



# Reviewing the Foundations

CSCI 601-771 (NLP: Advances in Self-Supervised Models)

# Language Modeling: Definitions and History

The

# The cat

The cat sat

The cat sat on

The cat sat on \_\_?\_\_

The cat sat on the mat.

$P(\text{mat} \mid \text{The cat sat on the})$ 

next word

context or prefix

# Probability of Upcoming Word

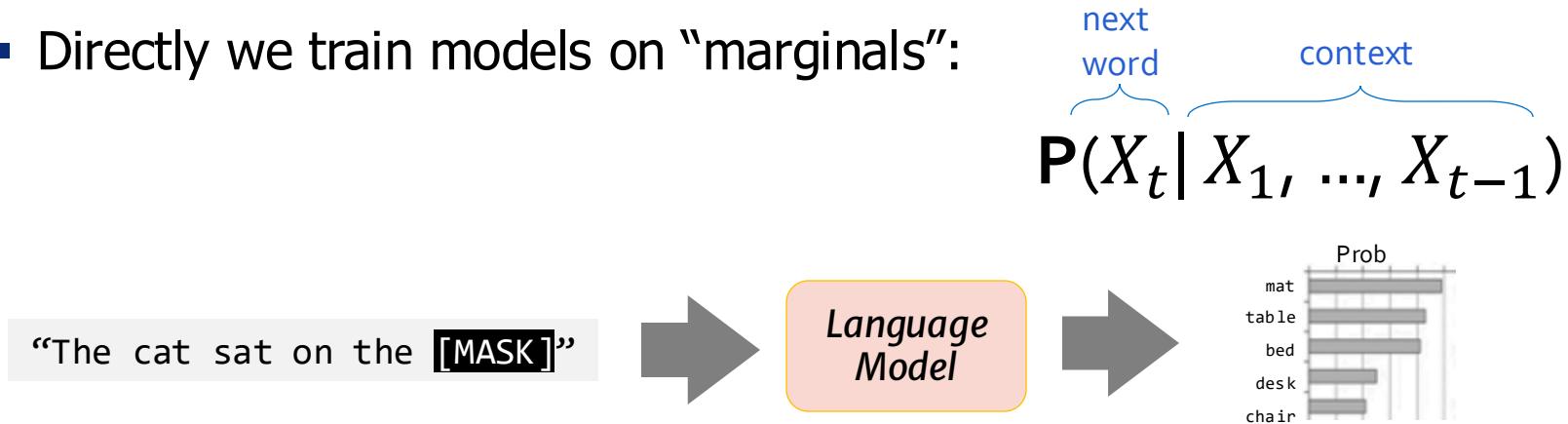
---

$$P(X_t | X_1, \dots, X_{t-1})$$



# LMs as a Marginal Distribution

- Directly we train models on “marginals”:



# LMs as Implicit Joint Distribution over Language

---

- While language modeling involves learning the marginals, we are implicitly learning the full/joint distribution of language.
  - Remember the chain rule:

$$P(X_1, \dots, X_t) = P(X_1) \prod_{i=1}^t P(X_i | X_1, X_2, \dots, X_{i-1})$$

- **Language Modeling**  $\triangleq$  learning prob distribution over language sequence.

# Doing Things with Language Model

---

- What is the probability of ....

“I like Johns Hopkins University”

“like Hopkins I University Johns”

# Doing Things with Language Model

---

- What is the probability of ....

“I like Johns Hopkins University”

“like Hopkins I University Johns”

- LMs assign a probability to every sentence (or any string of words).

$P(\text{"I like Johns Hopkins University EOS"}) = 10^{-5}$

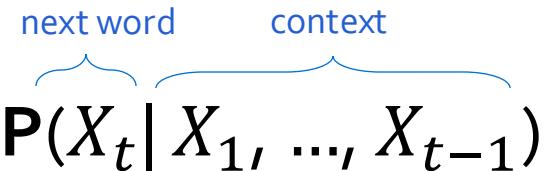
$P(\text{"like Hopkins I University Johns EOS"}) = 10^{-15}$

# Doing Things with Language Model (2)

- We can rank sentences.
- While LMs show “typicality”, this may be a proxy indicator to other properties:
  - Grammaticality, fluency, factuality, etc.

$$P(X_t | X_1, \dots, X_{t-1})$$

next word      context



$P("I \ like \ Johns \ Hopkins \ University. \ EOS") > P("I \ like \ John \ Hopkins \ University \ EOS")$

$P("I \ like \ Johns \ Hopkins \ University. \ EOS") > P("University. \ I \ Johns \ EOS \ Hopkins \ like")$

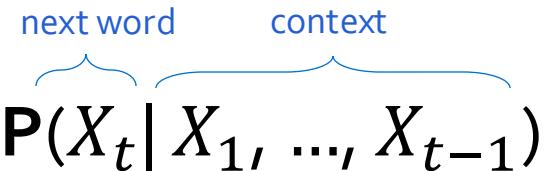
$P("JHU \ is \ located \ in \ Baltimore. \ EOS") > P("JHU \ is \ located \ in \ Virginia. \ EOS")$

# Doing Things with Language Model (3)

- Can also generate strings!

$$P(X_t | X_1, \dots, X_{t-1})$$

next word      context



- Let's say we start "*Johns Hopkins is*"
- Using this prompt as an initial condition, recursively sample from an LM:

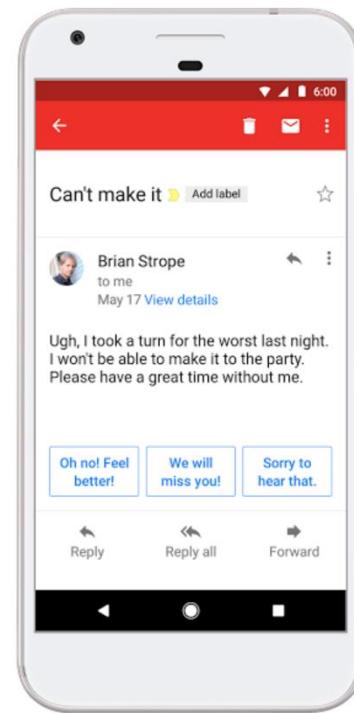
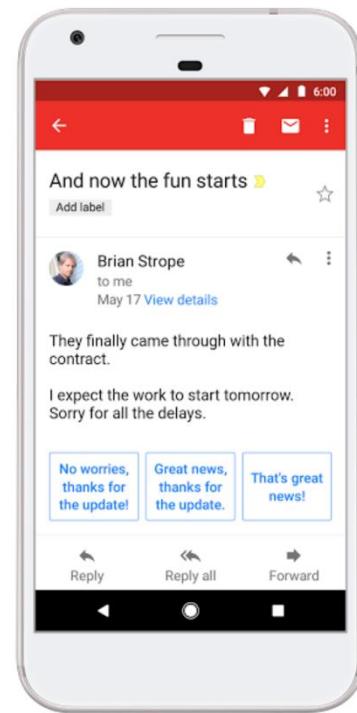
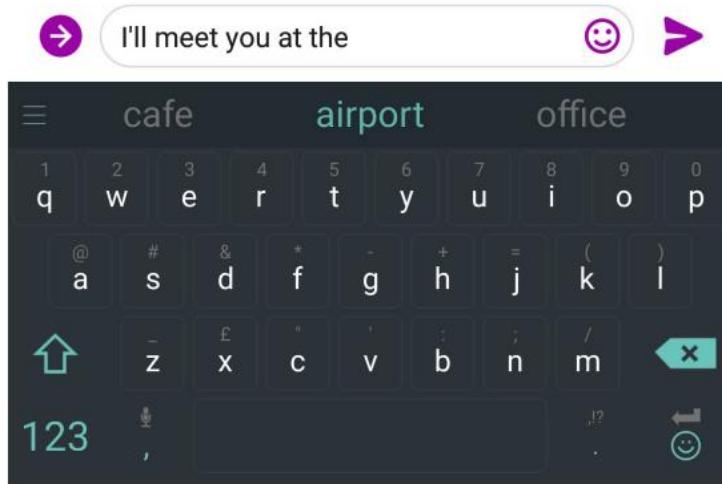
1. Sample from  $P(X | "Johns\ Hopkins\ is") \rightarrow "located"$
2. Sample from  $P(X | "Johns\ Hopkins\ is\ located") \rightarrow "at"$
3. Sample from  $P(X | "Johns\ Hopkins\ is\ located\ at") \rightarrow "the"$
4. Sample from  $P(X | "Johns\ Hopkins\ is\ located\ at\ the") \rightarrow "state"$
5. Sample from  $P(X | "Johns\ Hopkins\ is\ located\ at\ the\ state") \rightarrow "of"$
6. Sample from  $P(X | "Johns\ Hopkins\ is\ located\ at\ the\ state\ of") \rightarrow "Maryland"$
7. Sample from  $P(X | "Johns\ Hopkins\ is\ located\ at\ the\ state\ of\ Maryland") \rightarrow "EOS"$

# Why Care About Language Modeling?

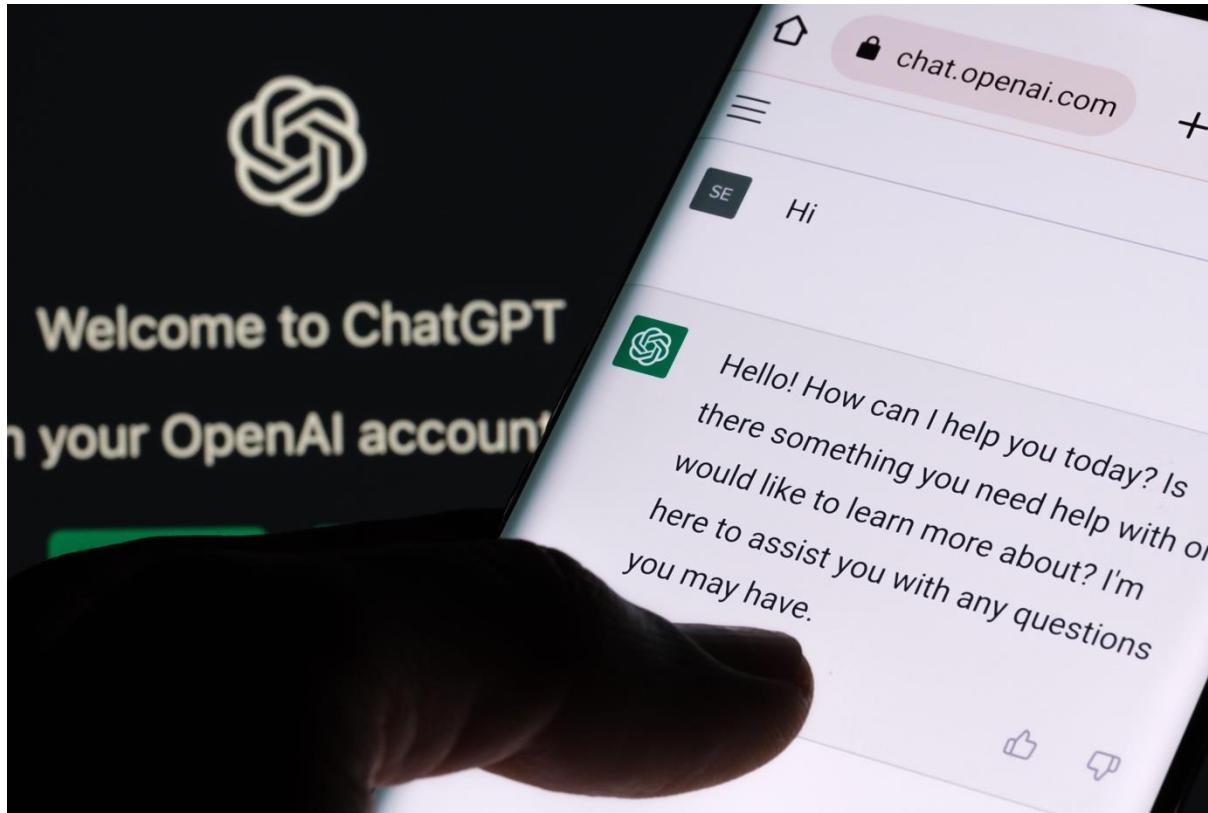
---

- Language Modeling is a **subcomponent superset** of many tasks:
  - Summarization
  - Machine translation
  - Spelling correction
  - Dialogue etc.
- Language Modeling is an effective proxy for **language understanding**.
  - **Effective ability to predict forthcoming words** requires **understanding of context/prefix**.

# You use Language Models every day!



# You use Language Models every day!



# Summary

---

- **Language modeling:** building probabilistic distribution over language.
- An accurate distribution of language enables us to solve many important tasks that involve language communication.
- **The remaining question:** how do you actually estimate this distribution?

# Language Modeling with Counting

# LMs as a Marginal Distribution

- Now the question is, how to estimate this distribution.

$$\mathbf{P}(X_t | X_1, \dots, X_{t-1})$$

How do we estimate these probabilities?

Let's just count!

$$P(\text{mat} | \text{the cat sat on the}) \approx \frac{\text{count}(\text{"the cat sat on the mat"})}{\text{count}(\text{"the cat sat on the"})}$$

Count how often  
“**the cat sat on the mat**”  
has appeared in the world (internet)!

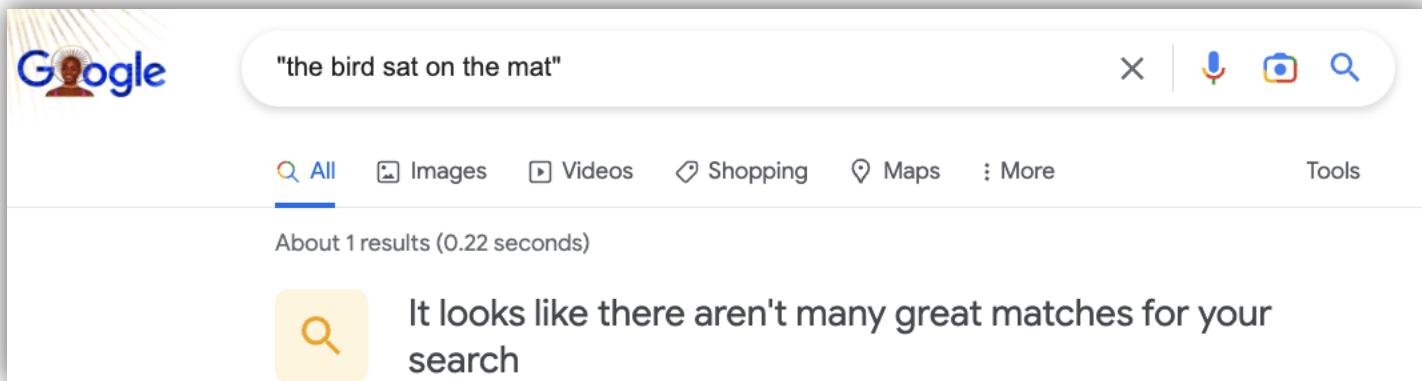
Divide that by, the count of  
“**the cat sat on the**”  
in the world (internet)!

$$\mathbf{P}(X_t | X_1, \dots, X_{t-1})$$

How do we estimate these probabilities?

Let's just count!

$$P(\text{mat} | \text{the cat sat on the}) \approx \frac{\text{count}(\text{"the cat sat on the mat"})}{\text{count}(\text{"the cat sat on the"})}$$



A screenshot of a Google search results page. The search bar contains the query "the bird sat on the mat". Below the search bar, there are navigation links for All, Images, Videos, Shopping, Maps, More, and Tools. A message indicates "About 1 results (0.22 seconds)". At the bottom, a yellow search icon is followed by the text: "It looks like there aren't many great matches for your search".

$$\mathbf{P}(X_t | X_1, \dots, X_{t-1})$$

How do we estimate these probabilities?

Let's just count!

$$P(\text{mat} | \text{the cat sat on the}) \approx \frac{\text{count}(\text{"the cat sat on the mat"})}{\text{count}(\text{"the cat sat on the"})}$$

Challenge: Increasing  $n$  makes **sparsity problems** worse.

Typically, we can't have  $n$  bigger than 5.

Some partial solutions (e.g., smoothing and backoffs)  
though still an open problem.

# Language Models: A History

- Shannon (1950): The redundancy and predictability (entropy) of English.



## Prediction and Entropy of Printed English

By C. E. SHANNON

(Manuscript Received Sept. 15, 1950)

A new method of estimating the entropy and redundancy of a language is described. This method exploits the knowledge of the language statistics possessed by those who speak the language, and depends on experimental results in prediction of the next letter when the preceding text is known. Results of experiments in prediction are given, and some properties of an ideal predictor are developed.



$$P(X_t | X_1, \dots, X_{t-1})$$

Shannon (1950) built an approximate language model with word co-occurrences.

**Markov assumptions:** every node in a Bayesian network is **conditionally independent** of its non-descendants, **given its parents**.

1<sup>st</sup> order approximation:

$$P(\text{mat} | \text{the cat sat on the}) \approx P(\text{mat} | \text{the})$$

1 element  
\_\_\_\_\_



$$P(X_t | X_1, \dots, X_{t-1})$$

Andrey Markov

Shannon (1950) built an approximate language model with word co-occurrences.

**Markov assumptions:** every node in a Bayesian network is **conditionally independent** of its non-descendants, **given its parents**.

2<sup>nd</sup> order approximation:

$$P(\text{mat} | \text{the cat sat on the}) \approx P(\text{mat} | \text{on the})$$

2 elements  
\_\_\_\_\_



$$P(X_t | X_1, \dots, X_{t-1})$$

Andrey Markov

Shannon (1950) built an approximate language model with word co-occurrences.

**Markov assumptions:** every node in a Bayesian network is **conditionally independent** of its non-descendants, **given its parents**.

3<sup>rd</sup> order approximation:

$$P(\text{mat} | \text{the cat sat on the}) \approx P(\text{mat} | \text{sat on the})$$

3 elements



$$P(X_t | X_1, \dots, X_{t-1})$$

Andrey Markov

Shannon (1950) built an approximate language model with word co-occurrences.

Then, we can use counts of approximate conditional probability.  
Using the 3<sup>rd</sup> order approximation, we can:

$$P(\text{mat} | \text{the cat sat on the}) \approx P(\text{mat} | \text{sat on the}) = \frac{\text{count("sat on the mat")}}{\text{count("on the mat")}}$$

# Understanding Sparsity: A Thought Experiment

---

- How common are zero-probabilities? 😐
- **Example:** Shakespeare as a text corpus
  - The size vocab used by Shakespeare:  $|V|=29,066$
  - Shakespeare produced:  $\sim 300,000$  bigrams
    - Out of  $|V|^2 = 844$  million possible bigrams
      - (some of them don't make sense, but ok!)
- So, **99.96%** of the possible bigrams are **never seen** (hence, have zero entries for bigram counts).

# N-gram Language Models

- **Terminology:** *n*-gram is a chunk of *n* consecutive words:
  - unigrams: “cat”, “mat”, “sat”, ...
  - bigrams: “the cat”, “cat sat”, “sat on”, ...
  - trigrams: “the cat sat”, “cat sat on”, “sat on the”, ...
  - four-grams: “the cat sat on”, “cat sat on the”, “sat on the mat”, ...
- *n*-gram language model:

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t | \underbrace{X_{t-n+1}, \dots, X_{t-1}}_{n-1 \text{ elements}})$$

# Generation from N-Gram Models

---

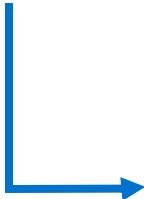
- You can build a simple **trigram** Language Model over a 1.7 million words corpus in a few seconds on your laptop\*

# Generation from N-Gram Models

- You can build a simple trigram Language Model over a 1.7 million words corpus in a few seconds on your laptop\*

today the   

get probability  
distribution



company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem: not  
much granularity in the  
probability distribution

Otherwise, seems reasonable!

# Generation from N-Gram Models

- Now we can sample from this mode:

today the   

get probability  
distribution



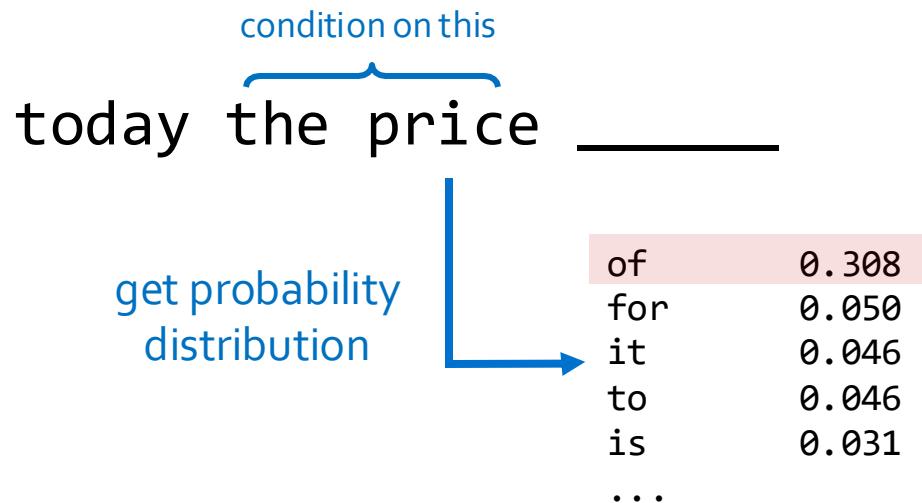
company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem: not  
much granularity in the  
probability distribution

Otherwise, seems reasonable!

# Generation from N-Gram Models

- Now we can sample from this mode:

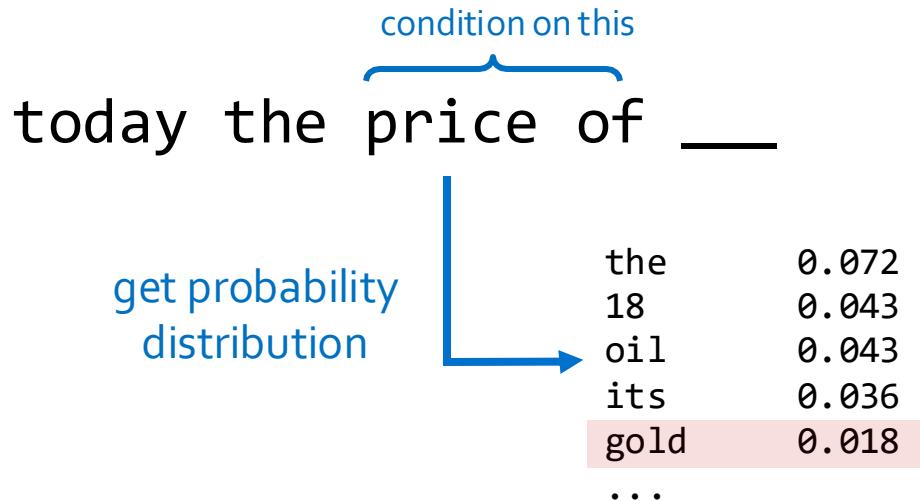


Sparsity problem: not much granularity in the probability distribution

Otherwise, seems reasonable!

# Generation from N-Gram Models

- Now we can sample from this mode:



Sparsity problem: not much granularity in the probability distribution

Otherwise, seems reasonable!

# N-Gram Models in Practice

- Now we can sample from this mode:

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 quite incoherent! To improve coherence, one may consider increasing larger than 3-grams, but that would worsen the sparsity problem!

# N-Grams LMs and Long-range Dependencies

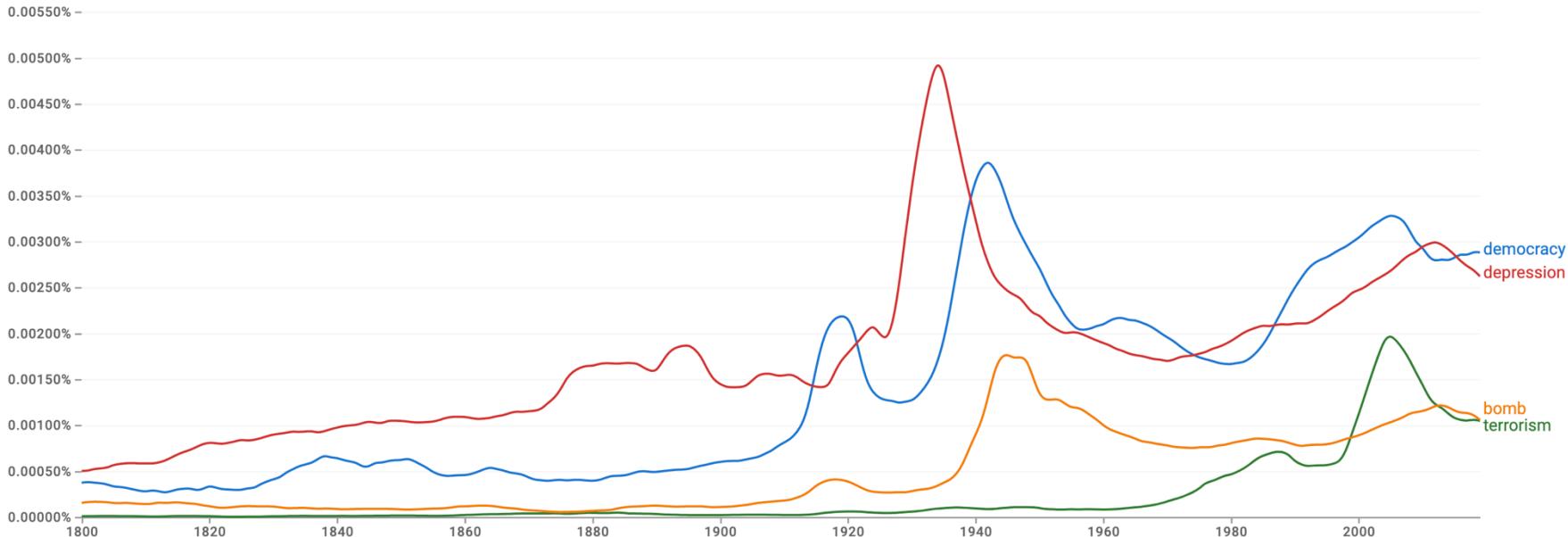
---

- In general, count-based LMs are insufficient models of language because language has **long-distance dependencies**:

“**The computer** which I had just put into the machine room on the fifth floor **crashed**.”

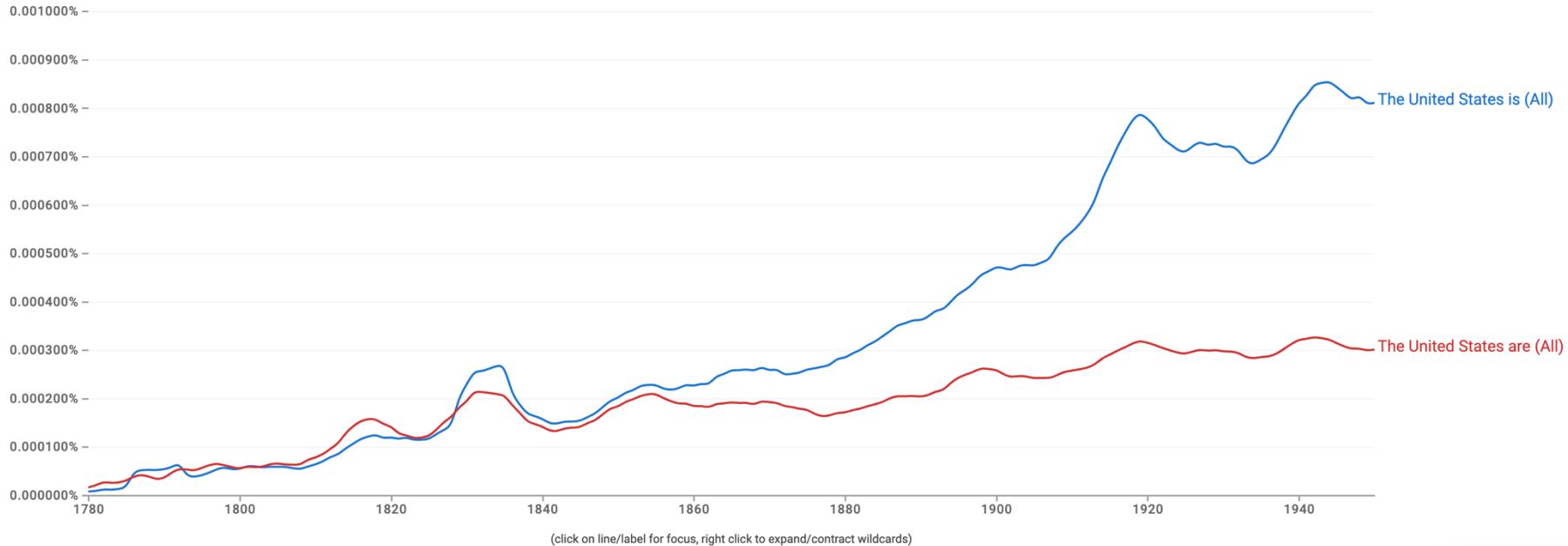
# Pre-Computed N-Grams

Google Books Ngram Viewer



# Pre-Computed N-Grams

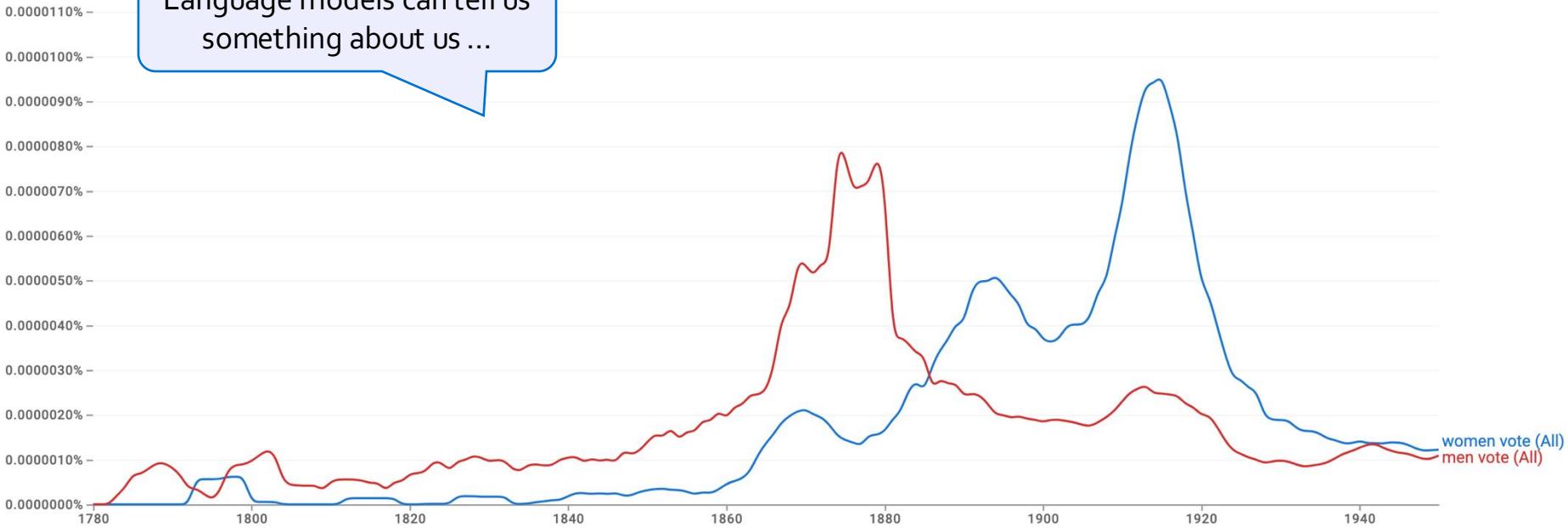
Google Books Ngram Viewer



# Pre-Computed N-Grams

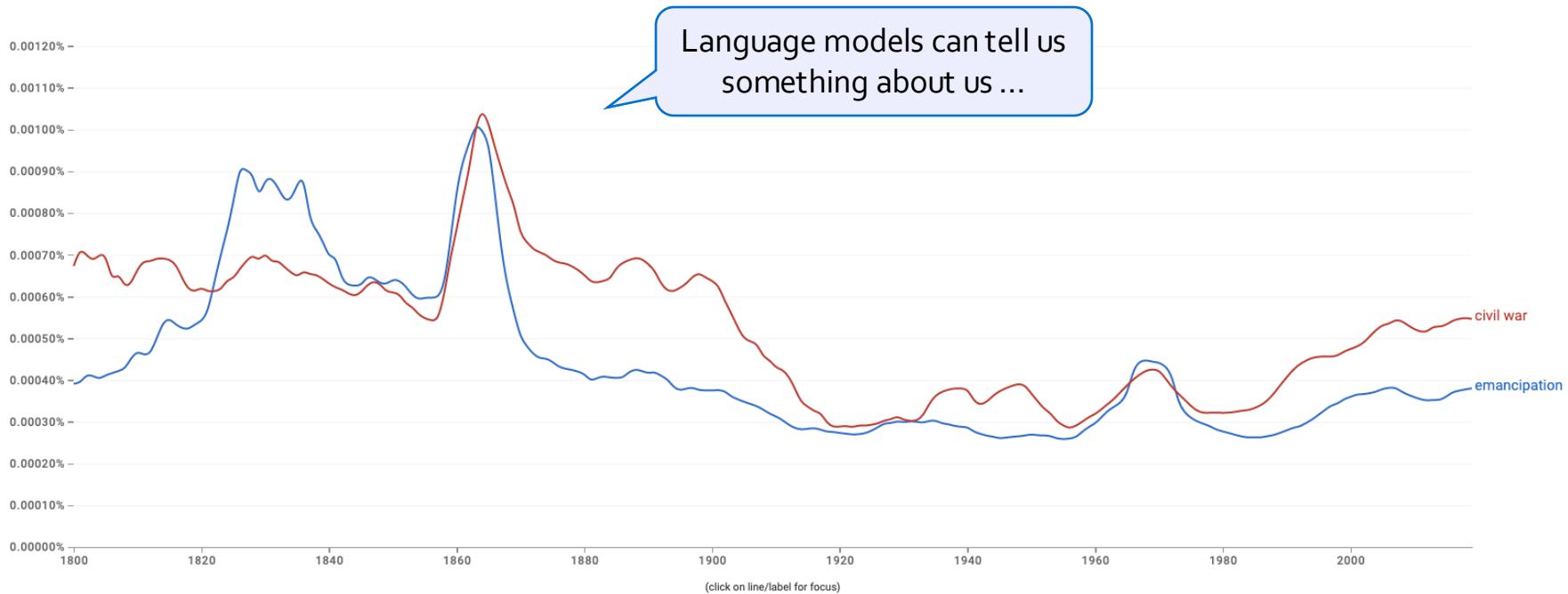
Google Books Ngram Viewer

Language models can tell us something about us ...



# Pre-Computed N-Grams

Google Books Ngram Viewer



# N-Gram Language Models, A Historical Highlight

"Every time I fire a linguist, the performance of the speech recognizer goes up"!!



Fred Jelinek  
(1932-2010)

- Probabilistic n-gram models of text generation [Jelinek+ 1980's, ...]
  - Applications: Speech Recognition, Machine Translation

532

PROCEEDINGS OF THE IEEE, VOL. 64, NO. 4, APRIL 1976

## Continuous Speech Recognition by Statistical Methods

FREDERICK JELINEK, FELLOW, IEEE

*Abstract—Statistical methods useful in automatic recognition of continuous speech are described. They concern modeling of a speaker and of an acoustic processor, extraction of the models' statistical parameters, and hypothesis search procedures and likelihood computations of linguistic decoding. Experimental results are presented that indicate the power of the methods.*

utterance models used will incorporate more grammatical features, and statistics will have been grafted onto grammatical models. Most methods presented here concern modeling of the speaker's and acoustic processor's performance and should, therefore, be universally useful.

Automatic recognition of continuous (English) speech is an

# Summary

---

- Learning a language model ~ learning **conditional probabilities** over language.
- One approach to estimating these probabilities: **counting word co-occurrences**.
- Challenges:
  - Word co-occurrences become **rare** for long sequences. (the sparsity issue)
  - But language understanding requires **long-range** dependencies.
- We need a better alternative! 😱
- **Next:** Measuring quality of language models.

# How Good are Language Models?

# Evaluating Language Models

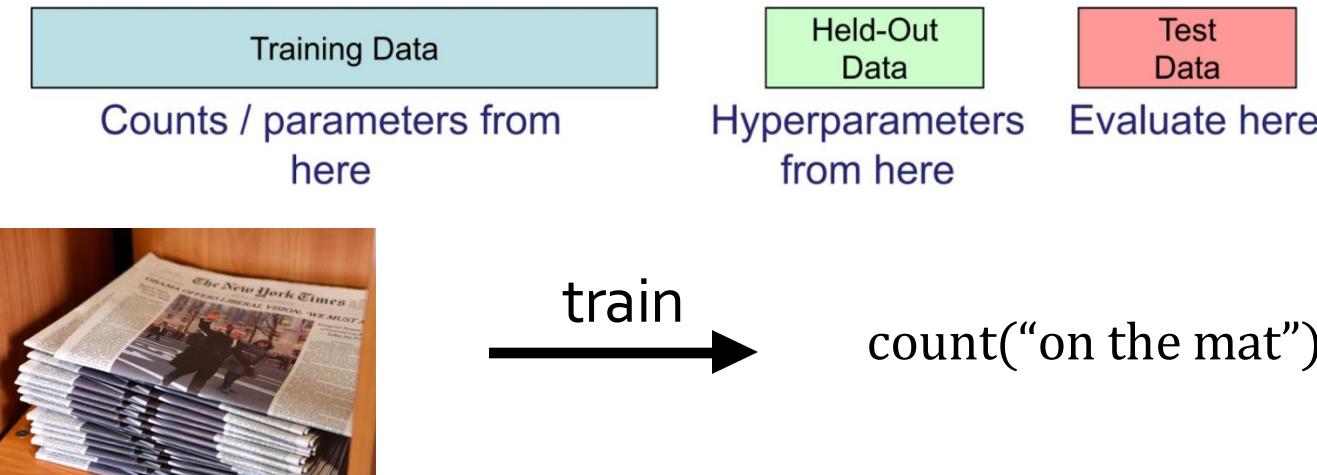
---

- Does our language model prefer good sentences to bad ones?
  - Assign higher probability to “real” or “frequently observed” sentences
  - Than “ungrammatical” or “rarely observed” sentences?
- We test the model’s performance on data we haven’t seen.

# Evaluating Language Models

Setup:

- Train it on a suitable training documents.
- Evaluate their predictions on different, unseen documents.
- An evaluation metric tells us how well our model does on the test set.



# Evaluating Language Models: Example

Setup:

- Train it on a suitable training documents.
- Evaluate their predictions on different, unseen documents.
- An evaluation metric tells us how well our model does on the test set.

Example: I use a bunch of New York Times articles to build a bigram probability table



A good language model  
should assign a high  
probability to held-out text!

train

count("on the mat")

eval



Now I'm going to evaluate the probability of some heldout data using our bigram table

# Be Careful About Data Leakage!

**Advice from a grandpa** 🧑:

- Don't allow test sentences to leak into training set.
- Otherwise, you will assign it an artificially high probability (==cheating).

Example: I use a bunch of New York Times articles to build a bigram probability table



A good language model  
should assign a high  
probability to held-out text!

train →

count("on the mat")

eval →



Now I'm going to evaluate the  
probability of some heldout  
data using our bigram table

# Evaluating Language Models: Intrinsic vs Extrinsic

- **Intrinsic:** measure how good we are at modeling language
- **Extrinsic:** build a new language model, use it for some task (MT, ASR, etc.)

Example: I use a bunch of New York Times articles to build a bigram probability table



train →



Now I'm going to evaluate the probability of some heldout data using our bigram table



# Evaluation Metric for Language Modeling: Perplexity

---

- **Perplexity** is the inverse probability of the test set, normalized by the number of words:

$$\text{ppl}(w_1, \dots, w_n) = \mathbf{P}(w_1, w_2, \dots, w_n)^{-\frac{1}{n}}$$

- A measure of **predictive quality** of a language model.
- **Minimizing** perplexity is the same as **maximizing** probability

# Evaluation Metric for Language Modeling: Perplexity

---

- **Perplexity** is the inverse probability of the test set, normalized by the number of words:

$$\text{ppl}(w_1, \dots, w_n) = \mathbf{P}(w_1, w_2, \dots, w_n)^{-\frac{1}{n}}$$

# Evaluation Metric for Language Modeling: Perplexity

- **Perplexity** is the inverse probability of the test set, normalized by the number of words:

$$\begin{aligned} \text{ppl}(w_1, \dots, w_n) &= \mathbf{P}(w_1, w_2, \dots, w_n)^{-\frac{1}{n}} \\ &= \sqrt[n]{\frac{1}{\mathbf{P}(w_1, w_2, \dots, w_n)}} = \sqrt[n]{\prod_{i=1}^n \frac{1}{\mathbf{P}(w_i | w_{<i})}} \quad \text{chain rule} \\ &= 2^H, \text{ where} \end{aligned}$$

$$H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

# Evaluation Metric for Language Modeling: Perplexity

- In practice, we prefer to use **log**-probabilities (also known as “logits”)
- We can rewrite perplexity formula in terms of log-probs:

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 P(w_i | w_1, \dots, w_{i-1})$$

**Recap:** Definition of cross-entropy between two distributions:

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x)$$

Can be interpreted as cross-entropy between LM prob and language prob. **Why?**

# Evaluation Metric for Language Modeling: Perplexity

---

- In practice, we prefer to use **log**-probabilities (also known as “logits”)
- We can rewrite perplexity formula in terms of log-probs:

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

- **Perplexity** for n-grams:
  - Unigrams:  $H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i)$
  - Bigrams:  $H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_{i-1})$
  - Trigrams:  $H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_{i-2}, w_{i-1})$
  - ...

# Intuition-building Quizzes (1)

- In practice, we prefer to use **log**-probabilities (also known as “logits”)
- We can rewrite perplexity formula in terms of log-probs:

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

- **Quiz:** let's suppose we have a sentence  $w_1, \dots, w_n$  and it's fixed. Our model will correctly guess each word with probability 1/5. What is perplexity of our model?

$$H = -\frac{1}{n} \left[ \log_2 \left( \frac{1}{5} \right) + \dots + \log_2 \left( \frac{1}{5} \right) \right] = -\log \left( \frac{1}{5} \right) \Rightarrow \text{ppl}(D) = 5$$

**Intuition:** the model is indecisive among 5 choices.

# Intuition-building Quizzes (2)

- In practice, we prefer to use **log**-probabilities (also known as “logits”)
- We can rewrite perplexity formula in terms of log-probs:

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

- **Quiz:** let's we evaluate an **exact** (!! model of language, i.e., our model always knows what exact word should follow a given context. What is the perplexity of this model?

$$\forall w \in V: \mathbf{P}(w_i | w_{1:i-1}) = 1 \Rightarrow \text{ppl}(D) = 2^{-\frac{1}{2} n \log_2 1} = 1$$

**Intuition:** the model is indecisive among 1 (the right!) choice!

# Intuition-building Quizzes (3)

- In practice, we prefer to use **log**-probabilities (also known as “logits”)
- We can rewrite perplexity formula in terms of log-probs:

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

- **Quiz:** let's evaluate a **confused** (!! ) model of language, i.e., our model has no idea what word should follow each context—it always chooses a uniformly random word. What is the perplexity of this model?

$$\forall w \in V: \mathbf{P}(w | w_{1:i-1}) = \frac{1}{|V|} \Rightarrow \text{ppl}(D) = 2^{-\frac{1}{n} n \log_2 \frac{1}{|V|}} = |V|$$

**Intuition:** the model is indecisive among all the vocabulary terms.

# Perplexity: Summary

$$\text{ppl}(w_1, \dots, w_n) = 2^H, \text{ where } H = -\frac{1}{n} \sum_{i=1}^n \log_2 \mathbf{P}(w_i | w_1, \dots, w_{i-1})$$

- Perplexity is a measure of model's uncertainty about next word (aka "average branching factor").
  - The larger the number of vocabulary, the more options there to choose from.
  - (the choice of atomic units of language impacts PPL — more on this later)
- Perplexity ranges between 1 and  $|V|$ .
- We prefer LMs with lower perplexity.

# Lower perplexity == Better Model

- Training on 38 million words, test 1.5 million words, Wall Street Journal

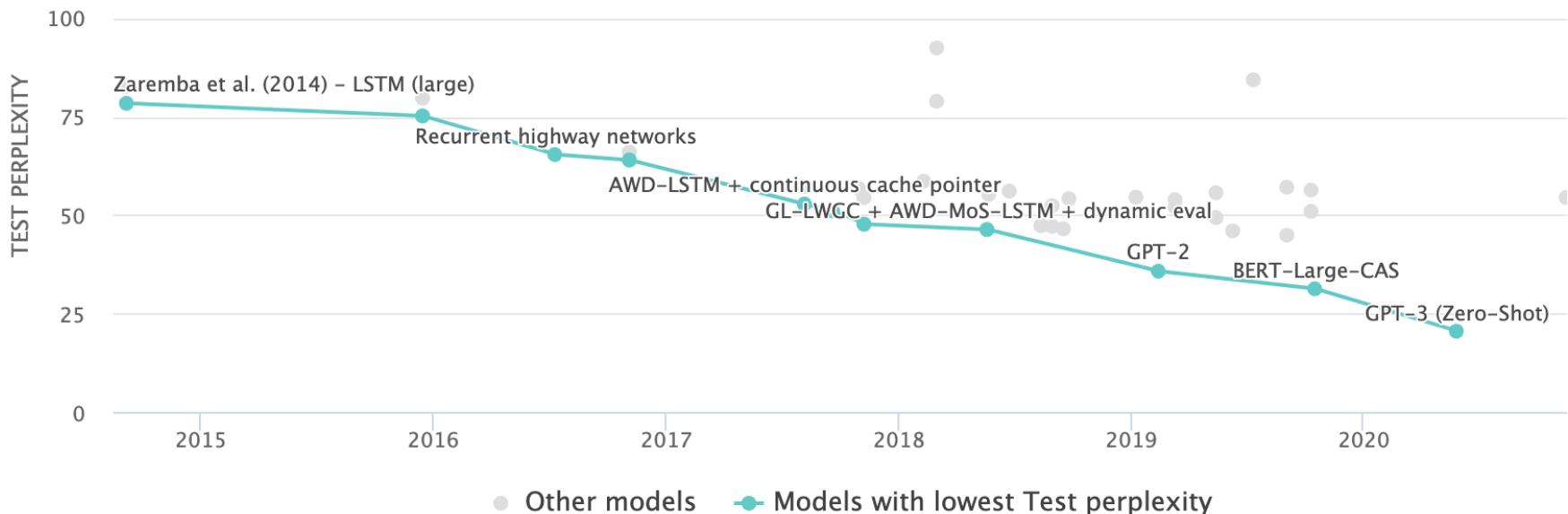
N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

Lower is better

Note these evaluations are done on data that was not used for “counting.” (no cheating!!)

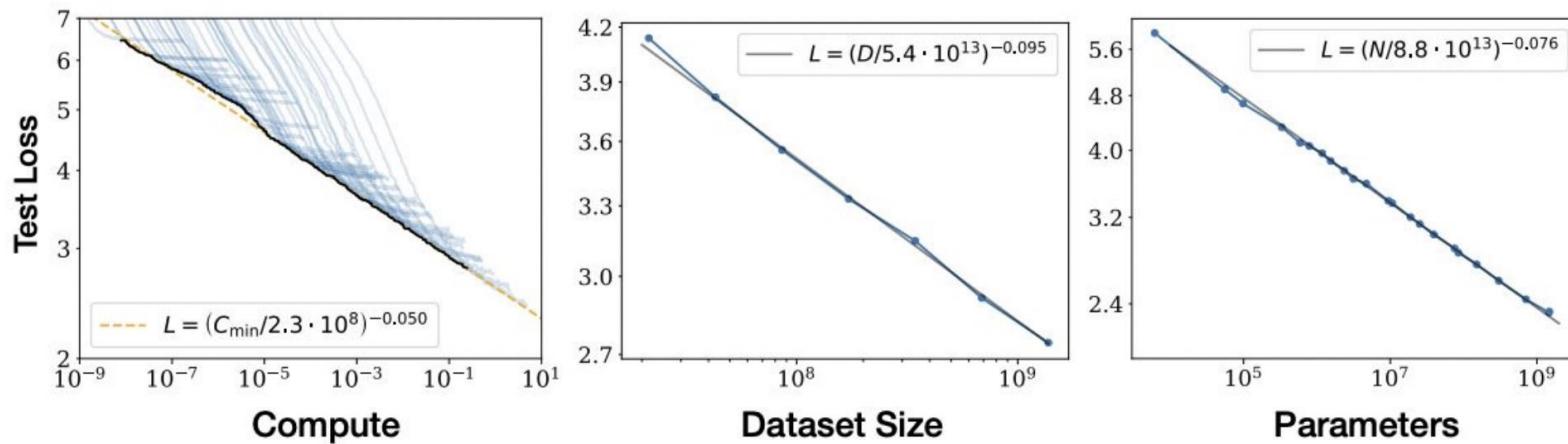
# Lower perplexity == Better Model

The PPL of modern language models have consistently been going down.



# Lower perplexity == Better Model

The PPL of modern language models have consistently been going down.



# Summary

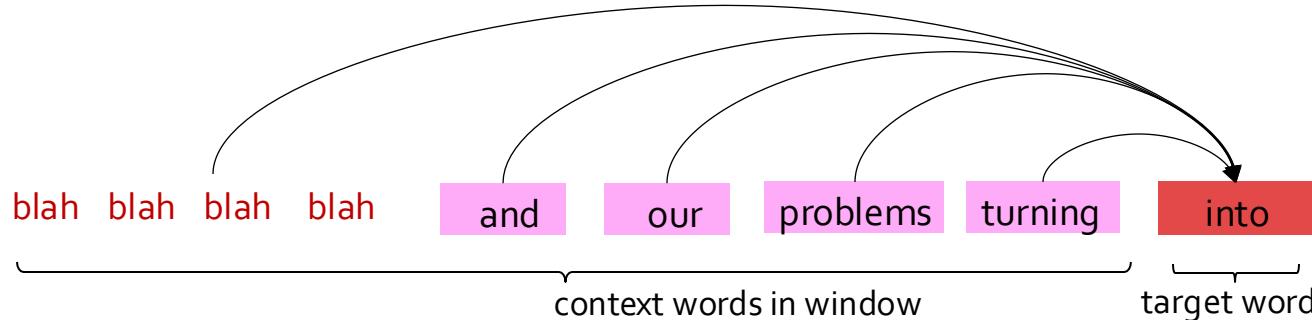
---

- Language Models (LM): distributions over language
- Measuring LM quality: use perplexity on held-out data.
- Count-based LMs have limitations.
  - Challenge with large N's: sparsity problem — many zero counts/probs.
  - Challenge with small N's: lack of long-range dependencies.
- Next: Rethinking language modeling as a statistical learning problem.

# Beyond Counting: Language Models as a Learning Problem

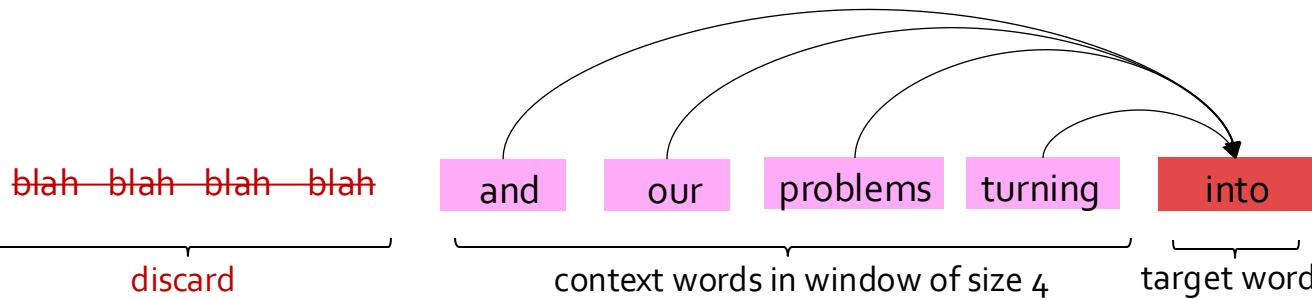
# LM as a Machine Learning Problem

- Given the embeddings of the context, predict the word on the right side.
  - Dropping the right context for simplicity -- not a fundamental limitation.
- Discard anything beyond its context window



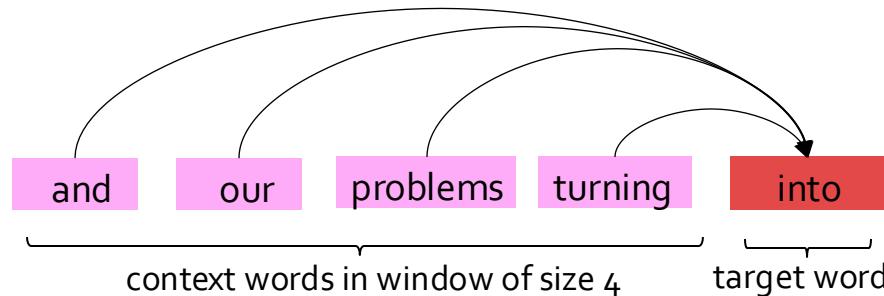
# LM as a Machine Learning Problem

- Given the embeddings of the context, predict the word on the right side.
  - Dropping the right context for simplicity -- not a fundamental limitation.
- Discard anything beyond its context window



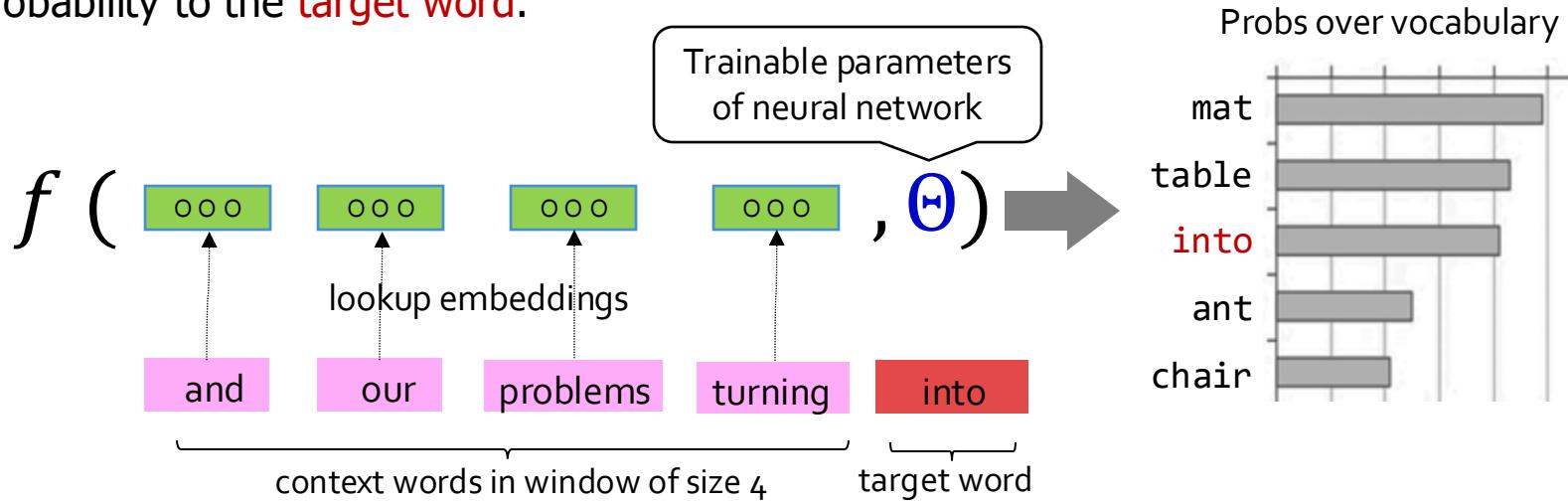
# LM as a Machine Learning Problem

- Given the embeddings of the context, predict the word on the right side.
  - Dropping the right context for simplicity -- not a fundamental limitation.
- Discard anything beyond its context window



# A Fixed-Window Neural LM

- Given the embeddings of the **context**, predict a **target word** on the right side.
  - Dropping the right context for simplicity -- not a fundamental limitation.
- Training this model is basically optimizing its parameters  $\Theta$  such that it assigns high probability to the **target word**.



# A Fixed-Window Neural LM

- It will also lay the foundation for the future models (recurrent nets, transformers, ...)
- But first we need to figure out how to train neural networks!

How do you build  
this function?

$$f($$



Trainable parameters  
of neural network



$$, \Theta)$$



Neural Networks  
for rescue!

lookup embeddings

our

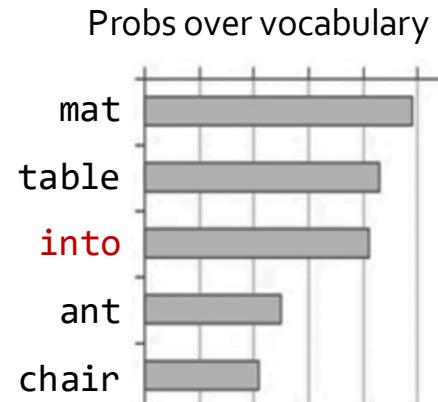
problems

turning

into

context words in window of size 4

target word



# From Counting (N-Gram) to Neural Models

---

- n-gram models of text generation [Jelinek+ 1980's, ...]
    - Applications: Speech Recognition, Machine Translation
  - “Shallow” statistical/neural language models (2000’s) [Bengio+ 1999 & 2001, ...]
- 

NeurIPS 2000

## A Neural Probabilistic Language Model

---

**Yoshua Bengio\*, Réjean Ducharme and Pascal Vincent**  
Département d’Informatique et Recherche Opérationnelle  
Centre de Recherche Mathématiques  
Université de Montréal  
Montréal, Québec, Canada, H3C 3J7  
*{bengioy,ducharme,vincentp}@iro.umontreal.ca*

# Summary

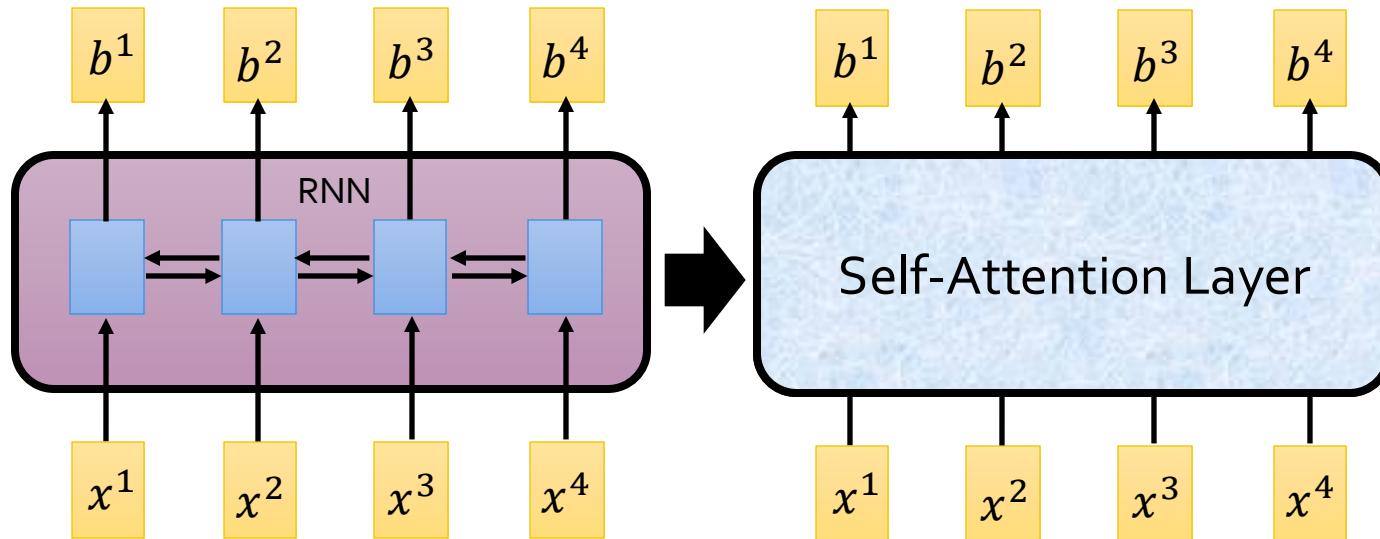
---

- **Language Modeling (LM)**, a useful predictive objective for language
- **Perplexity**, a measure of an LM's predictive ability
- **N-gram models** (~1980 to early 2000's),
  - Early instances of LMs
  - Difficult to scale to large window sizes
- **Shallow neural LMs** (early and mid-2000's),
  - We will need in coming sessions that one can build these models with neural networks.
  - These will be effective predictive models based on feed-forward networks

# Self-Attention

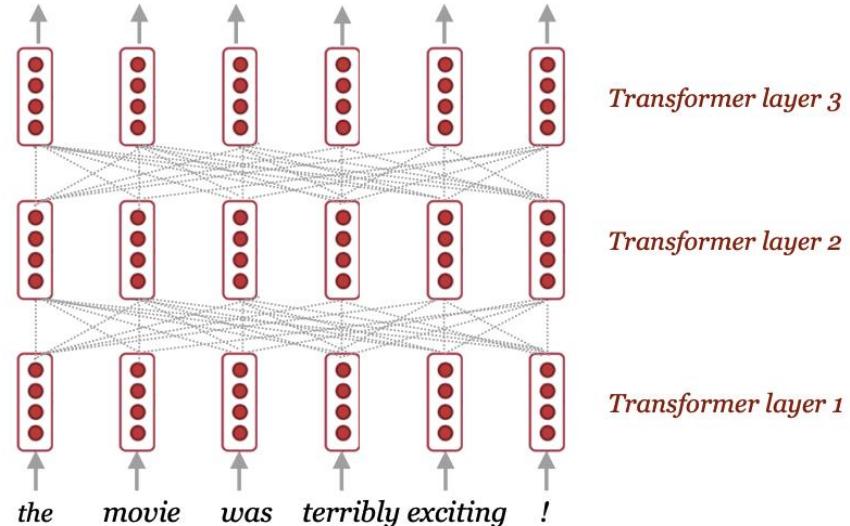
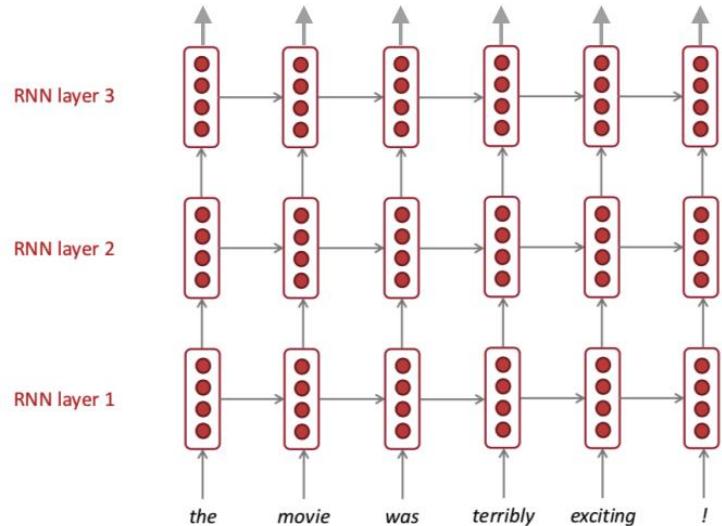
# Self-Attention

- $b^i$  is obtained based on the whole input sequence.
- can be parallelly computed.



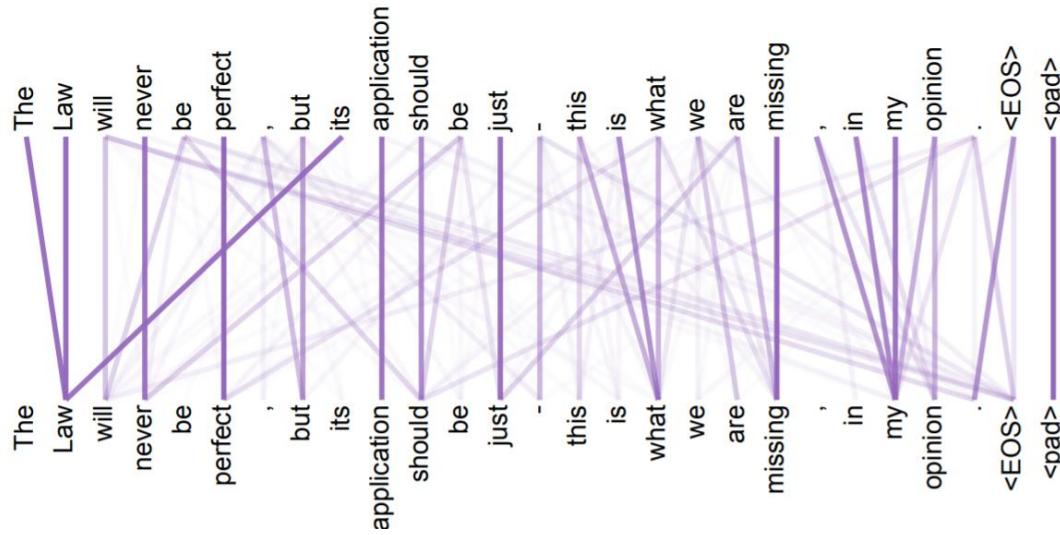
Idea: replace any thing done by RNN with **self-attention**.

# RNN vs Transformer



# Attention

- Core idea: build a mechanism to focus ("attend") on a particular part of the context.



# Defining Self-Attention

---

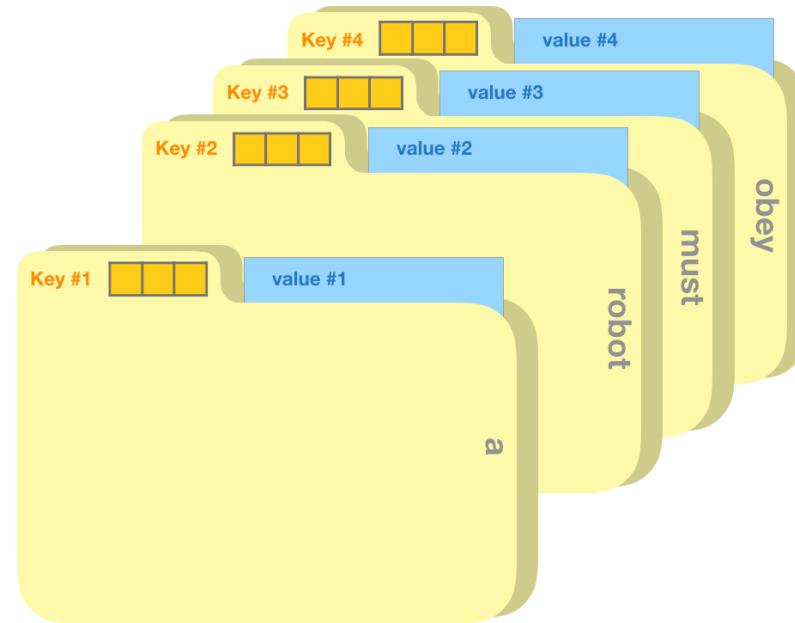
- Terminology:
  - **Query**: to match others
  - **Key**: to be matched
  - **Value**: information to be extracted

# Defining Self-Attention

- Terminology:
  - **Query**: to match others
  - **Key**: to be matched
  - **Value**: information to be

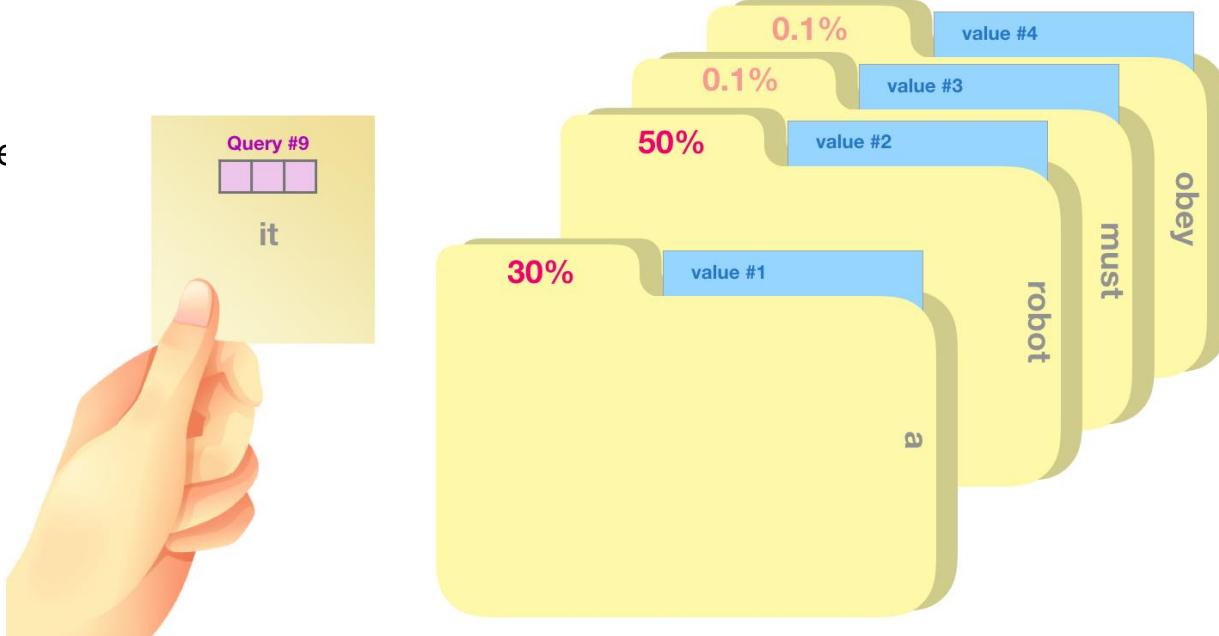


An analogy ....



# Defining Self-Attention

- Terminology:
  - **Query**: to match others
  - **Key**: to be matched
  - **Value**: information to be



*q*: query (to match others)

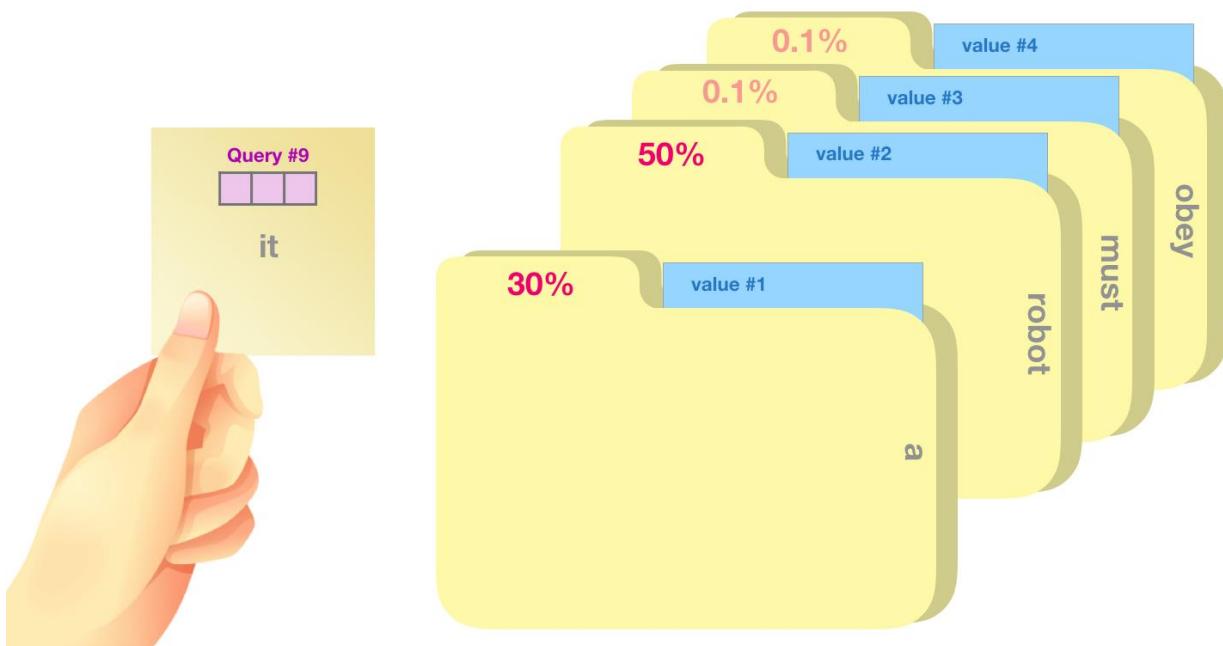
$$q_i = W^q x_i$$

*k*: key (to be matched)

$$k_i = W^k x_i$$

*v*: value (information to be extracted)

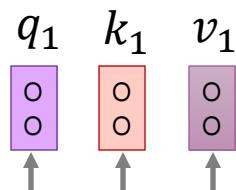
$$v_i = W^v x_i$$



$q$ : query (to match others)  
 $q_i = W^q x_i$

$k$ : key (to be matched)  
 $k_i = W^k x_i$

$v$ : value (information to be extracted)  
 $v_i = W^v x_i$



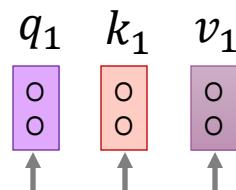
$x_1$

The

*q*: query (to match others)  
 $q_i = W^q x_i$

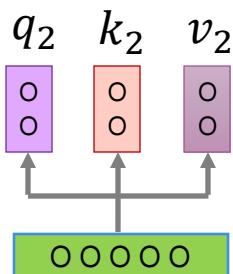
*k*: key (to be matched)  
 $k_i = W^k x_i$

*v*: value (information to be extracted)  
 $v_i = W^v x_i$



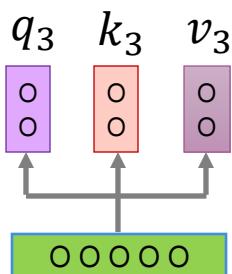
$x_1$

The



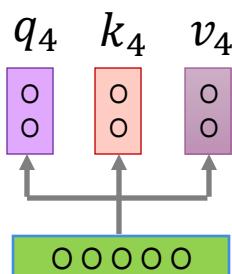
$x_2$

cat



$x_3$

sat



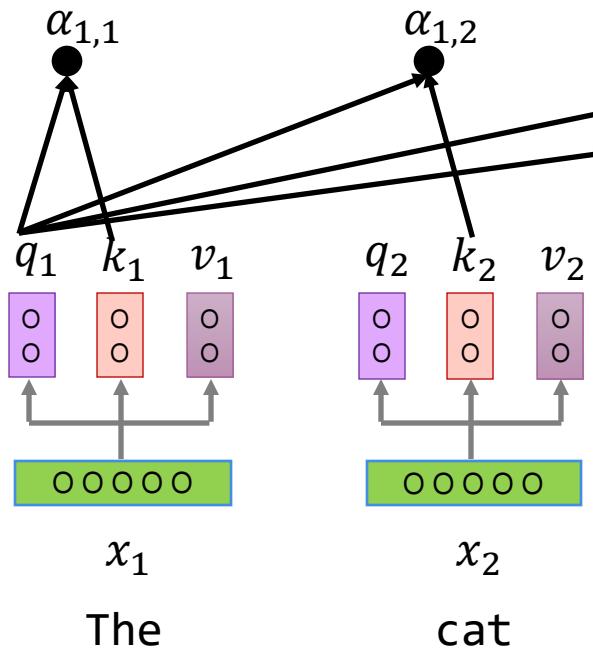
$x_4$

on

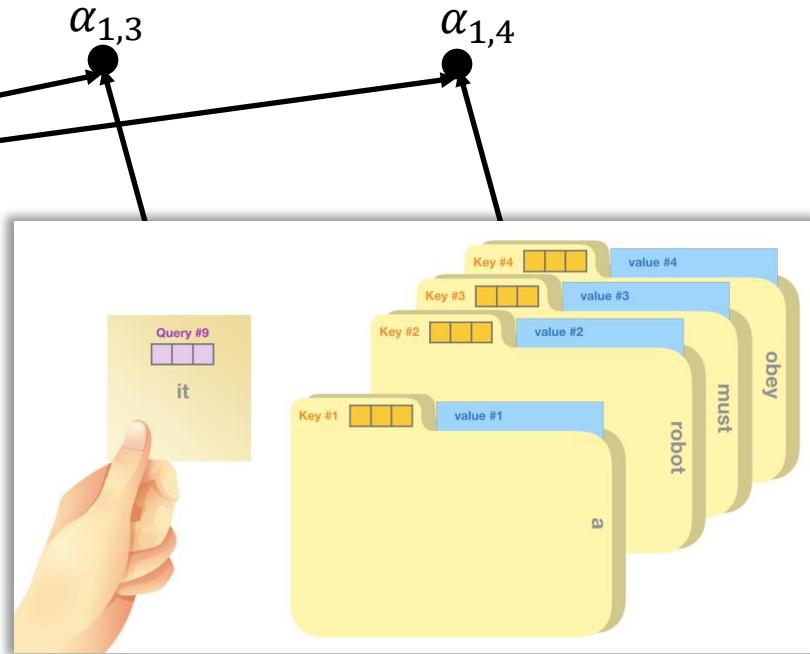
$$\alpha_{1,i} = \frac{q^1 \cdot k^i}{\sqrt{d}}$$

Scaled dot product

How much  
should "The"  
attend to other  
positions?

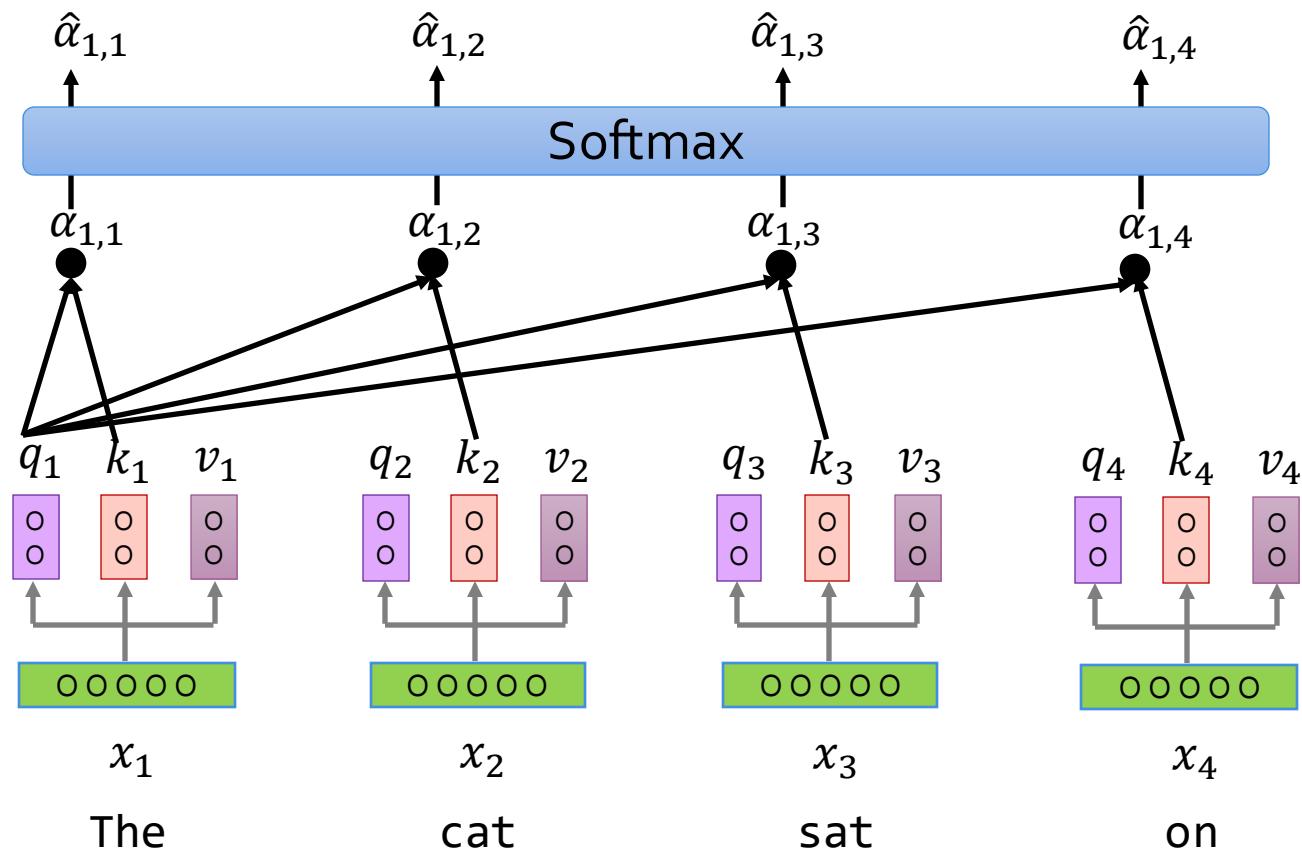


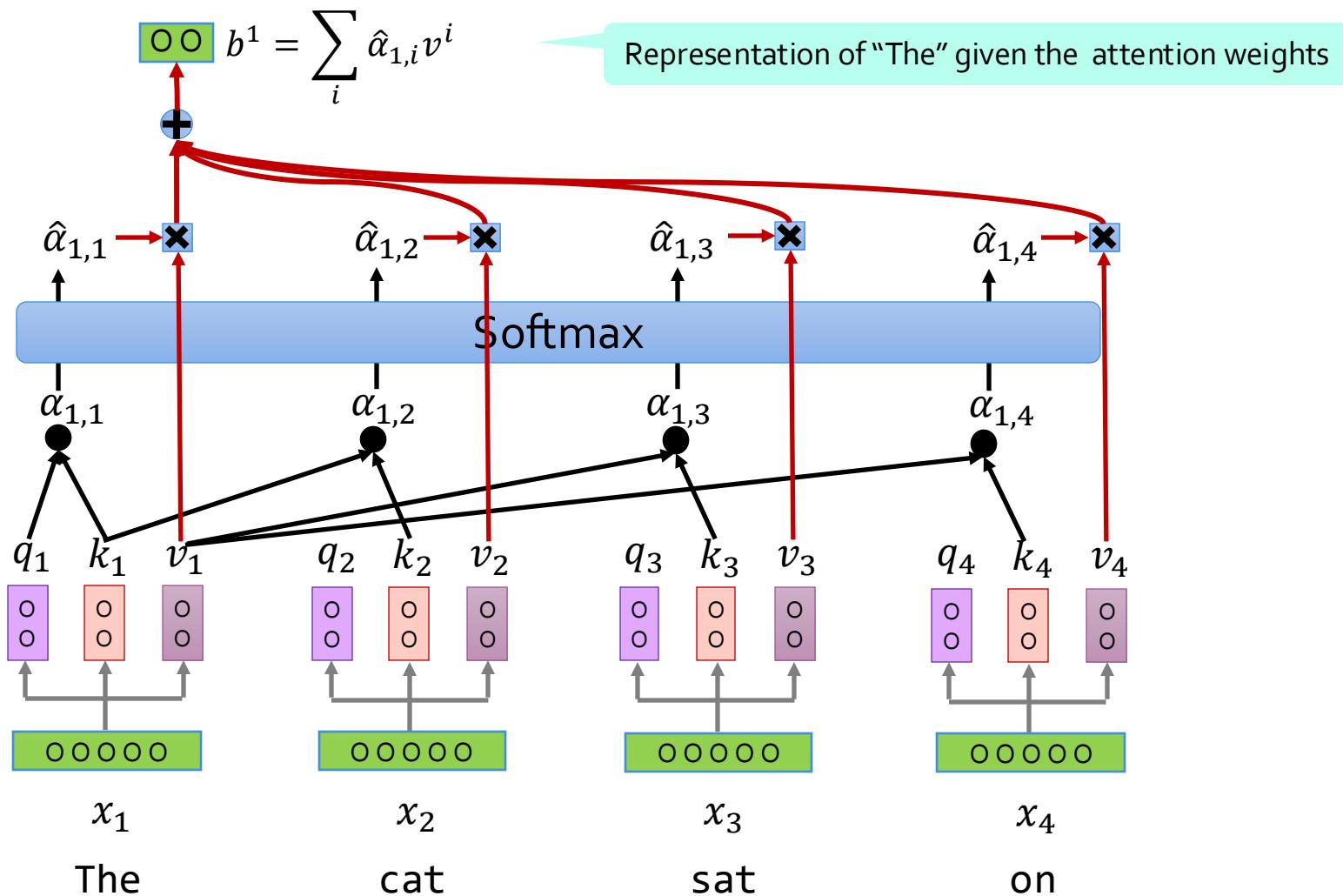
*q*: query (to match others)  
*k*: key (to be matched)  
*v*: value (information to be extracted)



$$\sigma(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

How much  
should "The"  
attend to other  
positions?

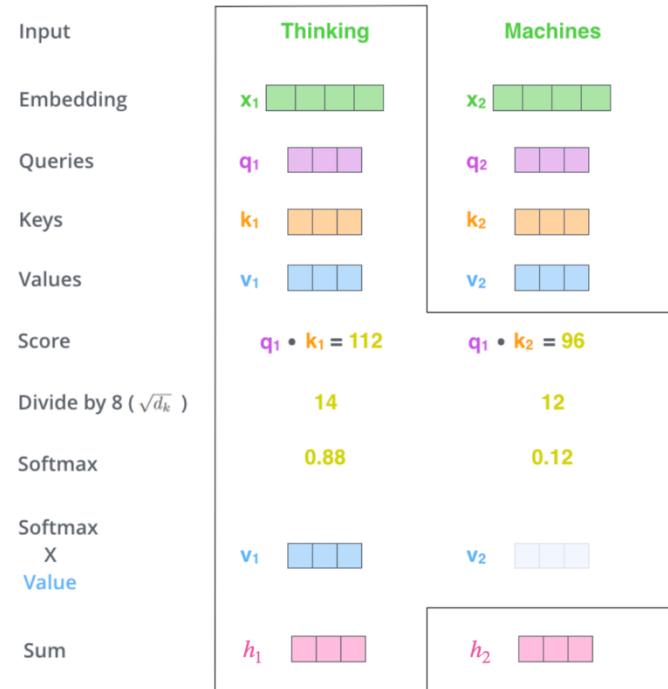




# Question

- What would be the output vector for the word “Thinking”?

- (a)  $0.5\mathbf{v}_1 + 0.5\mathbf{v}_2$
- (b)  $0.54\mathbf{v}_1 + 0.46\mathbf{v}_2$
- (c)  $0.88\mathbf{v}_1 + 0.12\mathbf{v}_2$
- (d)  $0.12\mathbf{v}_1 + 0.88\mathbf{v}_2$



# Self-Attention: Matrix Notation

$$X \in \mathbb{R}^{n \times d_1} \quad (n = \text{input length})$$

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

$$W^Q \in \mathbb{R}^{d_1 \times d_q}, W^K \in \mathbb{R}^{d_1 \times d_k}, W^V \in \mathbb{R}^{d_1 \times d_v}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$n \times d_q$        $d_k \times n$

$n$ :

Q: What is this softmax operation?

$$\begin{aligned} & \text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V} \\ &= \mathbf{H} \end{aligned}$$

# Self-Attention

- Can write it in matrix form:
- Given input  $\mathbf{x}$ :

$$Q = \mathbf{W}^q \mathbf{x}$$

$$K = \mathbf{W}^k \mathbf{x}$$

$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



hardmaru  
@hardmaru

...

The most important formula in deep learning after 2018

## Self-Attention

**What is self-attention?** Self-attention calculates a weighted average of feature representations with the weight proportional to a similarity score between pairs of representations. Formally, an input sequence of  $n$  tokens of dimensions  $d$ ,  $X \in \mathbf{R}^{n \times d}$ , is projected using three matrices  $W_Q \in \mathbf{R}^{d \times d_q}$ ,  $W_K \in \mathbf{R}^{d \times d_k}$ , and  $W_V \in \mathbf{R}^{d \times d_v}$  to extract feature representations  $Q$ ,  $K$ , and  $V$ , referred to as query, key, and value respectively with  $d_k = d_q$ . The outputs  $Q$ ,  $K$ ,  $V$  are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V. \quad (1)$$

So, self-attention can be written as,

$$S = D(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_q}}\right)V, \quad (2)$$

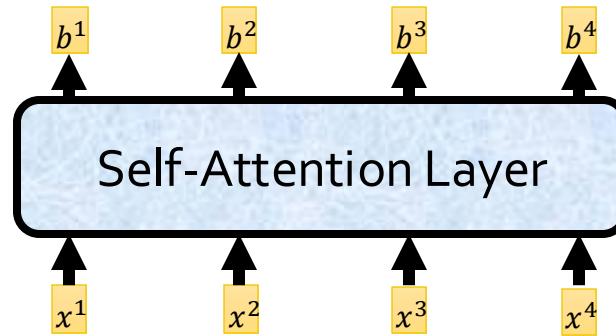
where softmax denotes a *row-wise* softmax normalization function. Thus, each element in  $S$  depends on all other elements in the same row.

9:08 PM · Feb 9, 2021 · Twitter Web App

553 Retweets 42 Quote Tweets 3,338 Likes

# Self-Attention: Back to Big Picture

- **Attention** is a powerful mechanism to create context-aware representations
- A way to focus on select parts of the input



- Better at maintaining **long-distance dependencies** in the context.

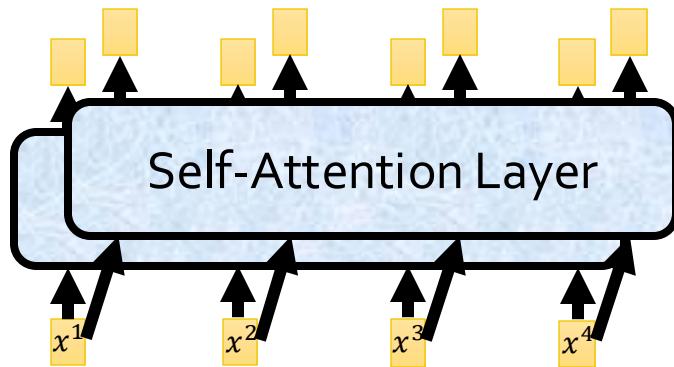
# Properties of Self-Attention

Layer Type	Complexity per Layer	Sequential Operations
Self-Attention	$O(n^2 \cdot d)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$

- $n$  = sequence length,  $d$  = hidden dimension
- Quadratic complexity, but:
  - $O(1)$  sequential operations (not linear like in RNN)
- Efficient implementations

# Multi-Headed Self-Attention

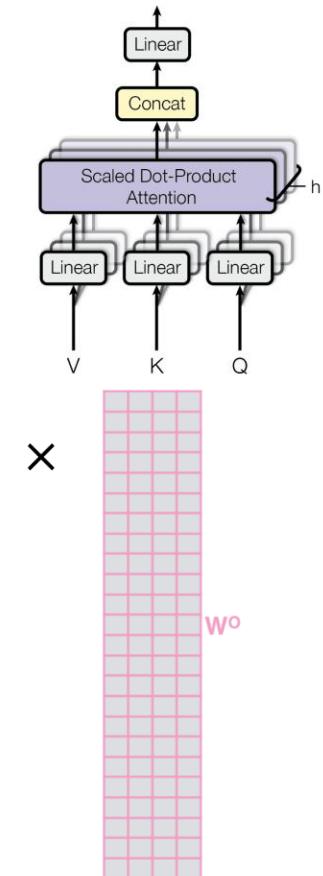
- Multiple parallel attention layers is quite common.
  - Each attention layer has its own parameters.
  - Concatenate the results and run them through a linear projection.



# Multi-Headed Self-Attention

- Just concatenate all the heads and apply an output projection matrix.

$$\begin{aligned}\text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{head}_i &= \text{Attention}(XW_i^Q, XW_i^K, XW_i^V)\end{aligned}$$



- In practice, we use a *reduced dimension* for each head.

$$W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$$

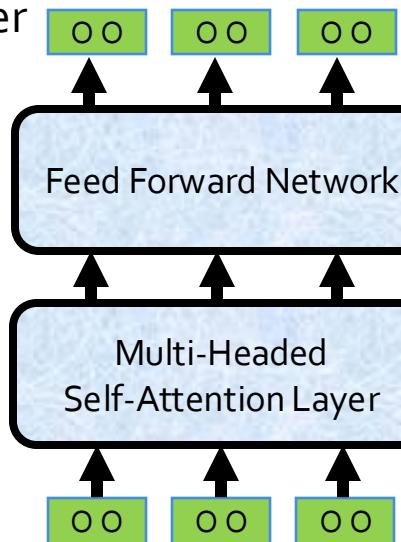
$$d_q = d_k = d_v = d/m \quad d = \text{hidden size, } m = \# \text{ of heads}$$

$$W^O \in \mathbb{R}^{d \times d_2} \quad \text{If we stack multiple layers, usually } d_1 = d_2 = d$$

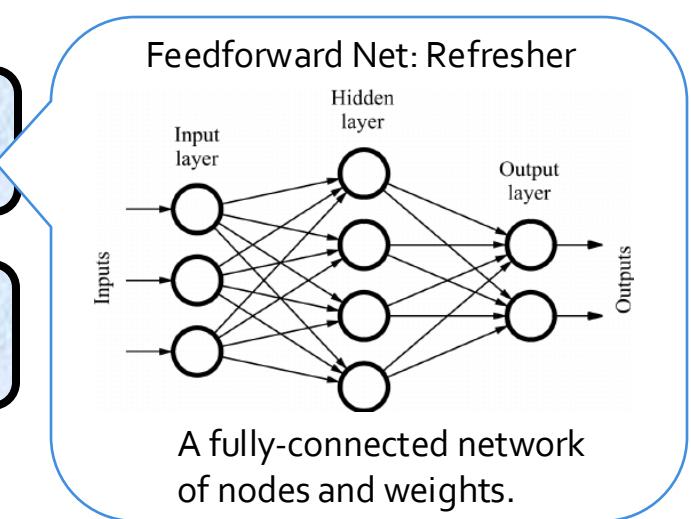
- The total computational cost is similar to that of single-head attention with full dimensionality.

# Combine with FFN

- Add a **feed-forward network** on top it to add more expressivity.
  - This allows the model to apply another transformation to the contextual representations (or “post-process” them).
  - Usually, the dimensionality of the hidden feedforward layer is 2-8 times larger than the input dimension.

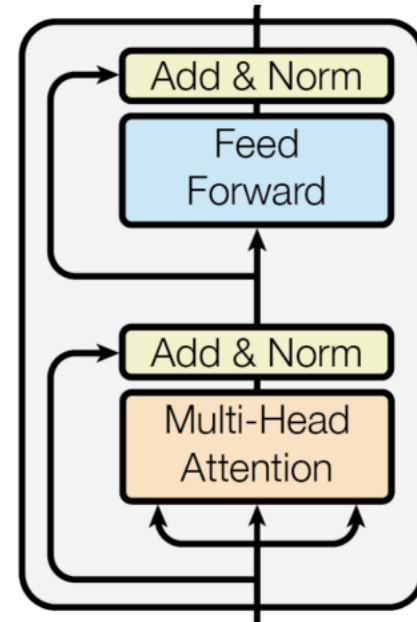


$$\text{FFN}(\mathbf{x}) = f(cW_1 + b_1)W_2 + b_2$$



# How Do We Prevent Vanishing Gradients?

- Residual connections let the model “skip” layers
  - These connections are particularly useful for training deep networks
- Use layer normalization to stabilize the network and allow for proper gradient flow

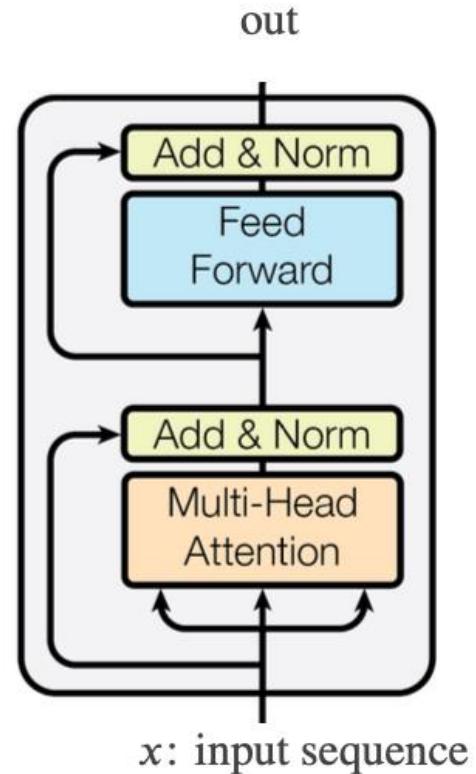


# Putting it Together: Self-Attention Block

Given input  $\mathbf{x}$ :

$$\begin{aligned}\text{out} &= LN(\tilde{\mathbf{c}} + \mathbf{c}') \\ \tilde{\mathbf{c}} &= \text{FFN}(\mathbf{c}') = f(\mathbf{c}'W_1 + b_1)W_2 + b_2\end{aligned}$$

$$\begin{aligned}\mathbf{c}' &= LN(\mathbf{c} + \mathbf{x}) \\ \mathbf{c} &= \text{MultiHeadedAttention}(\mathbf{x}; \mathbf{W}^q, \mathbf{W}^k, \mathbf{W}^v)\end{aligned}$$



# Summary: Self-Attention Block

---

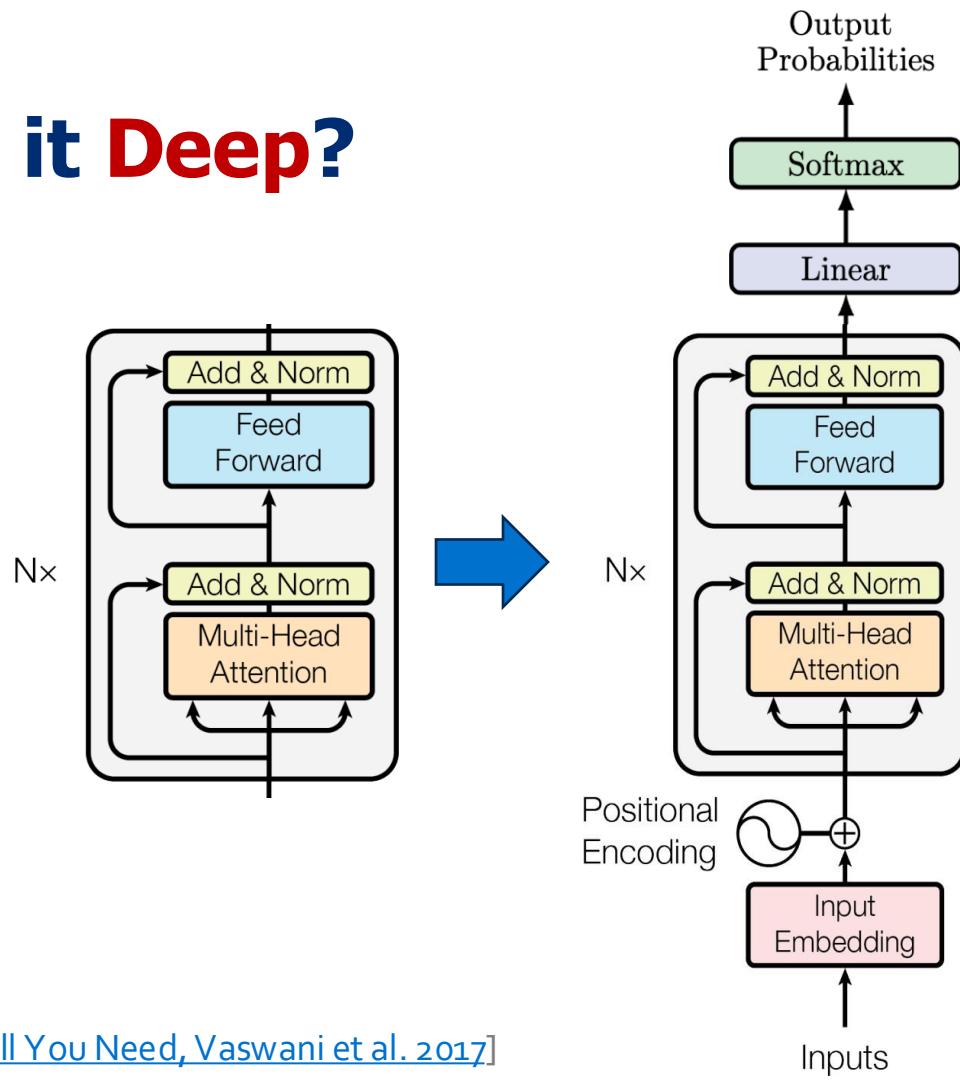
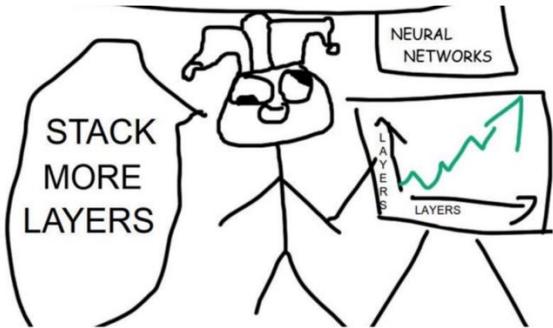
- **Self-Attention:** A critical building block of modern language models.
  - The idea is to compose meanings of words weighted according some similarity notion.
- **Next:** We will combine self-attention blocks to build various architectures known as Transformer.



# Transformer

# How Do We Make it Deep?

- Stack more layers!



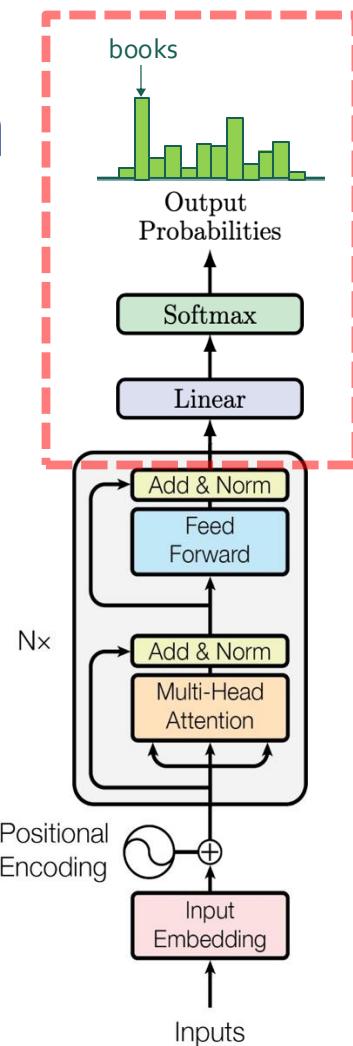
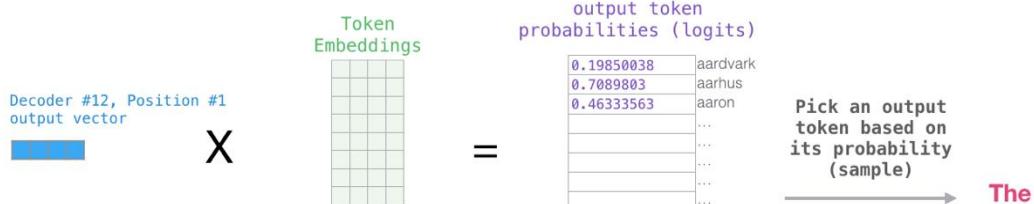
# From Representations to Prediction

- To perform prediction, add a classification head on top of the final layer of the transformer.
- This can be per token (Language modeling)
- Or can be for the entire sequence (only one token)

$\text{out} \in \mathbb{R}^{S \times d}$  ( $S$ : Sequence length)

$\text{logits} = \text{Linear}_{(d, V)}(\text{out}) = f(\text{out} \cdot W_V) \in \mathbb{R}^{S \times V}$

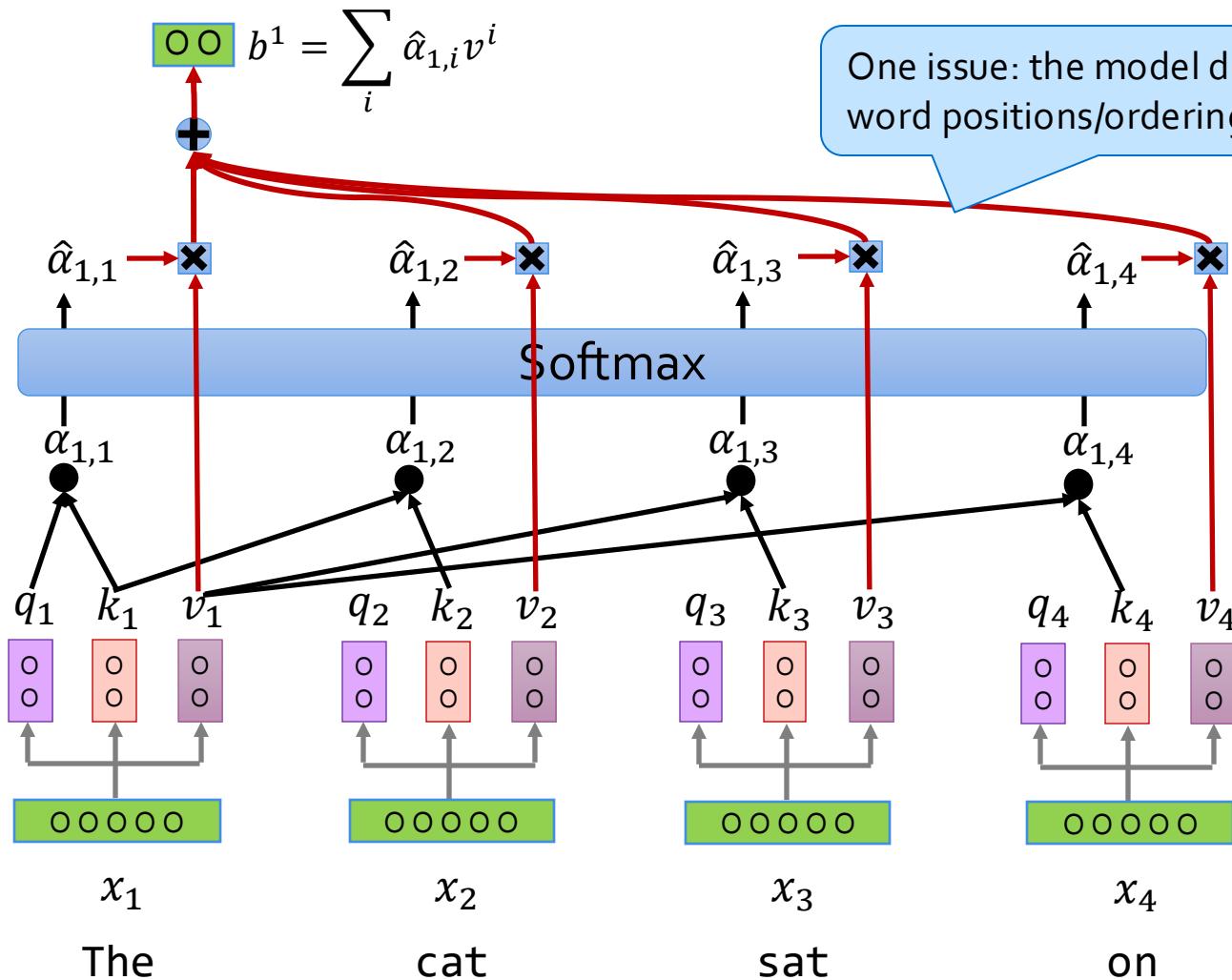
$\text{probabilities} = \text{softmax}(\text{logits}) \in \mathbb{R}^{S \times V}$





One last wrinkle though ...

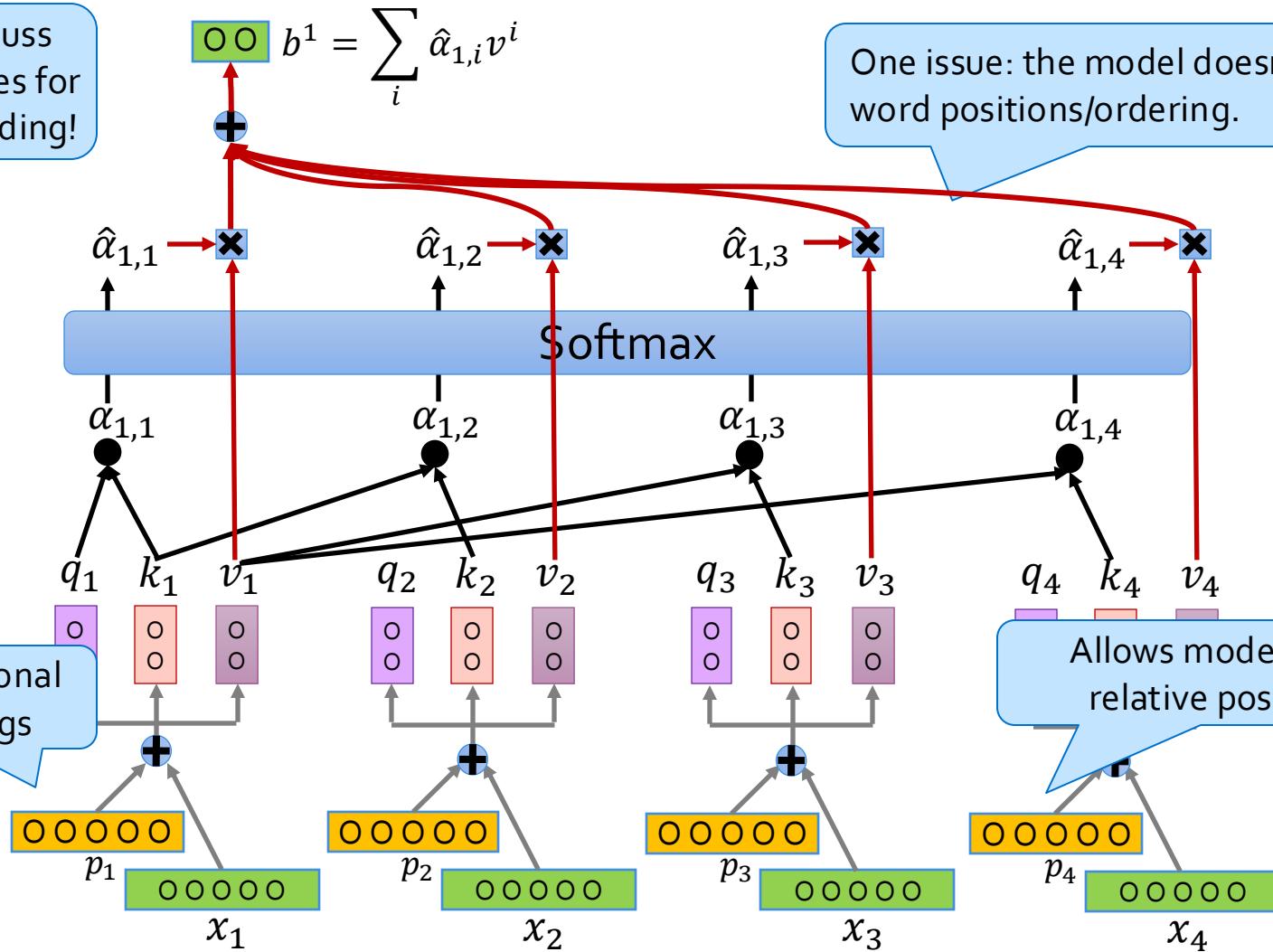




We will discuss various choices for these embedding!

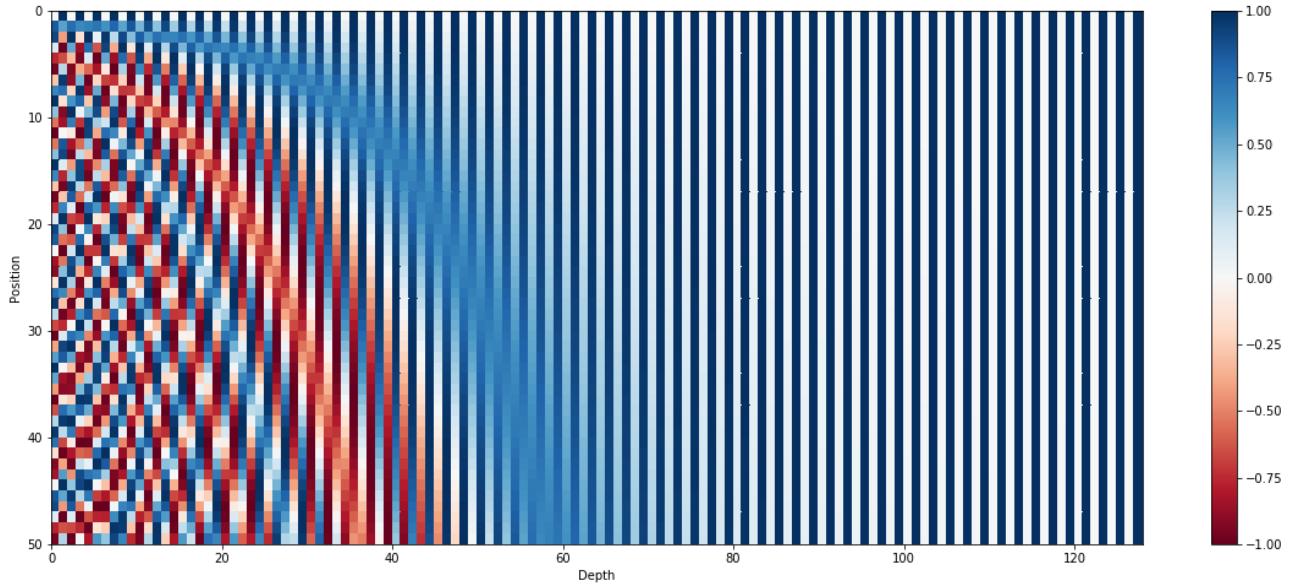
$$OO = \sum_i \hat{\alpha}_{1,i} v^i$$

One issue: the model doesn't know word positions/ordering.

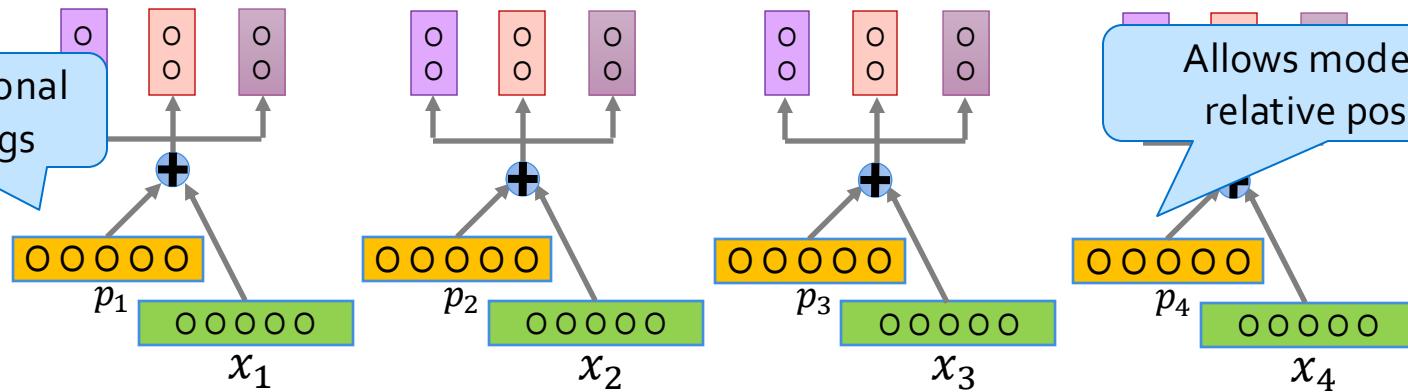


An approach:  
Sine/Cosine encoding

$$p_i = \begin{cases} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{cases}$$



$p_i$  are positional embeddings



# The Transformer Stack in PyTorch

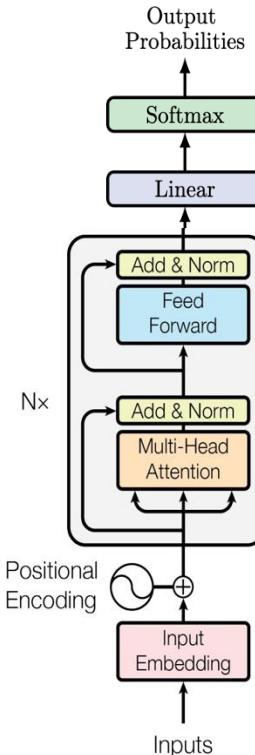


```
class Block(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.ln_1 = LayerNorm(config.n_embd, bias=config.bias)
        self.attn = CausalSelfAttention(config)
        self.ln_2 = LayerNorm(config.n_embd, bias=config.bias)
        self.mlp = MLP(config)

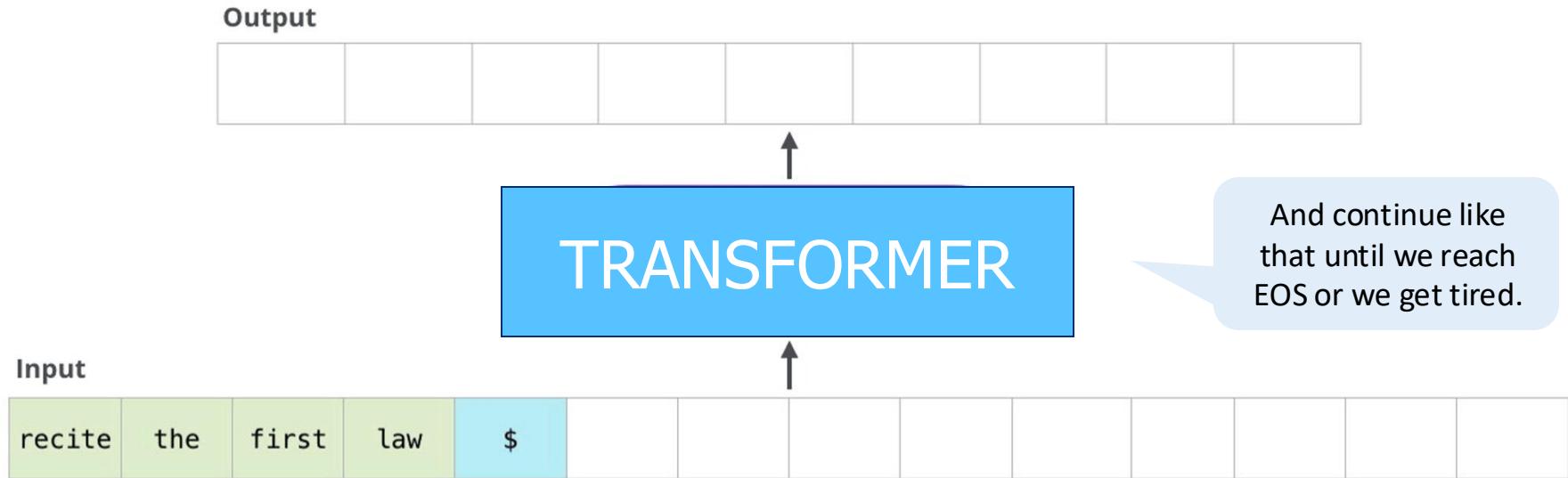
    def forward(self, x):
        x = x + self.attn(self.ln_1(x))
        x = x + self.mlp(self.ln_2(x))
        return x

    self.transformer = nn.ModuleDict(
        dict(
            wte=nn.Embedding(config.vocab_size, config.n_embd),
            wpe=nn.Embedding(config.block_size, config.n_embd),
            drop=nn.Dropout(config.dropout),
            h=nn.ModuleList([Block(config) for _ in range(config.n_layer)]),
            ln_f=LayerNorm(config.n_embd, bias=config.bias),
        )
    )

    self.lm_head = nn.Linear(config.n_embd, config.vocab_size, bias=False)
```



# Transformer-based Language Modeling



# Training a Transformer Language Model

- **Goal:** Train a Transformer for language modeling (i.e., predicting the next word).
- **Approach:** Train it so that each position is predictor of the next (right) token.
  - We just shift the input to right by one, and use as labels

(gold output)  $Y = \text{cat sat on the mat } </s>$

EOS special token



```
X = text[:, :-1]  
Y = text[:, 1:]
```

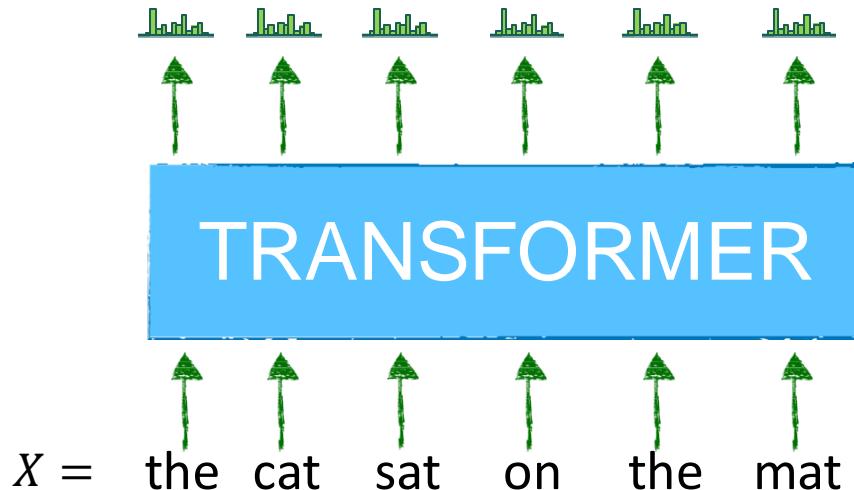
$X = \text{the cat sat on the mat}$

[Slide credit: Arman Cohan]

# Training a Transformer Language Model

- For each position, compute their corresponding **distribution** over the whole vocab.

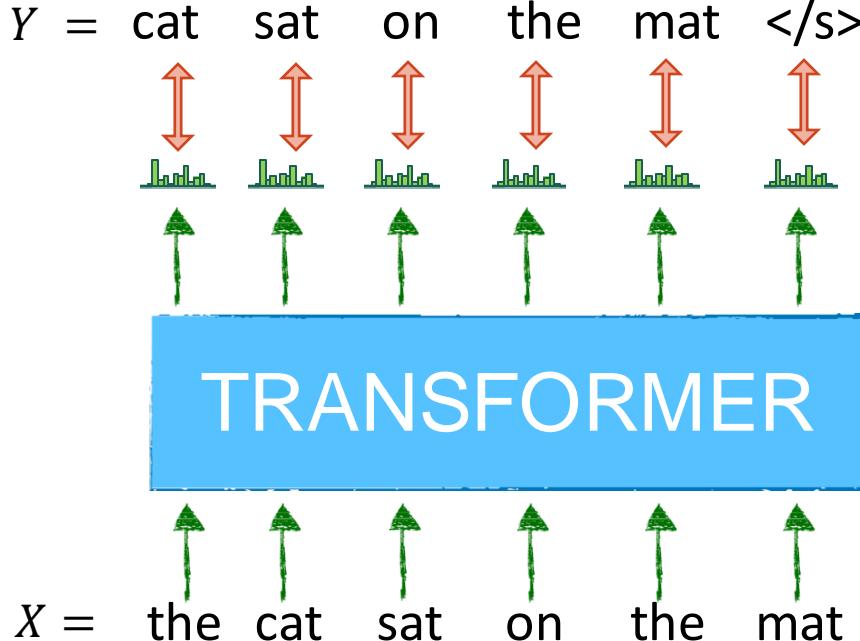
(gold output)  $Y = \text{cat sat on the mat } </s>$



# Training a Transformer Language Model

- For each position, compute the **loss** between the distribution and the gold output label.

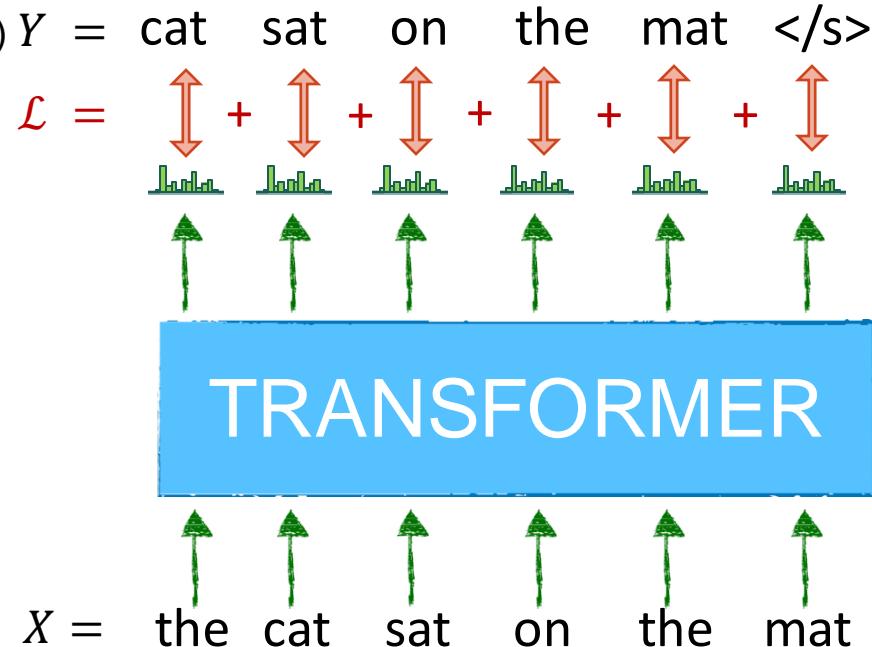
(gold output)  $Y = \text{cat sat on the mat } </s>$



# Training a Transformer Language Model

- Sum the position-wise loss values to obtain a **global loss**.

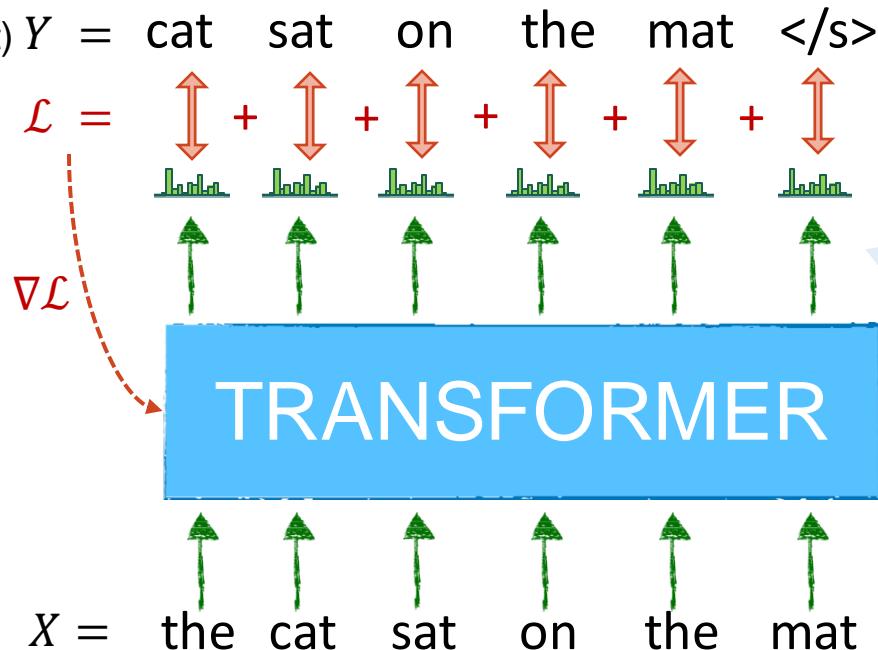
(gold output)  $Y = \text{cat sat on the mat } </s>$



# Training a Transformer Language Model

- Using this loss, do **Backprop** and **update** the Transformer parameters.

(gold output)  $Y = \text{cat sat on the mat } </s>$

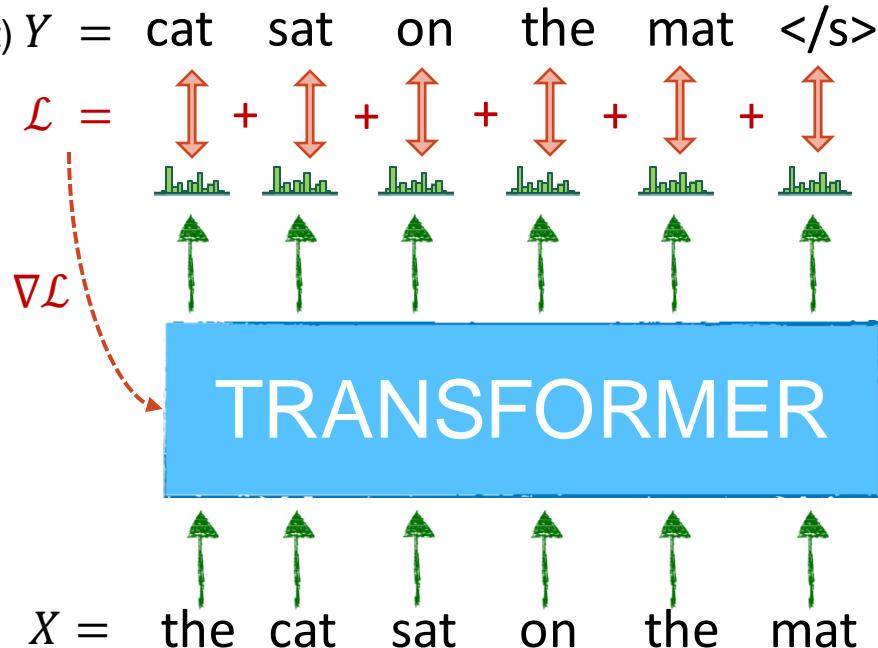


Well, this is not quite right 😊  
...  
what is the problem with this?

# Training a Transformer Language Model

- The model would solve the task by **copying** the next token to output (data leakage).
  - Does **not** learn anything useful

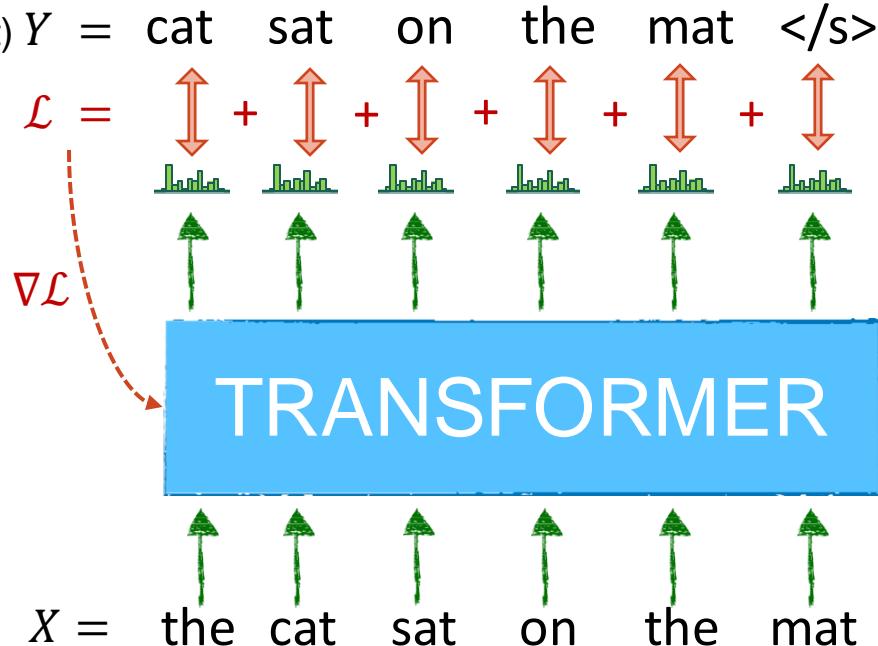
(gold output)  $Y = \text{cat sat on the mat } </s>$



# Training a Transformer Language Model

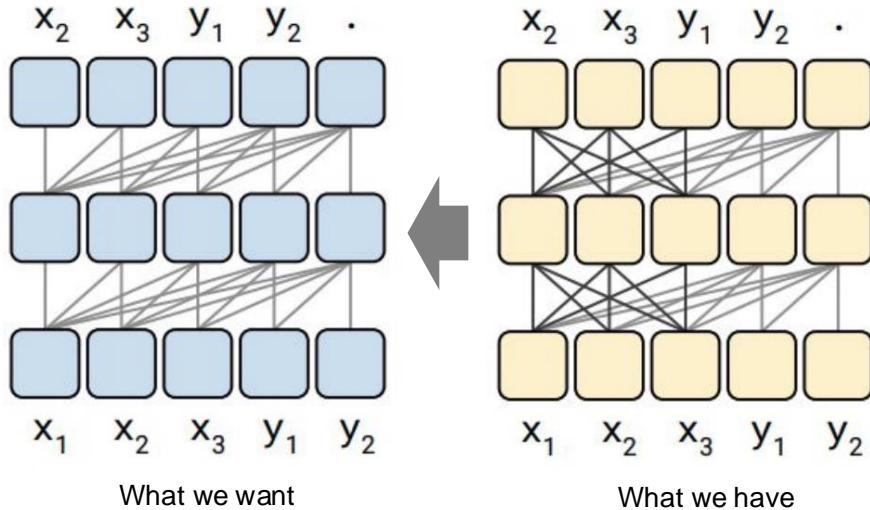
- We need to **prevent information leakage** from future tokens! How?

(gold output)  $Y = \text{cat sat on the mat } </s>$



# Attention mask

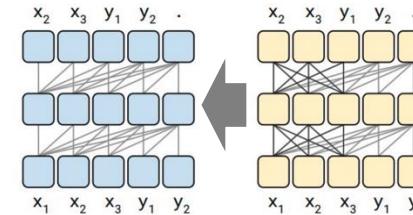
Attention raw scores												
0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11



# Attention mask

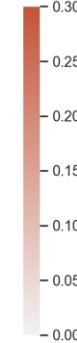
Attention raw scores

	1	2	3	4	5	6	7	8	9	10	11	12
0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11



Attention mask

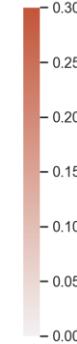
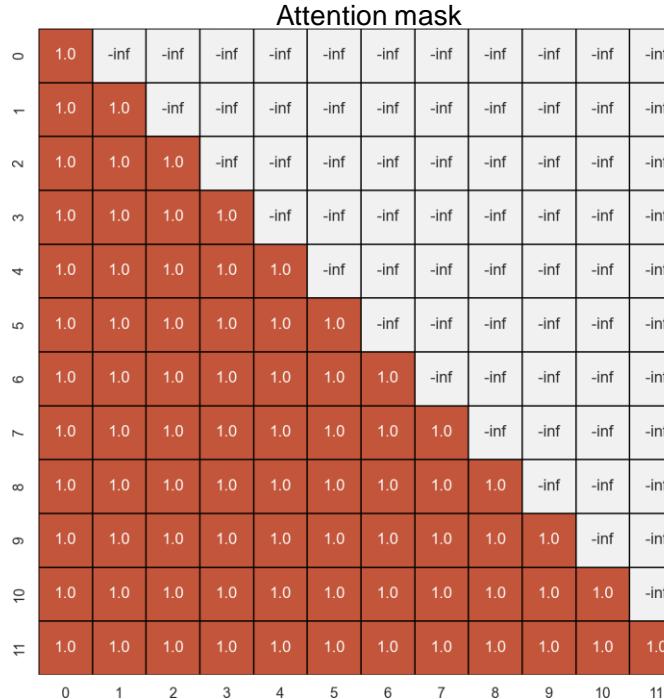
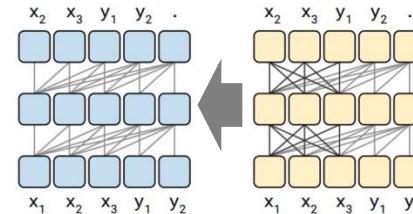
	0	1	2	3	4	5	6	7	8	9	10	11
0	1.0	-inf										
1	1.0	1.0	-inf									
2	1.0	1.0	1.0	-inf								
3	1.0	1.0	1.0	1.0	-inf							
4	1.0	1.0	1.0	1.0	1.0	-inf						
5	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf
6	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf
7	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf
8	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf
9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf
10	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf
11	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0



# Attention mask

Attention raw scores												
0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11

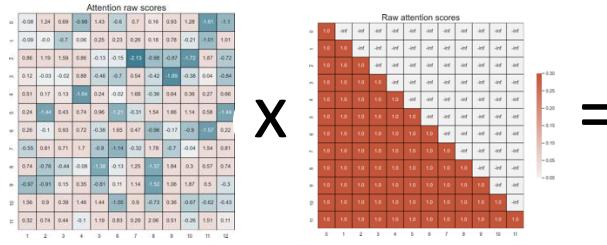
X



Note matrix multiplication is quite fast in GPUs.

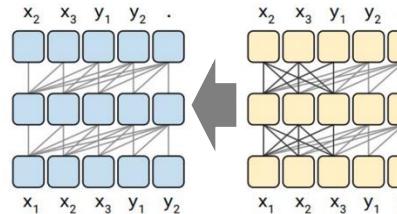
Arman Cohan

# Attention mask



### Masked attention raw scores

	0	-0.08	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
1	-0.09	-0.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
2	0.86	1.19	1.59	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
3	0.12	-0.03	-0.02	0.88	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
4	0.51	0.17	0.13	-1.64	0.24	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-inf	-inf	-inf	-inf	-inf	-inf	-inf
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-inf	-inf	-inf	-inf	-inf	-inf
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-inf	-inf	-inf	-inf	-inf
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	-inf	-inf	-inf	-inf
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	-inf	-inf	-inf
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-inf	-inf
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11	-inf
12	1	2	3	4	5	6	7	8	9	10	11	12	-inf

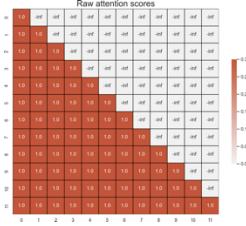


# Attention mask

The effect is more than just pruning out some of the wirings in self-attention block.

Attention raw scores											
-0.09	1.24	0.69	0.38	3.43	-0.1	0.7	1.16	0.93	1.28	1.81	-1.13
-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	1.01	1.01
-0.06	1.19	1.59	0.98	-0.15	0.15	1.11	0.96	-0.97	1.22	1.67	-0.72
-0.12	-0.03	-0.42	0.88	-0.46	-0.1	0.54	-0.42	1.86	-0.38	0.94	-0.84
-0.51	0.17	0.55	0.34	0.24	0.02	1.68	0.36	0.04	0.37	0.66	-0.01
-0.24	-1.64	0.43	0.78	0.96	1.2	-0.31	1.54	1.88	1.14	0.58	-1.44
-0.26	-0.1	0.63	0.72	0.36	1.84	0.47	1.96	-0.17	-0.5	0.67	0.22
-0.55	0.81	0.71	1.7	-0.18	0.32	1.78	-0.7	-0.94	1.54	0.81	-0.74
-0.74	-0.48	-0.48	0.08	-0.33	-0.13	1.25	1.37	1.84	0.3	0.57	0.74
-0.87	0.91	0.15	0.30	-0.81	0.11	1.14	-0.50	1.08	1.87	0.5	-0.37
-0.52	0.74	0.44	0.1	1.19	0.89	0.29	2.06	0.31	-0.28	1.51	0.11

X



Masked attention raw scores											
-0.08	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
-0.09	-0.0	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
-0.06	1.19	1.59	inf	inf	inf	inf	inf	inf	inf	inf	inf
-0.12	-0.03	0.88	inf	inf	inf	inf	inf	inf	inf	inf	inf
-0.51	0.17	0.55	-1.64	0.24	inf	inf	inf	inf	inf	inf	inf
-0.24	-1.64	0.43	0.78	0.96	1.2	inf	inf	inf	inf	inf	inf
-0.26	-0.1	0.63	0.72	0.36	1.84	0.47	inf	inf	inf	inf	inf
-0.55	0.81	0.71	1.7	-0.18	0.32	1.78	-0.7	inf	inf	inf	inf
-0.74	-0.48	-0.48	0.08	-0.33	-0.13	1.25	1.37	1.84	inf	inf	inf
-0.87	0.91	0.15	0.30	-0.81	0.11	1.14	-0.52	1.06	1.87	inf	inf
-0.52	0.74	0.44	0.1	1.19	0.89	0.29	2.06	0.31	-0.28	1.51	0.11

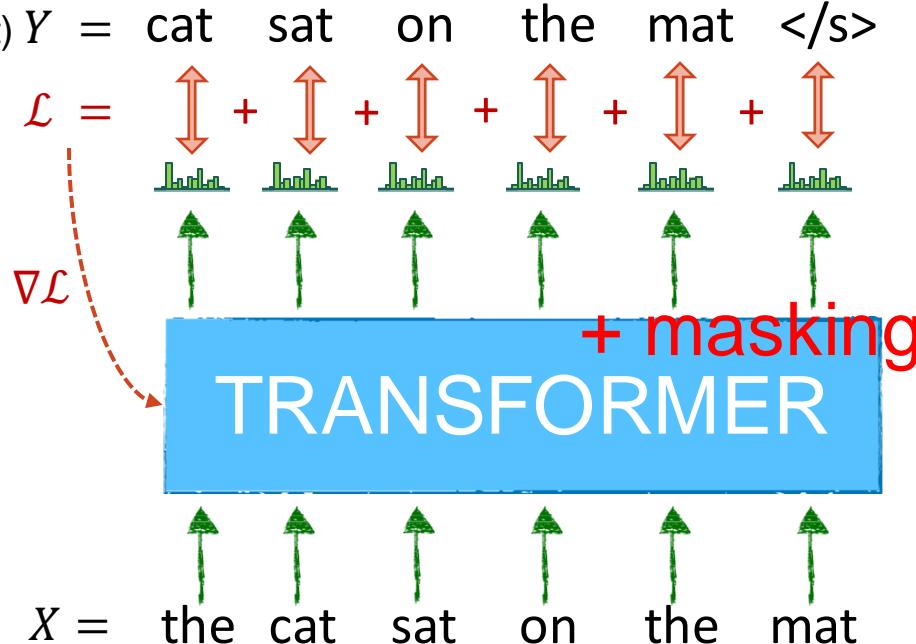
softmax



# Training a Transformer Language Model

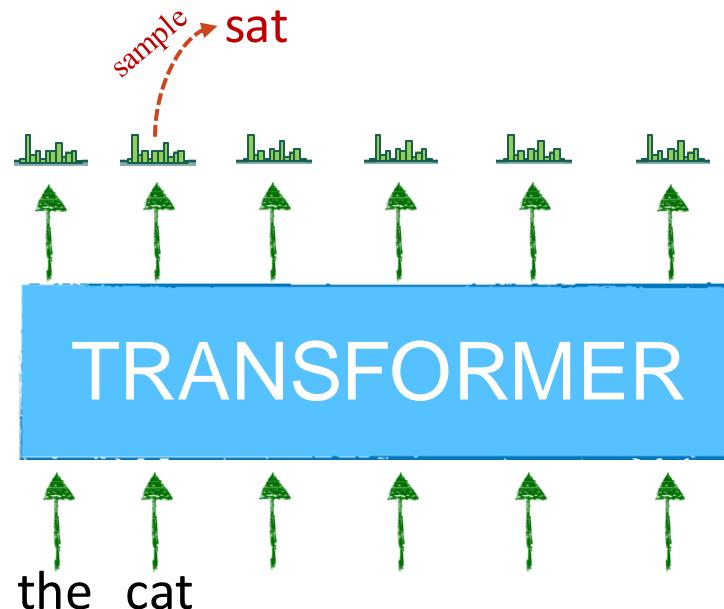
- We need to **prevent information leakage** from future tokens! How?

(gold output)  $Y = \text{cat sat on the mat } </s>$



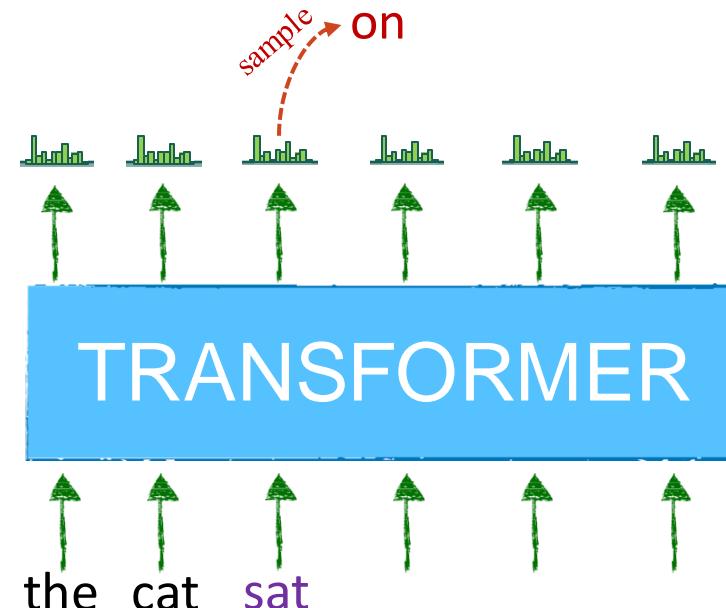
# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly



# How to use the model to generate text?

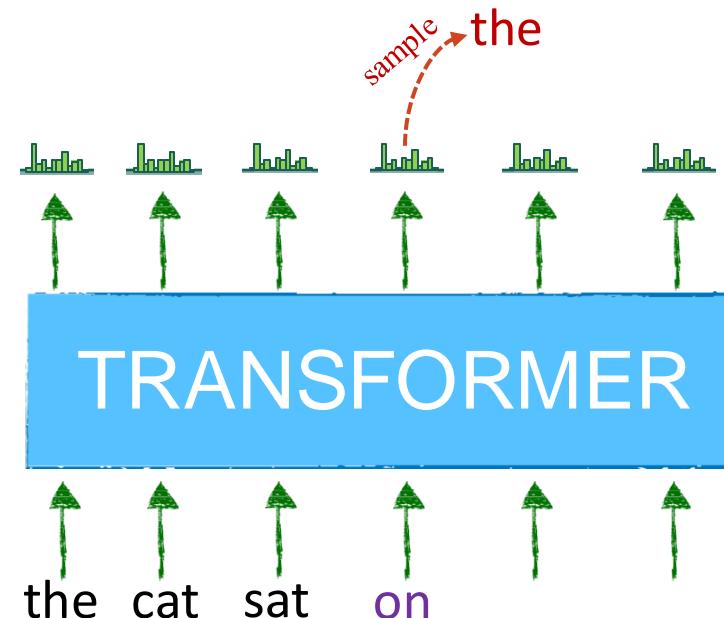
- Use the output of previous step as input to the next step repeatedly



The probabilities get revised upon adding a new token to the input.

# How to use the model to generate text?

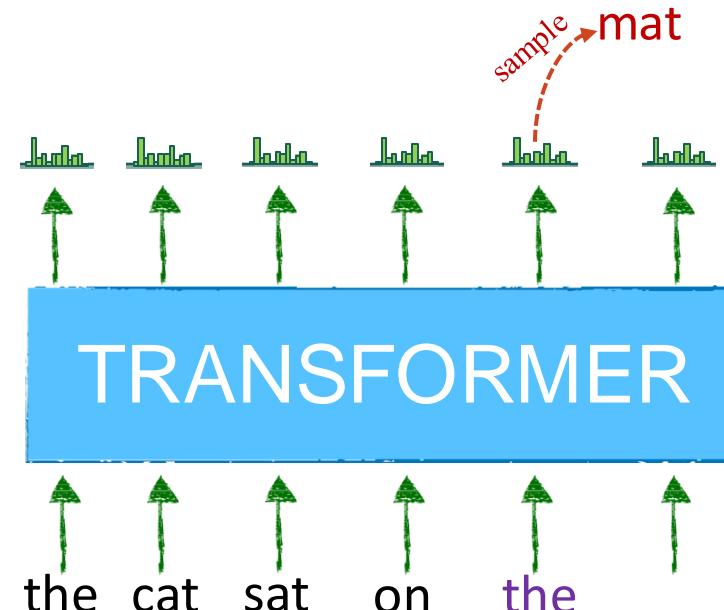
- Use the output of previous step as input to the next step repeatedly



The probabilities get revised upon adding a new token to the input.

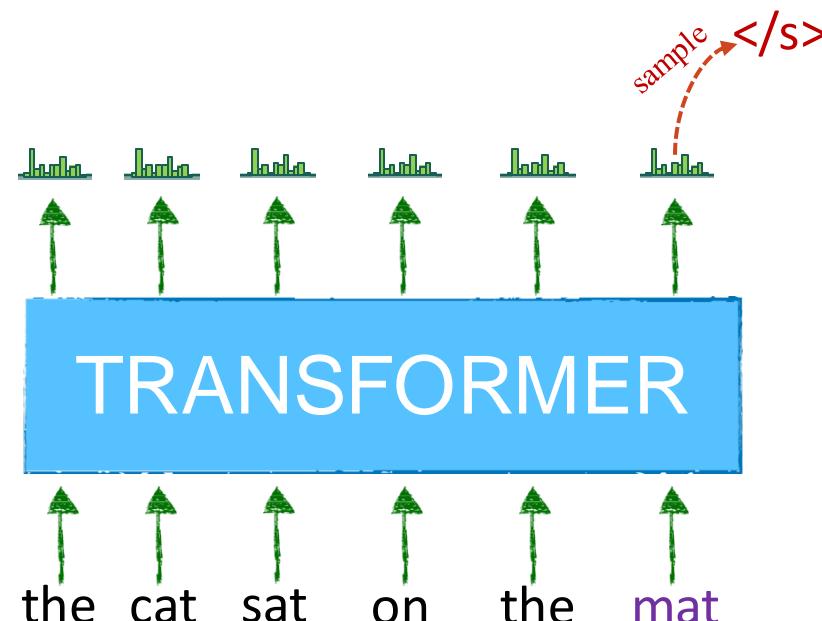
# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly



# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly



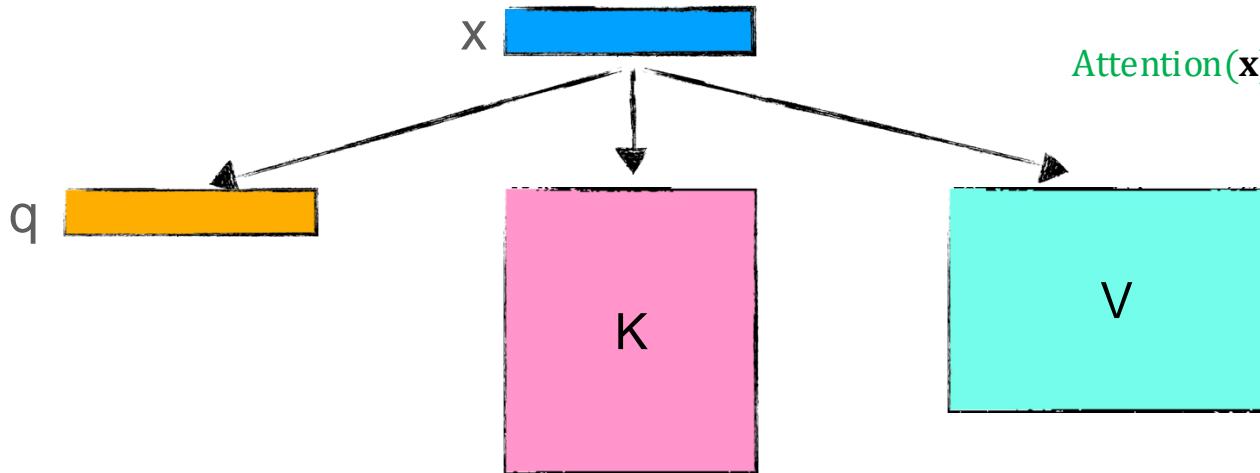
The probabilities get revised upon adding a new token to the input.

An important efficiency  
consideration about decoding!

# Making decoding more efficient

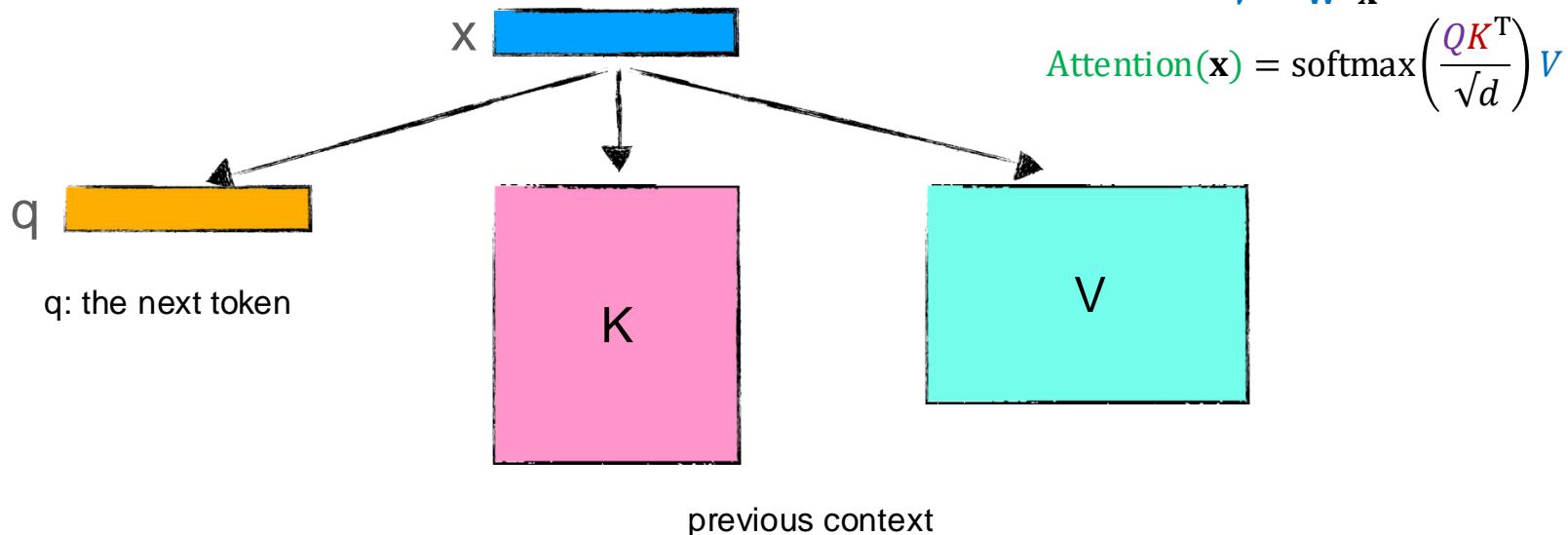
$$\begin{aligned} Q &= \mathbf{W}^q \mathbf{x} \\ K &= \mathbf{W}^k \mathbf{x} \\ V &= \mathbf{W}^v \mathbf{x} \end{aligned}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$



[Slide credit: Arman Cohan]

# Making decoding more efficient



[Slide credit: Arman Cohan]

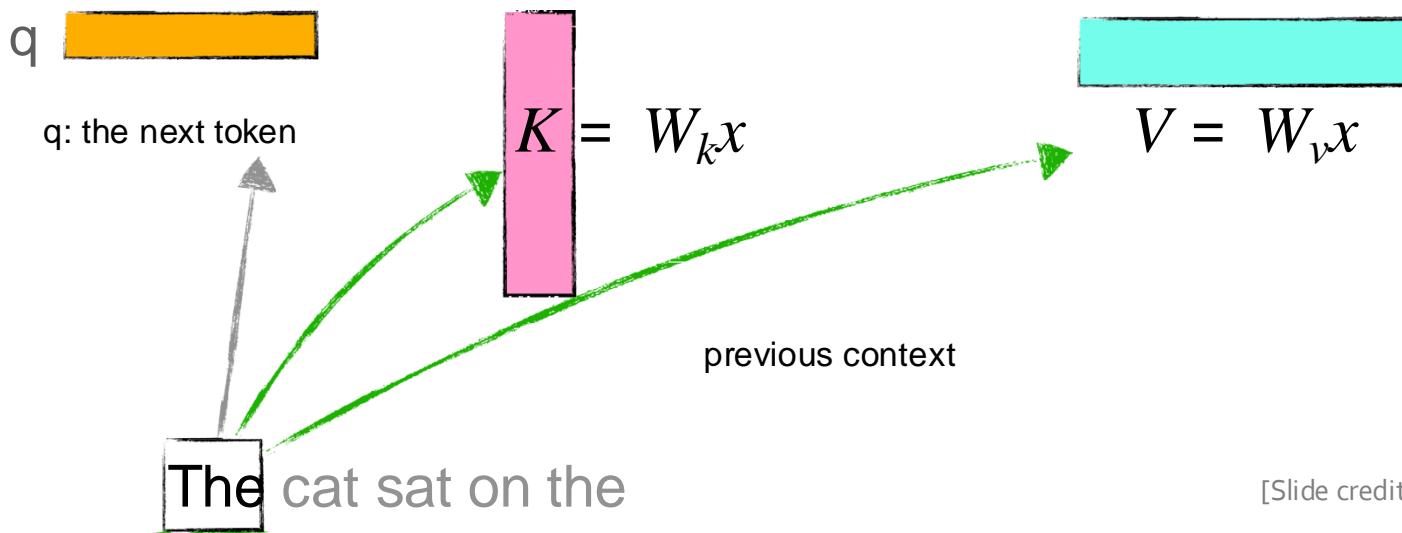
# Making decoding more efficient

$$Q = W^q \mathbf{x}$$

$$K = W^k \mathbf{x}$$

$$V = W^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

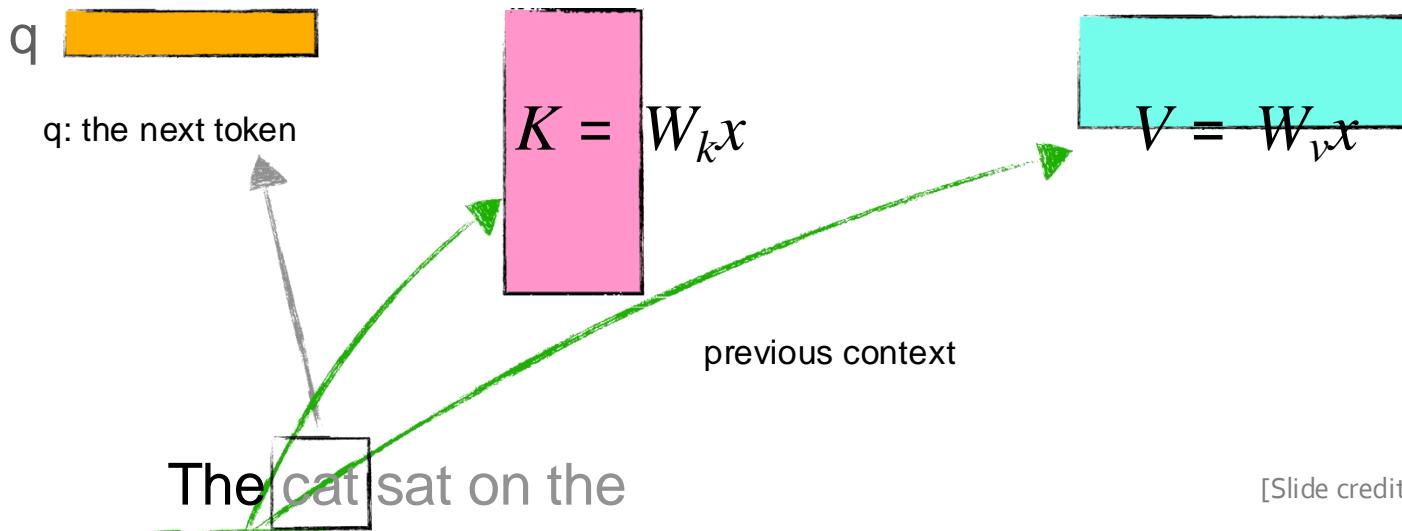
# Making decoding more efficient

$$Q = W^q x$$

$$K = W^k x$$

$$V = W^v x$$

$$\text{Attention}(x) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

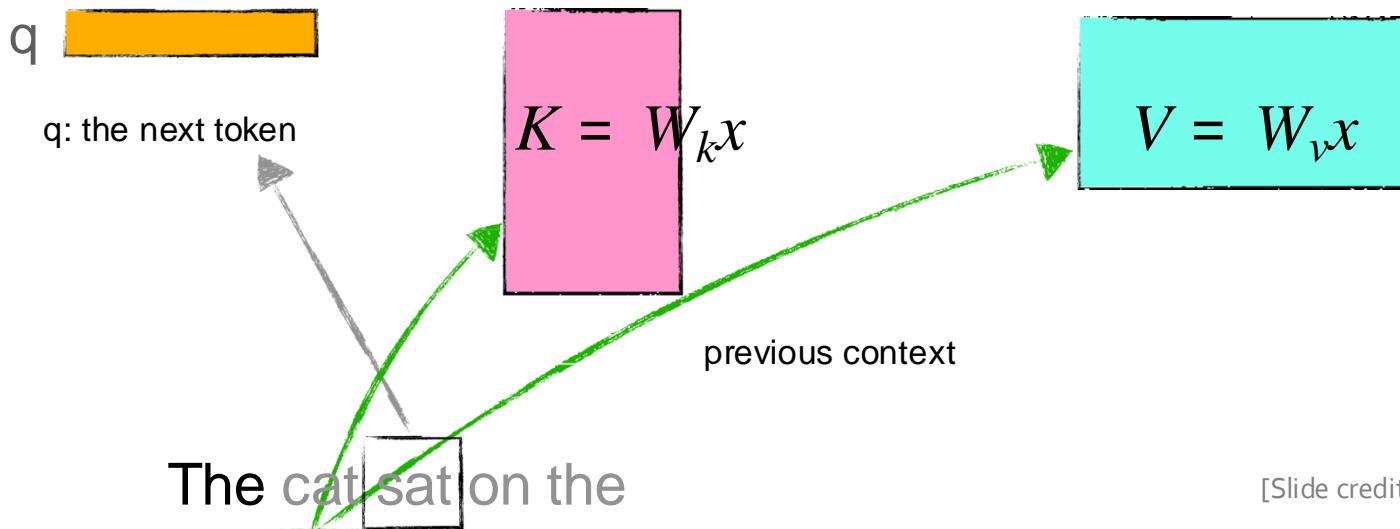
# Making decoding more efficient

$$Q = W^q x$$

$$K = W^k x$$

$$V = W^v x$$

$$\text{Attention}(x) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

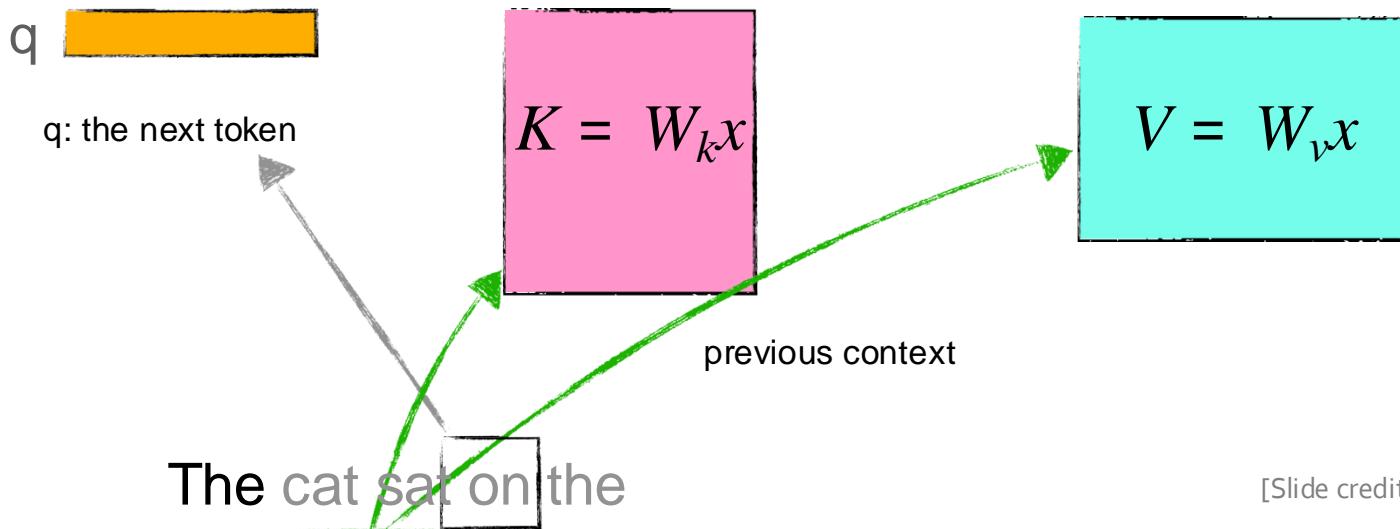
# Making decoding more efficient

$$Q = W^q \mathbf{x}$$

$$K = W^k \mathbf{x}$$

$$V = W^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

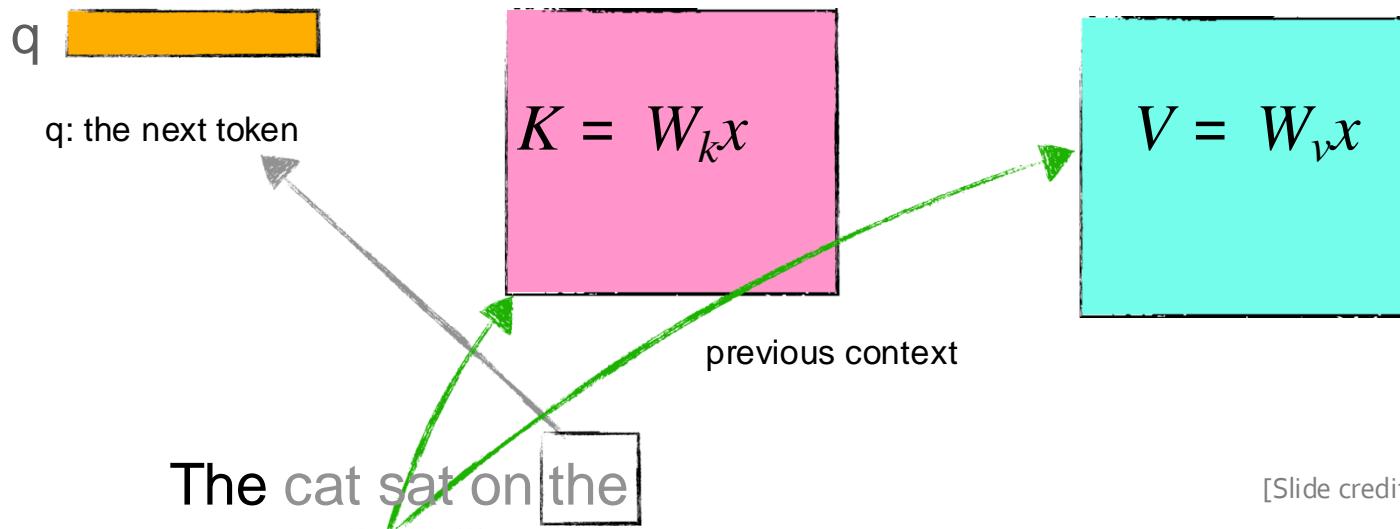
# Making decoding more efficient

$$Q = \mathbf{W}^q \mathbf{x}$$

$$K = \mathbf{W}^k \mathbf{x}$$

$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

# Making decoding more efficient

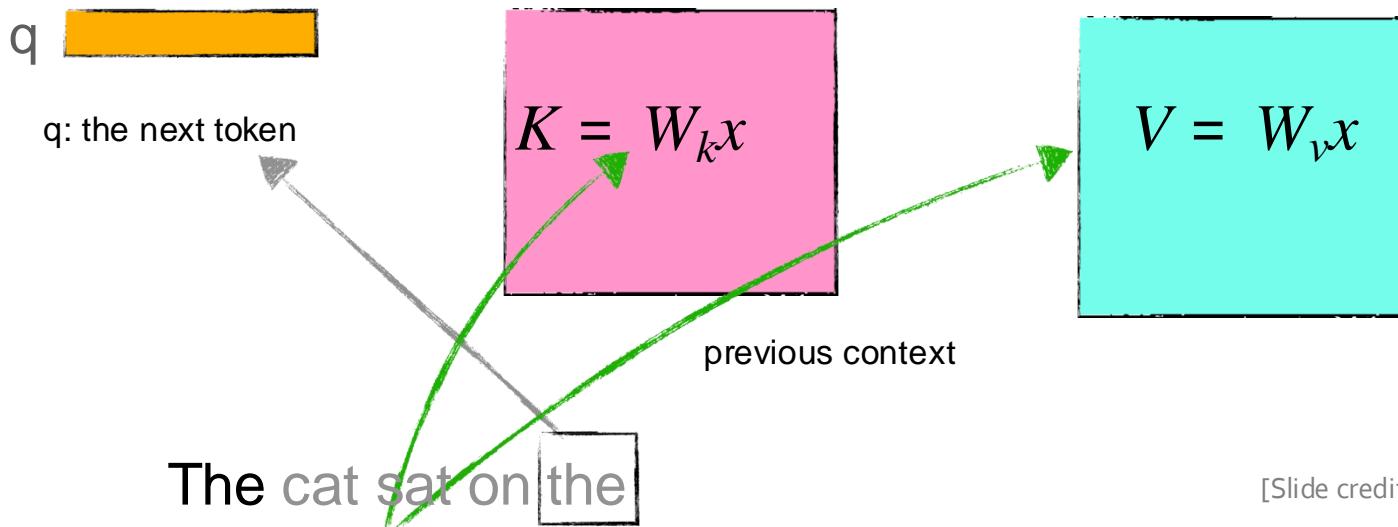
- We are computing the Keys and Values many times!
  - Let's reduce redundancy! 😊

$$Q = W^q \mathbf{x}$$

$$K = W^k \mathbf{x}$$

$$V = W^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



[Slide credit: Arman Cohan]

# Making decoding more efficient

- We are computing the Keys and Values many times!
  - Let's reduce redundancy! 😭

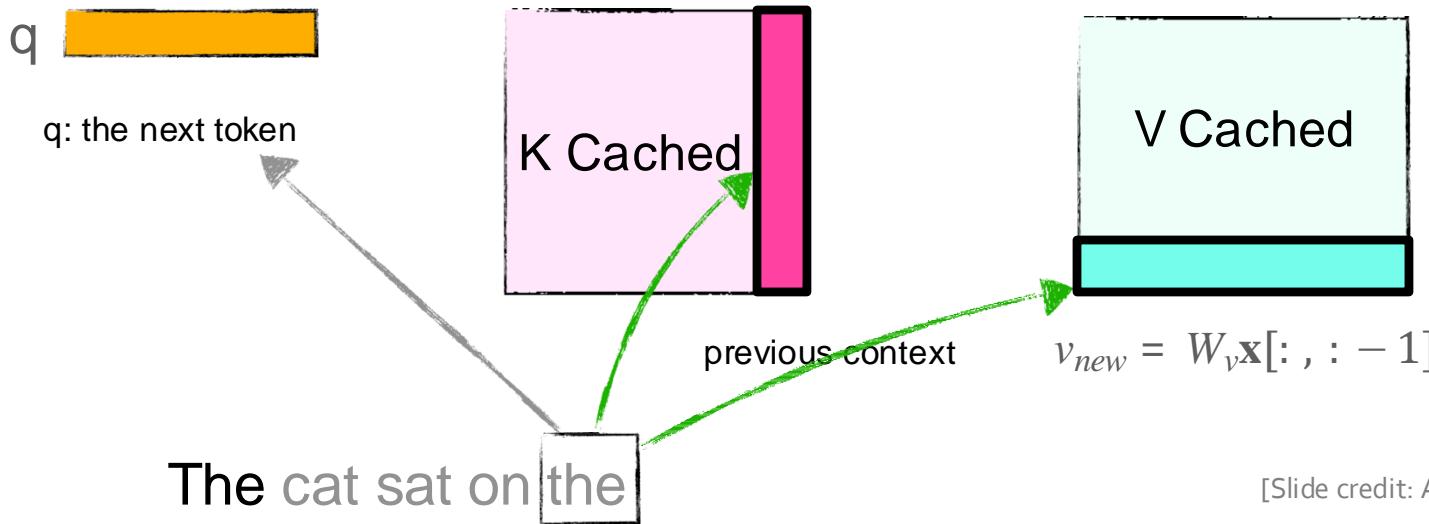
$$Q = \mathbf{W}^q \mathbf{x}$$

$$K = \mathbf{W}^k \mathbf{x}$$

$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

$$k_{new} = W_k \mathbf{x}[:, : - 1]$$



# Making decoding more efficient

- **Question:** How much memory does this K, V cache require?

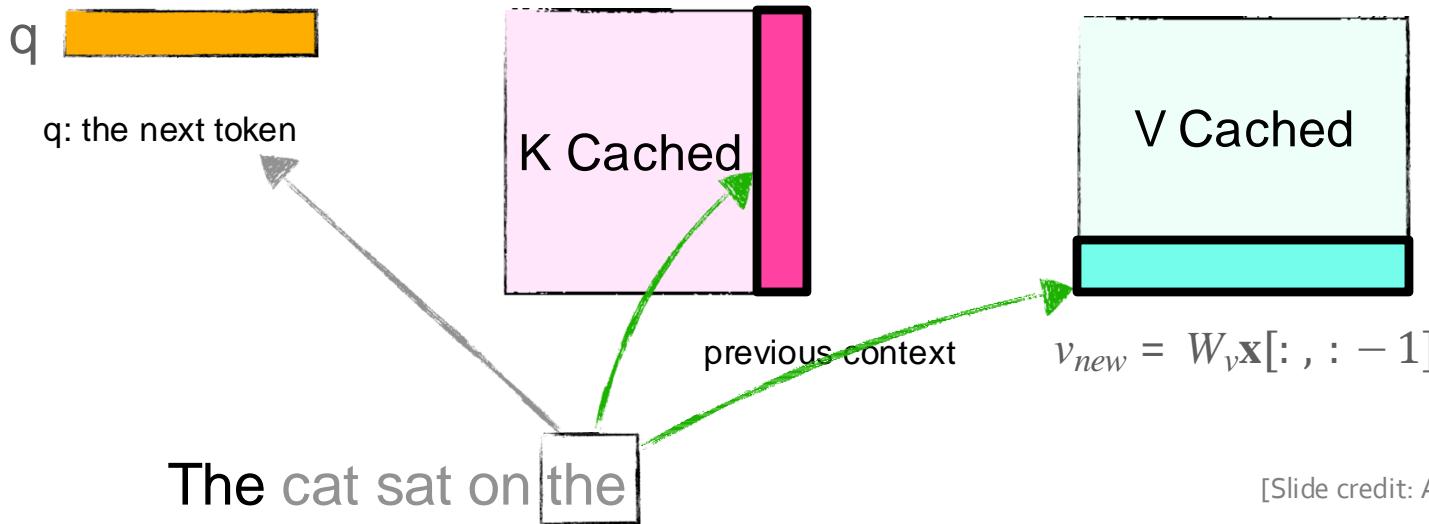
$$Q = \mathbf{W}^q \mathbf{x}$$

$$K = \mathbf{W}^k \mathbf{x}$$

$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

$$k_{new} = W_k \mathbf{x}[:, : - 1]$$



$$v_{new} = W_v \mathbf{x}[:, : - 1]$$

[Slide credit: Arman Cohan]

# Summary

---

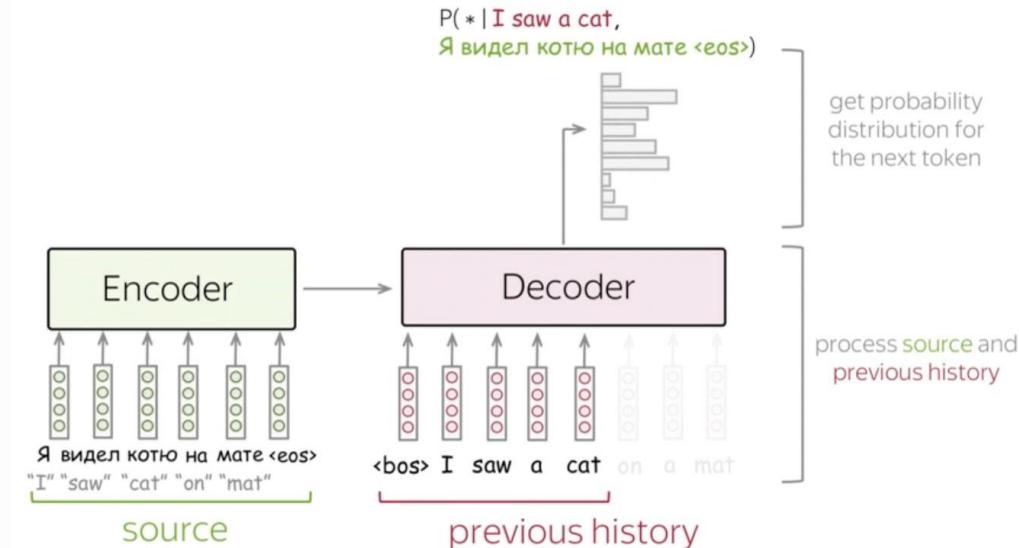
- This is a very generic Transformer!
- We will implement this in HW5 to build a simple Transformer Language Model!!
- **Next:**
  - Architectural variants
  - Efficiency issues.
  - ...



# Transformer Architectural Variants

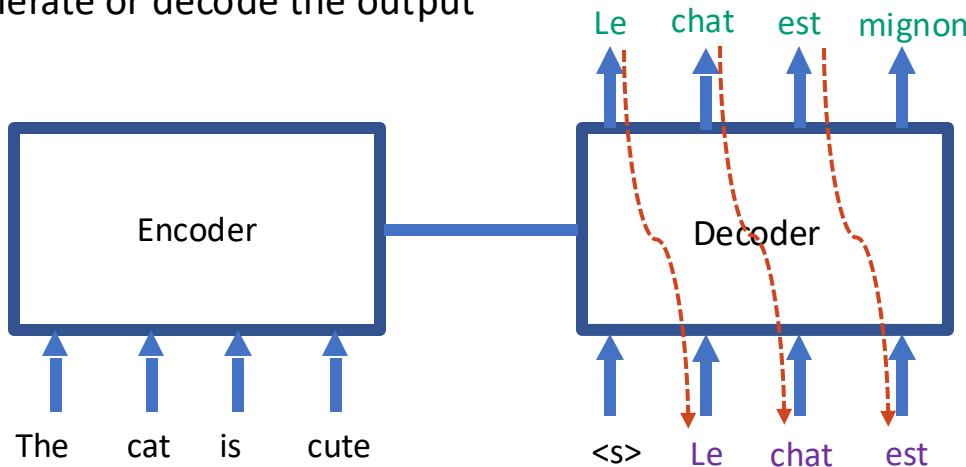
# Encoder-decoder

- It is possible to have two stacks of transformer layers
- The encoder is as we've seen
- We can also add a decoder layer that is identical to the encoder but we give it the ability to also attend to the input



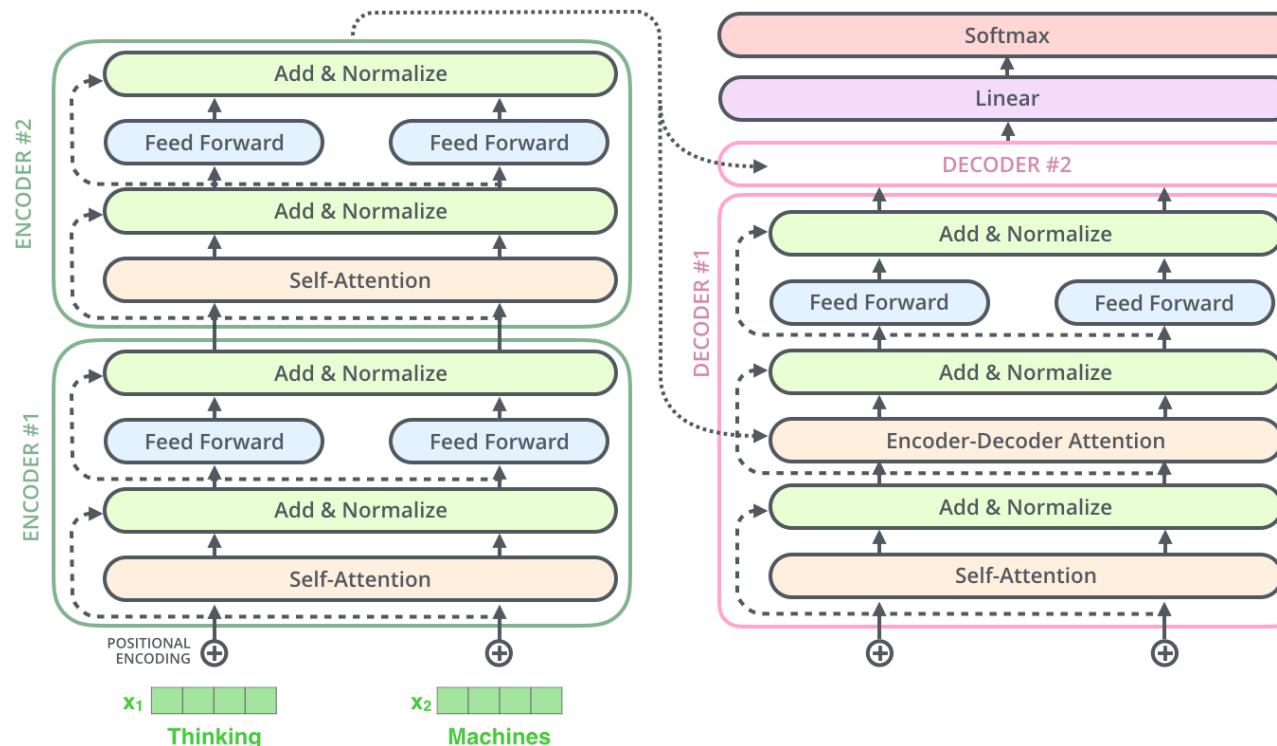
# Encoder-decoder models

- Encoder = read or encode the input,
- Decoder = generate or decode the output



# Transformer [Vaswani et al. 2017]

- An encoder-decoder architecture built with attention modules.



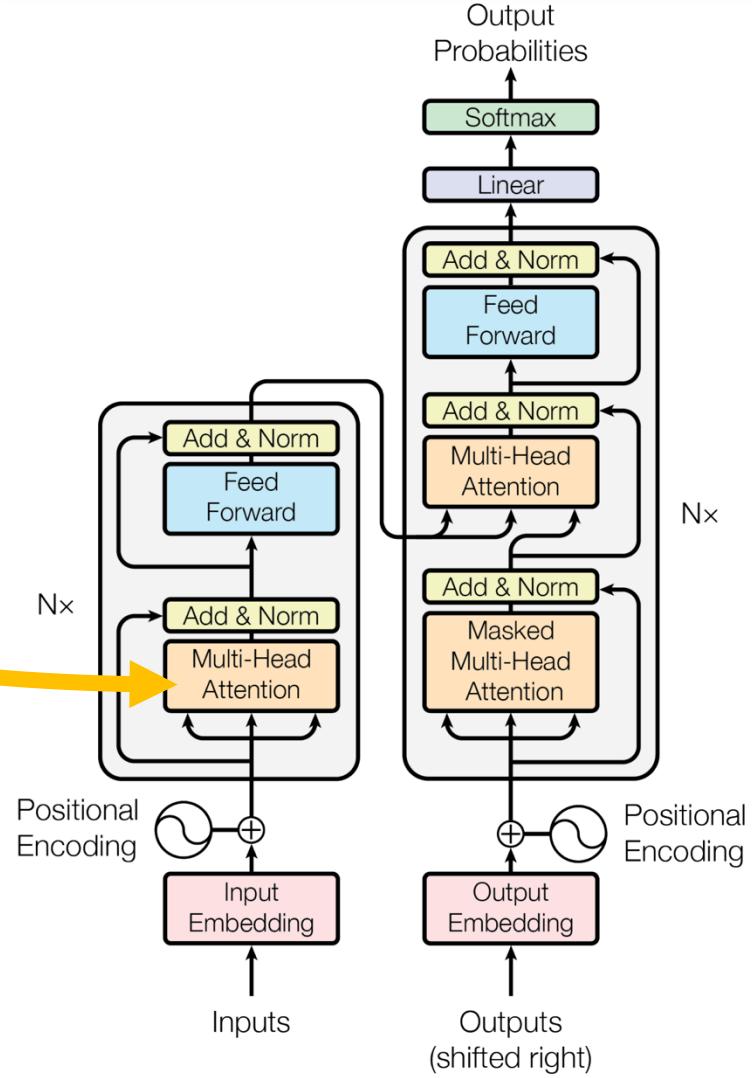
# Transformer

[Vaswani et al. 2017]

- Computation of **encoder** attends to both sides.



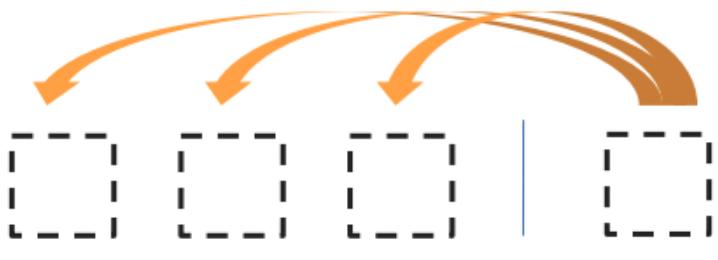
Encoder Self-Attention



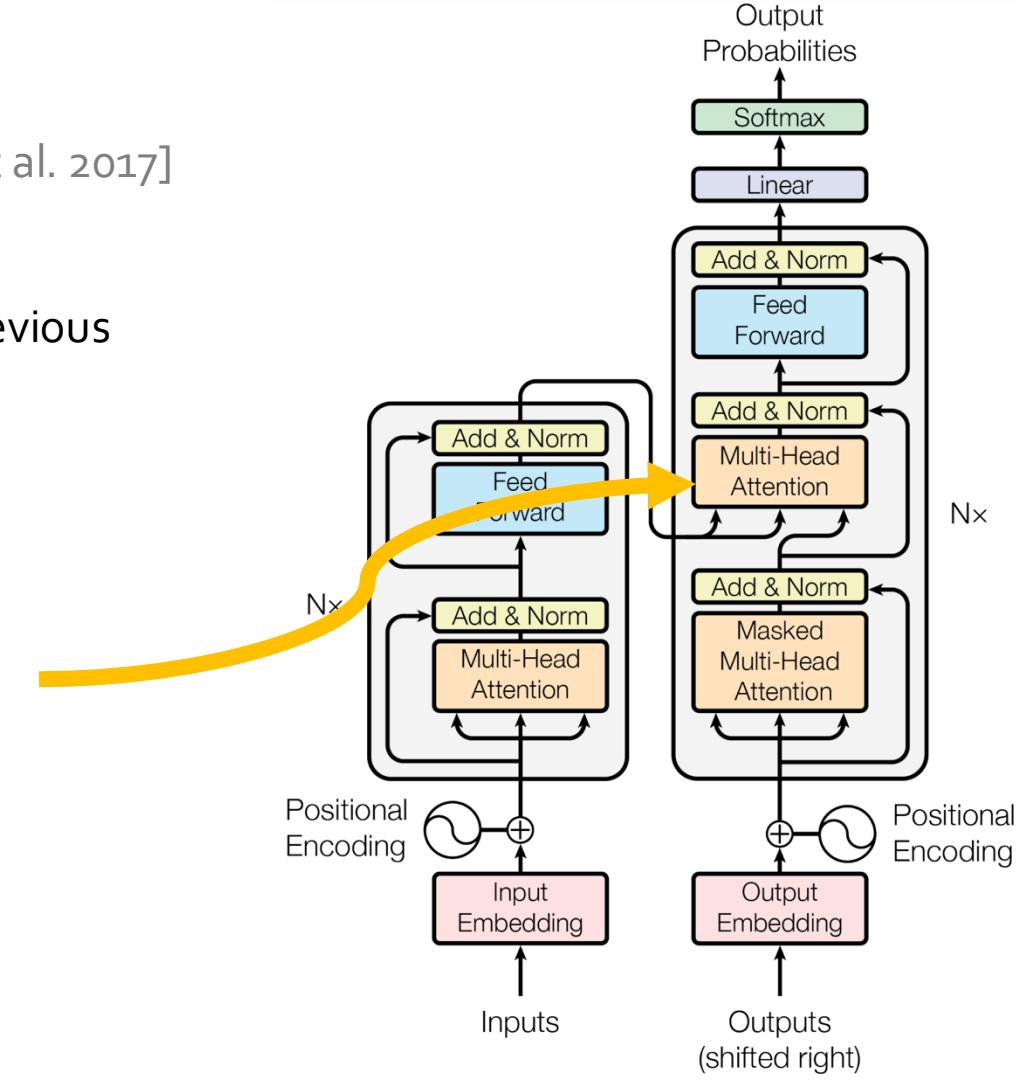
# Transformer

[Vaswani et al. 2017]

- At any step of **decoder**, it attends to previous computation of **encoder**



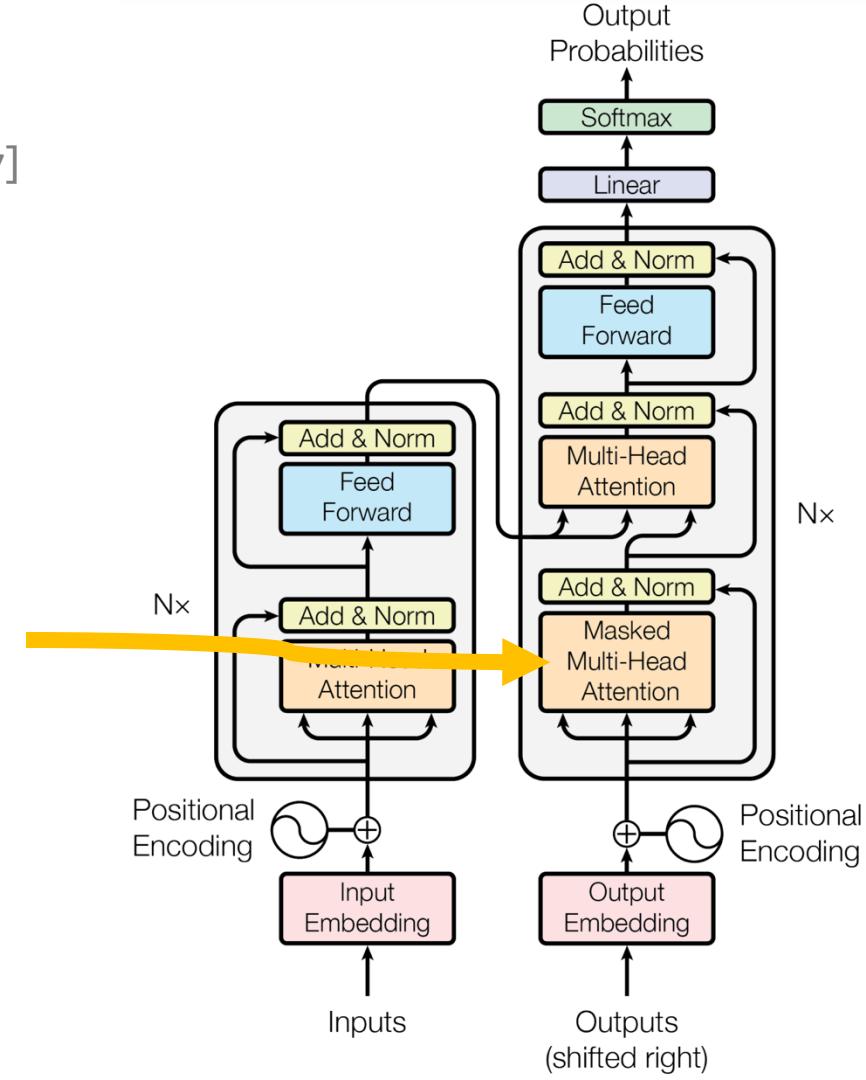
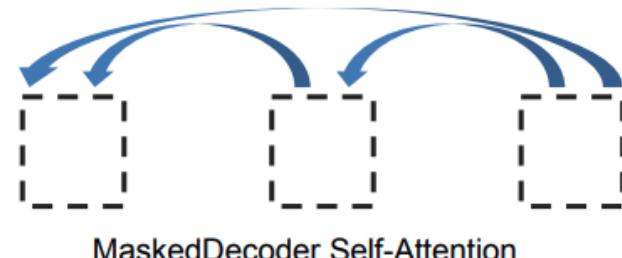
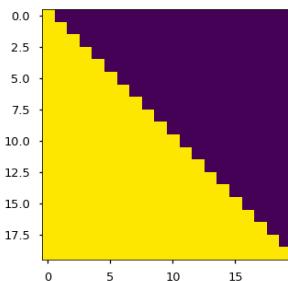
Encoder-Decoder Attention



# Transformer

 [Vaswani et al. 2017]

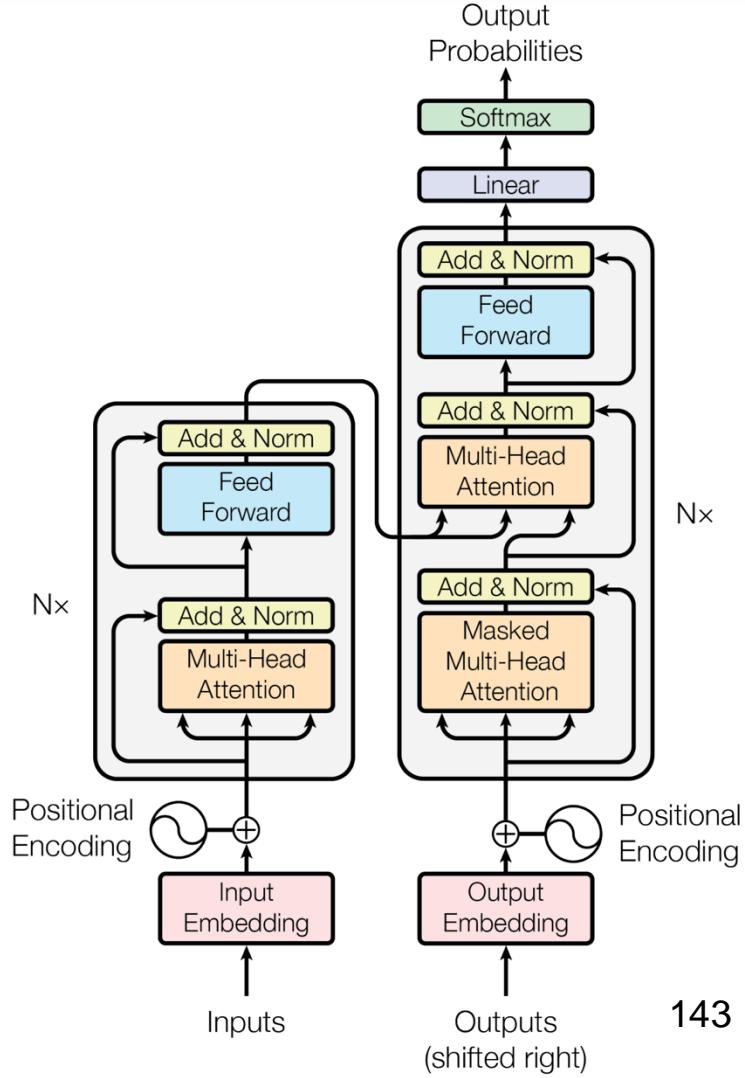
- At any step of **decoder**, it attends to previous computation of **encoder** as well as **decoder's** own generations



# Transformer

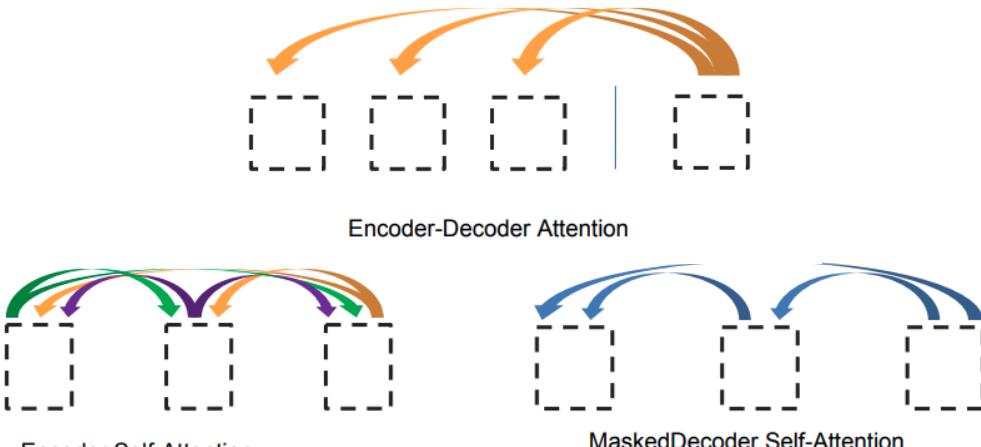
[Vaswani et al. 2017]

- At any step of **decoder**, it attends to previous computation of **encoder** as well as **decoder's** own generations
- At any step of **decoder**, **re-use** previous computation of **encoder**.
- Computation of **decoder** is **linear**, instead of quadratic.

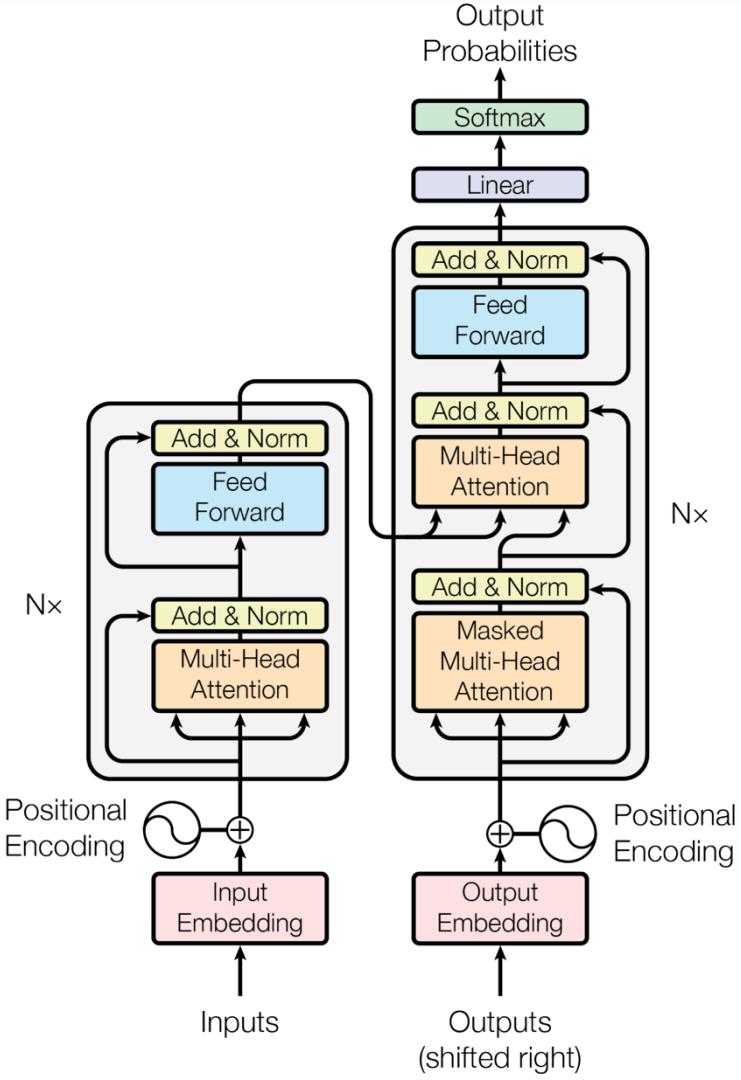


# Recap: Transformer

- Yaaay we know Transformers now! 🎉
- An **encoder-decoder** architecture
- 3 forms of attention



[[Attention Is All You Need, Vaswani et al. 2017](#)]



# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles,  $\sim 10$  web pages per query.
- For a passage of length  $N$  and a summary of length  $M$ , the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all  
and so quadratic.

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles,  $\sim 10$  web pages per query.
- For a passage of length  $N$  and a summary of length  $M$ , the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all  
and so quadratic in  $M$  and  $N$ .

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles,  $\sim 10$  web pages per query.
- For a passage of length  $N$  and a summary of length  $M$ , the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all  
and so quadratic in  $M$  and  $N$ .

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles,  $\sim 10$  web pages per query.
- For a passage of length  $N$  and a summary of length  $M$ , the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, cross attention is missing.

# Quiz: Enc-Dec Cost

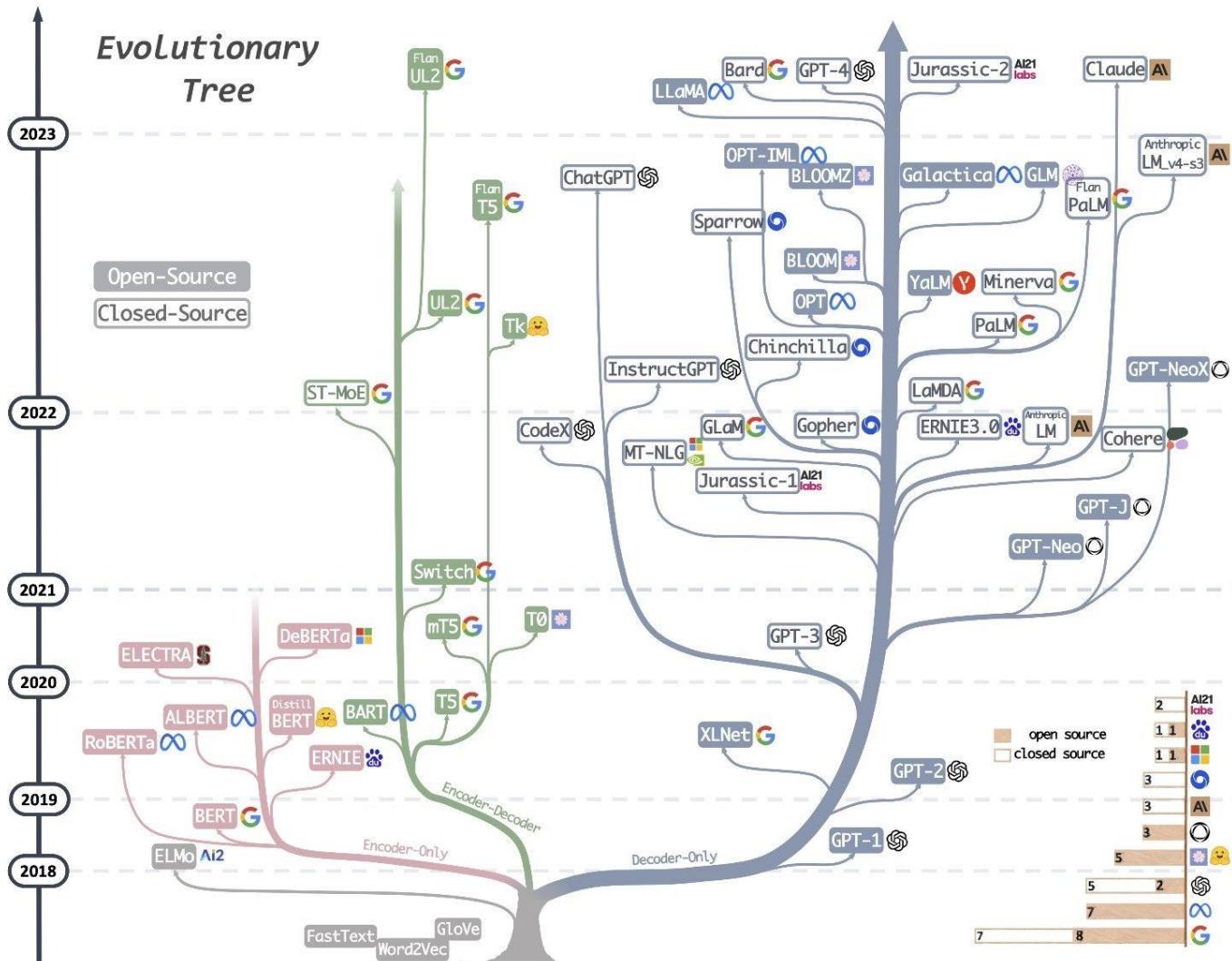
- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles,  $\sim 10$  web pages per query.
- For a passage of length  $N$  and a summary of length  $M$ , the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

Yes. The three terms are respectively the Encoder self-attention, Decoder self-attention, and Cross attention.

# After Transformer ...



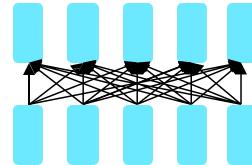
# *Evolutionary Tree*



Yang et al. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond, 2023

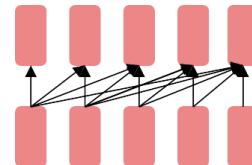
# Impact of Transformers

- A building block for a variety of LMs



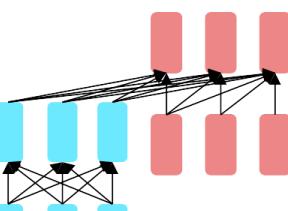
Encoders

- ❖ Examples: BERT, RoBERTa, SciBERT.
- ❖ Captures bidirectional context. Wait, how do we pretrain them?



Decoders

- ❖ Examples: GPT-2, GPT-3, LaMDA
- ❖ Other name: causal or auto-regressive language model
- ❖ Nice to generate from; can't condition on future words

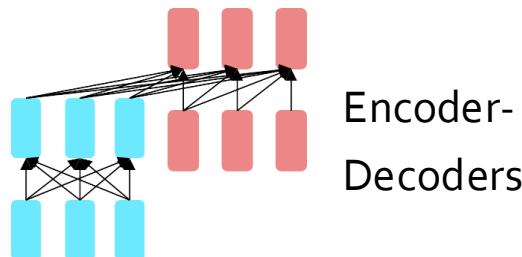


Encoder-  
Decoders

- ❖ Examples: Transformer, T5, Meena
- ❖ What's the best way to pretrain them?

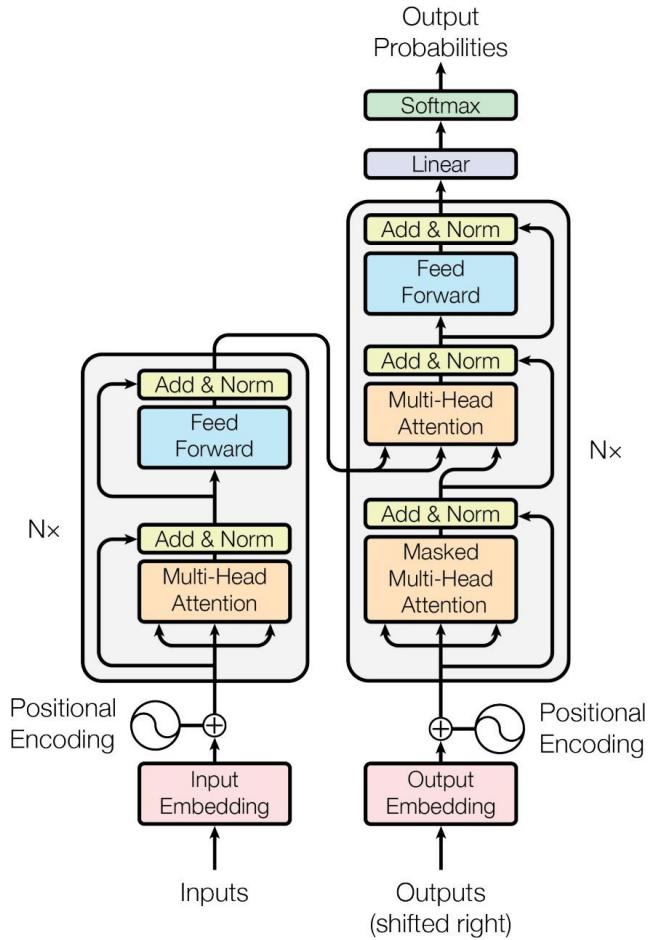
# Transformer Language Model Families

# Encoder-Decoder Family of Transformers



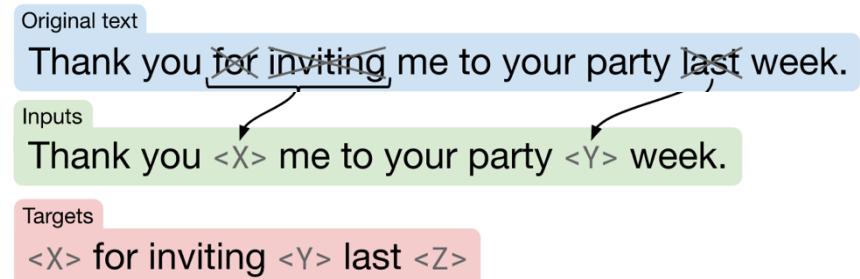
# Encoder-decoder Models

- The original transformer architecture was encoder decoder
- Encoder-decoder models are flexible in both generation and classification tasks
- How can we pretrain an encoder-decoder model like BERT to be a good general language pretrained LM?



# T5: Text-To-Text Transfer Transformer

- An encoder-decoder architecture
- Pre-training objective:  
corrupt and reconstruct objective



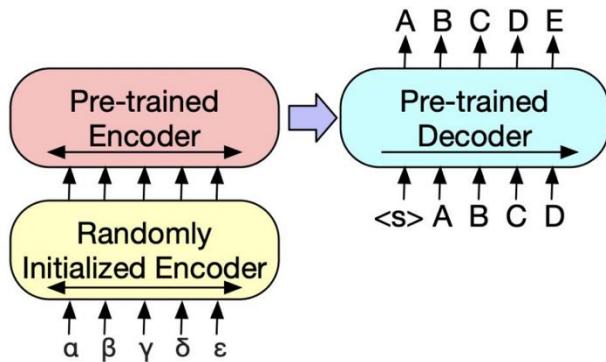
Model	Parameters	No. of layers	$d_{\text{model}}$	$d_{\text{ff}}$	$d_{\text{kv}}$	No. of heads
Small	60M	6	512	2048	64	8
Base	220M	12	768	3072	64	12
Large	770M	24	1024	4096	64	16
3B	3B	24	1024	16384	128	32
11B	11B	24	1024	65536	128	128

- The original paper is an excellent set of in-depth analysis of various parameters of model design. We discuss some of these results in other places.

<https://huggingface.co/t5-base>

# BART (Lewis et al. 2020)

- Similar Architecture as T5.
  - Corrupt the input -> ask the model to reconstruct the original input
  - Outperformed existing methods on generative tasks (question answering, and summarization).



**BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension**

Mike Lewis\*, Yinhan Liu\*, Naman Goyal\*, Marjan Ghazvininejad,  
Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, Luke Zettlemoyer  
Facebook AI

{mikelewis,yinhanliu,naman}@fb.com

# BART

The code might be outdated, but the logic is the same ...

```
from transformers import BartTokenizer, BartForConditionalGeneration

tokenizer = BartTokenizer.from_pretrained("facebook/bart-large")
model = BartForConditionalGeneration.from_pretrained("facebook/bart-large")

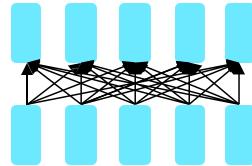
TXT = "The sun is <mask> ."
input_ids = tokenizer([TXT], return_tensors="pt")["input_ids"]
logits = model(input_ids).logits

masked_index = (input_ids[0] == tokenizer.mask_token_id).nonzero().item()
probs = logits[0, masked_index].softmax(dim=0)
values, predictions = probs.topk(5)

tokenizer.decode(predictions).split()
```

**Result:** ['located', 'at', 'approximately', 'also', 'about']

# Encoder-only Family of Transformers



BERT

# Bidirectional Encoder Representations from Transformers



# BERT

## Bidirectional Encoder Representations from Transformers

Like Bidirectional LSTMs (ELMo), let's look in both directions



# BERT

Bidirectional Encoder Representations from Transformers

Let's only use Transformer Encoders, no Decoders



# BERT

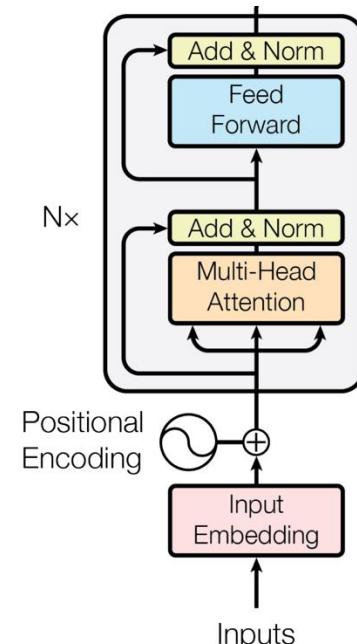
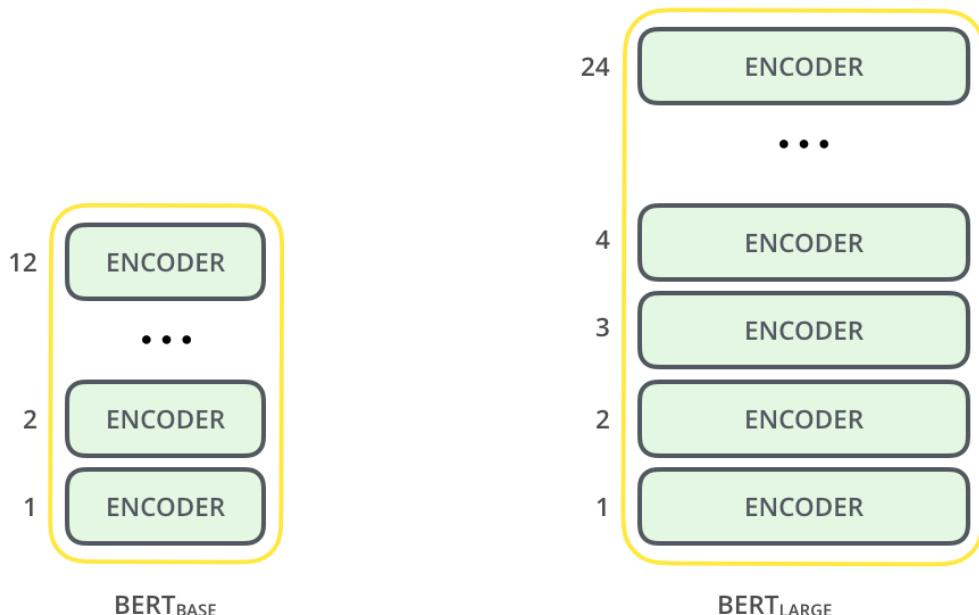
## Bidirectional Encoder Representations from Transformers

It's a language model that builds rich representations via self-supervised learning (pre-training)



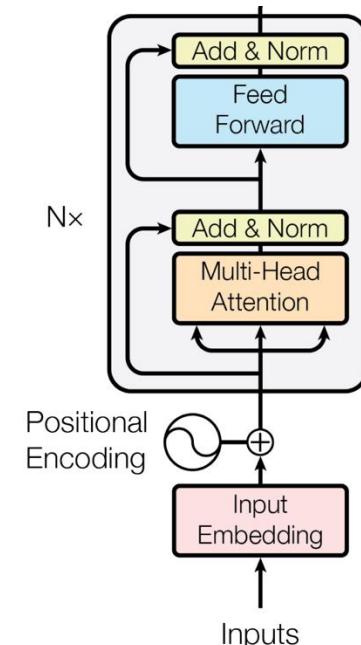
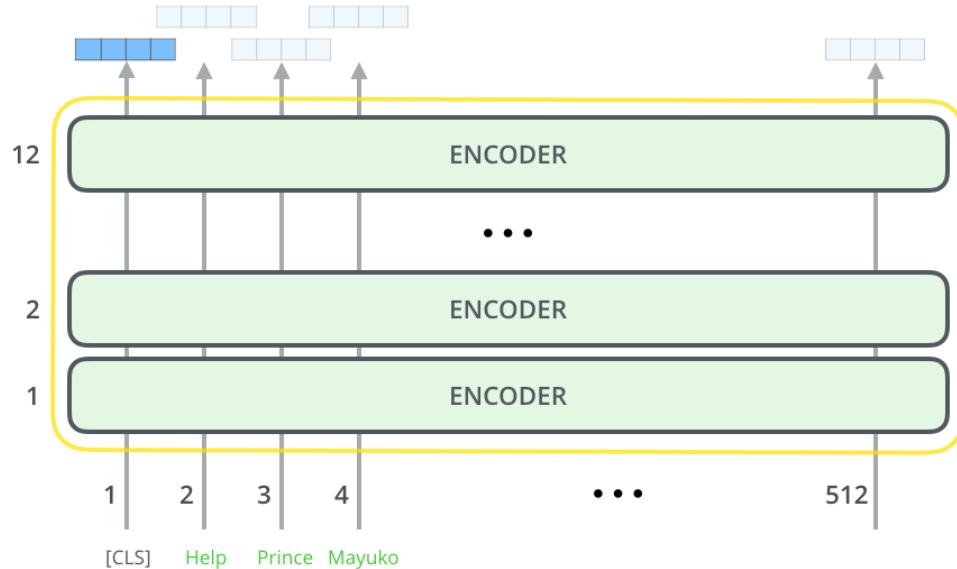
# BERT: Architecture

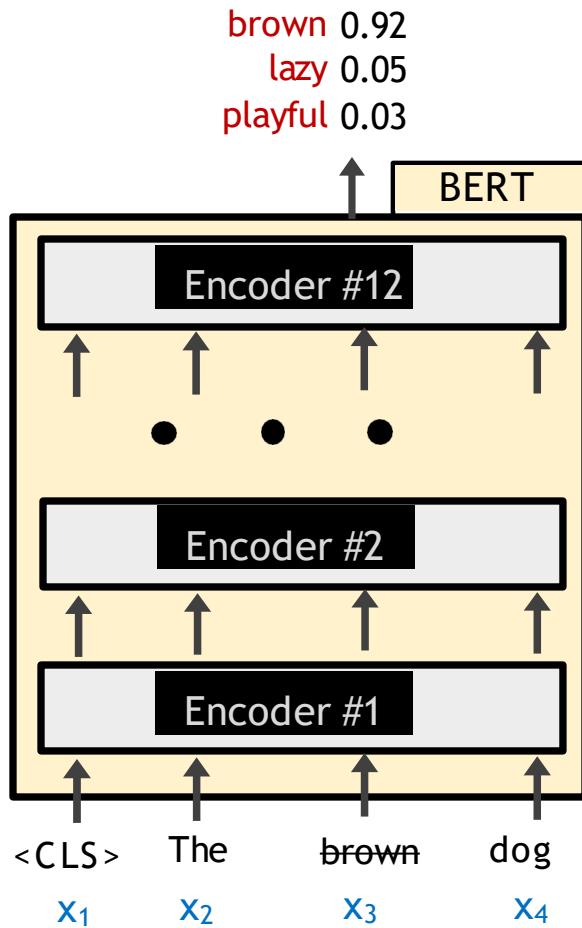
- Stacks of Transformer encoders



# BERT: Architecture

- Model output dimension: 512





BERT is trained to uncover masked tokens.

# Probing BERT Masked LM

- Masking words forces BERT to use context in both directions to predict the masked word.

Paris is the [MASK] of France.

Compute

Computation time on cpu: cached

capital	0.997
heart	0.001
center	0.000
centre	0.000
city	0.000

</> JSON Output

Maximize

# Probing BERT Masked LM

- Masking words forces BERT to use context in both directions to predict the masked word.

Today is Tuesday, so tomorrow is [MASK].

Compute

Computation time on cpu: cached



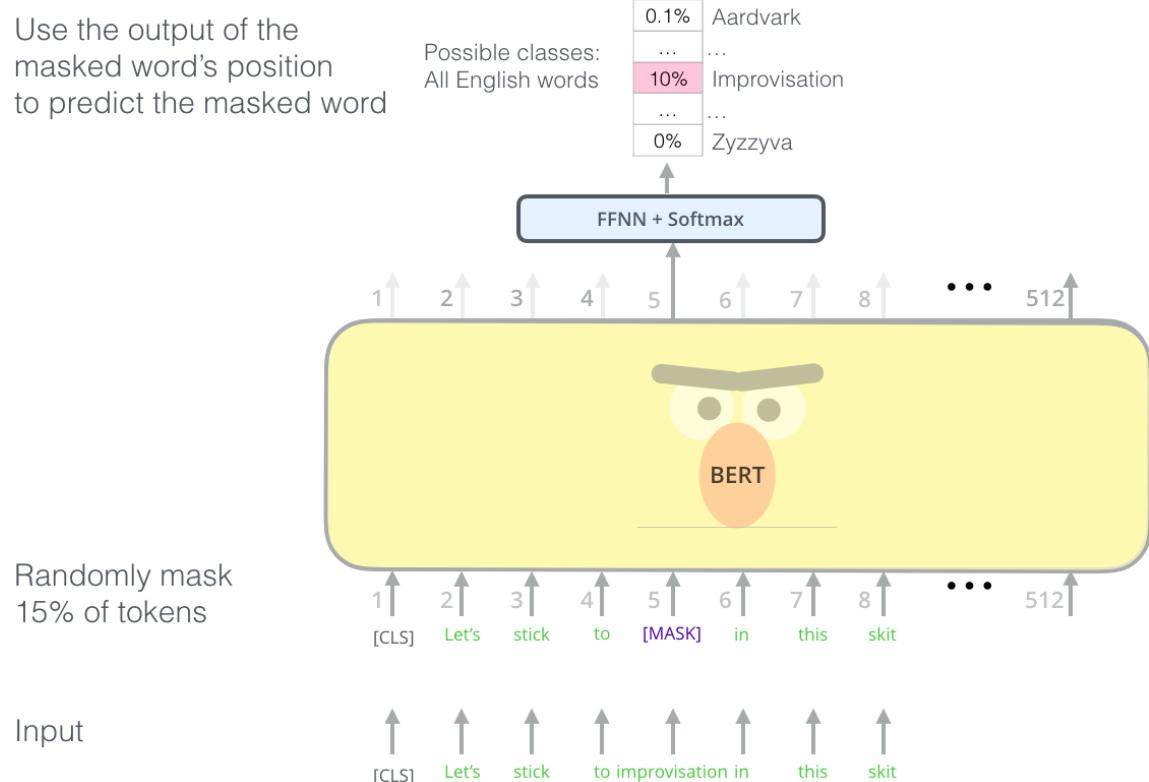
</> JSON Output

Maximize

# BERT: Pre-training Objective (1): Masked Tokens

- Randomly mask 15% of the tokens and train the model to predict them.

Use the output of the masked word's position to predict the masked word



# BERT: Pre-training Objective (1): Masked Tokens

store

Galon

the man went to the [MASK] to buy a [MASK] of milk

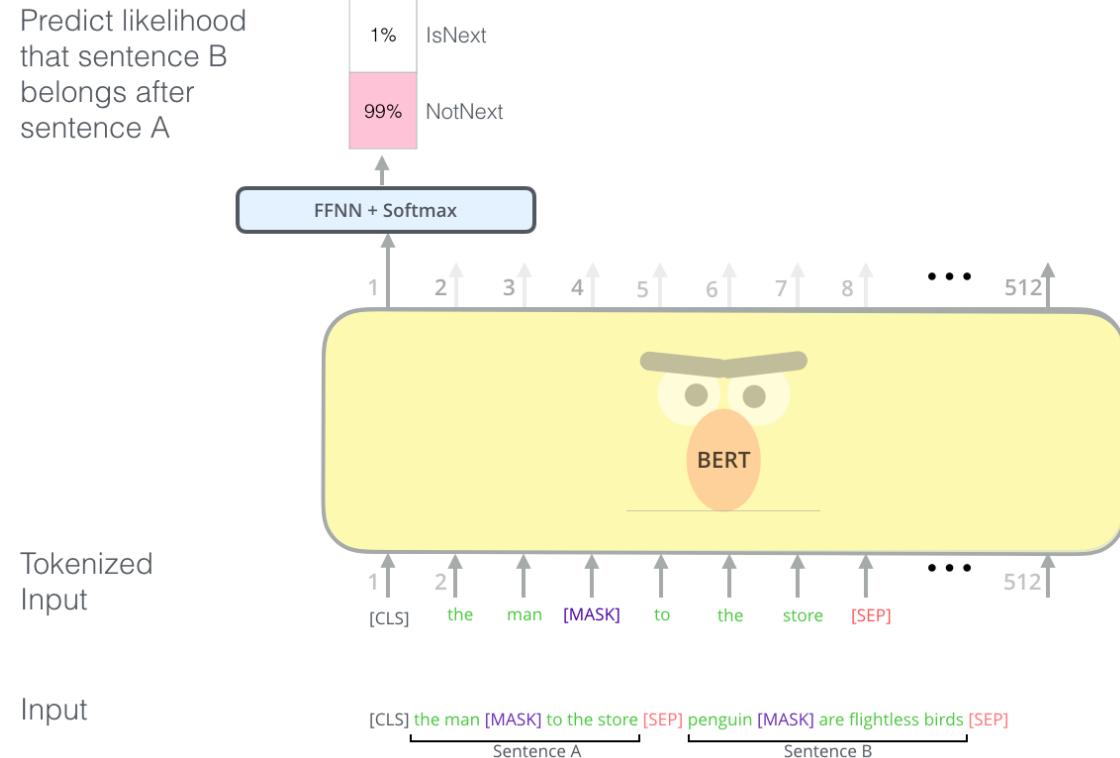
- Too little masking: Too **expensive** to train
- Too much masking: **Underdefined**
  - (not enough info for the model to recover the masked tokens)

Later work shows that more principled masking (instead of uniformly random) could benefit downstream task performance and result in faster training.

PMI Masking (Levine et al., 2021) <https://arxiv.org/pdf/2010.01825.pdf>  
SpanBERT (Joshi et al., 2020) <https://arxiv.org/pdf/1907.10529.pdf>

# BERT: Pre-training Objective (2): Sentence Ordering

- Predict sentence ordering
- 50% correct ordering, and 50% random incorrect ones



# BERT Pre-training Objective (2): Sentence Ordering

- Learn relationships between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

**Sentence A** = The man went to the store.

**Sentence B** = He bought a gallon of milk.

**Label** = IsNextSentence

**Sentence A** = The man went to the store.

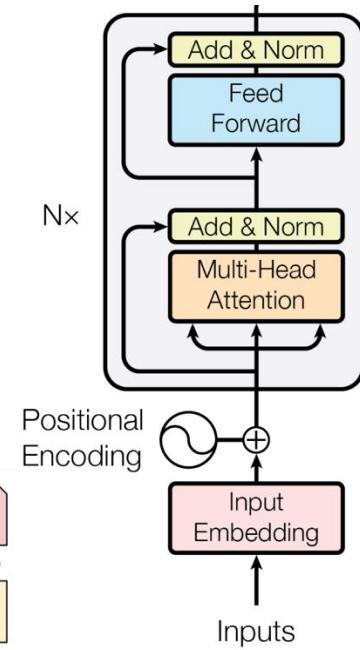
**Sentence B** = Penguins are flightless.

**Label** = NotNextSentence

# BERT: Input Representation

- Use 30,000 WordPiece vocabulary on input.
- Each token is sum of three embeddings
  - Addition to transformer encoder: sentence embedding

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	# #ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{\text{my}}$	$E_{\text{dog}}$	$E_{\text{is}}$	$E_{\text{cute}}$	$E_{[\text{SEP}]}$	$E_{\text{he}}$	$E_{\text{likes}}$	$E_{\text{play}}$	$E_{\#\#\text{ing}}$	$E_{[\text{SEP}]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$



# Training

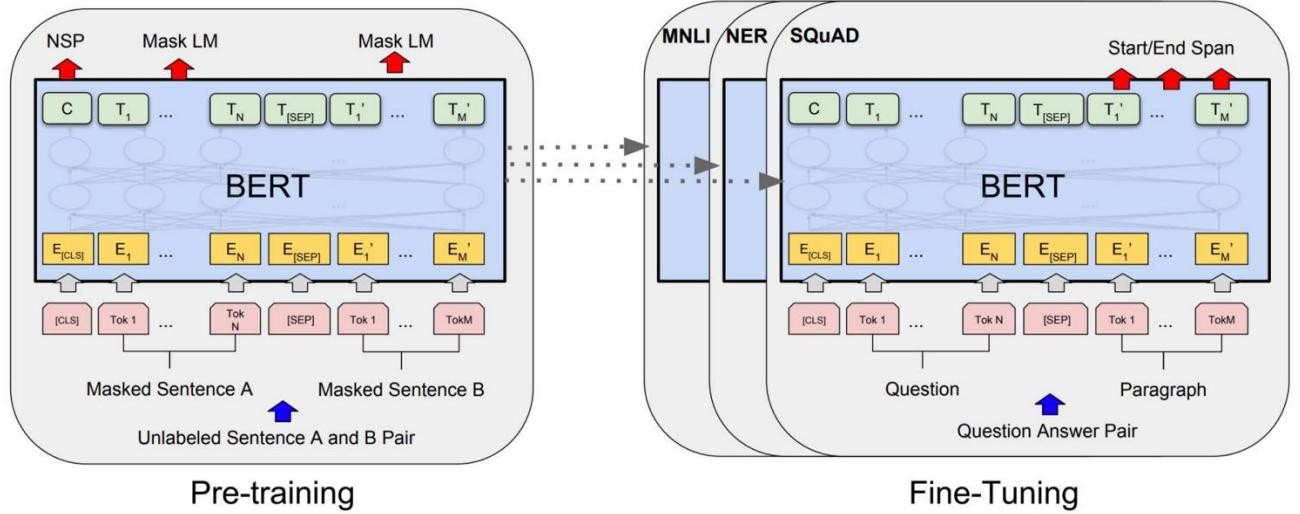
---

- Trains model on unlabeled data over different pre-training tasks (self-supervised learning)
- **Data:** Wikipedia (2.5B words) + BookCorpus (0.8B words)
- **Training Time:** 1M steps (~40 epochs)
- **Optimizer:** AdamW, 1e-4 learning rate, linear decay
- **BERT-Base:** 12-layer, 768-hidden, 12-head, sequence length of 512
- **BERT-Large:** 24-layer, 1024-hidden, 16-head, sequence length of 512
- Trained on 4x4 and 8x8 TPUs for 4 days (cost today using cloud TPU: \$1.3K and \$5K)

# Fine-tuning BERT

“Pretrain once, finetune many times.”

- **Idea:** Make pre-trained model **usable** in **downstream tasks**
- **Initialized** with pre-trained model parameters
- **Fine-tune** model parameters using labeled data from downstream tasks



## Leaderboard

- Human Performance (88.00%)
- Running Best
- Submissions

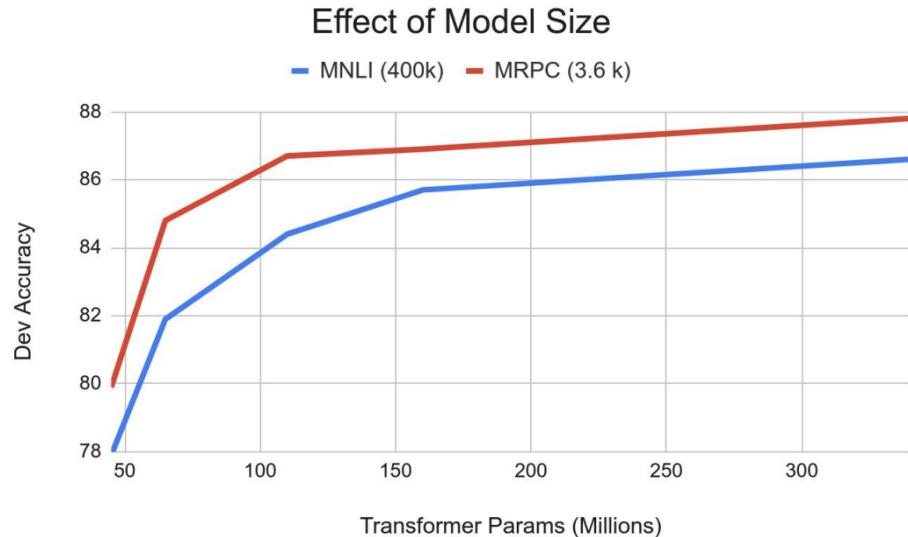
Rank	Model	Test Score
1	<b>BERT (Bidirectional Encoder Representations from Transfo...</b> <i>Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova</i> 10/11/2018	<b>86.28%</b>
2	<b>OpenAI Transformer Language Model</b> <i>Original work by Alec Radford, Karthik Narasimhan, Tim Salimans, ...</i> 10/11/2018	<b>77.97%</b>
3	<b>ESIM with ELMo</b> <i>Zellers, Rowan and Bisk, Yonatan and Schwartz, Roy and Choi, Yejin</i> 08/30/2018	<b>59.06%</b>
4	<b>ESIM with Glove</b> <i>Zellers, Rowan and Bisk, Yonatan and Schwartz, Roy and Choi, Yejin</i> 08/29/2018	<b>52.45%</b>

A girl is going across a set of monkey bars. She

- (i) jumps up across the monkey bars.
- (ii) struggles onto the bars to grab her head.
- (iii) gets to the end and stands on a wooden plank.
- (iv) jumps up and does a back flip.

- Run each Premise + Ending through BERT.
- Produce logit for each pair on token 0 ([CLS])

# Effect of Model Size



- Big models help a lot
- Going from 110M -> 340M params helps even on datasets with 3,600 labeled examples
- Improvements have **not** plateaued!

# Impact of BERT

---

- In order to have state-of-the-art performance on different tasks, there is no need for coming up with a novel model architecture
  - End of task-specific model architecture engineering
- An early sign that larger scales and self-supervised learning (language modeling) are the key for future performance improvements

# Why did no one think of this before?

---

- Why wasn't contextual pre-training popular before 2018 with ELMo?
- Good results on pre-training is  $>1,000x$  to 100,000 more expensive than supervised training.

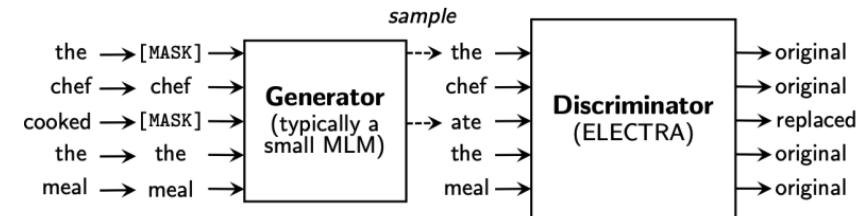
# What Happened After BERT?

- RoBERTa (Liu et al., 2019)
  - Exact same architecture as BERT
  - Drops the next sentence prediction loss!
  - Trained on 10x data (the original BERT was actually under-trained)
  - Much stronger performance than BERT (e.g., 94.6 vs 90.9 on SQuAD)

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

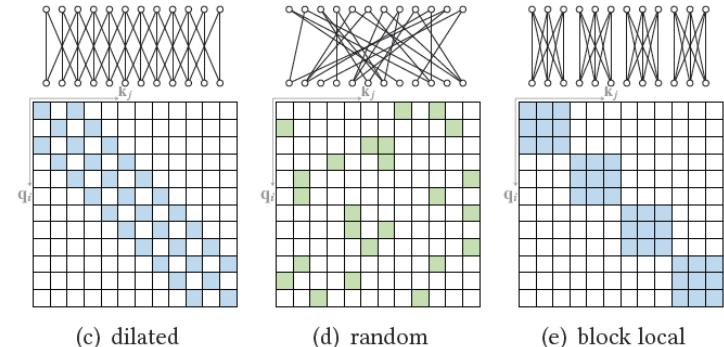
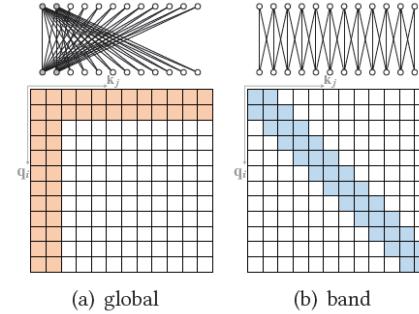
# What Happened After BERT?

- RoBERTa (Liu et al., 2019)
  - Exact same architecture as BERT
  - Drops the next sentence prediction loss!
  - Trained on 10x data (the original BERT was actually under-trained)
  - Much stronger performance than BERT (e.g., 94.6 vs 90.9 on SQuAD)
- ALBERT (Lan et al., 2020)
  - Increasing model sizes by sharing model parameters across layers
  - Less storage, much stronger performance but runs slower..
- ELECTRA (Clark et al., 2020)
  - Pre-training objective: replaced-token detection
  - Two models generator and discriminator (GAN-like)
  - It provides a more efficient training method



# What Happened After BERT?

- Models that handle long contexts
  - Longformer, Big Bird, ...
- Multilingual BERT
  - Trained single model on 104 languages from Wikipedia.
- BERT extended to different domains
  - SciBERT, BioBERT, FinBERT, ClinicalBERT, ...
- Making BERT smaller to use
  - DistillBERT, TinyBERT, ...



# Text generation using BERT

- Does not support generation or sequence-to-sequence tasks
  - Summarization, Translation, Text simplification, etc

**BERT has a Mouth, and It Must Speak:  
BERT as a Markov Random Field Language Model**

Alex Wang  
New York University  
alexwang@nyu.edu

Kyunghyun Cho  
New York University  
Facebook AI Research  
CIFAR Azrieli Global Scholar  
kyunghyun.cho@nyu.edu

**Mask-Predict: Parallel Decoding of  
Conditional Masked Language Models**

Marjan Ghazvininejad\*      Omer Levy\*      Yinhan Liu\*      Luke Zettlemoyer  
Facebook AI Research  
Seattle, WA

## Exposing the Implicit Energy Networks behind Masked Language Models via Metropolis–Hastings

Kartik Goyal, Chris Dyer, Taylor Berg-Kirkpatrick

## Leveraging Pre-trained Checkpoints for Sequence Generation Tasks

Sascha Rothe, Shashi Narayan, Aliaksei Severyn

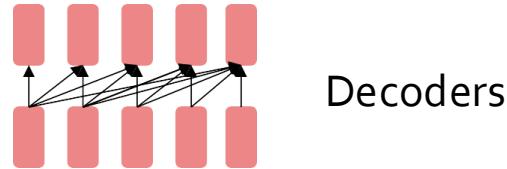
src	Der Abzug der franzsischen Kampftruppen wurde am 20. November abgeschlossen .
t = 0	The departure of the French combat completed completed on 20 November .
t = 1	The departure of French combat troops was completed on 20 November .
t = 2	The withdrawal of French combat troops was completed on November 20th .

# Summary Thus Far

---

- BERT and the family
- An encoder; Transformer-based networks trained on massive piles of data.
- Incredible for learning **contextualized** embeddings of words
- It's very useful to **pre-train** a large unsupervised/self-supervised LM then **fine-tune** on your particular task (replace the top layer, so that it can work)
- However, they were **not** designed to generate text.

# Decoder-only Family of Transformers



# GPT

Generative Pre-trained Transformer

GPT-2: A Big Language Model (2019)

---

Language Models are Unsupervised Multitask Learners

---

Alec Radford <sup>\*1</sup> Jeffrey Wu <sup>\*1</sup> Rewon Child <sup>1</sup> David Luan <sup>1</sup> Dario Amodei <sup>\*\*1</sup> Ilya Sutskever <sup>\*\*1</sup>

GPT: An Auto-Regressive LM (2018)

---

Improving Language Understanding  
by Generative Pre-Training

---

Alec Radford  
OpenAI  
alec@openai.com

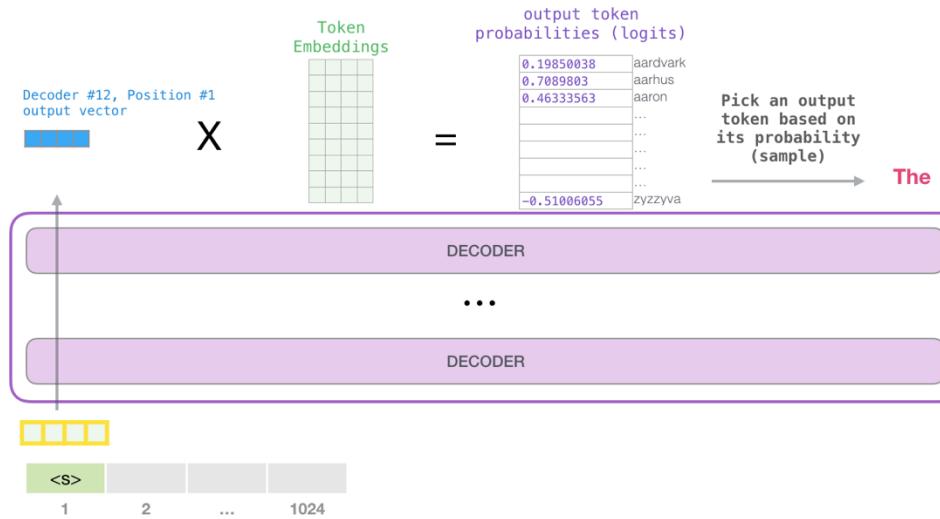
Karthik Narasimhan  
OpenAI  
karthikn@openai.com

Tim Salimans  
OpenAI  
tim@openai.com

Ilya Sutskever  
OpenAI  
ilyasu@openai.com

# GPT-2

- GPT-2 uses only Transformer Decoders (no Encoders) to generate new sequences from scratch or from a starting sequence
- As it processes each subword, it masks the “future” words and conditions on and attends to the previous words

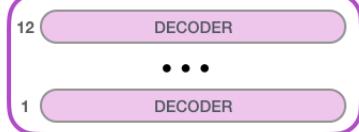


# GPT2: Model Sizes

Play with it here: <https://huggingface.co/gpt2>



GPT-2  
SMALL

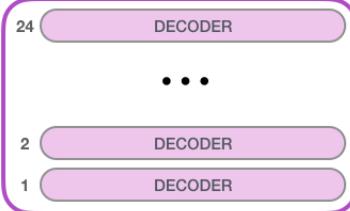


Model Dimensionality: 768

117M parameters



GPT-2  
MEDIUM

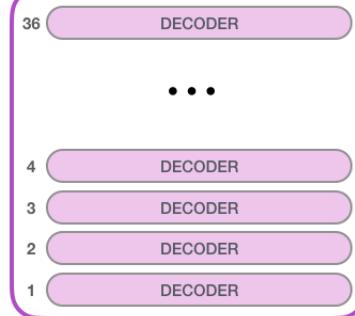


Model Dimensionality: 1024

345M



GPT-2  
LARGE

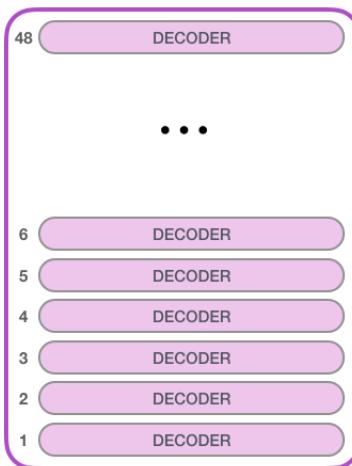


Model Dimensionality: 1280

762M



GPT-2  
EXTRA  
LARGE



Model Dimensionality: 1600

1542M

# GPT2: Some Results

Language Models are Unsupervised Multitask Learners

	LAMBADA (PPL)	LAMBADA (ACC)	CBT-CN (ACC)	CBT-NE (ACC)	WikiText2 (PPL)	PTB (PPL)	enwik8 (BPB)	text8 (BPC)	WikiText103 (PPL)	1BW (PPL)
SOTA	99.8	56.25	85.7	82.3	39.14	46.54	0.99	1.08	18.3	<b>21.8</b>
117M	<b>35.13</b>	45.99	<b>87.65</b>	<b>83.4</b>	<b>29.41</b>	65.85	1.16	1.17	37.50	75.20
345M	<b>15.60</b>	55.48	<b>92.35</b>	<b>87.1</b>	<b>22.76</b>	47.33	1.01	<b>1.06</b>	26.37	55.72
762M	<b>10.87</b>	<b>60.12</b>	<b>93.45</b>	<b>88.0</b>	<b>19.93</b>	<b>40.31</b>	<b>0.97</b>	<b>1.02</b>	22.05	44.575
1542M	<b>8.63</b>	<b>63.24</b>	<b>93.30</b>	<b>89.05</b>	<b>18.34</b>	<b>35.76</b>	<b>0.93</b>	<b>0.98</b>	<b>17.48</b>	42.16

Table 3. Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). Other language model results are from (Dai et al., 2019).

**Article:** Prehistoric man sketched an incredible array of prehistoric beasts on the rough limestone walls of a cave in modern day France 36,000 years ago.

Now, with the help of cutting-edge technology, those works of art in the Chauvet-Pont-d'Arc Cave have been reproduced to create the biggest replica cave in the world.

The manmade cavern named the Caverne du Pont-d'Arc has been built a few miles from the original site in Vallon-Pont-D'arc in Southern France and contains 1,000 painstakingly-reproduced drawings as well as around 450 bones and other features...

Cavemen and women sketched an incredible array of prehistoric beasts on the rough limestone walls of a cave 36,000 years ago and now a replica has been created (pictured)

...

**GPT-2:** The original site in Vallon-Pont-D'arc in Southern France is a Unesco World Heritage site and is the oldest known and the best preserved cave decorated by man. The replica cave was built a few miles from the original site in Vallon-Pont-D'Arc in Southern France. The cave contains images of 14 different species of animals including woolly rhinoceros, mammoths, and big cats.

**Reference:** Cave mimics famous Caverne du Pont-d'Arc in France, the oldest cave decorated by man and the best preserved. The replica contains all 1,000 paintings which include 425 such as a woolly rhinoceros and mammoths. Minute details were copied using 3D modelling and anamorphic techniques, often used to shoot widescreen images. The modern cave also includes replica paw prints of bears, bones and details preserved in the original cave.

# Impact of GPT2

---

- Zero-shot learning (no use of task-specific supervision) increasingly become a reality.

NMT: “Translate to french,” <English text>, <French text>.

QA: “Answer the question,” <Document>, <Question>, <Answer>.

SUMM: <Document> “TL; DR:” <Summarization>

# GPT-3: A Very Large Language Model (2020)

- More layers & parameters
- Bigger dataset
- Longer training
- Larger embedding/hidden dimension
- Larger context window



# Impact of GPT3

---

- Moving away from the fine-tuning paradigm
  - Zero/Few-shot learning and in-context learning
- Massive LM scale makes high zero/few-shot performance possible
- Start of closed source models
  - Not too many details about their model
  - No released code / model checkpoint
- Also revitalized open source efforts:
  - OPT, LLaMA by Meta, BLOOM by Huggingface, etc.

# GPT4

- Transformer-based
  - The rest is .... mystery! 😊
  - If we're going based on costs, GPT4 is ~15-30 times costlier than GPT3. That should give you an idea how its likely size!
- Note, these language models involve more than just pre-training.
  - Pre-training provides the foundation based on which we build the model.
  - We will discuss the later stages (post hoc alignment) in a 2-3 weeks.

Model	Usage	
davinci-002	\$0.0020 / 1K tokens	
Model	Input	Output
gpt-4	\$0.03 / 1K tokens	\$0.06 / 1K tokens

# Accessing API Models

```
import openai
```

```
openai.api_key = ("sk-[REDACTED]")
```

```
my_prompt = '''The sun is [MASK].
```

Replace [MASK] with the most probable 5 words to replace, and give me their probabilities.'''

```
# Here set parameters as you like
```

```
response = openai.Completion.create(
```

```
    engine="text-davinci-002",
```

```
    prompt=my_prompt,
```

```
    temperature=0,
```

```
    max_tokens=100,
```

```
)
```

```
print(response['choices'][0]['text'])
```

# Other Available [Decoder] LMs

---

EleutherAI: GPT-Neo (6.7B), GPT-J (6B), GPT-NeoX (20B)

<https://huggingface.co/EleutherAI>

<https://6b.eleuther.ai/>

LLaMA, 65B: <https://github.com/facebookresearch/llama>

Mistral and Mixtral:

<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

<https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

# Training Transformer LMs: Empirical Considerations

# Pre-training Transformer LMs

---

- You have learned about the basics of pre-training Transformer language models.
- There is so much empirical knowledge/experiences that goes into training these models.
- Various empirical issues about:
  - Preparation/pre-processing data
  - Efficient training of models
  - ...

# C4: The Data

---

- C4: Colossal Clean Crawled Corpus
  - Web-extracted text
  - English language only
  - 750GB

Data set	Size
★ C4	745GB
C4, unfiltered	6.1TB

# C4: The Data

Remove any:

- References to Javascript
- “Lorem ipsum” text — placeholder text commonly used to demonstrate the visual form of a document

Menu

Lemon

Introduction

The lemon, Citrus Limon (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae. The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a ph of around 2.2, giving it a sour taste.

Article

The origin of the le

vn, though

Retain:

- Sentences with terminal punctuation marks
- Pages with at least 5 sentences, sentences with at least 3 words

Please enable JavaScript to use our site.

Home

Products

Shipping

Contact

FAQ

Dried Lemons, \$3.59/pound

Organic dried lemons from our farm in California.

Lemons are harvested and sun-dried for maximum flavor.

Good in soups and on popcorn.

The lemon, Citrus Limon (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae.

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a ph of around 2.2, giving it a sour taste.

Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Curabitur in tempus quam. In mollis et ante at consectetur.

Aliquam erat volutpat.

Donec at lacinia est.

Duis semper, magna tempor interdum suscipit, ante elit molestie urna, eget efficitur risus nunc ac elit.

Fusce quis blandit lectus.

Mauris at mauris a turpis tristique lacinia at nec ante.

Aenean in scelerisque tellus, a efficitur ipsum:

Integer justo enim, ornare vitae sem non, mollis fermentum lectus.

Mauris ultrices nisl at libero porta sodales in ac orci.

```
function Ball(r) {  
    this.radius = r;  
    this.area = pi * r ** 2;  
    this.show = function(){  
        drawCircle(r);  
    }  
}
```

# Pre-training Data: Experiment

- Takeaway:
  - Clean and compact data is better than large, but noisy data.
  - Pre-training on in-domain data helps.

Data set	Size	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21

# Pre-training Data Duplicates

- There is a non-negligible number of duplicates in any pre-training data.

	% train examples with dup in train	% valid with dup in train	% valid with dup in valid
C4	3.04%	1.59%	4.60%
RealNews	13.63%	1.25%	14.35%
LM1B	4.86%	0.07%	4.92%
Wiki40B	0.39%	0.26%	0.72%

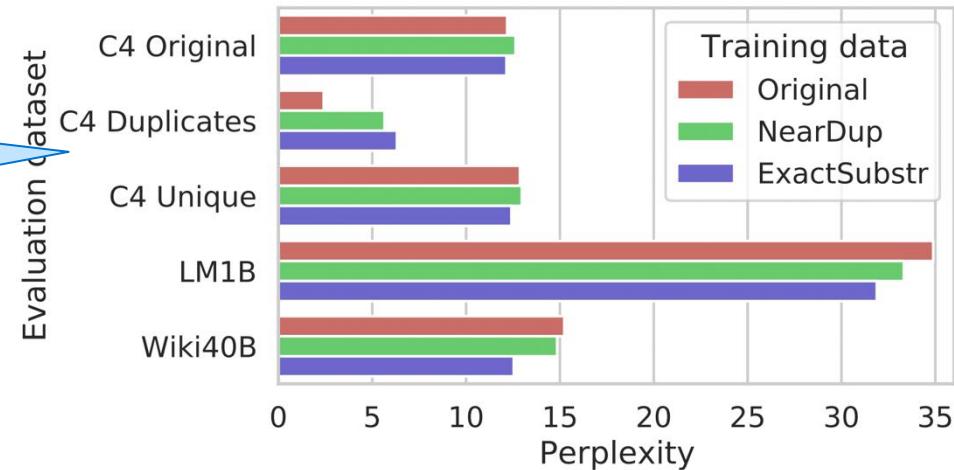
Dataset	Example	Near-Duplicate Example
Wiki-40B	\n_START_ARTICLE_\nHum Award for Most Impactful Character\n\n_START_SECTION_\nWinners and nominees\n\n_START_PARAGRAPH_\nIn the list below, winners are listed first in the colored row, followed by the other nominees. [...]	\n_START_ARTICLE_\nHum Award for Best Actor in a Negative Role\n\n_START_SECTION_\nWinners and nominees\n\n_START_PARAGRAPH_\nIn the list below, winners are listed first in the colored row, followed by the other nominees. [...]
LM1B	I left for California in 1979 and tracked Cleveland 's changes on trips back to visit my sisters .	I left for California in 1979 , and tracked Cleveland 's changes on trips back to visit my sisters .
C4	Affordable and convenient holiday flights take off from your departure country, "Canada". From May 2019 to October 2019, Condor flights to your dream destination will be roughly 6 a week! Book your Halifax (YHZ) - Basel (BSL) flight now, and look forward to your "Switzerland" destination!	Affordable and convenient holiday flights take off from your departure country, "USA". From April 2019 to October 2019, Condor flights to your dream destination will be roughly 7 a week! Book your Maui Kahului (OGG) - Dubrovnik (DBV) flight now, and look forward to your "Croatia" destination!

# Deduplicating Data Improves LMs

- Models: GPT-2-like (1.5B param) models
- On these datasets:
  - C4 : the original training data
  - C4-NearDup: C4 excluding exact duplicates
  - C4-ExactSubs: C4 excluding near-duplicates

Except when evaluated on duplicate evaluation data!

Training on deduplicated data always leads to lower PPL!

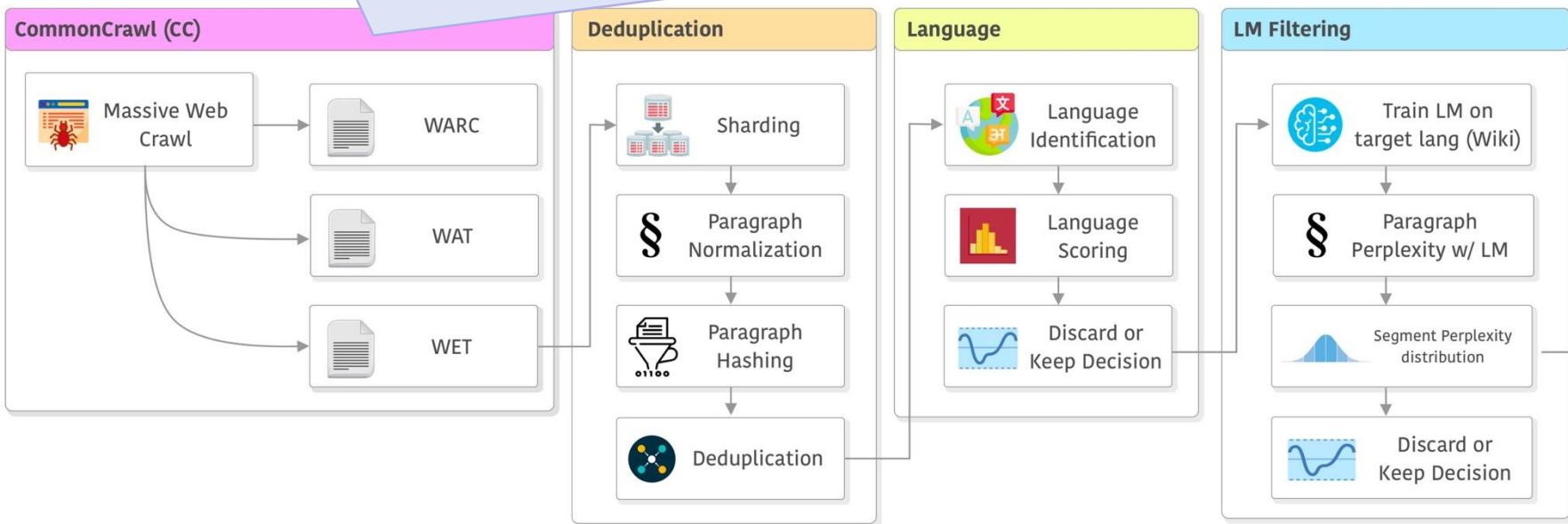


# LLaMA's Data Pipeline

Starts with the massive crawled data by CommonCrawl.

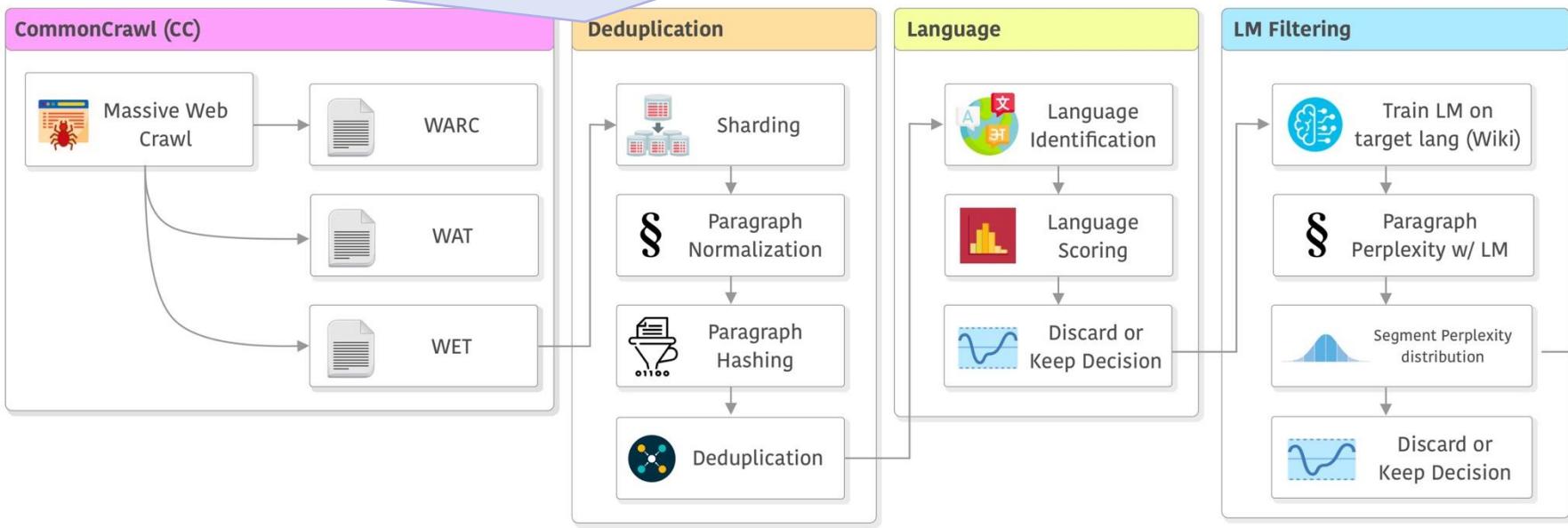
The WET format that contains textual information.

WARC is raw, WAT is metadata, WET is text+some metadata.



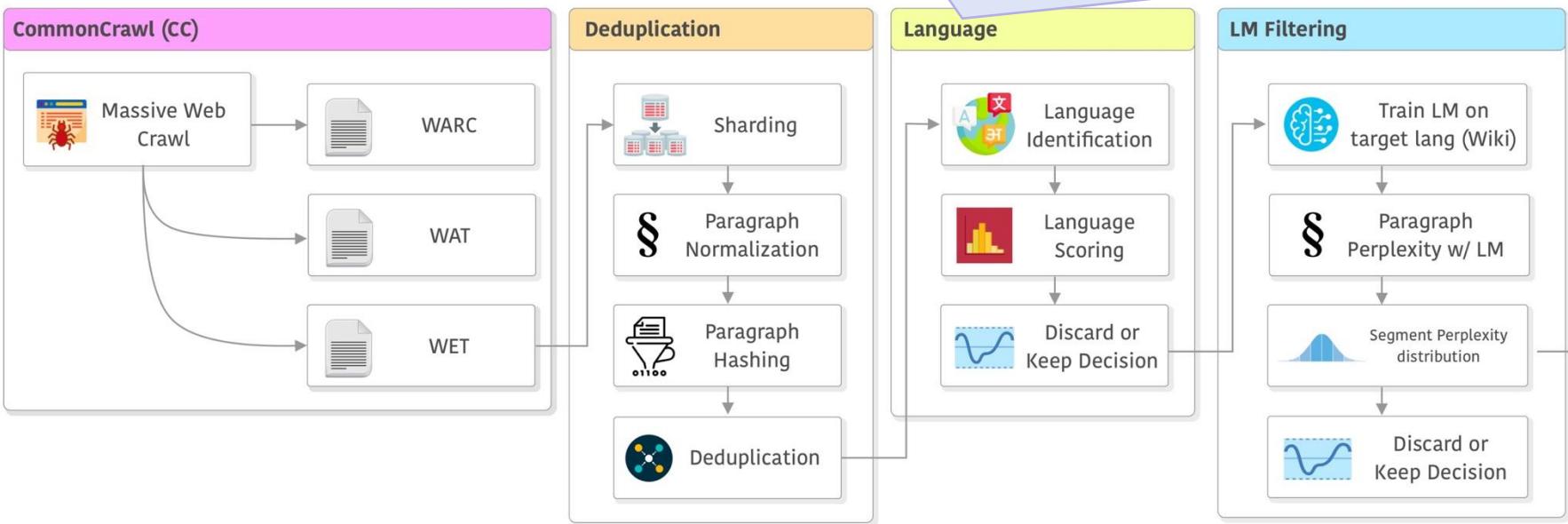
# LLaMA's Data Pipeline

Shard WET content into shards of 5GB each (one CC snapshot can have 30TB). Then you normalize paragraphs (lowercasing, numbers as placeholders, etc), compute per-paragraph hashes and then duplicate them.



# LLaMA's Data Pipeline

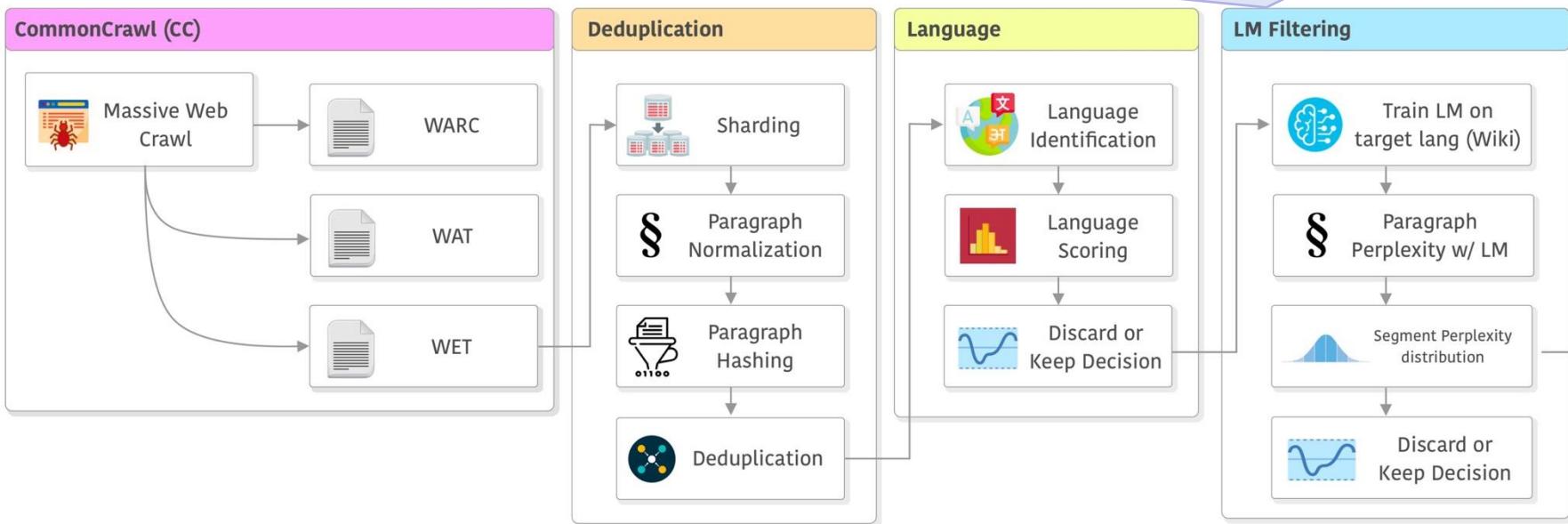
Perform language identification and decide whether to keep or discard languages.  
The order of when you do this in the pipeline can impact the language discrimination quality.



# LLaMA's Data Pipeline

Do further quality filtering: Train a simple LM (n-gram) on target languages using Wikipedia, then compute per-paragraph perplexity on the rest of the data:

- Very high PPL: Very different than Wiki and likely low-quality → Drop
- Very low PPL: Very similar or near duplicates to Wiki → Drop



# Grouped Query-Attention

- Used for training LLaMA 2.
- One key-value vector for each group of queries — an interpolation between “multi-head” attention and “multi-query” attention.

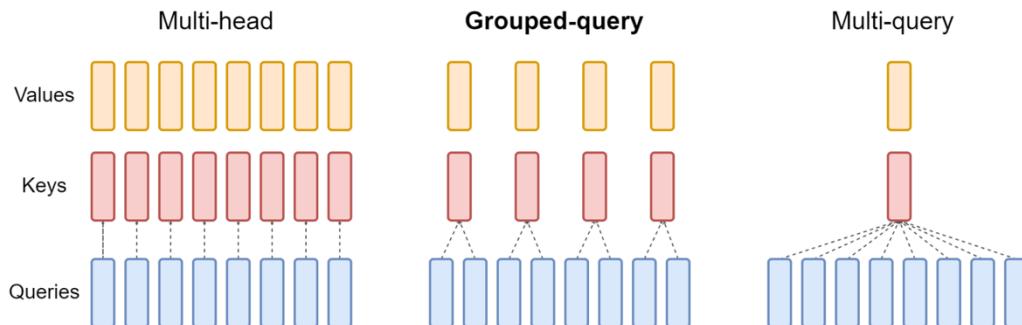


Figure 2: Overview of grouped-query method. Multi-head attention has  $H$  query, key, and value heads. Multi-query attention shares single key and value heads across all query heads. Grouped-query attention instead shares single key and value heads for each *group* of query heads, interpolating between multi-head and multi-query attention.

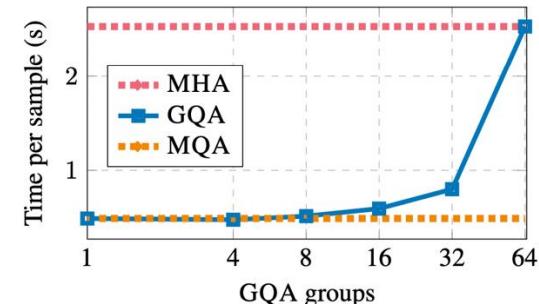
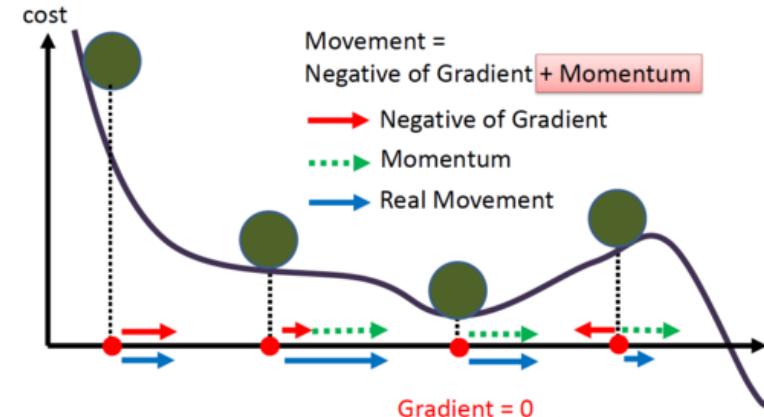


Figure 6: Time per sample for GQA-XXL as a function of the number of GQA groups with input length 2048 and output length 512. Going from 1 (MQA) to 8 groups adds modest inference overhead, with increasing cost to adding more groups.

- Improves inference scalability for our larger models

# Optimizers

- Most modern models use “AdamW” optimizer (not vanilla Gradient Descent).
  - Adam optimization is a stochastic gradient descent method that is based on adaptive estimation of first-order and second-order “momentums”.
  - “W” because it decouples “weight decay” from “learning rate”. (Details out of scope for us. See the cited paper.)



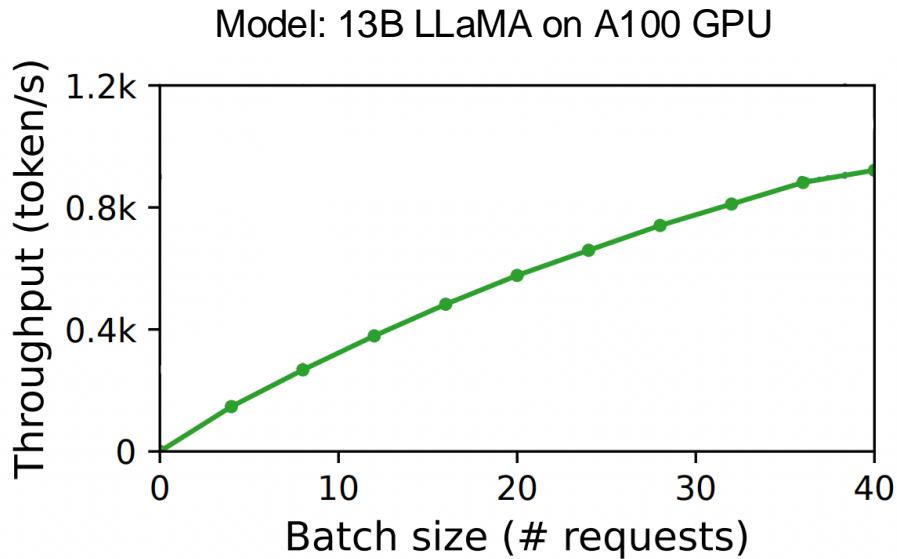
<https://pytorch.org/docs/stable/generated/torch.optim.AdamW.html>

<https://pytorch.org/docs/stable/generated/torch.optim.Adam.html>

[Decoupled Weight Decay Regularization, 2017]

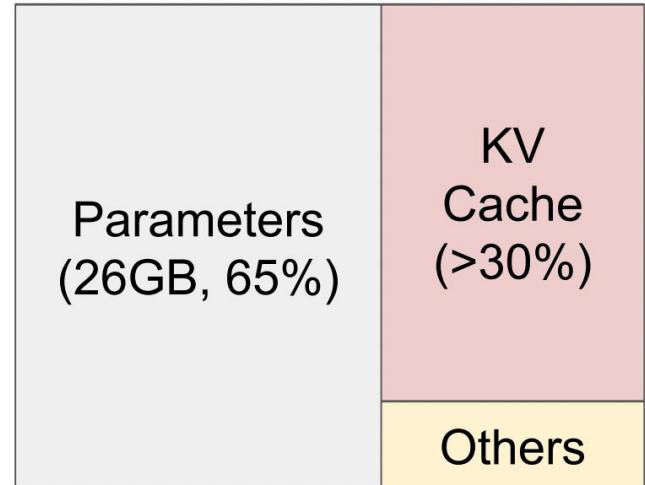
# Batching Data

- Previously we talked about the importance of batching data
- GPUs are faster at Tensor operations and hence, we want to do batch processing
- The larger batch of data, the faster they get processed.
- Alas, the speedup is often sub-linear (e.g., 2x larger batch leads to less than 2x speedup).



# The Memory Usage

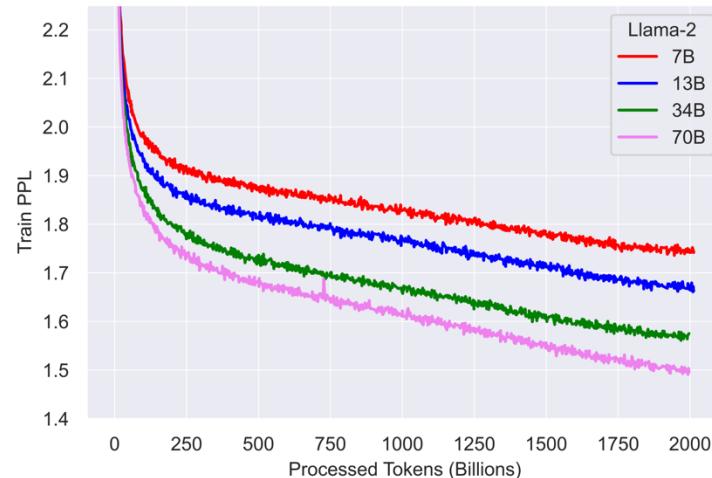
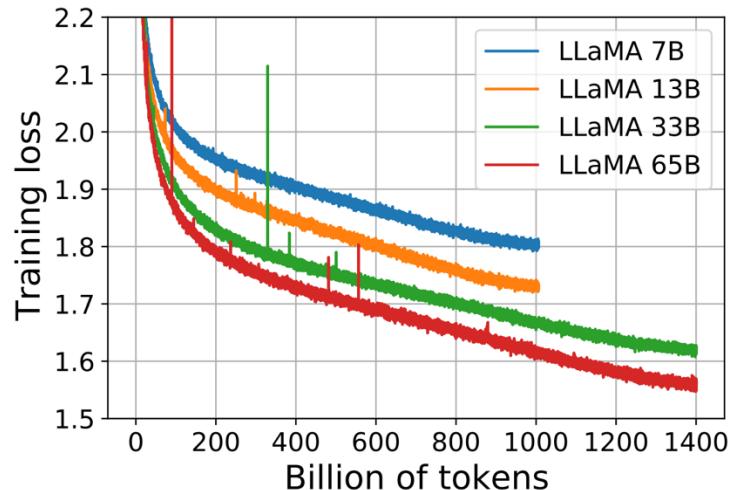
- Here is the memory usage of an NVIDIA A100 when serving (i.e., no training)
  - Model: 13B LLaMA
  - Batch size of 10
- ~65% of your GPU memory is the model parameters that **never change**
- ~32% of your memory are KV tensors that **change for each input.**
  - This KV cache will increase for larger batch sizes.
  - Managing this part of the memory is key for efficient training.



NVIDIA A100 40GB

# Convergence

- In practice, your model's loss should continue to go down with more training on more data.
- So, the real bottlenecks are:
  - (1) compute
  - (2) data
- Sometimes training diverges (spikes in the loss), at which point practitioners usually restart training from an earlier checkpoint.



# Summary

---

- There is much empirical knowledge that goes into engineering LMs.
- Here we covered basic topics about data and architecture engineering.
- Various topics are forthcoming: scaling laws, efficient training, etc.