

Today: improvements on BERT

1. training improvements \Rightarrow ROBERTa
more data
 2. longer sequences
Transformer XL (XLNet)
 3. more efficient pretraining objectives
ELECTRA
 4. smaller models
ALBERT
-

ROBERTa: very simple
Collection of modifications

1. train w/ bigger batches
 - \hookrightarrow smaller # of batches w/
larger batch size
 - \hookrightarrow gradient accumulation
to bypass GPU mem. limitations
2. no next sentence prediction
 - \hookrightarrow downstream perf. unaffected
 - \hookrightarrow [CLS] token gets no pretraining
3. pretrain on more data

↳ 16 GB → 160 GB

↳ common crawl,
URLs from Reddit

4. pretrain for longer (more batches/epochs)

↳ 500k steps

TransformerXL

BERT has a fixed token limit of 512 for its inputs. how can we model longer sequences?

↳ idea: add a recurrent mechanism that connects adjacent segments


↳ no gradient flow to previous segment
hidden states from prev. segment are cached

→ practical limit to this extended context window

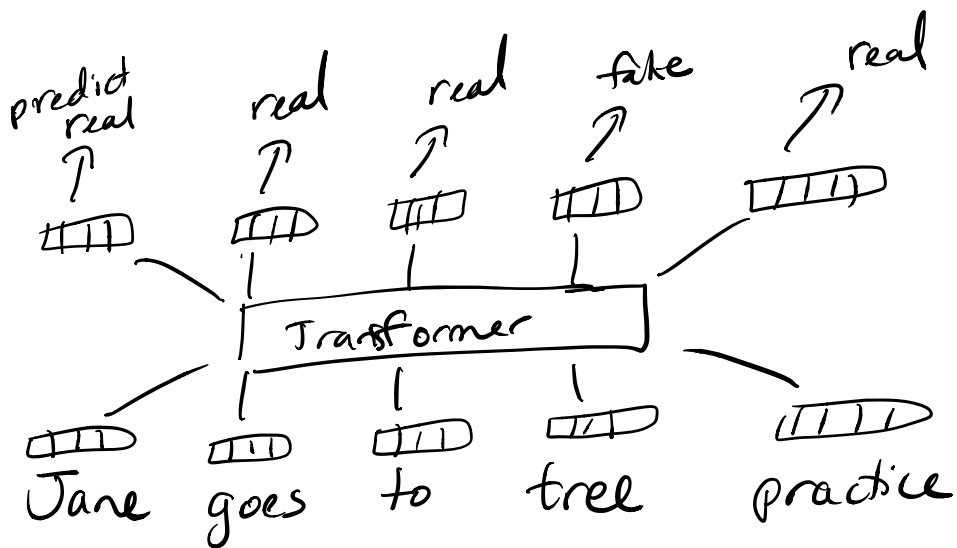
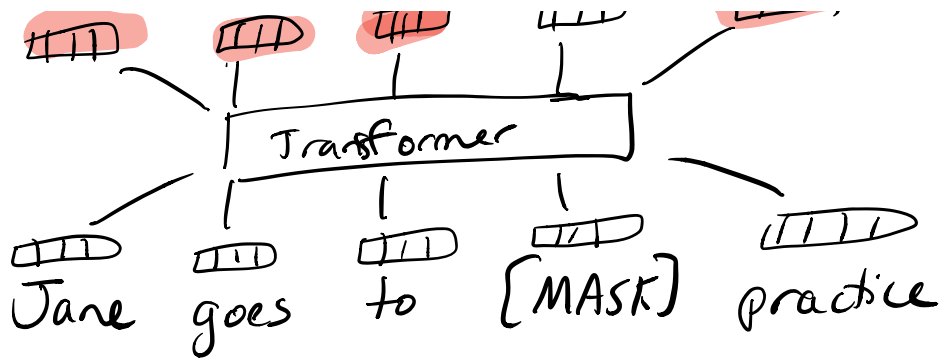
→ 900 words for TransformerXL

ELECTRA - cheaper obj. fn

predict "baseball"



The diagram shows a sequence of tokens represented by vertical bars. The first bar is red, the second is red, the third is red, and the fourth is red. An arrow points from the third bar to the text "predict 'baseball'". The fourth bar is circled in red.



↳ how do I decide which words to replace and with what?

↳ "generator" ⇒ coming up w/ fake words

↳ train a small BERT model

Jane goes to [MASK] practice

↓
 football
 basketball
 baseball

- ↳ sampled words from generator form fake words for ELECTRA
 - every single token is associated with a prediction of real/fake, not just 15% of words as in BERT
-

ALBERT - more params != better model

- cross-layer param sharing
 - Q, k, V projection matrices
 - W matrices in FF layers } shared across all layers

BERT - large: 334 M params

ALBERT-large: 18M params

↳ what if we make our shared set of params bigger?

ALBERT-XXL: 235 M params, 4096 d hidden state size

↳ outperforms BERT-large