Lec10. Transformers and Pretraining

1. subword modeling

We assume a fixed vocab of tens of thousands of words, built from the training set.
All *novel* words seen at test time are mapped to a single UNK.



Variation/misspellings/novel items -> UNK

--> The byte-pair encoding algorithm



2. Motivating model pretraining from word embeddings
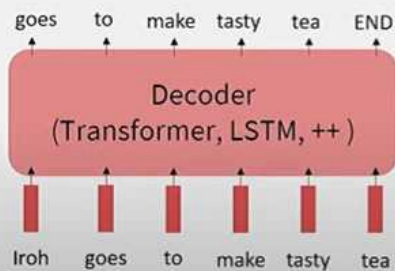
word2vec의 한계 ->

- The Pretraining/ Finetuning Paradigm



## The Pretraining / Finetuning Paradigm

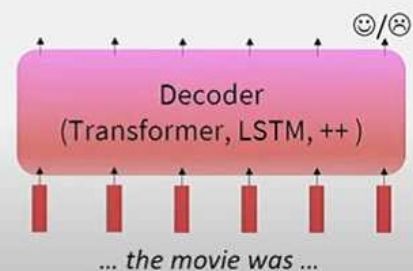Pretraining can improve NLP applications by serving as parameter initialization.

**Step 1: Pretrain (on language modeling)**
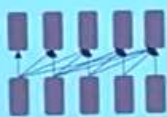Lots of text; learn general things!

goes    to    make    tasty    tea    END

Decoder
(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on your task)**
Not many labels; adapt to the task!

☺/☹

Decoder
(Transformer, LSTM, ++ )

... the movie was ...

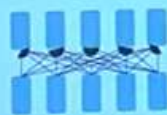3. Model pretraining three ways



## Pretraining for three types of architectures

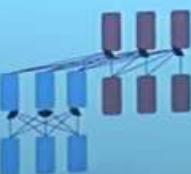The neural architecture influences the type of pretraining, and natural use cases.

**Decoders**
- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words

**Encoders**
- Gets bidirectional context – can condition on future!
- Wait, how do we pretrain them?

**Encoder-Decoders**
- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Stanford

1. Decoders



## Pretraining decoders

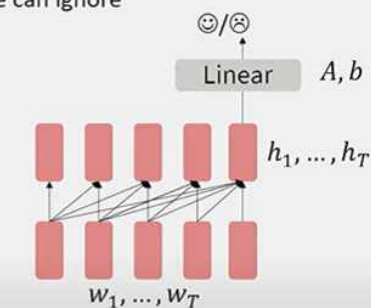When using language model pretrained decoders, we can ignore that they were trained to model $p(w_t|w_{1:t-1})$.

We can finetune them by training a classifier on the last word's hidden state.

$$h_1, ..., h_T = \text{Decoder}(w_1, ..., w_T)$$
$$y \sim Ah_t + b$$

Where $A$ and $b$ are randomly initialized and specified by the downstream task.

Gradients backpropagate through the whole network.

[Note how the linear layer hasn't been pretrained and must be learned from scratch.]

Stanford

eg. GPT(Generative Pretrained Transformer)
- Transformer decoder with 12 layers.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
- Trained on BooksCorpus

2. Encoders



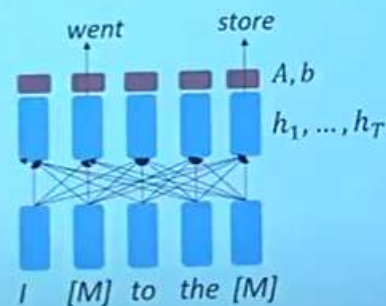## Pretraining encoders: what pretraining objective to use?

So far, we've looked at language model pretraining. But **encoders get bidirectional context**, so we can't do language modeling!

Idea: replace some fraction of words in the input with a special [MASK] token; predict these words.

$$h_1, ..., h_T = \text{Encoder}(w_1, ..., w_T)$$
$$y_i \sim Aw_i + b$$

Only add loss terms from words that are "masked out." If $\tilde{x}$ is the masked version of $x$, we're learning $p_\theta(x|\tilde{x})$. Called **Masked LM**.
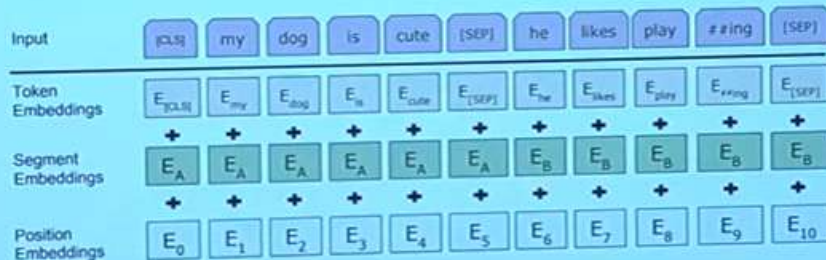
[Devlin et al., 2018]

eg. BERT(Bidirectional Encoder Representations from Transforms)
- "Masked LM" objective and released the weights of a pretrained Transformer
-

- Two models: (BERT-base/large)
- Extensions: RoBERTa, SpanBERT, +++
  3. Encoder-Decoder



-eg. T5: it can be finetuned to answer a wide range of questions, retrieving knowledge from its parameters.

4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning
- GPT3: has 175 billions parameters.
- -> without gradient steps