

COMP5046

Natural Language Processing

Lecture 12: Pretrained Model

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Dr. Caren Han

Semester 1, 2021

School of Computer Science,
University of Sydney

Deep Learning



Linguistics



Language



NLP



0 LECTURE PLAN

Lecture 12: Pretrained Model

1. The Rise of the Pre-trained Model
2. BERT
3. Post BERT
4. Multimodal Pretrained Model

<https://tutorcs.com>

Assignment Project Exam Help

WeChat: cstutorcs

The Rise of the Pre-trained Model

Pre-training and Transfer Learning

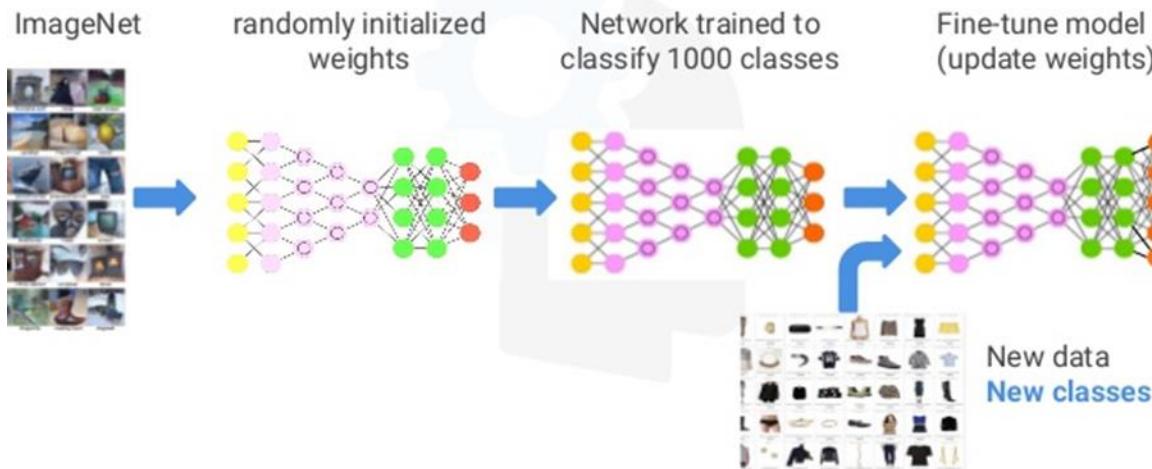
In computer vision, prove the value of transfer learning

- pre-training a neural network on a known task (i.e. ImageNet)
- performing fine-tuning
- using the trained neural network as the basis of a new purpose-specific model.

Assignment Project Exam Help

<https://tutorcs.com>

Transfer Learning
WeChat: cstutorcs



The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

Popular Pre-trained Model in NLP

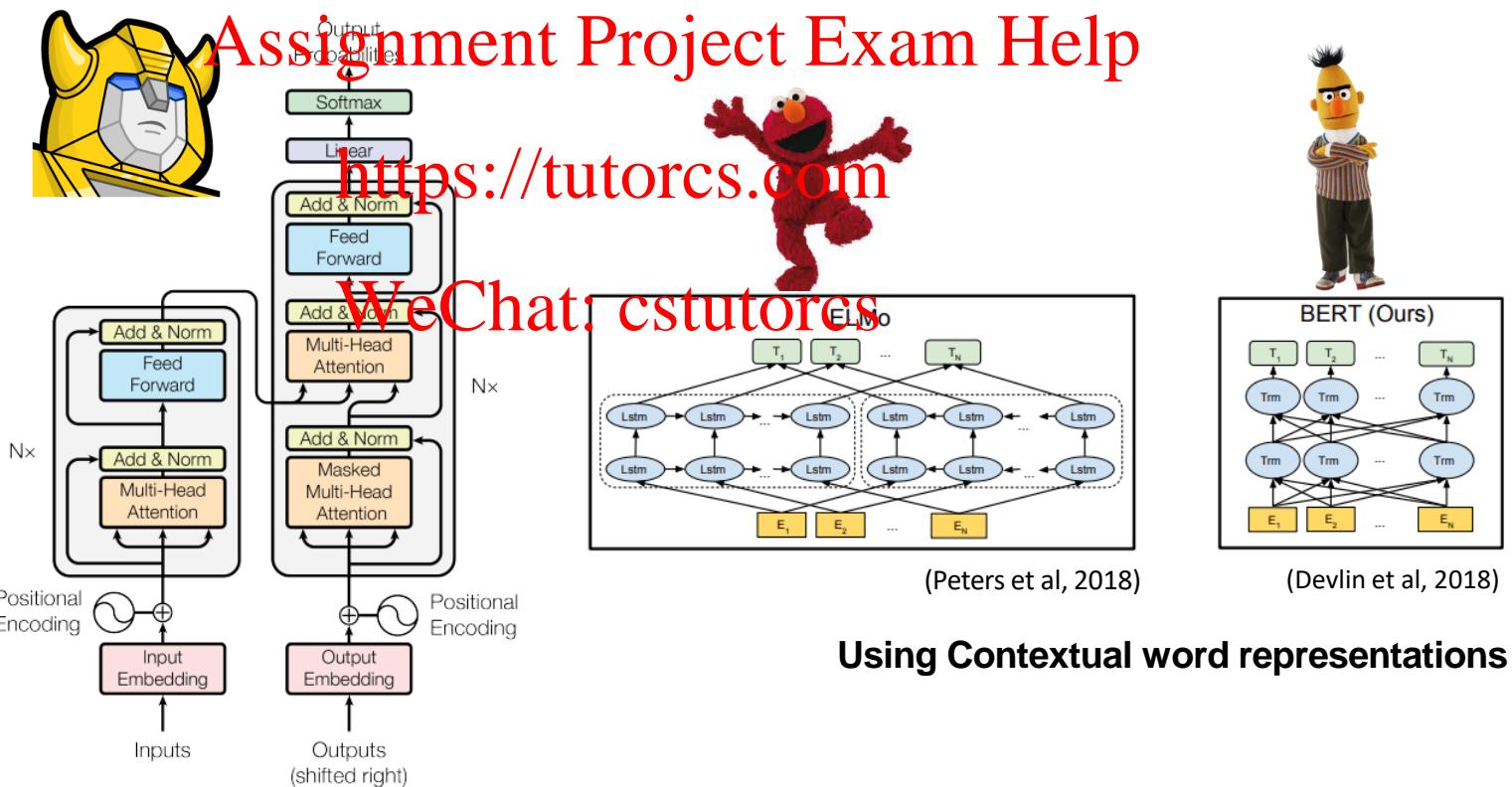


Figure 1: The Transformer - model architecture.

The Rise of the Pre-trained Model

Before we started...

Word Structure and subword models

We assume a fixed vocab of tens of thousands of words, built from the training set.

Assignment Project Exam Help

All novel words seen at test time are mapped to a single UNK.

<https://tutorcs.com>

	Word	Vocab mapping
Common words	computer	desk (index)
	play	cute (index)
Variations	coooooooooooooool	UNK (index)
misspellings	laern	UNK (index)
novel items	Transformerify	UNK (index)

The Rise of the Pre-trained Model

Before we started...

Word Structure and subword models

Many languages exhibit complex morphology, or word structure.

- The effect is more word types, each occurring fewer times.*

<https://tutorcs.com>

WeChat: cstutorcs

Conjugation of -ambia															[less ▲]				
Polarity	Persons			Persons			Classes			Classes			Classes			[less ▲]			
	1st Sg.	2nd Pl.	3rd Sg.	1st Pl.	2nd Pl.	3rd / M-wa	M-mi	4	5 Ma	6	7 Ki-vi	8	9 N	10 11/14	15/17	16 Pa	18 Mu		
Past																			
Positive	nillambia	tulambia	uliambia	milambia	mwallambia	aliambia	wallambia	uliambia	iliambia	illambia	yallambia	kiliambia	viliambia	iliambia	ziliambia	uliambia	kuliambia	pallambia	mullambia
Negative	sikuambia	hatukumbia	hukuambia	hamkuambi	a	hakuambia	hawkuambia	haukumbia	haikumbia	halikumbia	hayakumbia	hakiakumbia	havikumbia	haikumbia	hazikumbia	haukumbia	hakukumbia	hapakuambia	hamukumbia
Present																	[less ▲]		
Positive	ninambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	inaambia	linaambia	yanaambia	kinaambia	vinaambia	inaambia	zinaambia	unaambia	kunaambia	panaambia	munaambia
Negative	siambia	hatumambia	huambia	hamambia	haambia	hawaambia	haumambia	haambia	halambia	halaambia	hayaambia	hakambia	havambia	haambia	hazambia	hauambia	hakumambia	hapaambia	hamuambia
Future																	[less ▲]		
Positive	nitaambia	tutaambia	utaambia	mtaambia	ataambia	wataambia	utaambia	itaambia	itaambia	yataambia	kitaambia	vitaambia	itaambia	zitaambia	utaambia	kutaambia	pataambia	motaambia	hamutaambia
Negative	sitaambia	hatutambia	hutaambia	hamtaambia	a	hataambia	hawataambia	hautaambia	haltaambia	halytaambia	hayataambia	hakitaambia	havitaambia	hitaambia	hzitaambia	hautaambia	hakutaambia	hpataambia	hamutaambia
Subjunctive																	[less ▲]		
Positive	niambie	tuambie	uambie	mambie	aambie	waambie	uambie	lambie	lambie	yaambie	kiambie	viambie	iambie	ziambie	uambie	kuambie	paambie	muambie	musiambie
Negative	nisiambie	tusiambie	usambie	msiambie	asiambie	wasambie	usambie	isambie	isambie	ysambie	ksiambie	visambie	isambie	zisambie	usambie	kusambie	pasambie	musambie	musiambie
Present Conditional																	[less ▲]		
Positive	ningeambia	tungeambia	ungeambia	ingeambia	wangeambia	ungeambia	ingeambia	lingeambia	yingeambia	kingeambia	vingeambia	ingeambia	zingeambia	ungeambia	kungeambia	pangeambia	mungeambia		
Negative	nisingeambia	tusingeambia	usingeambia	isingeambia	msingeambia	asingeambia	wasingeambia	usingeambia	isingeambia	lyasingeambia	yasingeambia	visingeambia	isingeambia	zisingeambia	usingeambia	kusingeambia	passingeambia	musingeambia	hamungeambia
Past Conditional																	[less ▲]		
Positive	ningilambi	tungilambi	ungilambi	ngilambi	angilambi	wangilambi	ungilambi	ingilambi	ingilambi	yangilambi	kingilambi	vingilambi	ingilambi	zingilambi	unggilambi	kungilambi	pangilambi	mungilambi	
Negative	nisingilambi	tusingilambi	usingilambi	isiningilambi	msiningilambi	asingilambi	wasingilambi	usingilambi	isiningilambi	lyisingilambi	yasiningilambi	visiningilambi	isiningilambi	zisingilambi	usingilambi	kusingilambi	passiningilambi	musiningilambi	hamusingilambi
Conditional Contrary to Fact																	[less ▲]		
Positive	ningeliambi	tungettiambi	ungettiambi	ngeliambi	angeliambi	wangeliambi	ungettiambi	ingeliambi	ingeliambi	yangeliambi	kingeliambi	vingeliambi	ingeliambi	zingeliambi	ungeteliambi	kungeliambi	pangeliambi	mungeliambi	
Gnomic																	[less ▲]		
Positive	naambia	twamba	waambia	mwaambia	aambia	waambia	waambia	yaambia	yaambia	chaambia	vyambia	yaambia	zaambia	waambia	kwaambia	paambia	mwaambia		
Positive	nimeambia	tumeambia	umeambia	mmeambia	ameambia	wameambia	umeambia	imeambia	limeambia	yameambia	kimeambia	vimeambia	imeambia	zimeambia	umeambia	kumeambia	pameambia	mumeambia	

The Rise of the Pre-trained Model

Before we started...

The byte-pair encoding algorithm

Assignment Project Exam Help
Subword modeling in NLP encompasses a wide range of methods for reasoning about structure below the word level. (Parts of words, characters, bytes.)
<https://tutorcs.com>

WeChat: cstutorcs
Byte-pair encoding is a simple, effective strategy for defining a subword vocabulary.

1. *Start with a vocabulary containing only characters and an “end-of-word” symbol.*
2. *Using a corpus of text, find the most common adjacent characters “a,b”; add “ab” as a subword.*
3. *Replace instances of the character pair with the new subword; repeat until desired vocab size.*

The Rise of the Pre-trained Model

The byte-pair encoding algorithm (1994)

The byte-pair encoding algorithm: How it works?

A simple form of data compression in which the most common pair of consecutive bytes of data is replaced with a byte that does not occur within that data.

aaabdaaaabac <https://tutorcs.com>

aaabdaaaabac

Zabdzabac

ZYdZYac

XdXac

WeChat: cstutorcs

Replace $Z = aa$

Replace $Y = ab$

Replace $X = ZY$

Byte Pair	Replacement
ZY	X
ab	Y
aa	Z

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Traditional Word Encoding in NLP

Dictionary

#vocab: occurrence
low: 5,
lower: 2,
newest: 6,
widest: 3

Vocabulary

low, lower, newest, widest
<https://tutorcs.com>

*What if we have the word
'lowest' in the test set?*

OOV Issue

WeChat: cstutorcs

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Dictionary

Vocabulary

#vocab: occurrence

low: 5,

lower: 2,

newest: 6,

widest: 3

l, o, w, e, r, n, w, s, t, i, d
<https://tutorcs.com>

WeChat: cstutorcs

Character-based segmentation

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Dictionary

Vocabulary

#vocab: occurrence

low: 5,

lower: 2,

newest: 6,

widest: 3

i, o, w, e, r, n, w, s, t, i, d
<https://tutorcs.com>

WeChat: cstutorcs

Character-based segmentation

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation – 1st round

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Dictionary

Vocabulary

Assignment Project Exam Help

#vocab: occurrence

l o w: 5,

l o w e r: 2,

n e w e s t: 6,

w i d e s t: 3

Character-based segmentation

Dictionary

#vocab: occurrence

l o w: 5,

l o w e r: 2,

n e w e s t: 6,

w i d e s t: 3

https://tutorcs.com

WeChat: cstutorcs

Replace es = e s

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation – 2nd round

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Assignment Project Exam Help

Dictionary

Vocabulary

#vocab: occurrence

low: 5,

lower: 2,

newest: 6,

widest: 3

Character-based segmentation

Dictionary

#vocab: occurrence

low: 5,

lower: 2,

newest: 6,

widest: 3

https://tutorcs.com

WeChat: cstutorcs

Replace est = es t

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es, est

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation – 3rd round

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Assignment Project Exam Help

Dictionary

Vocabulary

#vocab: occurrence

l o w: 5,

l o w e r: 2,

n e w est: 6,

w i d est: 3

Character-based segmentation

Dictionary

#vocab: occurrence

l o w: 5,

l o w e r: 2,

n e w est: 6,

w i d est: 3

https://tutorcs.com

l, o, w, e, r, n, w, s, t, i, d, es, est

WeChat: cstutorcs

Replace lo = l o

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es, est, lo



The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation – 3rd round

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Assignment Project Exam Help

Dictionary

Vocabulary

#vocab: occurrence

l o w: 5,

l o w e r: 2,

n e w est: 6,

w i d est: 3

Character-based segmentation

Dictionary

#vocab: occurrence

lo w: 5,

lo w e r: 2,

n e w est: 6,

w i d est: 3

https://tutorcs.com

l, o, w, e, r, n, w, s, t, i, d, es, est

WeChat: cstutorcs

Replace lo = l o

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es, est, lo

Repeat this process until 10th round

The Rise of the Pre-trained Model

The byte-pair encoding algorithm in NLP

Subword Segmentation – Final (after 10th round)

Character/unicode → Vocabulary (Bottom up style)

→ The most common pair of consecutive bytes of data is replaced with a byte

Dictionary

Vocabulary

#vocab: occurrence

low: 5,

low e r: 2,

newest: 6,

widest: 3

l, o, w, e, r, n, w, s, t, i, d, es, est, lo,
low, ne, new, newest, wi, wid, widest

<https://tutorcs.com>

WeChat: cstutorcs

Assignment Project Exam Help
What if we have the word
'lowest' in the test set?

Training data vocabulary

lower, newest, widest, low

Vocabulary using BPE

l, o, w, e, r, n, s, t, i, d, es, est, lo,
low, ne, new, newest, wi, wid, widest

l, o, w, e, s, t

l, o, w, e, s, t

OOV : lowest

low, est

Repeat this process until 10th round

The Rise of the Pre-trained Model

Before we started...

Word Structure and subword models

Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.

<https://tutorcs.com>

In the worst case, words are split into as many subwords as they have characters

	Word	WeChat: cstutorcs	Vocab mapping
Common words	computer		Computer
	play		play
Variations	coooooool		coo##ooo#ool
misspellings	laern		la##ern##
novel items	Transformerify		Transformer##ify

The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

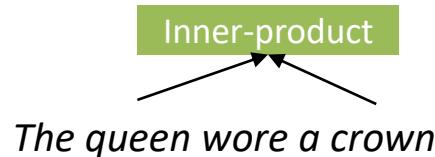
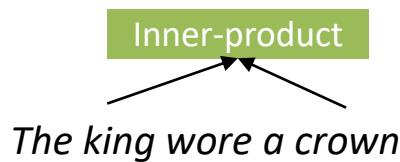
Popular Pre-trained Model: Word Embeddings

Word embeddings (e.g. word2vec) are the basis of deep learning for NLP

<https://tutorcs.com>
king [-0.5, -0.9, 1.4, ...] *queen* [0.7, 0.2, -0.5, 1.1, ...]

WeChat: cstutorcs

Word embeddings (word2vec, GloVe) are often pre-trained on text corpus from co-occurrence statistics



The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

Popular Pre-trained Model: Contextual Representations

Word embeddings (i.e. Word2Vec, fastText, GloVe) are applied in a context free manner

Assignment Project Exam Help

Step up to the **bat** — **bat** [0.7, 0.2, -0.5, 1.1, ...]

A vampire **bat** — **bat** [0.7, 0.2, -0.5, 1.1, ...]

WeChat: cstutorcs

Need to train *contextual representation* on text corpus

Step up to the **bat** — **bat** [1.1, -0.7, 0.8, 2.1, ...]

A vampire **bat** — **bat** [0.3, 0.5, -0.9, 1.3, ...]

The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

Early-Stage Pretraining in NLP

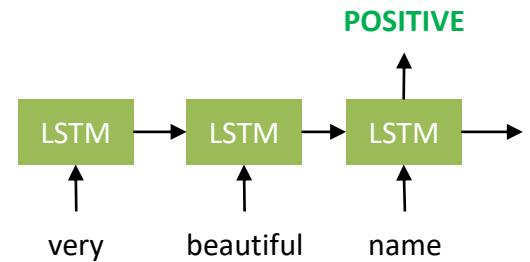
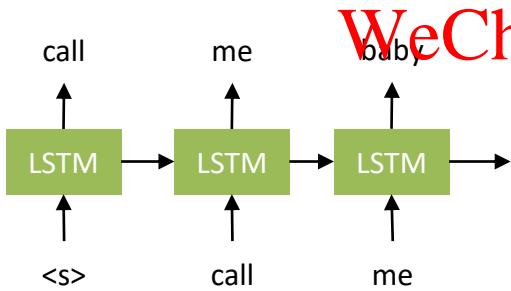
Semi-supervised Sequence Learning (Dai and Le, 2015)

Assignment Project Exam Help

Train LSTM Language Model

<https://tutorcs.com>

*Fine-tune on Classification Task
(e.g. sentiment analysis)*



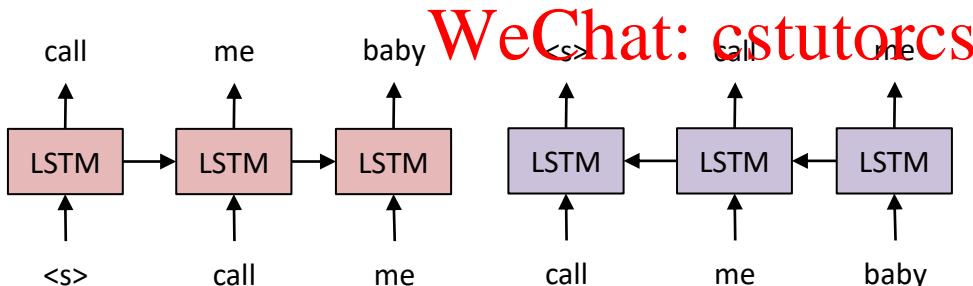
The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

ELMO: Deep Contextual Word Embeddings (Peters et al., 2018)

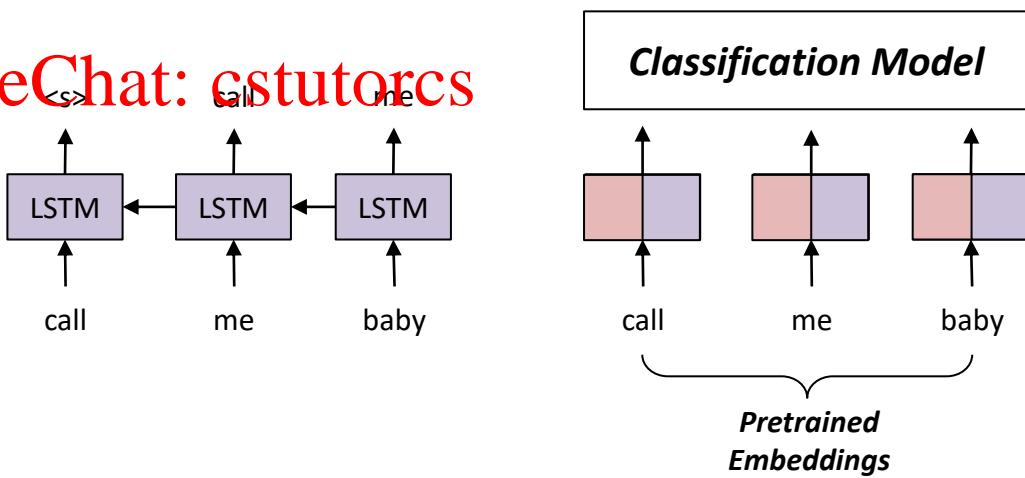
Assignment Project Exam Help

Train Separate Left-to-Right and Right-to-Left Language Models
<https://tutorcs.com>



WeChat: cstutorcs

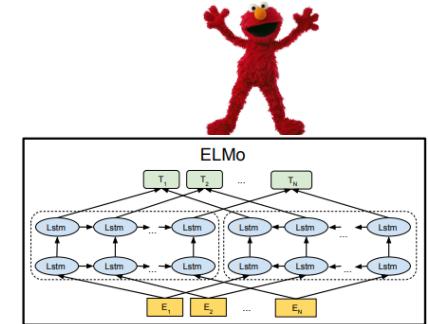
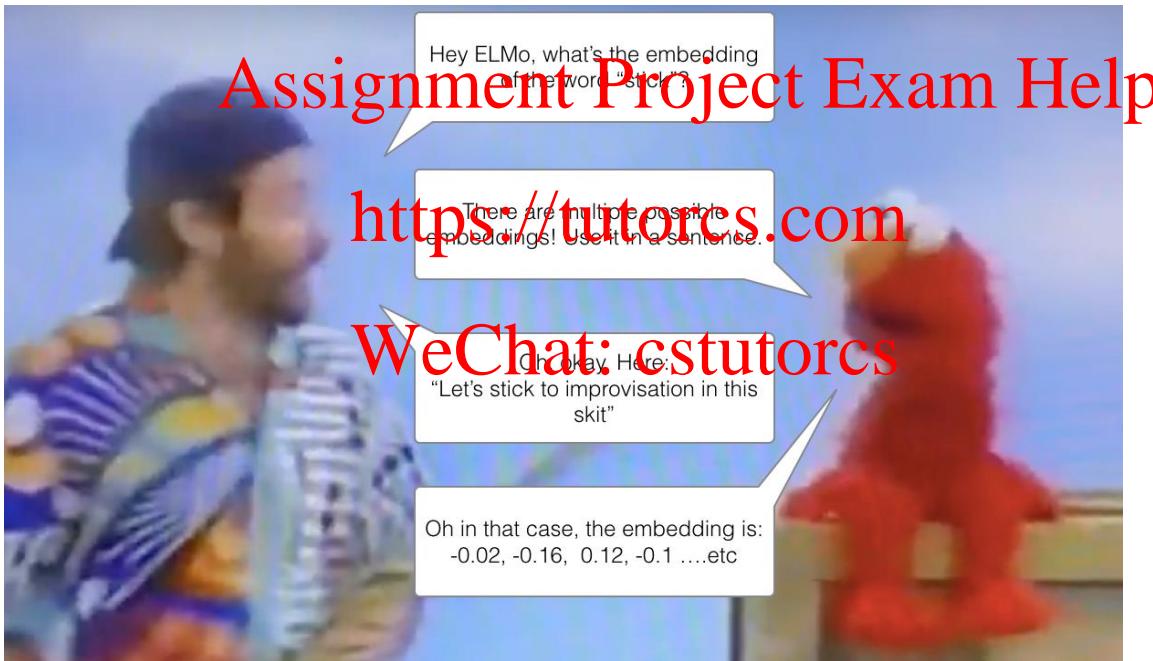
Apply as “Pretrained Embeddings”



The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

ELMo: Deep Contextual Word Embeddings (2017)



ELMo provided a significant step towards pre-training in the context of NLP. Let's dig in what the ELMo's big secret is!

The Rise of the Pre-trained Model



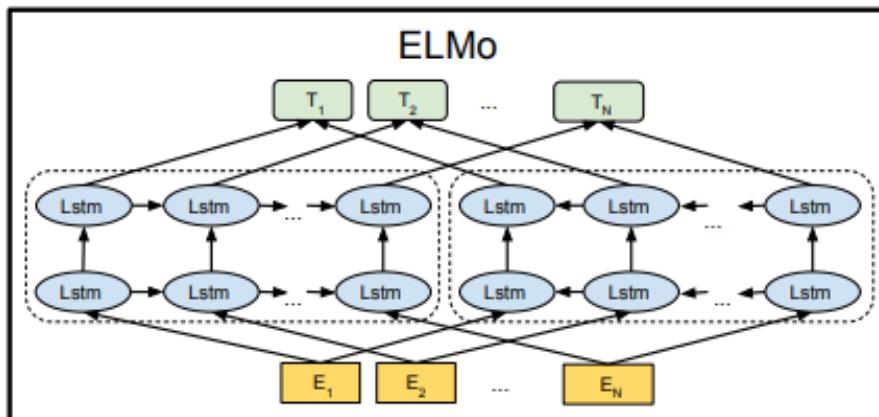
Pre-training and Transfer Learning in NLP

ELMo: Deep Contextual Word Embeddings (2017)

ELMo gained its ability to understand from being trained to predict the next word in a sequence of words, Language Modeling Tasks. This is convenient because we have vast amounts of text data that such a model can learn from without needing labels.

<https://tutorcs.com>

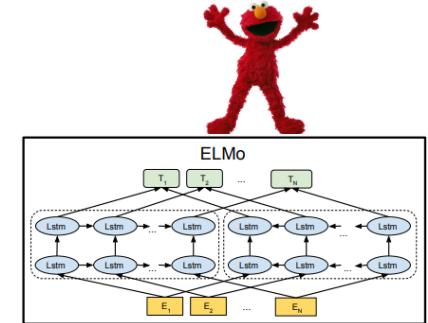
WeChat: cstutorcs



The Rise of the Pre-trained Model

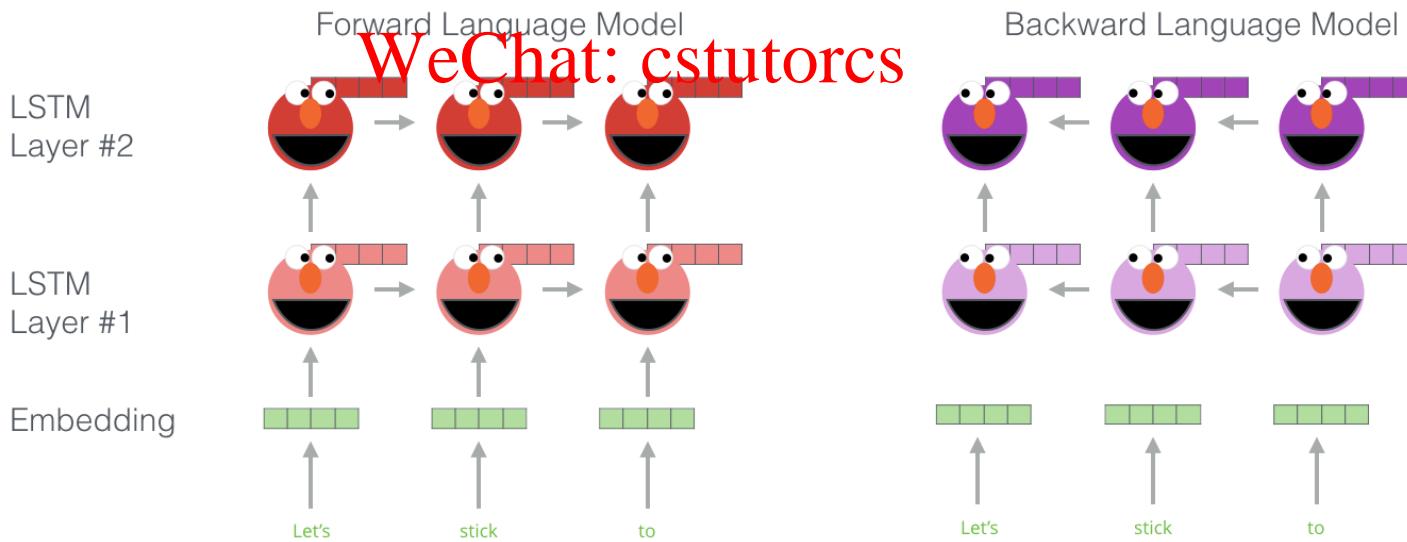
Pre-training and Transfer Learning in NLP

ELMo: Deep Contextual Word Embeddings (2017)



We can see the hidden state of each unrolled LSTM step popping out from behind ELMo's head. Those come in handy in the embedding process after this pre-training is done.

ELMo goes a step further and trains a bi-directional LSTM – so that its language model doesn't only have a sense of the next word, but also the previous word.



The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

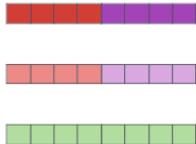
ELMo: Deep Contextualized Word Embeddings (2017)

ELMo comes with the contextualized embedding through grouping together the hidden states (and initial embedding) in a certain way (concatenation followed by weighted summation).

Assignment Project Exam Help

Embedding of “stick” in “Let’s stick to” - Step #2

1- Concatenate hidden layers



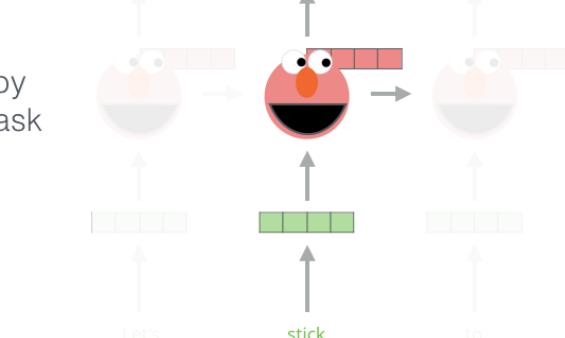
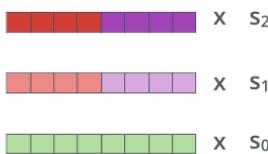
Forward Language Model



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

The Rise of the Pre-trained Model

Pre-training and Transfer Learning in NLP

Early-Stage Pretraining in NLP

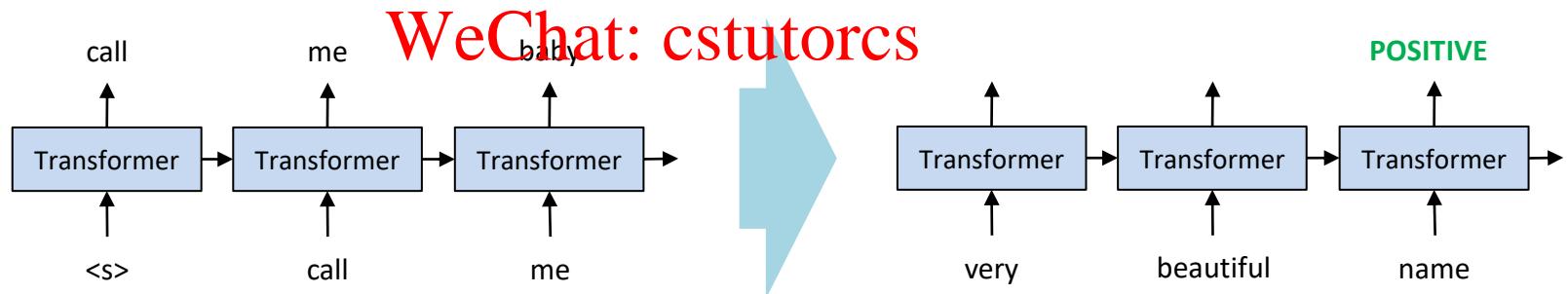
Improving Language Understanding by Generative Pre-Training (2018)

Assignment Project Exam Help

*Train Deep (12 layer) Transformer
Language Model*

<https://tutorcs.com>

*Fine-tune on Classification Task
(e.g. sentiment analysis)*



The Rise of the Pre-trained Model

Transformer (Recap)



1. Encoder

A stack of **N=6** identical layers.

Each layer with two sub-layers:

1. Multi-head self-attention mechanism

2. Position-wise fully connected feed-forward network

<https://tutorcs.com>

2. Decoder

WeChat: cstutorcs

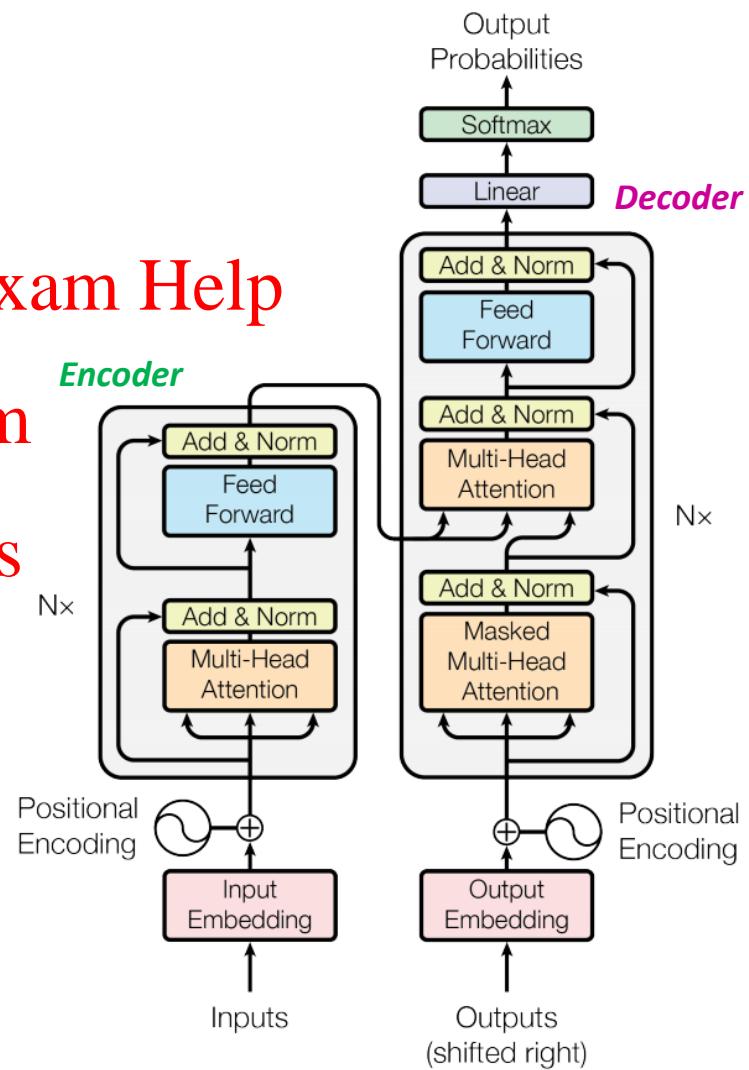
A stack of **N=6** identical layers.

Each layer with three sub-layers:

1. Multi-head self-attention mechanism

2. Position-wise fully connected feed-forward network

3. Masked Multi-head self-attention



The transformer – model architecture

The Rise of the Pre-trained Model

Transformer (Recap)



Multi-head Attention

- Models context

Assignment Project Exam Help

Feed-forward layers

- Computes non-linear hierarchical features

<https://tutorcs.com>

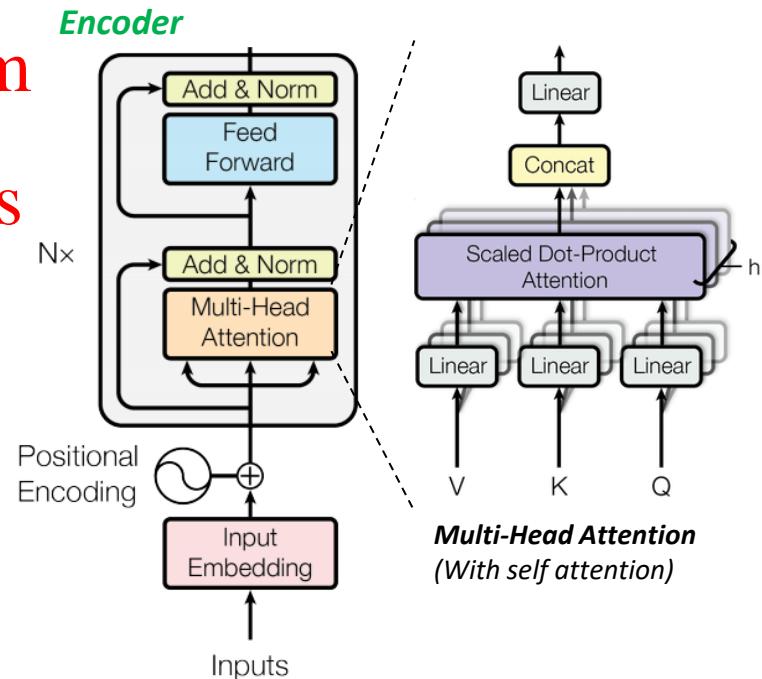
Layer norm and residuals

- Makes training deep networks healthy

WeChat: cstutorcs

Positional Embeddings

- Allows model to learn relative positioning



The transformer – Encoder

The Rise of the Pre-trained Model

Transformer (Recap)



Multi-head Attention

- Models context

Assignment Project Exam Help

Feed-forward layers

- Computes non-linear hierarchical features

<https://tutorcs.com>

Layer norm and residuals

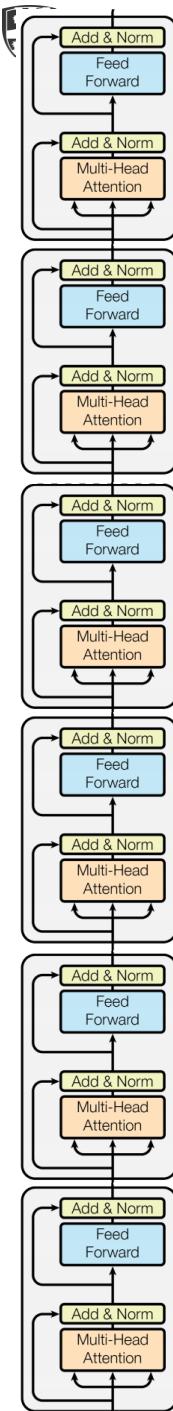
- Makes training deep networks healthy

WeChat: cstutorcs

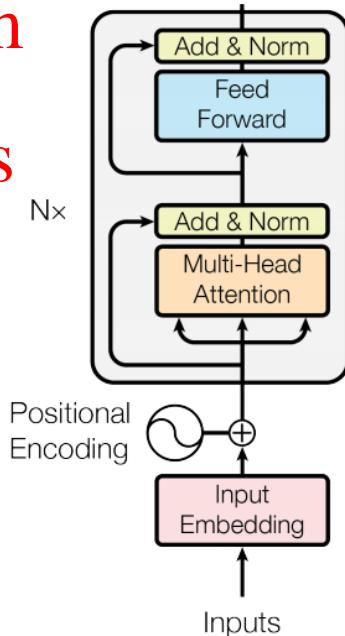
Positional Embeddings

- Allows model to learn relative positioning

Blocks are repeated 6 or more times (in vertical stack)



Encoder



The transformer – Encoder

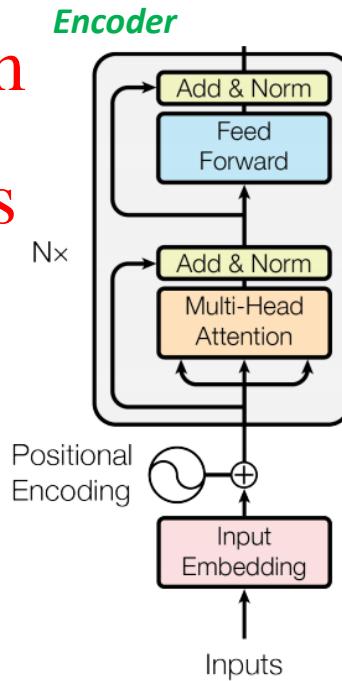
The Rise of the Pre-trained Model

Transformer VS LSTM

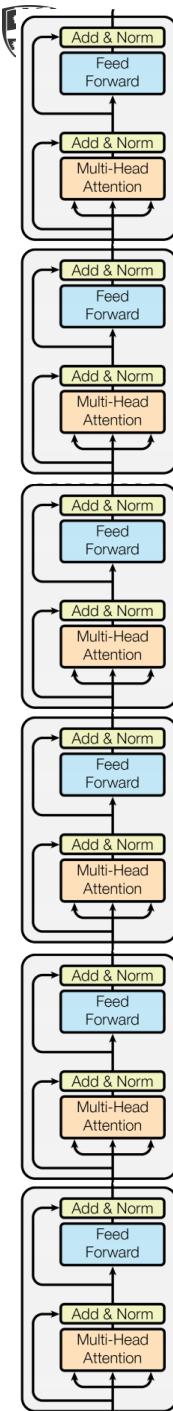
1. *Self-Attention == no locality bias*
 - *Long distance context has “equal opportunity”*
2. *Single multiplication per layer == efficiency on TPU*
 - *Effective batch size is number of words, not sequences.* <https://tutorcs.com>

Assignment Project Exam Help

WeChat: cstutorcs



The transformer – Encoder



Assignment Project Exam Help

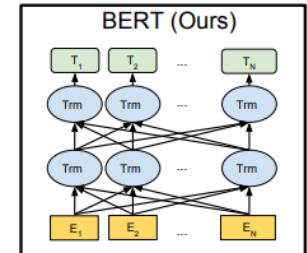
2 BERT <https://tutorcs.com>

WeChat: cstutorcs

BERT



BERT



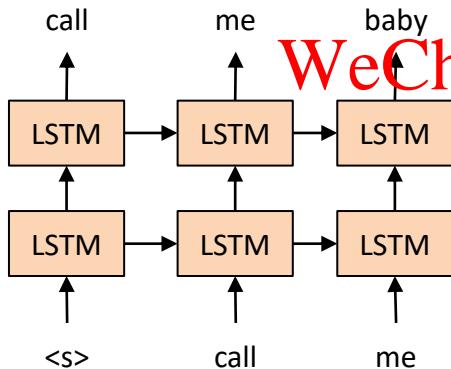
Problem with Previous Approaches

Problem: Language models only use left context or right context, but language understanding is bidirectional

Assignment Project Exam Help

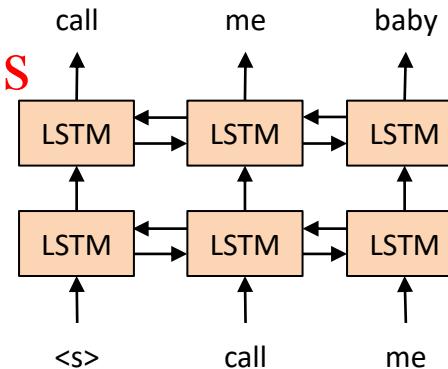
Unidirectional context
Build representation incrementally

<https://tutorcs.com>



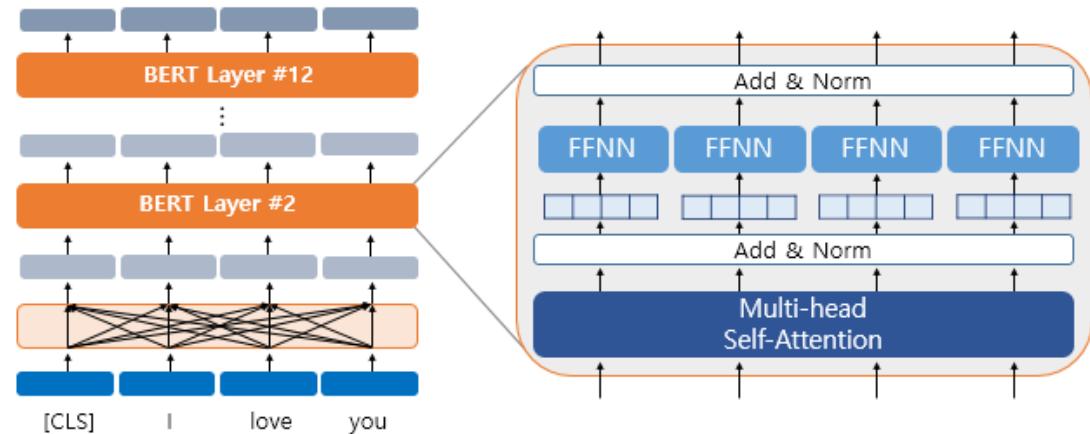
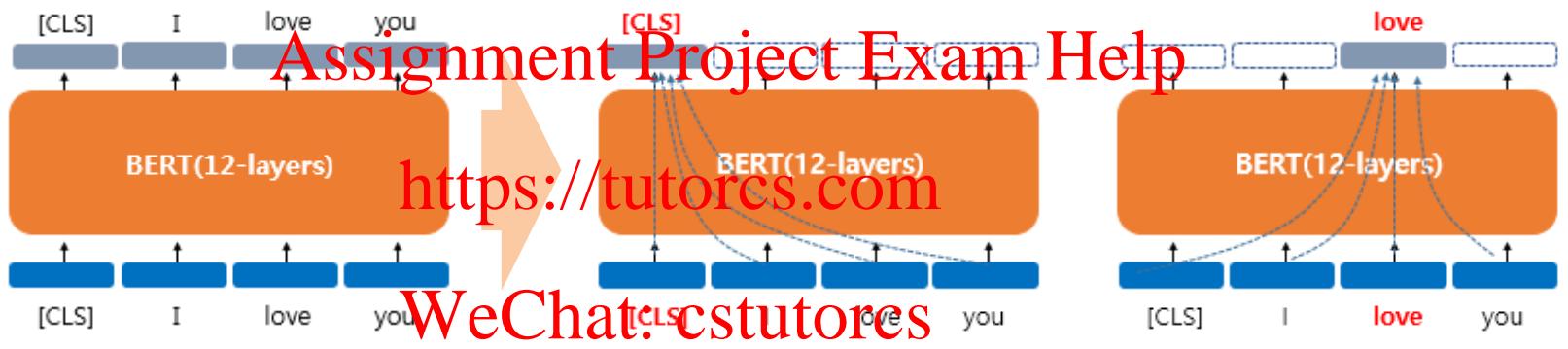
WeChat: cstutorcs

Bidirectional context
Word can "see themselves"



2 BERT

How the BERT is working (Brief)



Pre-training and Transfer Learning in NLP

BERT: Input Representation

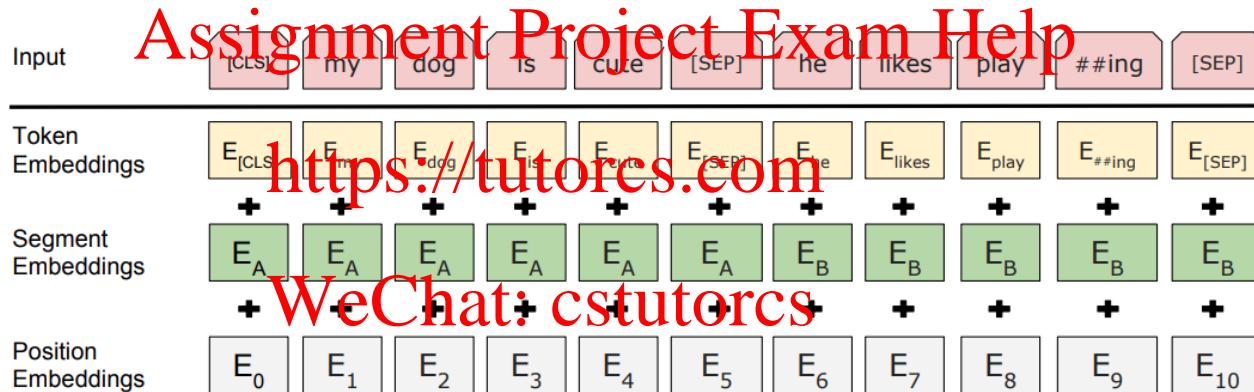


Figure 2: BERT input representation. The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

Word-piece Token Embedding

Word Structure and subword models

Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.

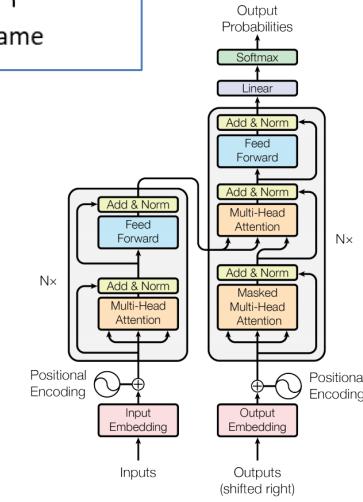
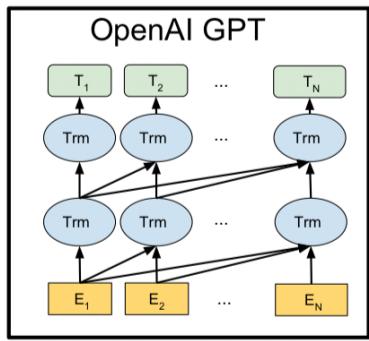
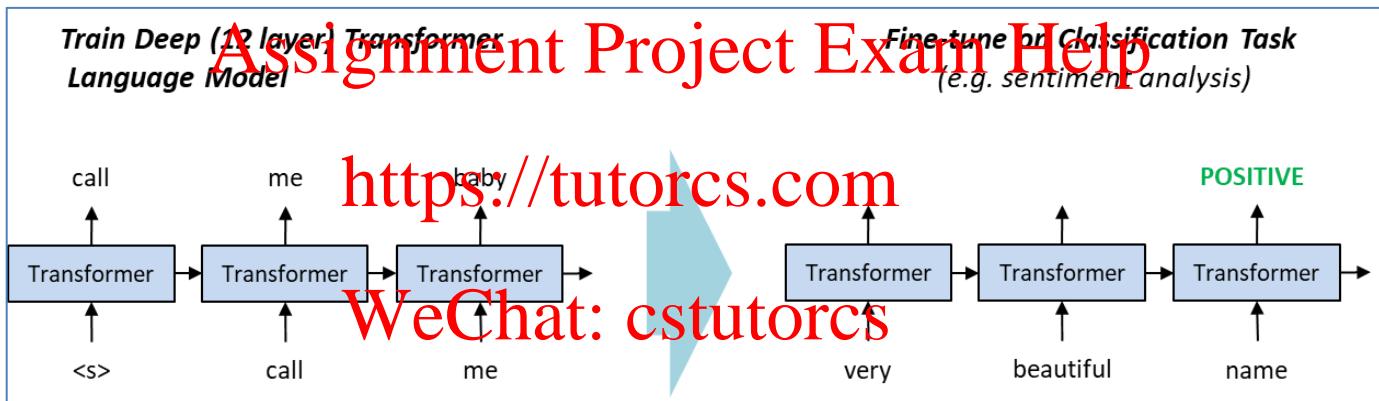
<https://tutorcs.com>
In the worst case, words are split into as many subwords as they have characters

	Word	Vocab mapping
Common words	computer	Computer
	play	play
Variations	coooooool	coo##ooo#ool
misspellings	laern	la##ern##
novel items	Transformerify	Transformer##ify

2 BERT

Pretraining for BERT

Remember GPT? Pretraining the Language Model



Pretraining: Masked Language Model (LM)

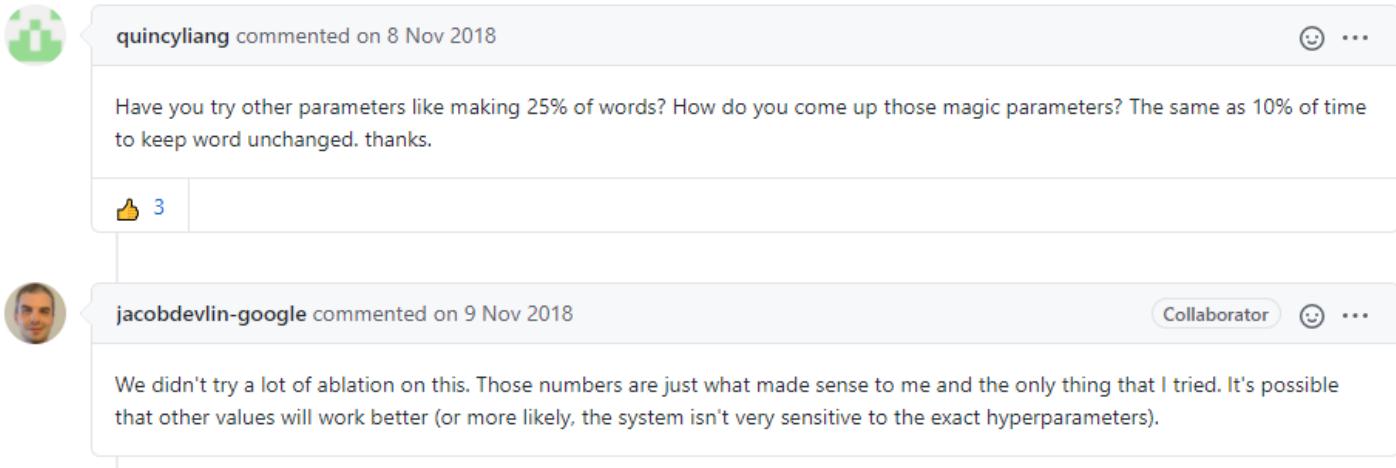
*Mask out $k\%$ of the input words, and then predict the masked words
(use $k=15\%$)*

Assignment Project Exam Help

store gallon
<https://tutorcs.com>

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

WeChat: cstutorcs



quincyliang commented on 8 Nov 2018

Have you try other parameters like making 25% of words? How do you come up those magic parameters? The same as 10% of time to keep word unchanged. thanks.

3

jacobdevlin-google commented on 9 Nov 2018

Collaborator

We didn't try a lot of ablation on this. Those numbers are just what made sense to me and the only thing that I tried. It's possible that other values will work better (or more likely, the system isn't very sensitive to the exact hyperparameters).

Pre-training and Transfer Learning in NLP

BERT: Masked language Model

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

Assignment Project Exam Help

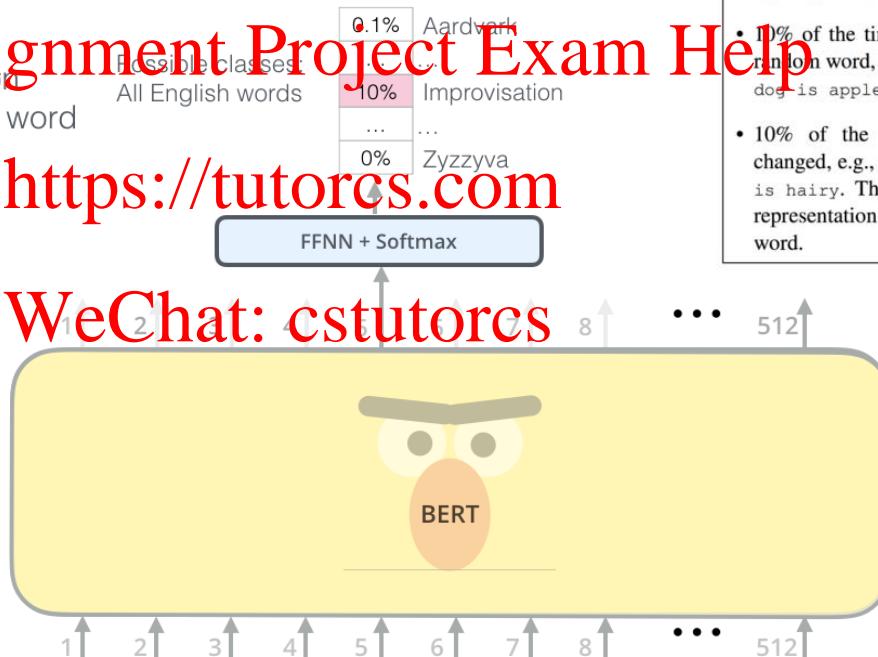
<https://tutorcs.com>

WeChat: cstutorcs

Randomly mask 15% of tokens

Input

[CLS] Let's stick to improvisation in this skit



With those 15%

- Rather than *always* replacing the chosen words with [MASK], the data generator will do the following:
- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this is to bias the representation towards the actual observed word.

[CLS] Let's stick to improvisation in this skit

Pre-training and Transfer Learning in NLP

BERT: Next Sentence Prediction

Predict likelihood
that sentence B
belongs after
sentence A

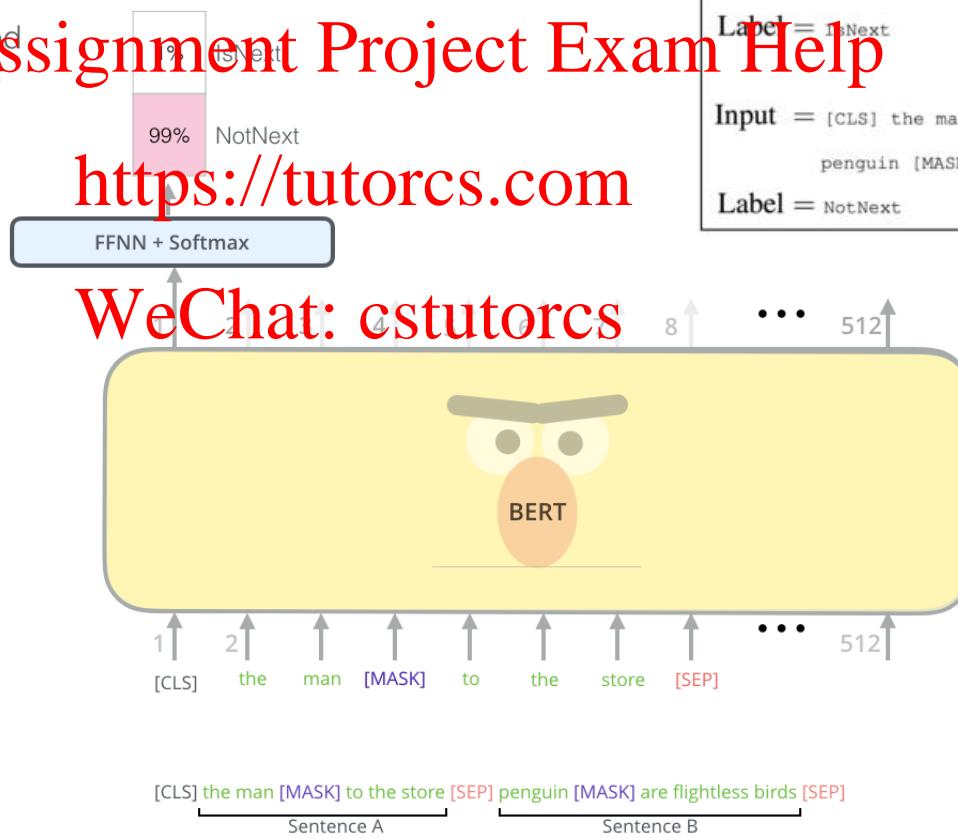
Assignment Project Exam Help

<https://tutorcs.com>

FFNN + Softmax

Tokenized Input

Input



Pre-training and Transfer Learning in NLP

BERT: Input Representation

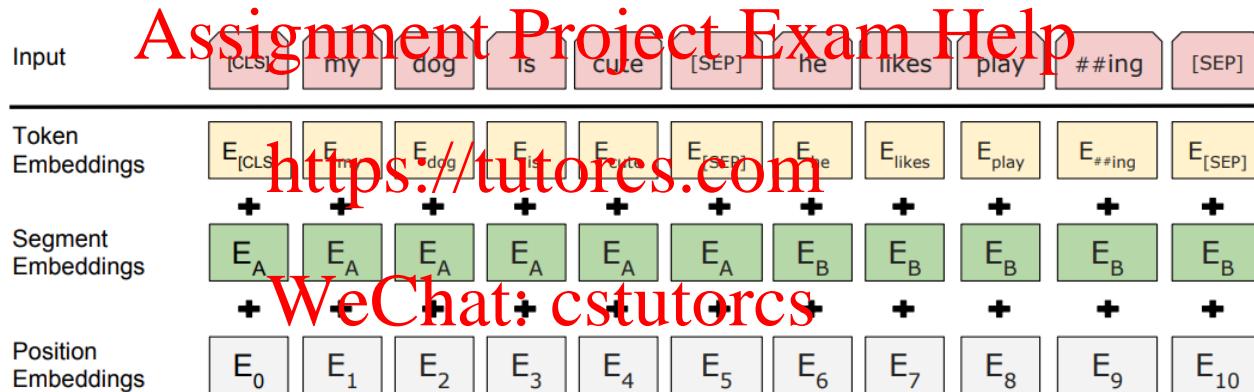


Figure 2: BERT input representation. The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

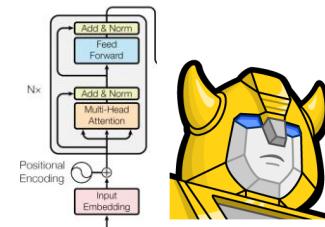
Model Details

- *Data: Wikipedia (2.5B words) + BookCorpus (800M words)*
- *Batch Size: 131,072 words (1024 sequences * 128 length or 256 sequences * 512 length)*
- *Training Time: 1M steps (~40 epochs)*
- *Optimizer: AdamW, 1e-4 learning rate, linear decay*

Assignment Project Exam Help

WeChat: cstutorcs

- *BERT-Base: 12-layer, 768-hidden, 12-head*
- *BERT-Large: 24-layer, 1024-hidden, 16-head*
- *Trained on 4x4 or 8x8 TPU slice for 4 days*



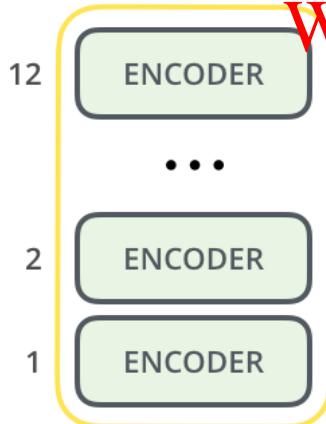
Pre-training and Transfer Learning in NLP

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

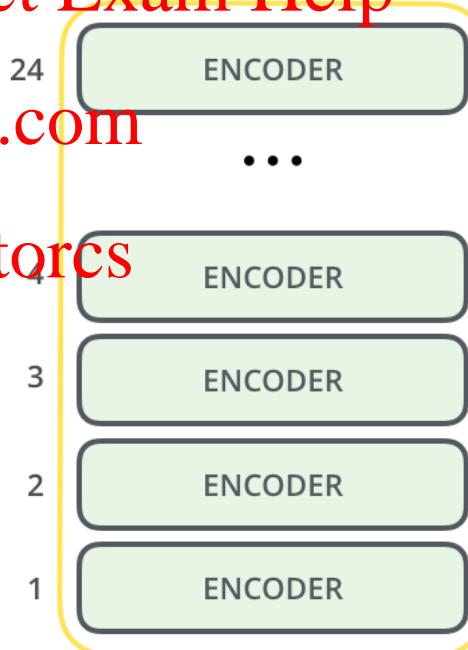
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



BERT_{BASE}



BERT_{LARGE}

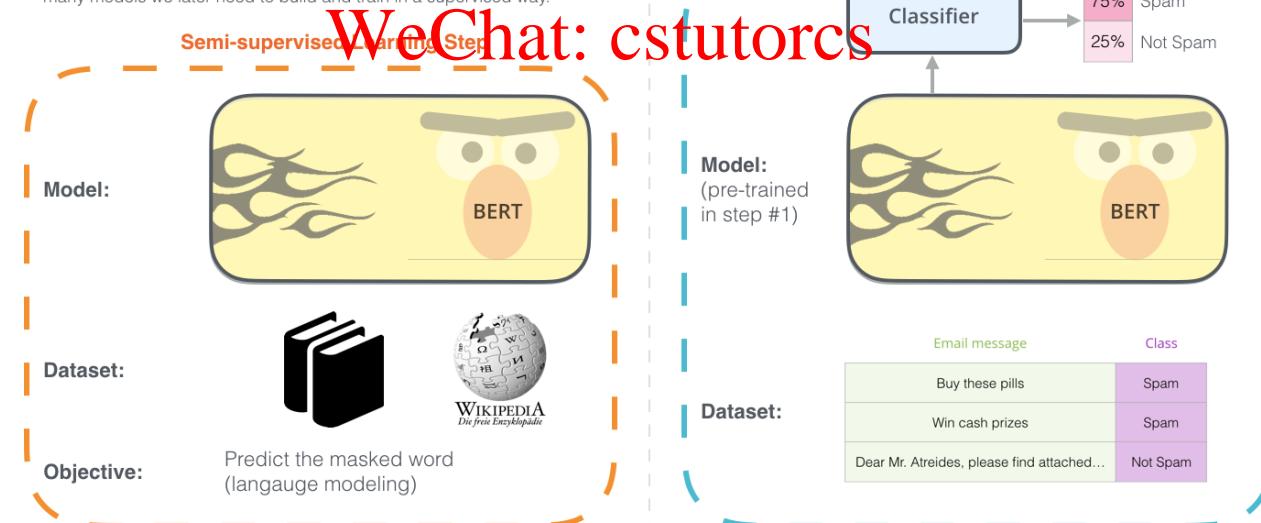
2 BERT

Pre-training and Transfer Learning in NLP

The two steps of how BERT is developed. Download the model pre-trained in step 1 (trained on un-annotated data), and only worry about fine-tuning it for step 2.

- 1 - **Semi-supervised** training on large amounts of text (books, wikipedia, etc).

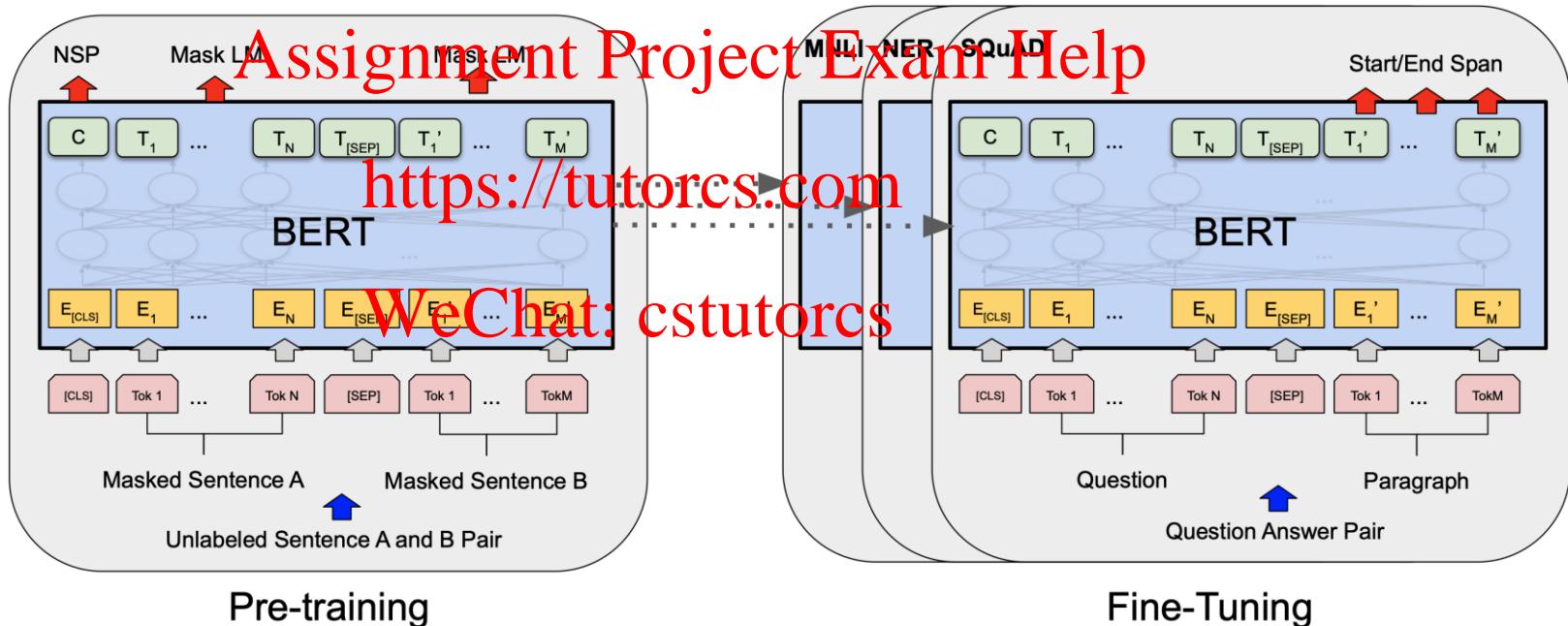
The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



Pretraining

Fine-Tuning

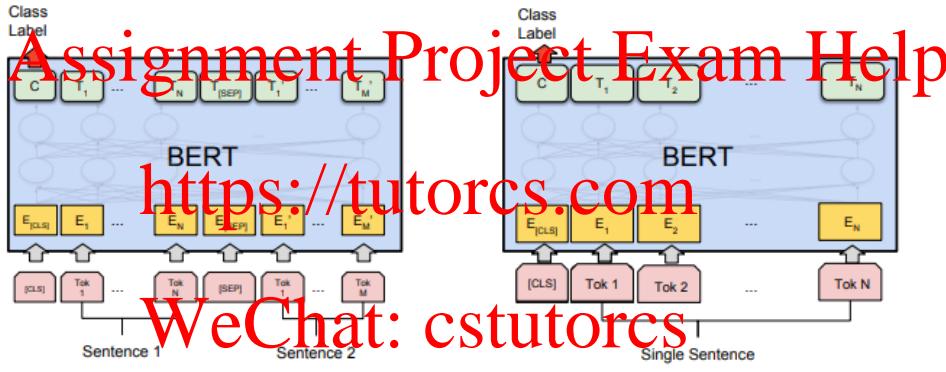
Fine-Tuning Procedure



2 BERT

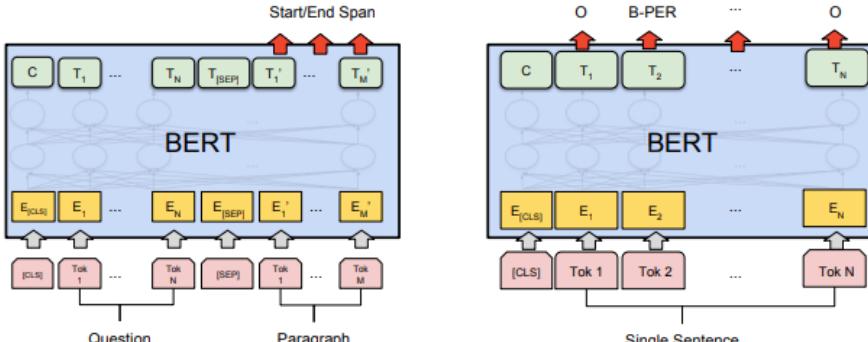
Fine-Tuning Procedure

The following shows a number of ways to use BERT for different 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

Assignment Project Exam Help
<https://tutorcs.com>
 WeChat: cstutorcs

2 BERT

Accuracy... Performance

System	MNLI (num 392k)	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attr	66.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.1	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT_{BASE} = (L=12, H=768, A=12); BERT_{LARGE} = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

Assignment Project Exam Help

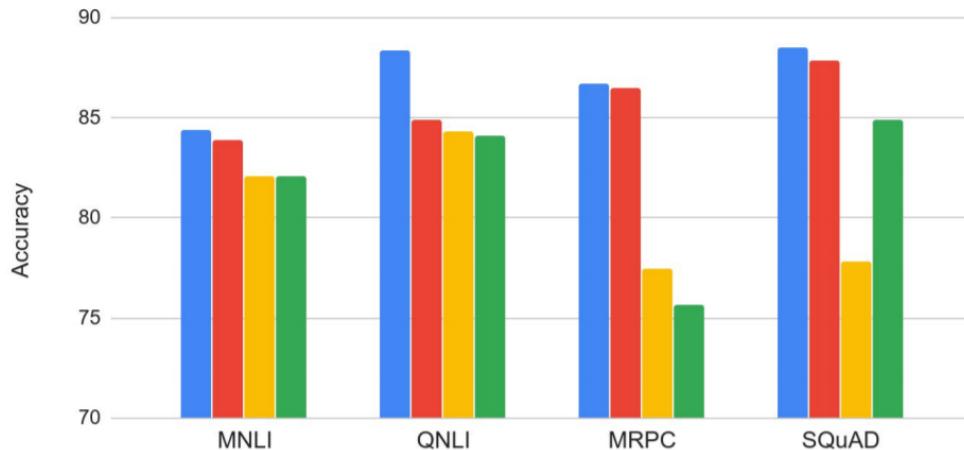
<https://tutorcs.com>
WeChat: cstutorcs

Effect of Pre-training Task

- *Masked LM (compared to left-to-right LM) is very important on some tasks, Next Sentence Prediction is important on other tasks.*
- *Left-to-right model does very poorly on word-level task (SQuAD), although this is mitigated by Bi-LSTM*

<https://tutorcs.com>

Effect of Pre-training Task
WeChat: cstutortores

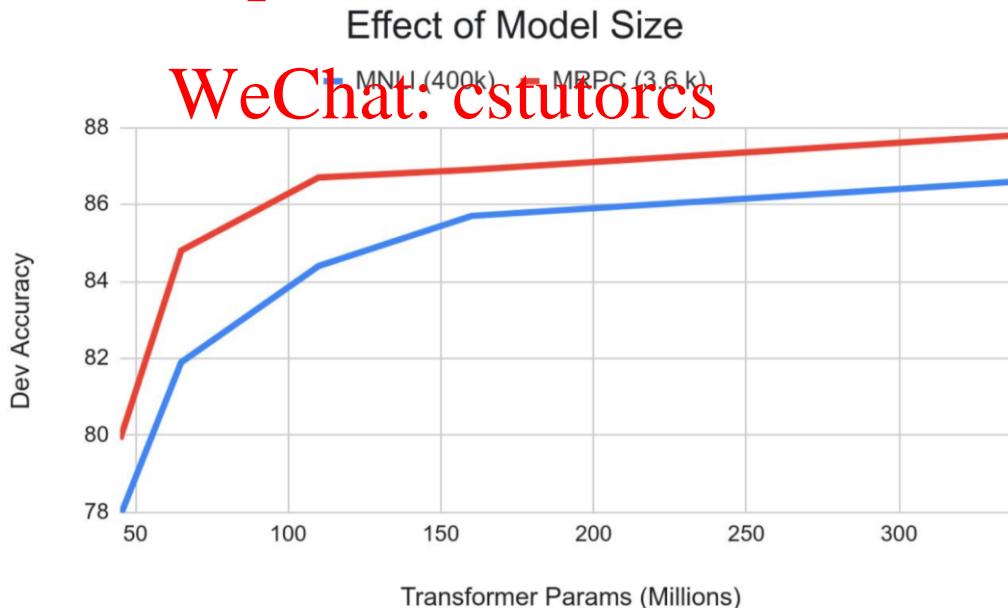


2 BERT

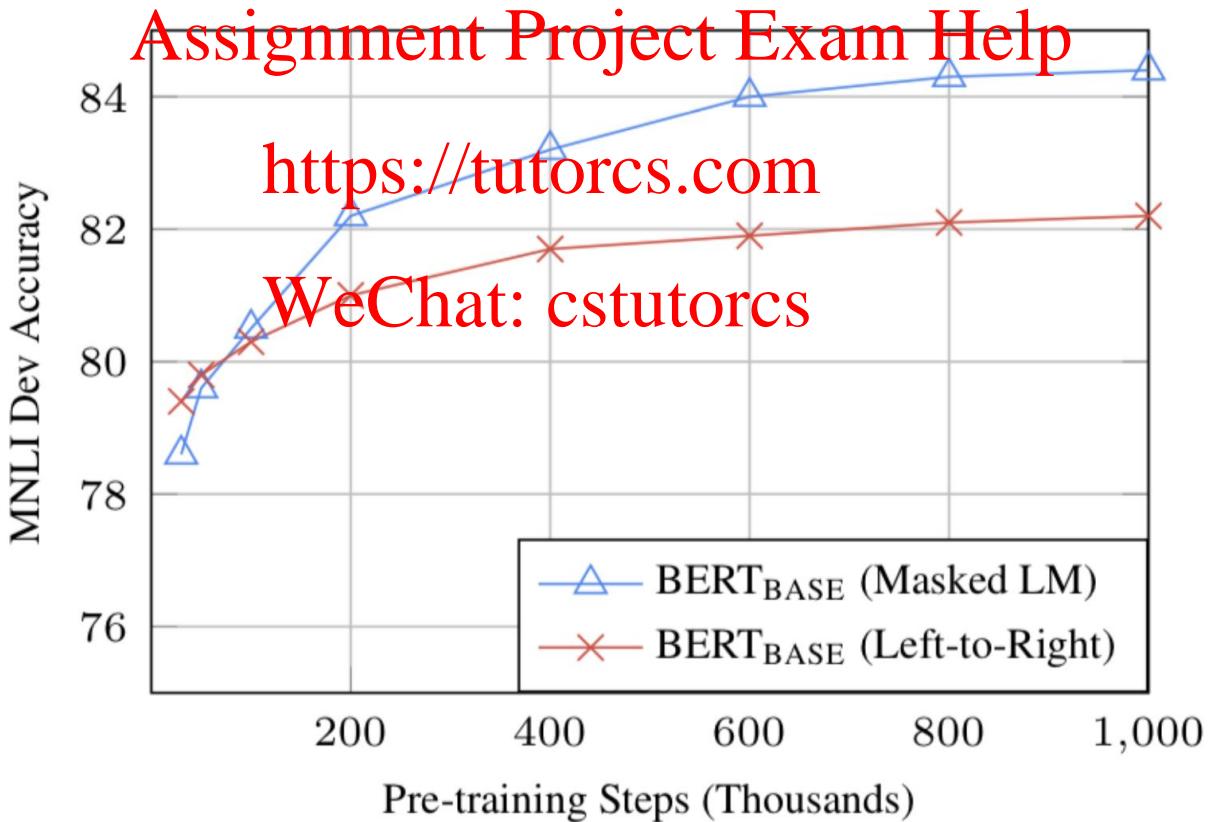
Effect of Model Size

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

<https://tutorcs.com>



2 BERT

Resource! Resource!*TPUs... and Resources..***BERT-Base:** 4 Cloud TPUs (16 TPU chips total) in 4 days**BERT-Large:** 16 Cloud TPUs (64 TPU chips total)

Questions

- *Why did no one think of this before?*
- *Better Question: Why wasn't contextual pre-training popular before 2018 with ELMo?*
- *Good Results on pre-training is > 1,000 x to 100,000 more expensive than supervised training.*

Assignment Project Exam Help

WeChat: cstutorcs

<https://tutorcs.com>

2 BERT

Questions

- *The model must be learning more than “contextual embeddings”*
- *Alternate interpretation: Predicting missing words (or next words) requires learning many types of language understanding features*
 - *Syntax, semantics, pragmatics, coreference, etc.*
- *Implication: Pre-trained model is much bigger than it needs to be to solve specific task*
- *Task-specific model distillation words very well*

2 BERT

More advanced pre-trained model

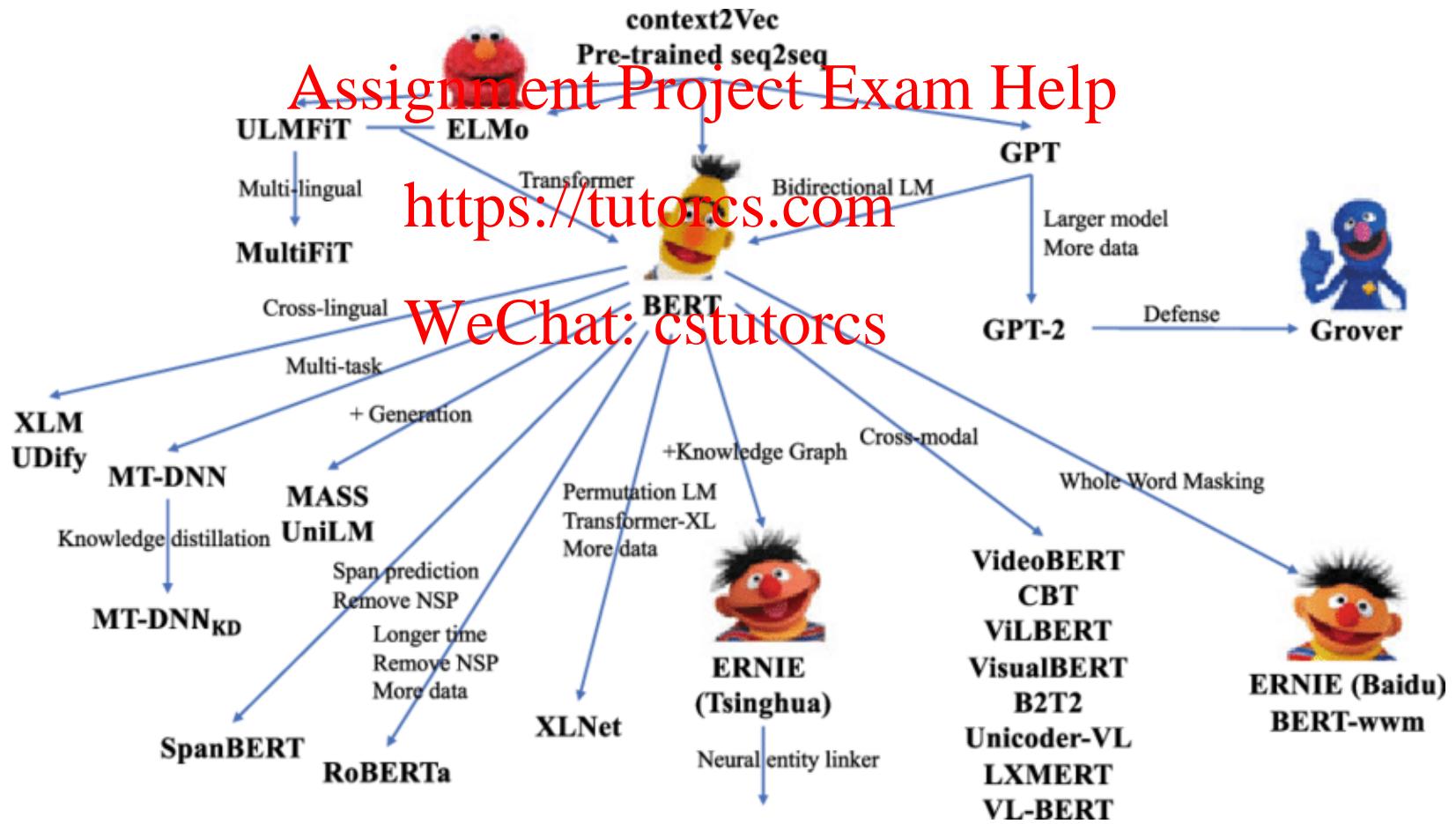
	BERT	RoBERTa	DistilBERT	XLNet
Size (millions)	Base: 110 Large: 340	Base: 110 Large: 340	Base: 66	Base: ~110 Large: ~340
Training Time	Base: 8 x V100 x 12 days* Large: 64 TPU Chips x 1 day (or 280 x V100 x 1 days*)	Large: 1024 x V100 x 1 day; 4-5 times more than BERT.	Base: 8 x V100 x 3.5 days; 4 times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

Assignment Project Exam Help

3 Post BERT
<https://tutorcs.com>

WeChat: cstutorcs

Pretrained Model Map



RoBERTa

A Robustly Optimized BERT Pretraining Approach (Liu et al, University of Washington and Facebook, 2019)

Assignment Project Exam Help

- *Trained BERT for more epochs and/or on more data*
 - *Showed that more epochs alone helps, even on same data*
 - *More data also helps*

WeChat: cstutorcs

- *Improved masking and pre-training data slightly*

	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT _{LARGE}	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet _{LARGE}	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	90.2/90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	91.3	-

XLNet

Generalized Autoregressive Pretraining for Language Understanding
(Yang et al, CMU and Google, 2019)

Assignment Project Exam Help

Innovation 1: Relative Position Embedding

- Sentence: ~~Caren ate a hot dog~~

WeChat: cstutorcs

- Absolute Attention: “How much should dog attend to hot (in any position), and how much should dog in position 4 attend to the word in position 3?”
- Relative Attention: “How much should dog attend to hot (in any position) and how much should dog attend to the previous word?”

XLNet

*Generalized Autoregressive Pretraining for Language Understanding
(Yang et al, CMU and Google, 2019)*

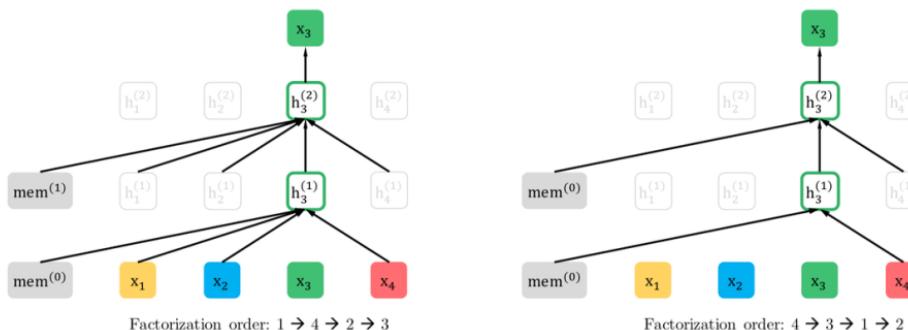
Assignment Project Exam Help

Innovation 2: Permutation Language Modeling

<https://tutorcs.com>

- In a left-to-right language model, every word is predicted based on all of the words to its left
- Instead: Randomly permute the order for every training sentence
- Equivalent to masking, but many more predictions per sentence
- Can be done efficiently with Transformers

WeChat: cstutorcs



XLNet

Generalized Autoregressive Pretraining for Language Understanding
(Yang et al, CMU and Google, 2019)

Assignment Project Exam Help

- Also used more data and bigger models, but showed that innovations improved on BERT even with same data and model size
- XLNet results:

WeChat: cstutorcs

Model	MNLI	QNLI	QQP	RTE	SST-2	MRPC	CoLA	STS-B
<i>Single-task single models on dev</i>								
BERT [2]	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0
RoBERTa [21]	90.2/90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4
XLNet	90.8/90.8	94.9	92.3	85.9	97.0	90.8	69.0	92.5

ALBERT

*A Lite BERT for Self-supervised Learning of Language Representations
(Lan et al, Google and TTI Chicago, 2019)*

Assignment Project Exam Help

Innovation 1: Factorized Embedding Parameterisation

- Use small embedding size (e.g., 128) and then project it to Transformer hidden size (e.g., 1024) with parameter matrix

WeChat: cstutorcs



ALBERT

Innovation #2: Cross-layer parameter sharing

- Share all parameters between Transformer layers

Assignment Project Exam Help

Results:

<https://tutorcs.com>

Models	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS
<i>Single-task single models on dev</i>								
BERT-large	88.6	92.3	91.3	70.4	93.2	88.0	60.6	90.0
XLNet-large	89.8	93.9	91.8	83.8	95.6	89.2	63.6	91.8
RoBERTa-large	90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4
ALBERT (1M)	90.4	95.2	92.0	88.1	96.8	90.2	68.7	92.7
ALBERT (1.5M)	90.8	95.3	92.2	89.2	96.9	90.9	71.4	93.0

ALBERT is light in terms of parameters, not speed

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
	xlarge	60M	92.5/86.1	86.1/83.1	86.4	92.4	74.8	85.5
	xxlarge	235M	94.1/88.3	88.1/85.1	88.0	95.2	82.3	88.7

T5

Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer (Raffel et al, Google, 2019)

Assignment Project Exam Help

Ablated many aspects of pre-training:

- *Model Size* <https://tutorcs.com>
- *Amount of Training data*
- *Domain/cleanness of training data* WeChat: cstutorcs
- *Pre-training objective details (e.g. span length of masked text)*
- *Ensembling*
- *Finetuning recipe (e.g. only allowing certain layers to finetune)*
- *Multi-task training*

T5

Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer (Raffel et al, Google, 2019)

Assignment Project Exam Help

Conclusions:

- *Scaling up model size and amount of training data helps a lot*
- *Best model is 11B parameters (BERT-Large is 330M), trained on 120B words of cleaned common crawl text*
- *Exact masking/corruptions strategy does not matter that much*
- *Mostly negative results for better finetuning and multi-task strategies*

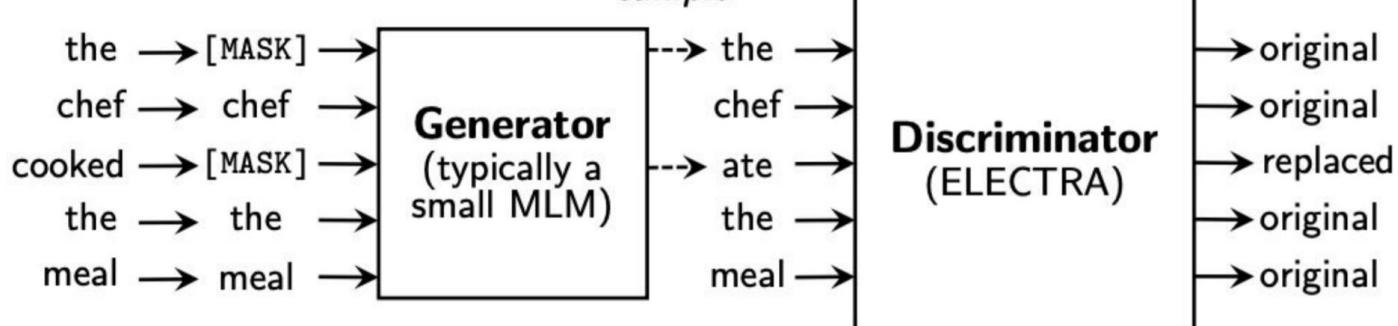
ELECTRA

Pre-training Text Encoders as Discriminators Rather Than Generators
(Clark et al, 2020)

Assignment Project Exam Help

Train model to discriminate locally plausible text from real text:
<https://tutorcs.com>

WeChat: cstutorcs



ELECTRA

*Pre-training Text Encoders as Discriminators Rather Than Generators
(Clark et al, 2020)*

Assignment Project Exam Help

Difficult to match SOTA results with less compute

<https://tutorcs.com>

Model	Train FLOPs	Params	SQuAD 1.1		SQuAD 2.0	
			EM	F1	EM	F1
BERT-Base	6.4e19 (0.09x)	110M	80.8	88.5	–	–
BERT	1.9e20 (0.27x)	335M	84.1	90.9	79.0	81.8
SpanBERT	7.1e20 (1x)	335M	88.8	94.6	85.7	88.7
XLNet-Base	6.6e19 (0.09x)	117M	81.3	–	78.5	–
XLNet	3.9e21 (5.4x)	360M	89.7	95.1	87.9	90.6
RoBERTa-100K	6.4e20 (0.90x)	356M	–	94.0	–	87.7
RoBERTa-500K	3.2e21 (4.5x)	356M	88.9	94.6	86.5	89.4
ALBERT	3.1e22 (44x)	235M	89.3	94.8	87.4	90.2
BERT (ours)	7.1e20 (1x)	335M	88.0	93.7	84.7	87.5
ELECTRA-Base	6.4e19 (0.09x)	110M	84.5	90.8	80.5	83.3
ELECTRA-400K	7.1e20 (1x)	335M	88.7	94.2	86.9	89.6
ELECTRA-1.75M	3.1e21 (4.4x)	335M	89.7	94.9	88.1	90.6

WeChat: cstutorcs

LongFormer

The Long-Document Transformer (Peters et al., 2020)

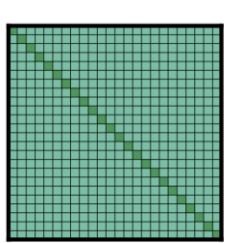
Why? Assignment Project Exam Help

- Traditional Transformer-based models are unable to process long sequences due to their self-attention operation, which scales quadratically with the sequence length
- To address this, Longformer uses an attention pattern that scales linearly with sequence length, making it easy to process documents of thousands of tokens or longer.

LongFormer

The Long-Document Transformer (Peters et al., 2020)

Assignment Project Exam Help

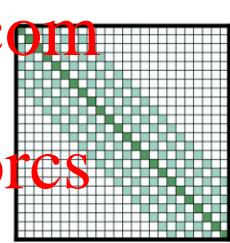


(a) Full n^2 attention

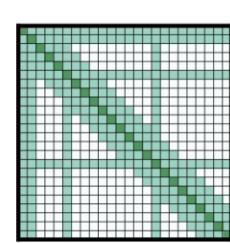


<https://tutorcs.com>
WeChat: cstutorcs

(b) Sliding window attention



(c) Dilated sliding window



(d) Global+sliding window

Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.

Applying Models to Production Services

- *BERT and other pre-trained language models are extremely large and expensive*
- *How are companies applying them to low-latency production services?*

<https://tutorcs.com>

GOOGLE \ TECH \ ARTIFICIAL INTELLIGENCE

Google is improving 10 percent of searches by
WeChat: cstutors
understanding language context

Say hello to BERT

By Dieter Bohn | @backlon | Oct 25, 2019, 3:01am EDT

Bing says it has been applying BERT since April

The natural language processing capabilities are now applied to all Bing queries globally.

[George Nguyen](#) on November 19, 2019 at 1:38 pm

Applying Models to Production Services

- *BERT and other pre-trained language models are extremely large and expensive*
- *How are companies applying them to low-latency production services?*

<https://tutorcs.com>

GOOGLE \ TECH \ ARTIFICIAL INTELLIGENCE

Google is improving 10 percent of searches by
WeChat: cstutorcs

The Answer is ‘Distillation’ (Model Compression)

The natural language processing capabilities are now applied to all Bing queries globally.

[George Nguyen](#) on November 19, 2019 at 1:38 pm

Distillation (Model Compression)

The idea has been around for a long time (from 2006)

- *Model Compression (Bucila et al. 2006)*
- *Distilling the Knowledge in a Neural Network (Hinton et al., 2015)*

<https://tutorcs.com>

WeChat: cstutorcs

Assignment Project Exam Help

4

Multimodal Pretrained Model

<https://tutorcs.com>

WeChat: cstutorcs

Multimodal Pretrained Model

Video Representation Learning

Supervised Learning: Large labelled data with CNN

- Expensive to collect labelled data
- Small corresponding label vocab not able to represent the nuances of actions (e.g. difference “sipping” - “drinking” - “gulping”)
- Represent short video clips (a few seconds long)

WeChat: cstutorcs *Unsupervised Learning: Learning density models from video*

- Single static stochastic variable, decoded into a sequence learning using RNN (VAE-style loss or GAN-style loss)
- Temporal stochastic variable (SV2P/SVCLP) or GAN-based (SAVP/MoCoGAN)
- What if not using explicit stochastic latent variable

Multimodal Pretrained Model



VideoBERT (ICCV 2019)

A Joint Model for Video and Language Representation Learning

Assignment Project Exam Help

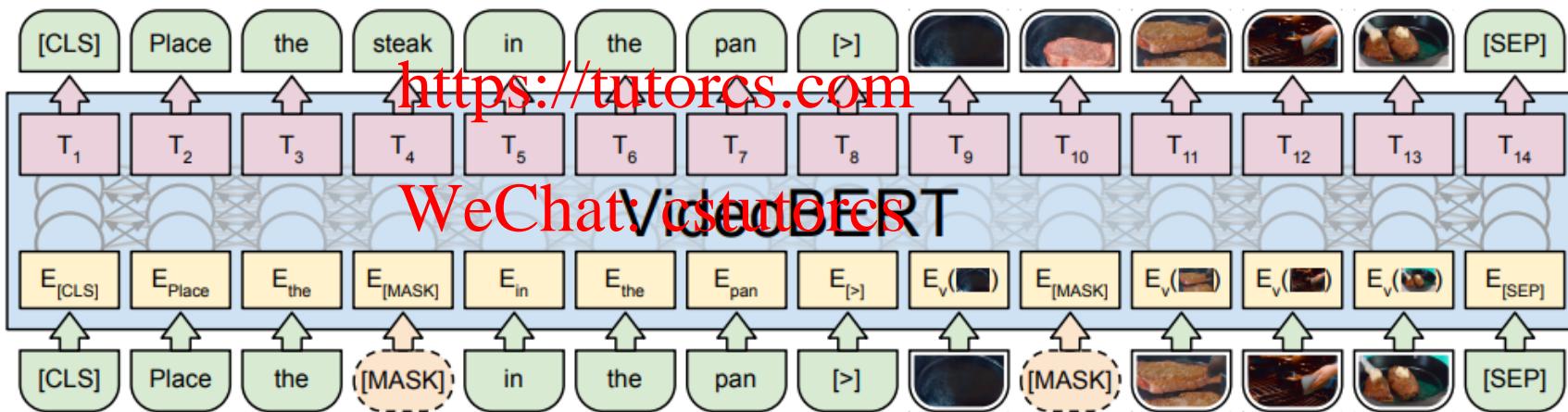


Figure 3: Illustration of VideoBERT in the context of a video and text masked token prediction, or *cloze*, task. This task also allows for training with text-only and video-only data, and VideoBERT can furthermore be trained using a linguistic-visual alignment classification objective (not shown here, see text for details).

Multimodal Pretrained Model



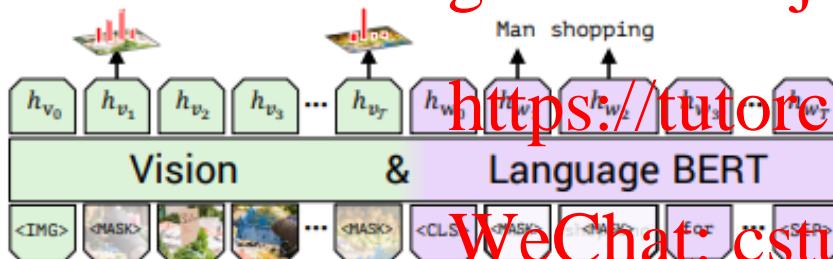
ViLBERT (NIPS 2019)

A Joint Model for Video and Language Representation Learning

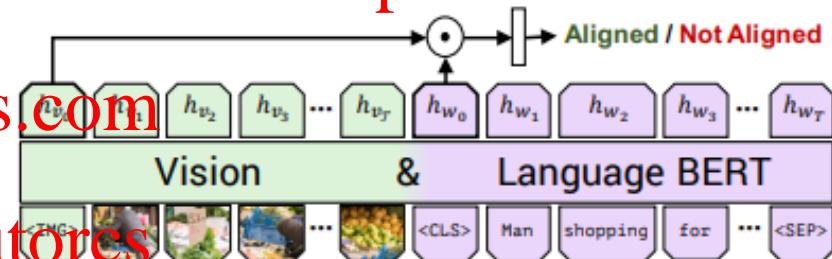
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



(a) Masked multi-modal learning



(b) Multi-modal alignment prediction

Figure 3: We train ViLBERT on the Conceptual Captions [24] dataset under two training tasks to learn visual grounding. In masked multi-modal learning, the model must reconstruct image region categories or words for masked inputs given the observed inputs. In multi-modal alignment prediction, the model must predict whether or not the caption describes the image content.

Multimodal Pretrained Model



ViLBERT (NIPS 2019)

Results across all transfer tasks

- *improves performance over a single-stream model*
- *result in improved vision/linguistic representations*
- *Finetuning from ViLBERT is a powerful strategy for vision and language tasks*

<https://tutorcs.com>

Table 1: Transfer task results for our ViLBERT model compared with existing state-of-the-art and sensible architectural ablations. [†] indicates models without pretraining on Conceptual Captions. For VCR and VQA which have private test sets, we report test results (in parentheses) only for our full model. Our full ViLBERT model outperforms task-specific state-of-the-art models across all tasks.

Method	VQA [3]			VCR [25]			RefCOCO+ [32]			Image Retrieval [26]			ZS Image Retrieval		
	test-dev (test-std)	Q→A	QA→R	Q→AR	val	testA	testB	R1	R5	R10	R1	R5	R10		
SOTA	DFAF [36]	70.22 (70.34)	-	-	-	-	-	-	-	-	-	-	-	-	-
	R2C [25]	-	63.8 (65.1)	67.2 (67.3)	43.1 (44.0)	-	-	-	-	-	-	-	-	-	-
	MAttNet [33]	-	-	-	-	65.33	71.62	56.02	-	-	-	-	-	-	-
	SCAN [35]	-	-	-	-	-	-	-	48.60	77.70	85.20	-	-	-	-
Ours	Single-Stream [†]	65.90	68.15	68.89	47.27	65.64	72.02	56.04	-	-	-	-	-	-	-
	Single-Stream	68.85	71.09	73.93	52.73	69.21	75.32	61.02	-	-	-	-	-	-	-
	ViLBERT [†]	68.93	69.26	71.01	49.48	68.61	75.97	58.44	45.50	76.78	85.02	0.00	0.00	0.00	0.00
	ViLBERT	70.55 (70.92)	72.42 (73.3)	74.47 (74.6)	54.04 (54.8)	72.34	78.52	62.61	58.20	84.90	91.52	31.86	61.12	72.80	

Assignment Project Exam Help

5

Before we finish this lecture...

<https://tutorcs.com>

WeChat: cstutorcs

0

Before we finish the lecture...

The current ML/DL-based NLP trends

- 1) Assume that we collect data from Royal Prince Alfred Hospital

Assignment Project Exam Help



ME

<https://tutorcs.com>

WeChat: cstutorcs



0

Before we finish the lecture...

The current ML/DL-based NLP trends

1) Assume that we collect data from Royal Prince Alfred Hospital

Assignment Project Exam Help



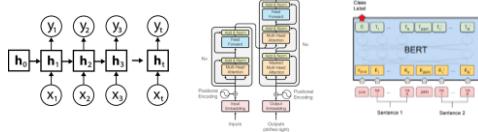
ME

<https://tutorcs.com>

WeChat: cstutorcs



***Training with RNN, Transformer,
or even pretrained model...***



2) Train and Test on data from the same hospital

0

Before we finish the lecture...

The current ML/DL-based NLP trends

1) Assume that we collect data from Royal Prince Alfred Hospital

Assignment Project Exam Help

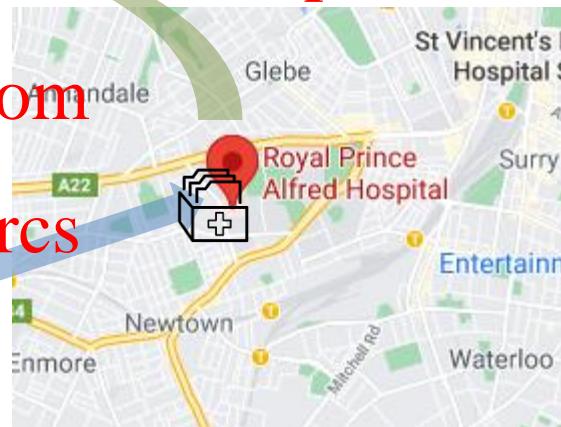


ME

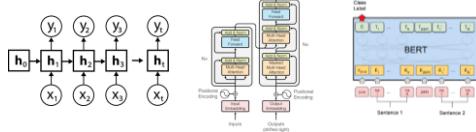


<https://tutorcs.com>

WeChat: cstutorcs



Training with RNN, Transformer,
or even pretrained model...



2) Train and Test on data from the same hospital

0 Before we finish the lecture...

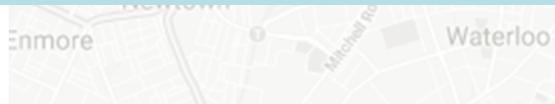
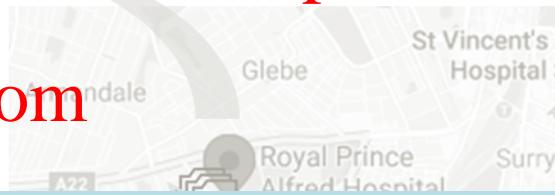
The current ML/DL-based NLP trends

1) Assume that we collect data from Royal Prince Alfred Hospital

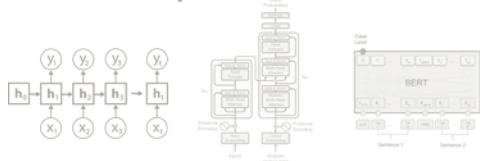
Assignment Project Exam Help

<https://tutorcs.com>

*"Indeed, WeChat postulates
comparable to human medical experts in spotting certain conditions"*



*Training with RNN, Transformer,
or even pretrained model...*



2) Train and Test on data from the same hospital

0 Before we finish the lecture...

The current ML/DL-based NLP trends

1) Assume that we collect data from Royal Prince Alfred Hospital

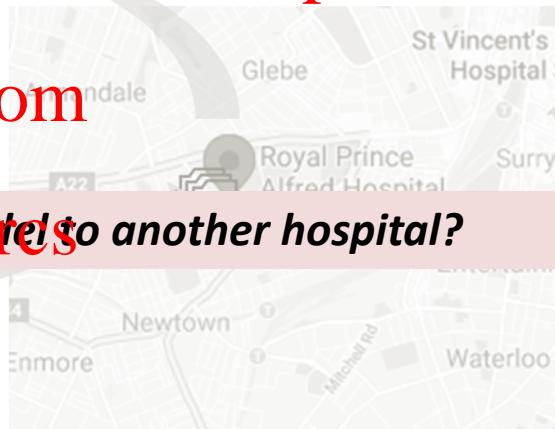
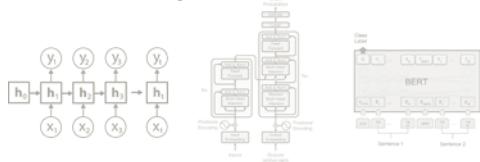
Assignment Project Exam Help

<https://tutorcs.com>

What if we use this same model to another hospital?



Training with RNN, Transformer,
or even pretrained model...



2) Train and Test on data from the same hospital

0

Before we finish the lecture...

The current ML/DL-based NLP trends

Assume you take that same DL-based NLP model, to St Vincent's Private Hospital, with an older testing machine, and the technician there uses a slightly different testing protocol

Assignment Project Exam Help



ME



<https://tutorcs.com>

WeChat: cstutorcs



0

Before we finish the lecture...

The current ML/DL-based NLP trends

Assume you take that same DL-based NLP model, to St Vincent's Private Hospital, with an older testing machine, and the technician there uses a slightly different testing protocol

Assignment Project Exam Help



ME



<https://tutorcs.com>

WeChat: cstutorcs



Data drifts to cause the performance of DL-based NLP model to degrade significantly

0 Before we finish the lecture...

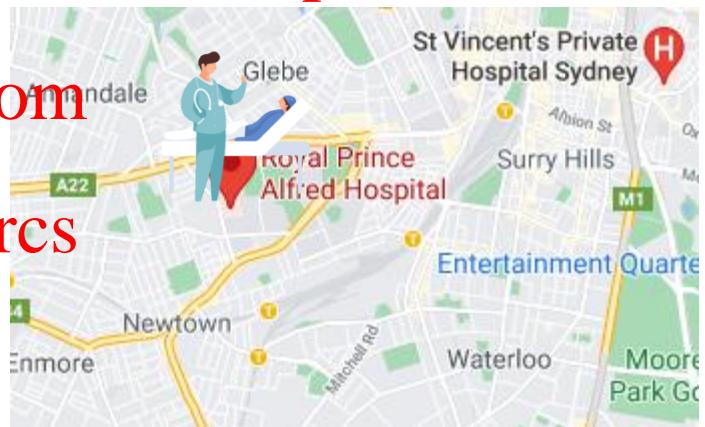
The current ML/DL-based NLP trends

In contrast, the Doctor just can walk down the street and diagnose the patient

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



0 Before we finish the lecture...

The current ML/DL-based NLP trends

*When a system is not performing well,
many teams instinctively try to improve the code
(try different model/component or change hyperparameters)*

Assignment Project Exam Help
<https://tutorcs.com>

*However, for many practical applications,
it is more effective instead to focus on improving the data*



0

Before we finish the lecture...

The current ML/DL-based NLP trends

*When a system is not performing well,
many teams instinctively try to improve the code
(try different model/component or change hyperparameters)*

Assignment Project Exam Help
<https://tutorcs.com>

*However, for many practical applications,
it is more effective instead to focus on improving the data*

*Everyone jokes about ML/DL is
80% data preparation, but no
one seems to care*

Prof. Andrew Ng (March, 2021)



0

Before we finish the lecture...

The current ML/DL-based NLP trends

Google facebook
Assignment Project Exam Help

 Microsoft  OpenAI

WeChat: cstutorcs

*There is unprecedented competition around beating the benchmarks.
If Google has BERT then OpenAI has GPT-3.*

*However, these fancy models take up only 20% of business problems.
What differentiates a good deployment is the quality of data.*

“Data Dispersion”

/ Reference

Reference for this lecture

- Deng, L., & Liu, Y. (Eds.). (2018). Deep Learning in Natural Language Processing. Springer.
- Rao, D., & McMahan, B. (2019). Natural Language Processing with PyTorch: Build Intelligent Language Applications Using Deep Learning. " O'Reilly Media, Inc.".
- Manning, C. D., Manning, C. D., & Schütze, H. (1999). Foundations of statistical natural language processing. MIT press.
- Manning, C 2018, Natural Language Processing with Deep Learning, lecture notes, Stanford University

<https://tutorcs.com>

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).
 - Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018) Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
 - Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep contextualized word representations. arXiv preprint arXiv:1802.05365.
 - Miller, A., Fisch, A., Dodge, J., Karimi, A. H., Bordes, A., & Weston, J. (2016). Key-value memory networks for directly reading documents. arXiv preprint arXiv:1606.03126.
- WeChat: estidores
- Drawings
 - <http://jalammar.github.io/illustrated-bert/>
 - <http://jalammar.github.io/illustrated-transformer/>