015
Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 77 University May 4, 2017

- There are some advanced techniques for Image Captioning which also use 'Attention'

Image Captioning with Attention



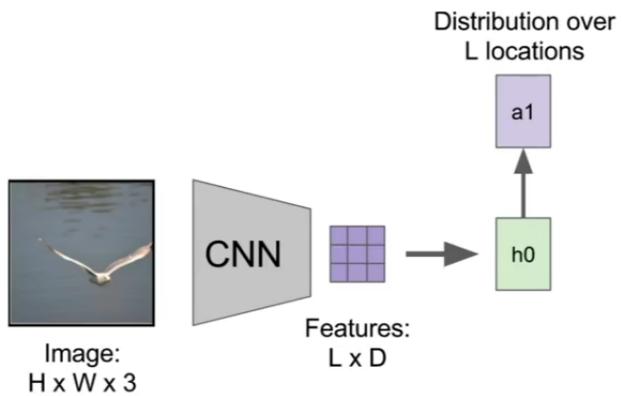
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 78 University May 4, 2017

- Instead of just passing a feature vector from the CNN, we will use a grid of vectors (attention vector?) which is fed to the RNN to generate more precise captions.

Image Captioning with Attention

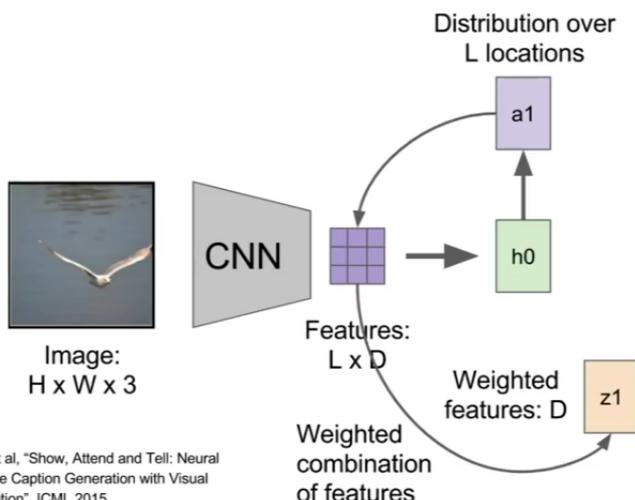


Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 79 University May 4, 2017

Image Captioning with Attention



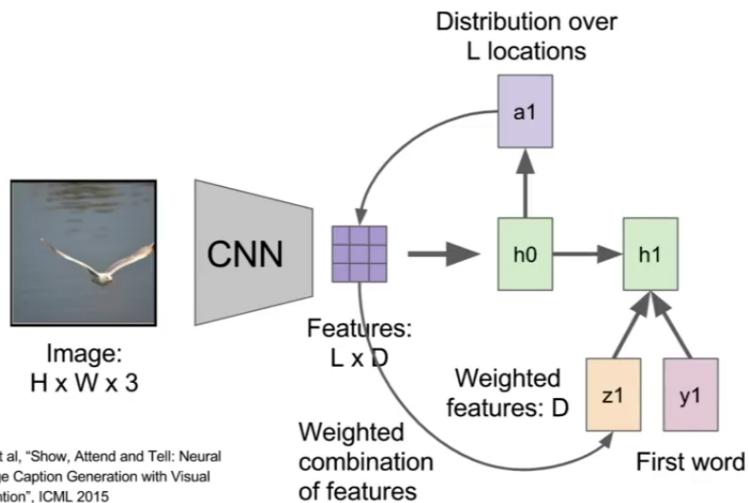
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

$$z = \sum_{i=1}^L p_i v_i$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 80 University May 4, 2017

Image Captioning with Attention

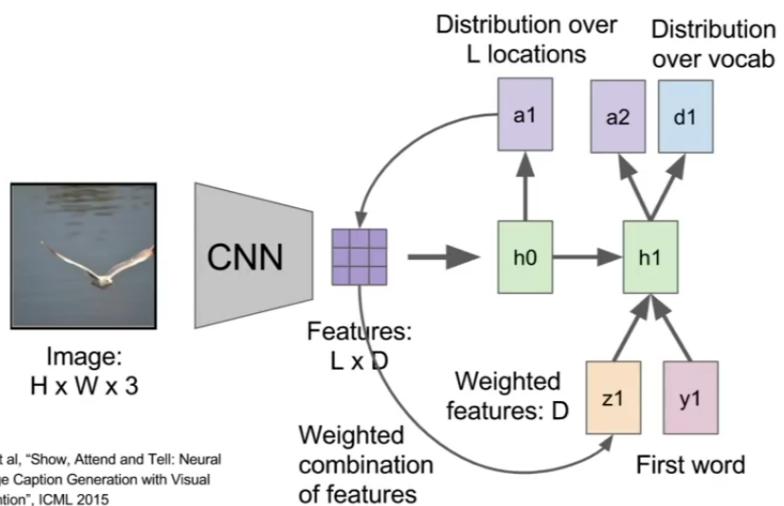


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 81 University May 4, 2017

- The summary vector z is a combination of vector for all L locations and is also fed as input the language model.

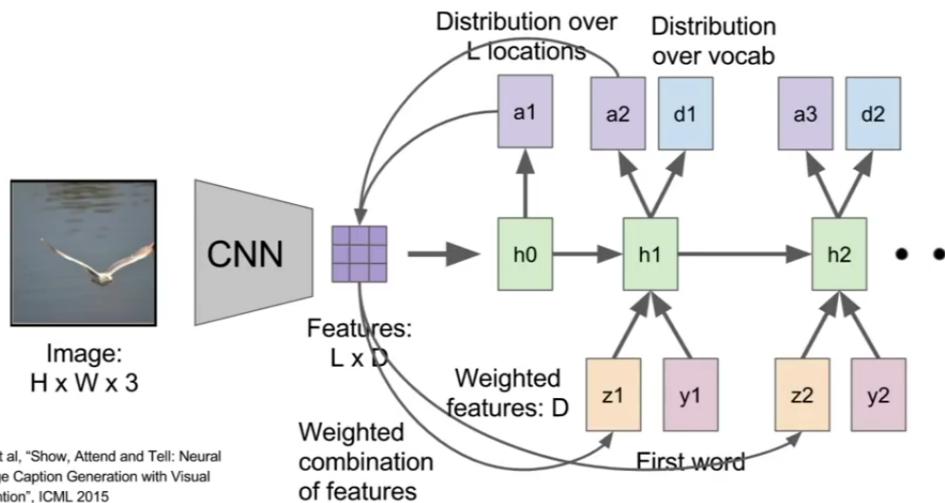
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 82 University May 4, 2017

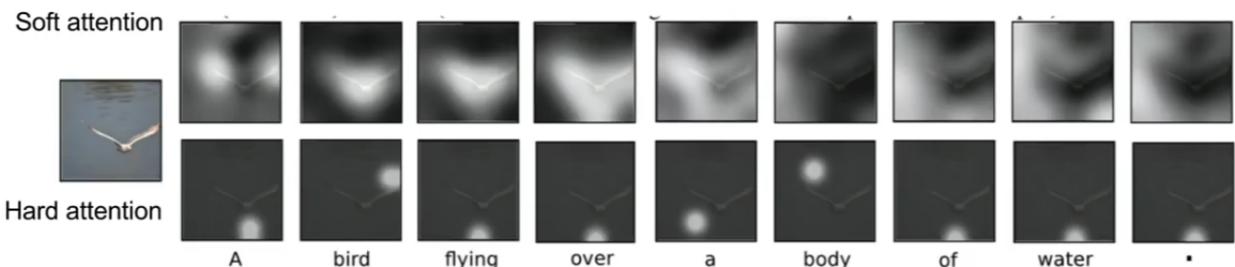
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 84 University May 4, 2017

Image Captioning with Attention



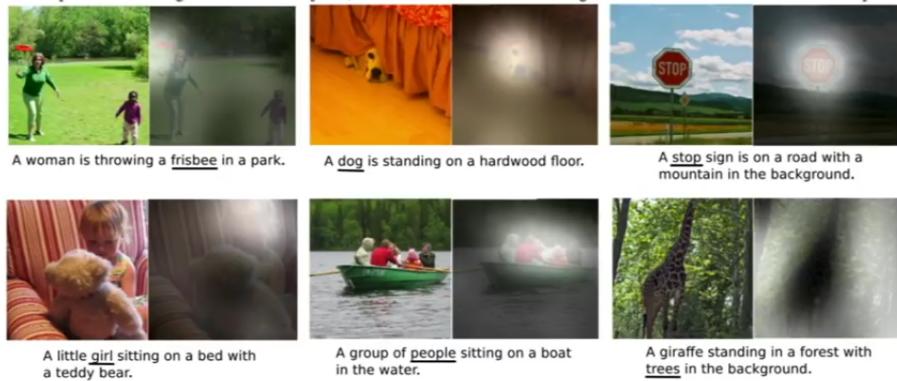
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015
Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 85 University May 4, 2017

- Soft attention is when we combine the values from all the locations during each time step.
- Hard attention is when we choose only one location to look in the image at each time step.
- Hard attention is tricky since it is not a differentiable function, so we cannot use just the vanilla backpropagation to train the model. (More details in Reinforcement Learning Lecture)

Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015
Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 86 University May 4, 2017

- We can notice that when the model generated a specific word (Eg. Frisbee) it had attention at the location where Frisbee was present.
- **The model had learnt where the attention is and used it to generate the caption.**
- Attention can also be used for various other tasks such as Visual Question Answering.

Visual Question Answering

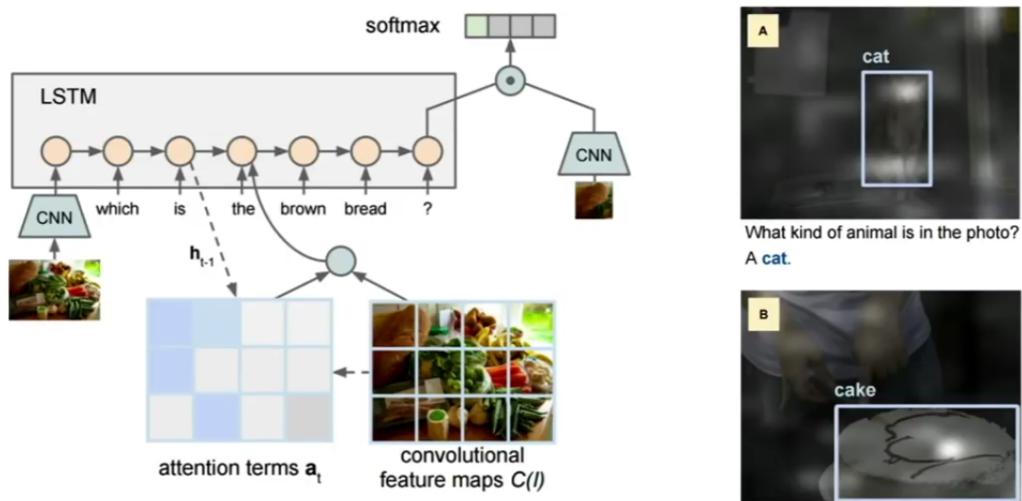


Agrawal et al, "VQA: Visual Question Answering", ICCV 2015
Zhu et al, "Visual7W: Grounded Question Answering in Images", CVPR 2016
Figure from Zhu et al. copyright IEEE 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 87 University May 4, 2017

Visual Question Answering: RNNs with Attention



Zhu et al., "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figures from Zhu et al., copyright IEEE 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 88 University May 4, 2017

Multilayer RNNs

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$ $W^l [n \times 2n]$

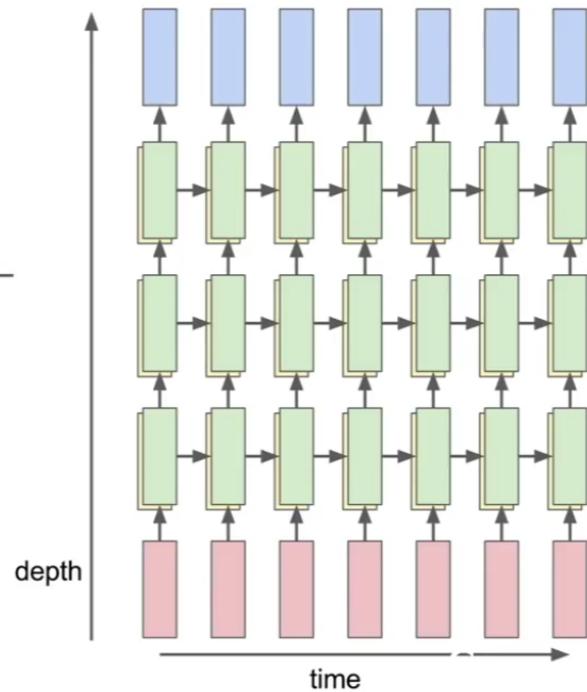
LSTM:

$$W^l [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_t^l \end{pmatrix}$$

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

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



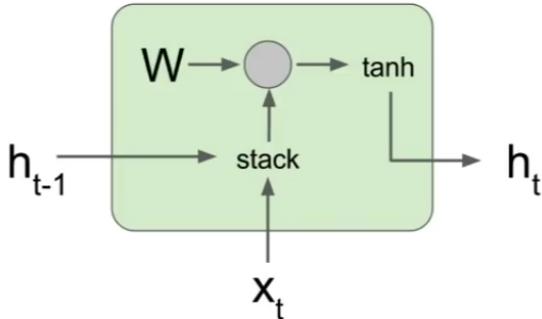
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 89 University May 4, 2017

- Here we have a depth of 3 layer RNN. Commonly we use 2 or 3 layer deep RNN.

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



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

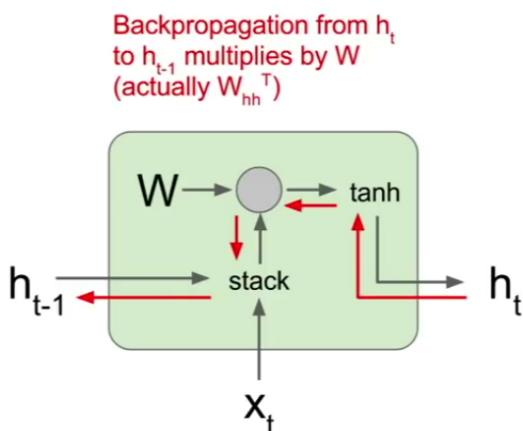
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 90 University May 4, 2017

- We are taking 2 vectors x_t and h_{t-1} and then stack them and multiply with W and then passed on to \tanh activation. This operation is one RNN cell.

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



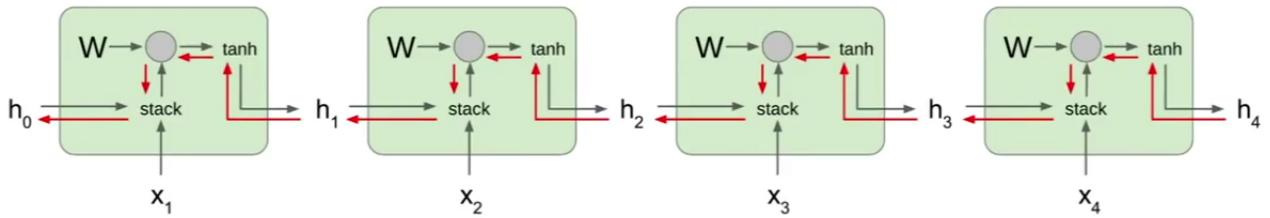
$$\begin{aligned}
 h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\
 &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\
 &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)
 \end{aligned}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 91 University May 4, 2017

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



Computing gradient of h_0 involves many factors of W (and repeated tanh)

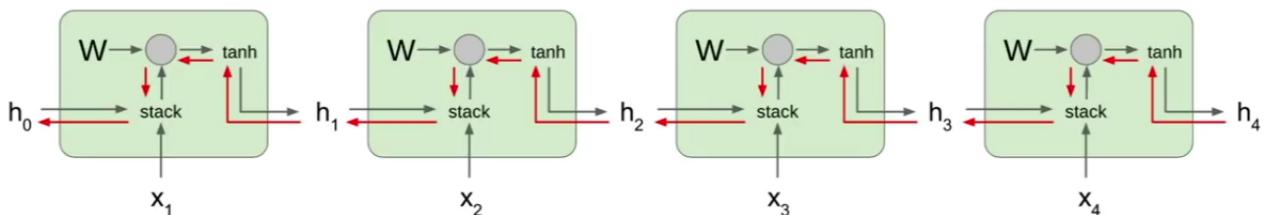
• Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 92 University, May 4, 2017

- So when we have many RNN cells the gradient for h_0 will have multiplication of multiple W^T .
- So we can either have exploding gradients or vanishing gradients problem.

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



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

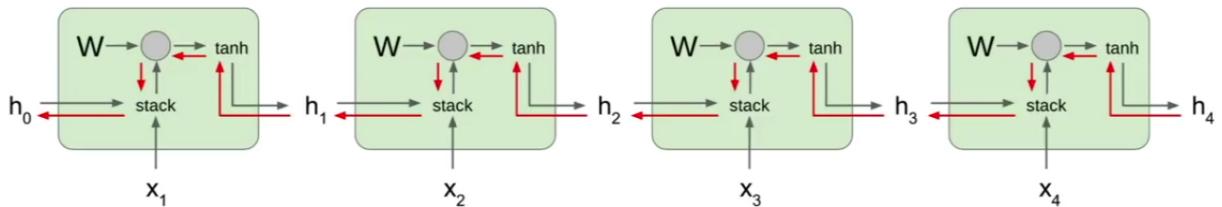
• Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 93 University, May 4, 2017

- Solution:

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



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

Gradient clipping: Scale gradient if its norm is too big

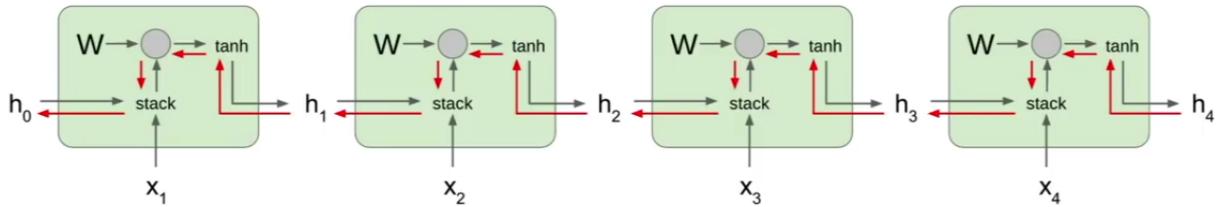
```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 94 University May 4, 2017

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



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

→ Change RNN architecture

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 95 University May 4, 2017

- Due to vanishing and exploding gradient problem, LSTM was made.

Long Short Term Memory (LSTM)

Vanilla RNN

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

LSTM

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

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

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

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

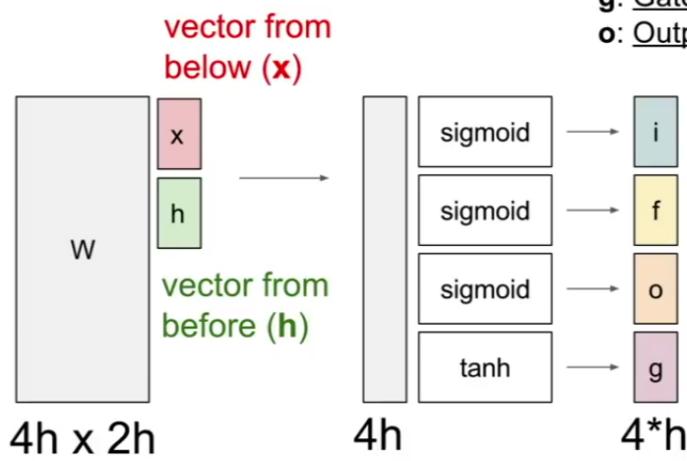
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 96 University May 4, 2017

- LSTM was conceptualized to alleviate the problem of exploding and vanishing gradients
- Instead of hacking on top of vanilla RNN, we can model the problem in a way which has better gradient flow property.
- LSTM was conceptualized in 1997.
- In Vanilla RNN we had 1 hidden state but in LSTM we have to maintain 2 hidden state.
- In LSTM we maintain the hidden state h_t and the state c_t which is the cell state.

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



f: Forget gate, Whether to erase cell
 i: Input gate, whether to write to cell
 g: Gate gate (?), How much to write to cell
 o: Output gate, How much to reveal cell

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

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

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

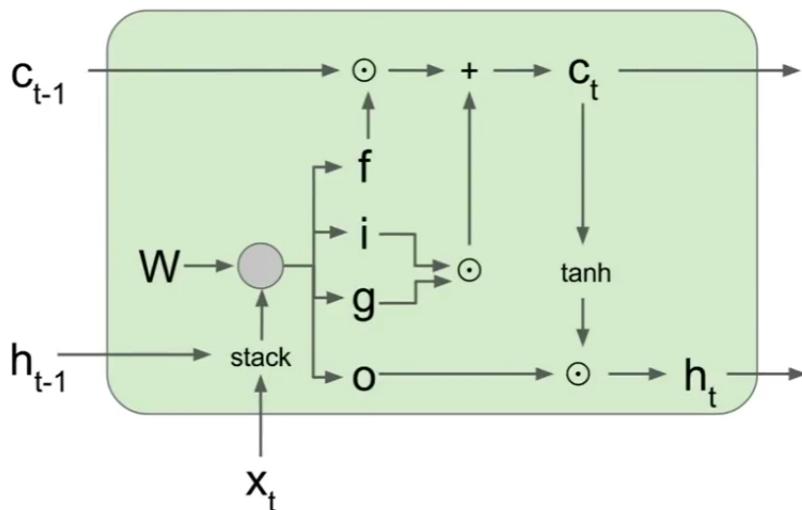
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 97 University May 4, 2017

- i, f, o use the sigmoid activation so their values will be between 0-1.
- $f \odot c_{t-1}$ is to know how much to forget from the previous cell state
-

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



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

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

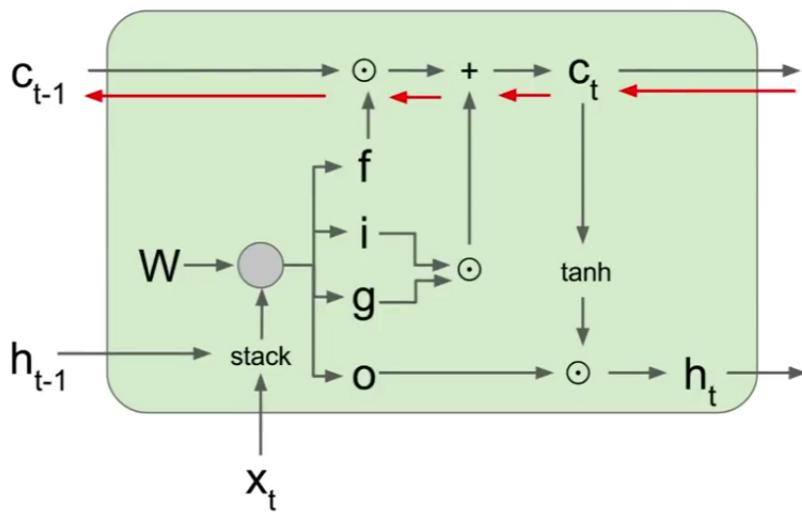
$$h_t = o \odot \tanh(c_t)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 98 University, May 4, 2017

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

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

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

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

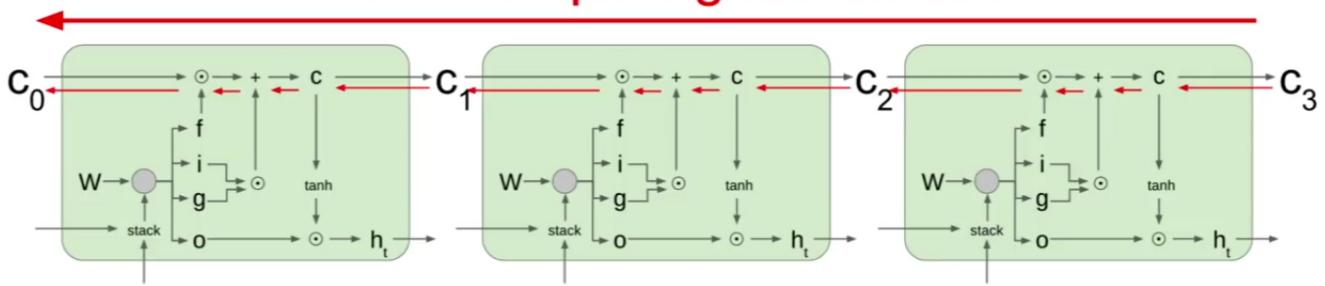
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 99 University, May 4, 2017

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Fei-Fei Li & Justin Johnson & Serena Yeung

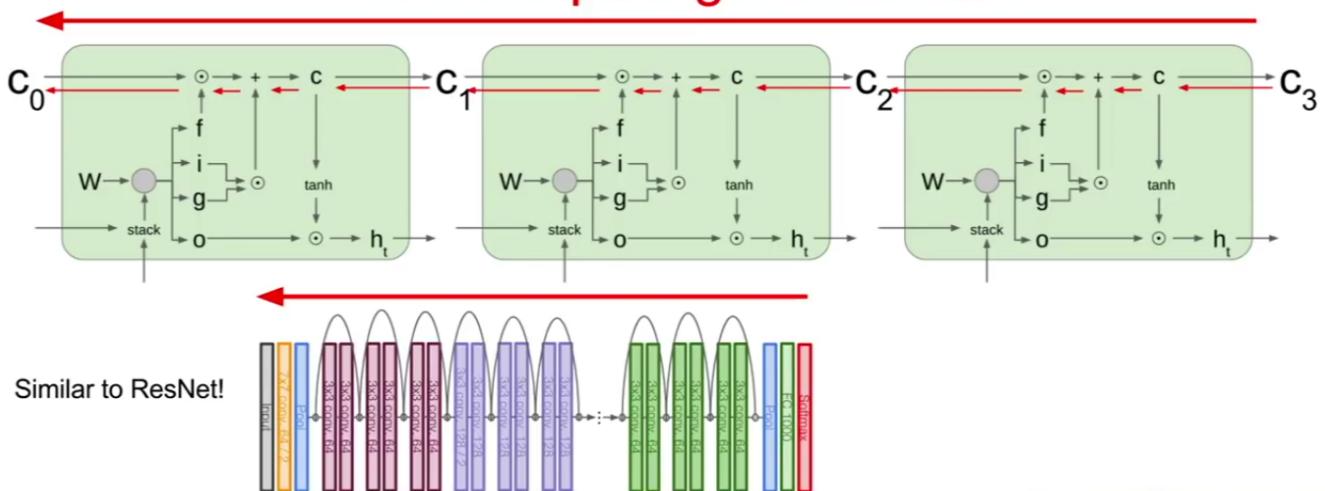
Lecture 10 - $\frac{10}{0}$ University, May 4, 2017

- Now the gradient flow will not have a recurrent multiplication of W_T to cause vanishing and exploding gradients.
- There can be problems of vanishing and exploding gradients but we can use biases which can be learnt during training and does not cause issues as Vanilla RNN.

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



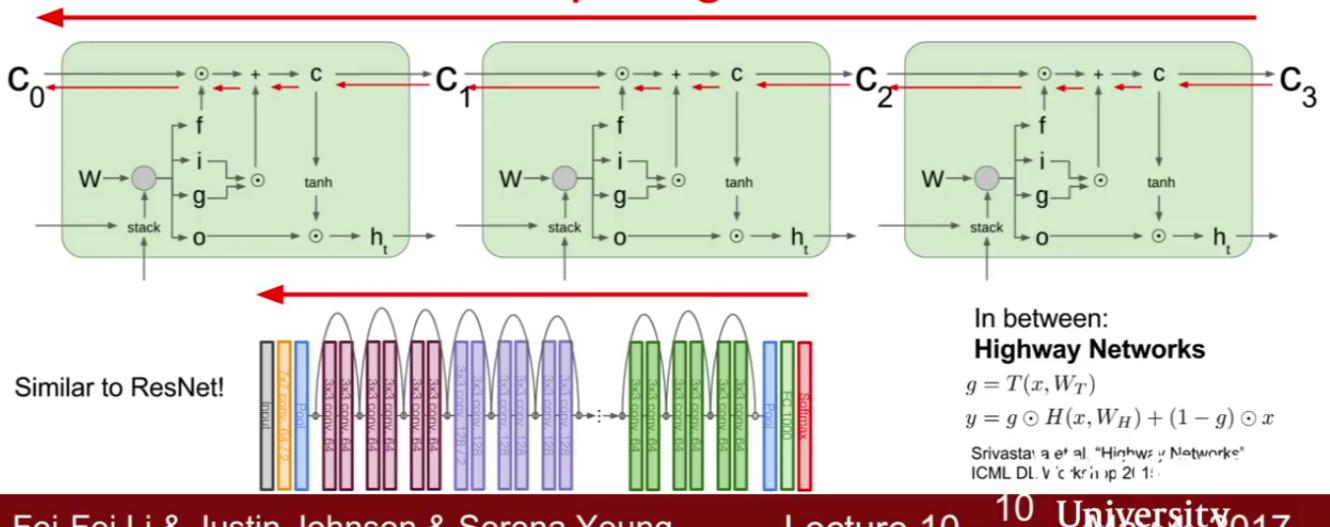
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - $\frac{10}{1}$ University, May 4, 2017

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - $\frac{10}{2}$ University, May 4, 2017

-

- Highway Networks paper was published before ResNet!

Other RNN Variants

GRU [Learning phrase representations using rnn encoder-decoder for statistical machine translation, Cho et al. 2014]

$$\begin{aligned} r_t &= \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\ z_t &= \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \end{aligned}$$

[LSTM: A Search Space Odyssey, Greff et al., 2015]

[An Empirical Exploration of Recurrent Network Architectures, Jozefowicz et al., 2015]

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_x) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hx}h_t + b_x) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hx}\tanh(h_t) + b_x) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - $\frac{10}{3}$ University, May 4, 2017

-

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
- Better understanding (both theoretical and empirical) is needed.