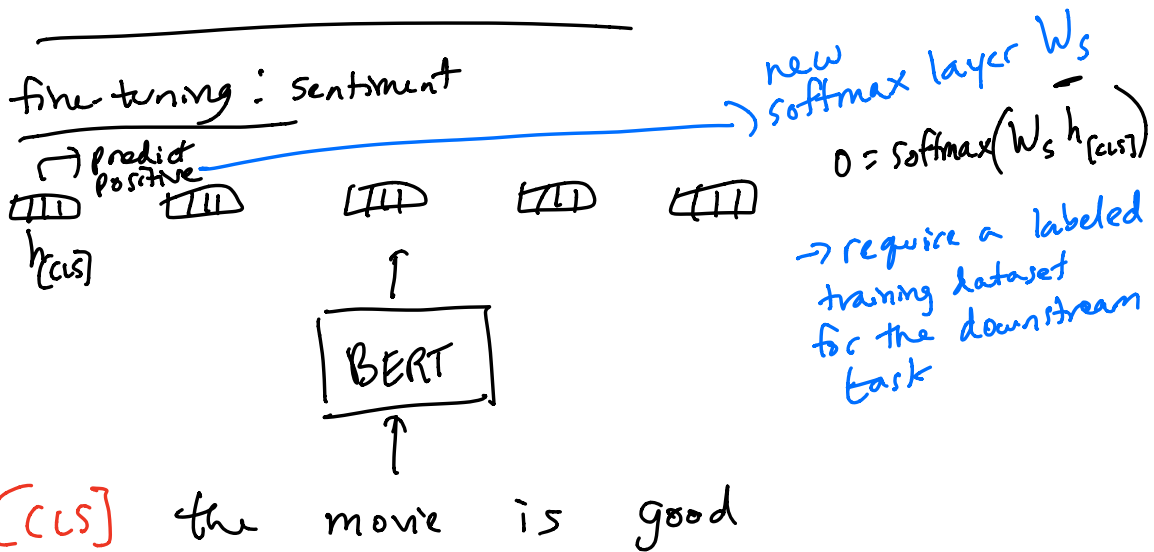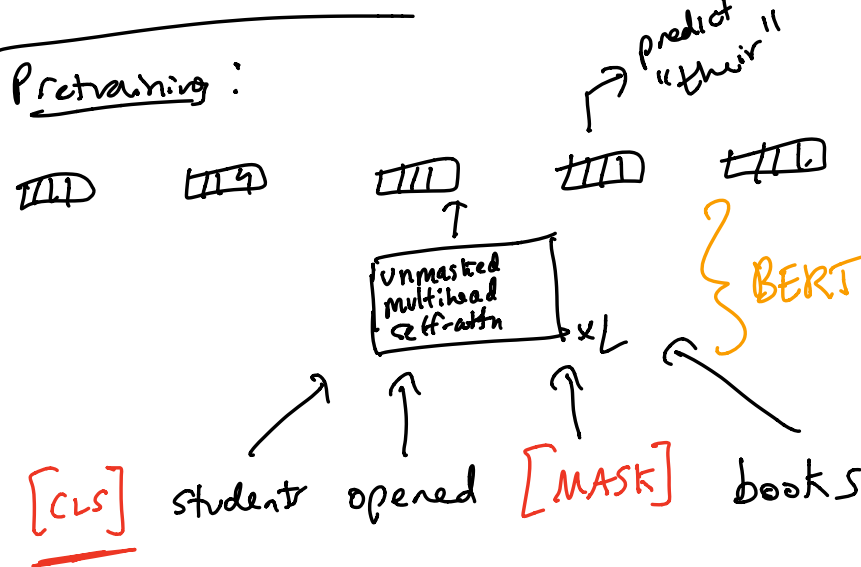# BERT:

↳ example of the encoder-only paradigm

↳ **Pretraining:** train w/ self-supervised obj called "masked LM"

↳ **Finetuning:** process of adapting the pretrained model to a particular downstream task

---

**Pretraining:**

→ predict "their"



```
unmasked
multihead
self-attn    × L
```

} BERT

[CLS]  students  opened  [MASK]  books

---

**fine-tuning : sentiment**

→ predict positive

new softmax layer $W_s$

$$O = softmax(W_s \, h_{[CLS]})$$

→ require a labeled training dataset for the downstream task

$h_{[CLS]}$

| BERT |

[CLS]  the  movie  is  good

## extractive QA

↳ inputs: + a question
         + a paragraph that contains the answer

↳ output: a span of the paragraph
         that answers the question

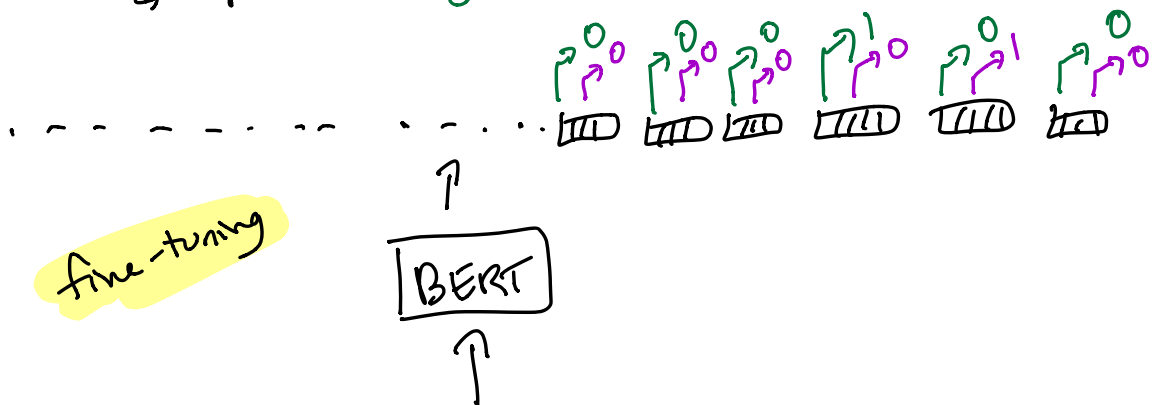↳ datasets: SQuAD v1,2 ; QuAC/CoQa,
           HotpotQA , ...

Q: Who starred in the Matrix as Neo?

P: $v_1$ $v_2$ $v_3$ ..... Neo was played by
   actor Keanu Reeves ... ..

A: (i, j)

---

How do we use BERT for extractive QA?
↳ 2 softmax layers on each token in passage
↳ predict beginning, end index of answer span

fine-tuning

BERT

[CLS] Who stored in the Matrix [SEP] $w_1$ $w_2$ $w_3$ ... Keanu Reeves ...

how do we select an answer span at test time?

→ find the span $W_{i...j}$ that maximizes

$$P_{START}(i) \cdot P_{END}(j)$$

↳ exclude spans where $j < i$

↳ exclude spans longer than a threshold

---

advanced variants of BERT:

↳ pretraining improvements  ⇒ RoBERTa
    ↳ more data

↳ longer sequences during pretraining
- BERT ⇒ 512 tokens max
- XLNet ⇒ 900 tokens

↳ more pretraining objectives

↳ ELECTRA

↳ → "students opened [MASK] books"

↳ "students opened Monster books"
    ↓ ↓     ↓     ↓     ↓
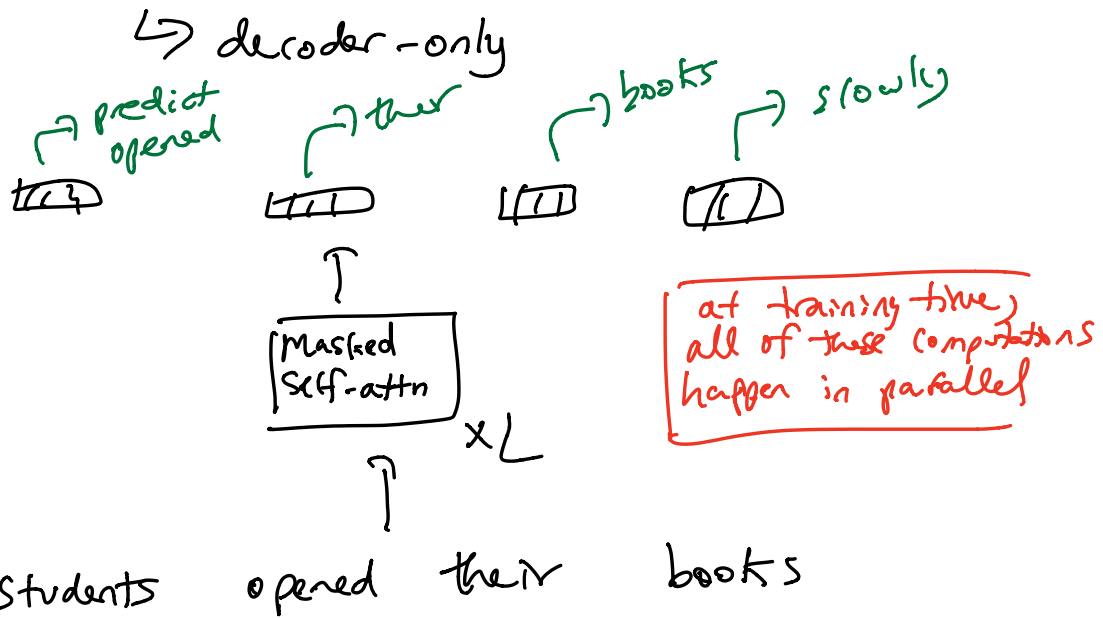  (real), corrupted?  real  Corrupted  real

↳ smaller models
  ↳ tinyBERT, distilBERT, ALBERT
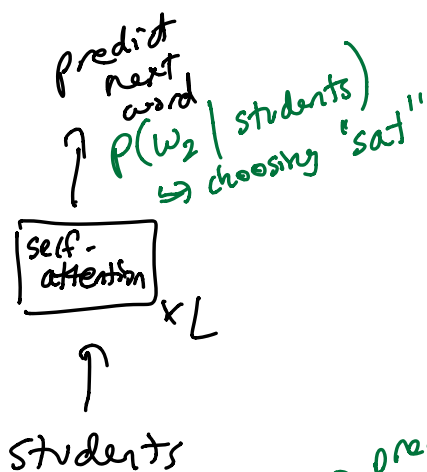
↳ distillation

↳ pruning

---

Transformer LMs at test-time:

---

↳ decoder-only

→ predict opened    → their    → books    → slowly



Masked Self-attn

×L

Students    opened    their    books

at training time,
all of these computations
happen in parallel

⇒ those four predictions can happen
simultaneously b/c we already know
the identity of the gold next tokens
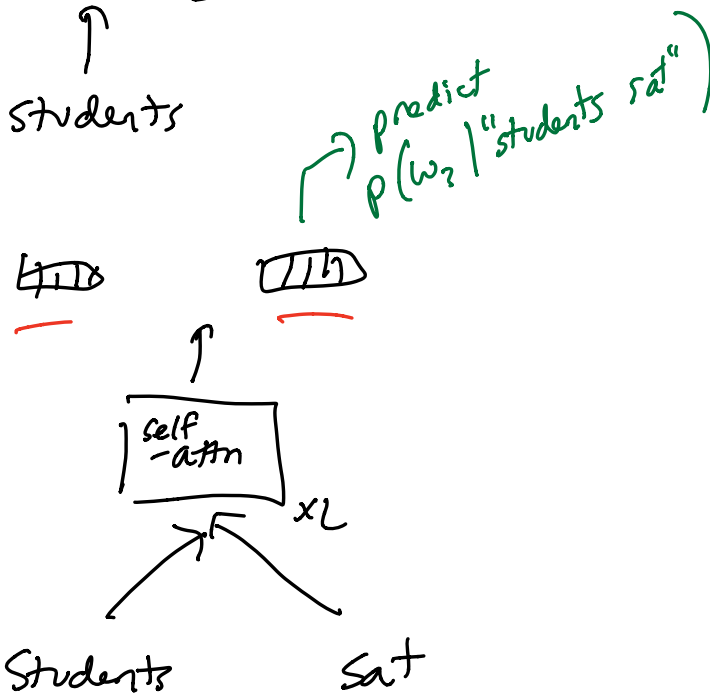                    ground-truth

test-time:

timestep $\mathbb{1}$:

predict next word

$\uparrow$ P($w_2$ | students)
$\hookrightarrow$ choosing "sat"

[self-attention] $\times L$

$\uparrow$

students

autoregressive decoding / LM

1. first have to choose a next word from P($w_2$ | students)
   $\hookrightarrow$ sample
   $\hookrightarrow$ argmax

2. concat this word to the input and then repeat step 1

predict P($w_2$ | "students sat")

[self -attn] $\times L$

Students          sat

at test-time, I have to **decode** the output word-by-word, because I don't have access to the gold next tokens

$\hookrightarrow$ decoding algorithms