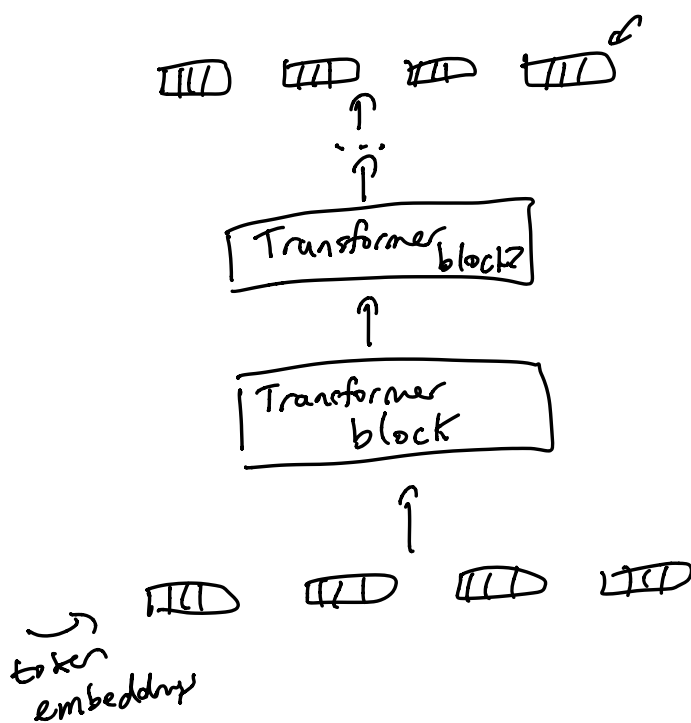


Today:

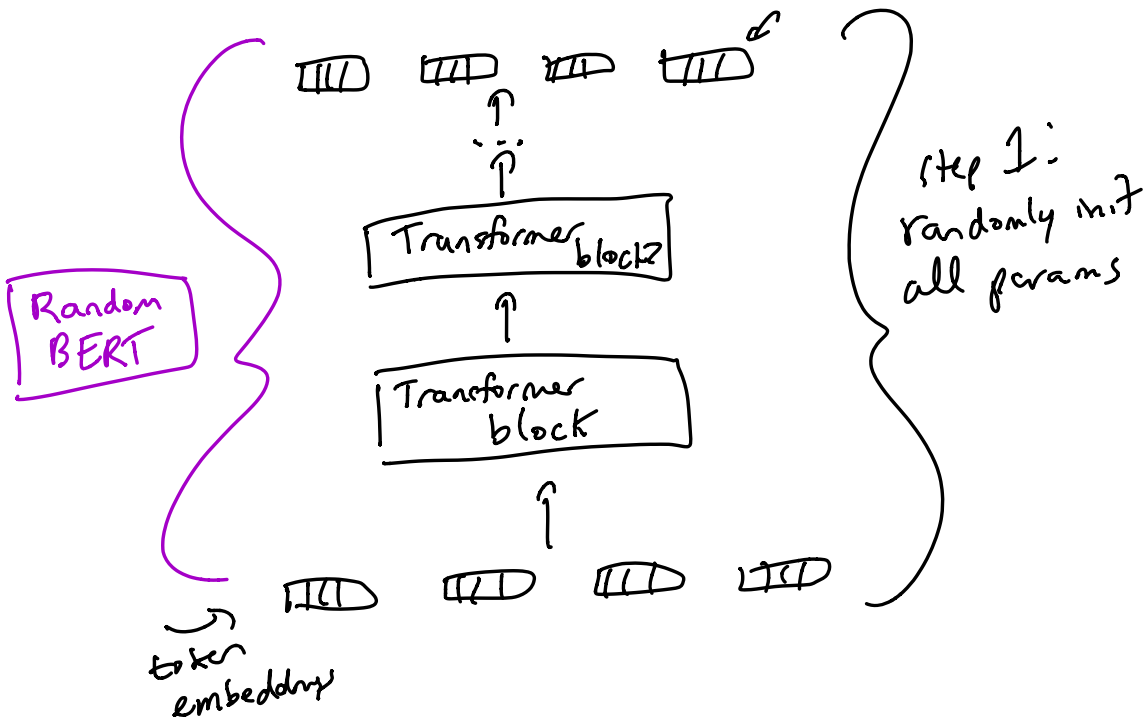
- "trained from scratch" vs. "fine-tuned"
- text-to-text transfer learning
- decoding algorithms to generate text

BERT for downstream tasks:

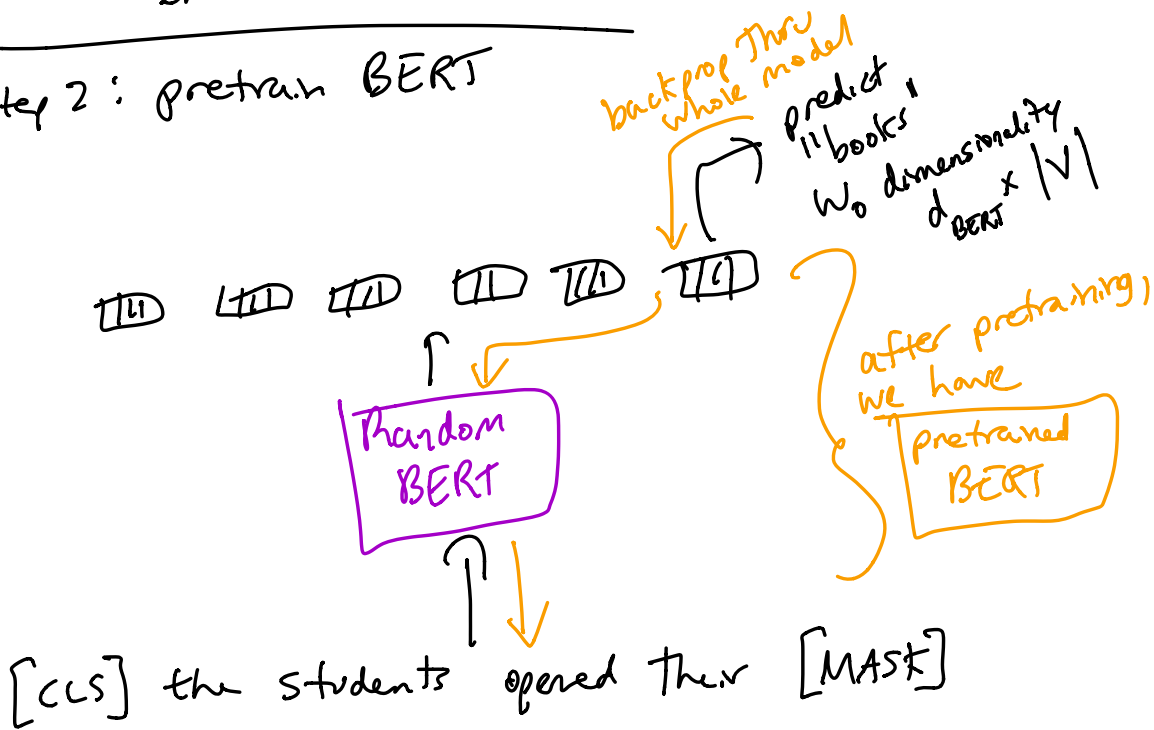
model: deep Transformer



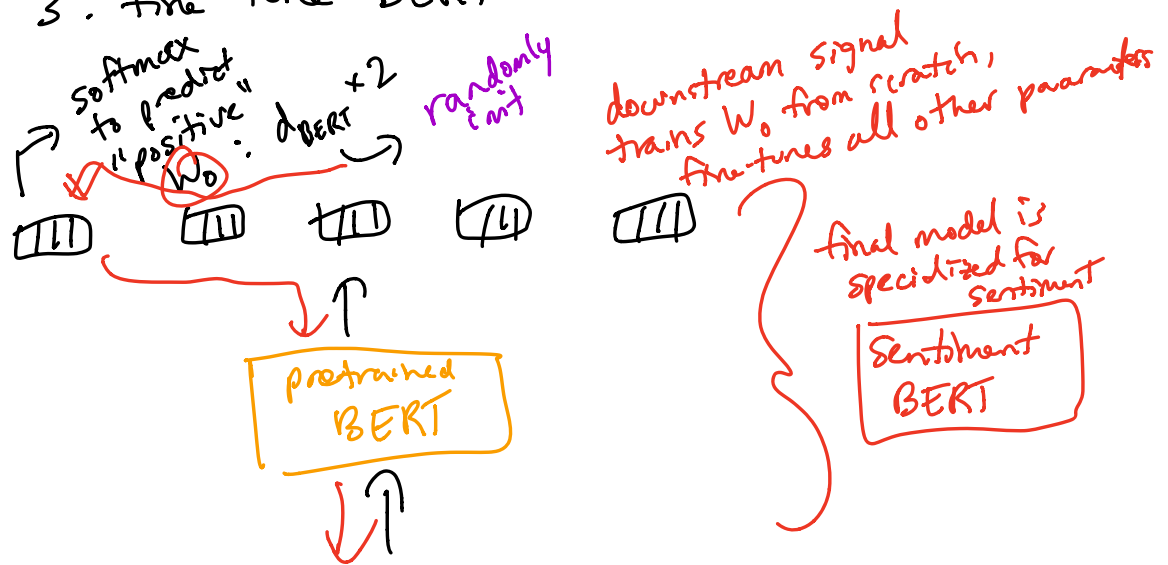
What are the model params?  
 $W_q, W_k, W_v$  } for each head,  
 $W_{ff}$  } feedforward matrix for each block  
word embeddings (maybe pos embeddings)  
softmax matrix  $W_o$



Step 2: pretrain BERT



### Step 3: fine-tune BERT



[CLS] this movie was awesome

### Limitations of BERT:

What tasks can BERT not solve?

↳ text generation

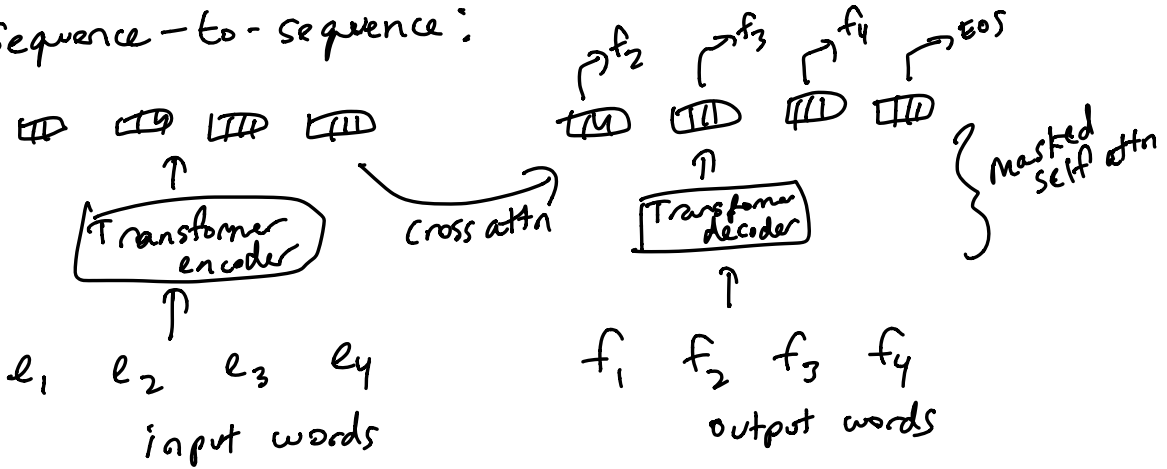
↳ translation, summarization

Can we develop self-supervised pretraining obj; that covers all types of NLP tasks

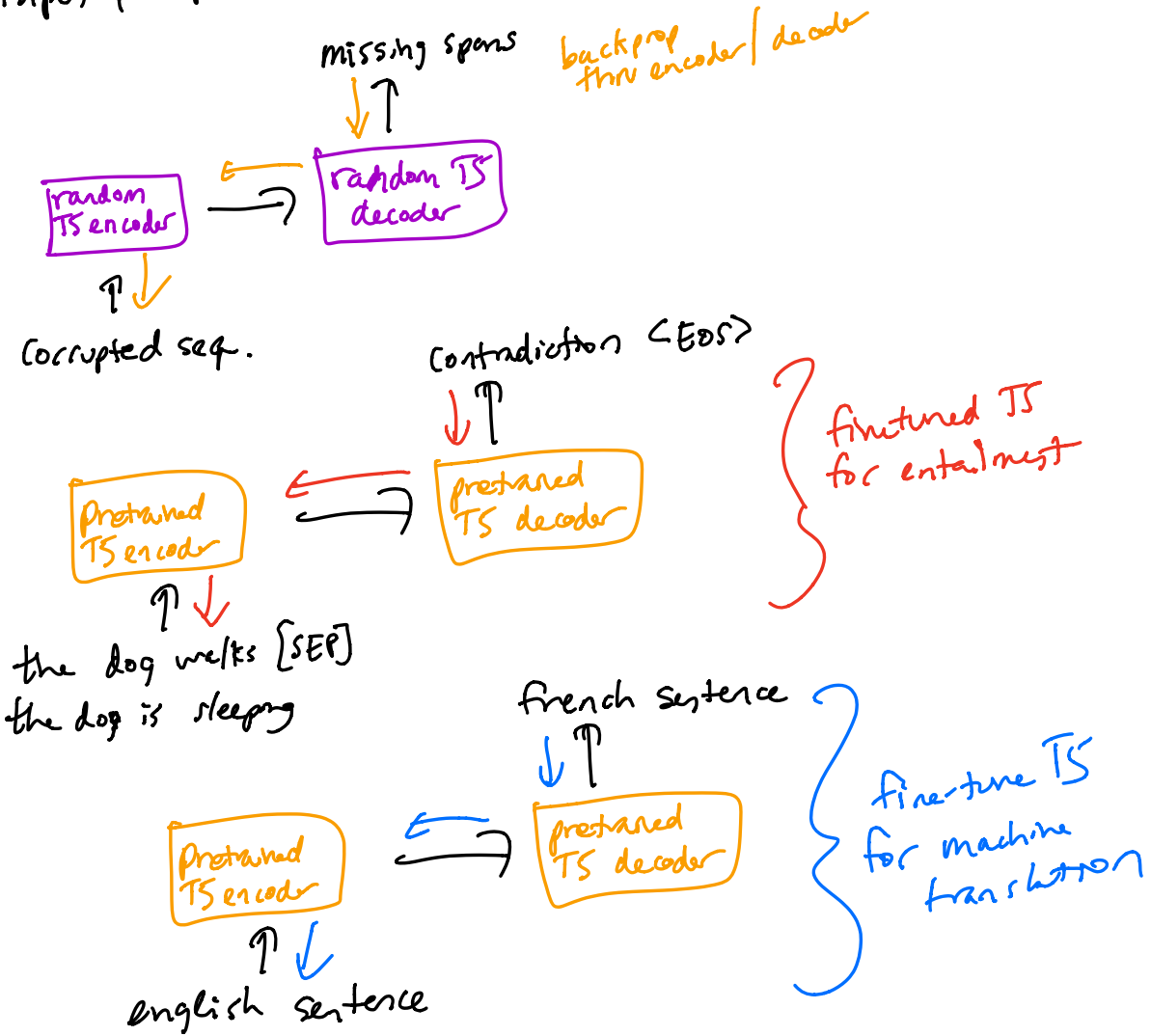
↳ classification, seq. labeling, generation

TS paper: reformulate every NLP task as a generation problem, "text-to-text"

Sequence-to-sequence:



input/output for TS:



decoder-only version of TS!

thanks for  $\langle X \rangle$  me to your party  $\langle X \rangle$  for inviting  $\langle \text{Eos} \rangle$   
marked self-attention, no cross attention