

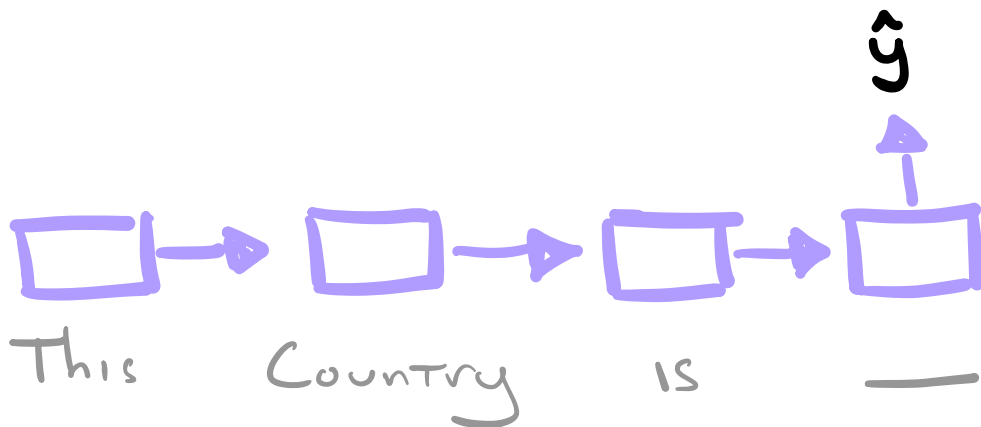
DS 4440 RNNs \rightarrow Contextual
embeddings
 \rightarrow Transformers

Autoregressive generative model:

$$P(x_T | x_{T-1} \dots x_1)$$

$$P(x) = \prod_T P(x_T | x_{<T})$$

e.g. Language Modeling.

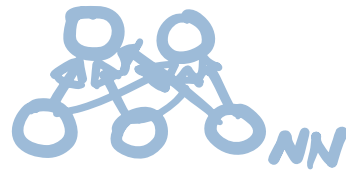


Language Modeling is a kind of
Self-Supervised learning, often
used for Pre-Training

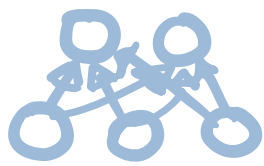
Unlabeled data

+ Self Supervised loss

PreTrain



e.g. The Internet



+ \mathcal{L} and (x, y) | Fine-Tune
OR

(x, y) | In Context learning (later)

Contextualized Word Vectors

Recall (9.30 lecture) that $W2V$ yields embeddings of words

"Cat" $\xrightarrow{W2V}$ $[x \ y \ z]$ | STATIC; always $x \ y \ z$

"Cute Cat"
 \downarrow
 $[x \ y \ z]$

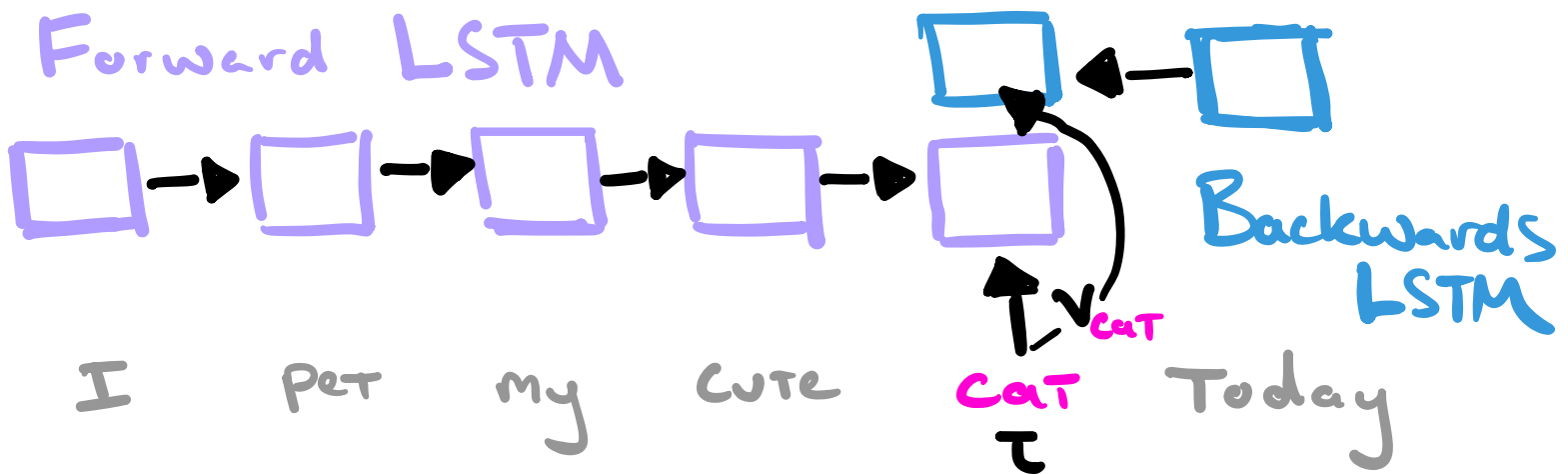
"fat Cat"
 \downarrow
 $[x \ y \ z]$

IDEA Train an RNN for language modeling; extract contextualized embeddings

ELMo

Embeddings from Language Models

[Peters et al., '18]



$$L = -P(\text{cat} | v_{\text{cat}}, \vec{h}_{\tau}) - P(\text{cat} | v_{\text{cat}}, \overleftarrow{h}_{\tau})$$

Self Supervised!
LM

$$\text{ELMo}_{\text{cat}} = [v_{\text{cat}}; \vec{h}_{\tau}; \overleftarrow{h}_{\tau}]$$

Can now use $ELMo_{CAT}$ VECTORS as features for downstream tasks

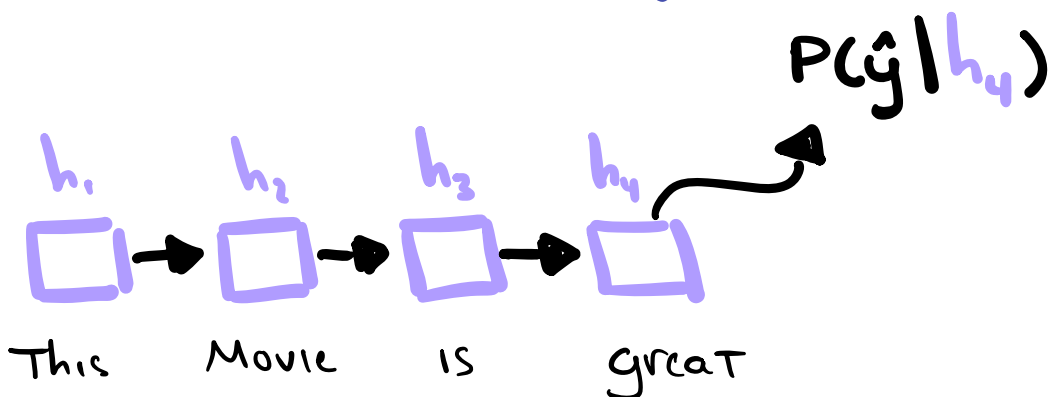
fine-Tuning!

This works pretty well! BUT, $ELMo$ is based on RNNs and so it's **Sloooooow** due to sequential processing

Enter Transformers which do away with recurrence using Self-attention.

(Vanilla)

ATTENTION!



Idea Weight hidden STATES

$$a_j = W_a \cdot h_j$$

$\langle 1 \times h \rangle \quad \langle h \rangle$

$$\bar{h} = \sum_{j=1}^n \alpha_j h_j$$

$$\alpha = \text{SoftMax}(a) \quad | \quad P(y|\bar{h})$$

Self-Attention

Generalizing: Assume keys k and queries q .

$$a_j = S(k_j, q) \quad \alpha = \text{SoftMax}(a)$$

$$\bar{h} = \sum_j \alpha_j v_j$$

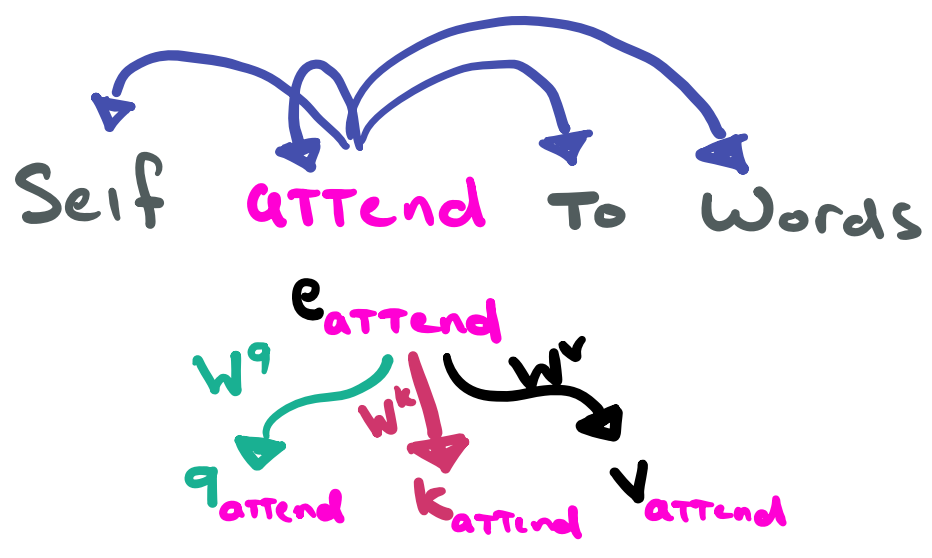
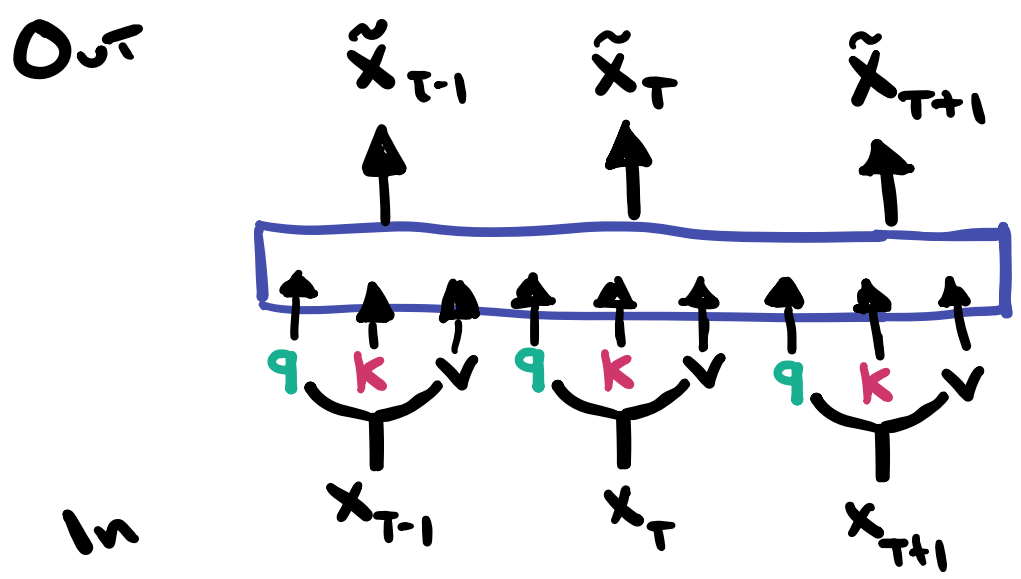
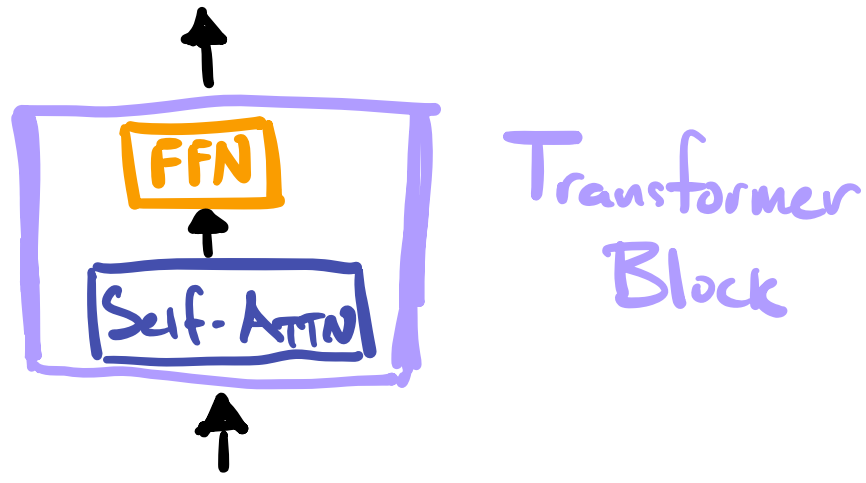
e.g. $(k_j \cdot q) / \sqrt{d}$ "dot product ATTN"
or
Scale by dims

$$v_j^T \text{Tanh}(W_k k_j + W_q q) \quad \text{"MLP ATTN"}$$

Transformers

[Vaswani et al '17]

Replace recurrence with repeated blocks of self-attention.



ATTEND

$$\begin{array}{l} \text{Self} \\ \text{ATTEND} \\ \text{to} \\ \text{Words} \end{array} \begin{array}{l} S(K_{\text{self}}, q_{\text{attend}}) \\ S(K_{\text{attend}}, q_{\text{attend}}) \\ S(K_{\text{to}}, q_{\text{attend}}) \\ S(K_{\text{words}}, q_{\text{attend}}) \end{array} \rightarrow \text{SM} \left(\begin{array}{c} \phantom{S(K_{\text{self}}, q_{\text{attend}})} \\ \phantom{S(K_{\text{attend}}, q_{\text{attend}})} \\ \phantom{S(K_{\text{to}}, q_{\text{attend}})} \\ \phantom{S(K_{\text{words}}, q_{\text{attend}})} \end{array} \right)$$

α_{attend}

$$\tilde{x}_{\text{attend}} = \sum_j \alpha_j^{\text{attend}} v_j$$

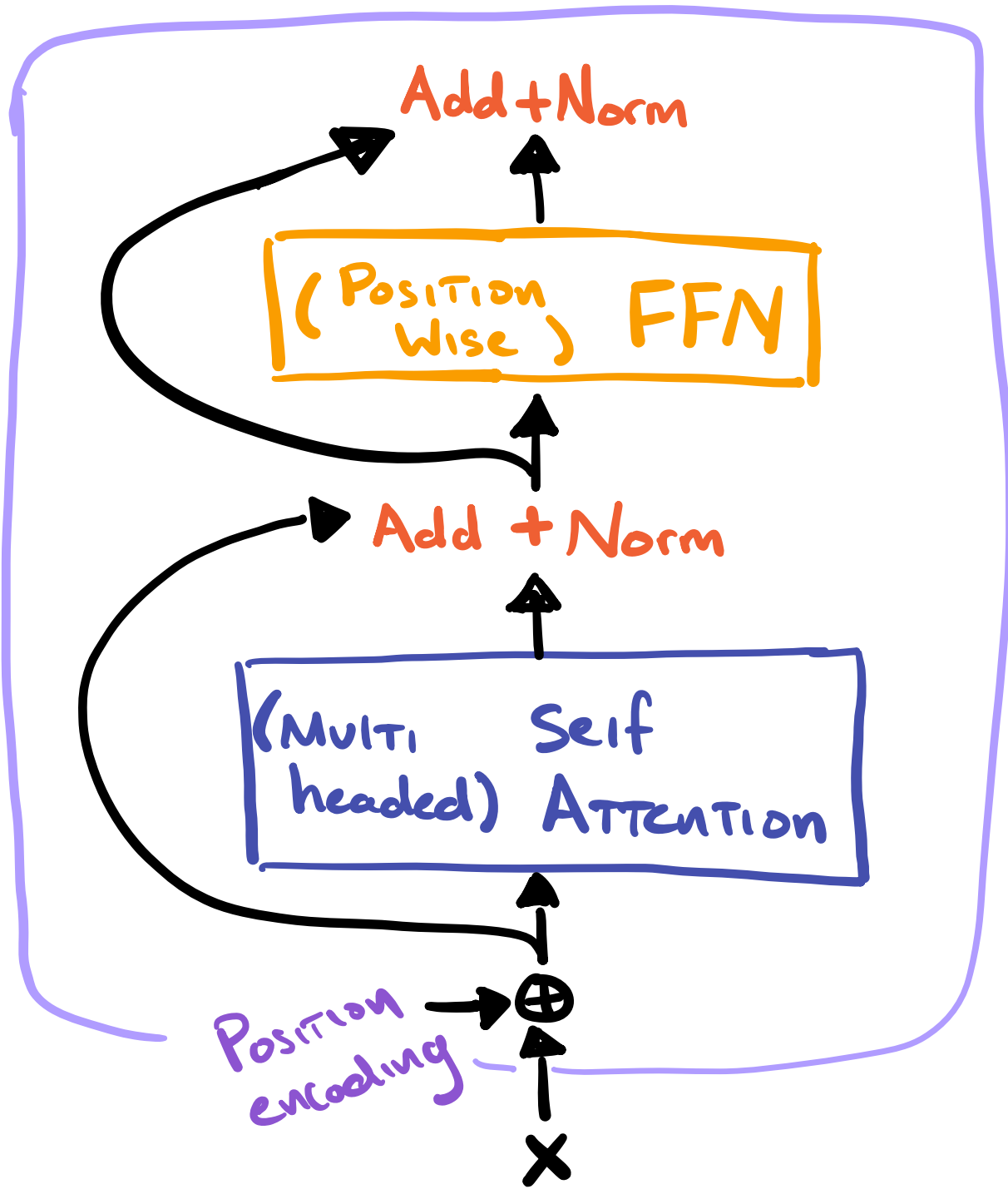
(See Colab exercise!)

In Multi-headed Attention we define L attention heads with independent params W_q^l, W_k^l, W_v^l .

Outputs are combined (e.g., \oplus).

After Self Attention, FFN (Position wise)

$$\text{FFN}(\tilde{x}_j)$$



Another detail is **add + Norm** which provides a residual shortcut and imposes a layer norm, which is like batch norm (10/7 lecture) but Transposed.

Positional encodings address the
issue of token positions (why
is this needed?)