

10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

RNN LMs



Transformer LMs

Matt Gormley, Henry Chai Lecture 18 Mar. 19, 2025

Reminders

- Homework 6: Learning Theory & Generative Models
 - Out: Sun, Mar-16
 - Due: Sat, Mar-22, 11:59pm
 - (only two grace/late days permitted)
- Exam 2: Wed, Mar-26, 7:00 pm 9:00 pm

BACKGROUND: N-GRAM LANGUAGE MODELS

Human Language Technologies



Machine Translation

기계 번역은 특히 영어와 한국어와 같은 언어 쌍의 경우 매우 어렵습니다.

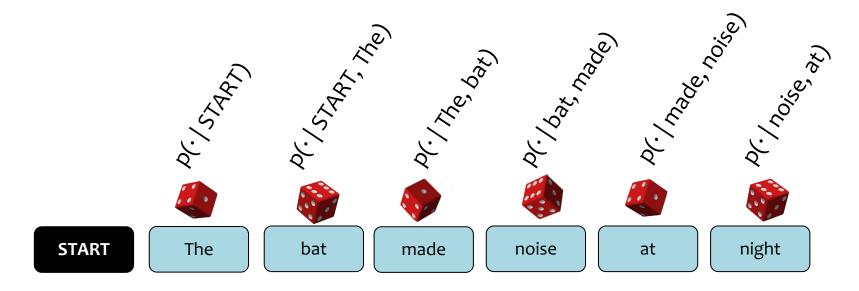
Summarization

```
Lorem ipsum dolor sit amet,

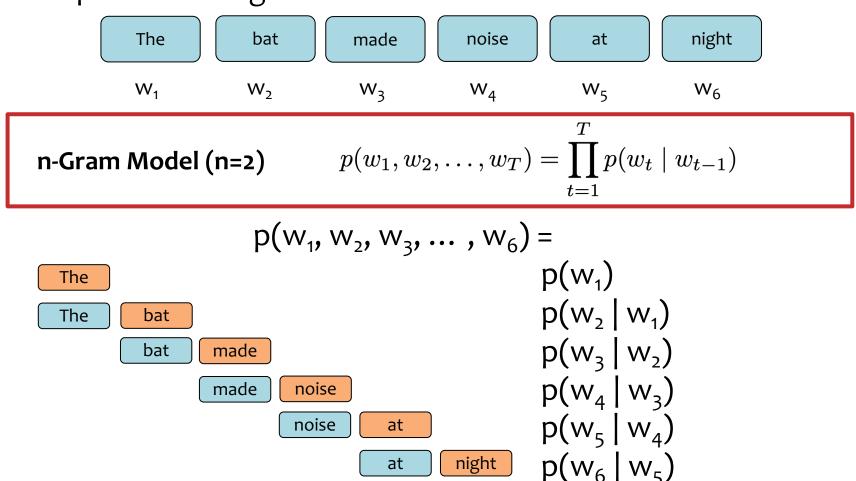
diu
lat Lorem ipsum dolor sit amet,

dil Gorini della d
```

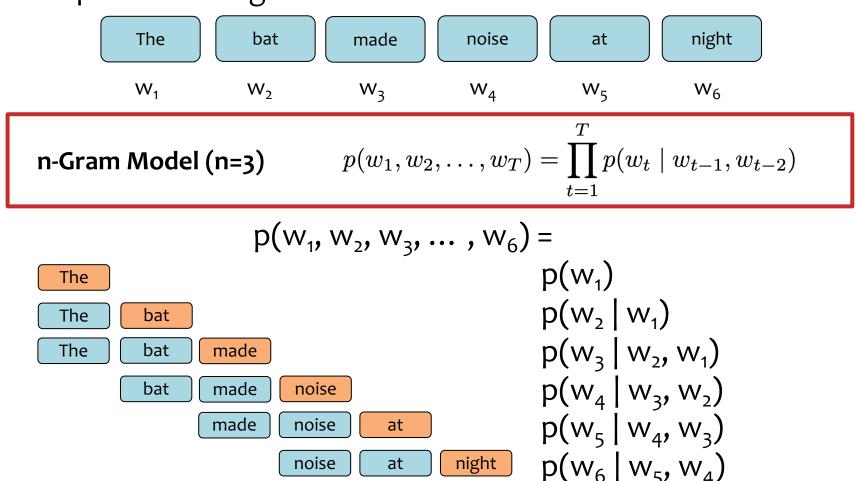
- <u>Goal</u>: Generate realistic looking sentences in a human language
- <u>Key Idea</u>: condition on the last n-1 words to sample the nth word



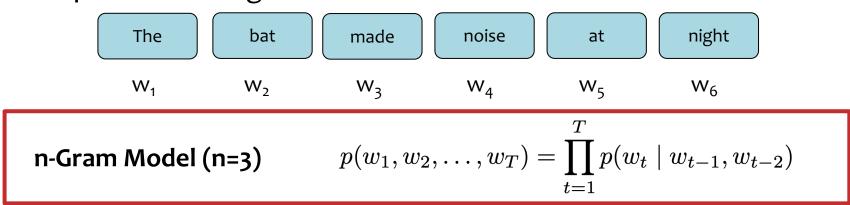
<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



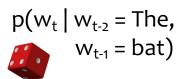
<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



$$p(w_1, w_2, ..., w_6) = p(w_1)$$
The Note: This is called a **model** because we made some **assumptions** about how many previous words to condition on (i.e. only n-1 words)

Learning an n-Gram Model

<u>Question</u>: How do we **learn** the probabilities for the n-Gram Model?



W _t	p(· ·,·)
ate	0.015
•••	

flies	0.046
•••	
zebra	0.000

$$p(w_t | w_{t-2} = made, w_{t-1} = noise)$$

W _t	p(· ·,·)
at	0.020

pollution	0.030
•••	
zebra	0.000

$$p(w_t | w_{t-2} = cows, w_{t-1} = eat)$$

w _t	p(· ·,·)
corn	0.420

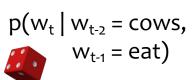
grass	0.510	
•••		

Learning an n-Gram Model

<u>Question</u>: How do we **learn** the probabilities for the n-Gram Model?

Answer: From data! Just count n-gram frequencies

```
... the cows eat grass...
... our cows eat hay daily...
... factory-farm cows eat corn...
... on an organic farm, cows eat hay and...
... do your cows eat grass or corn?...
... what do cows eat if they have...
... cows eat corn when there is no...
... which cows eat which foods depends...
... if cows eat grass...
... when cows eat corn their stomachs...
... should we let cows eat corn?...
```

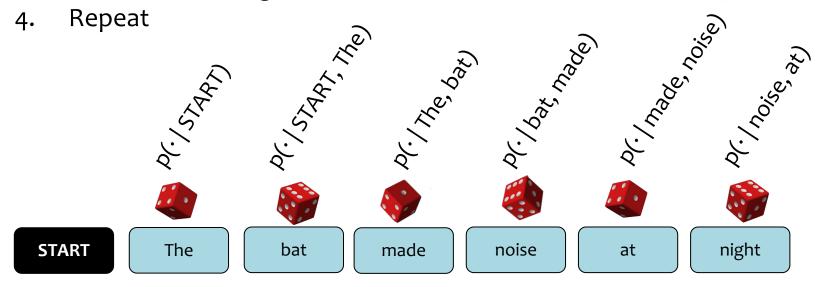


W _t	p(· ·,·)
corn	4/11
grass	3/11
hay	2/11
if	1/11
which	1/11

Sampling from a Language Model

<u>Question</u>: How do we sample from a Language Model? <u>Answer</u>:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up



Sampling from a Language Model

<u>Question</u>: How do we sample from a Language Model? Answer:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up
- 4. Repeat

Training Data (Shakespeaere)

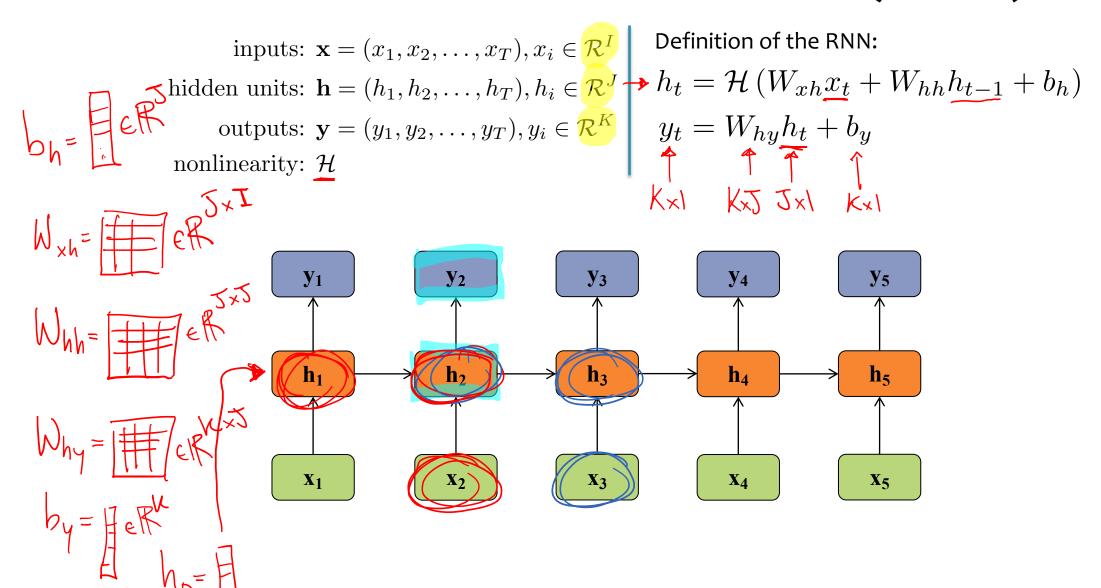
I tell you, friends, most charitable care ave the patricians of you. For your wants, Your suffering in this dearth, you may as well Strike at the heaven with your staves as lift them Against the Roman state, whose course will on The way it takes, cracking ten thousand curbs Of more strong link asunder than can ever Appear in your impediment. For the dearth, The gods, not the patricians, make it, and Your knees to them, not arms, must help.

5-Gram Model

Approacheth, denay. dungy
Thither! Julius think: grant,—O
Yead linens, sheep's Ancient,
Agreed: Petrarch plaguy Resolved
pear! observingly honourest
adulteries wherever scabbard
guess; affirmation—his monsieur;
died. jealousy, chequins me.
Daphne building. weakness: sun—
rise, cannot stays carry't,
unpurposed. prophet—like drink;
back—return 'gainst surmise
Bridget ships? wane; interim?
She's striving wet;

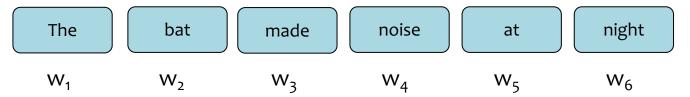
RECURRENT NEURAL NETWORK (RNN) LANGUAGE MODELS

Recurrent Neural Networks (RNNs)

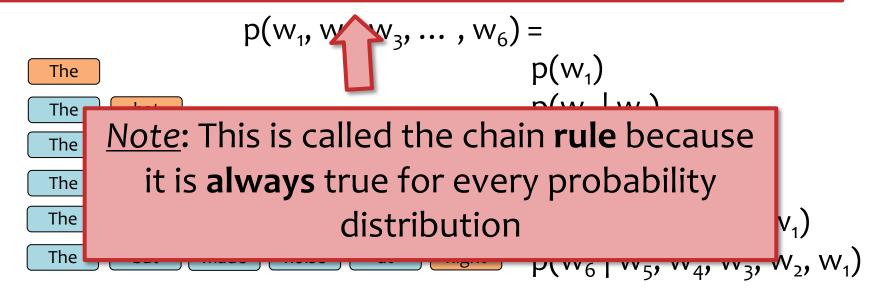


The Chain Rule of Probability

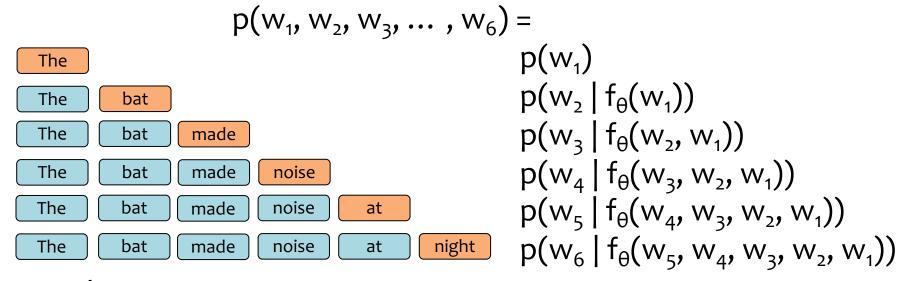
<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



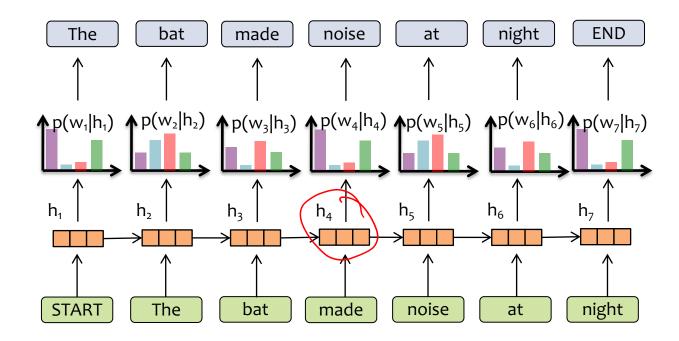
Chain rule of probability: $p(w_1, w_2, \dots, w_T) = \prod_{t=1}^T p(w_t \mid w_{t-1}, \dots, w_1)$



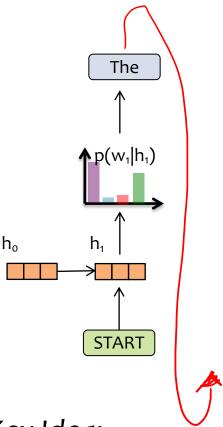
RNN Language Model:
$$p(w_1, w_2, \dots, w_T) = \prod_{t=1}^T p(w_t \mid f_{\boldsymbol{\theta}}(w_{t-1}, \dots, w_1))$$



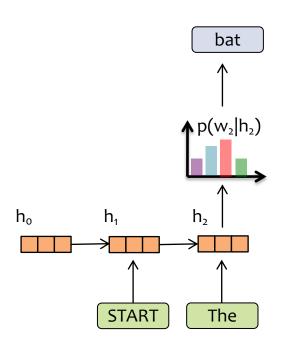
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector



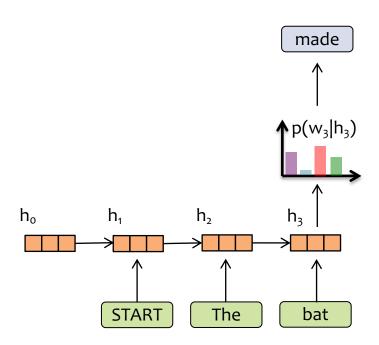
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



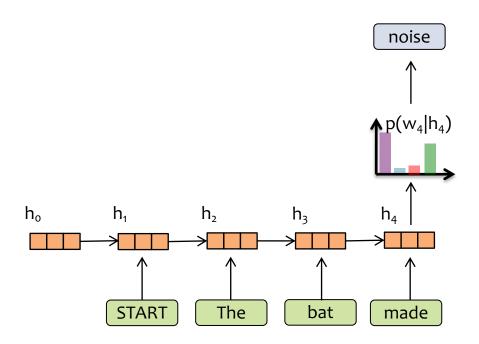
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



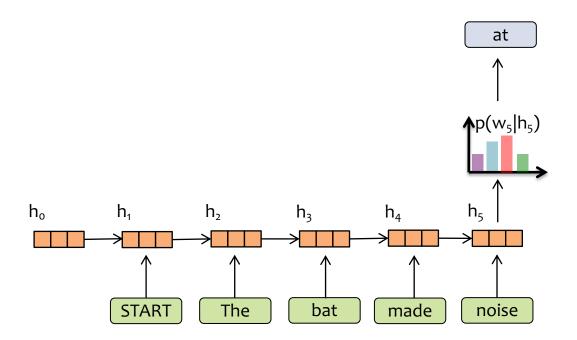
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



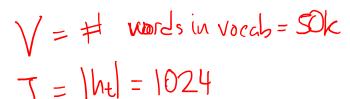
- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$

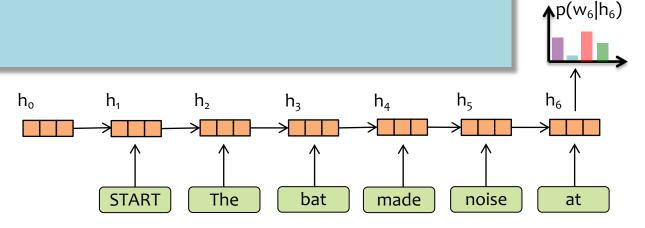
Poll 1

RNN Language Model V = # words in vocab = 50k

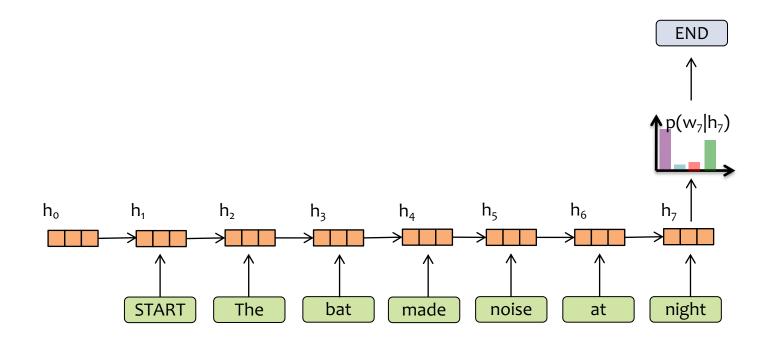


Question: How can we create a distribution $p(w_t|h_t)$ from h_t ?

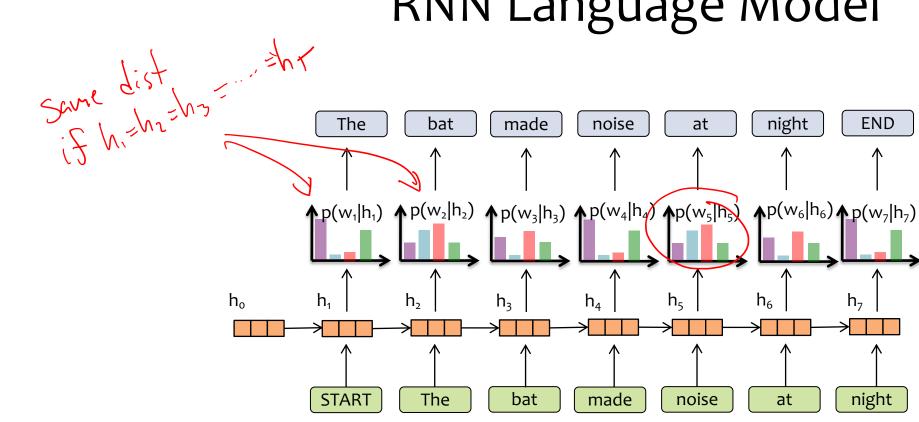
Answer:



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$



- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$

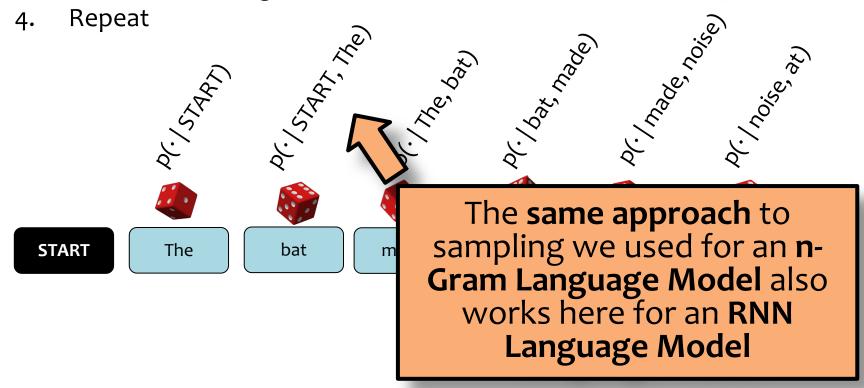


$$p(w_1, w_2, w_3, ..., w_T) = p(w_1 | h_1) p(w_2 | h_2) ... p(w_T | h_T)$$

Sampling from a Language Model

<u>Question</u>: How do we sample from a Language Model? <u>Answer</u>:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up



??

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 m
there My power to give thee but so much
service in the noble bondman here, Would
her wine.

KING LEAR: O, if you were a feeble state, 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.

??

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is ender; and, for your love, I would be as I must, for my own honour, if he re, out of my love to you, I came hither withal, that either you might stay him from his intend.

TOUCHSTONE: For my part, I had rather bear with you than bear you; yet I should bear no cross if I did bear you, for I think you have no money in your purse.

shall run into, in the is a thing of his own search and

altogether against my will.

Shakespeare's As You Like 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.

RNN-LM Sample

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is but young and tender; and, for your love, I would be loath to foil him, as I must, for my own honour, if he come in: therefore, out of my love to you, I came hither to acquaint you withal, that either you might stay him from his intendment or brook such disgrace well as he shall run into, in that it is a thing of his own search and altogether against my will.

TOUCHSTONE: For my part, I had rather bear with you than bear you; yet I should bear no cross if I did bear you, for I think you have no money in your purse.

RNN-LM Sample

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.

Shakespeare's As You Like It

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is but young and tender; and, for your love, I would be loath to foil him, as I must, for my own honour, if he come in: therefore, out of my love to you, I came hither to acquaint you withal, that either you might stay him from his intendment or brook such disgrace well as he shall run into, in that it is a thing of his own search and altogether against my will.

TOUCHSTONE: For my part, I had rather bear with you than bear you; yet I should bear no cross if I did bear you, for I think you have no money in your purse.

??

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 m
there My power to give thee but so much
service in the noble bondman here, Would
her wine.

KING LEAR: O, if you were a feeble state, 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.

??

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him well. Your brother is ender; and, for your love, I would be as I must, for my own honour, if he re, out of my love to you, I came hither withal, that either you might stay him from his intends or brook such disgrace well as he shall run into, in the

TOUCHSTONE: For my part, I had rather bear with you than bear you; yet I should bear no cross if I did bear you, for I think you have no money in your purse.

altogether against my will.

LEARNING AN RNN

Dataset for Supervised Part-of-Speech (POS) Tagging

Data: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$

Sample 1:	n	flies	p	d	$ \begin{array}{c} $
Sample 2:	n	n	v like	an	
Sample 3:	n	fly	with	n	$ \begin{array}{c c} $
Sample 4:	p with	n	you	will	

Recall

SGD and Mini-batch SGD

Algorithm 1 SGD

```
1: Initialize \theta^{(0)}
 2:
4: s = 0
 5: for t = 1, 2, ..., T do
      for i \in \mathsf{shuffle}(1, \ldots, N) do
              Select the next training point (x_i, y_i)
              Compute the gradient g^{(s)} = \nabla J_i(\theta^{(s-1)})
              Update parameters \theta^{(s)} = \theta^{(s-1)} - \eta g^{(s)}
 9:
              Increment time step s = s + 1
10:
         Evaluate average training loss J(\theta) = \frac{1}{n} \sum_{i=1}^{n} J_i(\theta)
11:
12: return \theta^{(s)}
```

Recall

SGD and Mini-batch SGD

Algorithm 1 Mini-Batch SGD

```
1: Initialize \theta^{(0)}
2: Divide examples \{1,\ldots,N\} randomly into batches \{I_1,\ldots,I_B\}
3: where \bigcup_{b=1}^{B} I_b = \{1, ..., N\} and \bigcap_{b=1}^{B} I_b = \emptyset
4: s = 0
 5: for t = 1, 2, ..., T do
      for b = 1, 2, ..., B do
              Select the next batch I_b, where m=|I_b|
              Compute the gradient g^{(s)} = \frac{1}{m} \sum_{i \in I_k} \nabla J_i(\theta^{(s)})
              Update parameters \theta^{(s)} = \theta^{(s-1)} - \eta q^{(s)}
9:
              Increment time step s = s + 1
10:
         Evaluate average training loss J(\theta) = \frac{1}{n} \sum_{i=1}^{n} J_i(\theta)
11:
12: return \theta^{(s)}
```


RNN

Algorithm 1 Elman RNN

```
1: procedure FORWARD(x_{1:T}, W_{ah}, W_{ax}, b_a, W_{yh}, b_y)
       Initialize the hidden state h_0 to zeros
2:
       for t in 1 to T do
                                                       Darams
3:
           Receive input data at time step t: x_t
4:
           Compute the hidden state update:
5:
              a_t = W_{ah} \cdot h_{t-1} + W_{ax} \cdot x_t + b_a
6:
              h_t = \sigma(a_t)
7:
           Compute the output at time step t:
8:
              y_t = W_{yh} \cdot h_t + b_y
9:
```

y_1 h_1 h_2 h_3 h_4 x_1 x_2 x_3 x_4

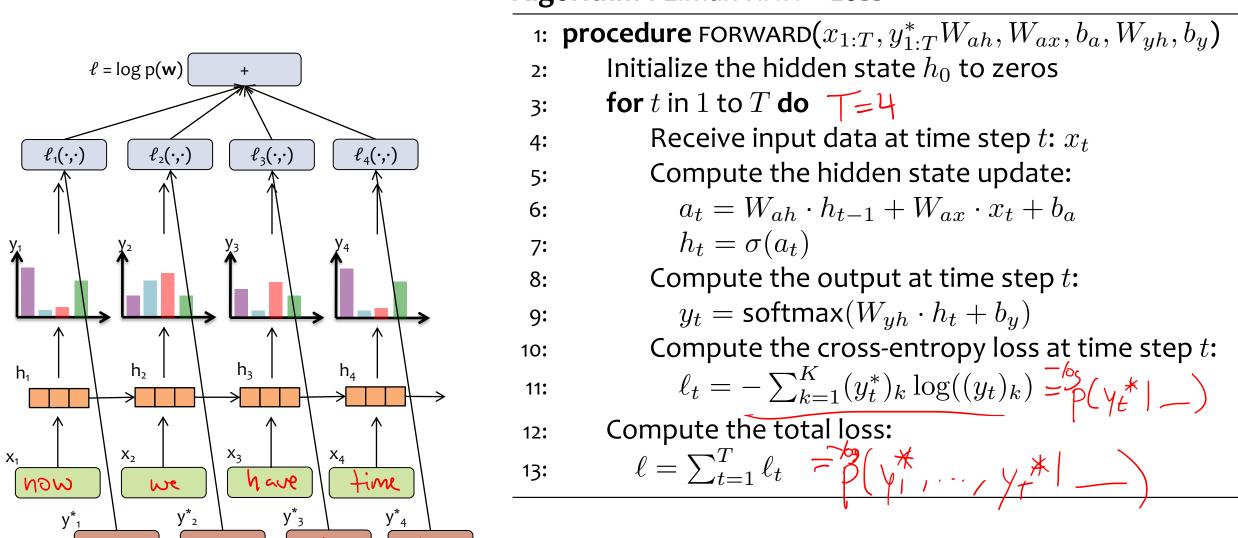
RNN

Algorithm 1 Elman RNN

```
1: procedure FORWARD(x_{1:T}, W_{ah}, W_{ax}, b_a, W_{yh}, b_y)
2: Initialize the hidden state h_0 to zeros
3: for t in 1 to T do
4: Receive input data at time step t: x_t
5: Compute the hidden state update:
6: a_t = W_{ah} \cdot h_{t-1} + W_{ax} \cdot x_t + b_a
7: h_t = \sigma(a_t)
8: Compute the output at time step t:
9: y_t = \operatorname{softmax}(W_{yh} \cdot h_t + b_y)
```

RNN + Loss

Algorithm 1 Elman RNN + Loss

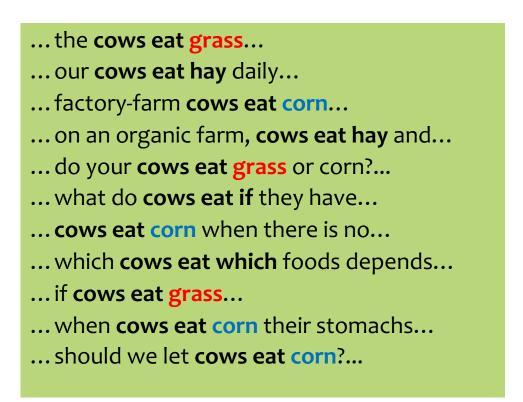


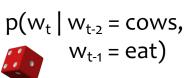
LEARNING AN RNN-LM

Learning a Language Model

<u>Question</u>: How do we **learn** the probabilities for the n-Gram Model?

Answer: From data! Just count n-gram frequencies





W _t	p(· ·,·)
corn	4/11
grass	3/11
hay	2/11
if	1/11
which	1/11

MLE for n-gram LM

- This counting method gives us the maximum likelihood estimate of the n-gram LM parameters
- We can derive it in the usual way:
 - Write the likelihood of the sentences under the n-gram LM
 - Set the gradient to zero
 and impose the constraint that the probabilities sumto-one
 - Solve for the MLE

Learning a Language Model

MLE for Deep Neural LM

- We can also use maximum likelihood estimation to learn the parameters of an RNN-LM or Transformer-LM too!
- But not in closed form instead we follow a different recipe:
 - Write the likelihood of the sentences under the Deep Neural LM model
 - Compute the gradient of the (batch) likelihood w.r.t.
 the parameters by AutoDiff
 - Follow the negative gradient using Mini-batch SGD (or your favorite optimizer)

MLE for n-gram LM

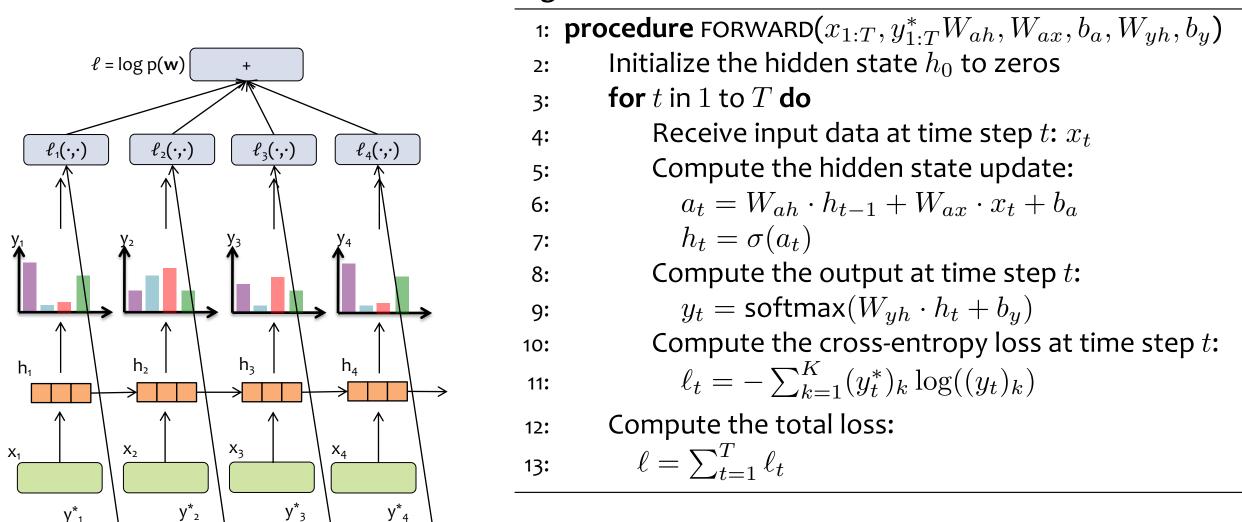
- This counting method gives us the maximum likelihood estimate of the n-gram LM parameters
- We can derive it in the usual way:
 - Write the likelihood of the sentences under the n-gram LM
 - Set the gradient to zero

 and impose the constraint that the probabilities sumto-one
 - Solve for the MLE

RNN + Loss

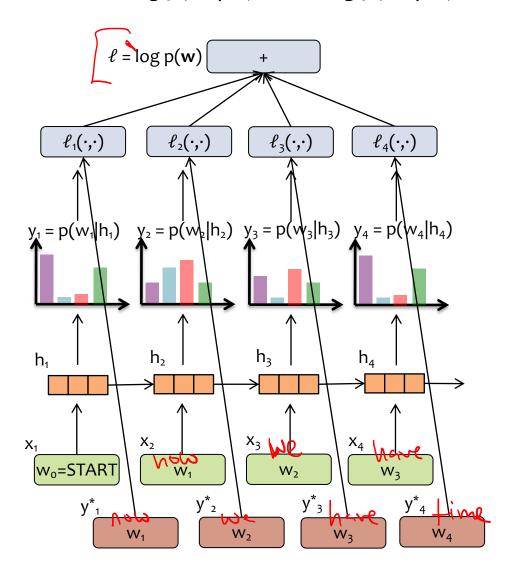
How can we use this to compute the loss for an RNN-LM?

Algorithm 1 Elman RNN + Loss



RNN-LM + Loss

How can we use this to compute the loss for an RNN-LM?



Algorithm 1 Elman RNN + Loss

1: procedure FORWARD $(x_{1:T}, y_{1:T}^* W_{ah}, W_{ax}, b_a, W_{yh}, b_y)$ Initialize the hidden state h_0 to zeros 2: **for** t in 1 to T **do** Receive input data at time step t: x_t Compute the hidden state update: 5: $a_t = W_{ah} \cdot h_{t-1} + W_{ax} \cdot x_t + b_a$ 6:

$$a_t = W_{ah} \cdot h_{t-1} + W_{ax} \cdot x_t + b_a$$

7:
$$h_t = \sigma(a_t)$$

Compute the output at time step *t*: 8:

9:
$$y_t = \operatorname{softmax}(W_{yh} \cdot h_t + b_y)$$

Compute the cross-entropy loss at time step t: 10:

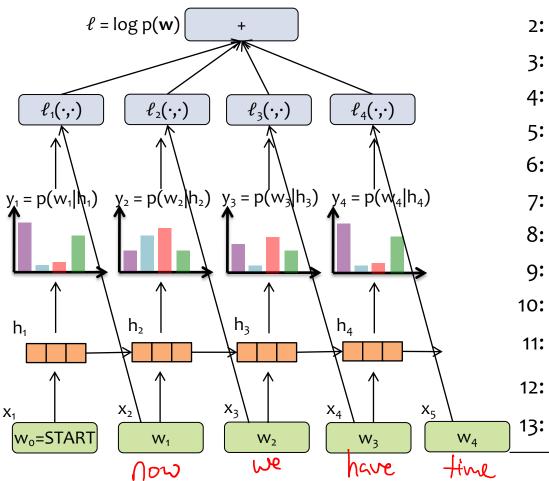
11:
$$\ell_t = -\sum_{k=1}^K (y_t^*)_k \log((y_t)_k)$$

Compute the total loss: 12:

13:
$$\ell = \sum_{t=1}^T \ell_t$$

RNN-LM + Loss

How can we use this to compute the loss for an RNN-LM?



Algorithm 1 Elman RNN + Loss

- 1: **procedure** FORWARD $(x_{1:T}, y_{1:T}^* W_{ah}, W_{ax}, b_a, W_{yh}, b_y)$
- 2: Initialize the hidden state h_0 to zeros
- 3: for t in 1 to T do
- 4: Receive input data at time step t: x_t
- 5: Compute the hidden state update:

$$a_t = W_{ah} \cdot h_{t-1} + W_{ax} \cdot x_t + b_a$$

$$h_t = \sigma(a_t)$$

Compute the output at time step *t*:

$$y_t = \mathsf{softmax}(W_{yh} \cdot h_t + b_y)$$

Compute the cross-entropy loss at time step t:

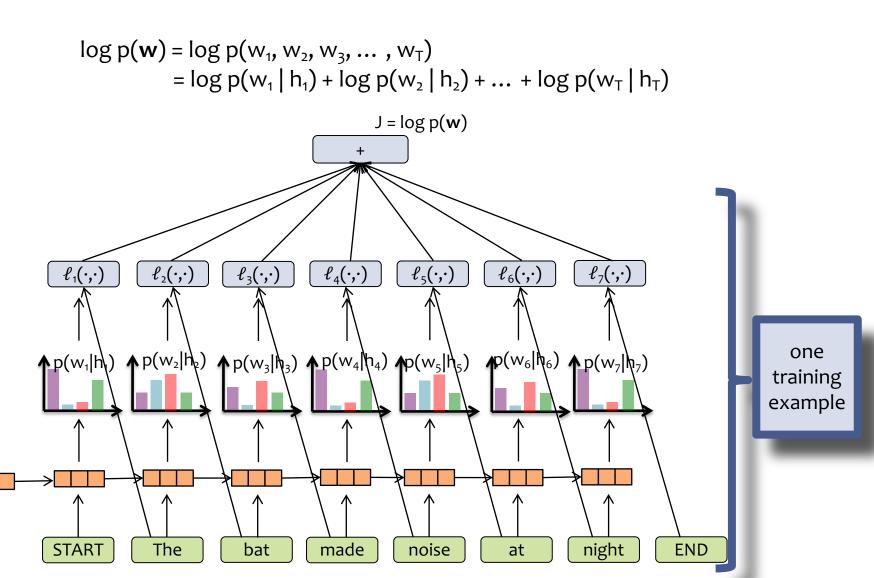
$$\ell_t = -\sum_{k=1}^{K} (y_t^*)_k \log((y_t)_k)$$

Compute the total loss:

$$\ell = \sum_{t=1}^{T} \ell_t$$

Learning an RNN-LM

- Each training example is a sequence (e.g. sentence), so we have training data D = {w⁽¹⁾, w⁽²⁾, ..., w^(N)}
- The objective function for a Deep LM (e.g. RNN-LM or Tranformer-LM) is typically the loglikelihood of the training examples: $J(\mathbf{\theta}) = \Sigma_i \log p_{\mathbf{\theta}}(\mathbf{w}^{(i)})$
- We train by mini-batch SGD (or your favorite flavor of mini-batch SGD)



LARGE LANGUAGE MODELS

How large are LLMs?

Comparison of some recent large language models (LLMs)

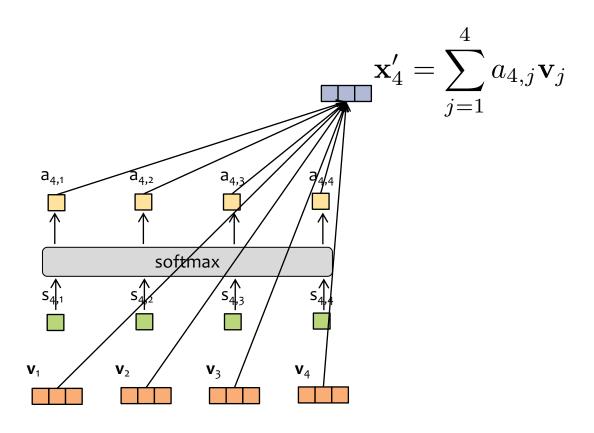
Model	Creators	Year of release	Training Data (# tokens)	Model Size (# parameters)
GPT-2	OpenAl	2019	~10 billion (40Gb)	1.5 billion
GPT-3	OpenAl	2020	300 billion	175 billion
PaLM	Google	2022	780 billion	540 billion
Chinchilla	DeepMind	2022	1.4 trillion	70 billion
LaMDA (cf. Bard)	Google	2022	1.56 trillion	137 billion
LLaMA	Meta	2023	1.4 trillion	65 billion
LLaMA-2	Meta	2023	2 trillion	70 billion
GPT-4	OpenAl	2023	?	? (1.76 trillion)
Gemini (Ultra)	Google	2023	?	? (1.5 trillion)
LLaMA-3	Meta	2024	15 trillion	405 billion

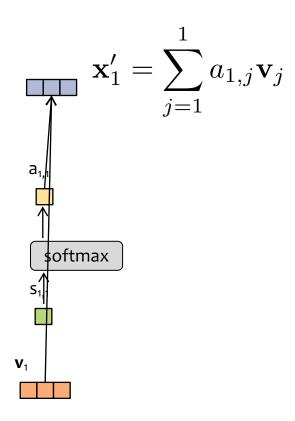
What is ChatGPT?

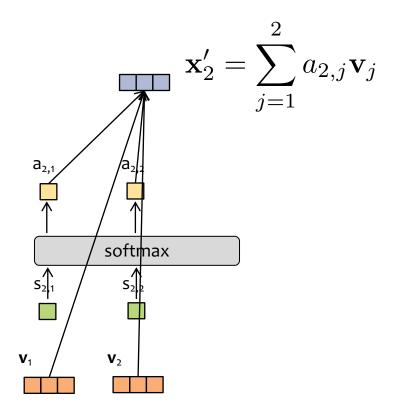
- ChatGPT is a large (in the sense of having many parameters) language model, fine-tuned to be a dialogue agent
- The base language model was originally GPT-3.5 which was trained on a large quantity of text

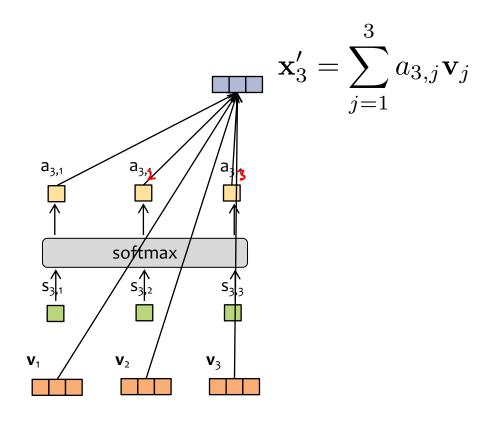
Transformer Language Models

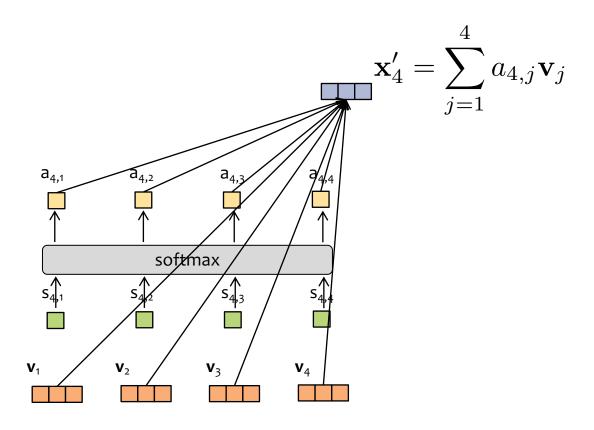
MODEL: GPT

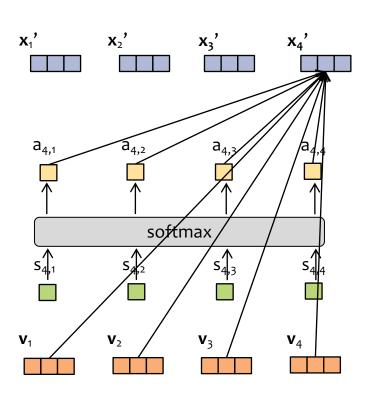










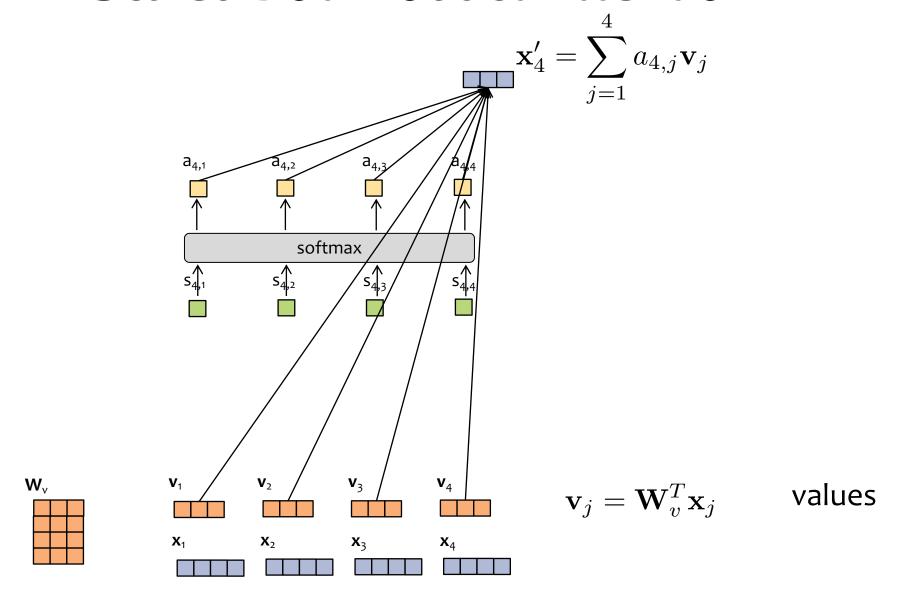


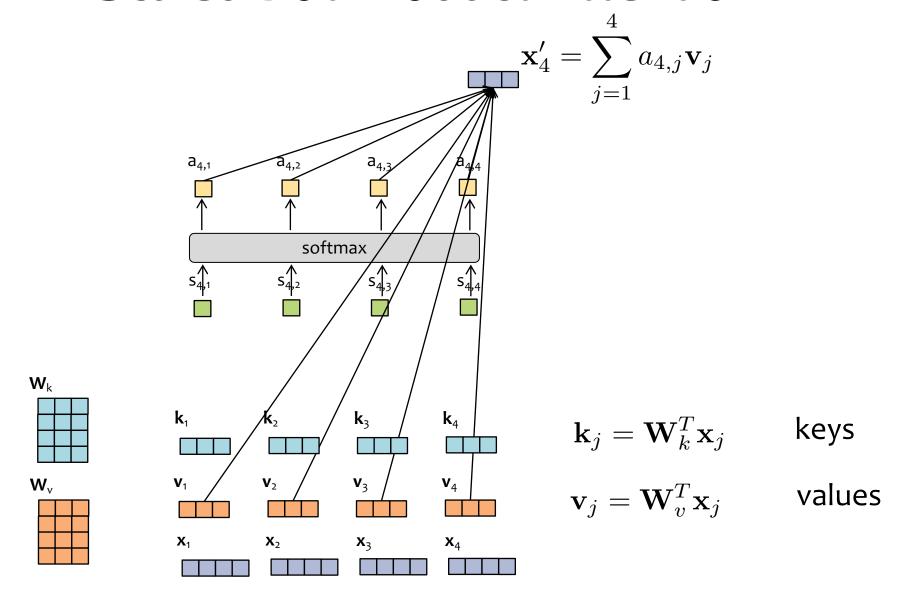
$$\mathbf{x}_t' = \sum_{j=1}^t a_{t,j} \mathbf{v}_j$$

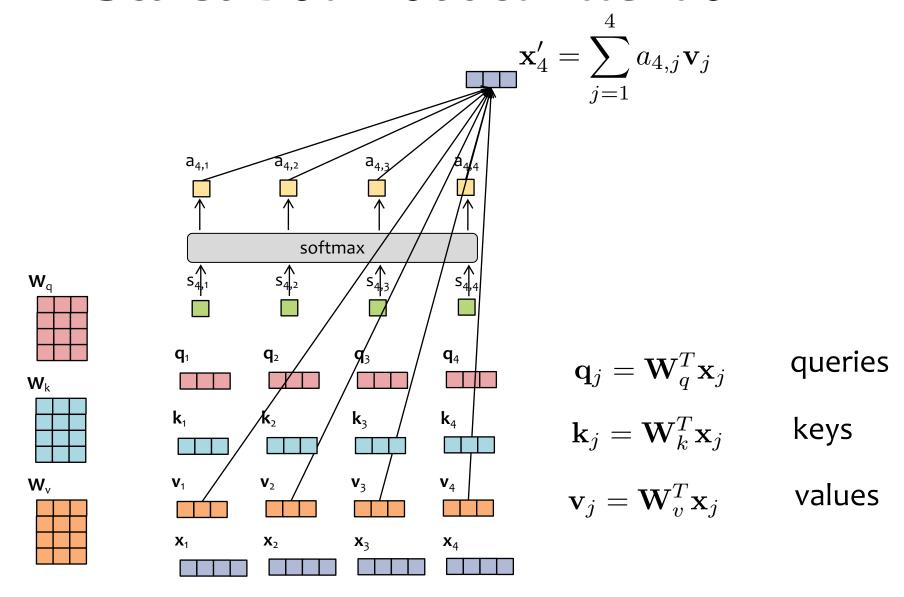
attention weights

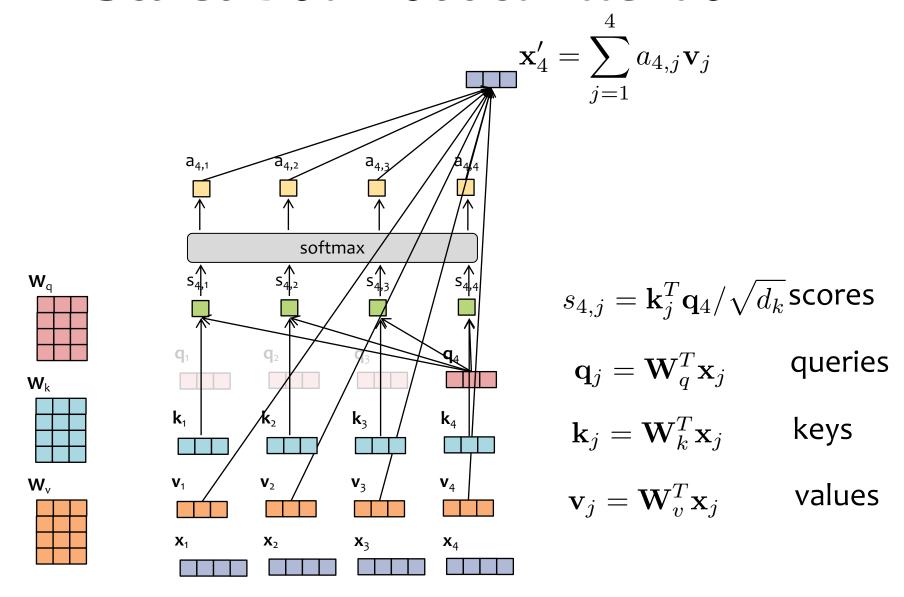
scores

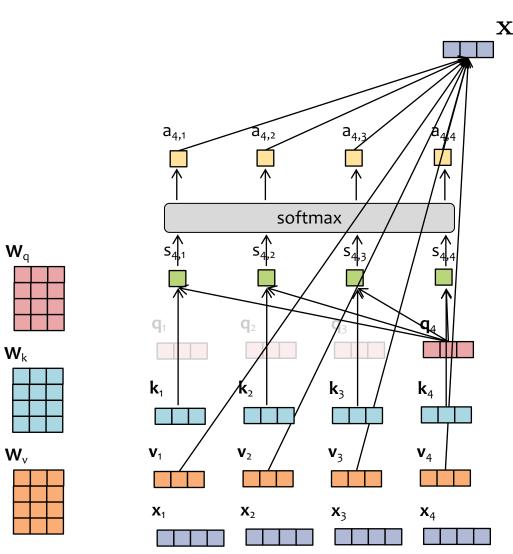
values











$$\mathbf{x}_4' = \sum_{j=1}^4 a_{4,j} \mathbf{v}_j$$

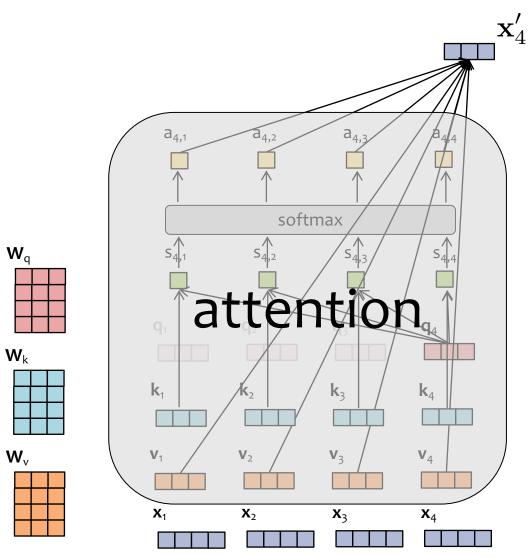
 $\mathbf{a}_4 = \mathsf{softmax}(\mathbf{s}_4)$ attention weights

$$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k}$$
 scores

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$$
 queries

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$$
 keys

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$$
 values



$$\mathbf{x}_4' = \sum_{j=1}^4 a_{4,j} \mathbf{v}_j$$

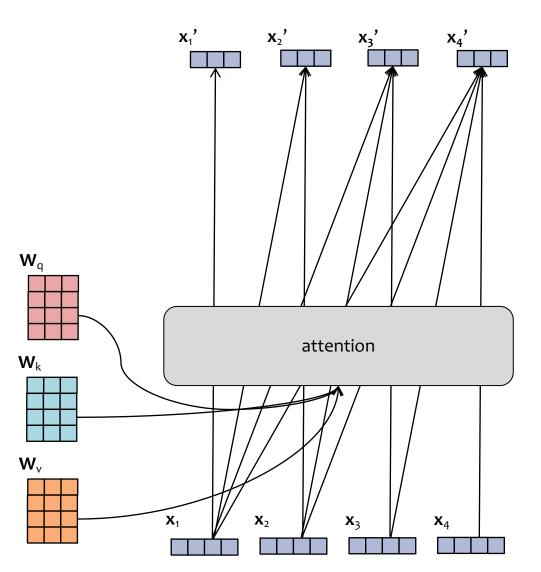
 $\mathbf{a}_4 = \mathsf{softmax}(\mathbf{s}_4)$ attention weights

$$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k}$$
 scores

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$$
 queries

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$$
 keys

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$$
 values



$$\mathbf{x}_t' = \sum_{j=1}^t a_{t,j} \mathbf{v}_j$$

 $\mathbf{a}_t = \mathsf{softmax}(\mathbf{s}_t)$ attention weights

$$s_{t,j} = \mathbf{k}_j^T \mathbf{q}_t / \sqrt{d_k}$$
 scores

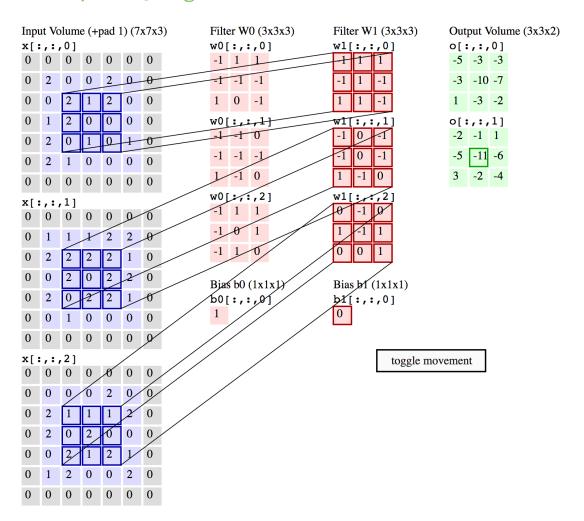
$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$$
 queries $\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$ keys

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$$
 keys

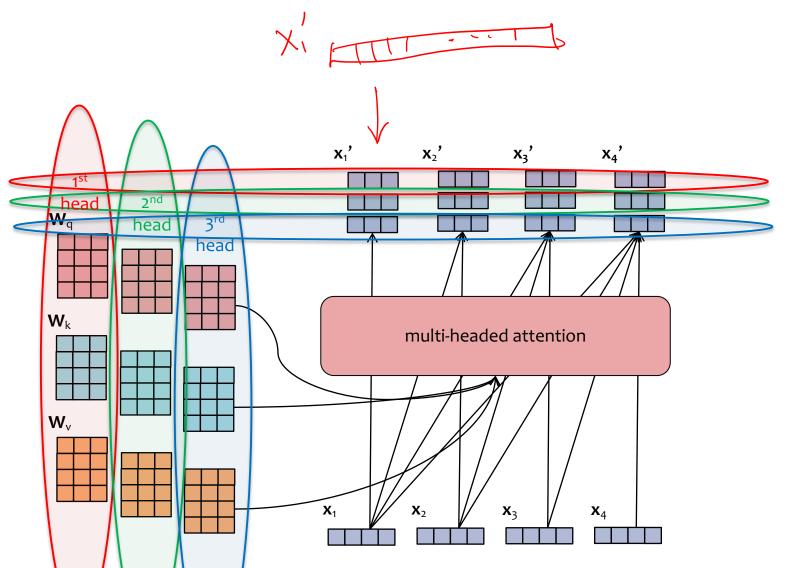
$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$$
 values

Animation of 3D Convolution

http://cs231n.github.io/convolutional-networks/



Multi-headed Attention

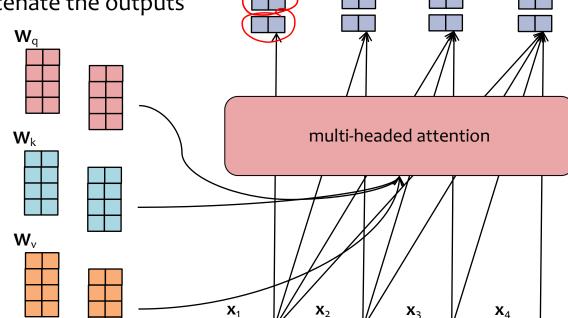


- Just as we can have multiple channels in a convolution layer, we can use multiple heads in an attention layer
- Each head gets its own parameters
- We can concatenate all the outputs to get a single vector for each time step

To ensure the dimension of the **input** embedding \mathbf{x}_t is the same as the **output** embedding \mathbf{x}_t , Transformers usually choose the embedding sizes and number of heads appropriately:

Multi-headed Attention

- $d_{model} = dim. of inputs = 4$
- $d_k = dim. of each output$
- h = # of heads = 2
- Choose $d_k = d_{model} / h = \frac{4}{2} = 2$ • Then concatenate the outputs



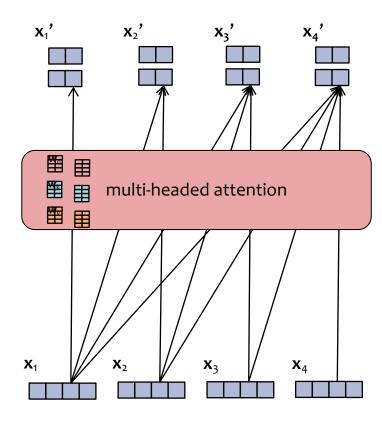
 \mathbf{x}_{2}

 X_3'

- Just as we can have multiple channels in a convolution layer, we can use multiple heads in an attention layer
- Each head gets its own parameters
- We can concatenate all the outputs to get a single vector for each time step

- To ensure the dimension of the input embedding x_t is the same as the output embedding x_t , Transformers usually choose the embedding sizes and number of heads appropriately:
 - $d_{model} = dim. of inputs$
 - $d_k = dim. of each output$
 - h = # of heads
 - Choose $d_k = d_{model} / h$
- Then concatenate the outputs

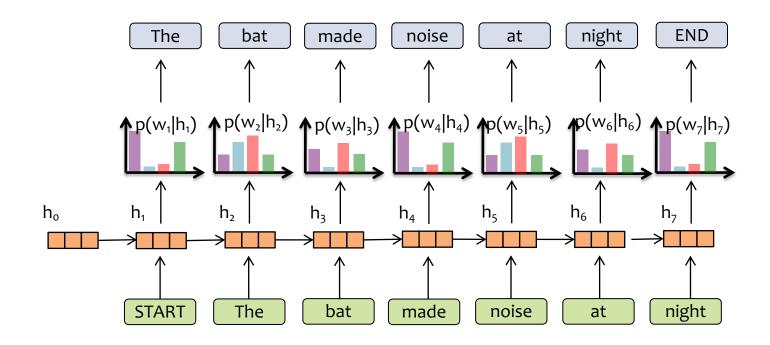
Multi-headed Attention



- Just as we can have multiple channels in a convolution layer, we can use multiple heads in an attention layer
- Each head gets its own parameters
- We can concatenate all the outputs to get a single vector for each time step

Recall

RNN Language Model



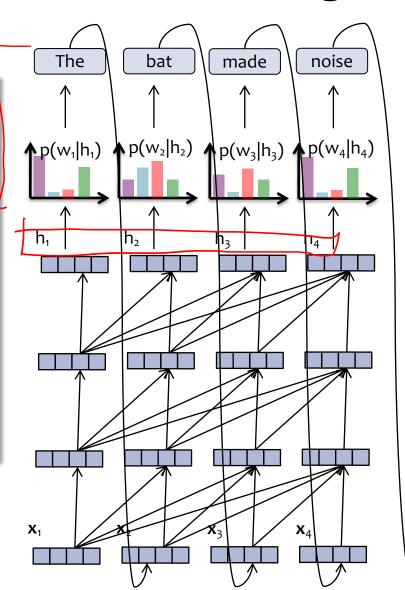
Key Idea:

- (1) convert all previous words to a fixed length vector
- (2) define distribution $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$ that conditions on the vector $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$

Transformer Language Model

Important!

- RNN computation graph grows linearly with the number of input tokens
- Transformer-LM computation graph grows quadratically with the number of input tokens



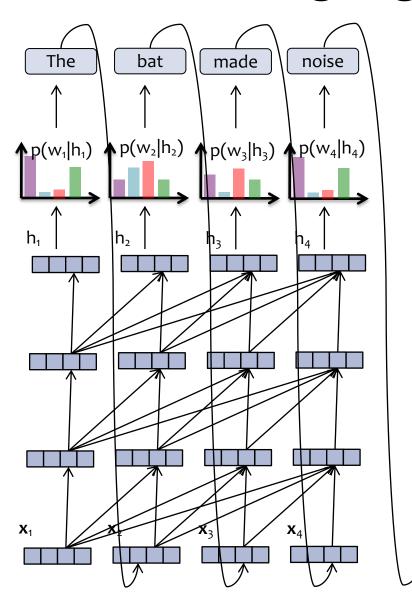
Each hidden vector looks back at the hidden vectors of the current and previous timesteps in the previous layer.

The language model part is just like an RNN-LM!

Transformer Language Model

Important!

- RNN computation graph grows linearly with the number of input tokens
- Transformer-LM computation graph grows quadratically with the number of input tokens



Each layer of a Transformer LM consists of several **sublayers**:

- 1. attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

Each hidden vector looks back at the hidden vectors of the current and previous timesteps in the previous layer.

The language model part is just like an RNN-LM!

Layer Normalization

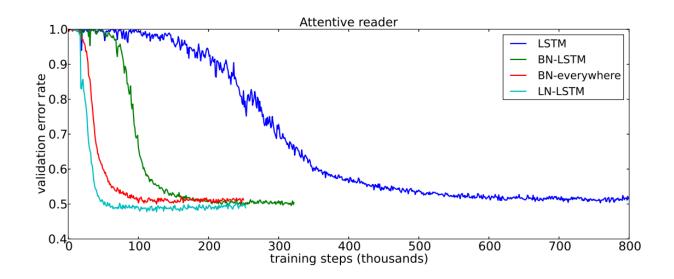
- The Problem: internal covariate shift occurs during training of a deep network when a small change in the low layers amplifies into a large change in the high layers
- One Solution: Layer
 normalization normalizes
 each layer and learns
 elementwise gain/bias
- Such normalization allows for higher learning rates (for faster convergence) without issues of diverging gradients

Given input $\mathbf{a} \in \mathbb{R}^K$, LayerNorm computes output $\mathbf{b} \in \mathbb{R}^K$:

$$\mathbf{b} = \boldsymbol{\gamma} \odot \frac{\mathbf{a} - \mu}{\sigma} \oplus \boldsymbol{\beta}$$

where we have mean $\mu = \frac{1}{K} \sum_{k=1}^{K} a_k$, standard deviation $\sigma = \sqrt{\frac{1}{K} \sum_{k=1}^{K} (a_k - \mu)^2}$, and parameters $\gamma \in \mathbb{R}^K$, $\beta \in \mathbb{R}^K$.

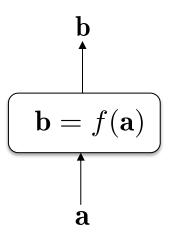
 \odot and \oplus denote elementwise multiplication and addition.

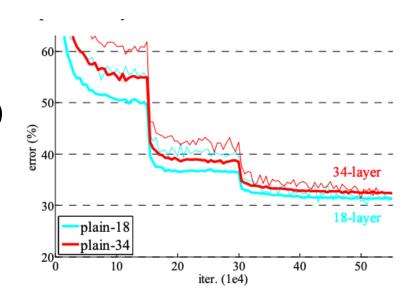


Residual Connections

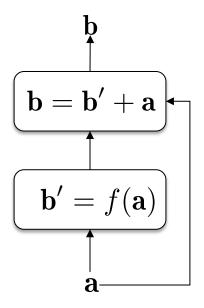
- The Problem: as network depth grows very large, a performance degradation occurs that is not explained by overfitting (i.e. train / test error both worsen)
- One Solution: Residual connections pass a copy of the input alongside another function so that information can flow more directly
- These residual connections allow for effective training of very deep networks that perform better than their shallower (though still deep) counterparts

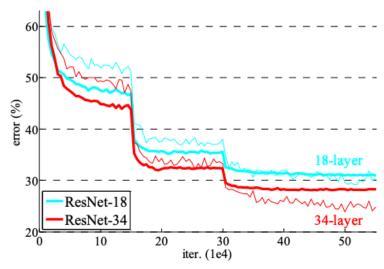
Plain Connection





Residual Connection

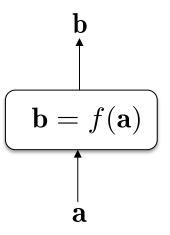




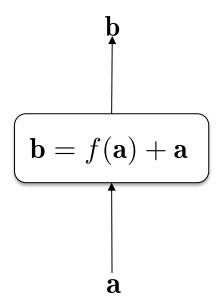
Residual Connections

- The Problem: as network depth grows very large, a performance degradation occurs that is not explained by overfitting (i.e. train / test error both worsen)
- One Solution: Residual connections pass a copy of the input alongside another function so that information can flow more directly
- These residual connections allow for **effective training of very deep networks** that perform better than their shallower (though still deep) counterparts

Plain Connection



Residual Connection



Why are residual connections helpful?

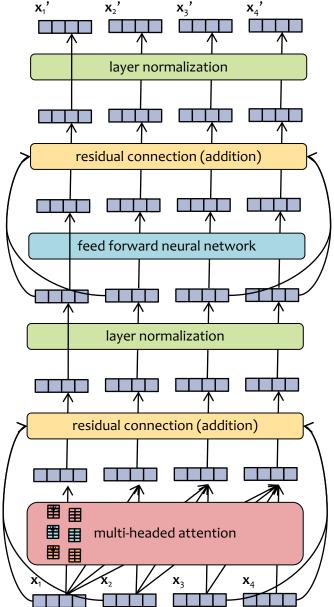
Instead of f(a) having to learn a full transformation of a, f(a) only needs to learn an additive modification of a (i.e. the residual).

Post-LN Version:

This is the version of the Transformer Layer that was introduced in the original paper in 2017.

The LayerNorm modules occur at the end of each set of 3 layers.

Transformer Layer



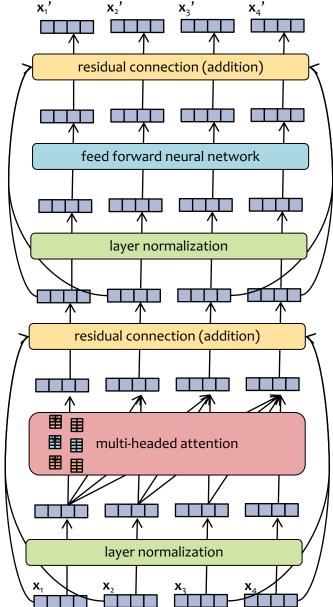
Each layer of a Transformer LM consists of several **sublayers**:

- 1. attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

Pre-LN Version:

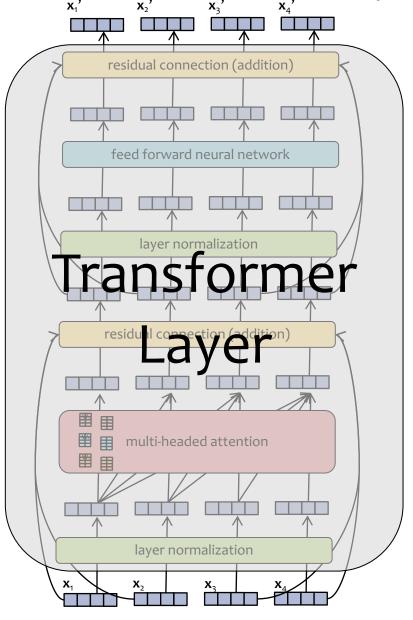
However, subsequent work found that reordering such that the LayerNorm's came at the beginning of each set of 3 layers, the multi-headed attention and feedforward NN layers tend to be better behaved (i.e. tricks like warm-up are less important).

Transformer Layer



- . attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

Transformer Layer



- 1. attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

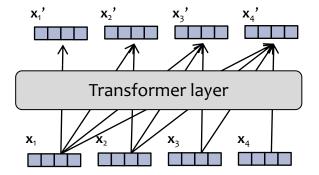
Transformer Layer



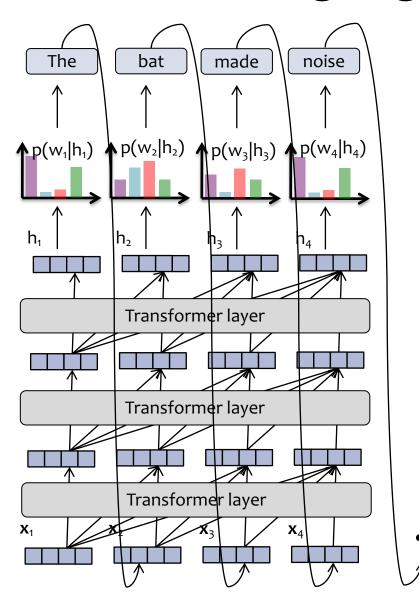
- . attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

Transformer Layer

- 1. attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections



Transformer Language Model



Each layer of a Transformer LM consists of several **sublayers**:

- 1. attention
- feed-forward neural network
- 3. layer normalization
- 4. residual connections

Each hidden vector looks back at the hidden vectors of the current and previous timesteps in the previous layer.

The language model part is just like an RNN-LM.

In-Class Poll



Suppose we have the following input embeddings and attention weights:

•
$$x_1 = [1,0,0,0] a_{4,1} = 0.1$$

•
$$x_2 = [0,1,0,0] a_{4,2} = 0.2$$

•
$$x_3 = [0,0,2,0] a_{4,3} = 0.6$$

•
$$x_4 = [0,0,0,1] a_{4,4} = 0.1$$

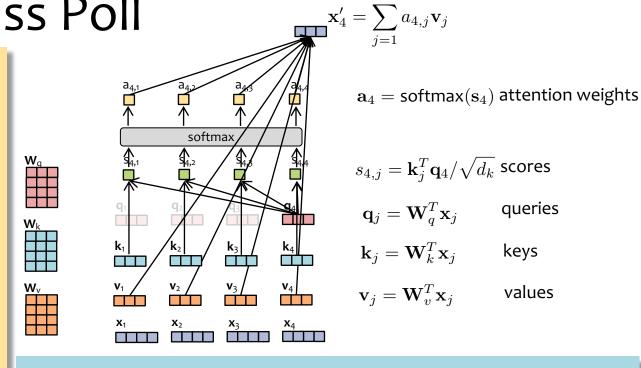
And $W_v = I$. Then we can compute x_4 .

Now suppose we swap the embeddings x_2 and x_3 such that

•
$$X_2 = [0,0,2,0]$$

•
$$X_3 = [0,1,0,0]$$

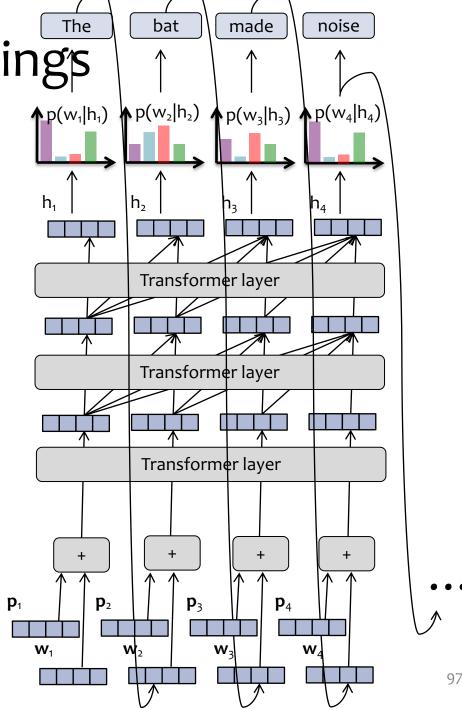
What is the new value of x_4 ?



Answer:

Position Embeddings

- The Problem: Because attention is position invariant, we need a way to learn about positions
- The Solution: Use (or learn) a collection of position specific embeddings: p_t represents what it means to be in position t. And add this to the word embedding w_t.
 - The **key idea** is that every word that appears in position t uses the same position embedding **p**_t
- There are a number of varieties of position embeddings:
 - Some are fixed (based on sine and cosine), whereas others are learned (like word embeddings)
 - Some are absolute (as described above) but we can also use relative position embeddings (i.e. relative to the position of the query vector)

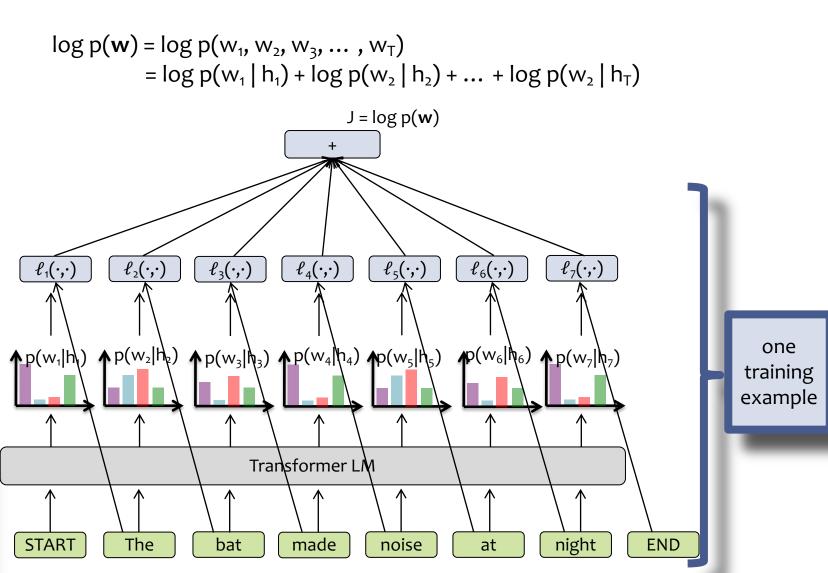


LEARNING A TRANSFORMER LM

Learning a Transformer LM

- Each training example is a sequence (e.g. sentence), so we have training data D = {w⁽¹⁾, w⁽²⁾,...,w^(N)}
- The objective function for a Deep LM (e.g. RNN-LM or Tranformer-LM) is typically the loglikelihood of the training examples: $J(\mathbf{\theta}) = \Sigma_i \log p_{\mathbf{\theta}}(\mathbf{w}^{(i)})$
- We train by mini-batch SGD (or your favorite flavor of mini-batch SGD)

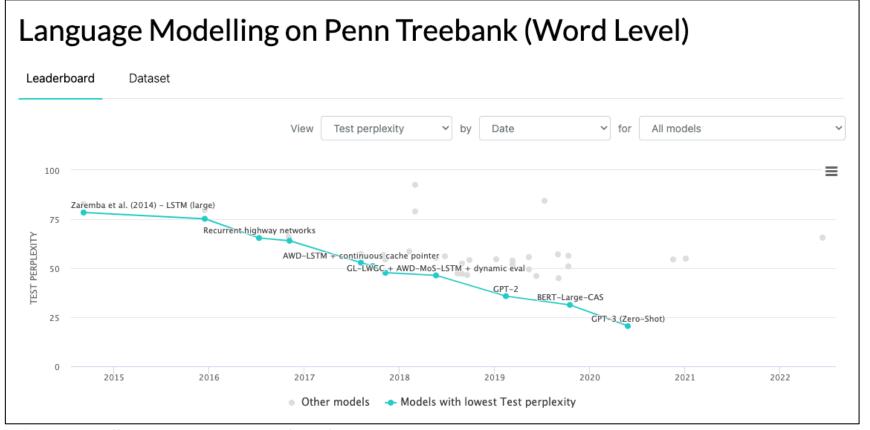
Training a Transformer-LM is the same, except we swap in a different deep language model.



Language Modeling

An aside:

- State-of-the-art language models currently tend to rely on transformer networks (e.g. GPT-3)
- RNN-LMs comprised most of the early neural LMs that led to current SOTA architectures



GPT-3

- GPT stands for Generative Pre-trained Transformer
- GPT is just a Transformer LM, but with a huge number of parameters

Model	# layers	dimension of states	dimension of inner states	# attention heads	# params
GPT (2018)	12	768	3072	12	117M
GPT-2 (2019)	48	1600			1542M
GPT-3 (2020)	96	12288	4*12288	96	175000M

Why does efficiency matter?

Case Study: GPT-3

- # of training tokens = 500 billion
- # of parameters = 175 billion
- # of cycles = 50
 petaflop/s-days
 (each of which
 are 8.64e+19
 flops)

Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens		
Common Crawl (filtered)	410 billion	60%	0.44		
WebText2	19 billion	22%	2.9		
Books1	12 billion	8%	1.9		
Books2	55 billion	8%	0.43		
Wikipedia	3 billion	3%	3.4		

Table 2.2: Datasets used to train GPT-3. "Weight in training mix" refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

Model Name	$n_{ m params}$	n_{layers}	d_{model}	$n_{ m heads}$	$d_{ m head}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

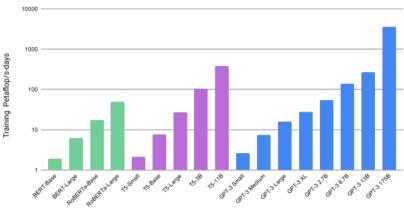


Figure 2.2: Total compute used during training. Based on the analysis in Scaling Laws For Neural Language Models [KMH+20] we train much larger models on many fewer tokens than is typical. As a consequence, although GPT-3 3B is almost 10x larger than RoBERTa-Large (355M params), both models took roughly 50 petaflop/s-days of compute during pre-training. Methodology for these calculations can be found in Appendix D.

Recap

Deep Learning

- AutoDiff
 - is a tool for computing gradients of a differentiable function, b = f(a)
 - the key building block is a module with a forward() and backward()
 - sometimes define f as code in forward() by chaining existing modules together
- Computation Graphs
 - are another way to define f (more conducive to slides)
 - we are considering various (deep) computation graphs: (1) CNN (2) RNN (3) RNN-LM
 (4) Transformer-LM
- Learning a Deep Network
 - deep networks (e.g. CNN/RNN) are trained by optimizing an objective function with SGD
 - compute gradients with AutoDiff

Language Modeling

- key idea: condition on previous words to sample the next word
- to define the probability of the next word...
 - ... n-gram LM uses collection of massive 50ksided dice
 - ... RNN-LM or Transformer-LM use a neural network
- Learning an LM
 - n-gram LMs are easy to learn: just count cooccurrences!
 - a RNN-LM / Transformer-LM is trained just like other deep neural networks