

# Natural Language Understanding, Generation, and Machine Translation

## Lecture 11: Contextualized Word Embeddings

---

Frank Keller

Week of 7 February 2022 (week 4)

School of Informatics

University of Edinburgh

[keller@inf.ed.ac.uk](mailto:keller@inf.ed.ac.uk)

The Story so Far

Bert Architecture

Masked Training

Pre-training and Finetuning Bert

Reading: Devlin et al. (2019).

## The Story so Far

---

## Self-attention

Self-attention is what we get when compute attention over the input sequence. Let  $x_1, \dots, x_t$  be the input vectors and  $y_1, \dots, y_t$  be the output vectors:

$$y_i = \sum_j w_{ij} x_j$$

We now compute the attention weight as the dot-product of each input token with every other token.

$$w'_{ij} = x_i \cdot x_j$$
$$w_i = \text{softmax}(w'_i)$$

Note that the attention weight is now called  $w'$  and the attention distribution  $w$  (rather than  $a$  and  $\alpha$ ).

# Self-attention

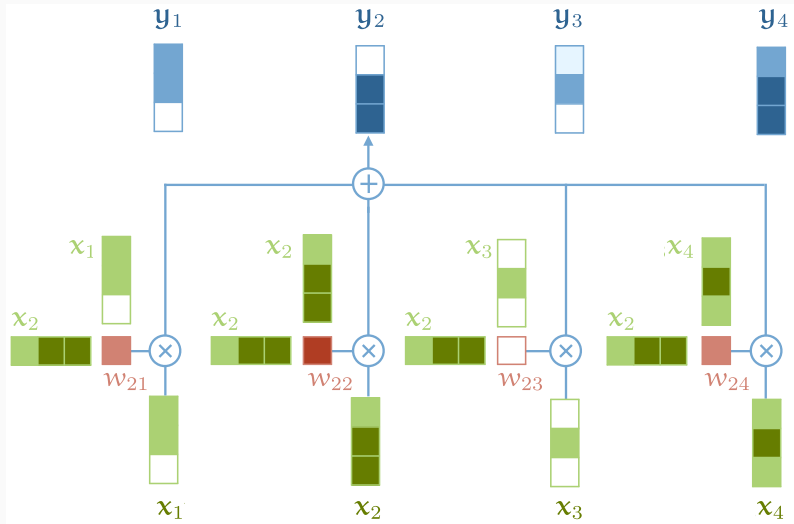
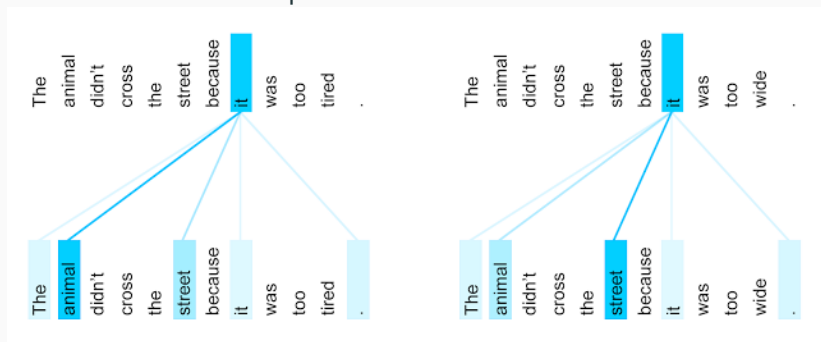


Figure from Bloem (2019).

# Multi-head Attention

*Multi-head attention* is able to jointly attend to different parts of the input. If we just have a single head, have to average.

For example, one head could attend to the subject of a verb, another one to its object. Or different heads could attend to different referents of a pronoun.



# The Transformer

The multi-head self attention mechanism needs to be integrated into a larger architecture to be useful:

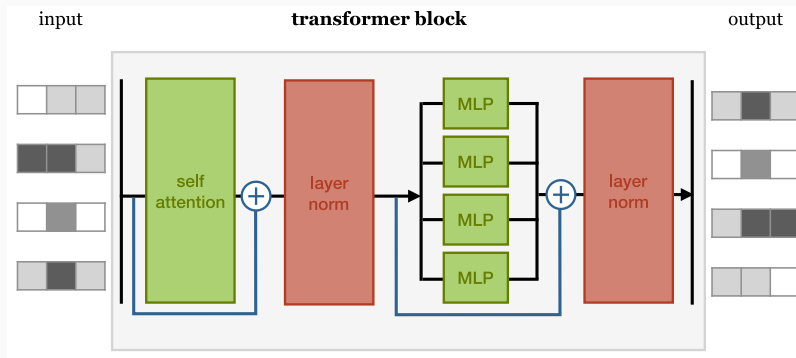


Figure from Bloem (2019).

Is the MLP same?



# From Transformers to Embeddings

- In lecture 9, we saw how we can stack transformers and use them for *classification*. We apply mean pooling to the output and map it onto a class vector.
- We can also use transformers for *sequence prediction*. For this we need to mask some of the input.
- In lecture 10, we saw how we can derive *embeddings* from neural language models.
- Transformers are language model that can represent context very well: they can learn *contextualized embeddings*.
- Contextualized language models can be used for feature extraction, but also in a *pre-training/finetuning* setup.



# Bert Architecture

---

# Bert Architecture

Bert (Bidirectional Encoder Representations from Transformers):

- designed for pre-training deep bidirectional representations from unlabeled text;
- conditions on left and right context in all layers;
- pre-trained model can be finetuned with one additional output layer for many tasks (e.g., NLI, QA, sentiment);
- for many tasks, no modifications to the Bert architecture are required;
- Devlin et al. (2019) report SotA results on 11 tasks using pre-training/finetuning approach.

# Bert Architecture

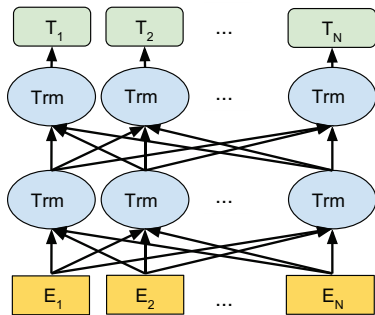
Bert uses bidirectional transformer representations. This is shown to be superior to competing architectures:

- GPT (Radford et al. 2018) also uses transformers, but is only trained left-to-right;
- Elmo (Peters et al. 2018) uses shallow concatenation of independently trained left-to-right and right-to-left LSTMs.

Main attraction of Bert: *pre-training/finetuning* approach (though this is also present in GPT).

# Bert Architecture

## BERT (Ours)



## OpenAI GPT

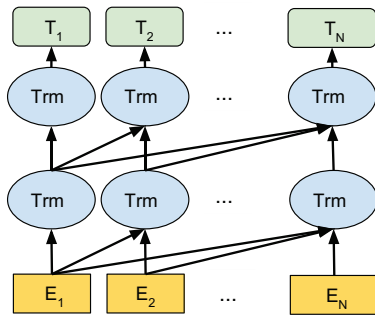


Figure from Devlin et al. (2019).

good at text generation

# Bert Architecture

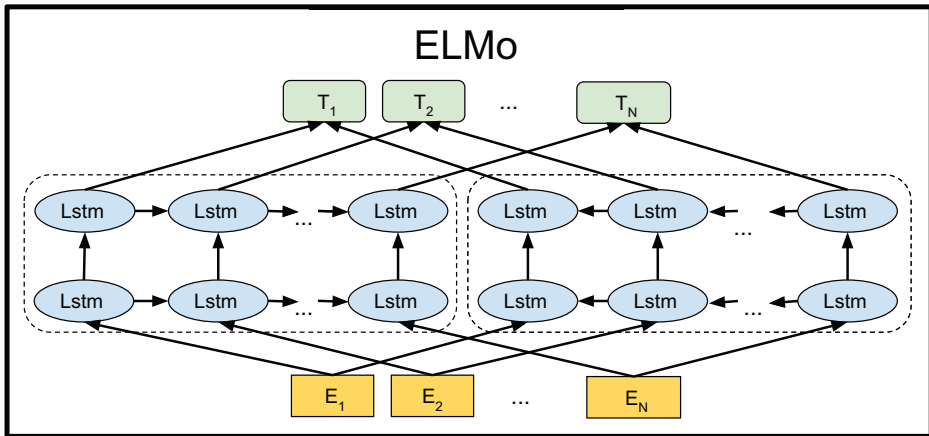


Figure from Devlin et al. (2019).

# Bert Architecture

Basic Bert architecture:

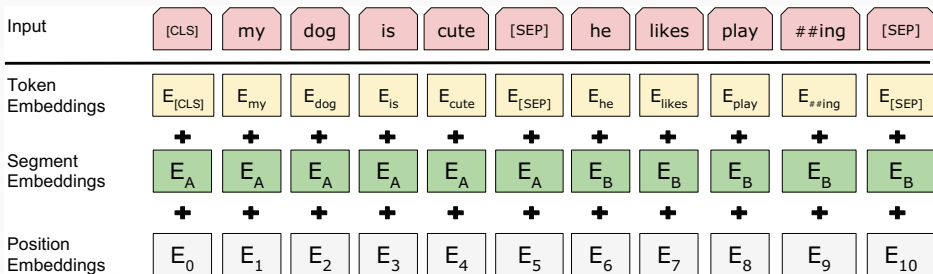
- multi-layer bidirectional transformer (Vaswani et al. 2017);
- $L$ : layers (transformer blocks);  $H$ : dimensionality of hidden layer,  $A$ : number of self-attention heads;
- Bert Base:  $L = 12$ ,  $H = 768$ ,  $A = 12$ , 110M parameters;
- Bert Large:  $L = 24$ ,  $H = 1024$ ,  $A = 16$ , 340M parameters.

# Input/Output Representation

- Input sequence: can be  $\langle \text{Question}, \text{Answer} \rangle$  pair, single sentence, or any other string of tokens;
- 30,000 token vocabulary, represented as WordPiece embeddings (handles OOV words);
- first token is always [CLS]: aggregate sentence representation for classification tasks;
- sentence pairs separated by [SEP] token; and by segment embeddings;
- token position represented by position embeddings (see lecture 10).

input is too short: use padding!

# Bert Architecture



$E$ : input embedding,  $C$ : final hidden vector of [CLS];  $T_i$ : final hidden vector for the  $i$ -th input token.

Figure from Devlin et al. (2019).



# Masked Training

---

# Masked Language Model

Important departure from previous embedding models (including GPT and Elmo):

*Don't train the model to predict the next word, but train it to predict the whole context.*

- Problem: how can we prevent *trivial copying* via the self-attention mechanism?
- Solution: mask 15% of the tokens in the input sequence; train the model to predict these.

# Masked Language Model

Problem: masking creates mismatch between pre-training and finetuning: [MASK] token is not seen during fine-tuning. Solution:

- do not always replace masked words with [MASK], instead choose 15% of token positions at random for prediction;
- if  $i$ -th token is chosen, we replace the  $i$ -th token with:
  1. the [MASK] token 80% of the time;
  2. a random token 10% of the time;
  3. the unchanged  $i$ -th token 10% of the time.
- Now use  $T_i$  to predict original token with cross entropy loss.

# Masked Language Model

Why this masking scheme?

- if we always use [MASK] token, the model would not have to learn good representation for other words;
- if we only use [MASK] token or random word, model would learn that observed word is never correct;
- if we only use [MASK] token or observed word, model would just learn to trivially copy.

## Secondary Task: Next Sentence Prediction

- Bert is designed to be used for tasks such as question answering (QA) and natural language inference (NLI) that require sentence pairs;
- Bert uses a special mechanism to capture this: pre-training on next sentence prediction task;
- generate training data: chose two sentences  $A$  and  $B$ , such that 50% of the time  $B$  is the actual next sentence of  $A$ , and 50% of the time a randomly selected sentence.

## Pre-training and Finetuning Bert

---

bert can do both word embedding (which contains the contextual information) and sentence embedding

# Pre-training and Finetuning Bert

pre-training do not need label, fine-tune needs label

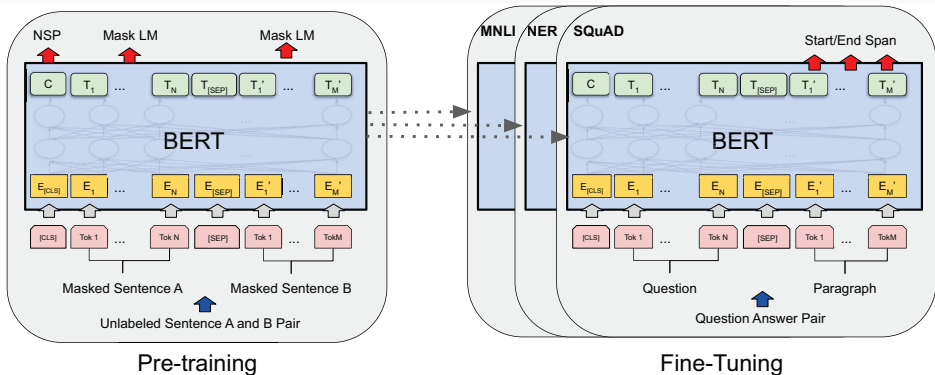


Figure from Devlin et al. (2019).

input and output are exactly same, maybe just position shift

# Pre-training and Finetuning Bert

- Pre-training on BooksCorpus (800M words) and English Wikipedia (2,500M words);
- to finetune for a new task, we plug task-specific inputs and outputs into Bert and re-train its parameters;
- input: designed to be two sentences:
  1. sentence pairs in paraphrasing;
  2. hypothesis-premise pairs in entailment;
  3. question-passage pairs in question answering;
  4. text- $\emptyset$  pair in text classification or sequence tagging;
- output:
  1. sequence of tokens in tasks such as QA (mark answer span);
  2. sequence of labels in tagging tasks such as NER;
  3. [CLS] representation is fed into an output layer for classification, such as entailment or sentiment.
- appendix of Devlin et al. (2019) gives examples.



## Bert Results: Glue Benchmark

System	MNLI-(m/mm) 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+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

Glue is a mixture of various natural language understanding tasks: MNLI, QNLI, WNLI: natural language inference; QQP: question equivalence; SST-2: sentiment; CoLA: linguistic acceptability; STS-B: semantic similarity; MRPC: paraphrasing; RTE: entailment.

how to evaluate the sentence paraphrasing?

language model can tell whether this is the paraphrasing or not, can not evaluate how good it is

bert just good at classification and tagging

## Bert Results: Squad 2.0 Question Answering

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	86.3	89.0	86.9	89.5
#1 Single - MIR-MRC (F-Net)	-	-	74.8	78.0
#2 Single - nlnet	-	-	74.2	77.1
Published				
unet (Ensemble)	-	-	71.4	74.9
SLQA+ (Single)	-		71.4	74.4
Ours				
BERT <sub>LARGE</sub> (Single)	78.7	81.9	80.0	83.1

## Feature Extraction vs. Finetuning: Named Entity Recognition

System	Dev F1	Test F1
ELMo (Peters et al., 2018a)	95.7	92.2
CVT (Clark et al., 2018)	-	92.6
CSE (Akbik et al., 2018)	-	<b>93.1</b>
Fine-tuning approach		
BERT <sub>LARGE</sub>	96.6	92.8
BERT <sub>BASE</sub>	96.4	92.4
Feature-based approach (BERT <sub>BASE</sub> )		
Embeddings	91.0	-
Second-to-Last Hidden	95.6	-
Last Hidden	94.9	-
Weighted Sum Last Four Hidden	95.9	-
Concat Last Four Hidden	96.1	-
Weighted Sum All 12 Layers	95.5	-

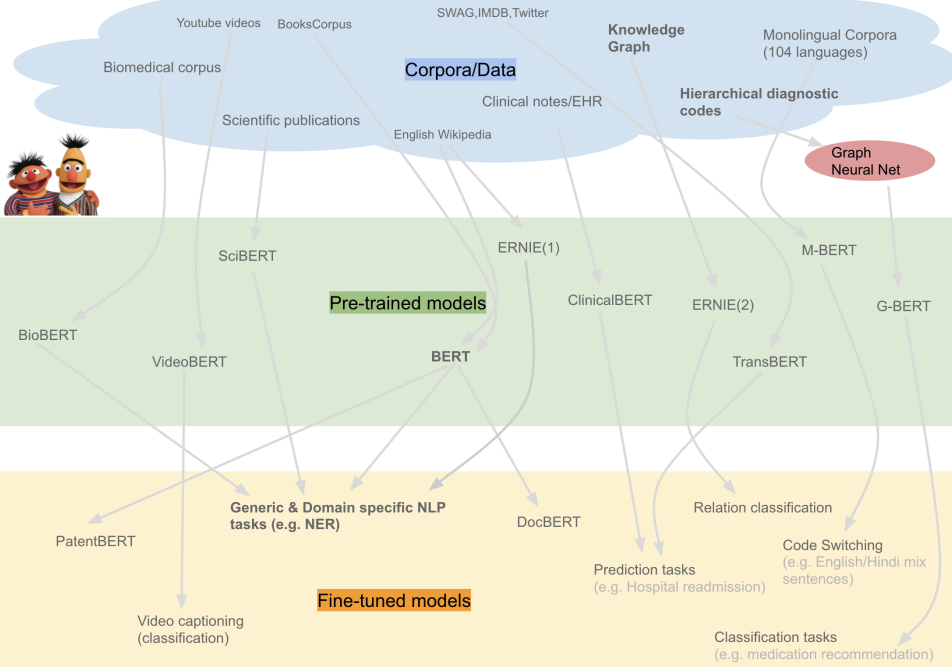
# Bert and the State of the Art

The current state of the art in NLP owes a lot to pre-training/finetuning of Bert (and similar models):

- pre-training a large-scale model like Bert, then fine-tuning on supervised data, is state of the art for many NLP tasks
- pre-training Bert requires lots of resources! Estimate (Tim Dettmers): 4 GPUs (RTX 2080 Ti) for 100 days
- fine-tuning existing model is relatively quick, but fitting model in GPU memory can be a challenge
- getting new state of the art results with ever-larger models and data is a game that will continue, but only few can play
- positive (for everybody else): we also need smart ideas for architectures, objective functions, evaluation, etc.

# Summary

- Bert is a contextualized language model that uses a deep, bidirectional transformer architecture;
- it is pre-trained on unlabeled text using masking and next sentence prediction;
- it is designed for finetuning with minimal architectural modifications;
- the input uses sentence pairs; the output can be sentences, labels, classification decisions, depending on task;
- Bert-based models are state of the art on many NLP tasks.
- There are many variants of Bert . . .



# References

- Bloem, Peter. 2019. Transformers from scratch. Blog, <http://www.peterbloem.nl/blog/transformers>.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*. Minneapolis, MN, pages 4171–4186.
- Peters, Matthew E., Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*. New Orleans, LA, pages 2227–2237.
- Radford, Alec, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning. Technical report, OpenAI.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*. Curran Associates, Red Hook, NY, pages 5998–6008.