

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Economics 2355: Transformer-Based Language Models

Melissa Dell

Harvard University

March 2021

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Contextualized Word Embeddings

- ▶ With word embeddings from a model like Word2Vec or GloVe, a given word will always have the same embedding, regardless of its context
- ▶ This is a problem, as many words have multiple meanings depending on the context in which they are used
- ▶ A transformer based model is not required to contextualize embeddings (although these models came to dominate for reasons we'll see)

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# ELMO

- ▶ ELMo was a prominent model in the development of contextualized word embeddings
- ▶ ELMo is a language model whose objective is to predict the next word in a sequence
- ▶ It also helped to popularize the approach of using a pre-trained language model to create representations for various downstream tasks

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

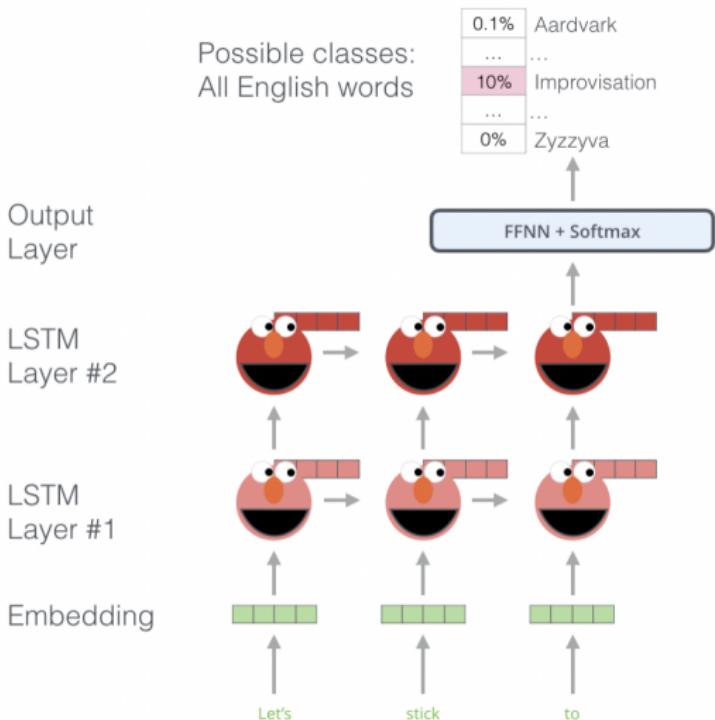
XLNet

Longformer

BigBird

Overview and  
What to Use

# ELMo is based on an LSTM



A step in the pre-training process of ELMo: Given "Let's stick to" as input, predict the next most likely word – a *language modeling* task. When trained on a large dataset, the model starts to pick up on language patterns. It's unlikely it'll accurately guess the next word in this example. More realistically, after a word such as "hang", it will assign a higher probability to a word like "out" (to spell "hang out") than to "camera".

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

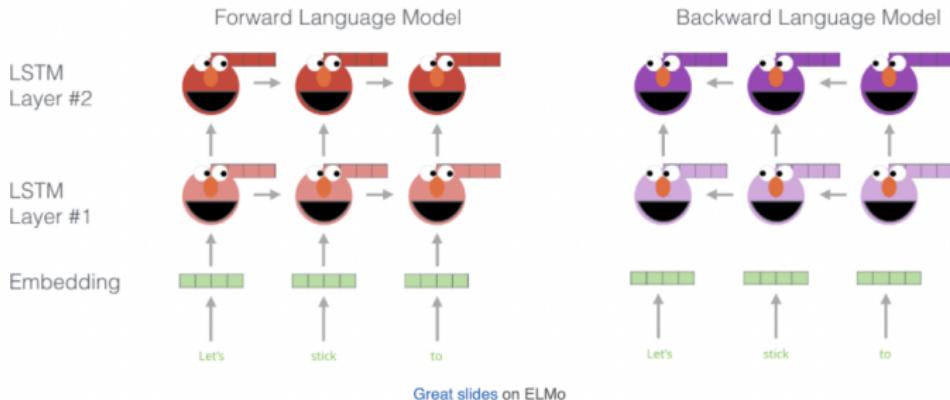
BigBird

Overview and  
What to Use

# ELMo is based on an LSTM

It is bi-directional, to allow contexts from both directions to matter

Embedding of "stick" in "Let's stick to" - Step #1



<http://jalammar.github.io/illustrated-bert/>

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Creating Contextualized Embeddings

- ▶ ELMo creates contextualized embeddings by concatenating the hidden states and initial embeddings across the forward and backward LSTMs
- ▶ Then it takes a weighted sum of the hidden states and initial embedding vectors

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

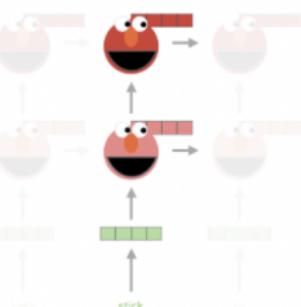
# Creating Contextualized Embeddings

Embedding of "stick" in "Let's stick to" - Step #2

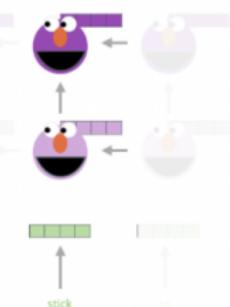
1- Concatenate hidden layers



Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task



3- Sum the (now weighted) vectors



ELMo embedding of "stick" for this task in this context

<http://jalammar.github.io/illustrated-bert/>

To create the initial embeddings, the model combines character level embeddings

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Open AI GPT Model

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

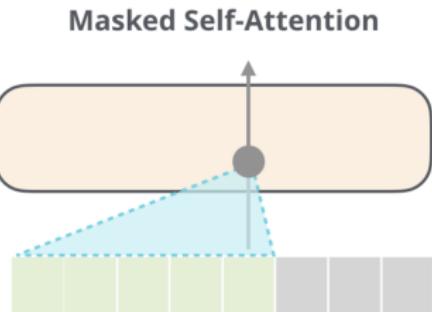
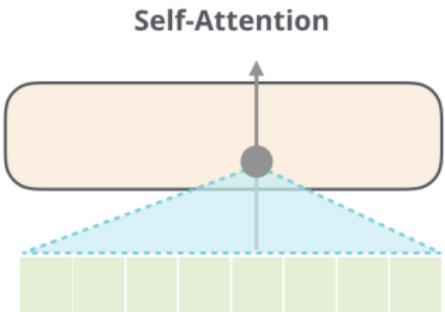
BigBird

Overview and  
What to Use

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# The Decoder Uses Masked Self-Attention

In the encoder self-attention blocks, a position can see tokens from the right. This is not the case for the decoder blocks. The idea here is akin to predicting the next word, given the previous contexts.



<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

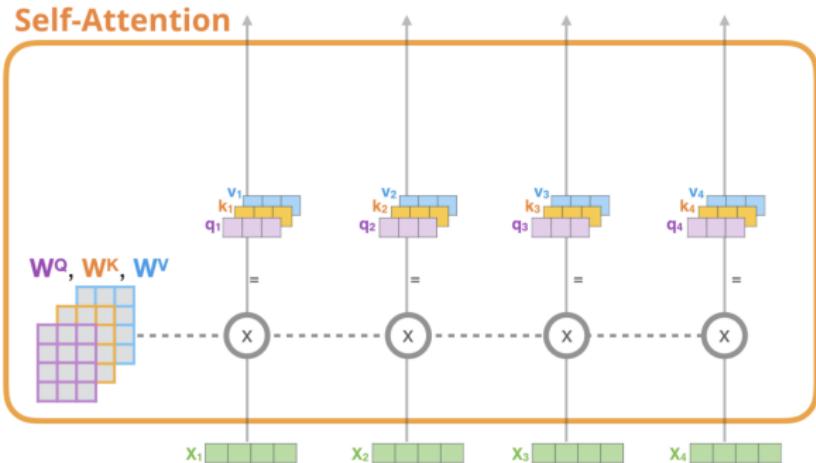
BigBird

Overview and  
What to Use

# A Recap of Self-Attention without Masking:

## Step 1

- 1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices  $W^Q$ ,  $W^K$ ,  $W^V$



<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

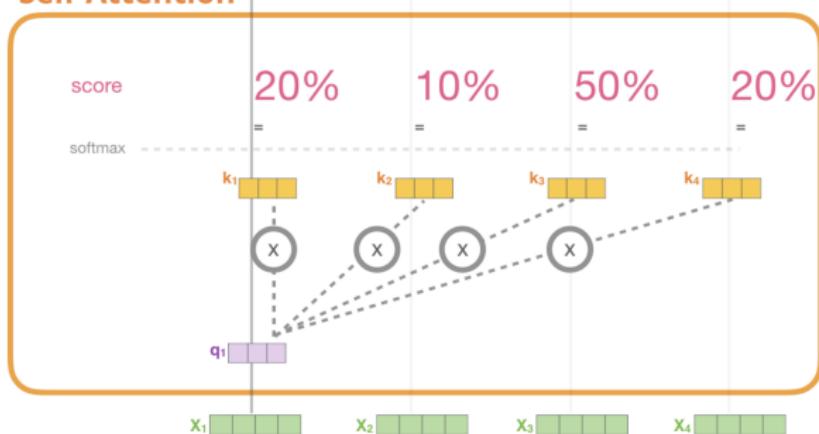
BigBird

Overview and  
What to Use

# A Recap of Self-Attention without Masking: Step 2

- 2) Multiply (dot product) the current **query vector**, by all the **key vectors**, to get a score of how well they match

## Self-Attention



<https://jalammar.github.io/illustrated-gpt2/>

# A Recap of Self-Attention without Masking: Step 3

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

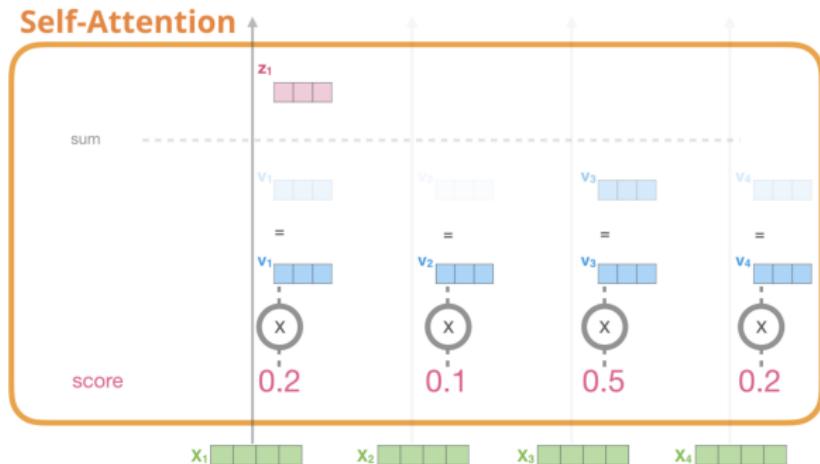
GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

The lower the score, the more transparent we're showing the value vector. That's to indicate how multiplying by a small number dilutes the values of the vector.

<https://jalammar.github.io/illustrated-gpt2/>

# A Recap of Self-Attention without Masking: Step 4

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

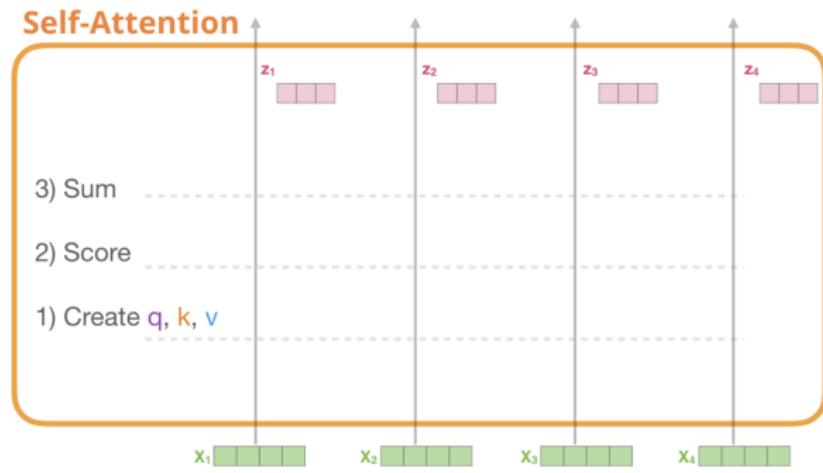
GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

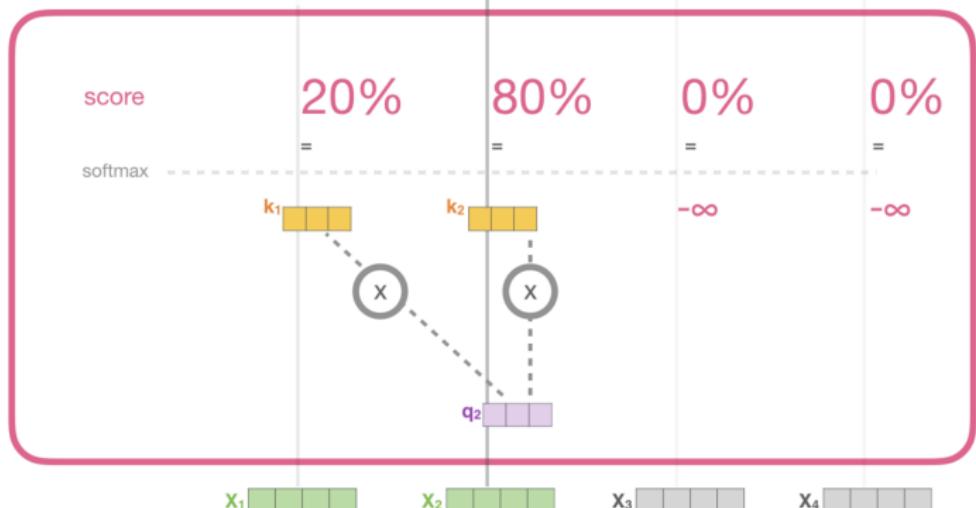
<https://jalammar.github.io/illustrated-gpt2/>

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Masked Self-Attention

Masked self-attention is identical to self-attention without masking, except at step 2.

## Masked Self-Attention



<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Masked Self-Attention

Consider a sentence with four words:

	Features				Labels	
position: 1	2	3	4			
Example:						
1	robot	must	obey	orders	must	
2	robot	must	obey	orders	obey	
3	robot	must	obey	orders	orders	
4	robot	must	obey	orders	<eos>	

<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Masked Self-Attention

Calculate the attention scores by multiplying the query matrix by the key matrix:

$$\begin{array}{c}
 \text{Queries} \\
 \begin{array}{cccc} \text{robot} & \text{must} & \text{obey} & \text{orders} \end{array} \times
 \begin{array}{c} \text{Keys} \\ \begin{array}{cccc} \text{robot} & \text{must} & \text{obey} & \text{orders} \\ \text{robot} & \text{must} & \text{obey} & \text{orders} \\ \text{robot} & \text{must} & \text{obey} & \text{orders} \\ \text{robot} & \text{must} & \text{obey} & \text{orders} \end{array} \end{array} = \begin{array}{c} \text{Scores} \\ \text{(before softmax)} \\ \begin{array}{ccccc} 0.11 & 0.00 & 0.81 & 0.79 \\ 0.19 & 0.50 & 0.30 & 0.48 \\ 0.53 & 0.98 & 0.95 & 0.14 \\ 0.81 & 0.86 & 0.38 & 0.90 \end{array} \end{array}
 \end{array}$$

<https://jalammar.github.io/illustrated-gpt2/>

In practice, each word is associated with a vector

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

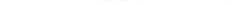
Overview and  
What to Use

# Masked Self-Attention

After creating the scores, multiply by an attention mask, which will be a triangle

Scores (before softmax)			
0.11	0.00	0.81	0.79
0.19	0.50	0.30	0.48
0.53	0.98	0.95	0.14
0.81	0.86	0.38	0.90

Apply Attention  
Mask



Masked Scores (before softmax)			
0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90

<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Masked Self-Attention

Then apply softmax

**Masked Scores**  
(before softmax)

0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90

Softmax  
(along rows)

**Scores**

1	0	0	0
0.48	0.52	0	0
0.31	0.35	0.34	0
0.25	0.26	0.23	0.26

<https://jalammar.github.io/illustrated-gpt2/>

# GPT

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

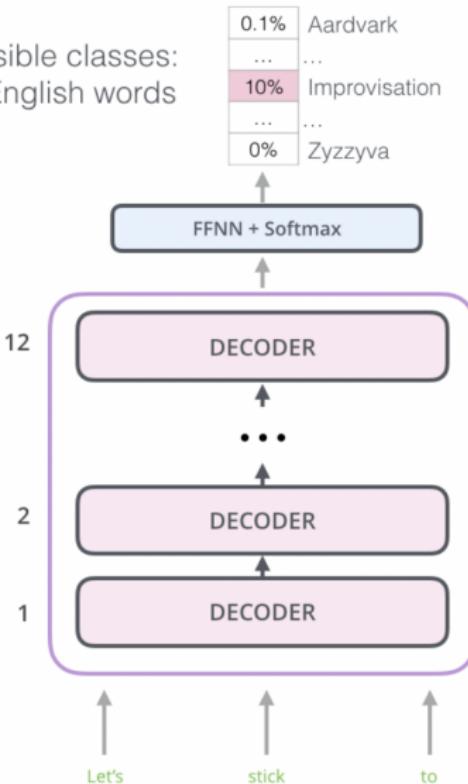
XLNet

Longformer

BigBird

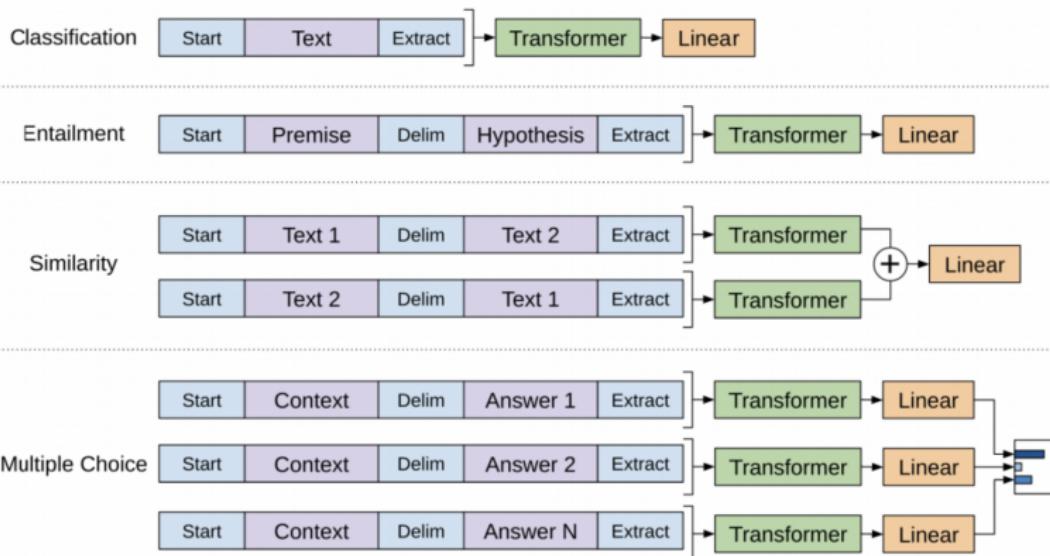
Overview and  
What to Use

Possible classes:  
All English words



<https://jalammar.github.io/illustrated-bert/>

Inputs can be transformed to configure the model to perform different downstream tasks



Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Transformers for Language Modeling

- ▶ At first it would seem more intuitive to stack decoder blocks, as the decode is a language model.
- ▶ But if you just use the decoder blocks, the language model is not bi-directional. You could feed the sequence into the decoder backwards, but then you would need a way to combine the forwards and backwards representations
- ▶ The encoder blocks are not a language model. To create representations for each word, the encoder can attend to all other positions in the sequence that was fed into the model
- ▶ We'll now see a very influential paper whose innovation was to realize that encoder blocks can be used to create a bidirectional language model through masking

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Outline

Overview

Contextualized Word Embeddings

GPT

**BERT**

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

**BERT**

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# BERT

- ▶ BERT, which stands for *Bidirectional Encoder Representations for Transformers*, marked a major leap forward in NLP, developing a Transformer-based language model architecture
- ▶ While it continues the Sesame Street theme, the name is also instructive
  - ▶ Bidirectional refers to the fact that it considers contexts from both directions (as with the bidirectional LSTM);
  - ▶ Encoder refers to the fact that it stacks encoder transformer blocks
  - ▶ Its most remarkable feature is creating representations that can be used for a very diverse set of downstream tasks

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Bidirectional Encoding

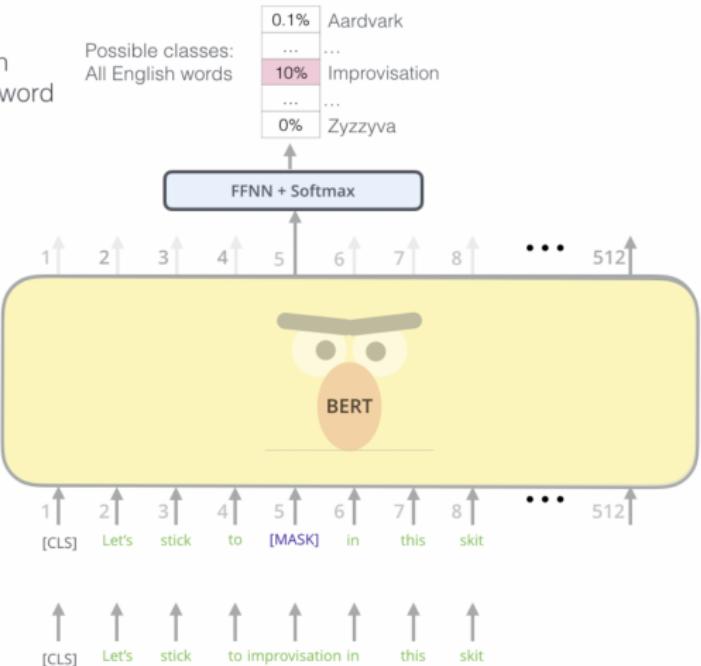
- ▶ The transformer encoder reads the sequence all at once, making it bidirectional
- ▶ BERT trains the encoder as a language model by using an approach called Masked Language Model (MLM): 15% of tokens are sampled at random. Of those, 80% are replaced with the MASK token. 10% are replaced with a random word, and 10% are kept the same
- ▶ The model is trained by predicting these 15% of tokens. The numbers are a bit arbitrary, and there weren't experiments reported in the paper on this

# Masking

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

Randomly mask 15% of tokens

Input



<http://jalammar.github.io/illustrated-bert/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

Economics 2355

Melissa Dell

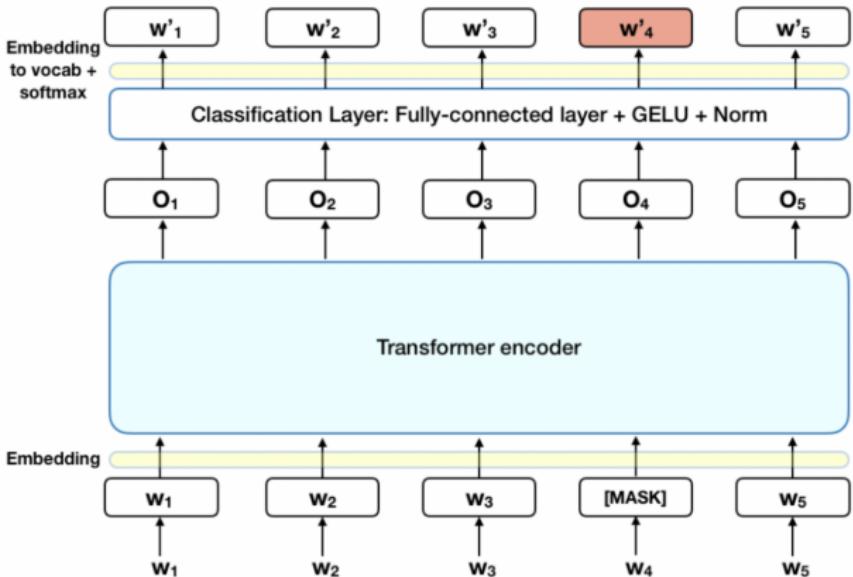
BERT

T5

XI Net

BigBird

## Overview and What to Use



[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Next Sentence Prediction

- ▶ A second key part of the training objective is next sentence prediction
- ▶ The model receives pairs of sentences as inputs and predicts whether the second sentence comes directly after the first in the original corpus
- ▶ Half the time it is the next sentence, the other half the time it is chosen randomly

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

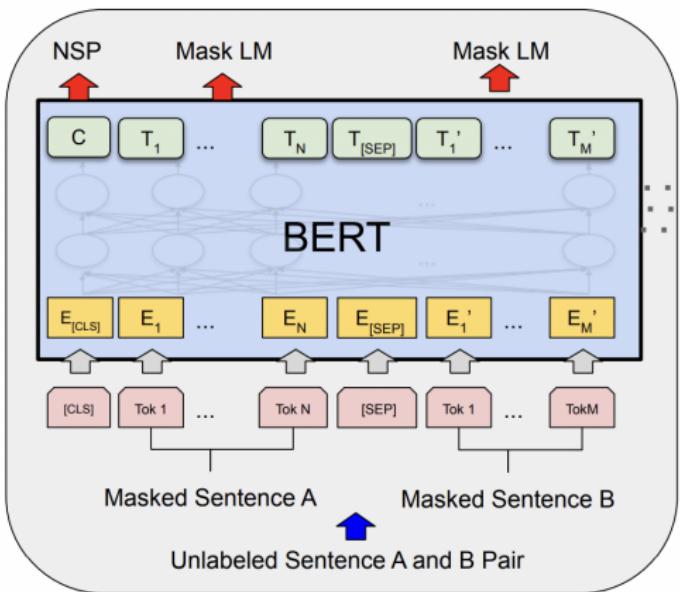
Longformer

BigBird

Overview and  
What to Use

# Next Sentence Prediction

It does this by inserting a special *CLS* token at the begin of the first sentence:



Pre-training

## Next Sentence Prediction

Economics 2355

Melissa Dell

BERT

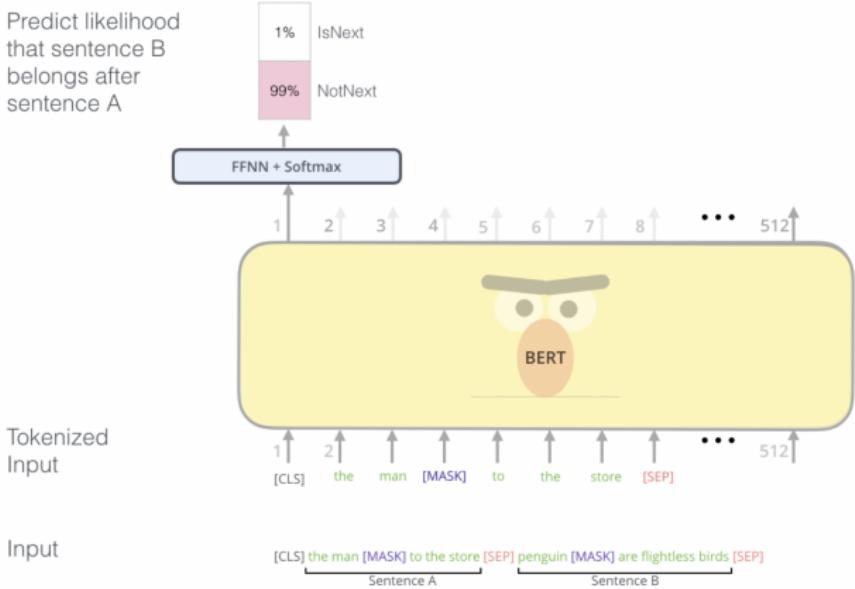
T5

XLNet

## Longformer

BigBird

## Overview and What to Use



Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

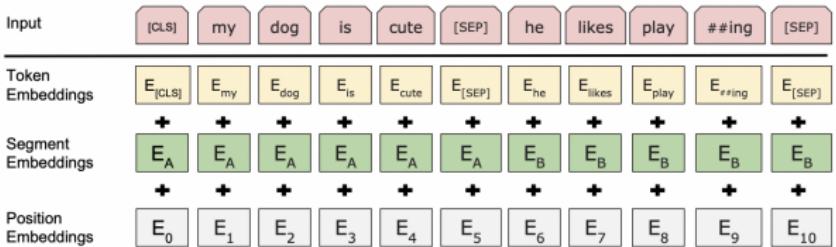
Longformer

BigBird

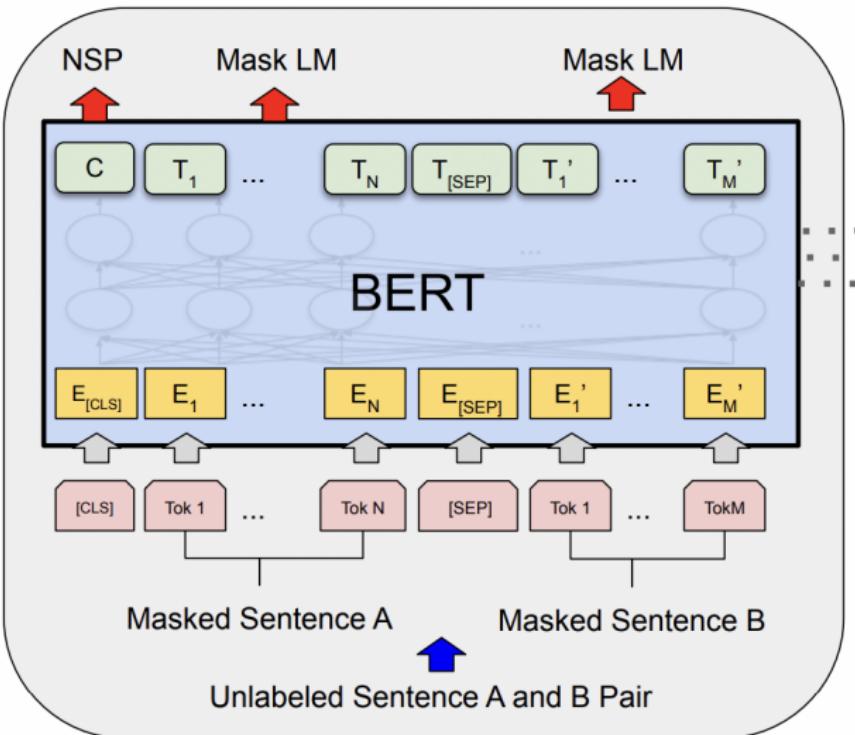
Overview and  
What to Use

# Training Input

Here is how tokens are input at training:



Devlin et al.

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

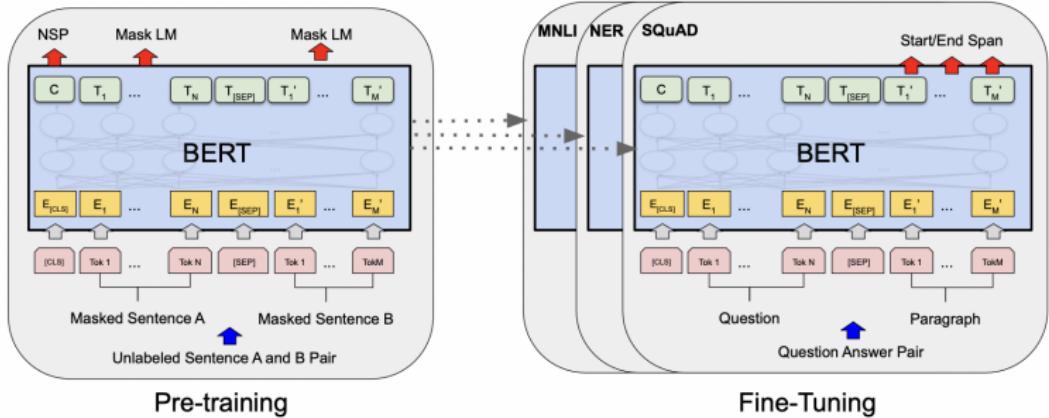
XLNet

Longformer

BigBird

Overview and  
What to Use

# Adapt to Fine-Tuning on Specific Tasks



Devlin et al.

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

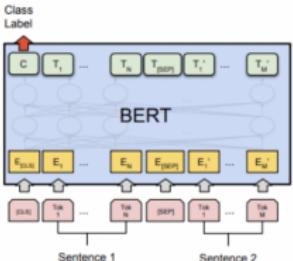
XLNet

Longformer

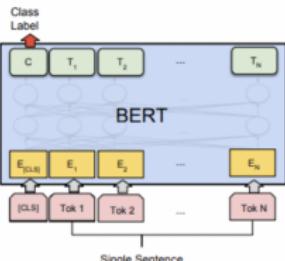
BigBird

Overview and  
What to Use

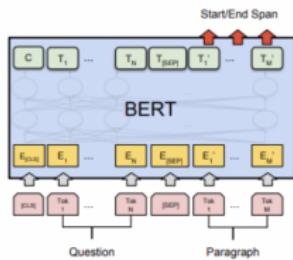
# Adapt to Fine-Tuning on Specific Tasks



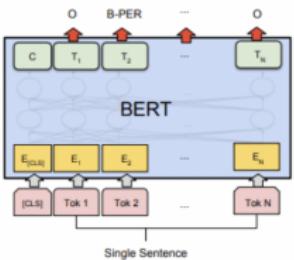
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# An Alternative to Fine-Tuning

- ▶ Take the pre-trained BERT embeddings
- ▶ Feed them into your own existing downstream model
- ▶ The paper shows, with an example from named entity recognition, that this works almost as well as fine-tuning
- ▶ Your mileage will probably vary

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

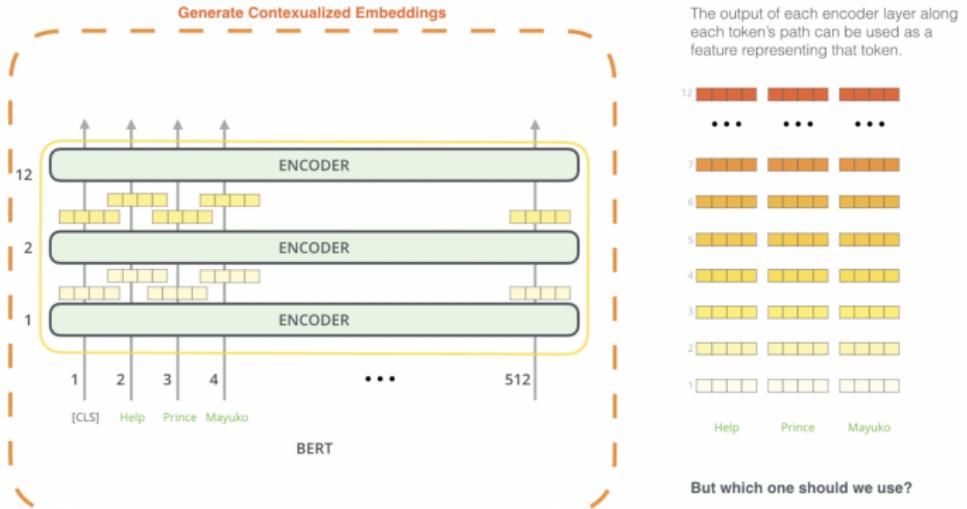
XLNet

Longformer

BigBird

Overview and What to Use

# Which Layer to Use?



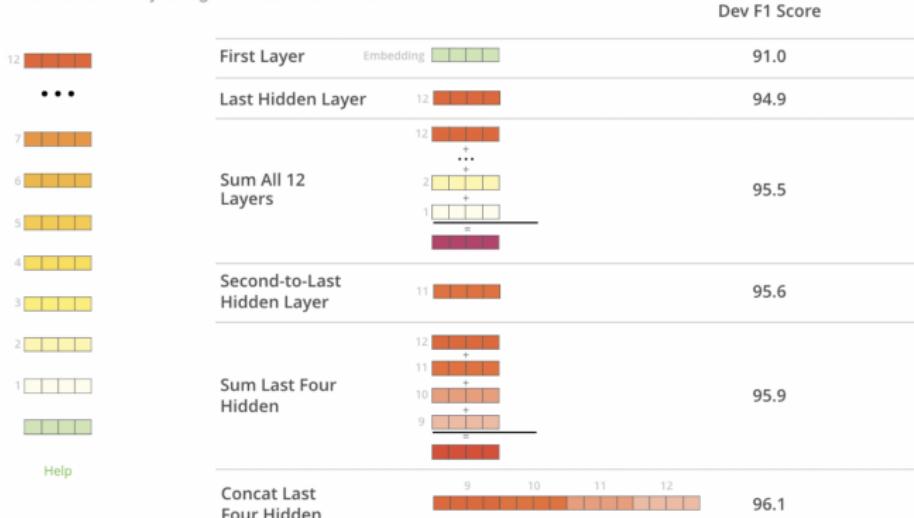
<http://jalammar.github.io/illustrated-bert/>

# What Layer to Use

The fine-tuned model achieves a score of 96.4

What is the best contextualized embedding for "Help" in that context?

For named-entity recognition task CoNLL-2003 NER



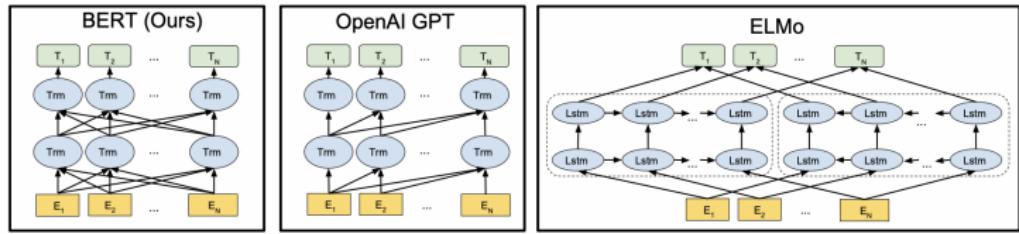
<http://jalammar.github.io/illustrated-bert/>

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Which Layer to Use

- ▶ Likely depends on the task
- ▶ General consensus that upper layers are more task specific
- ▶ A paper “How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings” finds that representations of the same word in different contexts are less similar in the upper layers, suggesting that the upper layers produce more context specific representations

## BERT Summary



Devlin et al.

# BERT Options

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# RoBERTa

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# RoBERTa

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# RoBERTa Performance

	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT <sub>LARGE</sub>	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet <sub>LARGE</sub>	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	<b>90.2/90.2</b>	<b>94.7</b>	<b>92.2</b>	<b>86.6</b>	<b>96.4</b>	<b>90.9</b>	<b>68.0</b>	<b>92.4</b>	<b>91.3</b>	-
<i>Ensembles on test (from leaderboard as of July 25, 2019)</i>										
ALICE	88.2/87.9	95.7	<b>90.7</b>	83.5	95.2	92.6	<b>68.6</b>	91.1	80.8	86.3
MT-DNN	87.9/87.4	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2/89.8	98.6	90.3	86.3	<b>96.8</b>	<b>93.0</b>	67.8	91.6	<b>90.4</b>	88.4
RoBERTa	<b>90.8/90.2</b>	<b>98.9</b>	90.2	<b>88.2</b>	96.7	92.3	67.8	<b>92.2</b>	89.0	<b>88.5</b>

# RoBERTa Pre-Trained Models

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

## Pre-trained models

Model	Description	# params	Download
<code>roberta.base</code>	RoBERTa using the BERT-base architecture	125M	<a href="#">roberta.base.tar.gz</a>
<code>roberta.large</code>	RoBERTa using the BERT-large architecture	355M	<a href="#">roberta.large.tar.gz</a>
<code>roberta.large.mnli</code>	<code>roberta.large</code> finetuned on <a href="#">MNLI</a>	355M	<a href="#">roberta.large.mnli.tar.gz</a>
<code>roberta.large.wsc</code>	<code>roberta.large</code> finetuned on <a href="#">WSC</a>	355M	<a href="#">roberta.large.wsc.tar.gz</a>

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# RoBERTa Pre-Trained Models

- ▶ A Winograd schema is a pair of sentences that differ in only one or two words and that contain an ambiguity that is resolved in opposite ways in the two sentences and requires the use of world knowledge and reasoning
  - ▶ “The city councilmen refused the demonstrators a permit because they [feared/advocated] violence.”
  - ▶ If the word is “feared”, then “they” presumably refers to the city council; if it is “advocated” then “they” presumably refers to the demonstrators.
- ▶ The Multi-Genre Natural Language Inference (MultiNLI) corpus is a crowd-sourced collection of 433k sentence pairs annotated with textual entailment information.

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

## Examples

### Premise

#### Fiction

The Old One always comforted Ca'daan, except today.

### Label

*neutral*

### Hypothesis

Ca'daan knew the Old One very well.

#### Letters

Your gift is appreciated by each and every student who will benefit from your generosity.

*neutral*

Hundreds of students will benefit from your generosity.

#### Telephone Speech

yes now you know if everybody like in August when everybody's on vacation or something we can dress a little more casual or

*contradiction*

August is a black out month for vacations in the company.

#### 9/11 Report

At the other end of Pennsylvania Avenue, people began to line up for a White House tour.

*entailment*

People formed a line at the end of Pennsylvania Avenue.

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# A Lightweight BERT

- ▶ While RoBERTa gives the best performance, you might instead want a lightweight model to reduce the computational burden
- ▶ DistilBERT has half the number of parameters of BERT, but retains 97% of the accuracy on benchmark tasks
- ▶ RoBERTa took 4 times as long to train as BERT, whereas DistilBERT trained 4 times faster
- ▶ Trained on the same data as BERT (16GB text, 3.3 billion words)

# DistilBERT

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

**ALBERT**

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

**ALBERT**

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# ALBERT

Melissa Dell

- ▶ Another lightweight option is “A Light BERT for Self-Supervised Learning of Language Representations”
- ▶ ALBERT shares parameters between transformer layers (reducing parameters by 70%)
- ▶ It also uses small embeddings for the input vectors, which are context independent (128 d versus 768 d for the hidden state representations); intuition is that it requires larger embeddings to learn context-dependent representation
- ▶ ALBERT base reduces parameters by 89%
- ▶ Can scale up the size of the hidden representations (while still sharing parameters); ALBERT XXL scales up the hidden representations to 4096, still has a 30% parameter reduction relative to BERT

# ALBERT

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

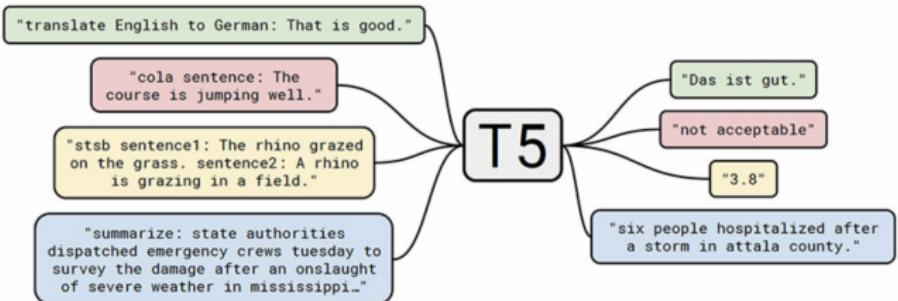
BigBird

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Going the Other Direction: T5

- ▶ A recent enormous model by Google: Text-to-Text Transformer
- ▶ Suggests every NLP model be treated as text-to-text: perspective on transfer learning
- ▶ Produces text as output for all tasks, even those like classification that don't normally have a text output (engineering convenience)
- ▶ The model understands what task is to be performed with a task specific prefix added to the original sentence (i.e. translate French to English, summarize, etc.)
- ▶ SOTA on more than 20 tasks; on many benchmarks, doing as well as humans

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Problems

- ▶ The model is huge: thirty times larger than BERT
- ▶ Too much for commodity hardware to handle, would need to use Google's Cloud TPU platform
- ▶ Will likely be smaller versions, akin to ALBERT

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Problems

- ▶ Pushing human performance on benchmarks, but if you work with these systems, know that they can often fail in unhuman-like ways; you would know if you were talking with a T5 chatbot versus a human customer service representative
- ▶ Language understanding is not solved. Suggests we need more challenging and realistic text datasets (far from obvious *ex ante* how well the off-the-shelf model would perform on our NLP tasks)
- ▶ Biases: Winogender Schemas show that can often correctly classify sentences in stereotypical situations (i.e. male doctor, female nurse) but not vice versa

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

**GPT2/GPT3**

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

**GPT2/GPT3**

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# GPT2/GPT3

- ▶ Speaking of large models, we would be remiss not to mention Open AI's GPT2/GPT3, the successors to the GPT model we saw earlier
- ▶ If you understand how the decoder transformer blocks work, you understand GPT2/3 at a general level
- ▶ Since it is a language model that generates text one word at a time, it is good for text generation tasks: give it a prompt and it will continue the text
- ▶ Unlike the other models we are discussing, it is proprietary

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

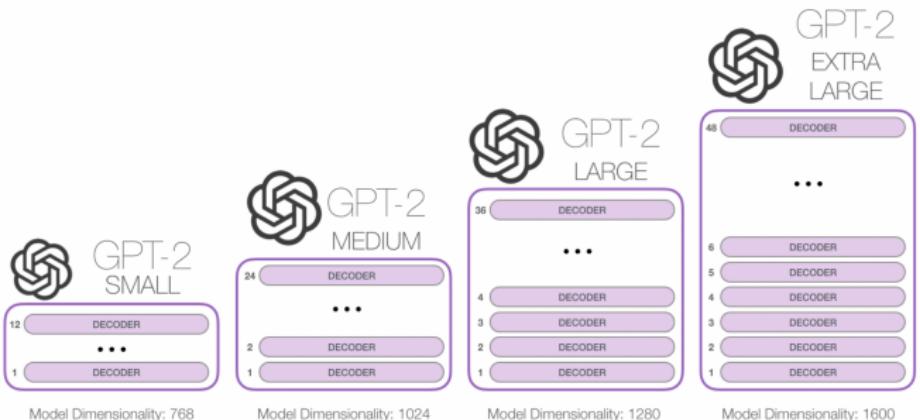
GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

<https://jalammar.github.io/illustrated-gpt2/>

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

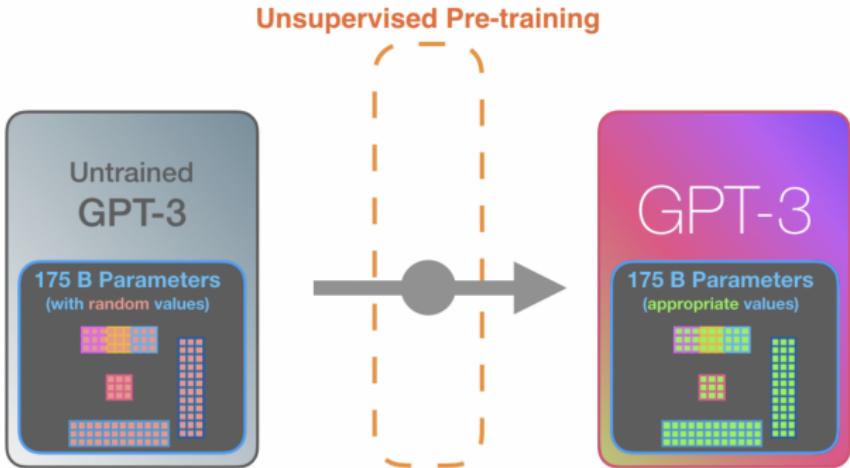
Longformer

BigBird

Overview and  
What to Use

# GPT3

GPT3 blew this out of the water with 175B parameters



<https://jalammar.github.io/how-gpt3-works-visualizations-animations/>

Context window is 2048 tokens and 96 layers deep. 355 GPU years to train, at a cost of \$4.6 million

Also Trained for Downstream Tasks Like  
Summarization

Economics 2355

Melissa Dell

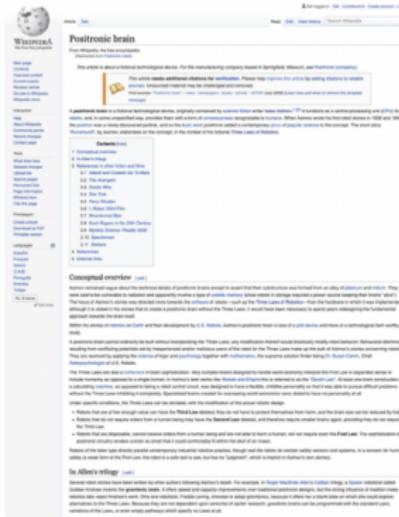
T5

GPT2/GPT3

YI Net

BioPixel

## Overview and What to Use



# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Attention

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

## Standard Attention Computation

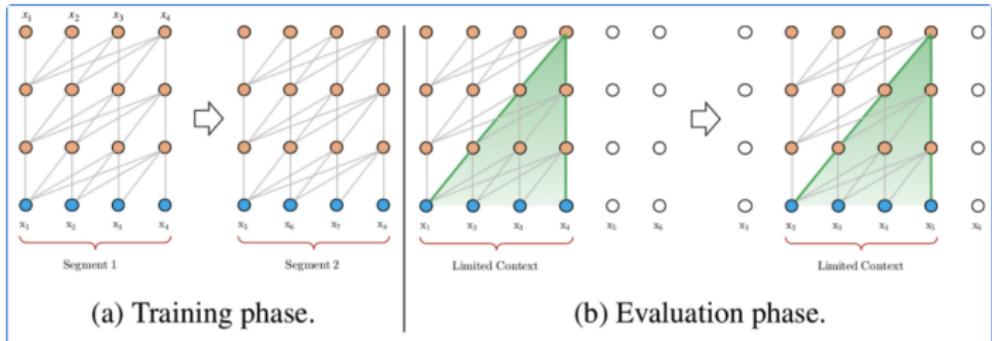
T5

Transformers XL

# XLNet

## BigBird

## Overview and What to Use



Training and Evaluation of the vanilla Transformer language model. Source: [Transformer-XL](#)

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

Introduces two features to overcome the above  
shortcoming:

- ▶ Recurrence Mechanism
- ▶ Relative Positional Encoding

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Recurrence Mechanism

To integrate long term dependencies, each hidden layer receives two inputs

- ▶ The output of the previous hidden layer of that segment, as in Bert
- ▶ The output of the previous hidden layer from the previous segment to create long-term dependencies.

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

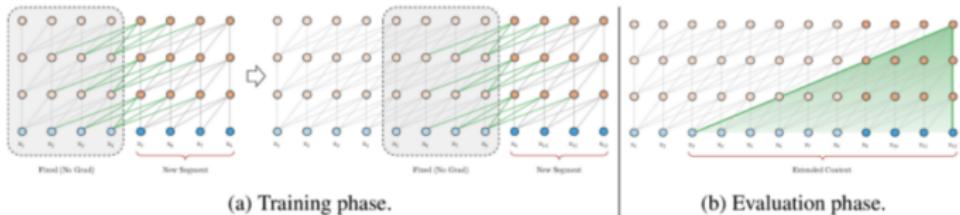
XLNet

Longformer

BigBird

Overview and  
What to Use

# Recurrence Mechanism



Training and Evaluation of the Transformer-XL language model. Source: [Transformer-XL](#)

The two inputs are concatenated and then used to calculate the Key and the Value matrices of the current Head of the current layer of the current segment

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Relative Positional Encoding

- ▶ The original positional encoding handles each segment separately and, as a result, tokens from different segments have the same positional encoding.
- ▶ For example, the first token of the first and the second segments will have the same encoding, although their position and importance are different
- ▶ The paper presents a new positional encoding that is part of each attention module, as opposed to encoding position only before the first layer, and is based on the relative distance between tokens and not their absolute position (incorporates several learned vectors)

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# XLNet

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

- ▶ XLNet uses TransformerXL, with its main contribution being a novel training objective
- ▶ BERT and ELMo improved on SOTA by incorporating left and right contexts
- ▶ XLNet takes this further by predicting each word in a sequence using any combination of other words in the sequence
- ▶ Presented difficult, and at times ambiguous contexts from which to infer whether or not a word is in a sentence; intuition is that this allows it to squeeze more info from the training data

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

- ▶ An autoregressive language model predicts  $P(x_i|x_{<i})$ ; ELMo also looks at  $P(x_i|x_{>i})$
- ▶ Consider a sequence with  $T=4$  tokens. Now consider the set of all  $4!$  permutations $Z = [1, 2, 3, 4], [1, 2, 4, 3], \dots, [4, 3, 2, 1]$ . The XLNet model is auto-regressive over all such permutations
- ▶ Samples from all possible permutations, since there are many

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Example

- ▶ Consider the permutation [3, 2, 4, 1].
- ▶ When calculating the probability of the 1st element in that order (i.e., token 3), the model has no context as the other tokens have not yet been seen.
- ▶ So the mask would be [0, 0, 0, 0].
- ▶ For the 2nd element (token 2), the mask is [0, 0, 1, 0] as its only context is token 3.
- ▶ Following that logic, the 3rd and 4th elements (tokens 4 and 1) have masks [0, 1, 1, 0] and [0, 1, 1, 1].
- ▶ Stacking all those in the token order gives the attention matrix

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Takes A Long Time to Train

ELMo	40 GPU days
BERT	450 GPU days
XLNet	2000 GPU days

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Takes A Long Time to Train

	BERT	RoBERTa	DistilBERT	XLNet
<b>Size (millions)</b>	<b>Base:</b> 110 <b>Large:</b> 340	<b>Base:</b> 110 <b>Large:</b> 340	<b>Base:</b> 66	<b>Base:</b> ~110 <b>Large:</b> ~340
<b>Training Time</b>	<b>Base:</b> 8 x V100 x 12 days* <b>Large:</b> 64 TPU Chips x 4 days (or 280 x V100 x 1 days*)	<b>Large:</b> 1024 x V100 x 1 day; 4-5 times more than BERT.	<b>Base:</b> 8 x V100 x 3.5 days; 4 times less than BERT.	<b>Large:</b> 512 TPU Chips x 2.5 days; 5 times more than BERT.
<b>Performance</b>	Outperforms state-of-the-art in Oct 2018	2-20% improvement over BERT	3% degradation from BERT	2-15% improvement over BERT
<b>Data</b>	16 GB BERT data (Books Corpus + Wikipedia). 3.3 Billion words.	<b>160 GB</b> (16 GB BERT data + 144 GB additional)	16 GB BERT data. 3.3 Billion words.	<b>Base:</b> 16 GB BERT data <b>Large:</b> 113 GB (16 GB BERT data + 97 GB additional). 33 Billion words.
<b>Method</b>	BERT (Bidirectional Transformer with MLM and NSP)	BERT without NSP**	BERT Distillation	Bidirectional Transformer with Permutation based modeling

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Longformer

- ▶ To process a sequence of  $n$  tokens requires  $n^2$  attention calculations for each attention head during the forward pass
- ▶ BERT addresses this by enforcing a hard limit of 512 (sub-word) tokens, more than enough to process the overwhelming majority of sequences in most benchmark datasets
- ▶ However, this is not sufficient if you need to attend within longer texts, like entire documents
- ▶ Still know the original word order through the Transformer's positional encodings

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Basic Idea of Longformer

- ▶ Use an attention sliding window of size  $w$
- ▶ This allows a multilayer transformer model to have a receptive field that covers the entire document
- ▶ Same intuition as the receptive field within CNNs, as convolutions are sliding windows

# Longformer

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

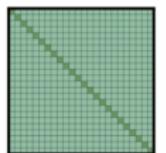
Transformers XL

XLNet

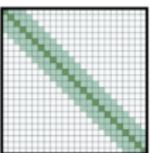
Longformer

BigBird

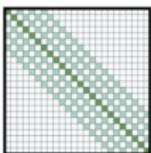
Overview and  
What to Use



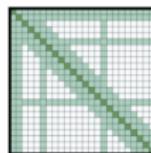
(a) Full  $n^2$  attention



(b) Sliding window attention



(c) Dilated sliding window

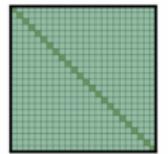
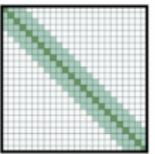


(d) Global+sliding window

Dilated sliding window allows for an even larger receptive field

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

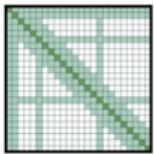
# Longformer

(a) Full  $n^2$  attention

(b) Sliding window attention



(c) Dilated sliding window

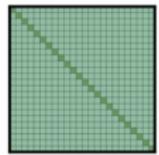
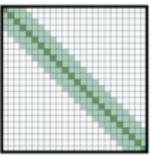


(d) Global+sliding window

The CLS token is supposed to be able to aggregate the entire sequence into a single representation to allow for classification.

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

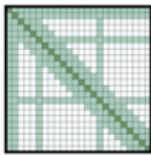
# Longformer

(a) Full  $n^2$  attention

(b) Sliding window attention



(c) Dilated sliding window



(d) Global+sliding window

Task-specific, global attention on special tokens. This attention is symmetric in that every token in the sequence can attend to the special token, as it can attend all of them.

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

**BigBird**

Overview and What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

**BigBird**

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# BigBird: Transformers for Long Text

- ▶ BigBird is an even more recent model by Google that has a similar intuition
- ▶ If we want to attend throughout long text, we need a way around exponential global attention
- ▶ Can process sequences 8x longer than BERT, because attention is  $O(n)$

# BigBird

T5

XLNet

## BigBird

## Overview and What to Use

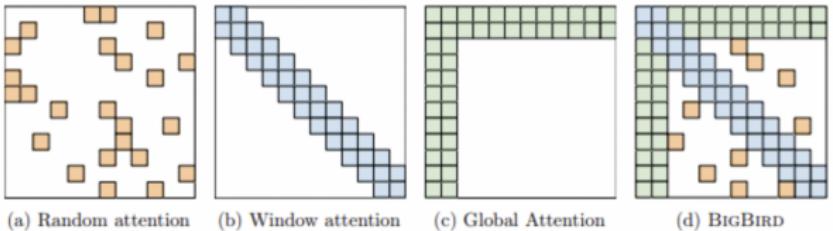


Figure 1: Building blocks of the attention mechanism used in BiGBiRD. White color indicates absence of attention. (a) random attention with  $r = 2$ , (b) sliding window attention with  $w = 3$  (c) global attention with  $g = 2$ . (d) the combined BiGBiRD model.

# Outline

Overview

Contextualized Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and What to Use

Melissa Dell

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Why Have Transformers Come to Dominate Language Modeling?

## 1. Self-attention removes locality bias

- ▶ By design, an LSTM pays the most attention to the closest contexts
- ▶ As long as it is within the sequence window, long-distance context has equal opportunity to be important in the transformer

## 2. The entire input can be processed in parallel, whereas an LSTM requires sequential processing

- ▶ The parallel compute power of GPUs can be fully harnessed, allowing for massively larger models to be trained

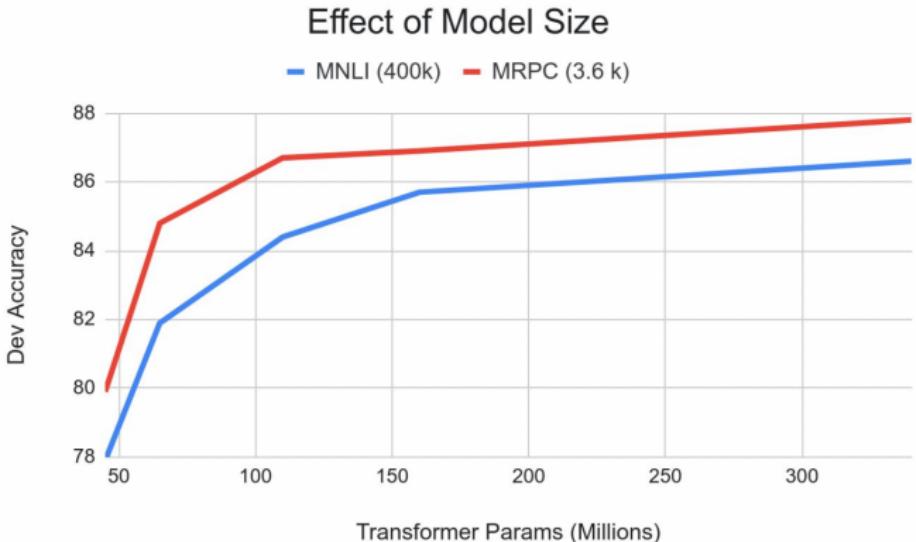
[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Key Theme: Bigger is Better

- ▶ One of the clearest lessons to emerge from this literature is that bigger models, trained on more data for longer, achieve SoTA, on a wide variety of downstream tasks
- ▶ Performance gains have not yet asymptoted, even when looking at very large models (i.e. T5 has 11 B parameters)
- ▶ This is true even when fine-tuning models to tasks with a small number of labels

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Bigger is Better



Remarkable the degree to which DL architectures have been able to solve overfitting and vanishing gradient problems

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Bigger is Better

T5 did a bunch of ablation studies to understand which features matter, i.e.

- ▶ Model size
- ▶ Amount of training data
- ▶ Domain/noise in training data
- ▶ Pre-training objectives
- ▶ Ensembling
- ▶ Details of fine-tuning
- ▶ Multi-task training

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Bigger is Better

- ▶ The only things that made much of a difference were model size and amount of training data
- ▶ Best model has 11B parameters (BERT large has 330M) trained on 120B words of cleaned Common Crawl Text
- ▶ ALBERT is light in terms of parameters, not speed of pre-training. The version that beats BERT has fewer parameters but is trained for significantly longer

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

# Bigger is Better

Model	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg	Speedup
BERT	base	108M	90.4/83.2	80.4/77.6	84.5	92.8	68.2	82.3
	large	334M	92.2/85.5	85.0/82.2	86.6	93.0	73.9	85.2
ALBERT	base	12M	89.3/82.3	80.0/77.1	81.6	90.3	64.0	80.1
	large	18M	90.6/83.9	82.3/79.4	83.5	91.7	68.5	82.4
ALBERT	xlarge	60M	92.5/86.1	86.1/83.1	86.4	92.4	74.8	85.5
	xxlarge	235M	<b>94.1/88.3</b>	<b>88.1/85.1</b>	<b>88.0</b>	<b>95.2</b>	<b>82.3</b>	<b>88.7</b>

Being light in terms of parameters is still very beneficial in terms of being able to deploy or fine-tune it on your local machine

[Overview](#)[Contextualized Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and What to Use](#)

# Language Modeling and Downstream Tasks

- ▶ Before the past few years, it was common to design different customized architectures for different types of downstream language tasks. These could be highly specialized to particular tasks
- ▶ The Transformer-based language models obviated this, providing a way to fine-tune most downstream tasks you'd want to do within the architecture of the language model (other models like ELMO had a similar idea, but post-BERT we've really seen this take off)
- ▶ This has made them enormously influential, underscoring just how much performance on downstream tasks is about having good representations (achieved through training massive models on massive data), versus about engineered domain knowledge

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

GPT2/GPT3

Transformers XL

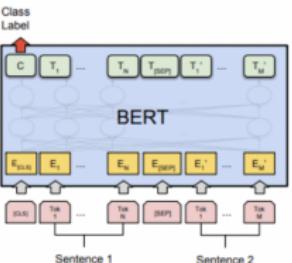
XLNet

Longformer

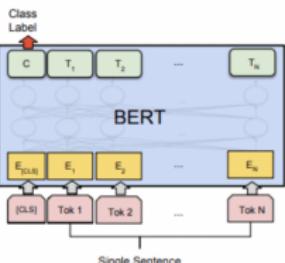
BigBird

Overview and  
What to Use

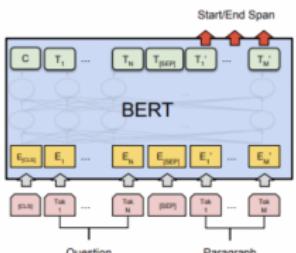
# BERT and Fine-Tuning



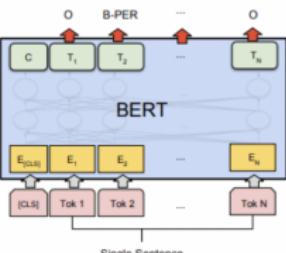
(a) Sentence Pair Classification Tasks:  
MNLI, QQP,QNLI,STS-B,MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

Overview

Contextualized  
Word Embeddings

GPT

BERT

RoBERTa

DistilBERT

ALBERT

T5

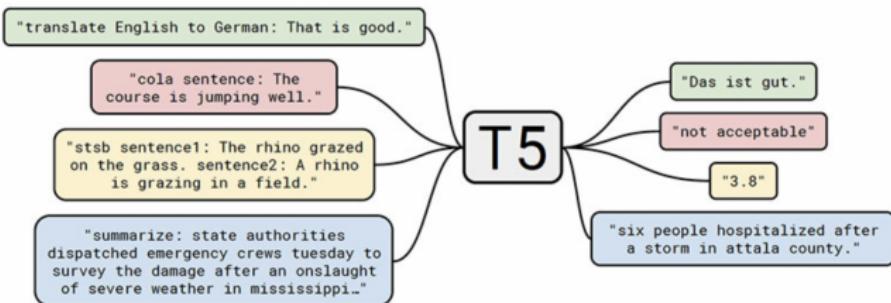
GPT2/GPT3

Transformers XL

XLNet

Longformer

BigBird

Overview and  
What to Use

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# Summary

- ▶ Bottom line: pre-trained language models work extremely well for downstream tasks
- ▶ The main thing that makes them work better is building bigger models trained for longer on more data. This is very expensive
- ▶ Fortunately, Google/Facebook/Microsoft have a lot of resources and a business incentive to pre-train them and have for the most part open-sourced their innovations (Open AI/Microsoft excepted)
- ▶ Distillation and parameter sharing (ALBERT) make it possible to do inference and potentially fine-tune on commercial hardware
- ▶ These models remove the need to design lots of customized architectures for downstream tasks

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# What to Use

- ▶ Different models have different architectures, training objectives, and training data
- ▶ Papers will say something about differences with other models and what matters but far from comprehensive (expensive to do a bunch of ablations)
- ▶ Hence, the influence of these different factors - which will matter differentially for different downstream tasks - is conflated
- ▶ You need to experiment (i.e. with pre-trained embeddings)

[Overview](#)[Contextualized  
Word Embeddings](#)[GPT](#)[BERT](#)[RoBERTa](#)[DistilBERT](#)[ALBERT](#)[T5](#)[GPT2/GPT3](#)[Transformers XL](#)[XLNet](#)[Longformer](#)[BigBird](#)[Overview and  
What to Use](#)

# What to Use?

- ▶ Is there a pre-trained version that approximates your task (i.e. RoBERTa News; BERTweet)
- ▶ ALBERT, RoBERTa, or XLNet are good places to start. If you are heavily computationally constrained, try inference with one of the lightweight ALBERT models
- ▶ Long text/documents: BigBird or perhaps Longformer
- ▶ This literature moves fast