# CS 378 Lecture 24
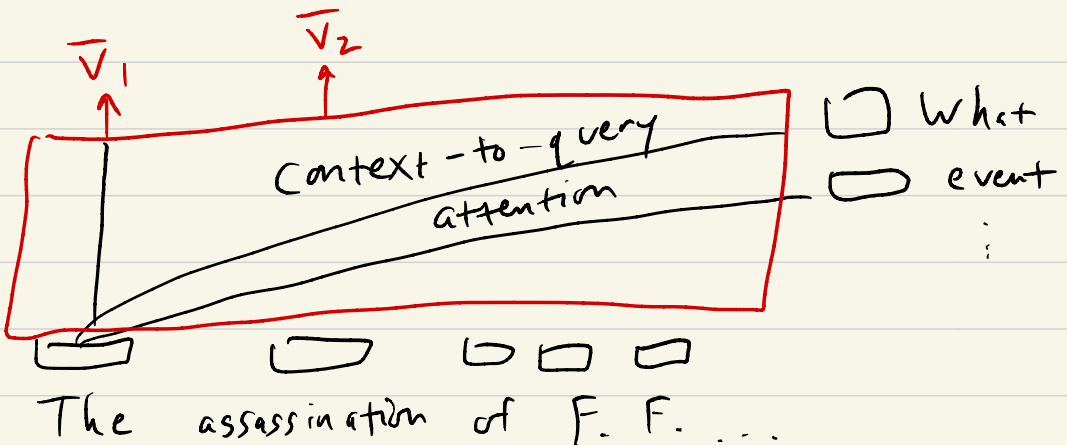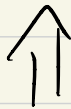
Today
1. Self-attention for language modeling
2. Transformers
3. BERT
4. Analysis + results of BERT

Recap    QA        rest of the model

↑

$\overline{V_1}$        $\overline{V_2}$

context-to-query
attention

☐ What
☐ event
⋮

The assassination of F. F. ...

ELMo: train a RNN LM on lots
of data, use it to produce
"contextualized" embeddings

Self-attention

Lang modeling: $P(\overline{w}) = P(w_1) P(w_2 | w_1)$
$$P(w_3 | w_1 w_2) \cdots$$

n-grams: look at past $n-1$ words only
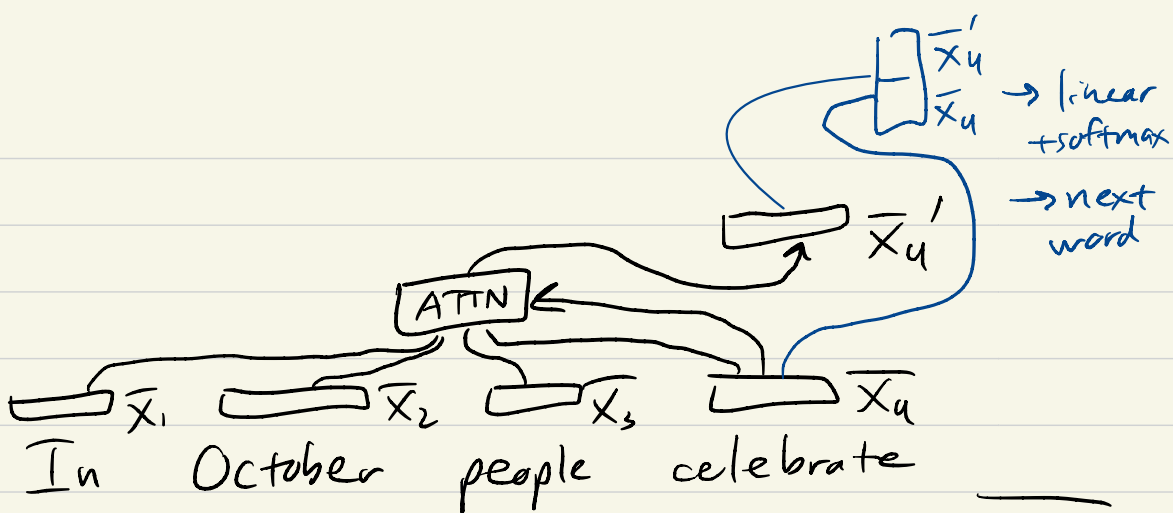
RNNs: look at everything, but they can
forget stuff

In October, people in the US
celebrate _____

Halloween

Predicting the next word requires
looking back a long way, but
sparsely

Alice really likes to go to the movies
with me. She likes horror movies,
I'm good friends with _____.

her
Alice

Self-attention: look back at the sequence
so far to predict the next word

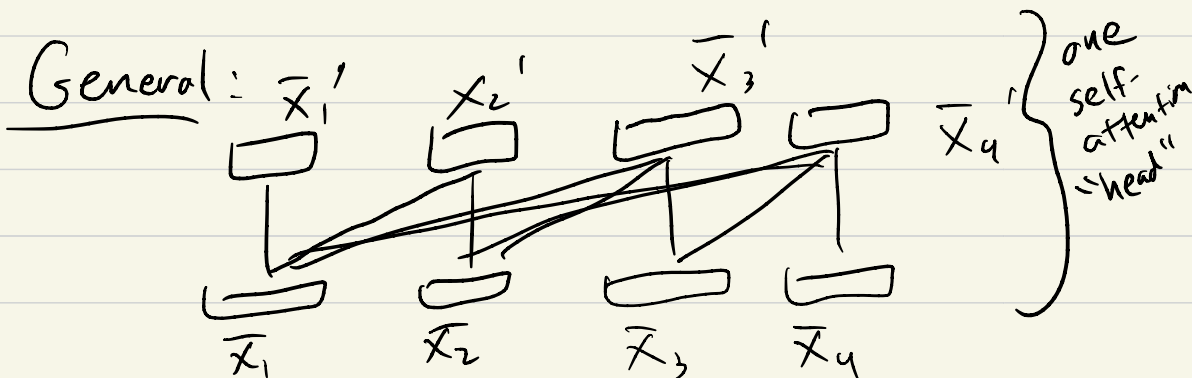$$\alpha_4 = \text{softmax}_i \left( \overline{x}_4^{\top} W \overline{x}_i \right)$$

$\overline{x}_4$ "key"

$\overline{x}_1 \cdots \overline{x}_4$ "values" the attention is over

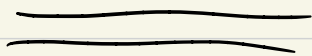$$\overline{x}_4' = \sum_i \overline{\alpha}_{4,i} \overline{x}_i$$

General:

Follows same abstraction as RNN:

Sequence of vectors $x_1, \ldots, x_n$
$\Rightarrow$ new sequence of vectors where
$x_i'$ "knows about" $x_1, \ldots, x_{i-1}$

Advantages: easy access to past words
    parallelizable

Disadvantages: not as powerful as LSTMs
            (so far)

———————

We want to look back at lots of
things in the context

Multi-head self-attention: $K$ "heads"
    which each do an attn computation

Alice likes going ... Movies with me

Combine (average)

$$\bar{x}_u^{(2)}$$
$$\bar{x}_u^{(1)}$$

ATTN

$\bar{x}^u$

$\alpha^{(1)}$
$\alpha^{(2)}$

In October people celebrate

$$\alpha_u^{(k)} = \text{softmax}_i \left( \bar{x}_u^T W^{(k)} \bar{x}_i \right)$$

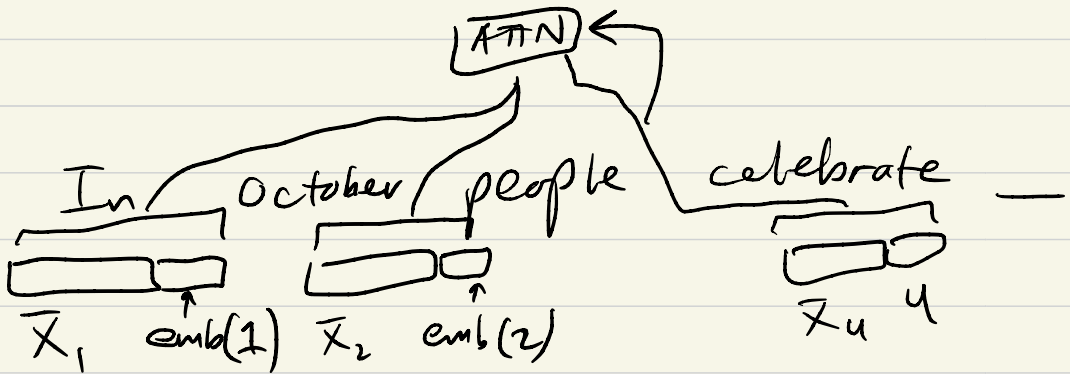$$\bar{x}_u'^{(k)} = \sum_i \alpha_{u,i}^{(k)} V^{(k)} \bar{x}_i$$

↖ new param matrix
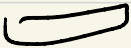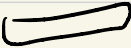
$K = 1 \ldots K$, do independent
copies of the computation

$\left( W^{(k)}, V^{(k)} \right)$ is a head

# Positional encoding

Attention doesn't know the order of the words

Solution: encode position into $\bar{X}_i$



$\bar{X}_1$  emb(1)  $\bar{X}_2$  emb(2)        $\bar{X}_u$  y

1  ▭
2  ▭
3  ▭

50-dim embs, trained with the rest of the model

# Transformer

more
parans

$K = 16$ heads

$16 \times (W, V)$
matrices

Feedforward

Multi-head
self-attn

many layers

word embs → ▭ ▭ ← posn embs