

CS388: Natural Language Processing

Lecture 8: Pre-trained Encoders

Greg Durrett



Announcements

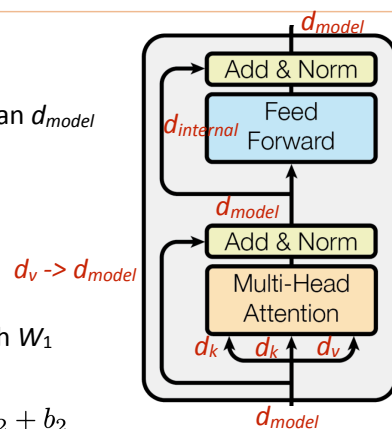
- P2 due Tuesday
- Final project released, proposals due Feb 23



Recall: Transformers

- Vectors: d_{model}
- Queries/keys: d_k , always smaller than d_{model}
- Values: separate dimension d_v , output is multiplied by W^O which is $d_v \times d_{model}$ so we can get back to d_{model} before the residual
- FFN can explode the dimension with W_1 and collapse it back with W_2

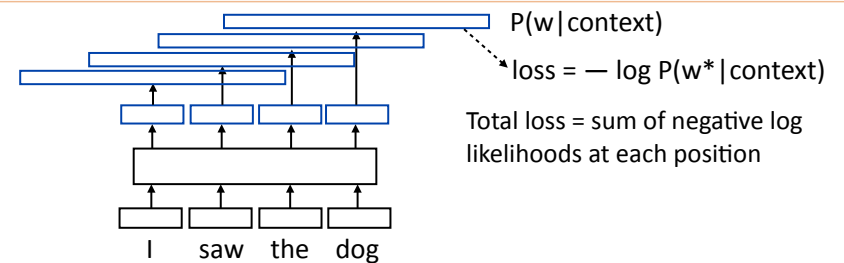
$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



Vaswani et al. (2017)



Recall: Training Transformer LMs



```
loss_fcn = nn.NLLLoss()
loss += loss_fcn(log_probs, ex.output_tensor)
```

[seq len, num output classes] [seq len]

- Batching is a little tricky with NLLLoss: need to collapse [batch, seq len, num classes] to [batch * seq len, num classes]. You do not need to batch



Today

- ELMo
- BERT
- Subword tokenization
- BERT results, BERT variants
- Applying BERT

ELMo



What is pre-training?

- “Pre-train” a model on a large dataset for task X, then “fine-tune” it on a dataset for task Y
- Key idea: X is somewhat related to Y, so a model that can do X will have some good neural representations for Y as well
- ImageNet pre-training is huge in computer vision: learn generic visual features for recognizing objects
- GloVe can be seen as pre-training: learn vectors with the skip-gram objective on large data (task X), then fine-tune them as part of a neural network for sentiment/any other task (task Y)



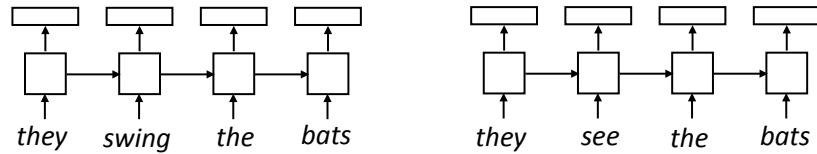
GloVe is insufficient

- GloVe uses a lot of data but in a weak way
- GloVe gives a single embedding for each word is wrong

they swing the bats they see the bats
- Identifying discrete word senses is hard, doesn't scale. Hard to identify how many senses each word has
- How can we make our word embeddings more *context-dependent*?
Use language model pretraining!



Context-dependent Embeddings



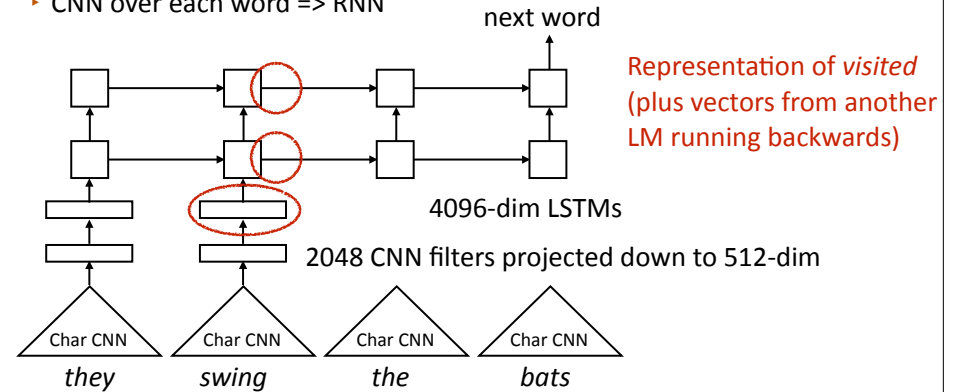
- ▶ Train a neural language model to predict the next word given previous words in the sentence, use the hidden states (output) at each step as *word embeddings*
- ▶ This is the key idea behind ELMo: language models can allow us to form useful word representations in the same way word2vec did

Peters et al. (2018)



ELMo

- ▶ CNN over each word => RNN



Peters et al. (2018)



ELMo

- ▶ Use the embeddings as a drop-in replacement for GloVe
- ▶ Huge gains across many high-profile tasks: NER, question answering, semantic role labeling (similar to parsing), etc.
- ▶ But what if the pre-training **isn't just for the embeddings?**

BERT



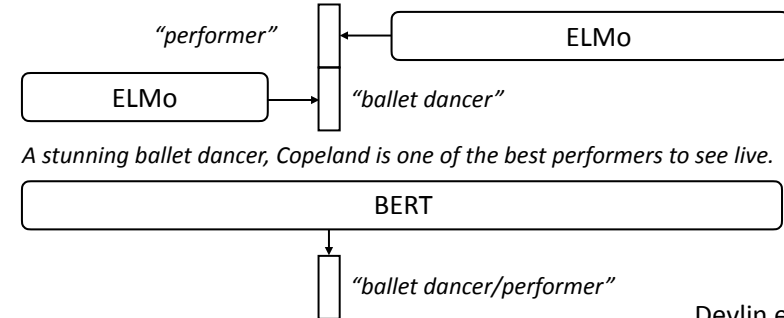
BERT

- AI2 made ELMo in spring 2018, GPT (transformer-based ELMo) was released in summer 2018, BERT came out October 2018
- Four major changes compared to ELMo:
 - Transformers instead of LSTMs
 - Bidirectional model with “Masked LM” objective instead of standard LM
 - Fine-tune instead of freeze at test time
 - Operates over word pieces (byte pair encoding)



BERT

- ELMo is a unidirectional model (as is GPT): we can concatenate two unidirectional models, but is this the right thing to do?
- ELMo reprs look at each direction in isolation; BERT looks at them jointly

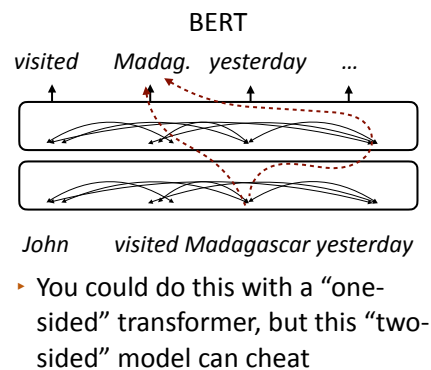
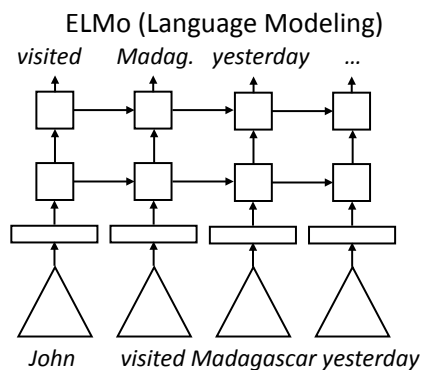


Devlin et al. (2019)



BERT

- How to learn a “deeply bidirectional” model? What happens if we just replace an LSTM with a transformer?



- You could do this with a “one-sided” transformer, but this “two-sided” model can cheat



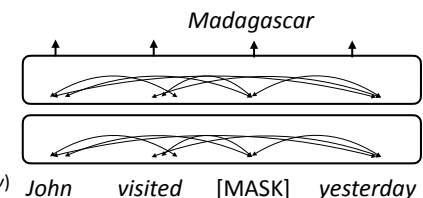
Masked Language Modeling

- How to prevent cheating? Next word prediction fundamentally doesn't work for bidirectional models, instead do *masked language modeling*

- BERT formula: take a chunk of text, mask out 15% of the tokens, and try to predict them

- Optimize

$P(\text{Madagascar} \mid \text{John visited } [\text{MASK}] \text{ yesterday})$

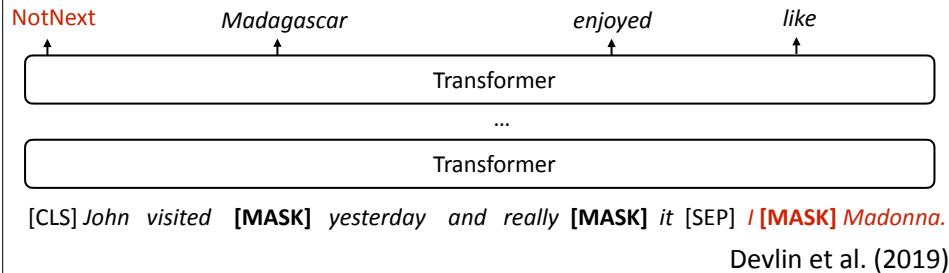


Devlin et al. (2019)



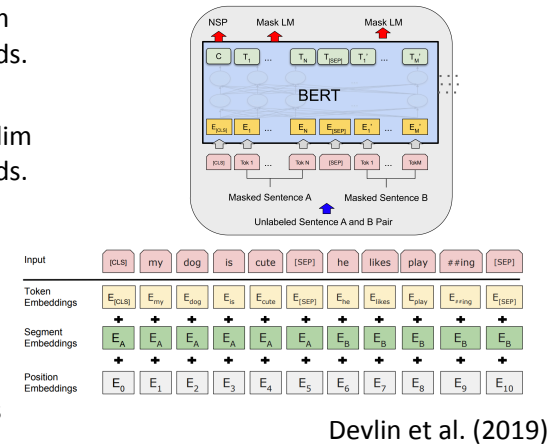
Next “Sentence” Prediction

- Input: [CLS] Text chunk 1 [SEP] Text chunk 2
- 50% of the time, take the true next chunk of text, 50% of the time take a random other chunk. Predict whether the next chunk is the “true” next
- BERT objective: masked LM + next sentence prediction

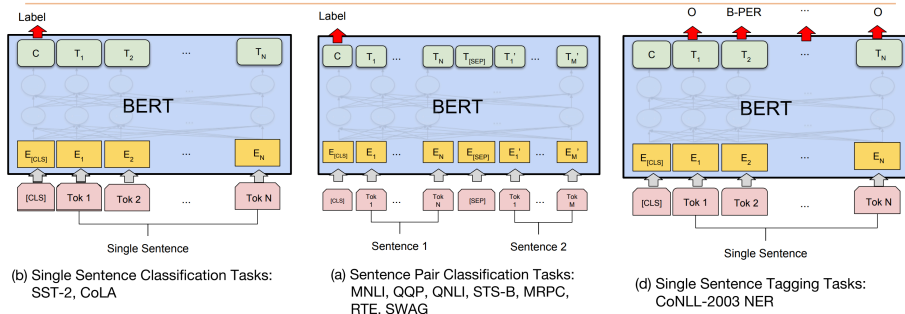


BERT Architecture

- BERT Base: 12 layers, 768-dim per wordpiece token, 12 heads. Total params = 110M
- BERT Large: 24 layers, 1024-dim per wordpiece token, 16 heads. Total params = 340M
- Positional embeddings and segment embeddings, 30k word pieces
- This is the model that gets **pre-trained** on a large corpus



What can BERT do?



- Artificial [CLS] token is used as the vector to do classification from
 - Sentence pair tasks (entailment): feed both sentences into BERT
 - BERT can also do tagging by predicting tags at each word piece
- Devlin et al. (2019)



Natural Language Inference

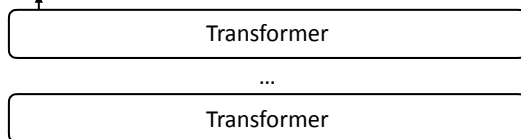
Premise		Hypothesis
A boy plays in the snow	<i>entails</i>	A boy is outside
A man inspects the uniform of a figure	<i>contradicts</i>	The man is sleeping
An older and younger man smiling	<i>neutral</i>	Two men are smiling and laughing at cats playing

- Long history of this task: “Recognizing Textual Entailment” challenge in 2006 (Dagan, Glickman, Magnini)
- Early datasets: small (hundreds of pairs), very ambitious (lots of world knowledge, temporal reasoning, etc.)



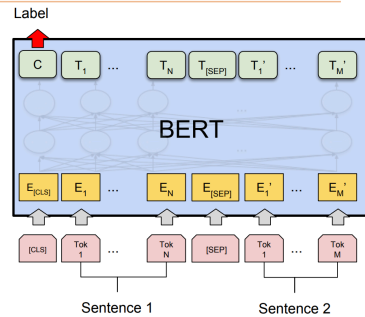
What can BERT do?

Entails (first sentence implies second is true)



[CLS] A boy plays in the snow [SEP] A boy is outside

- How does BERT model sentence pairs?
- Transformers can capture interactions between the two sentences, even though the NSP objective doesn't really cause this to happen



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



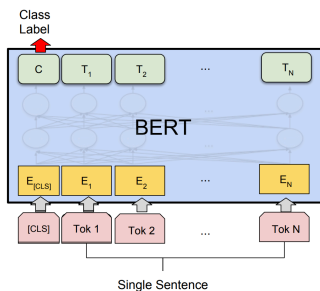
What can BERT NOT do?

- BERT **cannot** generate text (at least not in an obvious way)
- Can fill in MASK tokens, but can't generate left-to-right (well, you could put MASK at the end repeatedly, but this is slow)
- Masked language models are intended to be used primarily for "analysis" tasks



Fine-tuning BERT

- Fine-tune for 1-3 epochs, batch size 2-32, learning rate 2e-5 - 5e-5



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

- Large changes to weights up here (particularly in last layer to route the right information to [CLS])
- Smaller changes to weights lower down in the transformer
- Small LR and short fine-tuning schedule mean weights don't change much
- Often requires tricky learning rate schedules ("triangular" learning rates with warmup periods)

Subword Tokenization



Handling Rare Words

- Words are a difficult unit to work with. Why?
 - When you have 100,000+ words, the final matrix multiply and softmax start to dominate the computation, many params, still some words you haven't seen...
- Character-level models were explored extensively in 2016-2018 but simply don't work well — becomes very expensive to represent sequences



Subword Tokenization

- Subword tokenization: wide range of schemes that use tokens that are **between characters and words** in terms of granularity
- These “word pieces” may be full words or parts of words

_the _eco tax _port i co _in _Po nt - de - Bu is ...
- _ indicates the word piece starting a word (can think of it as the space character).

Sennrich et al. (2016)



Subword Tokenization

- Subword tokenization: wide range of schemes that use tokens that are **between characters and words** in terms of granularity
- These “word pieces” may be full words or parts of words

_the _eco tax _port i co _in _Po nt - de - Bu is...

Output: _le _port ique _éco taxe _de: _Pont - de - Bui s

- Can achieve transliteration with this, subword structure makes some translations easier to achieve

Sennrich et al. (2016)



Byte Pair Encoding (BPE)

- Start with every individual byte (basically character) as its own symbol
- ```
for i in range(num_merges):
 pairs = get_stats(vocab)
 best = max(pairs, key=pairs.get)
 vocab = merge_vocab(best, vocab)
```
- Count bigram character cooccurrences
  - Merge the most frequent pair of adjacent characters
  - Doing 8k merges => vocabulary of around 8000 word pieces. Includes many whole words
  - Most SOTA NMT systems use this on both source + target

Sennrich et al. (2016)



## Byte Pair Encoding (BPE)

|                    |                   |                    |                     |
|--------------------|-------------------|--------------------|---------------------|
| <b>Original:</b>   | furiously         | <b>Original:</b>   | tricycles           |
| <b>BPE:</b>        | _fur   iously (b) | <b>BPE:</b>        | _t   ric   y   cles |
| <b>Unigram LM:</b> | _fur   ious   ly  | <b>Unigram LM:</b> | _tri   cycle   s    |

|                    |                                                           |
|--------------------|-----------------------------------------------------------|
| <b>Original:</b>   | Completely preposterous suggestions                       |
| <b>BPE:</b>        | _Comple   t   ely   _prep   ost   erous   _suggest   ions |
| <b>Unigram LM:</b> | _Complete   ly   _pre   post   er   ous   _suggestion   s |

- ▶ BPE produces less linguistically plausible units than another technique based on a unigram language model: rather than greedily merge, find chunks which make the sequence look likely under a unigram LM
- ▶ Unigram LM tokenizer leads to slightly better BERT

Bostrom and Durrett (2020)

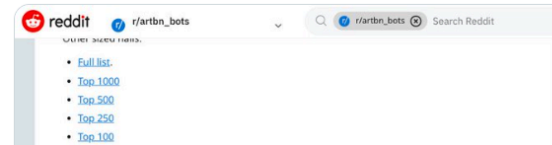


## What's in the token vocab?



Matthew Watkins  
@SoC\_trilogy

I've just found out that several of the anomalous GPT tokens ("TheNitromeFan", " SolidGoldMagikarp", " davidjl", " Smartstocks", " RandomRedditorWithNo", ) are handles of people who are (competitively? collaboratively?) counting to infinity on a Reddit forum. I kid you not.



| Rank | User                                    | Counts |
|------|-----------------------------------------|--------|
| 1    | <a href="#">/u/davidjl123</a>           | 163477 |
| 2    | <a href="#">/u/Smartstocks</a>          | 113829 |
| 3    | <a href="#">/u/atomicimploder</a>       | 103178 |
| 4    | <a href="#">/u/TheNitromeFan</a>        | 84581  |
| 5    | <a href="#">/u/SolidGoldMagikarp</a>    | 65753  |
| 6    | <a href="#">/u/RandomRedditorWithNo</a> | 63434  |
| 7    | <a href="#">/u/rideride</a>             | 59024  |
| 8    | <a href="#">/u/MooraeII</a>             | 57785  |
| 9    | <a href="#">/u/Removedpixel</a>         | 55080  |
| 10   | <a href="#">/u/Adinida</a>              | 48415  |
| 11   | <a href="#">/u/rschaosid</a>            | 47132  |



## Tokenization Today

- ▶ **All pre-trained** models use some kind of subword tokenization with a tuned vocabulary; usually between 50k and 250k pieces (larger number of pieces for multilingual models)
- ▶ As a result, classical word embeddings like GloVe **are not used**. All subword representations are randomly initialized and learned in the Transformer models

BERT results, BERT variants





## Evaluation: GLUE

| Corpus                          | Train | Test | Task                | Metrics                      | Domain              |
|---------------------------------|-------|------|---------------------|------------------------------|---------------------|
| Single-Sentence Tasks           |       |      |                     |                              |                     |
| CoLA                            | 8.5k  | 1k   | acceptability       | Matthews corr.               | misc.               |
| SST-2                           | 67k   | 1.8k | sentiment           | acc.                         | movie reviews       |
| Similarity and Paraphrase Tasks |       |      |                     |                              |                     |
| MRPC                            | 3.7k  | 1.7k | paraphrase          | acc./F1                      | news                |
| STS-B                           | 7k    | 1.4k | sentence similarity | Pearson/Spearman corr.       | misc.               |
| QQP                             | 364k  | 391k | paraphrase          | acc./F1                      | social QA questions |
| Inference Tasks                 |       |      |                     |                              |                     |
| MNLI                            | 393k  | 20k  | NLI                 | matched acc./mismatched acc. | misc.               |
| QNLI                            | 105k  | 5.4k | QA/NLI              | acc.                         | Wikipedia           |
| RTE                             | 2.5k  | 3k   | NLI                 | acc.                         | news, Wikipedia     |
| WNLI                            | 634   | 146  | coreference/NLI     | acc.                         | fiction books       |

Wang et al. (2019)



## Results

| System                | MNLI-(m/mm)      | QQP         | QNLI        | SST-2       | CoLA        | STS-B       | MRPC        | RTE         | Average     |
|-----------------------|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                       | 392k             | 363k        | 108k        | 67k         | 8.5k        | 5.7k        | 3.5k        | 2.5k        | -           |
| 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        | 79.9        | 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 <sub>BASE</sub>  | 84.6/83.4        | 71.2        | 90.1        | 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>91.1</b> | <b>94.9</b> | <b>60.5</b> | <b>86.5</b> | <b>89.3</b> | <b>70.1</b> | <b>81.9</b> |

- Huge improvements over prior work (even compared to ELMo)
- Effective at “sentence pair” tasks: textual entailment (does sentence A imply sentence B), paraphrase detection

Devlin et al. (2018)



## RoBERTa

- “Robustly optimized BERT”

- 160GB of data instead of 16 GB

- Dynamic masking: standard BERT uses the same MASK scheme for every epoch, RoBERTa recomputes them

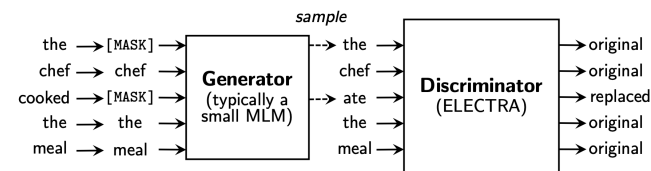
- New training + more data = better performance

| Model                                   | data  | bsz | steps | SQuAD (v1.1/2.0) | MNLI-m      | SST-2       |
|-----------------------------------------|-------|-----|-------|------------------|-------------|-------------|
| RoBERTa                                 |       |     |       |                  |             |             |
| with BOOKS + WIKI                       | 16GB  | 8K  | 100K  | 93.6/87.3        | 89.0        | 95.3        |
| + additional data (§3.2)                | 160GB | 8K  | 100K  | 94.0/87.7        | 89.3        | 95.6        |
| + pretrain longer                       | 160GB | 8K  | 300K  | 94.4/88.7        | 90.0        | 96.1        |
| + pretrain even longer                  | 160GB | 8K  | 500K  | <b>94.6/89.4</b> | <b>90.2</b> | <b>96.4</b> |
| BERT <sub>LARGE</sub> with BOOKS + WIKI |       |     |       |                  |             |             |
|                                         | 13GB  | 256 | 1M    | 90.9/81.8        | 86.6        | 93.7        |

Liu et al. (2019)

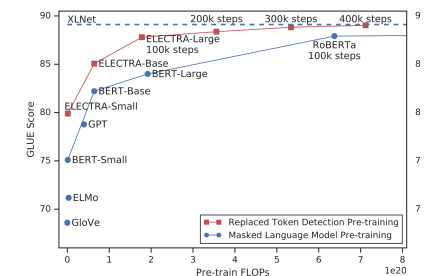


## ELECTRA



Clark et al. (2020)

- Discriminator to *detect* replaced tokens rather than a generator to actually *predict* what those tokens are
- More efficient, strong performance





## DeBERTa

- Slightly better variant

He et al. (2021)

$$A_{i,j} = \{H_i, P_{i|j}\} \times \{H_j, P_{j|i}\}^\top$$

$$= H_i H_j^\top + H_i P_{j|i}^\top + P_{i|j} H_j^\top + P_{i|j} P_{j|i}^\top \quad (2)$$

That is, the attention weight of a word pair can be computed as a sum of four attention scores using disentangled matrices on their contents and positions as *content-to-content*, *content-to-position*, *position-to-content*, and *position-to-position*<sup>2</sup>.

| Model                    | CoLA Mcc    | QQP Acc     | MNLI-m/mm Acc    | SST-2 Acc   | STS-B Corr  | QNLI Acc    | RTE Acc     | MRPC Acc    | Avg.         |
|--------------------------|-------------|-------------|------------------|-------------|-------------|-------------|-------------|-------------|--------------|
| BERT <sub>large</sub>    | 60.6        | 91.3        | 86.6/-           | 93.2        | 90.0        | 92.3        | 70.4        | 88.0        | 84.05        |
| RoBERTa <sub>large</sub> | 68.0        | 92.2        | 90.2/90.2        | 96.4        | 92.4        | 93.9        | 86.6        | 90.9        | 88.82        |
| XLNet <sub>large</sub>   | 69.0        | 92.3        | 90.8/90.8        | <b>97.0</b> | 92.5        | 94.9        | 85.9        | 90.8        | 89.15        |
| ELECTRA <sub>large</sub> | 69.1        | <b>92.4</b> | 90.9/-           | 96.9        | 92.6        | 95.0        | 88.0        | 90.8        | 89.46        |
| DeBERTa <sub>large</sub> | <b>70.5</b> | 92.3        | <b>91.1/91.1</b> | 96.8        | <b>92.8</b> | <b>95.3</b> | <b>88.3</b> | <b>91.9</b> | <b>90.00</b> |



## Using BERT

- HuggingFace Transformers: big open-source library with most pre-trained architectures implemented, weights available

- Lots of standard models...

and “community models”

Model architectures

Transformers currently provides the following NLU/NLG architectures:

1. BERT (from Google) released with the paper [BERT: Pre-training of Deep Understanding](#) by Jacob Devlin, Ming-Wei Chang, Kenton Lee and Krist
2. GPT (from OpenAI) released with the paper [Improving Language Understanding](#) by Radford, Karthik Narasimhan, Tim Salimans and Ilya Sutskever.
3. GPT-2 (from OpenAI) released with the paper [Language Models are Un](#) Jeffrey Wu\*, Rewon Child, David Luarn, Dario Amodei\*\* and Ilya Sutske
4. Transformer-XL (from Google/CMU) released with the paper [Transform](#) Fixed-Length Context by Zihang Dai\*, Zhilin Yang\*, Yiming Yang, Jaime
5. XLNet (from Google/CMU) released with the paper [XLNet: Generalized Understanding](#) by Zhilin Yang\*, Zihang Dai\*, Yiming Yang, Jaime Carbon
6. XLM (from Facebook) released together with the paper [Cross-lingual L](#) and Alexis Conneau.
7. RoBERTa (from Facebook), released together with the paper [A Robustly](#)

mrm8488/spanbert-large-finetuned-tacred ★

mrm8488/xlm-multi-finetuned-xquadv1 ★

nlpauub/bert-base-greek-uncased-v1 ★

nlpotown/bert-base-multilingual-uncased-sentiment ★

patrickvonplaten/reformer-crime-and-punish ★

redewiedergabe/bert-base-historical-german-rw-cased ★

roberta-base ★

severinsimmler/literary-german-bert ★

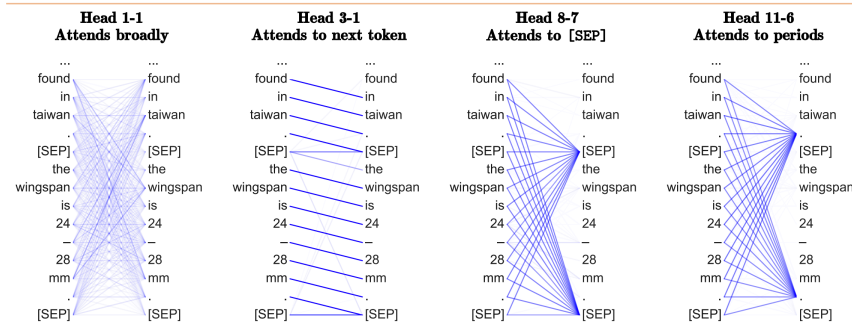
seyonec/ChemBERTa-zinc-base-v1 ★

...

...



## What does BERT learn?

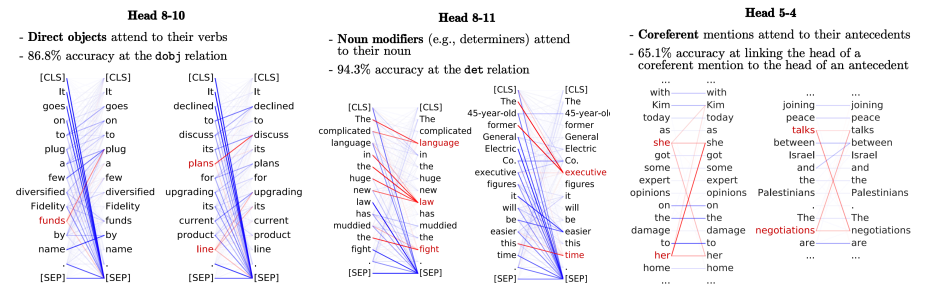


- Heads on transformers learn interesting and diverse things: content heads (attend based on content), positional heads (based on position), etc.

Clark et al. (2019)



## What does BERT learn?



- Still way worse than what supervised systems can do, but interesting that this is learned organically

Clark et al. (2019)

## Applying BERT



## Two Tasks

- Compared to ELMo, BERT is very good at **sentence-pair** tasks
  - Paraphrase detection
  - Semantic textual similarity
  - Textual entailment**
  - Question answering** (not really a sentence pair, but it's a pair of inputs)



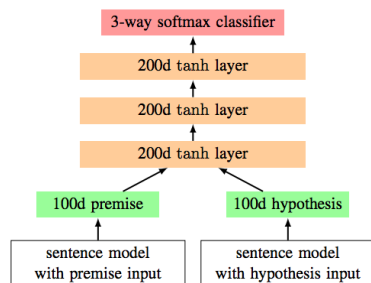
## SNLI Dataset

- Show people captions for (unseen) images and solicit entailed / neural / contradictory statements

- >500,000 sentence pairs
- One possible architecture:

300D BiLSTM: 83% accuracy  
(Liu et al., 2016)

- One of the first big successes of LSTM-based classifiers (sentiment results were more marginal)



Bowman et al. (2015)



## MNLI Dataset

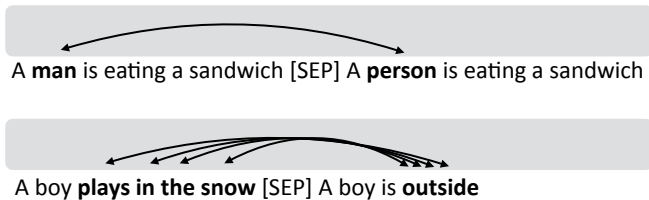
- Drawn from multiple genres of text

| Premise                                                                                                                        | Label                | Hypothesis                                                |
|--------------------------------------------------------------------------------------------------------------------------------|----------------------|-----------------------------------------------------------|
| <b>Fiction</b>                                                                                                                 |                      |                                                           |
| The Old One always comforted Ca'daan, except today.                                                                            | <i>neutral</i>       | Ca'daan knew the Old One very well.                       |
| <b>Letters</b>                                                                                                                 |                      |                                                           |
| Your gift is appreciated by each and every student who will benefit from your generosity.                                      | <i>neutral</i>       | Hundreds of students will benefit from your generosity.   |
| <b>Telephone Speech</b>                                                                                                        |                      |                                                           |
| yes now you know if if everybody like in August when everybody's on vacation or something we can dress a little more casual or | <i>contradiction</i> | August is a black out month for vacations in the company. |
| <b>9/11 Report</b>                                                                                                             |                      |                                                           |
| At the other end of Pennsylvania Avenue, people began to line up for a White House tour.                                       | <i>entailment</i>    | People formed a line at the end of Pennsylvania Avenue.   |

Williams et al. (2018)



## How do models do it?



- Transformers can easily learn to spot words or short phrases that are transformed
- But**, models are often overly sensitive to lexical overlap

Williams et al. (2018)



## Question Answering

- Many types of QA:
- We'll focus on **factoid questions** being answered **from text**
  - E.g., "What was Marie Curie the first female recipient of?" — unlikely you would have this answer in a database
  - Not appropriate: "When was Marie Curie born?" — probably answered in a DB
  - Not appropriate: "Why did World War II start?" — no simple answer



## SQuAD

Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the **Nobel Prize**. Famous musicians include Władysław Szpilman and Frédéric Chopin. Though Chopin was born in the village of Żelazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.

Answer = Nobel Prize

- Assume we know a passage that contains the answer. More recent work has shown how to retrieve these effectively (will discuss when we get to QA)



## SQuAD

Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the **Nobel Prize**. ...

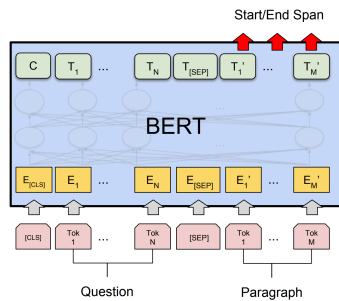
- Predict answer as a pair of (start, end) indices given question  $q$  and passage  $p$ ; compute a score for each word and softmax those

$$P(\text{start} \mid q, p) = \begin{matrix} 0.01 & 0.01 & 0.01 & 0.85 & 0.01 \\ \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ \text{recipient of the } \mathbf{Nobel Prize} . \end{matrix}$$

$P(\text{end} \mid q, p)$  = same computation but different params



## QA with BERT



What was Marie Curie the first female recipient of ? [SEP] One of the most famous people born in Warsaw was Marie ...

- In a couple lectures, we will look at what BERT learns when trained on this kind of data

Devlin et al. (2019)



## Takeaways

- Pre-trained models and BERT are very powerful for a range of NLP tasks
- These models have enabled big advances in NLI and QA specifically
- Next time: pre-trained decoders (GPT-3) and encoder-decoder models (T5)