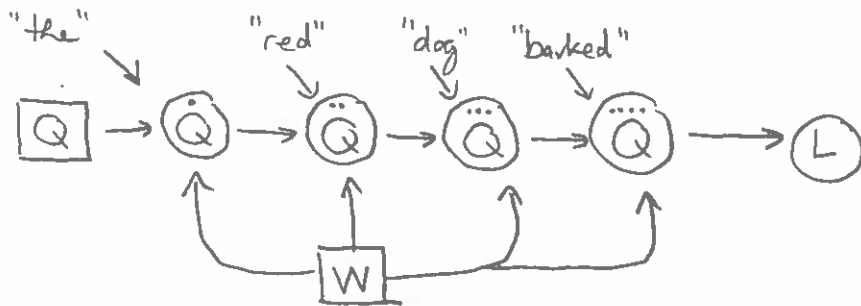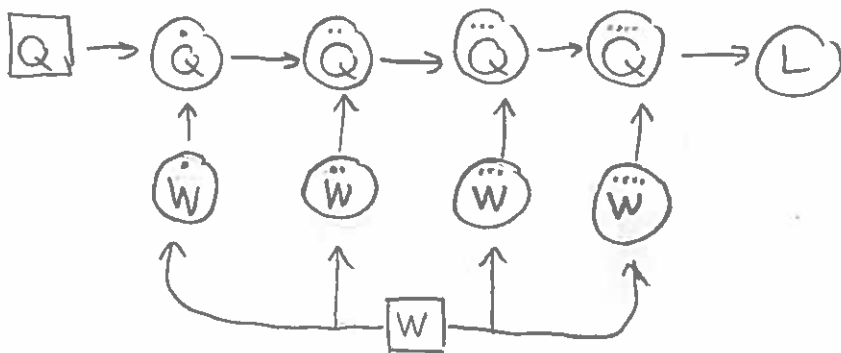# Transformer Networks

① A dominant approach to processing sequences is the RNN:



But one downside to an RNN is what happens when we need to compute partial derivatives for the model weights. Let's first create copies of $W$ and apply the Chain Rule:



So:
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \dot{W}} \frac{\partial \dot{W}}{\partial W} + \frac{\partial L}{\partial \ddot{W}} \frac{\partial \ddot{W}}{\partial W} + \frac{\partial L}{\partial \dddot{W}} \frac{\partial \dddot{W}}{\partial W} + \frac{\partial L}{\partial \ddddot{W}} \frac{\partial \ddddot{W}}{\partial W}$$

$$= \frac{\partial L}{\partial \dot{W}} + \frac{\partial L}{\partial \ddot{W}} + \frac{\partial L}{\partial \dddot{W}} + \frac{\partial L}{\partial \ddddot{W}}$$

$$= \frac{\partial L}{\partial \ddddot{Q}} \frac{\partial \ddddot{Q}}{\partial \dddot{Q}} \frac{\partial \dddot{Q}}{\partial \ddot{Q}} \frac{\partial \ddot{Q}}{\partial \dot{Q}} \frac{\partial \dot{Q}}{\partial \dot{W}} + \frac{\partial L}{\partial \ddddot{Q}} \frac{\partial \ddddot{Q}}{\partial \dddot{Q}} \frac{\partial \dddot{Q}}{\partial \ddot{Q}} \frac{\partial \ddot{Q}}{\partial \ddot{W}}$$

$$+ \frac{\partial L}{\partial \ddddot{Q}} \frac{\partial \ddddot{Q}}{\partial \dddot{Q}} \frac{\partial \dddot{Q}}{\partial \dddot{W}} + \frac{\partial L}{\partial \ddddot{Q}} \frac{\partial \ddddot{Q}}{\partial \ddddot{W}}$$

② Thus, it takes linear time to compute (linear in the length of the sequence).

I mean, that seems ok, since processing a sequence always takes time linear in the sequence, so what's the big deal?

③ Well, the downside isn't that it takes linear time to compute $\frac{\partial L}{\partial W^{(k)}}$ for all k.
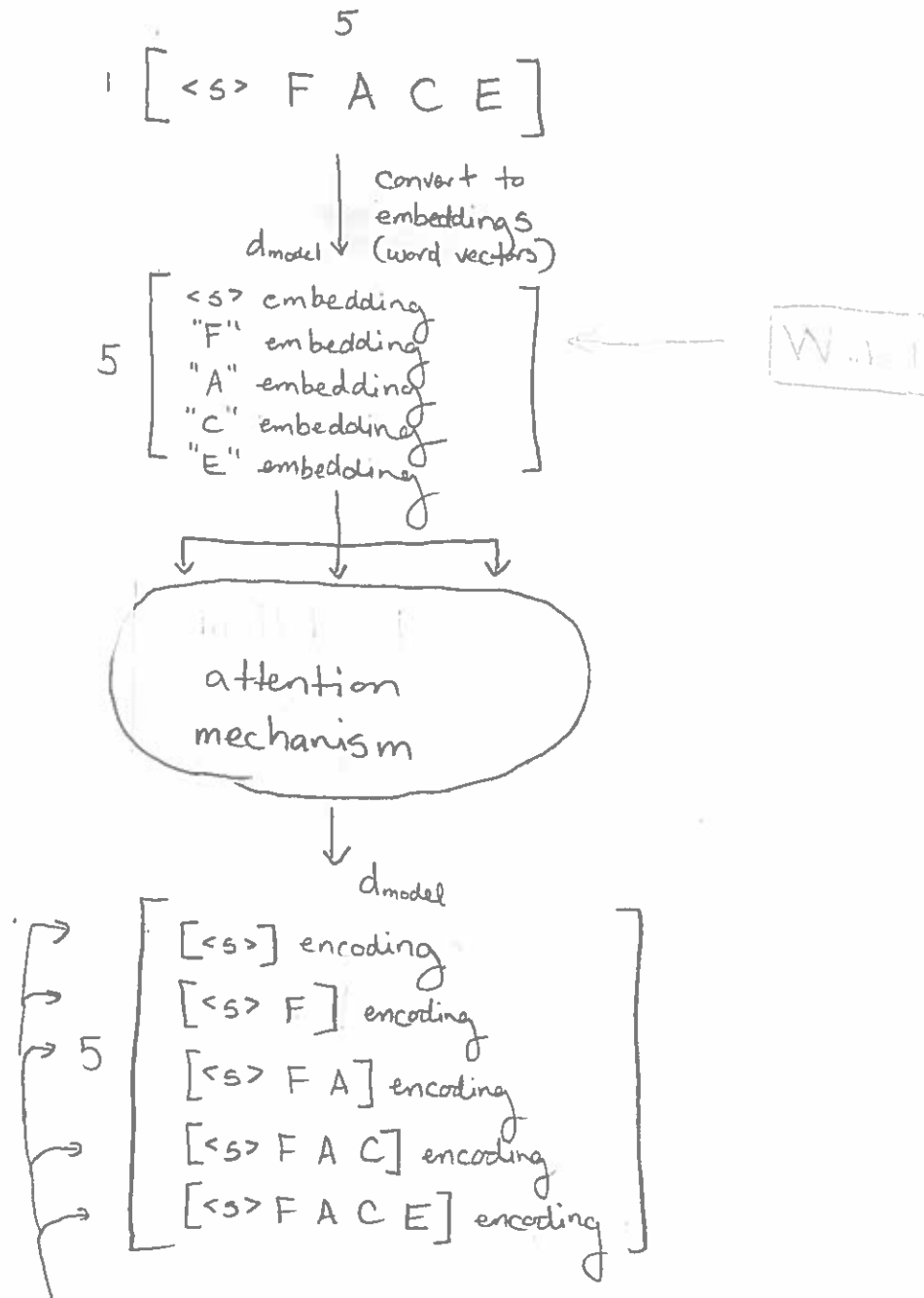
The downside is that it takes linear time to compute <u>just</u> $\frac{\partial L}{\partial W}$. This means that if we want to parallelize the computation, and compute $\frac{\partial L}{\partial W^{(k)}}$ on k processors, it'll still take linear time.

# TRANSFORMER NETWORKS

④ A "transformer" network is a way to encode a variable-length sequence $x$ of length $n$ such that we do not need to sequentially encode each prefix of $x$ in order to encode $x$.

⑤ At a high level, it looks as follows:

$$1 \begin{bmatrix} & & \overset{5}{\phantom{x}} & & \\ \text{<s>} & F & A & C & E \end{bmatrix}$$

$\downarrow$ convert to embeddings (word vectors)

$d_{model} \downarrow$

$$5 \begin{bmatrix} \text{<s> embedding} \\ \text{"F" embedding} \\ \text{"A" embedding} \\ \text{"C" embedding} \\ \text{"E" embedding} \end{bmatrix} \longleftarrow \boxed{W...}$$

↓ ↓ ↓

( attention mechanism )

$\downarrow$ $d_{model}$

$$5 \begin{bmatrix} [\text{<s>}] \text{ encoding} \\ [\text{<s> F}] \text{ encoding} \\ [\text{<s> F A}] \text{ encoding} \\ [\text{<s> F A C}] \text{ encoding} \\ [\text{<s> F A C E}] \text{ encoding} \end{bmatrix}$$
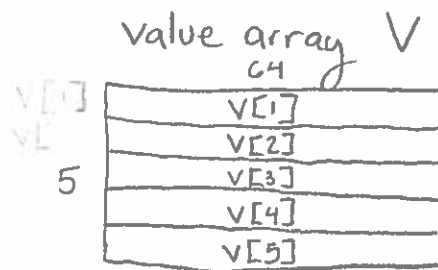
All sequence prefixes are encoded in a single forward pass!

⑥ Well, what does this magical "attention mechanism" do that manages to encode the prefixes in a parallelizable way?

At its core, it tries to emulate a dictionary (or hash table). Imagine we have an array of 5 "values", where each value is a length-64 real vector:

value array V

| | | 
|---|---|
| V[1] | |
| V[2] | |
| V[3] | |
| V[4] | |
| V[5] | |

5

64

This array can be stored as a $5 \times 64$ matrix.

⑦ In addition, suppose that each array position is associated with a "key", which is a length-32 real vector.

keys K

5

32  $k_1$ $k_2$ $k_3$ $k_4$ $k_5$

We can store the 5 keys as a $32 \times 5$ matrix.

⑧ To "look up" a value that corresponds to key $k_j$ (i.e $V[j]$), we can do it with matrix operations. Suppose a query $q = k_4^T$.

(a) multiply $q$ by $K$:

$$
\underset{1}{\boxed{\quad\underset{k_4}{\overset{32}{\phantom{xxx}}}\quad}}
\quad
\underset{32}{\overset{5}{\boxed{k_1\;k_2\;k_3\;k_4\;k_5}}}
\;=\;
\boxed{k_4\cdot k_1 \mid k_4\cdot k_2 \mid k_4\cdot k_3 \mid k_1\, k_4 \mid k_4\cdot k_5}
$$

(b) $k_4 \cdot k_4 = \|k_4\|^2$ is the maximal value of the resulting vector. By applying softmax, we can further magnify the differences to obtain an "attention vector" $A$:

$$
\text{softmax}\left(\boxed{k_4\cdot k_1 \mid k_4\cdot k_2 \mid k_4\cdot k_3 \mid k_4\cdot k_4 \mid k_4\cdot k_5}\right) \approx \boxed{0 \mid 0 \mid 0 \mid 1 \mid 0}
$$

(c) now we can retrieve $V[4]$ by multiplying $A$ by $V$:

$$
\underset{1}{\overset{5}{\boxed{0 \mid 0 \mid 0 \mid 1 \mid 0}}}
\quad
\underset{5}{\overset{64}{\boxed{\begin{array}{c} V[1] \\ V[2] \\ V[3] \\ V[4] \\ V[5] \end{array}}}}
\;=\;
\boxed{\quad V[4] \quad}
$$

⑨ This mechanism also allows us to "hash" any query $q$ to a corresponding array position. Observe:

(a)



$$\boxed{q} \; (32) \qquad \begin{array}{|c|c|c|c|c|} \hline k_1 & k_2 & k_3 & k_4 & k_5 \\ \hline \end{array} \; = \; \begin{array}{|c|c|c|c|c|} \hline q \cdot k_1 & q \cdot k_2 & q \cdot k_3 & q \cdot k_4 & q \cdot k_5 \\ \hline \end{array}$$

(b)

$$\text{softmax}\left( \begin{array}{|c|c|c|c|c|} \hline q \cdot k_1 & q \cdot k_2 & q \cdot k_3 & q \cdot k_4 & q \cdot k_5 \\ \hline \end{array} \right) \approx \begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 0 & 0 \\ \hline \end{array}$$

if $q$ is closer to $k_2$ than the other keys

(c)

$$\begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 0 & 0 \\ \hline \end{array} \; (5) \qquad \begin{array}{|c|} \hline V[1] \\ \hline V[2] \\ \hline V[3] \\ \hline V[4] \\ \hline V[5] \\ \hline \end{array} \; (64) \; = \; \boxed{V[2]}$$

⑩ Note also that some queries will return, rather than a single array element, a linear combination of array elements, e.g. suppose

$$q \cdot k_1 = 2.0 \qquad q \cdot k_4 = 1.8$$
$$q \cdot k_2 = 2.4 \qquad q \cdot k_5 = -12.5$$
$$q \cdot k_3 = -10$$

Then:

$$\text{softmax}\left( \begin{array}{|c|c|c|c|c|} \hline q \cdot k_1 & q \cdot k_2 & q \cdot k_3 & q \cdot k_4 & q \cdot k_5 \\ \hline \end{array} \right) \approx \begin{array}{|c|c|c|c|c|} \hline .30 & .45 & 0 & .25 & 0 \\ \hline \end{array}$$

and

$$\begin{array}{|c|c|c|c|c|} \hline .3 & .45 & 0 & .25 & 0 \\ \hline \end{array} \begin{array}{|c|} \hline V[1] \\ \hline V[2] \\ \hline V[3] \\ \hline V[4] \\ \hline V[5] \\ \hline \end{array} \; = \; \boxed{.3\,V[1] + .45\,V[2] + .25\,V[4]}$$

# TRANSFORMER NETWORKS

11) It's straightforward to extend this mechanism to handle multiple queries at once by stacking query vectors into a query matrix:

(a)

$$
\underset{3}{\begin{array}{c} 32 \\ \begin{array}{|c|} \hline q_1 \\ \hline q_2 \\ \hline q_3 \\ \hline \end{array} \end{array}}
\quad
\underset{32}{\begin{array}{c} 5 \\ \begin{array}{|c|c|c|c|c|} \hline k_1 & k_2 & k_3 & k_4 & k_5 \\ \hline \end{array} \end{array}}
\;=\;
\underset{3}{\begin{array}{c} 5 \\ \begin{array}{|c|c|c|c|c|} \hline q_1 \cdot k_1 & q_1 \cdot k_2 & q_1 \cdot k_3 & q_1 \cdot k_4 & q_1 \cdot k_5 \\ \hline q_2 \cdot k_1 & q_2 \cdot k_2 & q_2 \cdot k_3 & q_2 \cdot k_4 & q_2 \cdot k_5 \\ \hline q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 & q_3 \cdot k_4 & q_3 \cdot k_5 \\ \hline \end{array} \end{array}}
$$

(b)

$$
\begin{bmatrix}
\text{softmax}\left( \begin{array}{|c|c|c|c|c|} \hline q_1 \cdot k_1 & q_1 \cdot k_2 & q_1 \cdot k_3 & q_1 \cdot k_4 & q_1 \cdot k_5 \\ \hline \end{array} \right) \\
\text{softmax}\left( \begin{array}{|c|c|c|c|c|} \hline q_2 \cdot k_1 & q_2 \cdot k_2 & q_2 \cdot k_3 & q_2 \cdot k_4 & q_2 \cdot k_5 \\ \hline \end{array} \right) \\
\text{softmax}\left( \begin{array}{|c|c|c|c|c|} \hline q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 & q_3 \cdot k_4 & q_3 \cdot k_5 \\ \hline \end{array} \right)
\end{bmatrix}
\quad \text{e.g.} \approx \quad
\begin{bmatrix}
\begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}
\end{bmatrix}
$$

(c)

$$
\begin{bmatrix}
\begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}
\end{bmatrix}
\quad
\begin{array}{|c|} \hline V[1] \\ \hline V[2] \\ \hline V[3] \\ \hline V[4] \\ \hline V[5] \\ \hline \end{array}
\;=\;
\begin{bmatrix}
\begin{array}{|c|} \hline V[2] \\ \hline V[1] \\ \hline V[5] \\ \hline \end{array}
\end{bmatrix}
$$

12) Expressed as a causal diagram, the attention mechanism is:



$$A = \text{softmax}(RK, \; \dim = 1)$$

$$O = AV$$

this means we apply softmax to each row

(13) Let's see what happens if we put this attention mechanism into our encoder:

$1 \begin{bmatrix} \text{<s>} & F & A & C & E \end{bmatrix}$

$\downarrow$ 512

| <s> | embedding |
|-----|-----------|
| F | |
| A | |
| C | |
| E | |

5

X

$\downarrow$

E

$W^{(R)} \rightarrow$ R $\quad W^{(K)} \rightarrow$ K $\quad W^{(V)} \rightarrow$ V

A

O

$W \rightarrow$ Q

32

| $r_1$ |
| $r_2$ |
| $r_3$ |
| $r_4$ |
| $r_5$ |

5

32 | $k_1$ | $k_2$ | $k_3$ | $k_4$ | $k_5$ |

64

| $v_1$ |
| $v_2$ |
| $v_3$ |
| $v_4$ |
| $v_5$ |

5

5

A

5

64

| $o_1$ |
| $o_2$ |
| $o_3$ |
| $o_4$ |
| $o_5$ |

5

$\downarrow$ 64

| $q_1$ |
| $q_2$ |
| $q_3$ |
| $q_4$ |
| $q_5$ |

5

There are four "dimensionalities" in the encoder:

$|X|$ : length of the input seq

$d_{model}$ : length of the embedding

$d_{key}$ : length of each "key"

$d_{value}$ : length of each "value"

⑭ In the causal diagram, no matter what the size of the input sequence, the longest path has only 6 nodes: _____

$$\boxed{X} \rightarrow \textcircled{E} \rightarrow \textcircled{R} \rightarrow \textcircled{A} \rightarrow \textcircled{O} \rightarrow \textcircled{Q}$$

Contrast this to the RNN encoder, whose longest path was linear in the length of the input, e.g.

$$\boxed{Q} \rightarrow \overset{①}{\textcircled{Q}} \rightarrow \overset{②}{\textcircled{Q}} \rightarrow \cdots \rightarrow \overset{ⓝ}{\textcircled{Q}}$$

Hence, we can compute any quantity of the transformer encoder in constant time, while computing an RNN state variable requires linear time.

⑮ A generalized version of the transformer's attention mechanism splits the embedding matrix into equal-sized chunks and runs a separate attention mechanism on each chunk; e.g.:

(16) This "multi-head" attention mechanism can be efficiently implemented with tensor operations:

$$X$$

192
$$5 \begin{bmatrix} E \end{bmatrix}$$

$N^{(R)} \rightarrow 5 \begin{bmatrix} & 192 & \\ & R & \\ & & \end{bmatrix}$
$W^{(K)} \rightarrow 192 \begin{bmatrix} & 5 & \\ & K & \\ & & \end{bmatrix}$
$W^{(V)} \rightarrow 5 \begin{bmatrix} & 192 & \\ & & \\ & & \end{bmatrix}$

↓ reshape
↓ reshape
↓ reshape

$5 \begin{bmatrix} \overset{64}{R_1} & \overset{64}{R_2} & \overset{64}{R_3} \end{bmatrix}$
$64 \begin{bmatrix} \overset{5}{K_1} & \overset{5}{K_2} & \overset{5}{K_3} \end{bmatrix}$
$5 \begin{bmatrix} \overset{64}{V_1} & \overset{64}{V_2} & \overset{64}{V_3} \end{bmatrix}$

5
$5 \begin{bmatrix} A_1 \end{bmatrix}$

$5 \begin{bmatrix} \overset{64}{Q_1} & \overset{64}{Q_2} & \overset{64}{Q_3} \end{bmatrix}$

↓ reshape
192
$$5 \begin{bmatrix} Q \end{bmatrix}$$

(17) Next, let's take a closer look at the embedding process:

$$[ <s> \ F \ A \ C \ E ]$$

$$\downarrow$$

$$\begin{bmatrix} <s> \text{ embedding} \\ \text{"F" embedding} \\ \text{"A" embedding} \\ \text{"C" embedding} \\ \text{"E" embedding} \end{bmatrix}$$

(18) One straightforward approach is to just use pretrained word vectors from software like word2vec. However, the transformer network does not encode anything about the position of the words, thus:

$$\begin{bmatrix} <s> \text{ embedding} \\ F \\ A \\ C \\ E \end{bmatrix} \quad \text{is treated just like} \quad \begin{bmatrix} <s> \text{ embedding} \\ C \\ A \\ F \\ E \end{bmatrix}$$

This "bag-of-words" assumption is problematic for applications for which word order matters.



the dog bit the man

the man bit the dog

(19) Looking at the transpose of the embedding matrix:

<s> embedding ↓

"E" embedding ↓

$$\begin{bmatrix} 0.8 & 2.2 & -0.1 & -2.8 & 0.0 \\ -0.4 & 2.1 & -0.5 & -2.4 & 0.0 \\ 1.2 & 0.0 & 3.0 & 0.0 & -1.8 \end{bmatrix}$$

we could try to encode "position" by adding some function of the position to each dimension:

$$\begin{bmatrix} 0.8 & 2.2 & -0.1 & -2.8 & 0.0 \\ +f_1(1) & +f_1(2) & +f_1(3) & +f_1(4) & +f_1(5) \\ & & & & \\ -0.4 & 2.1 & -0.5 & -2.4 & 0.0 \\ +f_2(1) & +f_2(2) & +f_2(3) & +f_2(4) & +f_2(5) \\ & & & & \\ 1.2 & 0.0 & 3.0 & 0.0 & -1.8 \\ +f_3(1) & +f_3(2) & +f_3(3) & +f_3(4) & +f_3(5) \end{bmatrix}$$

(20) A common way to implement this "positional encoding" is to use sinusoids of different frequencies and offsets:



$$\begin{bmatrix} 0.8 & 2.2 & -0.1 & -2.8 & 0.0 \\ +0 & +1 & +0 & -1 & +0 \\ & & & & \\ & & & & \\ & & & & \\ -0.4 & 2.1 & -0.5 & -2.4 & 0.0 \\ +1 & +0 & -1 & 0 & +1 \end{bmatrix}$$

etc.

# TRANSFORMER NETWORKS

21) Adding positional encodings to our overall picture, we have:

(22) Now suppose we want to use a transformer for translation. When "decoding" (producing the output language), we want to take into consideration both the input encodings and the output generated so far:

"input"

$[\text{<s> F A C E}]$

↓

```
┌──────────────┐
│   encoder    │
└──────────────┘
```

↓

```
┌──────────────┐
│      q₁      │
│      q₂      │
│      q₃      │
│      q₄      │
│      q₅      │
└──────────────┘
```

$[\text{<s> F A C E}]$ → encoder → $q_1, q_2, q_3, q_4, q_5$

"output so far"

$[\text{<s> C A}]$

↓

```
┌──────────────┐
│   encoder    │
└──────────────┘
```

↓

```
┌──────────────┐
│     q'₁      │
│     q'₂      │
│     q'₃      │
└──────────────┘
```

$q'_1, q'_2, q'_3$

we would like to use the most recent state $q'_3$ (which tells us where we are in the decoding process) to direct the decoder's attention to the relevant encodings of the input

(23) Again we use the multi-head attention mechanism, but rather than use the input as our request, key, and value array, we use the "output so far" as the request:

"output so far"

$$[<s> \ C \ A]$$

self-attn encoder

| $q_1'$ |
|---|
| $q_2'$ |
| $q_3'$ |

"input"

$$[<s> \ F \ A \ C \ E]$$

self-attn encoder

| $q_1$ |
|---|
| $q_2$ |
| $q_3$ |
| $q_4$ |
| $q_5$ |

R  K  V

multi-head attention

| $q_1''$ |
|---|
| $q_2''$ |
| $q_3''$ |

these states have already been used for prediction, and are sort of trash at this point

← this final row is then used to predict the next letter of the translation (using a linear layer + softmax)

24) During training, we know the entire output in advance, so for efficiency we could do the predictions all at once:

"output"

$[$ <s> C A F E $]$

"input"

$[$ <s> F A C E $]$

self-attn encoder

self-attn encoder

R   K   V

multi-head attention

| $q_1''$ | → use to predict "C" |
| $q_2''$ | → use to predict "A" |
| $q_3''$ | → use to predict "F" |
| $q_4''$ | → use to predict "E" |
| $q_5''$ | → use to predict "</s>" |

25) This requires a small tweak to the output's self-attn encoder, so that earlier output states cannot "see" future output states.

26) This modification is a "mask" matrix that zeroes out information coming from future output (compare w 16):

$$\boxed{X}$$

$\downarrow 192$

$$5 \left[ \quad E \quad \right]$$

$192$       $5$       $192$

$$5 \left[ \quad R \quad \right] \qquad 192 \left[ \quad k \quad \right] \qquad 5 \left[ \quad V \quad \right]$$

$\downarrow$ reshape      $\downarrow$ reshape      $\downarrow$ reshape

$$5 \left[ \; \overset{64}{R_1} \; \overset{64}{R_2} \; \overset{64}{R_3} \; \right] \qquad 64 \left[ \; \overset{5}{K_1} \; \overset{5}{K_2} \; \overset{5}{K_3} \; \right] \qquad 5 \left[ \; \overset{64}{V_1} \; \overset{64}{V_2} \; \overset{64}{V_3} \; \right]$$

mask

$$\begin{array}{cccccc} 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 \end{array}$$

$\downarrow 5$

$$5 \; \square \; \leftarrow$$

$\downarrow 5$

$$5 \; \boxed{A_1}$$

$$5 \left[ \; \overset{64}{Q_1} \quad \overset{64}{Q_2} \quad \overset{64}{Q_3} \; \right]$$

$\downarrow$ reshape
$\quad 192$

$$5 \left[ \qquad Q \qquad \right]$$