

Natural Language Processing

CSC413/2516



- Central to human intelligence.
- Tremendous practical value.
- Colossal developments recently.

Goal of This Tutorial

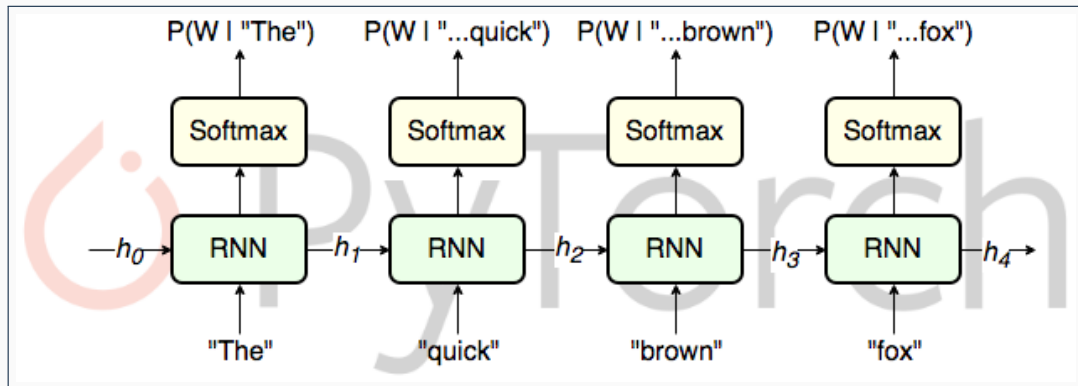
- Basic NLP literacy.
- Getting up to date with recent developments.
 - ▶ Architectures
 - ▶ Language tasks
 - ▶ Tremendous developments in the field recently!
- Know where to look at if you're starting an NLP project

- A statistical model that assigns probabilities to the words in a sentences.
- **Most commonly:** Given previous words, what should the next one be?
- **Neural language model:** Model the probability of words given others using neural networks.

Which architecture is most suitable?

Recurrent Architectures

- We can use recurrent architectures.
- LSTM, GRU ...
- Great for variable length inputs, like sentences.



What are some of the problems with recurrent architectures?

What are some of the problems with recurrent architectures?

- Not parallelizable across instances.
- Cannot model long dependences.
- Optimization difficulties (vanishing gradients).

We'd like an architectural primitive that is:

- Ideally feed-forward
- Can facilitate between-token interactions
- Can model long dependences easily.

Attention to the rescue!

- There are many forms of attention. Today we'll focus on **scaled dot product attention**.

- Three inputs: queries, keys and values.
- "Return a combination of the values based on the similarities between keys and queries".
- Dimensionalities:
 - ▶ $Q \in \mathbb{R}^{n_q \times d_{kq}}$
 - ▶ $K \in \mathbb{R}^{n_{kv} \times d_{kq}}$
 - ▶ $V \in \mathbb{R}^{n_{kv} \times d_v}$

- $$\mathbb{A}(Q, K, V) = V(\text{softmax}(\frac{K^T Q}{d_{kq}})) \quad (1)$$

Quiz!

- True or False: The dimensionality of queries have to match the dimensionality of the keys.

Quiz!

- True or False: The dimensionality of queries have to match the dimensionality of the keys. TRUE
- True or False: The number of keys have to match the number of values.

Quiz!

- True or False: The dimensionality of queries have to match the dimensionality of the keys. TRUE
- True or False: The number of keys have to match the number of values. TRUE
- True or False: The number of keys have to match the number of queries.

Quiz!

- True or False: The dimensionality of queries have to match the dimensionality of the keys. TRUE
- True or False: The number of keys have to match the number of values. TRUE
- True or False: The number of keys have to match the number of queries. FALSE

Quiz!

- True or False: The dimensionality of queries have to match the dimensionality of the keys. TRUE
- True or False: The number of keys have to match the number of values. TRUE
- True or False: The number of keys have to match the number of queries. FALSE
- Dimensionalities:
 - ▶ $Q \in \mathbb{R}^{n_q \times d_{kq}}$
 - ▶ $K \in \mathbb{R}^{n_{kv} \times d_{kq}}$
 - ▶ $V \in \mathbb{R}^{n_{kv} \times d_v}$

We need to make two central decisions:

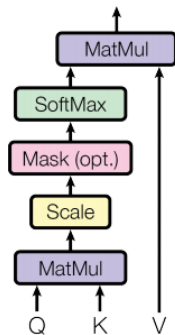
- How do we compute similarity?
- How do we 'normalize' the similarity scores amongst values?

- **Similarity:** Dot product between keys and queries.
- **Interesting theorem:** In high dimensions, two randomly sampled ¹ vectors are almost always approximately perpendicular to each other.
- **Normalization:** Softmax along the keys/values!
- **Result:** Scaled dot product attention.
- We get the following attention mechanism:

$$\mathbb{A}(Q, K, V) = V(\text{softmax}(\frac{K^T Q}{d_{kq}})) \quad (2)$$

¹From, lets say, a isotropic multivariate Gaussian distribution.

Scaled Dot-Product Attention



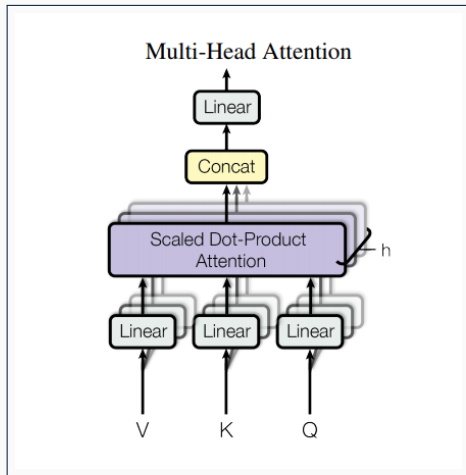
Covered in the lecture

<https://csc413-uoft.github.io/2021/assets/slides/lec08.pdf>.

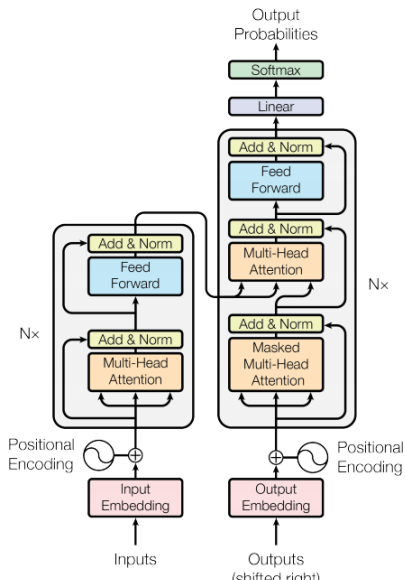
- What is self-attention?
- Use the same tensor to compute keys, values and queries!
- What is position embedding?
- Encode token's position in the input text into a vector.

- Lingering question: What is learned in an attention layer?
 - ▶ The space in which the similarities are computed.
 - ▶ The transformations on the values.
- What if we'd like to have different notions of similarity on the same set of tokens?
- **Multi head attention** to the rescue!

Multi head attention



Transformers



Properties of the transformer architecture:

- Fully feed forward.
- Equivariance properties of scaled dot product attention (important):
 - ▶ How does the output change if we permute the order of queries? (equivariance)
 - ▶ How does the output change if we permute the key-value pairs in unison? (invariance)

Exercise: How about self-attention?

Performance Comparison

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

What's next?

What's next?

- Brief intro to self supervised learning
- GPT (i.e. Self-supervised training of language model)
- BERT (i.e. Self-supervised training of language model)

Three types of learning:

- Supervised learning
- Reinforcement learning
- **Unsupervised/self-supervised learning:**
 - ▶ When the label is in the data itself!
 - ▶ Possible to make use of large amounts of data with no additional labelling efforts.

Examples of self-supervised learning:

- Predict next frame in a video.
- Image completion.
- Auto-encoding tasks.
- Rotation prediction.
- **Predicting next word from previous ones.**

- Can we use large amounts of text data to pretrain language models?

- Can we use large amounts of text data to pretrain language models?
- Considerations:
 - ▶ How can we fuse both left-right and right-left context?
 - ▶ How can we facilitate non-trivial interactions between input tokens?

- Can we use large amounts of text data to pretrain language models?
- Considerations:
 - ▶ How can we fuse both left-right and right-left context?
 - ▶ How can we facilitate non-trivial interactions between input tokens?
- Previous approaches:
 - ▶ ELMO (Peters. et. al., 2017): Bidirectional, but shallow.
 - ▶ GPT (Radford et. al., 2018), GPT-2 (Radford et. al., 2019): Deep, unidirectional.
 - ▶ BERT (Devlin et. al., 2018): Deep and bidirectional!

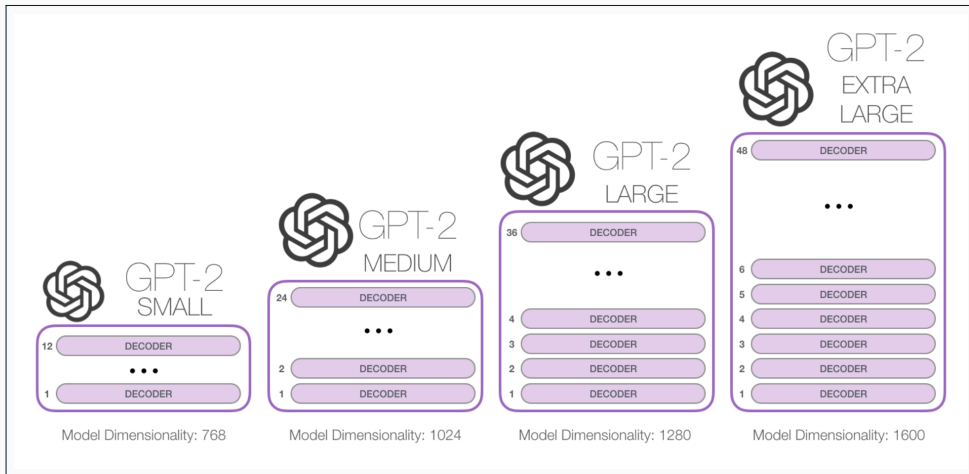
GPT and GPT-2 have very similar model architecture, we use GPT-2 as the example from now on.

The GPT-2 wasn't a particularly novel architecture – very similar to the decoder-only transformer, but a very large model trained on massive datasets.

- The workflow includes:
 - ▶ Pretrain on generic, self-supervised tasks, using large amounts of data (like all of Wikipedia)
 - ▶ Fine-tune on specific tasks with limited, labelled data.
- The pretraining tasks (will talk about this in more detail later):
 - ▶ Predict next token (i.e. word/character) - to learn contextualized token representations

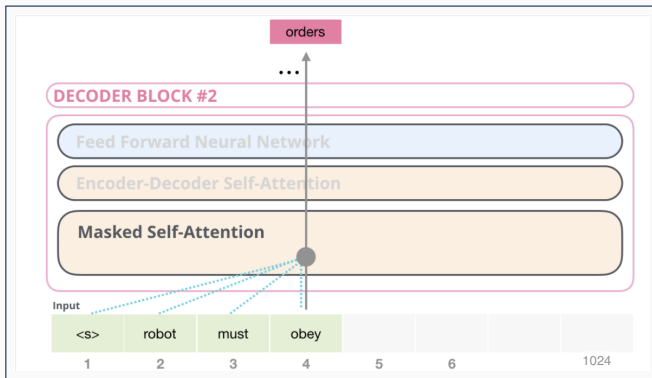
Variants of GPT-2 of different sizes

Credits of following figures to Jay Alammar. He has a series of high quality blogs on transformers <http://jalammar.github.io/illustrated-gpt2/>.



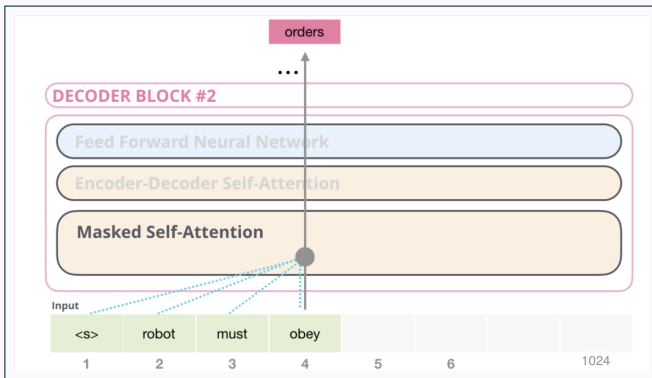
How GPT-2 works

Recall the transformer decoder. It generates the next token in an auto-regressive way.



How GPT-2 works

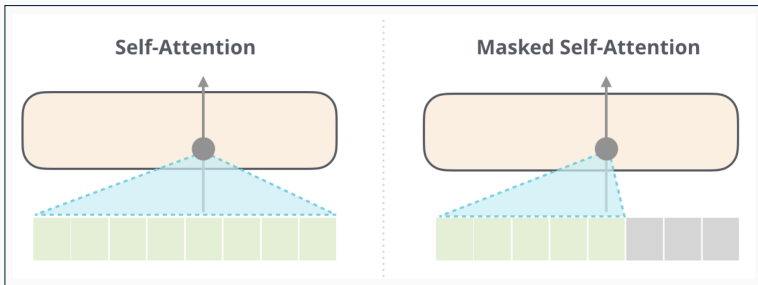
Recall the transformer decoder. It generates the next token in an auto-regressive way.



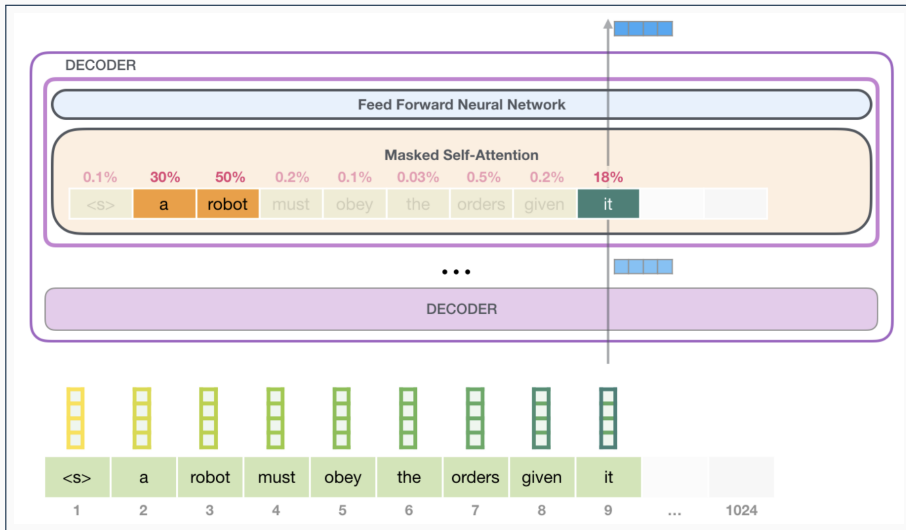
To make this work, future tokens need to be masked out using the **Masked Self-Attention**.

Masked Self-Attention

A normal self-attention block allows a position to peak at tokens to its right. Masked self-attention prevents that from happening. This is one distinction between GPT-2 and BERT.



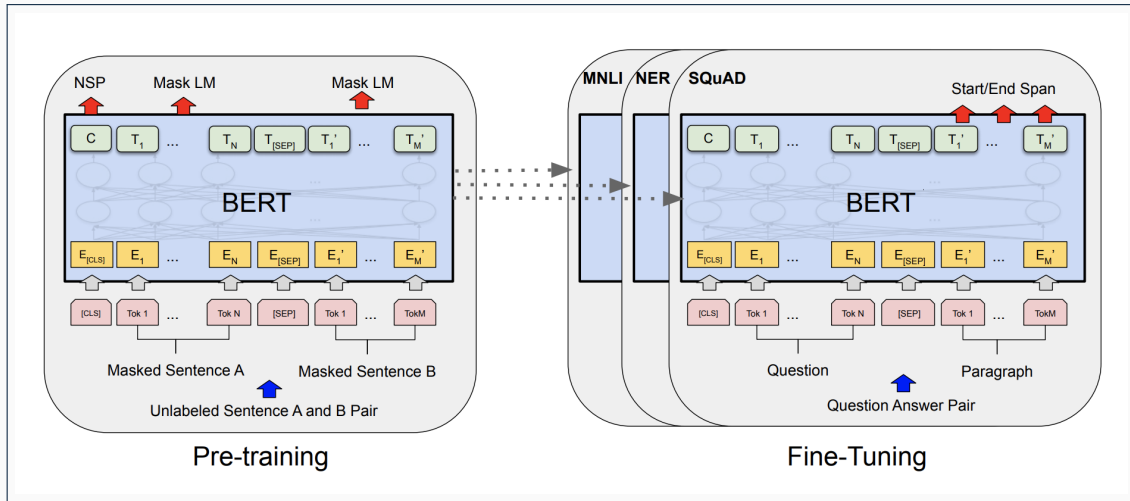
Under the hood



Next, BERT

- The BERT workflow includes:
 - ▶ Pretrain on generic, self-supervised tasks, using large amounts of data (like all of Wikipedia)
 - ▶ Fine-tune on specific tasks with limited, labelled data.
- The pretraining tasks (will talk about this in more detail later):
 - ▶ Masked Language Modelling (to learn contextualized token representations)
 - ▶ Next Sentence Prediction (summary vector for the whole input)

BERT Architecture



Properties

- Two input sequences.
 - ▶ Many NLP tasks have two inputs (question answering, paraphrase detection, entailment detection etc.)
- Computes embeddings
 - ▶ Both token, position and segment embeddings.
 - ▶ Special start and separation tokens.
- Architecture
 - ▶ Basically the same as transformer encoder.
- Outputs:
 - ▶ Contextualized token representations.
 - ▶ Special tokens for context.

BERT Embeddings

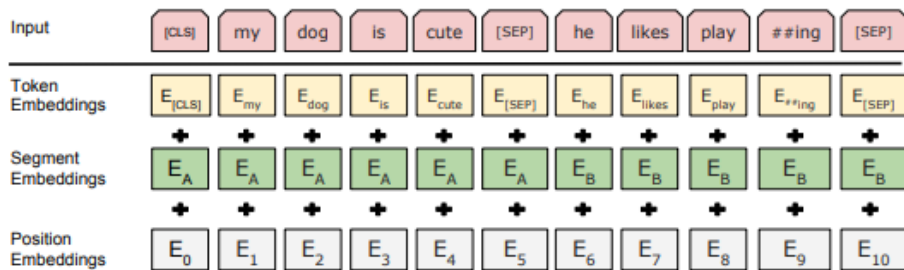


Figure 2: BERT input representation. The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

- How we tokenize the inputs is very important!
- BERT uses the WordPiece tokenizer (Wu et. al. 2016)

(Aside) Tokenizers

- Tokenizers have to balance the following:
 - ▶ Being comprehensive (rare words? translation to different languages)
 - ▶ Total number of tokens
 - ▶ How semantically meaningful each token is.
- This is an activate area of research.

- Masked Language Modelling (i.e. Cloze Task (Taylor, 1953))
- Next sentence prediction

- Mask 15% of the input tokens. (i.e. replace with a dummy masking token)
- Run the model, obtain the embeddings for the masked tokens.
- Using these embeddings, try to predict the missing token.
- "I love to eat peanut ___ and jam. " Can you guess what's missing?

This procedure forces the model to encode context information in the features of all of the tokens.

Next Sentence Prediction

- Goal is to summarize the complete context (i.e. the two segments) in a single feature vector.
- Procedure for generating data
 - ▶ Pick a sentence from the training corpus and feed it as "segment A".
 - ▶ With 50% probability, pick the following sentence and feed that as "segment B".
 - ▶ With 50% probability, pick the a random sentence and feed it as "segment B".
- Using the features for the context token, predict whether segment B is the following sentence of segment A.
- Turns out to be a very effective pretraining technique!

Procedure:

- Add a final layer on top of BERT representations.
- Train the whole network on the fine-tuning dataset.
- Pre-training time: In the order of days on TPUs.
- Fine tuning task: Takes only a few hours max.

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 _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). 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.⁸ BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

The BERT section of this tutorial is influenced by the fantastic talk by Lukasz Kaiser on transformers: <https://www.youtube.com/watch?v=rBCqOTefxvgt=1704s>

The GPT-2 section is influenced by the illustrated transformer blogs <http://jalammar.github.io/>.

Hugging Face and tutorial notebooks:

<https://huggingface.co/transformers/notebooks.html>