# LLMs

## Intro + Language modeling + Pretraining

Tatsunori Hashimoto

# Roadmap

1. LLMs today

2. Language modeling

3. Neural architectures – RNN/Transformer

4. Pretraining (part 1)

# Language models in the spotlight

Language models growing in relevance and capabilities

# Plotting

**Prompt:** Can you generate a pyplot for the following data: $x = [1, 3, 5, 6, 8], y1 = [2, 3, 5, 18, 1], y2 = [3, 5, 6, 8, 1], y3 = [5, 1, 2, 3, 4], y4 = [9, 7, 2, 3, 1]$. I want $y1, y2$ to be in the same plot, but $y3$ is in another plot next to that plot, $y4$ is in below. I want the legend of $y1$ to be "bob", $y2$ to be "alice", $y3$ to be "bilbo", $y4$ to be "allie". I want the $x$-axis to be labeled with "time" and $y$ axis to be labeled with "money". I want to add a 10 %-40% random error bar to each curve, through all times (including non-integers). I want smoothed curves to show the plot, and smoothed error bar. Do not use linear interpolation, use smooth interpolation! I want to also add some small zig-zag to the smoothed curve to make it look more real. I want to put a baseline as the mean of each line. I want to put a pie chart below indicating the fr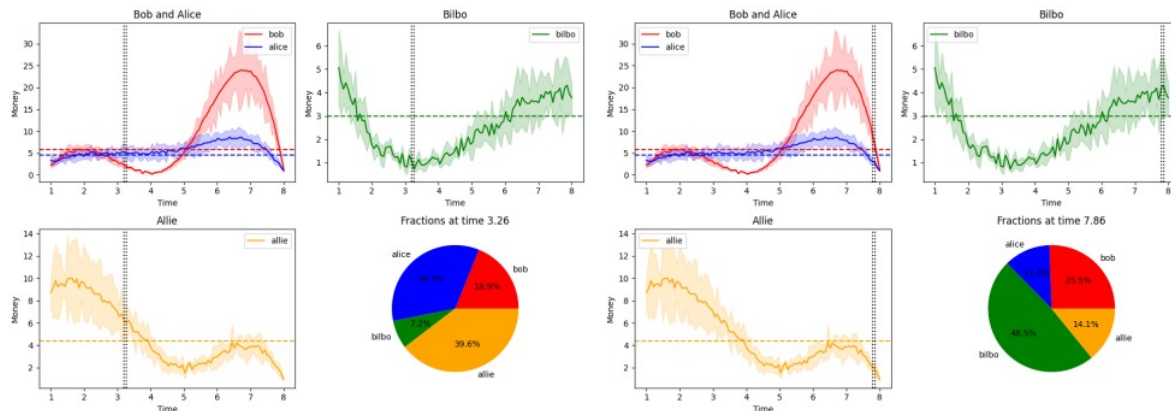action of the four people in each time step. I also want an animation to show how the fractions are changing in the pie chart in **continuous time**. Interpolate the missing fractions! I also want vertical line animation in other three plots to match the pie chart. I want the fanciest plot. Please add as many fancy things as possible.

Here is a possible code to generate a pyplot for the following data, with some comments to explain the steps:

**GPT-4:** [We give two snapshots from the animation resulting from the produced code]



Text example from [Bubeck et al 2023]

# Math

**Prompt:** Suppose $g(x) = f^{-1}(x), g(0) = 5, g(4) = 7, g(3) = 2, g(7) = 9, g(9) = 6$ what is $f(f(f(6)))$?

**GPT-4:** There are different ways to approach this problem, but one possible method is to use the fact that $g(x) = f^{-1}(x)$ implies that $f(g(x)) = x$ for any $x$ in the domain of $g$. That is, applying $f$ to both sides of the equation $g(x) = f^{-1}(x)$ gives $f(g(x)) = f(f^{-1}(x)) = x$. This means that we can undo the effect of $f$ by applying $g$, and vice versa. Using this idea, we can try to find $f(6)$ by looking for the value of $x$ such that $g(x) = 6$. From the given information, we see that $g(9) = 6$, so $f(6) = f(g(9)) = 9$. Similarly, we can find $f(9)$ by looking for the value of $x$ such that $g(x) = 9$. We see that $g(7) = 9$, so $f(9) = f(g(7)) = 7$. Finally, we can find $f(7)$ by looking for the value of $x$ such that $g(x) = 7$. We see that $g(4) = 7$, so $f(7) = f(g(4)) = 4$. Therefore, $f(f(f(6))) = f(f(9)) = f(7) = 4$.

Text example from [Bubeck et al 2023]

# Underlying this: rapid progress on benchmarks

It's an exciting time: extremely rapid, sudden improvements in NLP



Kiela et al

# Major driver: foundation models and internet data

What is the big change these last 10 years? What enabled recent advances in GenAI?

**'Foundation models'** that can effectively leverage vast amounts of internet data



**Intuition**: learn 'task-agnostic' structures (syntax of language, common sense) from freely available internet data.

[Bommasani et al 2021]

# How do LLMs work?

**Step 1 - Pretraining**: learn to autocomplete text on the internet



Maybe we don't just want to mimic users on the internet..

# How do LLMs work?

**Step 2 – Post-training:** explicitly reinforce desired behaviors identified by annotators



[Ouyang 2020]

# Language Modeling

**Language Modeling** is the task of predicting what word comes next

*the students opened their _____*

book
s
laptops
exams
mind

More formally: given a sequence of words $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(t)}$ compute the probability distribution of the next word $\boldsymbol{x}^{(t+1)}$ :

$$P(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})$$

where $\boldsymbol{x}^{(t+1)}$ can be any word in the vocabular $V = \{\boldsymbol{w}_1, ..., \boldsymbol{w}_{|V|}\}$

A system that does this is called a **Language Model**

[Slide from CS224n]

# Language Modeling

You can also think of a Language Model as a system that
assigns a probability to a piece of text

For example, if we have some text $x^{(1)}, \ldots, x^{(T)}$ , then the probability of this text (according to the Language Model) is:

$$P(\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)}|\ \boldsymbol{x}^{(1)}) \times \cdots \times P(\boldsymbol{x}^{(T)}|\ \boldsymbol{x}^{(T-1)}, \ldots, \boldsymbol{x}^{(1)})$$

$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)}|\ \boldsymbol{x}^{(t-1)}, \ldots, \boldsymbol{x}^{(1)})$$

This is what our LM provides

[Slide from CS224n]

# n-gram Language Models

*the students opened their  _____*

**Question**: How to learn a Language Model?

**Answer** (pre- Deep Learning): learn an *n*-gram Language Model!

**Definition:** An *n*-gram is a chunk of *n* consecutive words.
- unigrams: "the", "students", "opened", "their"
- bigrams: "the students", "students opened", "opened their"
- trigrams: "the students opened", "students opened their"
- four-grams: "the students opened their"

**Idea:** Collect statistics about how frequent different n-grams are and use these to predict next word.

[Slide from CS224n]

# n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

~~as the proctor started the clock, the~~ *students opened their* _____

discard

condition on this

$$P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w})}{\text{count}(\text{students opened their})}$$

For example, suppose that in the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

[Slide from CS224n]

# Generating text with a n-gram Language Model

You can also use a Language Model to generate text

*today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

…but **incoherent.** We need to consider more than three words at a time if we want to model language well.

But increasing *n* worsens sparsity problem, and increases model size…

[Slide from CS224n]

# How to build a *neural* language model?

Recall the Language Modeling task:

- Input: sequence of words $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(t)}$
- Output: prob. dist. of the next word $P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})$

How about a window-based neural model?

LOCATION

$\boldsymbol{U}$

$\boldsymbol{W}$

*museums*   *in*   *Paris*   *are*   *amazing*

[Slide from CS224n]

# A fixed-window neural Language Model

output distribution
$$\hat{\boldsymbol{y}} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden laver
$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word
$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot
$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

# A fixed-window neural Language Model

Approximately: Y. Bengio, et al. (2000/2003): A Neural Probabilistic Language Model

**Improvements** over *n*-gram LM:

- No sparsity problem
- Don't need to store all observed *n*-grams

Remaining **problems**:

- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in $W$. No symmetry in how the inputs are processed.

We need a neural architecture that can process *any length input*



books

laptops

a                    zoo

$U$

$W$

the        student       opened       their
$\boldsymbol{x}^{(1)}$      $\boldsymbol{x}^{(2)}$      $\boldsymbol{x}^{(3)}$      $\boldsymbol{x}^{(4)}$

[Slide from CS224n]

# A Simple RNN Language Model

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

output distribution

$$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

word

$$e^{(t)} = Ex^{(t)};$$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

*books*   *laptops*

a    zo o

$U$

$h^{(0)}$  $h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

$W_h$  $W_h$  $W_h$  $W_h$

$W_e$  $W_e$  $W_e$  $W_e$

$e^{(1)}$  $e^{(2)}$  $e^{(3)}$  $e^{(4)}$

$E$  $E$  $E$  $E$

*the*   *student*   *opened*   *their*
$x^{(1)}$  $x^{(2)}$  $x^{(3)}$  $x^{(4)}$

**Note**: *this input sequence could be much longer now!*

18

[Slide from CS224n]

# RNN Language Models

RNN **Advantages**:
- Can process any length input
- Computation for step $t$ can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

RNN **Disadvantages**:
- Recurrent computation is slow
- In practice, difficult to access information from many steps

More on these later

$$\hat{\boldsymbol{y}}^{(4)} = P(\boldsymbol{x}^{(5)}|\text{the students opened their})$$

books

laptops

a                    zo
                       o

$\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$   $\boldsymbol{h}^{(1)}$   $\boldsymbol{h}^{(2)}$   $\boldsymbol{h}^{(3)}$   $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$

$\boldsymbol{W}_e$   $\boldsymbol{W}_e$   $\boldsymbol{W}_e$   $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$   $\boldsymbol{e}^{(2)}$   $\boldsymbol{e}^{(3)}$   $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$

*the*        *student*     *opened*      *their*
$\boldsymbol{x}^{(1)}$    $\boldsymbol{x}^{(2)}$    $\boldsymbol{x}^{(3)}$    $\boldsymbol{x}^{(4)}$

[Slide from CS224n]

# Training an RNN Language Model

Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$

Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ for *every step t.*

- i.e., predict probability dist of *every word*, given words so far

Loss function on step *t* is cross-entropy between predicted probability distributic $\hat{y}^{(t)}$ , and the true next wo $y^{(t)}$ (one-hot f $x^{(t+1)}$ ):

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

[Slide from CS224n]

# Evaluating Language Models

The standard evaluation metric for Language Models is perplexity.

$$\text{perplexity} = \prod_{t=1}^{T} \left( \frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})} \right)^{1/T}$$

Normalized by number of words

Inverse probability of corpus, according to Language Model

This is equal to the exponential of the cross-entropy loss     : $J(\theta)$

$$= \prod_{t=1}^{T} \left( \frac{1}{\hat{\boldsymbol{y}}^{(t)}_{\boldsymbol{x}_{t+1}}} \right)^{1/T} = \exp\left( \frac{1}{T} \sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}^{(t)}_{\boldsymbol{x}_{t+1}} \right) = \exp(J(\theta))$$

**Lower** perplexity is better!

[Slide from CS224n]

# RNNs greatly improved perplexity over what came before

*n*-gram model

Increasingly complex RNNs

| Model | Perplexity |
|---|---|
| Interpolated Kneser-Ney 5-gram (Chelba et al., 2013) | 67.6 |
| RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013) | 51.3 |
| RNN-2048 + BlackOut sampling (Ji et al., 2015) | 68.3 |
| Sparse Non-negative Matrix factorization (Shazeer et al., 2015) | 52.9 |
| LSTM-2048 (Jozefowicz et al., 2016) | 43.7 |
| 2-layer LSTM-8192 (Jozefowicz et al., 2016) | 30 |
| Ours small (LSTM-2048) | 43.9 |
| Ours large (2-layer LSTM-2048) | 39.8 |

Perplexity improves
(lower is better)

[Slide from CS224n]

# Problems with RNNs: Vanishing and Exploding Gradients

[Slide from CS224n]

# Vanishing gradient intuition

$J^{(4)}(\theta)$

$\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}$  $\boldsymbol{W}$  $\boldsymbol{W}$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$

What happens if these are small?

**Vanishing gradient problem:**
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

[Slide from CS224n]

# Vanishing gradient proof sketch (linear case)

Recall:
$$\boldsymbol{h}^{(t)} = \sigma\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)$$

What if $\sigma$ were the identity function, $\sigma(x) = x$ ?

$$\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} = \text{diag}\left(\sigma'\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)\right) \boldsymbol{W}_h \qquad \text{(chain rule)}$$

$$= \boldsymbol{I} \; \boldsymbol{W}_h = \boldsymbol{W}_h$$

Consider the gradient of the loss $J^{(i)}(\theta)$ on step $i$, with respect to the hidden state $\boldsymbol{h}^{(j)}$ on some previous step $j$. Let $\ell = i - j$

$$\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \le i} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \qquad \text{(chain rule)}$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \le i} \boldsymbol{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \boxed{\boldsymbol{W}_h^\ell} \qquad \text{(value of } \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \text{)}$$

If $W_h$ is "small", then this term gets exponentially problematic as $\ell$ becomes large

[Slide from CS224n]

**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf
(and supplemental materials), at http://proceedings.mlr.press/v28/pascanu13-supp.pdf

# Why is vanishing gradient a problem?



Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are updated only with respect to near effects, not long-term effects.

[Slide from CS224n]

# Issues with recurrent models: Linear interaction distance

**O(sequence length)** steps for distant word pairs to interact means:

- Hard to learn long-distance dependencies (because gradient problems!)
- Linear order of words is "baked in"; we already know linear order isn't the right way to think about sentences…



*The* **chef** *who* …          ***was***

Info of **chef** has gone through O(sequence length) many layers!

# Issues with recurrent models: Lack of parallelizability

Forward and backward passes have **O(sequence length)** unparallelizable operations

- GPUs can perform a bunch of independent computations at once!
- But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
- Inhibits training on very large datasets!



$h_1$     $h_2$                                                    $h_T$

[Slide from CS224n]

Numbers indicate min # of steps before a state can be computed

# Attention

**Attention** provides a solution to the bottleneck problem.

**Core idea**: on each step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence

First, we will show via diagram (no equations), then we will show with equations

[Slide from CS224n]

# The starting point: mean-pooling for RNNs

**positive**

Sentence encoding

How to compute sentence encoding?

**Usually better**: Take element-wise max or mean of all hidden states

*overall*    *I*    *enjoyed*    *the*    *movie*    *a*    *lot*

Starting point: a *very* basic way of 'passing information from the encoder' is to *average*

[Slide from CS224n]

# Attention is *weighted* averaging, which lets you do lookups!

Attention is just a **weighted** average – this is very powerful if the weights are learned!

In **attention**, the **query** matches all **keys** *softly*, to a weight between 0 and 1. The keys' **values** are multiplied by the weights and summed.

In a **lookup table**, we have a table of **keys** that map to **values**. The **query** matches one of the keys, returning its value.

[Slide from CS224n]

# Sequence-to-sequence with attention

**Core idea**: on each step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence

dot product

Attention scores {

Encoder RNN {

Decoder RNN

Source sentence (input)

il    a    m'    entarté        <START>

[Slide from CS224n]

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

*il*   *a*   *m'*   *entarté*      *<START>*

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

*il*   *a*   *m'*   *entarté*   *<START>*

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention

dot product

Attention scores

Encoder RNN

Decoder RNN

*il*   *a*   *m'*   *entarté*   *&lt;START&gt;*

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



On this decoder timestep, we're mostly focusing on the first encoder hidden state (*"he"*)

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il     a     m'    entarté          <START>

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



Source sentence (input)

# Sequence-to-sequence with attention



Sometimes we take the attention output from the previous step, and also feed it into the decoder (along with the usual decoder input). We do this in Assignment 4.

[Slide from CS224n]

# Sequence-to-sequence with attention

Source sentence (input)

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*with*

$\hat{y}_4$

*il*  *a*  *m'*  *entarté*  *<START >*  *he*  *hit*  *me*

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

$a$

$\hat{y}_5$

*il    a    m'    entarté*    *<START>    he    hit    me    with*

Source sentence (input)

[Slide from CS224n]

# Sequence-to-sequence with attention



Source sentence (input)

# Attention: in equations

We have encoder hidden states $h_1, \ldots, h_N \in \mathbb{R}^h$

On timestep $t$, we have decoder hidden state $s_t \in \mathbb{R}^h$

We get the attention scores $e^t$ for this step:

$$e^t = [s_t^T h_1, \ldots, s_t^T h_N] \in \mathbb{R}^N$$

We take softmax to get the attention distribution $\alpha^t$ or this step (this is a probability distribution and sums to 1)

$$\alpha^t = \mathrm{softmax}(e^t) \in \mathbb{R}^N$$

We use $\alpha^t$ to take a weighted sum of the encoder hidden states to get the attention out $a_t$

$$a_t = \sum_{i=1}^{N} \alpha_i^t h_i \in \mathbb{R}^h$$

Finally we concatenate the attention output $a_t$ with the decoder hidden state $s_t$ and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

# Attention is great!

Attention solves the bottleneck problem

- Attention allows decoder to look directly at source; bypass bottleneck

Attention helps with the vanishing gradient problem

- Provides shortcut to faraway states

Attention provides some interpretability

- By inspecting attention distribution, we see what the decoder was focusing on
- We get (soft) alignment for free!
- This is cool because we never explicitly trained an alignment system
- The network just learned alignment by itself



[Slide from CS224n]

# Do we even need recurrence at all?

Abstractly: Attention is a way to pass information from a sequence ($x$) to a neural network input. ($h_t$)

- This is also *exactly* what RNNs are used for – to pass information!
- **Can we just get rid of the RNN entirely?** Maybe attention is just a better way to pass information!



**2014-2017ish Recurrence**

Lots of trial and error →

**2021 ??????**

[Slide from CS224n]

# The building block we need: *self* attention

What we talked about – **Cross** attention: paying attention to the input x to generate $y_t$



- What we need – **Self** attention: to generate $y_t$, we need to pay attention to $y_{<t}$

[Slide from CS224n]

# Self-Attention Hypothetical Example

attention
weights
for
**"learned"**

q

| v | v | v | v | v | v | v | v |
|---|---|---|---|---|---|---|---|
| k | k | k | k | k | k | k | k |

I    went    to    Stanford    CS    224n    and    learned

[Slide from CS224n]

## Self-Attention: keys, queries, values from the same sequence

Let $\boldsymbol{w}_{1:n}$ be a sequence of words in vocabulary $V$, like *Zuko made his uncle tea*.

For each $\boldsymbol{w}_i$, let $\boldsymbol{x}_i = E\boldsymbol{w}_i$, where $E \in \mathbb{R}^{d \times |V|}$ is an embedding matrix.

1. Transform each word embedding with weight matrices Q, K, V, each in $\mathbb{R}^{d \times d}$

$$\boldsymbol{q}_i = Q\boldsymbol{x}_i \text{ (queries)} \qquad \boldsymbol{k}_i = K\boldsymbol{x}_i \text{ (keys)} \qquad \boldsymbol{v}_i = V\boldsymbol{x}_i \text{ (values)}$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

$$\boldsymbol{e}_{ij} = \boldsymbol{q}_i^\top \boldsymbol{k}_j \qquad \boldsymbol{\alpha}_{ij} = \frac{\exp(\boldsymbol{e}_{ij})}{\sum_{j'} \exp(\boldsymbol{e}_{ij'})}$$

3. Compute output for each word as weighted sum of values

$$\boldsymbol{o}_i = \sum_j \boldsymbol{\alpha}_{ij} \, v_i$$

[Slide from CS224n]

# Barriers and solutions for Self-Attention as a building block

**Barriers**

- Doesn't have an inherent notion of order!

**Solutions**

# Fixing the first self-attention problem: sequence order

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values.

- Consider representing each **sequence index** as a **vector**

$$\boldsymbol{p}_i \in \mathbb{R}^d, \text{ for } i \in \{1, 2, \dots, n\} \text{ are position vectors}$$

- Don't worry about what the $p_i$ are made of yet!

- Easy to incorporate this info into our self-attention block: just add the $\boldsymbol{p}_i$ to our inputs!

- Recall that $\boldsymbol{x}_i$ is the embedding of the word at index $i$. The positioned embedding is:

$$\widetilde{\boldsymbol{x}}_i = \boldsymbol{x}_i + \boldsymbol{p}_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add…

[Slide from CS224n]

# Position representation vectors learned from scratch

- **Learned absolute position representations:** Let all $p_i$ be learnable parameters!
  Learn a matrix $\boldsymbol{p} \in \mathbb{R}^{d \times n}$, and let each $\boldsymbol{p}_i$ be a column of that matrix!

- Pros:
  - Flexibility: each position gets to be learned to fit the data
- Cons:
  - Definitely can't extrapolate to indices outside $1, \dots, n$.
- Many systems use this!

- Sometimes people try more flexible representations of position:
  - Relative linear position attention [Shaw et al., 2018]
  - Dependency syntax-based position [Wang et al., 2019]

[Slide from CS224n]

# Barriers and solutions for Self-Attention as a building block

**Barriers**

- Doesn't have an inherent notion of order!

- No nonlinearities for deep learning! It's all just weighted averages

**Solutions**

- Add position representations to the inputs

[Slide from CS224n]

# Adding nonlinearities in self-attention

- Note that there are no elementwise nonlinearities in self-attention; stacking more self-attention layers just re-averages **value** vectors (Why? Look at the notes!)

- Easy fix: add a **feed-forward network** to post-process each output vector.

$$m_i = MLP(\text{output}_i)$$
$$\quad = W_2 * \text{ReLU}(W_1\,\text{output}_i + b_1) + b_2$$



$w_1$     $w_2$     $w_3$        $w_n$

*The*     *chef*     *who*      *food*

Intuition: the FF network processes the result of attention

# Barriers and solutions for Self-Attention as a building block

| **Barriers** | **Solutions** |
|---|---|

- Doesn't have an inherent notion of order! ⟶ • Add position representations to the inputs

- No nonlinearities for deep learning magic! It's all just weighted averages ⟶ • Easy fix: apply the same feedforward network to each self-attention output.

- Need to ensure we don't "look at the future" when predicting a sequence ⟶
  - Like in machine translation
  - Or language modeling

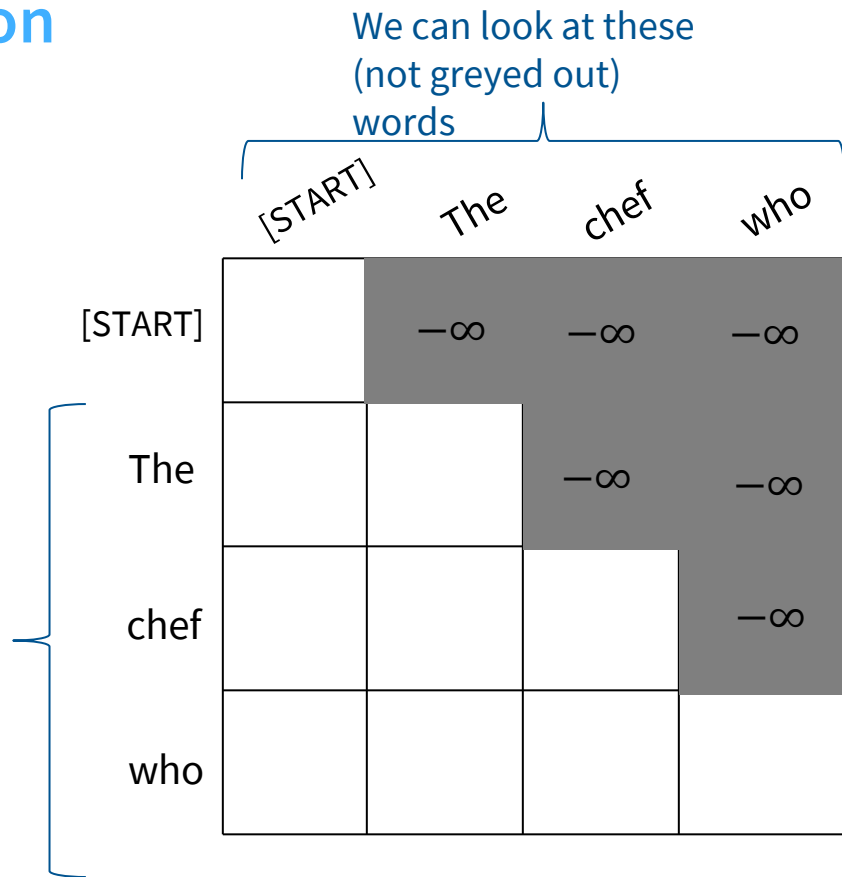# Masking the future in self-attention

- To use self-attention in **decoders**, we need to ensure we can't peek at the future.

- At every timestep, we could change the set of **keys and queries** to include only past words. (Inefficient!)

- To enable parallelization, we **mask out attention** to future words by setting attention scores to $-\infty$.

We can look at these (not greyed out) words

For encoding these words

$$e_{ij} = \begin{cases} q_i^\top k_j, & j \leq i \\ -\infty, & j > i \end{cases}$$

|  | [START] | The | chef | who |
|---|---|---|---|---|
| [START] |  | $-\infty$ | $-\infty$ | $-\infty$ |
| The |  |  | $-\infty$ | $-\infty$ |
| chef |  |  |  | $-\infty$ |
| who |  |  |  |  |

[Slide from CS224n]

# Barriers and solutions for Self-Attention as a building block

| **Barriers** | **Solutions** |
|---|---|

- Doesn't have an inherent notion of order!  ⟶  • Add position representations to the inputs

- No nonlinearities for deep learning magic! It's all just weighted averages  ⟶  • Easy fix: apply the same feedforward network to each self-attention output.

- Need to ensure we don't "look at the future" when predicting a sequence  ⟶  • Mask out the future by artificially setting attention weights to 0!
  - Like in machine translation
  - Or language modeling

# Necessities for a self-attention building block:

- **Self-attention**:
  - the basis of the method.
- **Position representations**:
  - Specify the sequence order, since self-attention is an unordered function of its inputs.
- **Nonlinearities**:
  - At the output of the self-attention block
  - Frequently implemented as a simple feed-forward network.
- **Masking**:
  - In order to parallelize operations while not looking at the future.
  - Keeps information about the future from "leaking" to the past.

Probabilities

Softmax

Linear

Repeat for number of encoder blocks

Feed-Forward

Masked Self-Attention

Block

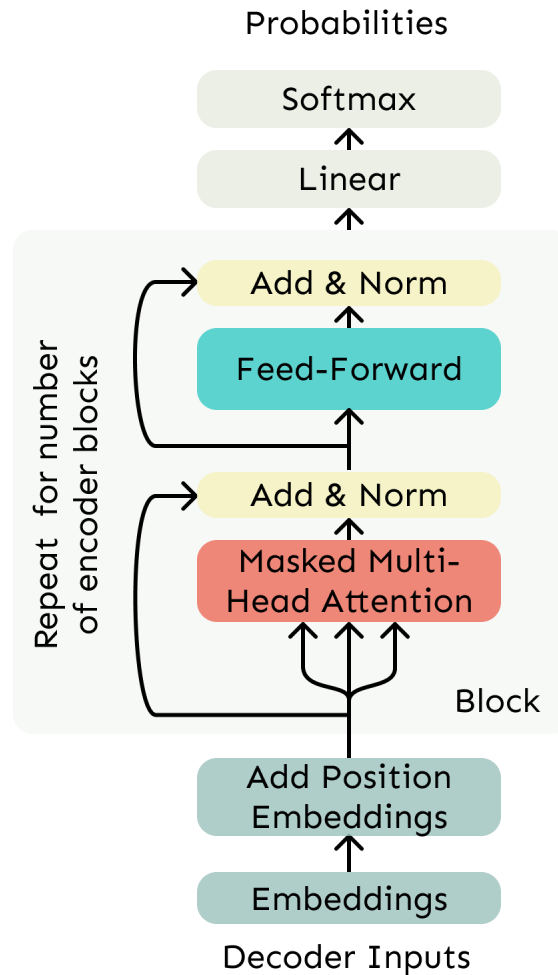Add Position Embeddings

Embeddings

Inputs

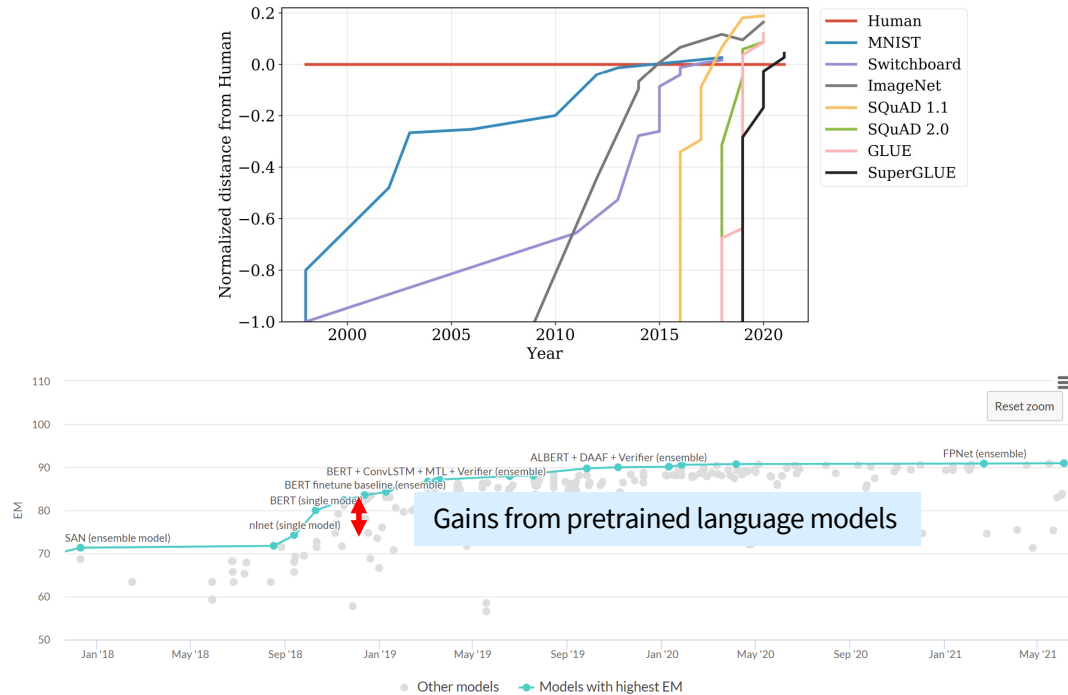[Slide from CS224n]

# Other components we wont cover

- **Multi-head attention**
  - Have multiple, but smaller attention heads and linearly mix the outputs

- **Residual connections**
  - Have residual connections around the attention and feed forward (addresses vanishing gradients)

- **Layer norm**
  - Standardize the activation across the hidden dimension coordinate.

# The Transformer Decoder

- The Transformer Decoder is a stack of Transformer Decoder **Blocks**.
- Each Block consists of:
  - Self-attention
  - Add & Norm
  - Feed-Forward
  - Add & Norm
- That's it! We've gone through the Transformer Decoder.



Probabilities

Softmax

Linear

Repeat for number of encoder blocks

Add & Norm

Feed-Forward

Add & Norm

Masked Multi-Head Attention

Block

Add Position Embeddings

Embeddings

Decoder Inputs

60

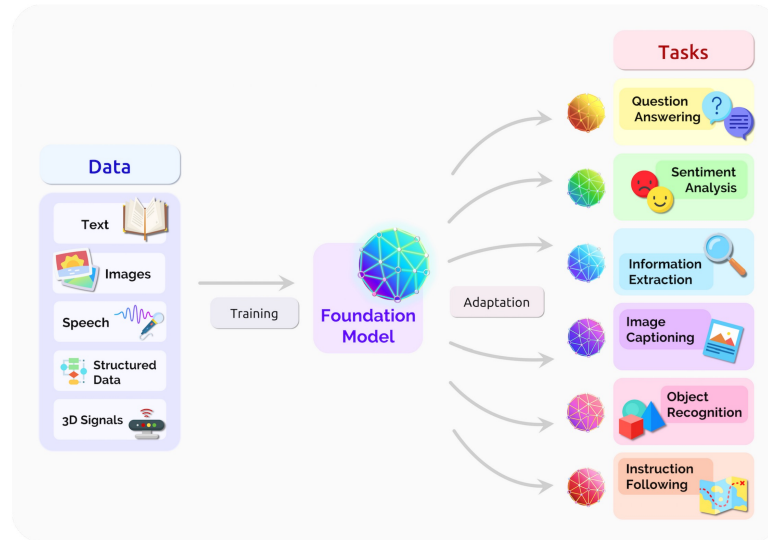[Slide from CS224n]

# The pretraining revolution



Gains from pretrained language models

Pretraining has had a major, tangible impact on how well NLP systems work

# Pretraining – scaling unsupervised learning on the internet



**Key ideas in pretraining**
- Make sure your model can process large-scale, diverse datasets
- Don't use labeled data (otherwise you can't scale!)
- Compute-aware scaling

# What kinds of things does pretraining teach?

There's increasing evidence that pretrained models learn a wide variety of things about the statistical properties of language.

*Stanford University is located in _____, California.* [Trivia]

*I put ___ fork down on the table.* [syntax]

*The woman walked across the street, checking for traffic over ___ shoulder.* [coreference]

*I went to the ocean to see the fish, turtles, seals, and _____.* [lexical semantics/topic]

*Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ___.* [sentiment]

Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the _____. [some reasoning – this is harder]

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, ____ [some basic arithmetic; they don't learn the Fibonnaci sequence]

Models also learn – and can exacerbate racism, sexism, all manner of bad biases.

[Slide from CS224n]

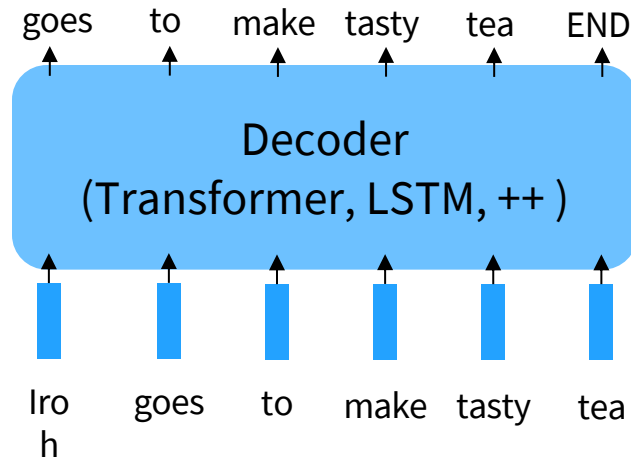# Pretraining through language modeling [Dai and Le, 2015]

Recall the **language modeling** task:

- Model $p_\theta(w_t|w_{1:t-1})$, the probability distribution over words given their past contexts.

- There's lots of data for this! (In English.)

**Pretraining through language modeling:**

Train a neural network to perform language modeling on a large amount of text.

Save the network parameters.

goes    to    make   tasty    tea    END

Decoder
(Transformer, LSTM, ++ )

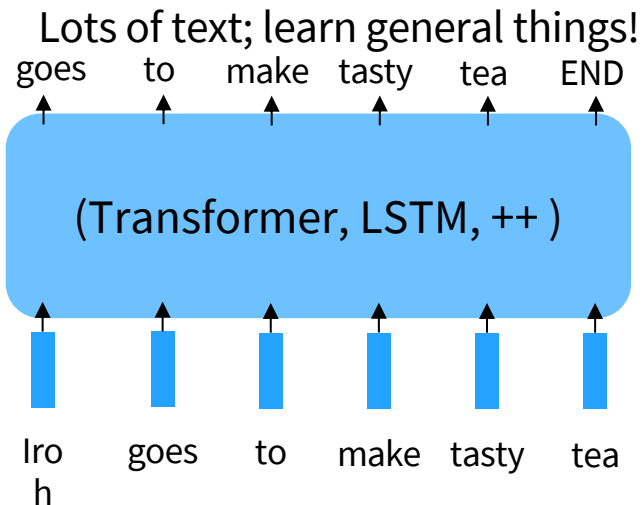Iroh   goes    to    make   tasty    tea
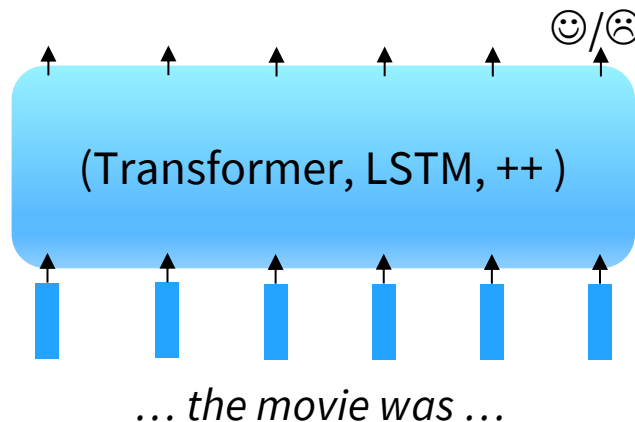
[Slide from CS224n]

# The Pretraining / Finetuning Paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!

goes    to    make    tasty    tea    END

(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on your task)**

Not many labels; adapt to the task!

☺/☹

(Transformer, LSTM, ++ )

*… the movie was …*

[Slide from CS224n]

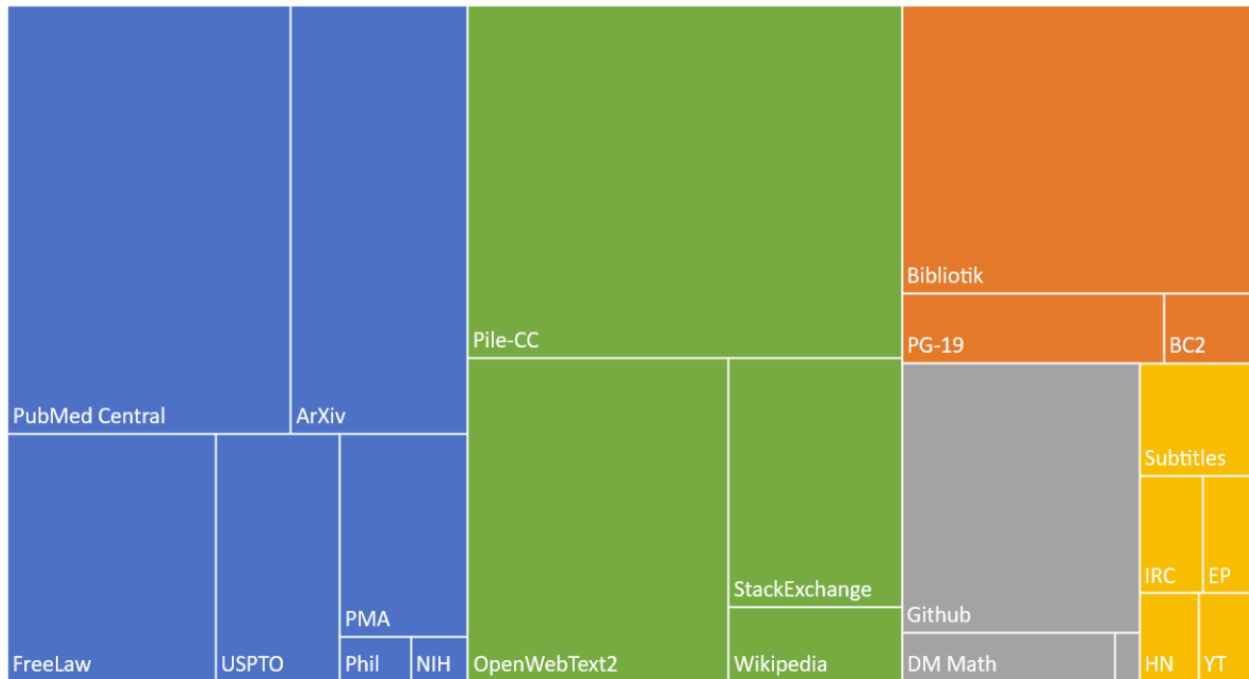# Stochastic gradient descent and pretrain/finetune

Why should pretraining and finetuning help, from a "training neural nets" perspective?

- Consider, provides parameters $\hat{\theta}$ by approximating $\min_{\theta} \mathcal{L}_{\mathrm{pretrain}}(\theta)$.

     (The pretraining loss.)

- Then, finetuning approximates $\min_{\theta} \mathcal{L}_{\mathrm{finetune}}(\theta)$, starting at $\hat{\theta}$.

     (The finetuning loss)

- The pretraining may matter because stochastic gradient descent sticks (relatively) close to $\hat{\theta}$ during finetuning.

     So, maybe the finetuning local minima near $\hat{\theta}$ tend to generalize well!

     And/or, maybe the gradients of finetuning loss near $\hat{\theta}$ propagate nicely!

[Slide from CS224n]

# Where does this data come from?


Composition of the Pile by Category

| Model | Training Data |
|-------|---------------|
| BERT | BookCorpus, English Wikipedia |
| GPT-1 | BookCorpus |
| GPT-3 | CommonCrawl, WebText, English Wikipedia, and 2 book databases ("Books 1" and "Books 2") |
| GPT-3.5+ | Undisclosed |

# Recap

1.  **Language modeling:** generative learning of natural language (often autoregressive)

2.  **Transformers:** scalable, easy to to train, parallelizable architecture for sequence modeling

3.  **Pretraining:** language modeling over the internet in a task-agnostic manner.