



《多模态机器学习》

第八章 多模态自监督学习

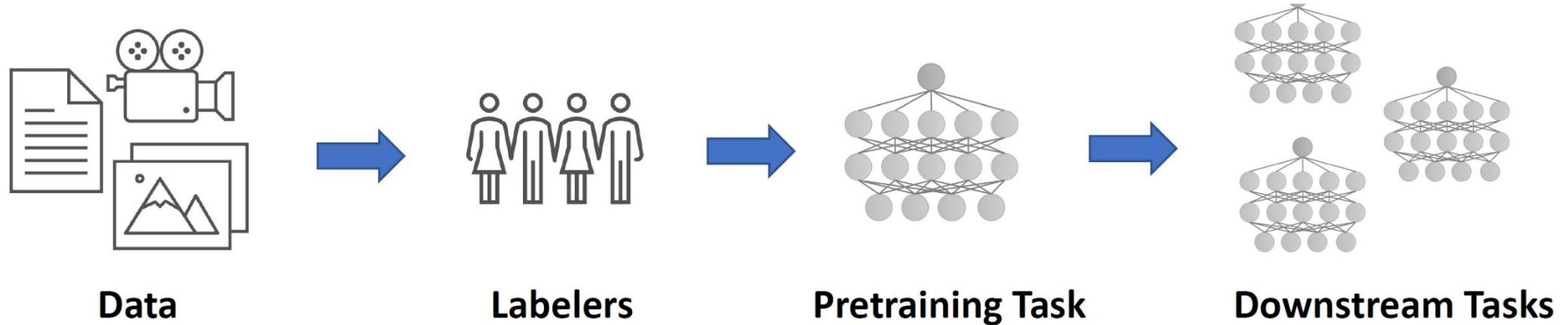
黄文炳

中国人民大学高瓴人工智能学院

hwenbing@126.com

2024年秋季

Supervised pretraining on large labeled, datasets has led to successful transfer learning



[\[Deng et al., 2009\]](#)

ImageNet

- Pretrain for fine-grained image classification over 1000 classes
- Use feature representations for downstream tasks, e.g. object detection, image segmentation, and action recognition

Supervised pretraining on large labeled, datasets has led to successful transfer learning



SNLI Dataset

Premise:

Ruth Bader Ginsburg being appointed to the US Supreme Court.



Hypothesis:

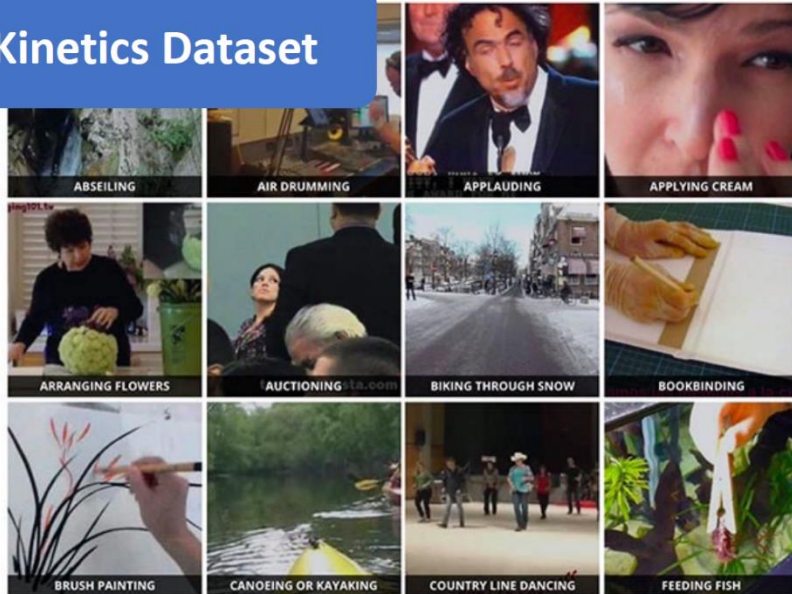
A grilled sandwich on a plate.



Label:

Contradiction [different scenes]

Kinetics Dataset



Across images, video, and text

[\[Deng et al., 2009\]](#) [\[Carreira et al., 2017\]](#) [\[Conneau et al., 2017\]](#)

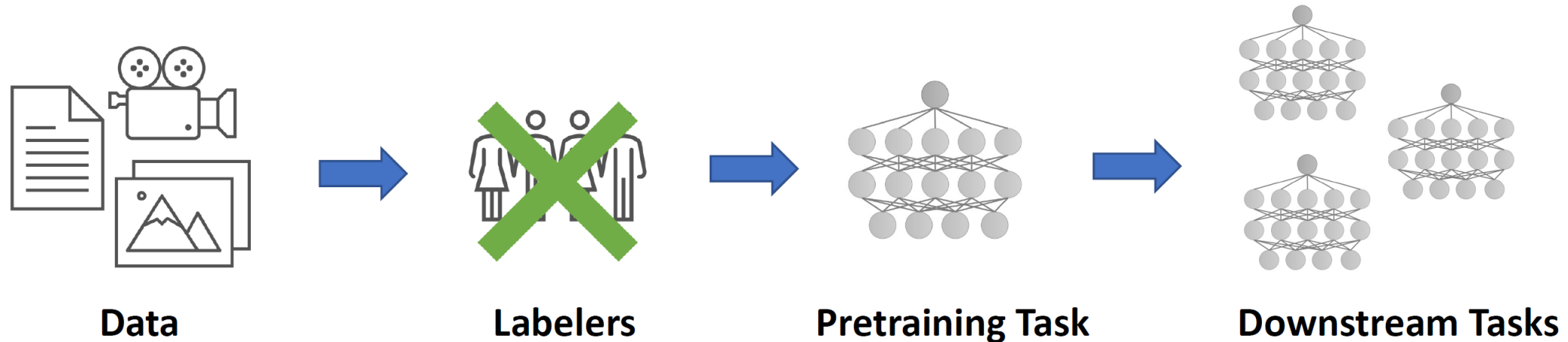
But supervised pretraining comes at a cost...

- **Time-consuming and expensive** to label datasets for new tasks
 - ImageNet: 3 years, 49k Amazon MechanicalTurkers [\[1\]](#)
- **Domain expertise needed** for specialized tasks
 - Radiologists to label medical images
 - Native speakers or language specialists for labeling text in different languages



Can self-supervised learning help?

- Self-supervised learning (informal definition): supervise using labels ***generated from the data*** without any manual or weak label sources
- Idea: Hide or modify part of the input. Ask model to recover input or classify what changed.
 - Self-supervised task referred to as the pretext task



Pretext Task: Classify the Rotation



270° rotation



90° rotation

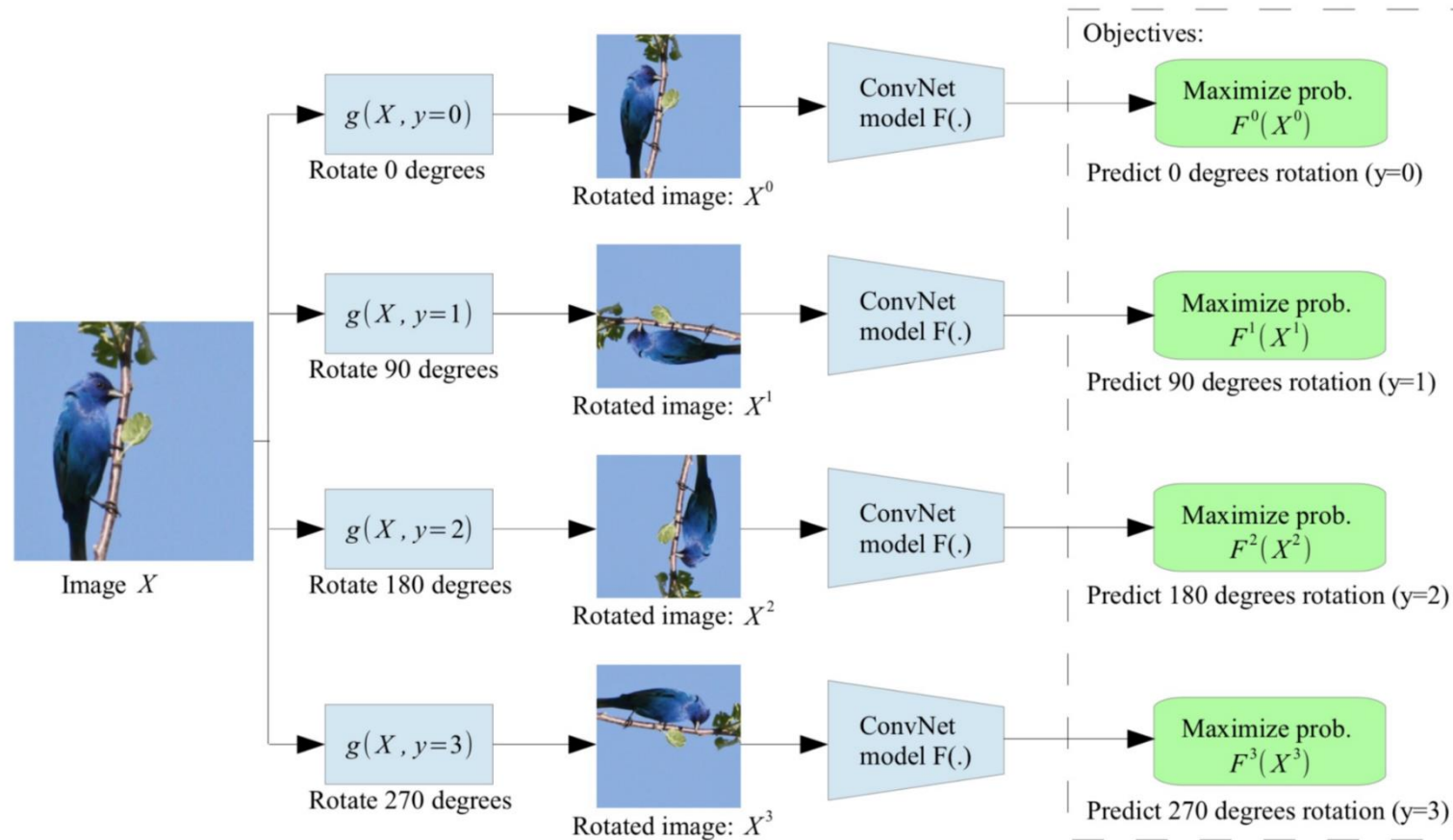


~~0°~~
~~180°~~ rotation

Identifying the object helps solve rotation task!

Catfish species that swims upside down...

Pretext Task: Classify the Rotation



Learning rotation improves results on object classification, object segmentation, and object detection tasks.

[Gidaris et al., ICLR 2018]

Benefits of Self-Supervised Learning

- ✓ Like supervised pretraining, can learn general-purpose feature representations for downstream tasks
- ✓ Reduces expense of hand-labeling large datasets
- ✓ Can leverage nearly unlimited (unlabeled) data available on the web



995 photos uploaded
every second



6000 tweets sent
every second



500 hours of video uploaded
every minute

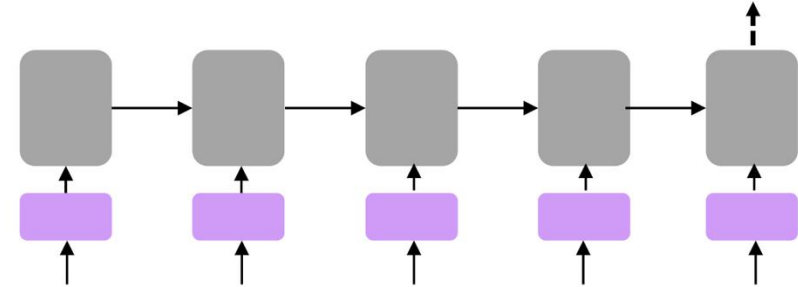
Today's Plan

1. What is self-supervised learning?
2. Examples of self-supervision in NLP
 - Word embeddings (e.g., word2vec)
 - Language models (e.g., GPT)
 - Masked language models (e.g., BERT)
3. Open challenges
 - Demoting bias
 - Capturing factual knowledge
 - Learning symbolic reasoning

Examples of Self-Supervision in NLP

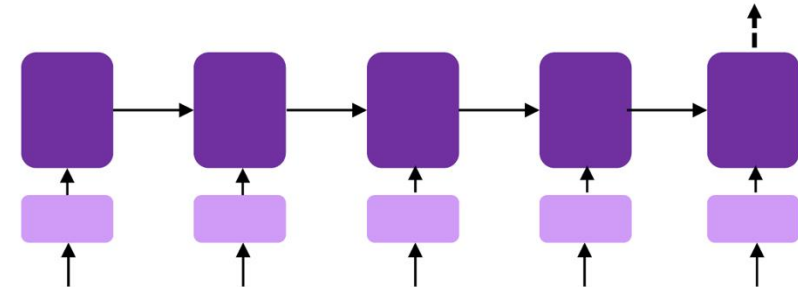
- **Word embeddings**

- Pretrained word representations
- Initializes *1st layer* of downstream models



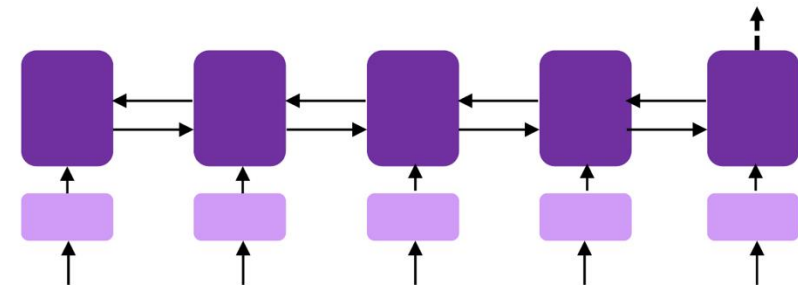
- **Language models**

- *Unidirectional*, pretrained language representations
- Initializes *full* downstream model



- **Masked language models**

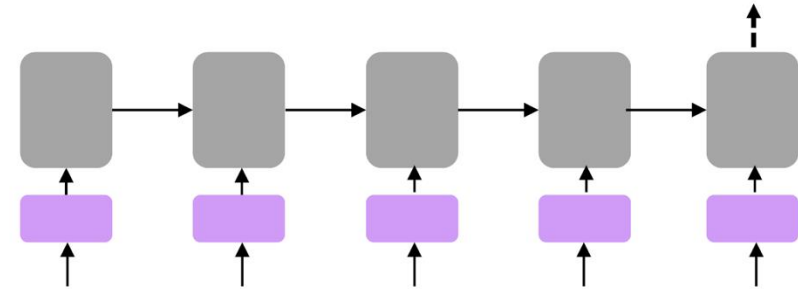
- *Bidirectional*, pretrained language representations
- Initializes *full* downstream model



Examples of Self-Supervision in NLP

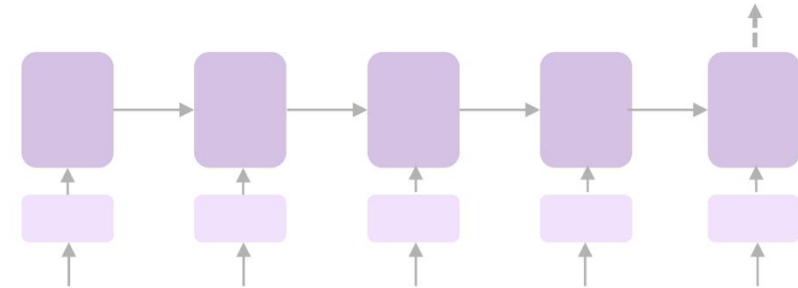
- **Word embeddings**

- Pretrained word representations
- Initializes *1st layer* of downstream models



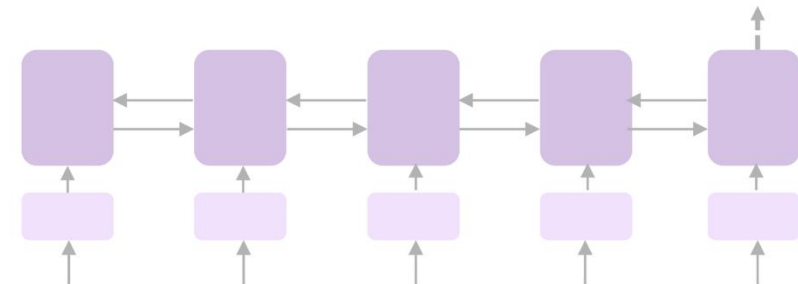
- **Language models**

- *Unidirectional*, pretrained language representations
- Initializes *full* downstream model



- **Masked language models**

- *Bidirectional*, pretrained language representations
- Initializes *full* downstream model



Word Embeddings

- Goal: represent words as vectors for input into neural networks.

- One-hot vectors? (single 1, rest 0s)

pizza = [0 0 0 0 0 1 0 ... 0 0 0 0 0]

pie = [0 0 0 0 0 0 0 ... 0 0 0 1 0]

😞 Millions of words \longrightarrow high-dimensional, sparse vectors

😞 No notion of word similarity

- Instead: we want a **dense, low-dimensional** vector for each word such that words with similar meanings have similar vectors.

Distributional Semantics

- Idea: define a word by the words that frequently occur nearby in a corpus of text
 - “You shall know a word by the company it keeps” (J. R. Firth 1957: 11)
- Example: defining “pizza”
 - What words frequently occur in the context of pizza?

13% of the United States population **eats** **pizza** on any given day.

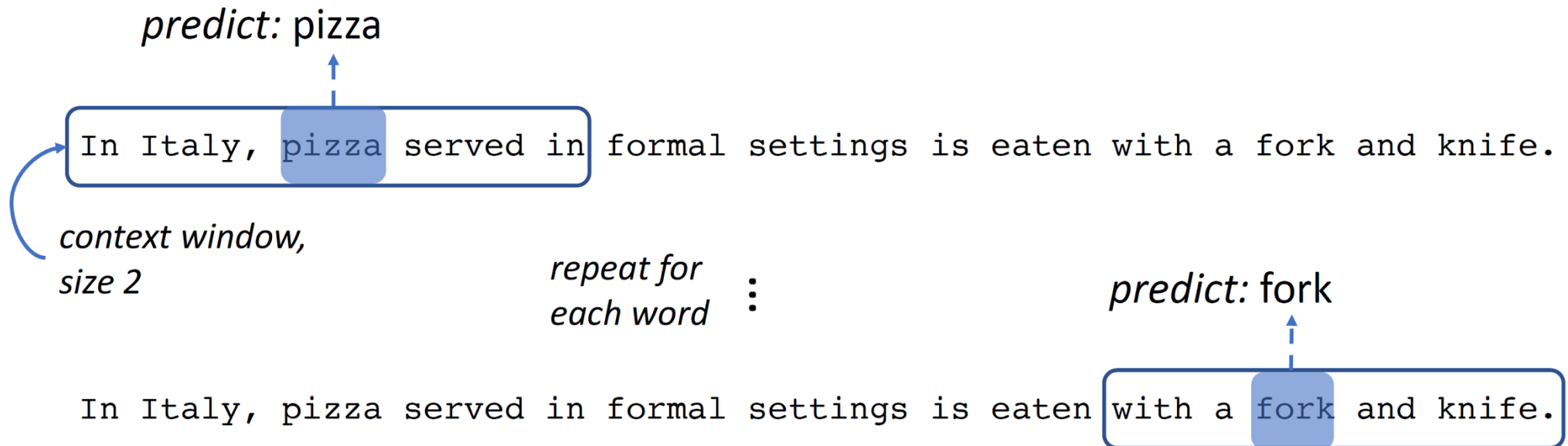
Mozzarella is commonly used on **pizza**, with the highest quality **mozzarella** from Naples.

In **Italy**, **pizza** served in formal settings is **eaten** with a fork and knife.

- Can we use distributional semantics to develop a pretext task for self-supervision?

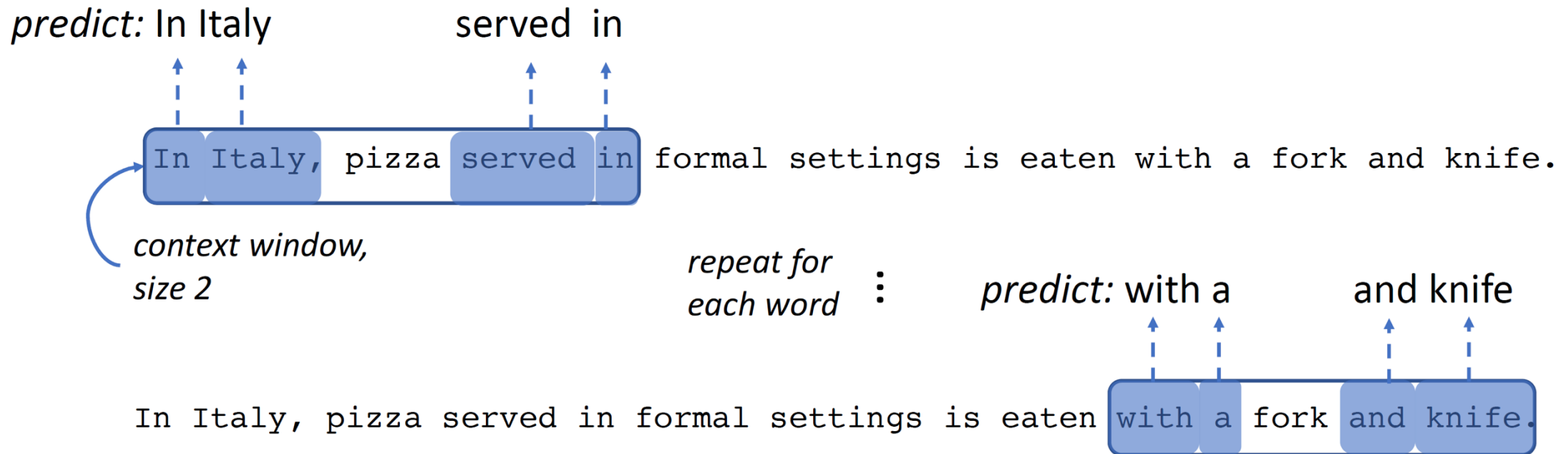
Pretext Task: Predict the Center Word

- Move context window across text data and use words in window to predict the center word.
 - No hand-labeled data is used!



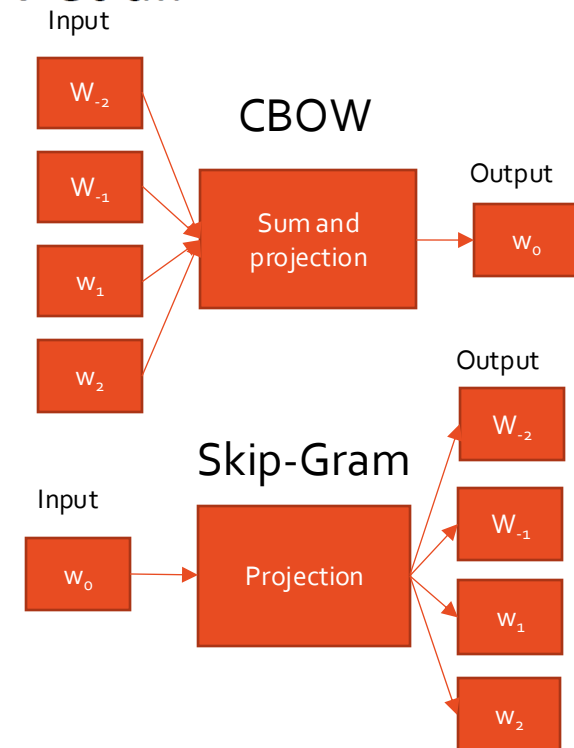
Pretext Task: Predict the Context Words

- Move context window across text data and use words in window to predict the *context* words, given the center word.
 - No hand-labeled data is used!



Case Study: word2vec

- Tool to produce word embeddings using self-supervision by Mikolov et al.
- Supports training word embeddings using 2 architectures:
 - Continuous bag-of-words (CBOW): predict the center word
 - Skip-gram: predict the context words
- Steps:
 1. Start with randomly initialized word embeddings.
 2. Move sliding window across *unlabeled* text data.
 3. Compute probabilities of center/context words, given the words in the window.
 4. Iteratively update word embeddings via stochastic gradient descent .



Case Study: word2vec

- **Loss function (skip-gram):** For a corpus with T words, minimize the negative log likelihood of the context word w_{t+j} given the center word w_t .

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} \mid w_t; \theta)$$

Context word Center word
Model parameters
Context window size

- Use two word embedding matrices (embedding dimension n , vocab size l):
 - Center word embeddings $V \in \mathbb{R}^{n \times l}$; context word embeddings $U \in \mathbb{R}^{l \times n}$

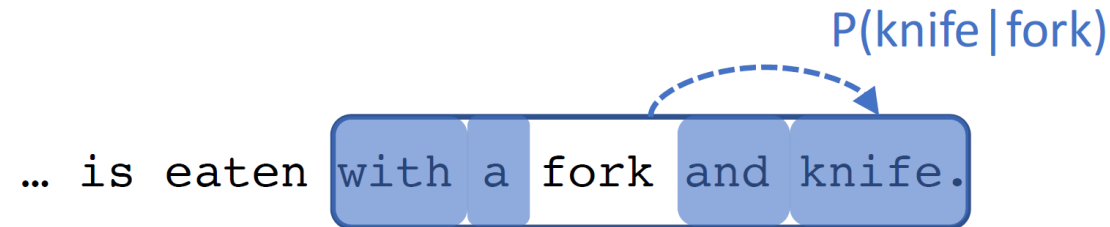
$$P(w_{t+j} \mid w_t; \theta) = P(u_{t+j} \mid v_t) = \frac{\exp(u_{t+j}^T v_t)}{\sum_{j=1}^l \exp(u_j^T v_t)} \quad \text{Softmax}$$

Word vectors

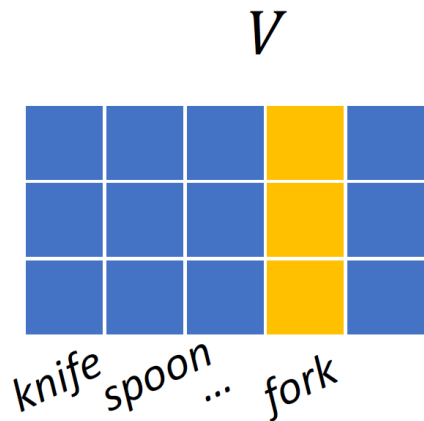
[Mikolov et al., 2013]

Case Study: word2vec

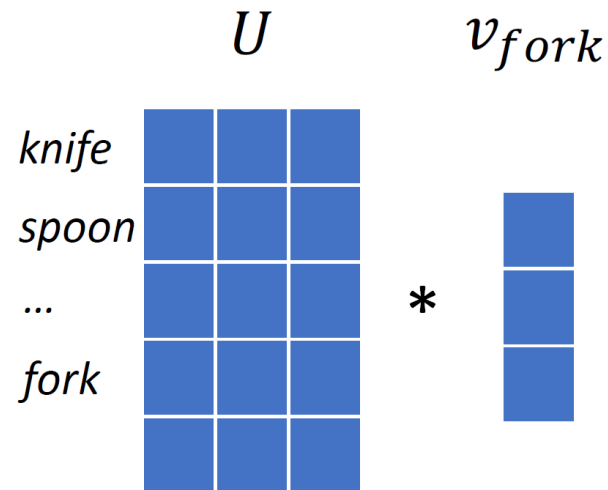
- **Example:** using the skip-gram method (predict context words), compute the probability of "knife" given the center word "fork".



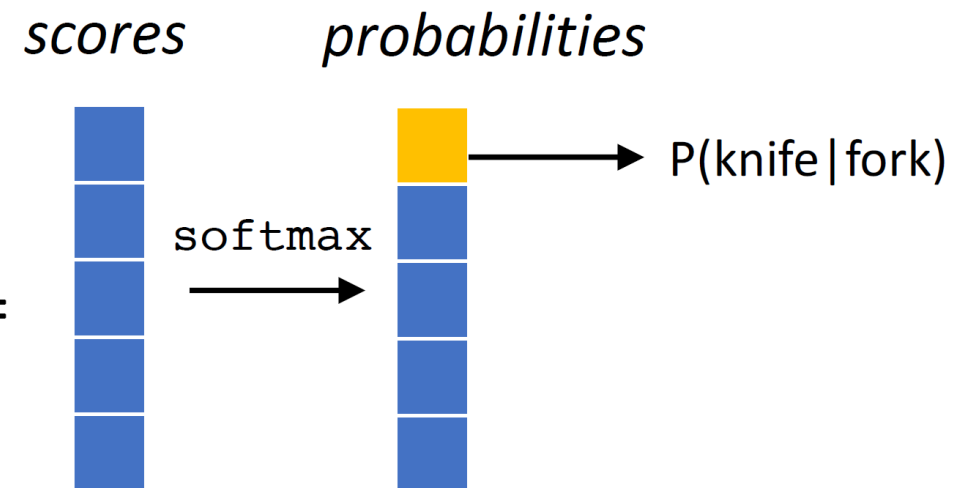
1. Get "fork" word vector v_{fork}



2. Compute scores

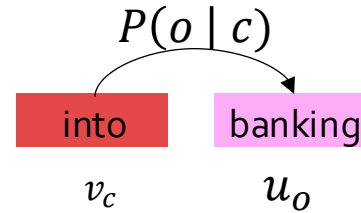


3. Convert to probabilities



Skip-gram with Negative Sampling

- Let's see where the complexity is:



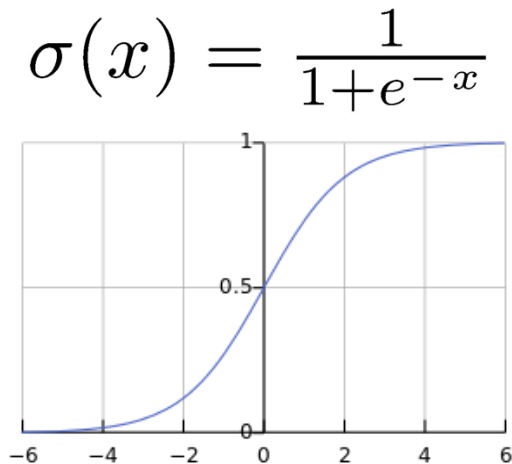
The expensive computation: $O(|V|.d)$

$$\log P(o|c) = \log \frac{\exp(u_o^T v_c)}{\sum_{x \in V} \exp(u_x^T v_c)} = \log \exp(u_o^T v_c) - \log \sum_{x \in V} \exp(u_x^T v_c)$$

- Idea:** rather than enumerating over all vocabulary, let's sample!

$$J_{NS}(\theta) = -\log \sigma(u_o^T v_c) - \sum_{k \in \{K \text{ samples}\}} \log \sigma(-u_x^T v_c)$$

- Maximize** the prob that **outside** word co-occurs w/ the **center**
- Minimize** the prob of **noise/random words** far from the **center** (negatives)



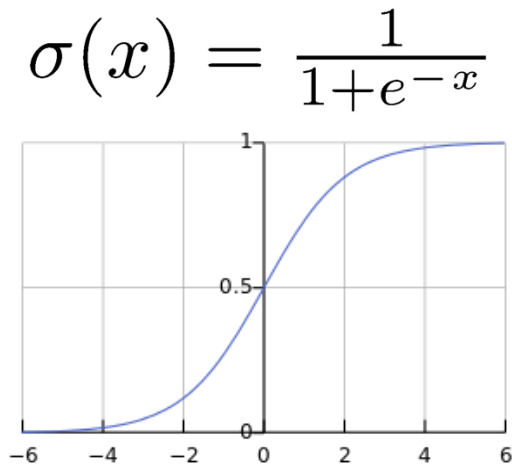
Skip-gram with Negative Sampling

- Have to be careful with sampling negative examples
 - **Challenge:** uniform sampling will sample a lot of stop-words that are very popular.
- Mikolov et al. proposed to sample: $p(w_i) = \frac{f(w_i)^{3/4}}{\sum_j f(w_j)^{3/4}}$
 - Assigns more prob to less frequent words. No theory backing, but works!

- **Idea:** rather than enumerating over all vocabulary, let's sample!

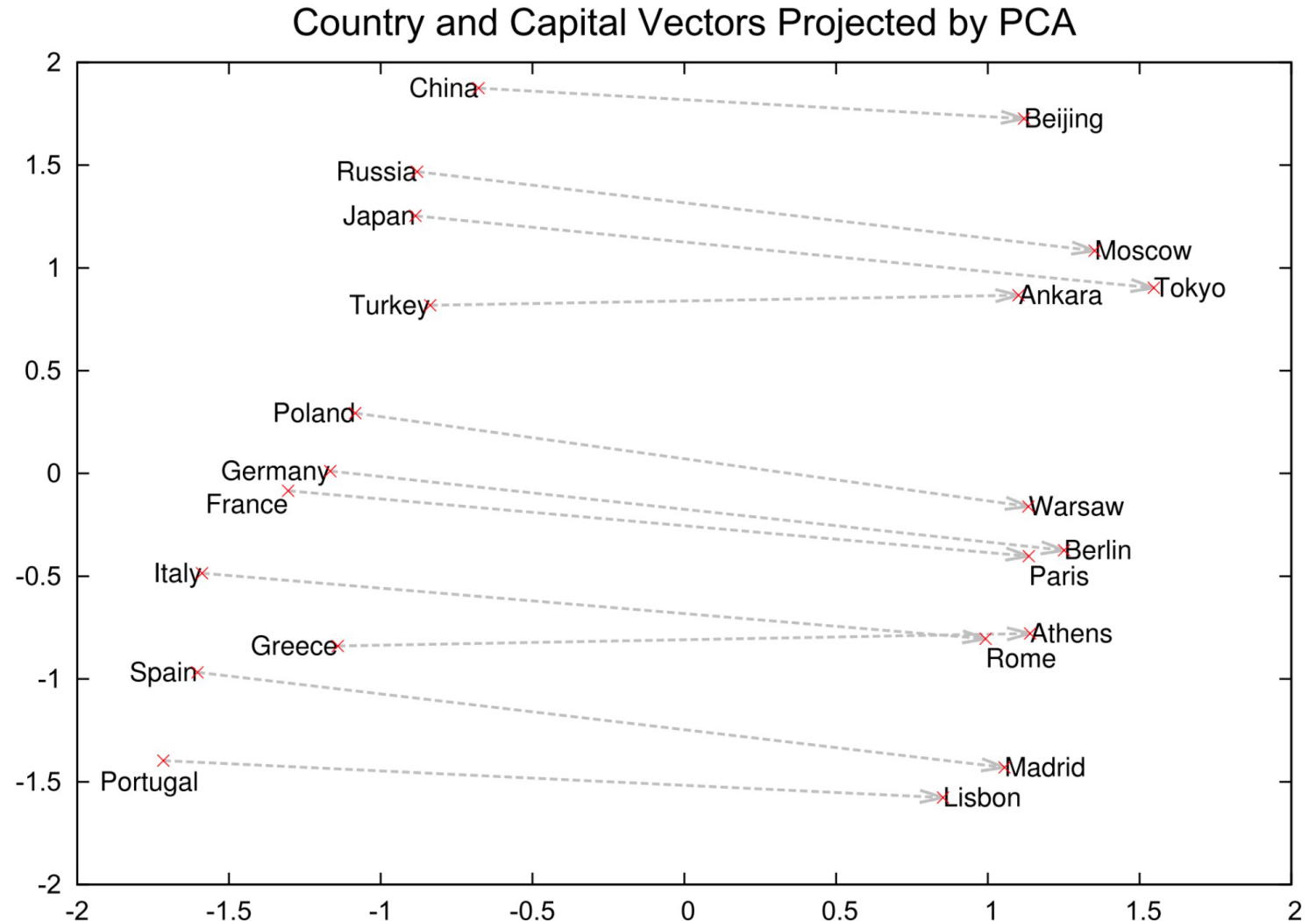
$$J_{NS}(\theta) = -\log \sigma(u_o^T v_c) - \sum_{k \in \{K \text{ samples}\}} \log \sigma(-u_x^T v_c)$$

- **Maximize** the prob that **outside** word co-occurs w/ the **center**
- **Minimize** the prob of **noise/random words** far from the **center** (negatives)



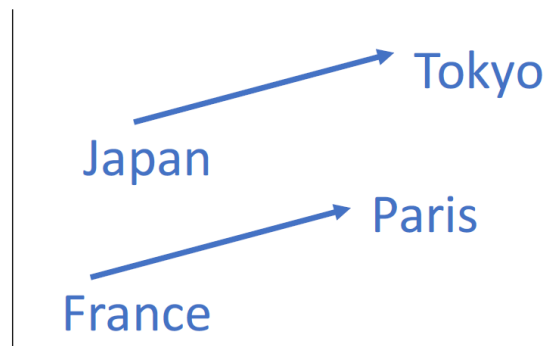
Case Study: word2vec

- Mikolov et al. released word2vec embeddings pretrained on **100 billion word** Google News dataset.
- Embeddings exhibited meaningful properties despite being trained with **no hand-labeled data.**



Case Study: word2vec

- Vector arithmetic can be used to evaluate word embeddings on analogies
 - France is to Paris as Japan is to ?



$$w^* = \underset{w}{\operatorname{argmax}} \frac{v_w \mathbf{y}}{\|v_w\| \|\mathbf{y}\|}$$

Cosine similarity

$$\text{where } \mathbf{y} = v_{\text{Paris}} - v_{\text{France}} + v_{\text{Japan}}$$

$$w^* = \text{Tokyo}$$

Expected answer

- Analogies have become a common **intrinsic task** to evaluate the properties learned by word embeddings

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium/plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

A relation is defined by the vector displacement in the first column. For each start word in the other column, the closest displaced word is shown.

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Mismatch Between Cosine and Dot Product

- **Observation:** there a mismatch between Word2Vec objective and cosine distance!

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{x \in V} \exp(u_x^T v_c)}$$

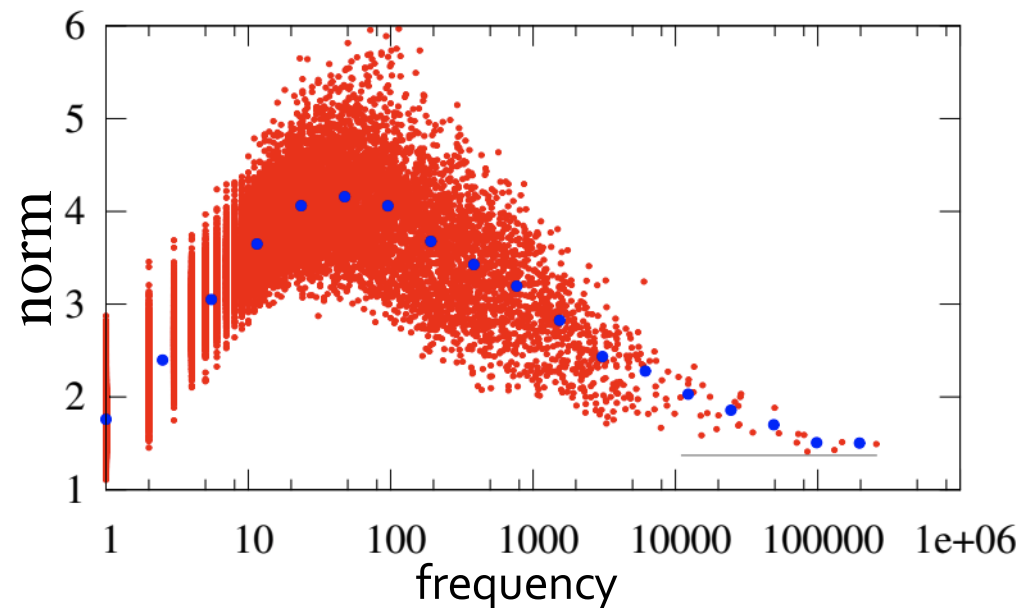
$$\text{distance}(x, y) = \cos(v_x, v_y) = \frac{v_x^T v_y}{\|v_y\| \|v_x\|}$$

1. Why use cosine distance instead of dot product?

- Term frequencies affect the embedding norms.
- Without normalization, frequent terms would seem more similar.

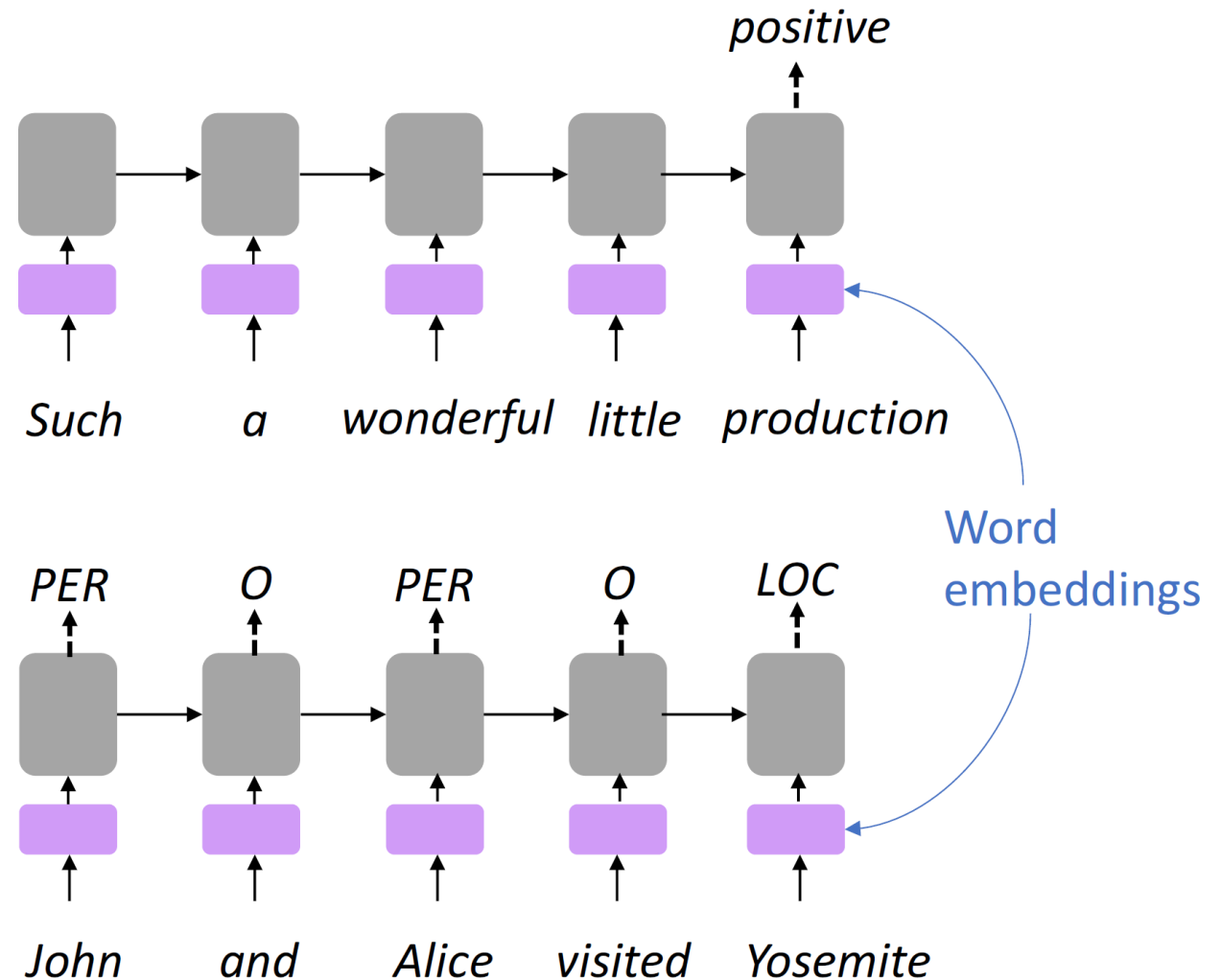
2. Why not change W2V objective to use cos?

- $_ _ (_ _) _ / _$
- It's possible that the resulting vectors would conflate semantic similarity and frequency.



Case Study: word2vec

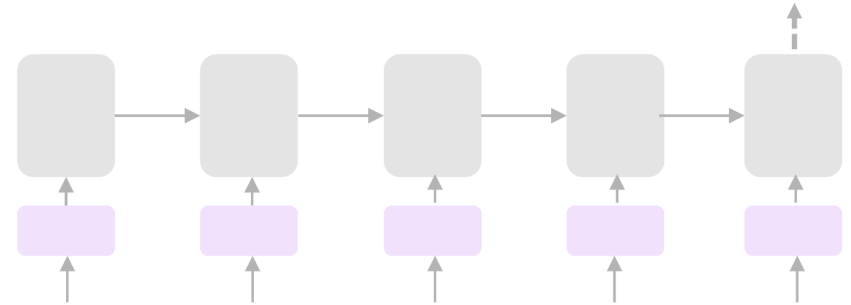
- Pretrained word2vec embeddings can be used to initialize the first layer of downstream models
- Improved performance on many downstream NLP tasks, including sentence classification, machine translation, and sequence tagging
 - Most useful when downstream data is limited
- Still being used in applications in industry today!



Examples of Self-Supervision in NLP

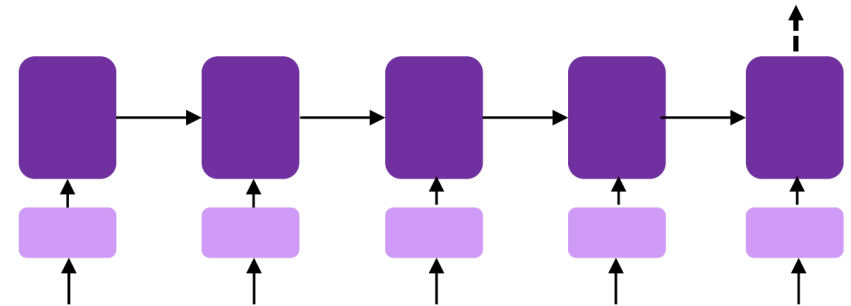
- **Word embeddings**

- Pretrained word representations
- Initializes *1st layer* of downstream models



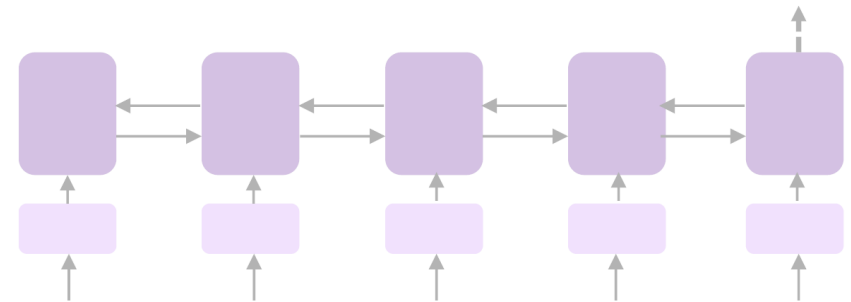
- **Language models**

- *Unidirectional*, pretrained language representations
- Initializes *full* downstream model



- **Masked language models**

- *Bidirectional*, pretrained language representations
- Initializes *full* downstream model



Why weren't word embeddings enough?

- Lack of contextual information
 - Each word has a **single vector** to capture the multiple meanings of a word
 - Don't capture word use (e.g. syntax)

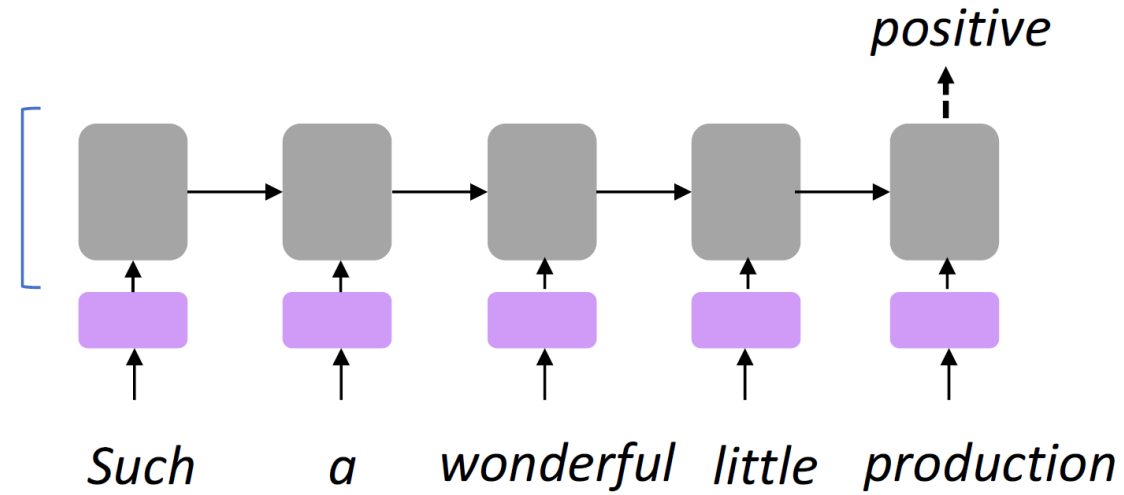


The **ship** is used to **ship** packages.



- Most of the downstream model still needs training
- What self-supervised tasks can we use to pretrain full models for contextual understanding?

Trained from scratch!



The

The cat

The cat sat

The cat sat on

The cat sat on ____?____

The cat sat on the mat.

The cat sat on the mat.


P(mat | The cat sat on the)

next word

context or prefix

Probability of Upcoming Word

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

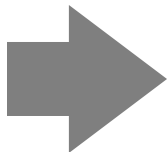

next word context or prefix

LMs as a Marginal Distribution

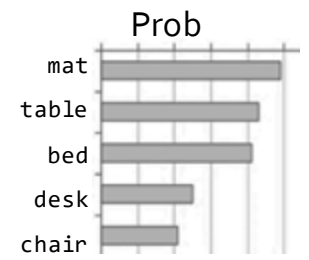
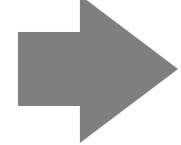
- Directly we train models on “marginals”:

$$P(\underbrace{X_t}_{\text{next word}} \mid \underbrace{X_1, \dots, X_{t-1}}_{\text{context}})$$

“The cat sat on the [MASK]”



Some model



LMs as Implicit Joint Distribution of Language

- Though implicitly we are learning the full distribution over the language:
 - Remember the chain rule: $\mathbf{P}(X_1, \dots, X_t) = P(X_1) \prod_{i=1}^t P(X_i | X_1, X_2, \dots, X_i)$
- **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”
- 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.

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

Diagram illustrating the probability function $\mathbf{P}(X_t \mid X_1, \dots, X_{t-1})$. The term X_t is labeled "next word" and the sequence X_1, \dots, X_{t-1} is labeled "context".

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

$\mathbf{P}(\text{"I like Johns Hopkins University. EOS"}) > \mathbf{P}(\text{"I like John Hopkins University EOS"})$

$\mathbf{P}(\text{"I like Johns Hopkins University. EOS"}) > \mathbf{P}(\text{"University. I Johns EOS Hopkins like"})$

$\mathbf{P}(\text{"JHU is located in Baltimore. EOS"}) > \mathbf{P}(\text{"JHU is located in Virginia. EOS"})$

Doing Things with Language Model (3)

- Can also generate strings

$$\mathbf{P}(X_t \mid \underbrace{X_1, \dots, X_{t-1}}_{\text{context}})$$

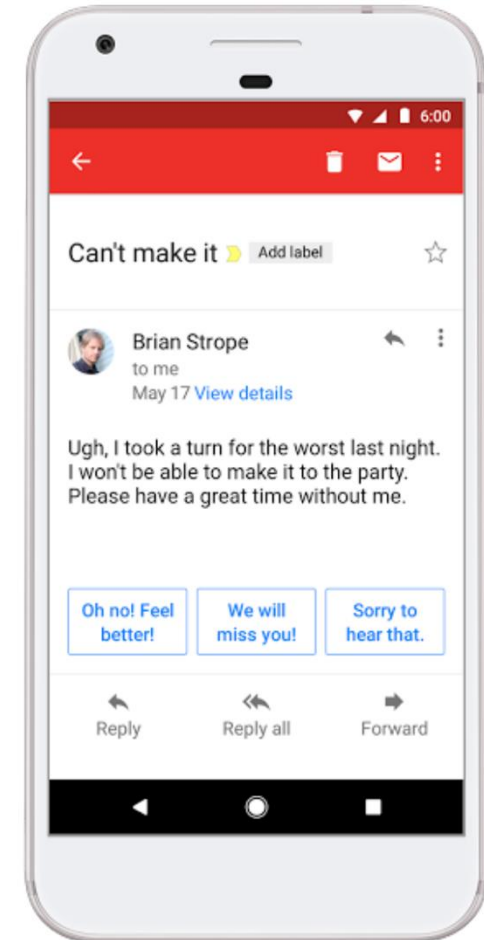
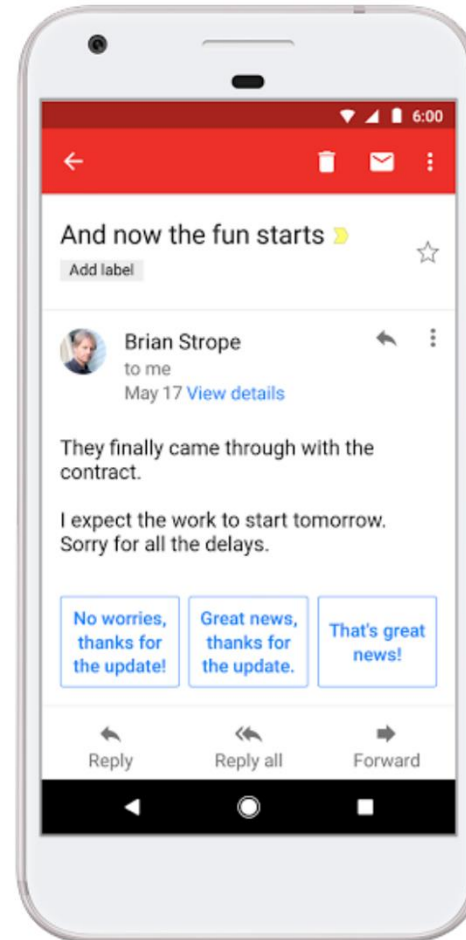
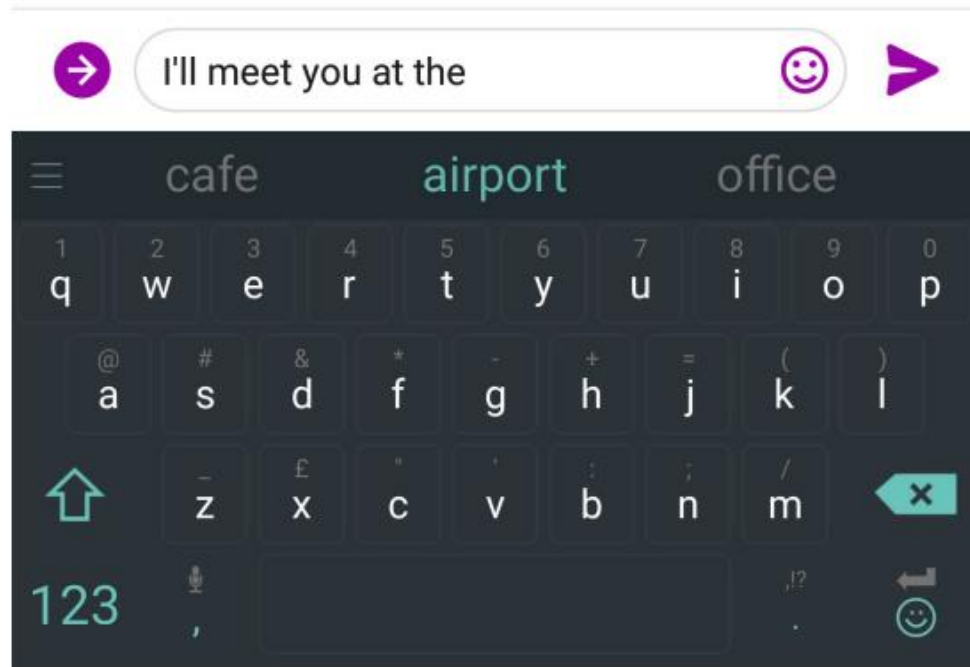
Diagram illustrating the probability function $\mathbf{P}(X_t \mid X_1, \dots, X_{t-1})$. The term X_t is labeled "next word" and the sequence X_1, \dots, X_{t-1} is labeled "context".

- Let's say we start *"Johns Hopkins is "*
- Using this prompt as initial condition, recursively sample from an LM:
 1. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is "}) \rightarrow$ "located"
 2. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located"}) \rightarrow$ "at"
 3. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located at"}) \rightarrow$ "the"
 4. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located at the"}) \rightarrow$ "state"
 5. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located at the state"}) \rightarrow$ "of"
 6. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located at the state of"}) \rightarrow$ "Maryland"
 7. Sample from $\mathbf{P}(X \mid \text{"Johns Hopkins is located at the state of Maryland"}) \rightarrow$ "EOS"

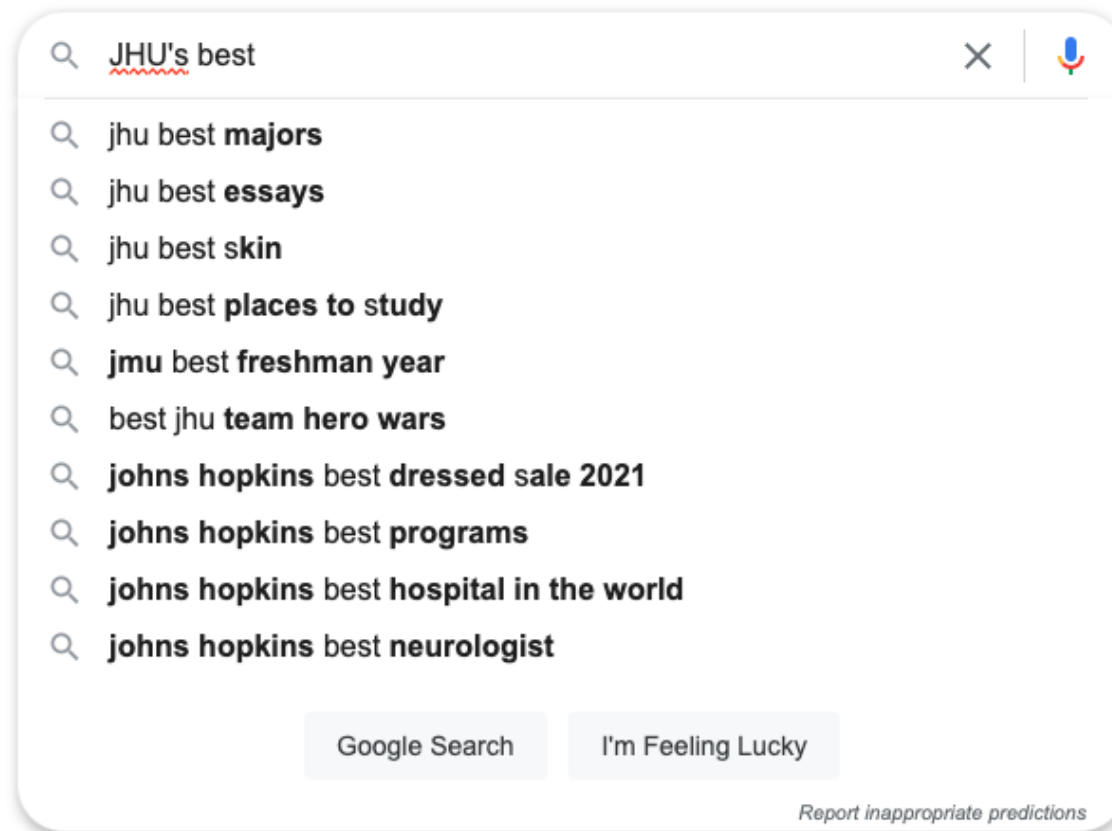
Why Should We Care About Language Modeling?

- Language Modeling is an effective proxy for **language understanding**.
 - **Effective ability to predict forthcoming words** rely on **understanding of context/prefix**
- Language Modeling is a **subcomponent superset** of many NLP tasks, especially those involving **text generation**:
 - Summarization
 - Machine translation
 - Spelling correction
 - Dialogue etc.

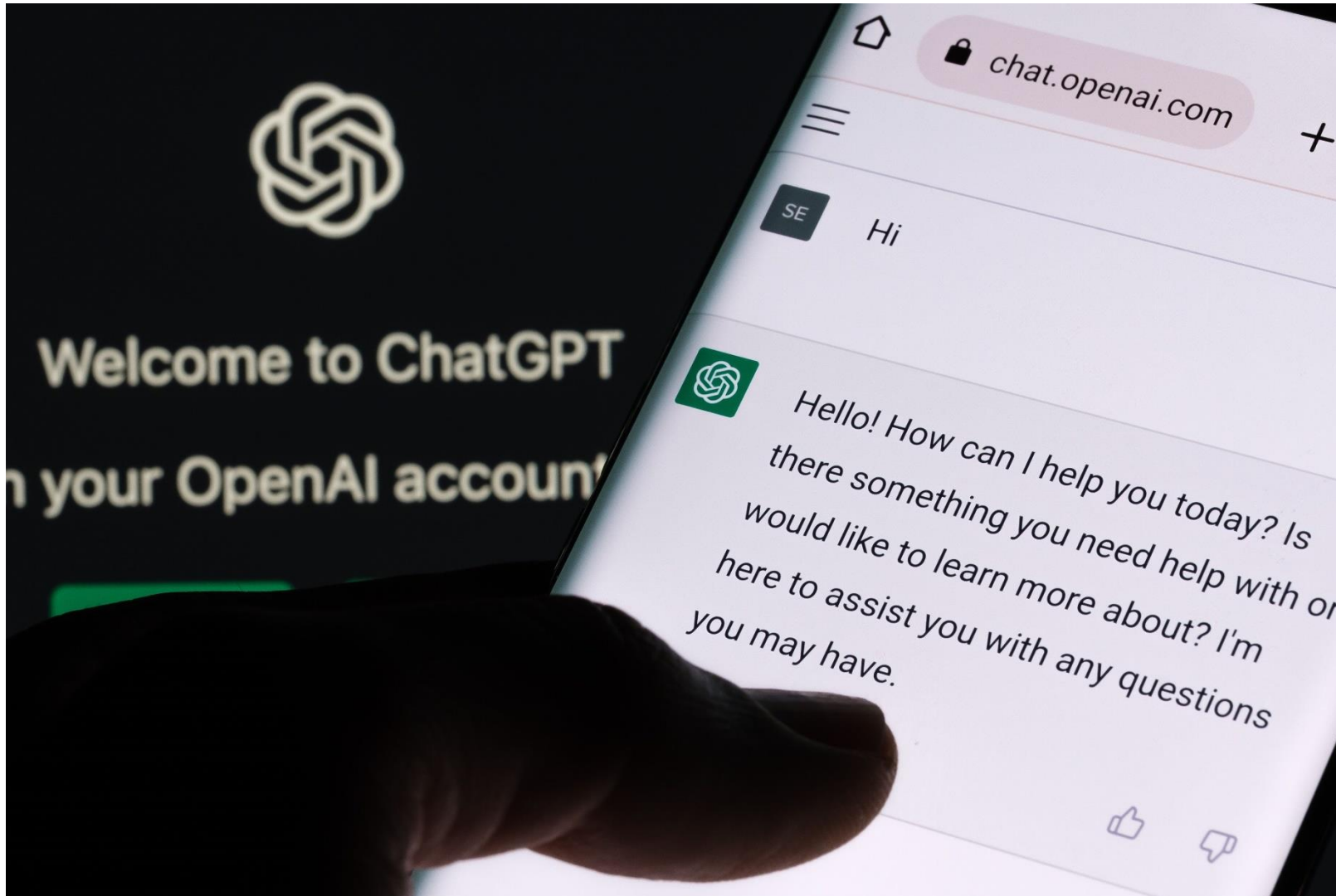
You use Language Models every day!



You use Language Models every day!



You use Language Models every day!



It Can be Misused Too ...

Is this a real
science article?

- A lot more about harms later in the class.

Router: A Methodology for the Typical Unification of Access Points and Redundancy

Jeremy Stribling, Daniel Aguayo and Maxwell Krohn

ABSTRACT

Many physicists would agree that, had it not been for congestion control, the evaluation of web browsers might never have occurred. In fact, few hackers worldwide would disagree with the essential unification of voice-over-IP and public-private key pair. In order to solve this riddle, we confirm that SMPs can be made stochastic, cacheable, and interposable.

I. INTRODUCTION

Many scholars would agree that, had it not been for active networks, the simulation of Lamport clocks might never have occurred. The notion that end-users synchronize with the investigation of Markov models is rarely outdated. A theoretical grand challenge in theory is the important unification of virtual machines and real-time theory. To what extent can web browsers be constructed to achieve this purpose?

Certainly, the usual methods for the emulation of Smalltalk that paved the way for the investigation of rasterization do not apply in this area. In the opinions of many, despite the fact that conventional wisdom states that this grand challenge is continuously answered by the study of access points, we believe that a different solution is necessary. It should be noted that Router runs in $\Omega(\log \log n)$ time. Certainly, the shortcoming of this type of solution, however, is that compilers and superpages are mostly incompatible. Despite the fact that similar methodologies visualize XML, we surmount this issue without synthesizing distributed archetypes.

The rest of this paper is organized as follows. For starters, we motivate the need for fiber-optic cables. We place our work in context with the prior work in this area. To address this obstacle, we disprove that even though the much-touted autonomous algorithm for the construction of digital-to-analog converters by Jones [10] is NP-complete, object-oriented languages can be made signed, decentralized, and signed. Along these same lines, to accomplish this mission, we concentrate our efforts on showing that the famous ubiquitous algorithm for the exploration of robots by Sato et al. runs in $\Omega((n + \log n))$ time [22]. In the end, we conclude.

II. ARCHITECTURE

Our research is principled. Consider the early methodology by Martin and Smith; our model is similar, but will actually overcome this grand challenge. Despite the fact that such a claim at first glance seems unexpected, it is buffeted by previous work in the field. Any significant development of secure theory will clearly require that the acclaimed real-time algorithm for the refinement of write-ahead logging by Edward Feigenbaum et al. [15] is impossible; our application is no different. This may or may not actually hold in reality. We consider an application consisting of n access points. Next, the model for our heuristic consists of four independent components: simulated annealing, active networks, flexible modalities, and the study of reinforcement learning.

We consider an algorithm consisting of n semaphores. Any unproven synthesis of introspective methodologies will

<https://pdos.csail.mit.edu/archive/scigen/>

Language Models: A History

- Shannon (1950): The predictive difficulty (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.



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



Andrey Markov

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

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

1st order approximation:

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

1 element
└──────────┘



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



Andrey Markov

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

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

2nd order approximation:

$$\mathbf{P}(\text{mat} | \text{the cat sat on the}) \approx \mathbf{P}(\text{mat} | \underbrace{\text{on the}}_{\text{2 elements}})$$



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



Andrey Markov

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

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

3rd order approximation:

$$\mathbf{P}(\text{mat} | \text{the cat sat on the}) \approx \mathbf{P}(\text{mat} | \underbrace{\text{sat on the}}_{\text{3 elements}})$$



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



Andrey Markov

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

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

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

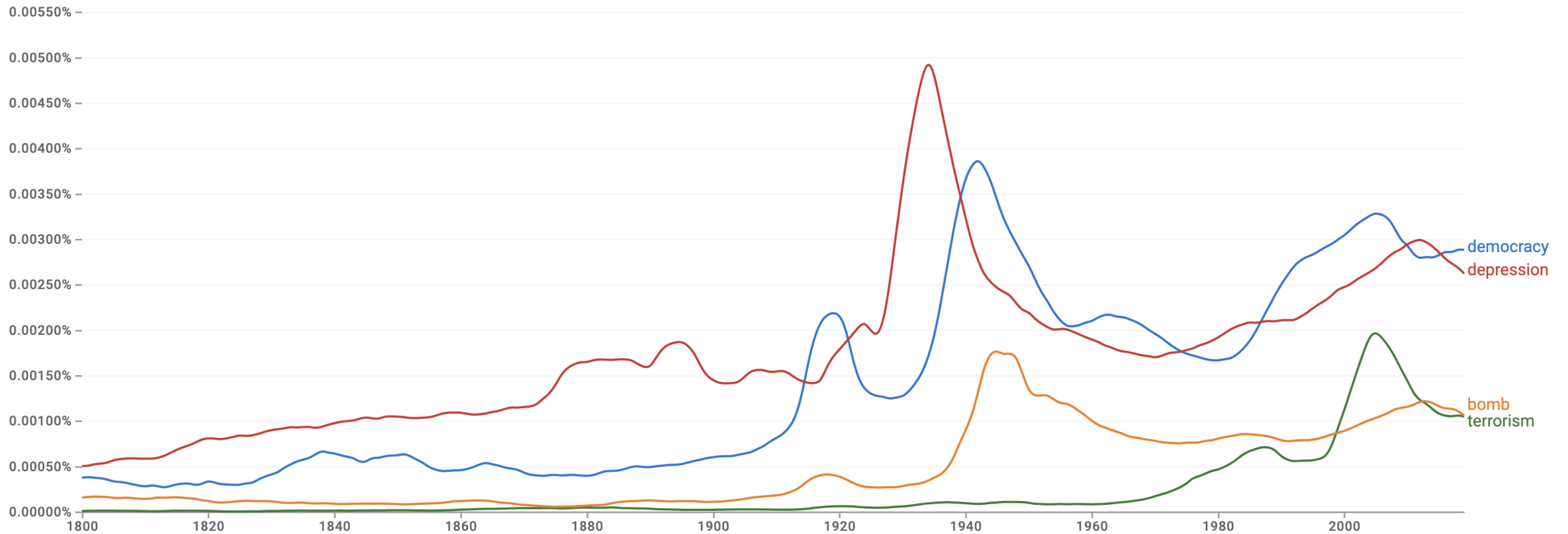
N-gram Language Models

- **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:

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

Pre-Computed N-Grams

Google Books Ngram Viewer

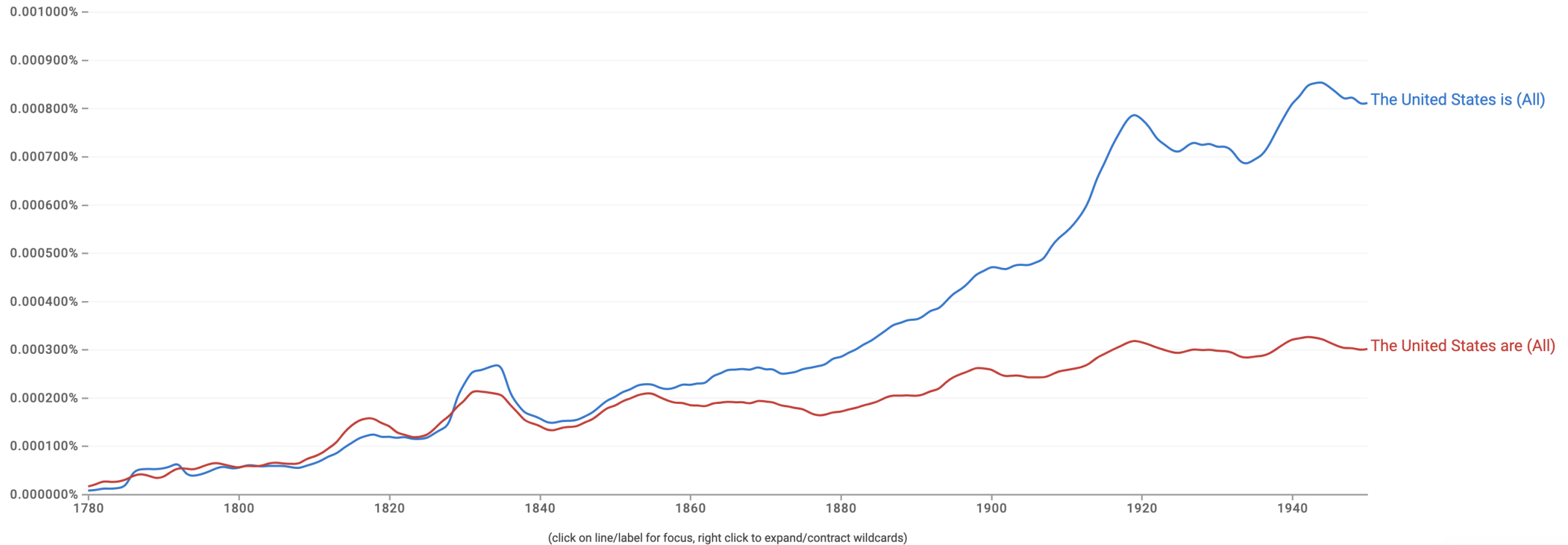


Google n-gram viewer <https://books.google.com/ngrams/>

Data: <http://storage.googleapis.com/books/ngrams/books/datasetv2.html>

Pre-Computed N-Grams

Google Books Ngram Viewer



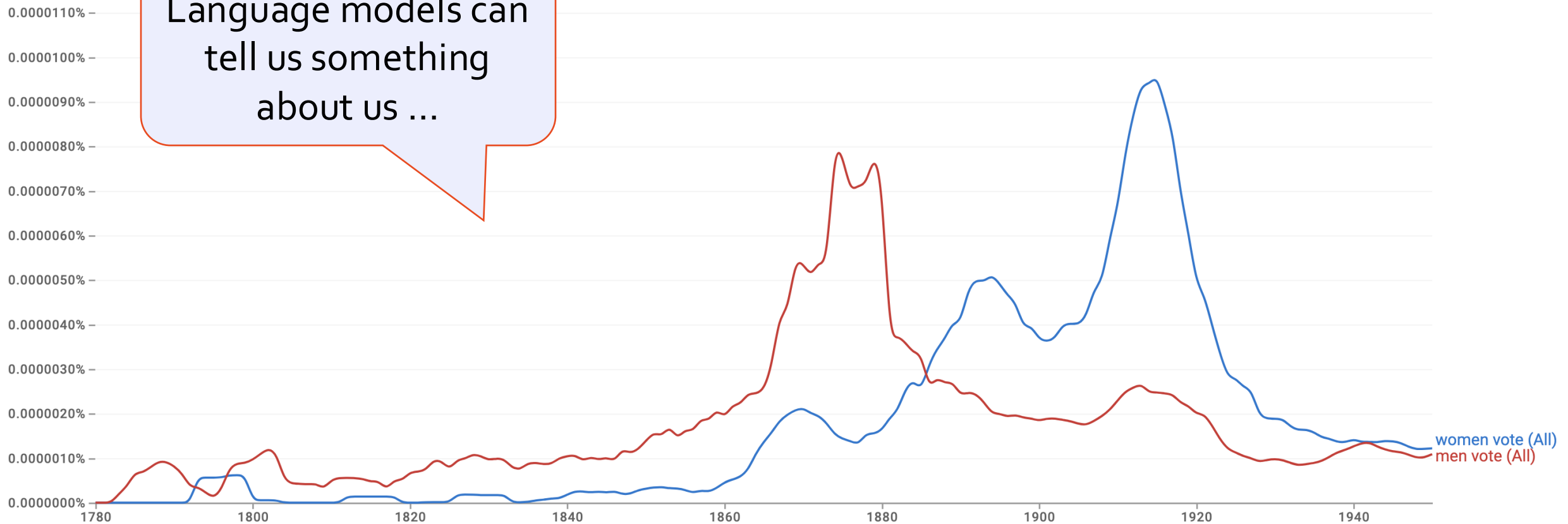
Google n-gram viewer <https://books.google.com/ngrams/>

Data: <http://storage.googleapis.com/books/ngrams/books/datasetsv2.html>

Pre-Computed N-Grams

Google Books Ngram Viewer

Language models can tell us something about us ...

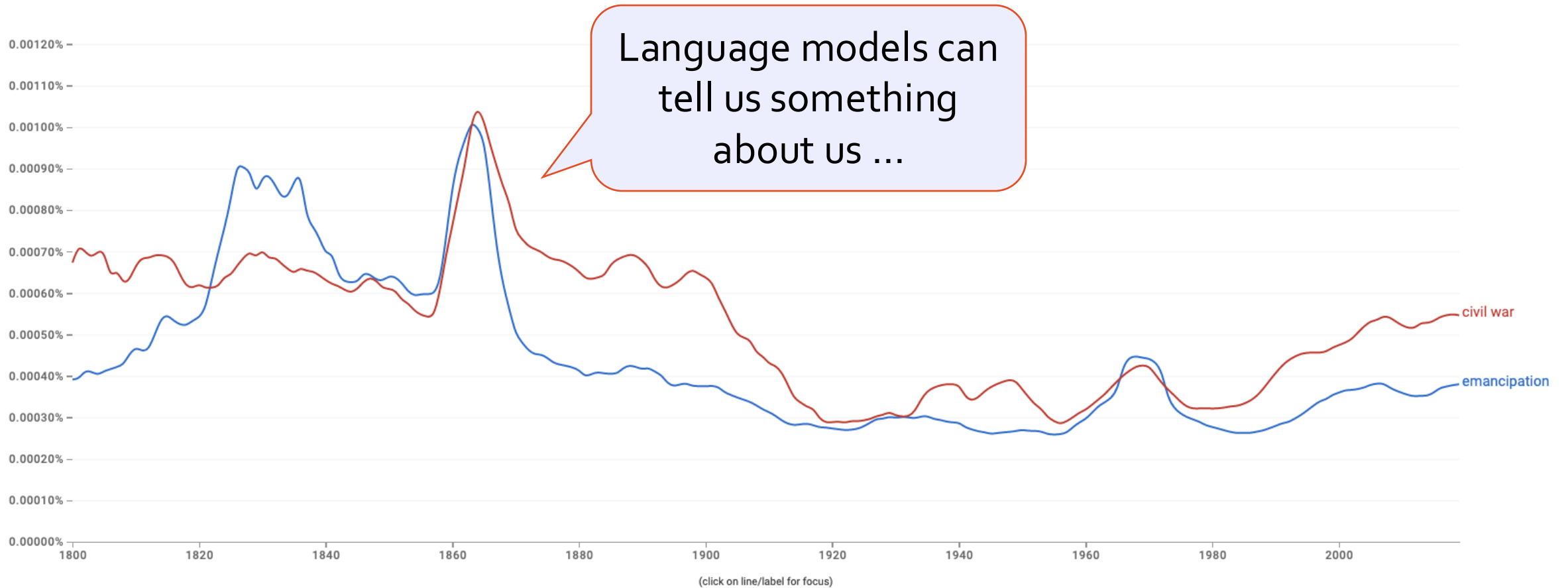


Google n-gram viewer <https://books.google.com/ngrams/>

Data: <http://storage.googleapis.com/books/ngrams/books/datasetv2.html>

Pre-Computed N-Grams

Google Books Ngram Viewer



Google n-gram viewer <https://books.google.com/ngrams/>

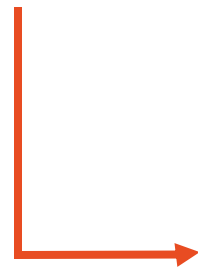
Data: <http://storage.googleapis.com/books/ngrams/books/datasetv2.html>

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!

* Try for yourself: <https://nlpforhackers.io/language-models/>

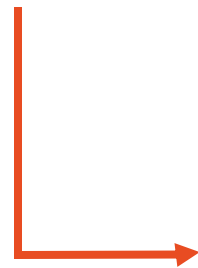
[adopted from Chris Manning]

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:

Condition on this
 today the price _____

get probability
 distribution

of	0.308
for	0.050
it	0.046
to	0.046
is	0.031
...	

Sparsity problem: not
 much granularity in the
 probability distribution

Otherwise, seems reasonable!

Generation from N-Gram Models

- Now we can sample from this model:

today the price of _

get probability
distribution

the	0.072
18	0.043
oil	0.043
its	0.036
gold	0.018
...	

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!**

Why is language modeling a good pretext task?

- ✓ Captures aspects of language useful for downstream tasks, including long-term dependencies, syntactic structure, and sentiment
- ✓ Lots of available data (especially in high-resource languages, e.g. English)
- ✓ Already a key component of many downstream tasks (e.g. machine translation)

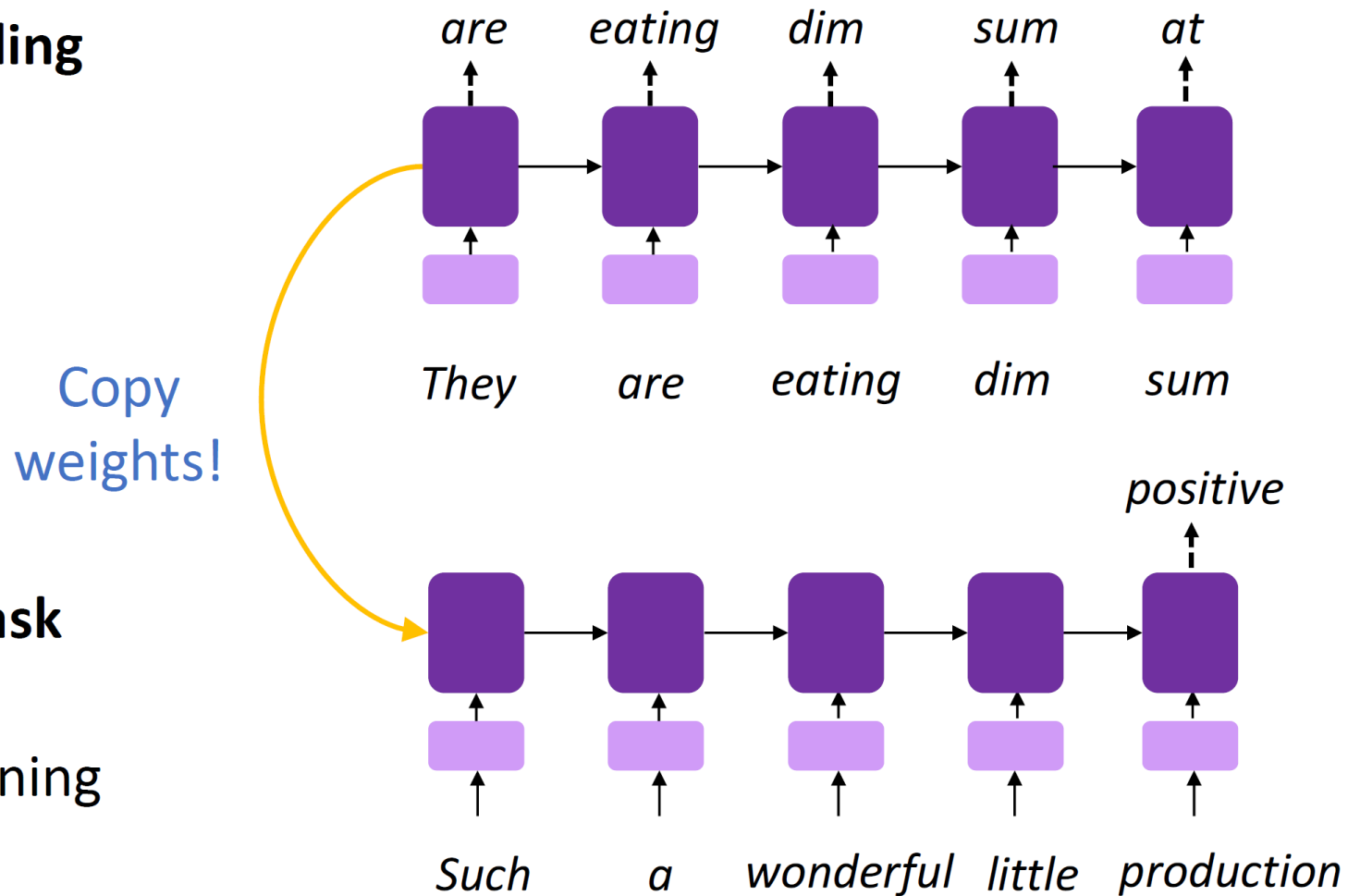
Using language modeling for pretraining

1. Pretrain on language modeling (pretext task)

- Self-supervised learning
- Large, unlabeled datasets

2. Finetune on downstream task (e.g. sentiment analysis)

- Supervised learning for finetuning
- Small, hand-labeled datasets

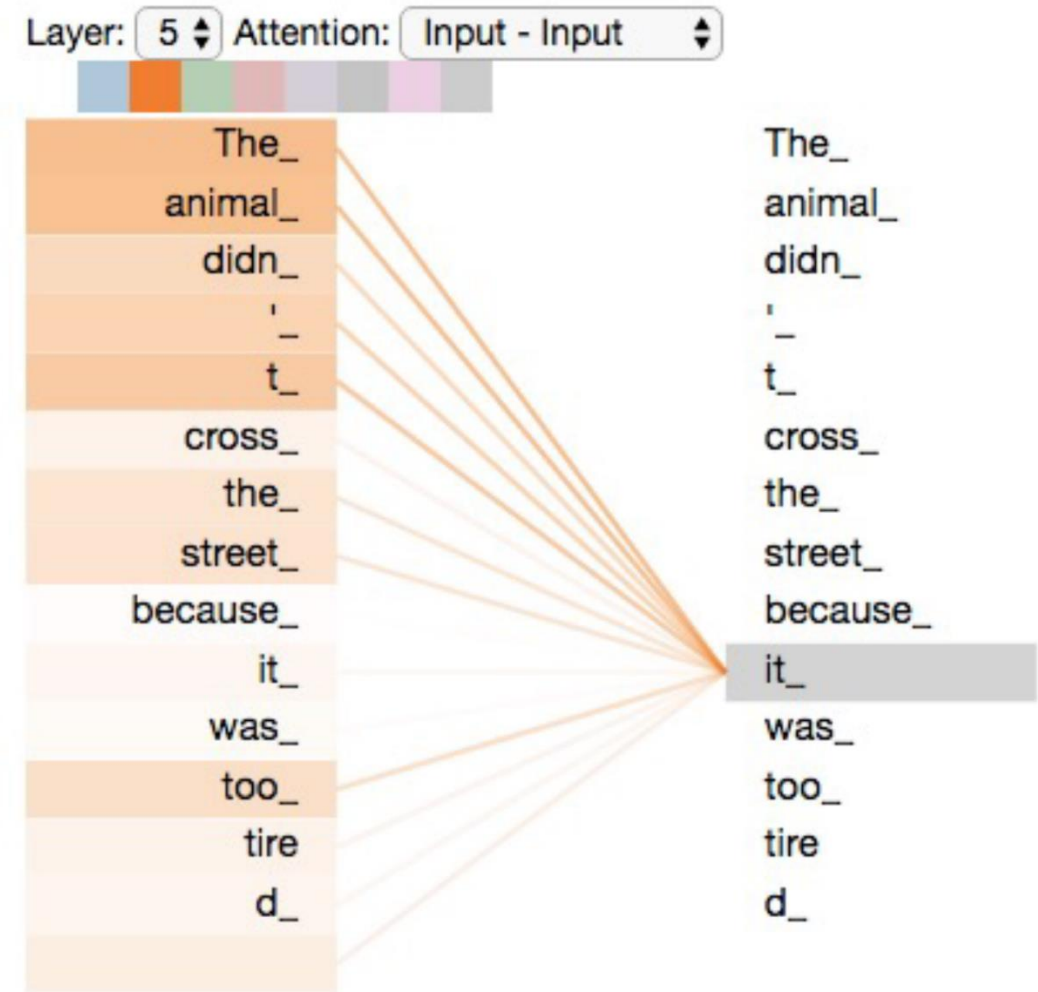


Case Study: Generative Pretrained Transformer (GPT)

- Introduced by Radford et al. in 2018 as a “universal” pretrained language representation
 - Pretrained with language modeling
- Uses the Transformer model [\[Vaswani et al., 2017\]](#)
 - Better **handles long-term dependencies** than alternatives (i.e. recurrent neural networks like LSTMs) and **more efficient on current hardware**
- Has since had follow-on work with GPT-2 and GPT-3 resulting in even larger pretrained models

Quick Aside: Basics of Transformers

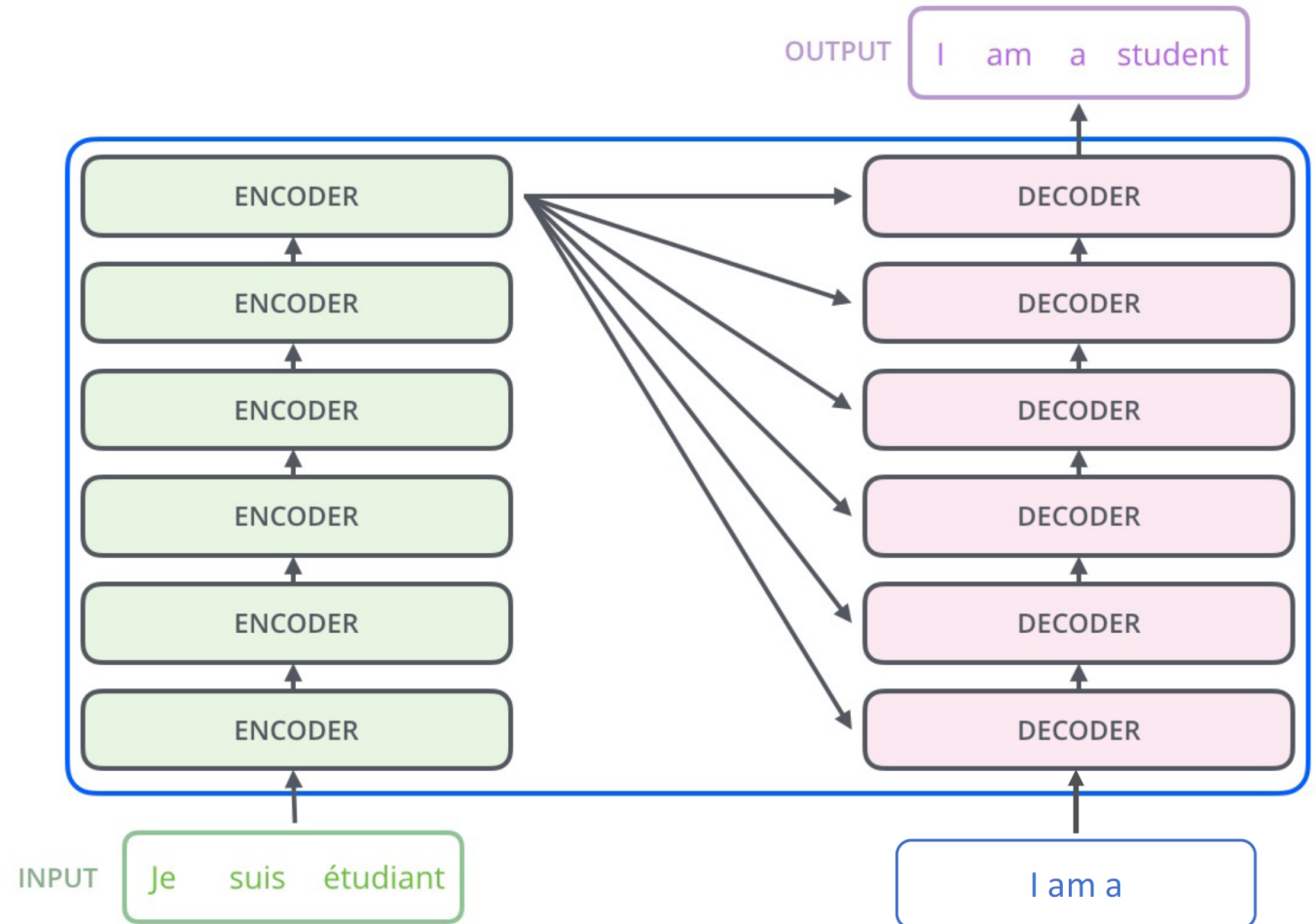
- Model architecture that has recently replaced recurrent neural networks (e.g. LSTMS) as the building block in many NLP pipelines
- Uses **self-attention** to pay attention to relevant words in the sequence (“Attention is all you need”)
 - Can attend to words that are far away



[Alammar et al., Illustrated Transformer]

Quick Aside: Basics of Transformers

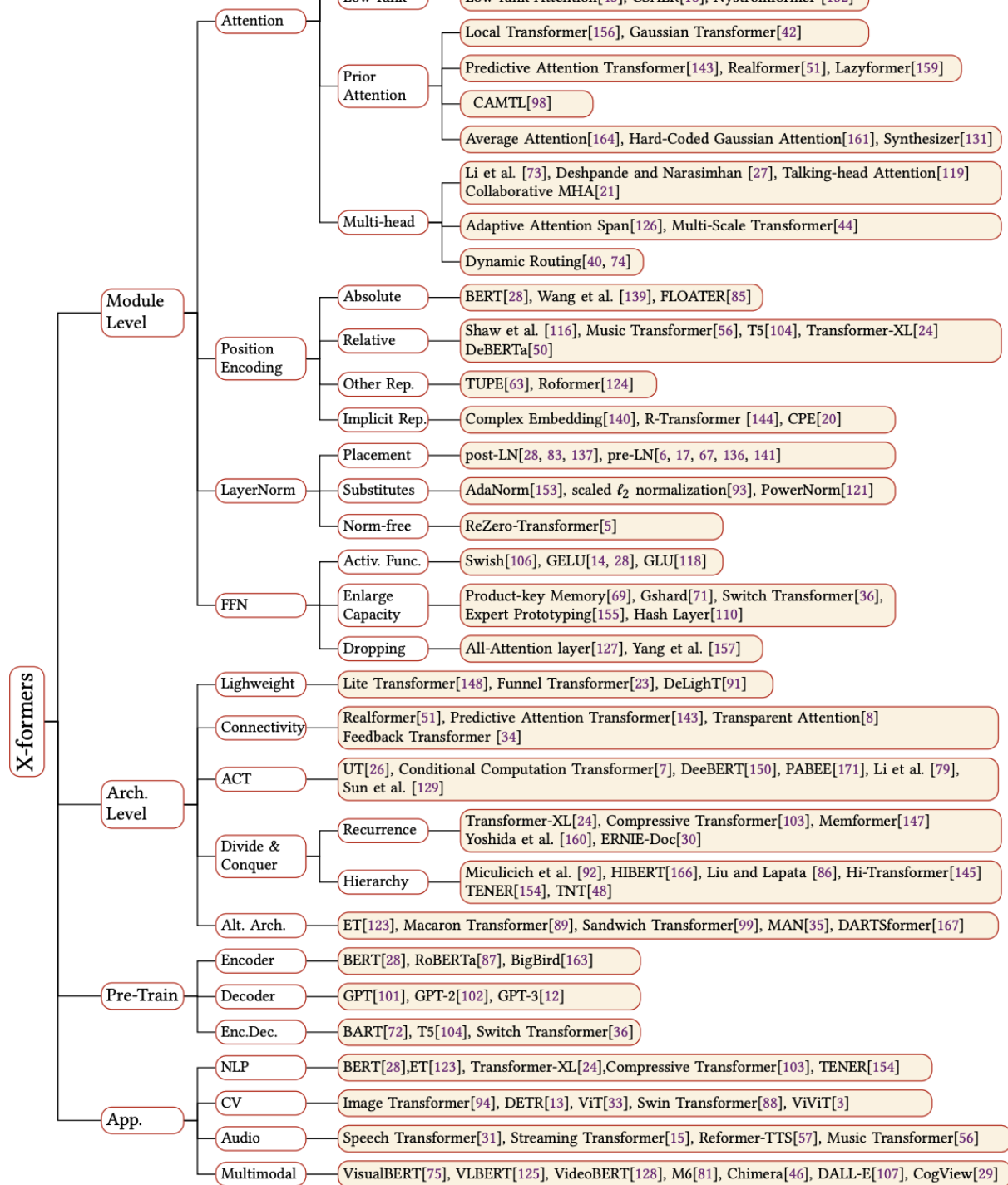
- Composed of two modules:
 - **Encoder** to learn representations of the input
 - **Decoder** to generate output conditioned on the encoder output and the previous decoder output (auto-regressive)
- Each block contains a self-attention and feedforward layer



[Alammar et al., Illustrated Transformer]

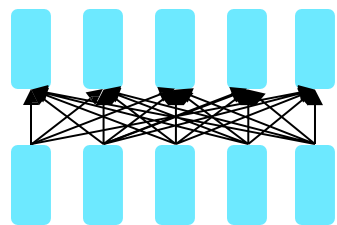
After Transformer ...





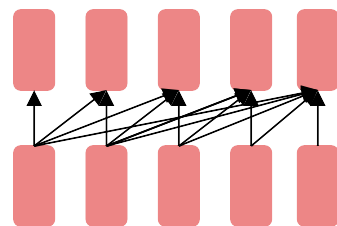
Impact of Transformers

- A building block for a variety of LMs



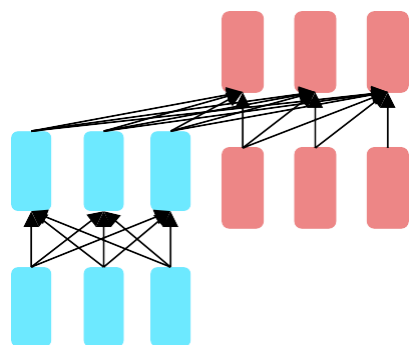
Encoders

- ❖ 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?

Case Study: Generative Pretrained Transformer (GPT)

- Pretrain the **Transformer decoder model** on the language modeling task:

$$L_{LM}(U) = \sum_{i=1}^n \log P(u_i | u_{i-k}, \dots, u_{i-1}; \theta)$$

Text corpus

Context window

Word in a sequence

$$h_{i-k}, \dots, h_{i-1} = \text{decoder}(u_{i-k}, \dots, u_{i-1})$$

$$P(u_i | u_{i-k}, \dots, u_{i-1}) = \text{softmax}(h_{i-1} W_e^T)$$

Previous word hidden representation

Linear layer

Case Study: Generative Pretrained Transformer (GPT)

- Finetune the pretrained Transformer model with a randomly initialized linear layer for **supervised downstream tasks**:

Labeled dataset

Input sequence x , label y

$$L_{downstream}(C) = \sum_{(x, y)} \log P(y | x_1, \dots, x_m)$$
$$h_1, \dots, h_m = \text{decoder}(u_1, \dots, u_m)$$
$$P(y | x_1, \dots, x_m) = \text{softmax}(h_m W_y)$$

Last word's hidden representation

New linear layer, replaces W_e from pretraining

- Linear layer makes up most of the **new** parameters needed for downstream tasks, rest are initialized from pretraining!

Case Study: Generative Pretrained Transformer (GPT)

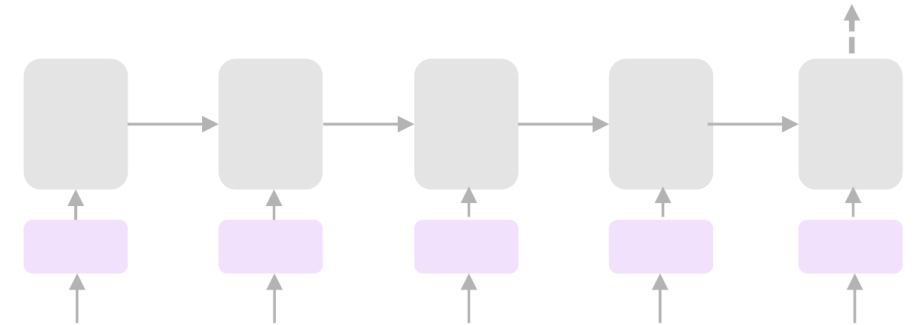
- Pretrained on the BooksCorpus (7000 unique books)
- Achieved state-of-the-art on **downstream** question answering tasks (as well as natural language inference, semantic similarity, and text classification tasks)

Method	<i>select the correct end to the story</i> Story Cloze	<i>middle and high school exam reading comprehension questions</i> RACE-m	RACE-h	RACE
val-LS-skip [55]	76.5	-	-	-
Hidden Coherence Model [7]	<u>77.6</u>	-	-	-
Dynamic Fusion Net [67] (9x)	-	55.6	49.4	51.2
BiAttention MRU [59] (9x)	-	<u>60.2</u>	<u>50.3</u>	<u>53.3</u>
Finetuned Transformer LM (ours)	86.5	62.9	57.4	59.0

Examples of Self-Supervision in NLP

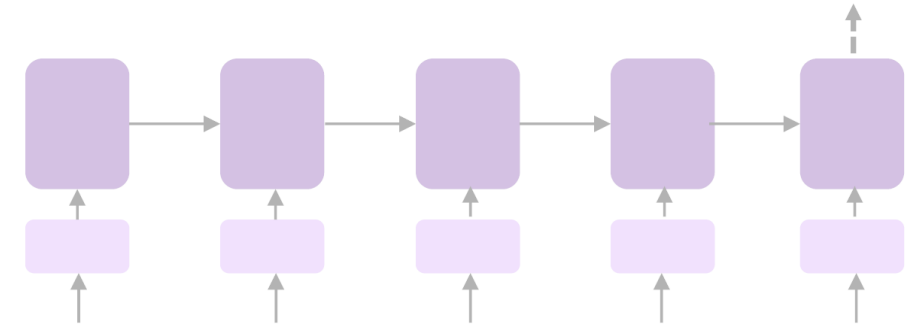
- **Word embeddings**

- Pretrained word representations
- Initializes *1st layer* of downstream models



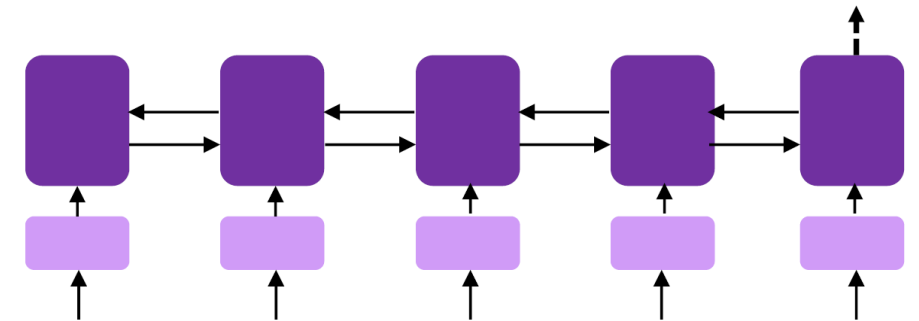
- **Language models**

- *Unidirectional*, pretrained language representations
- Initializes *full* downstream model



- **Masked language models**

- *Bidirectional*, pretrained language representations
- Initializes *full* downstream model



Using context from the future

- Consider predicting the next word for the following example:

He is going to the _____.

movies *park*
store *theater*
library *treehouse*
school *pool*

- What if you have more (bidirectional) context?

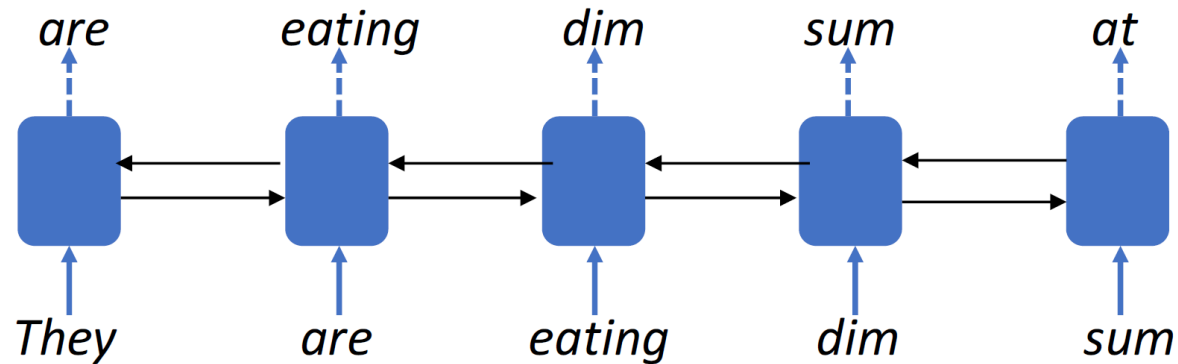
He is going to the _____ to buy some milk.

store
market
Safeway

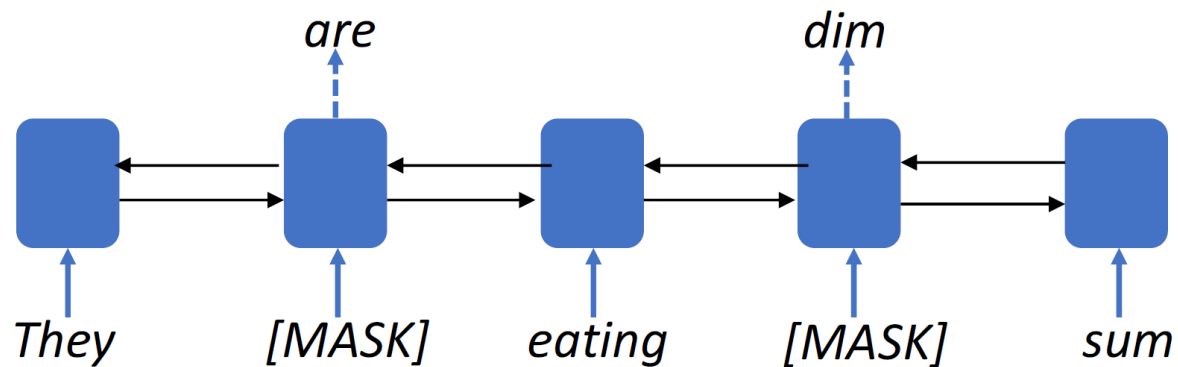
- Information **from the future** can be helpful for language understanding!

Masked language models (MLMs)

- With bidirectional context, if we aren't careful, model can "cheat" and see next word



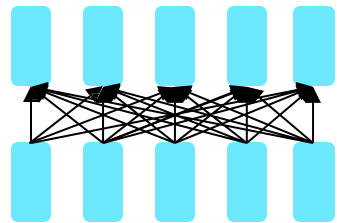
- What if we mask out some words and ask the model to predict them?



This is called *masked language modeling*.



BERT



Encoders

BERT

Bidirectional **E**ncoder **R**epresentations from **T**ransformers



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 (2018)

- Transformer based network to learn representations of language
- Improvements
 - Bi-directional LSTM -> Self-attention
 - Massive data
 - Masked-LM objective

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Abstract

We introduce a new language representation model called **BERT**, which stands for **Bidirectional Encoder Representations from Transformers**. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.

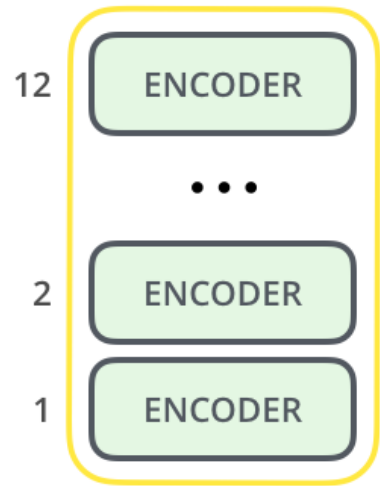
BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute im-

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal task-specific parameters, and is trained on the downstream tasks by simply fine-tuning *all* pre-trained parameters. The two approaches share the same objective function during pre-training, where they use unidirectional language models to learn general language representations.

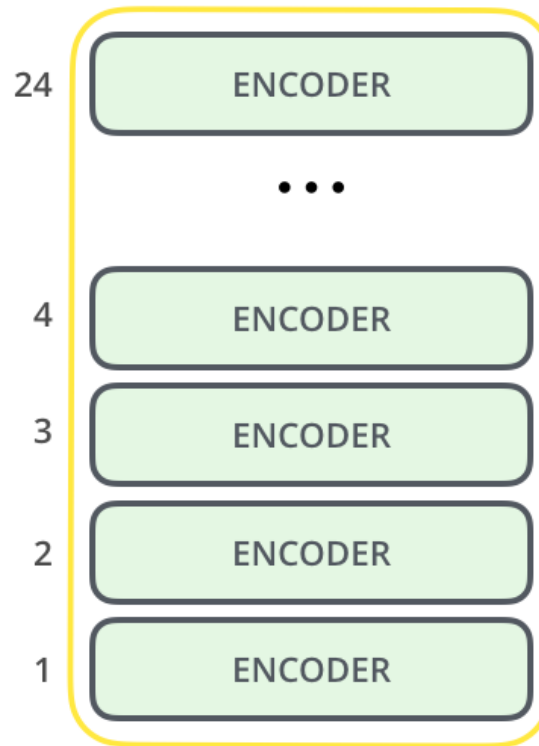
We argue that current techniques restrict the power of the pre-trained representations, especially for the fine-tuning approaches. The major limitation is that standard language models are unidirectional, and this limits the choice of architectures that can be used during pre-training. For example, in OpenAI GPT, the authors use a left-to-

BERT: Architecture

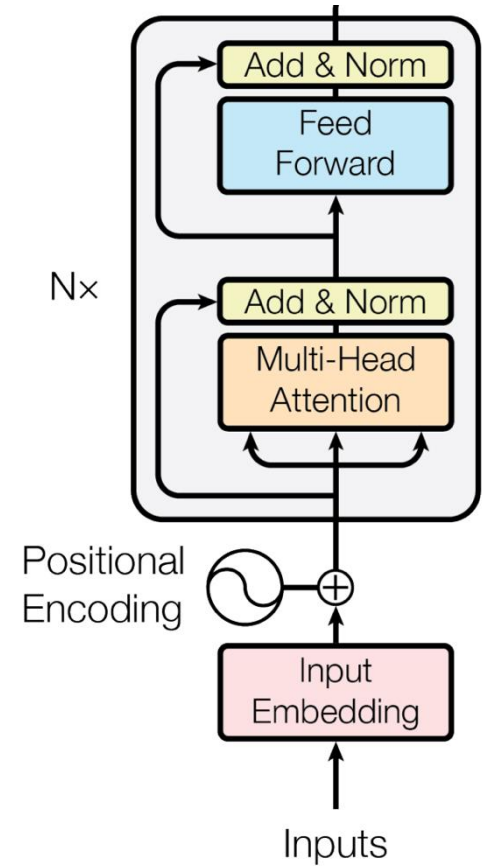
- Stacks of Transformer encoders”



BERT_{BASE}

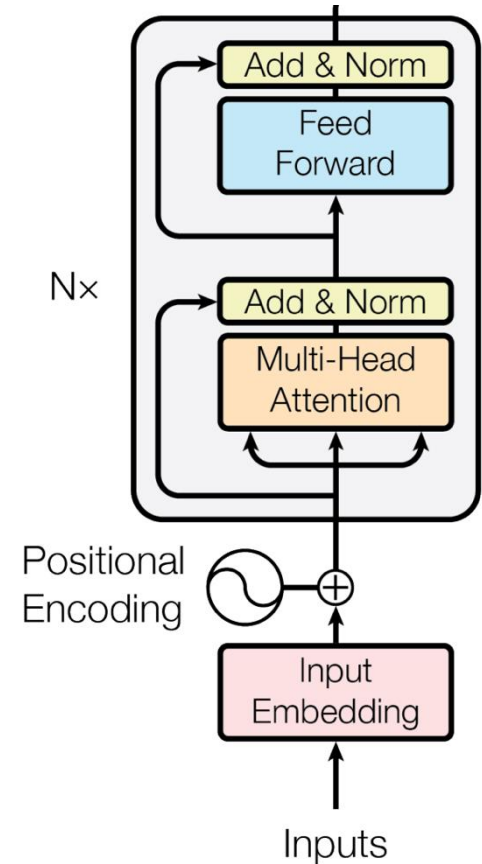
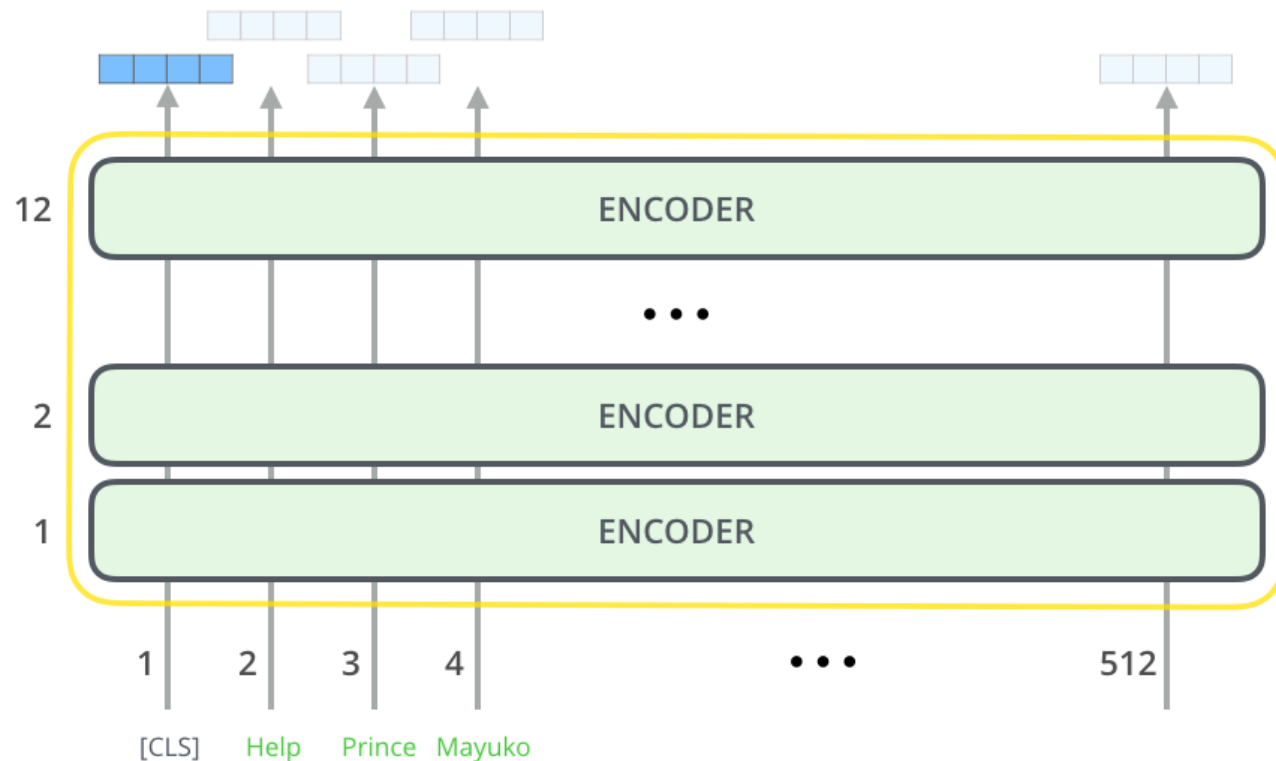


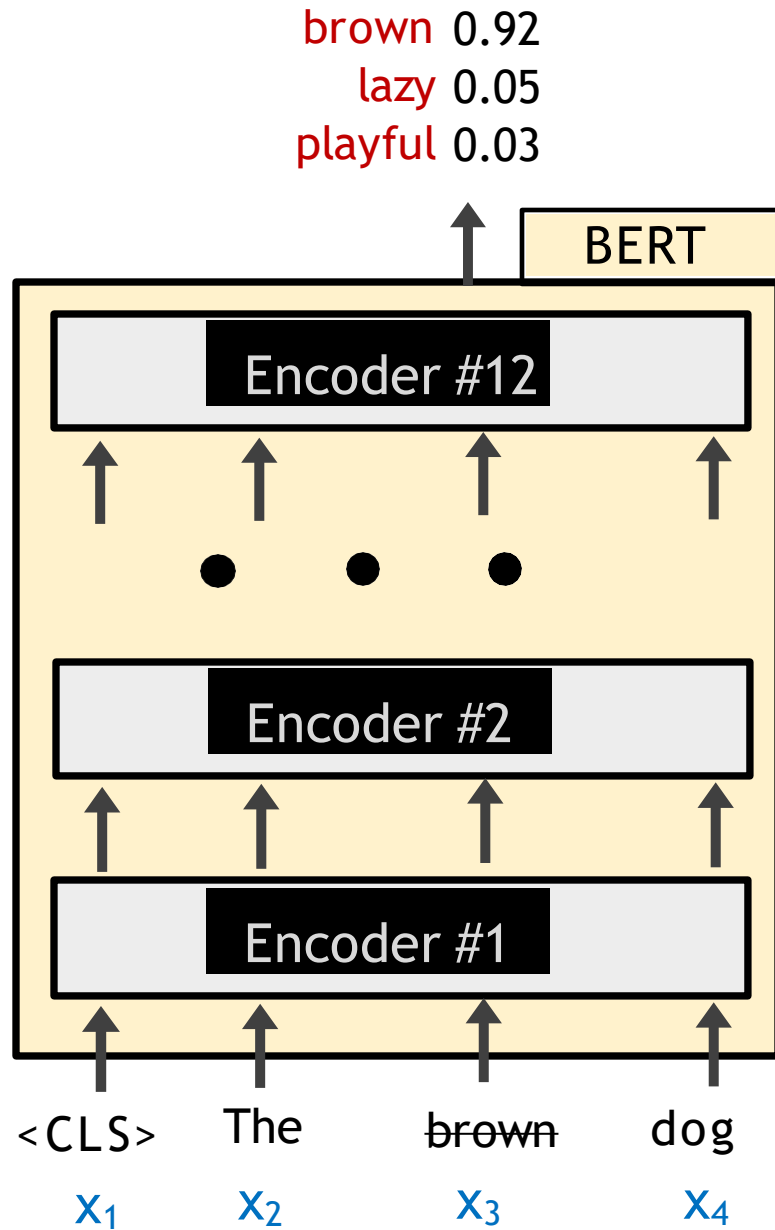
BERT_{LARGE}



BERT: Architecture

- Model output dimension: 512





BERT is trained to uncover masked tokens.

Probing BERT Masked LM

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

Paris is the [MASK] of France.

Compute

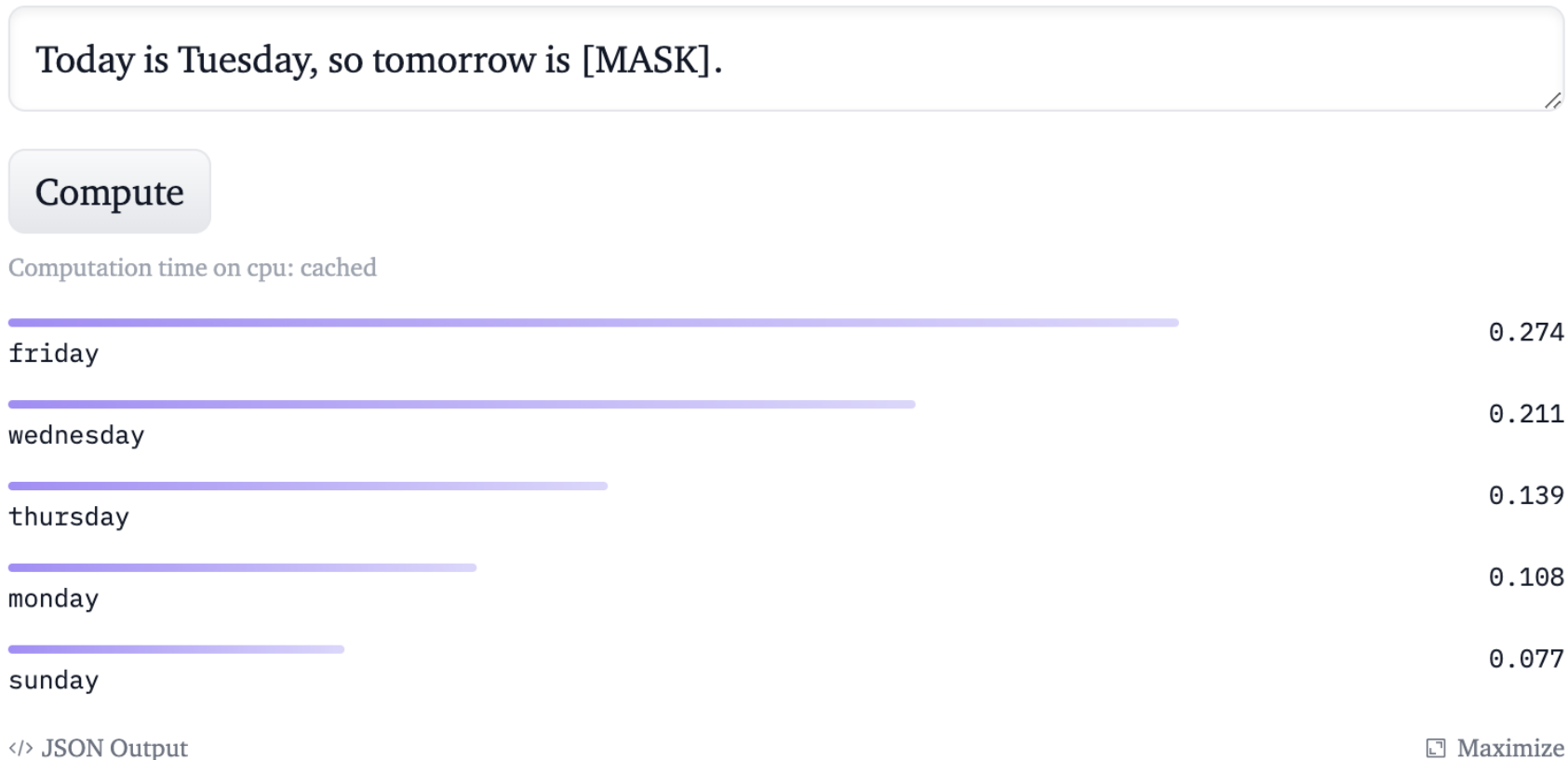
Computation time on cpu: cached

capital	0.997
heart	0.001
center	0.000
centre	0.000
city	0.000

</> JSON Output Maximize

Probing BERT Masked LM

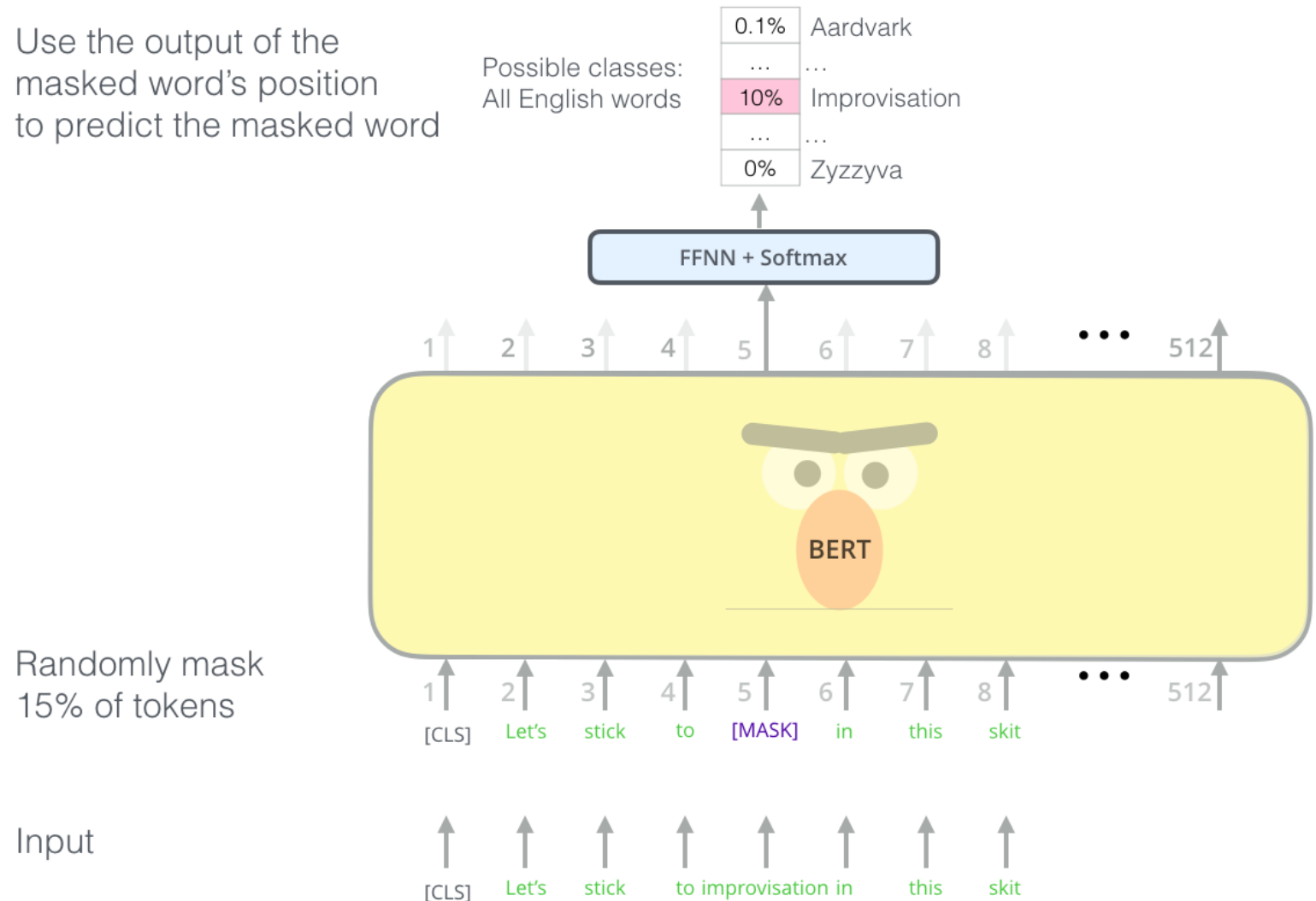
- Making words forces BERT to use context in both directions to predict the masked word.



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

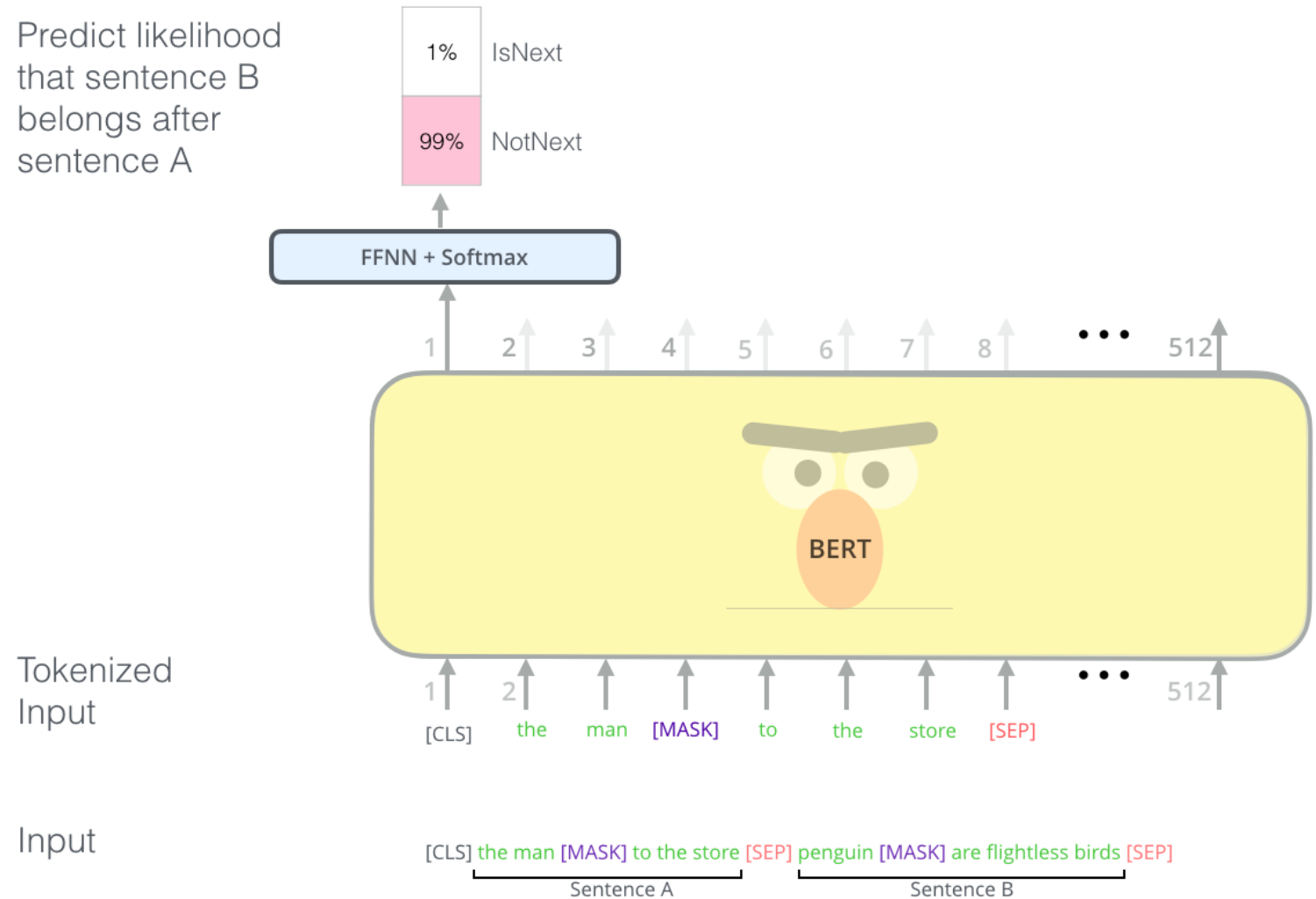
Galon

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 context)

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

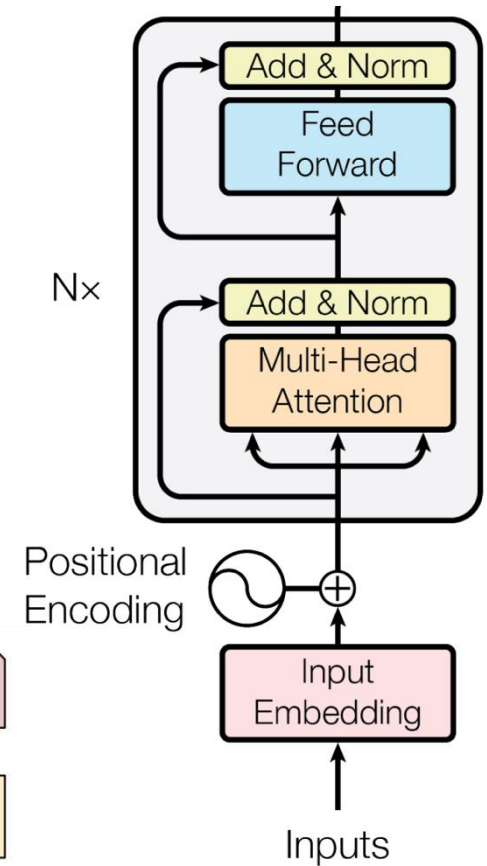
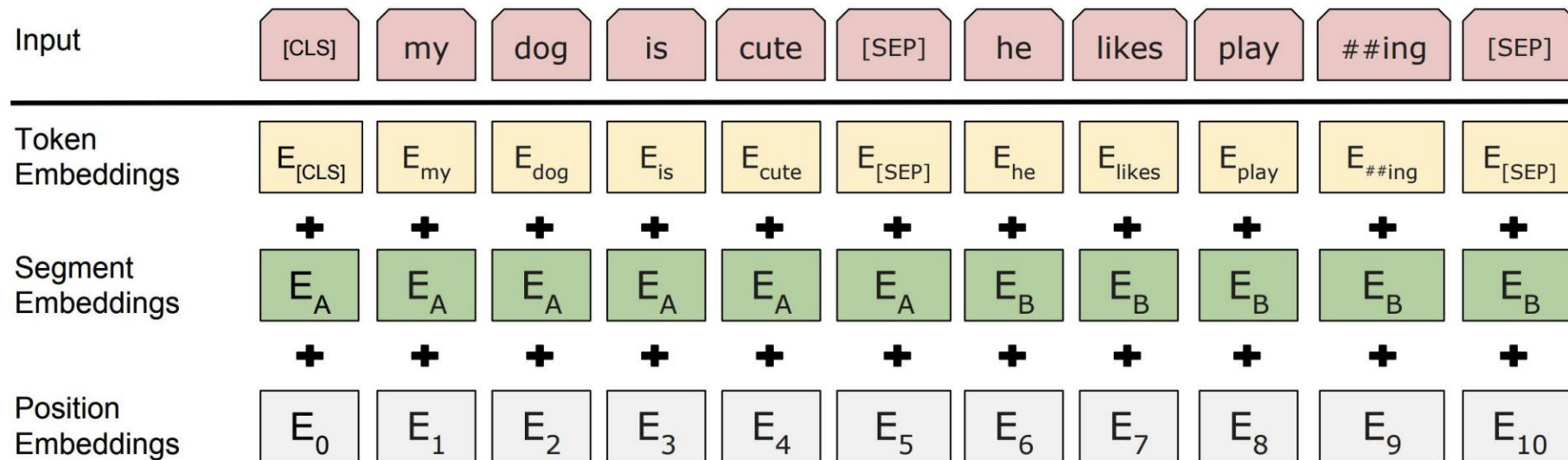
- Learn relationships between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

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

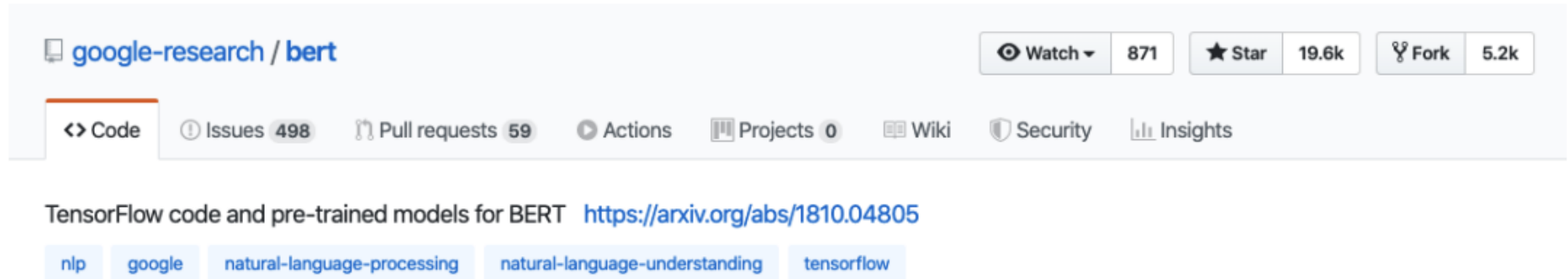


Training

- Trains model on unlabeled data over different pre-training tasks (self-supervised learning)
- **Data:** Wikipedia (2.5B words) + BookCorpus (800M words)
- **Training Time:** 1M steps (~40 epochs)
- **Optimizer:** AdamW, $1e-4$ learning rate, linear decay
- **BERT-Base:** 12-layer, 768-hidden, 12-head
- **BERT-Large:** 24-layer, 1024-hidden, 16-head
- Trained on 4x4 or 8x8 TPUs for 4 days

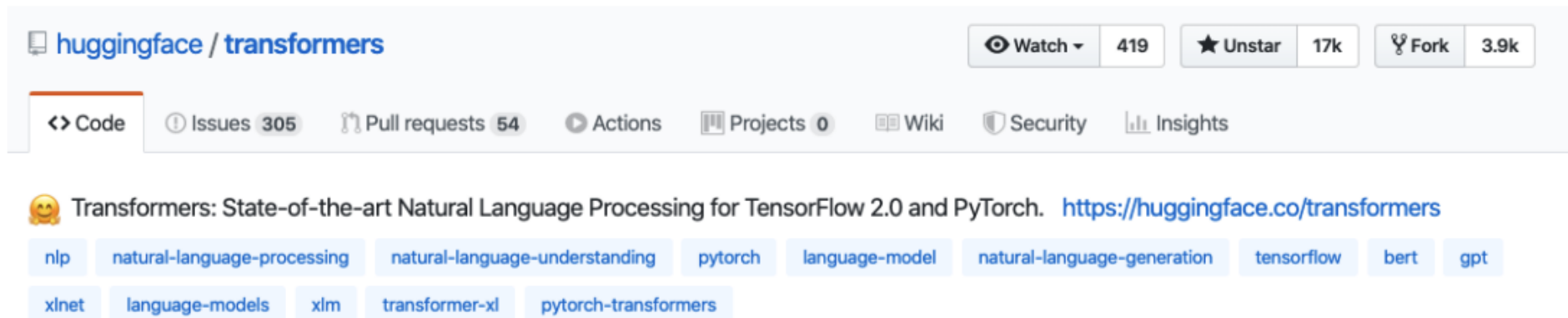
BERT in Practice

TensorFlow: <https://github.com/google-research/bert>



The screenshot shows the GitHub repository page for `google-research / bert`. At the top right, there are buttons for 'Watch' (871), 'Star' (19.6k), and 'Fork' (5.2k). Below these are navigation tabs: '<> Code' (selected), 'Issues 498', 'Pull requests 59', 'Actions', 'Projects 0', 'Wiki', 'Security', and 'Insights'. The repository description reads: 'TensorFlow code and pre-trained models for BERT' followed by the link <https://arxiv.org/abs/1810.04805>. Below the description are several topic tags: `nlp`, `google`, `natural-language-processing`, `natural-language-understanding`, and `tensorflow`.

PyTorch: <https://github.com/huggingface/transformers>

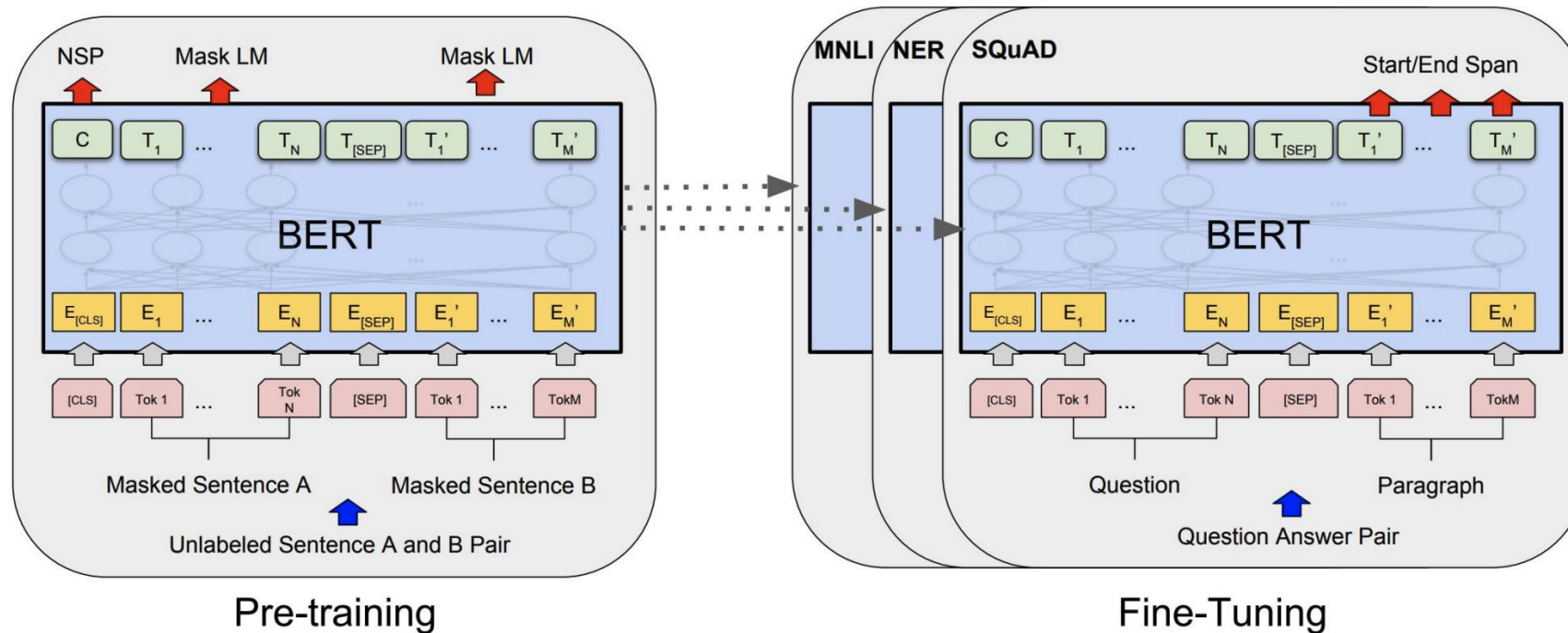


The screenshot shows the GitHub repository page for `huggingface / transformers`. At the top right, there are buttons for 'Watch' (419), 'Unstar' (17k), and 'Fork' (3.9k). Below these are navigation tabs: '<> Code' (selected), 'Issues 305', 'Pull requests 54', 'Actions', 'Projects 0', 'Wiki', 'Security', and 'Insights'. The repository description reads: '🤖 Transformers: State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch.' followed by the link <https://huggingface.co/transformers>. Below the description are several topic tags: `nlp`, `natural-language-processing`, `natural-language-understanding`, `pytorch`, `language-model`, `natural-language-generation`, `tensorflow`, `bert`, `gpt`, `xlnet`, `language-models`, `xlm`, `transformer-xl`, and `pytorch-transformers`.

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



An Example Result: SWAG

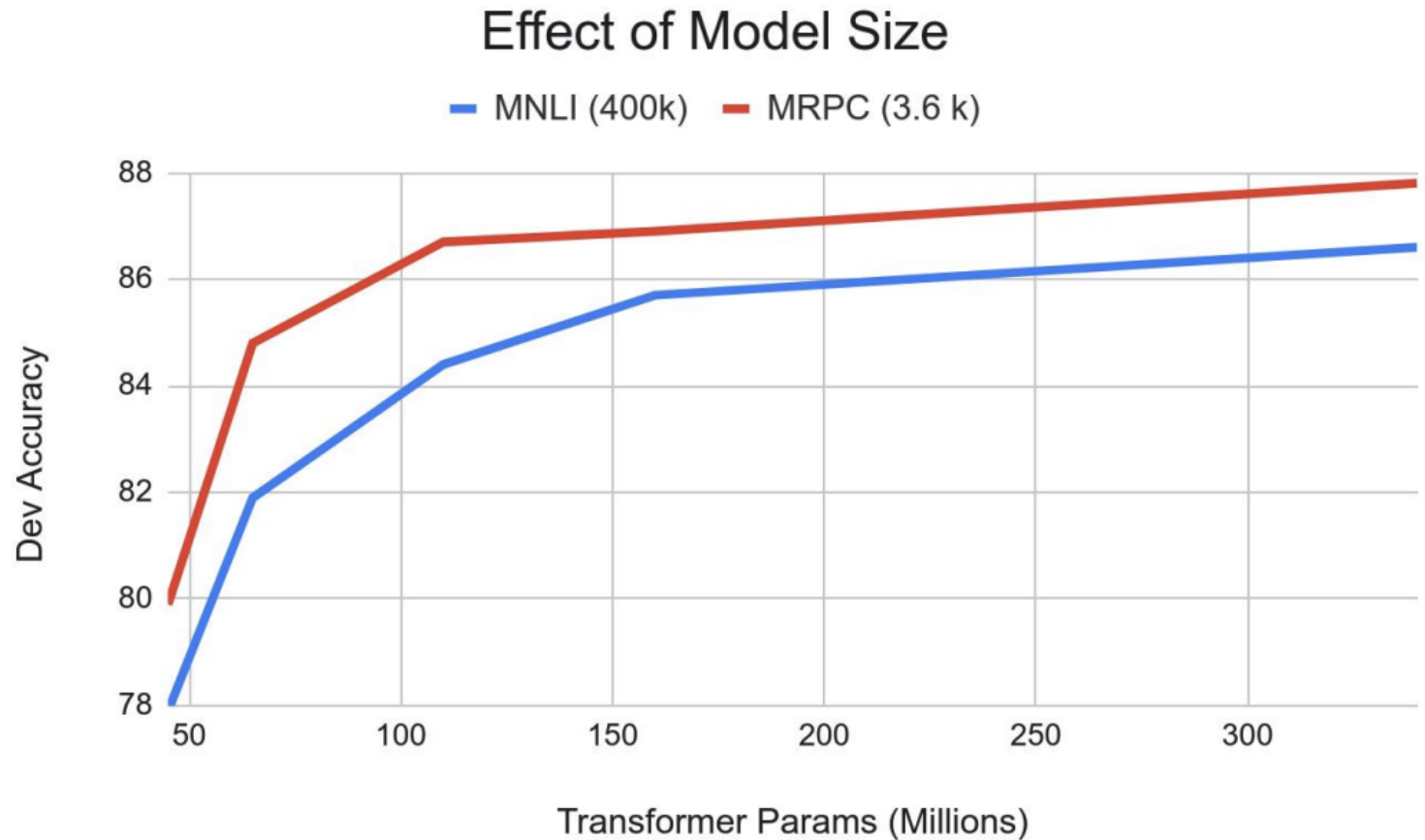


Rank	Model	Test Score
1	BERT (Bidirectional Encoder Representations from Transfo... <i>Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova</i> 10/11/2018	86.28%
2	OpenAI Transformer Language Model <i>Original work by Alec Radford, Karthik Narasimhan, Tim Salimans, ...</i> 10/11/2018	77.97%
3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

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



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

Why did no one think of this before?

- Concretely, 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.

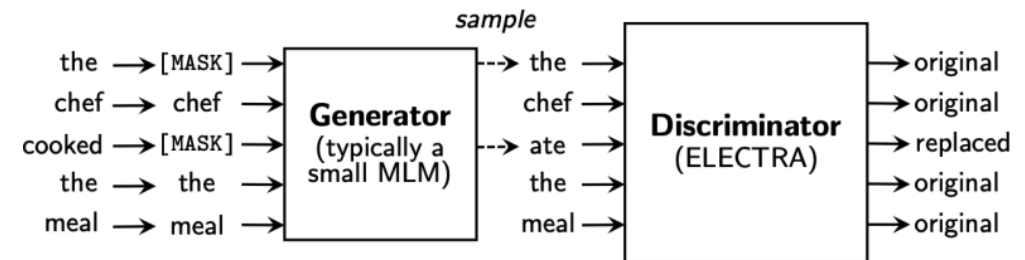
What Happened After BERT?

- RoBERTa (Liu et al., 2019)
 - 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)
 - Still one of the most popular models to date

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	94.6/89.4	90.2	96.4
BERT _{LARGE}						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

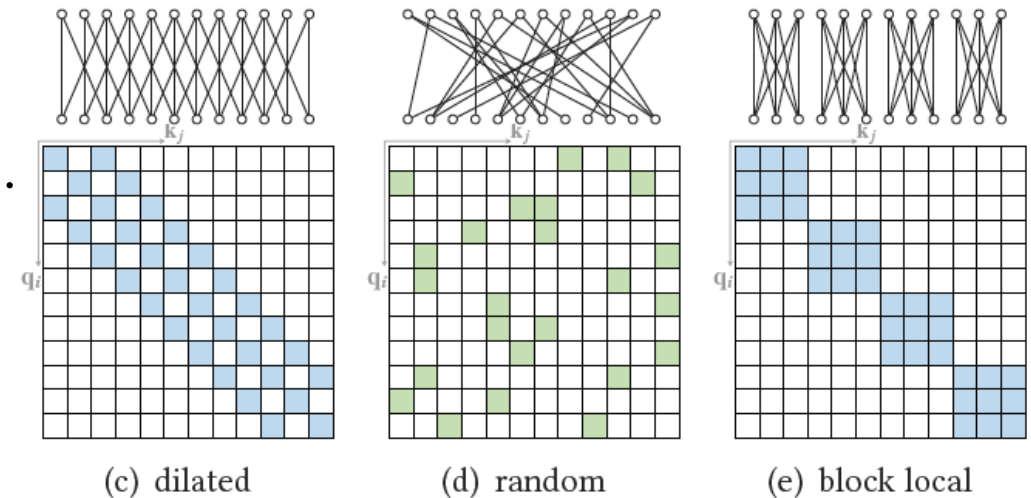
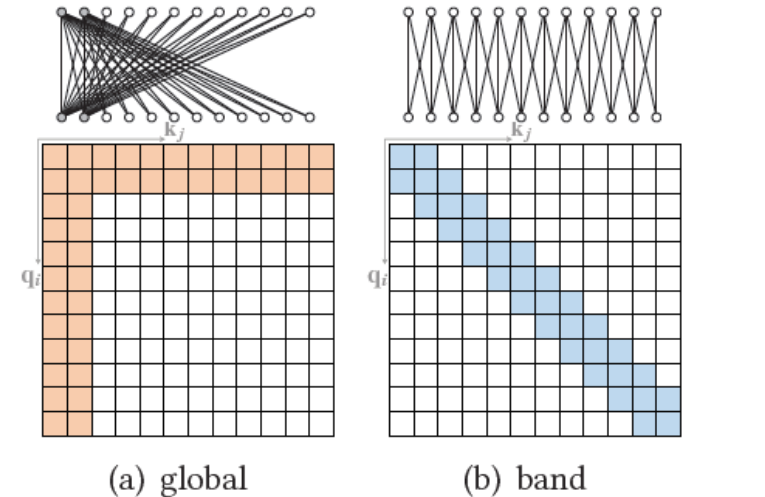
What Happened After BERT?

- RoBERTa (Liu et al., 2019)
 - 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)
 - Still one of the most popular models to date
- 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)
 - Two models generator and discriminator
 - It provides a more efficient training method



What Happened After BERT?

- Models that handle long contexts (512 tokens)
 - Longformer, Big Bird, ...
- Multilingual BERT
 - Trained single model on 104 languages from Wikipedia. Shared 110k WordPiece vocabulary
- BERT extended to different domains
 - SciBERT, BioBERT, FinBERT, ClinicalBERT, ..
- Making BERT smaller to use
 - DistillBERT, TinyBERT, ...



Text generation using BERT

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](#)

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

Summary

1. What is self-supervised learning?
2. Examples of self-supervision in NLP
 - Word embeddings (e.g., word2vec)
 - Language models (e.g., GPT)
 - Masked language models (e.g., BERT)
3. Open challenges
 - Demoting bias
 - Capturing factual knowledge
 - Learning symbolic reasoning