$$\text{softmax}\left(\left[x_1, x_2, x_3, x_4, x_5\right]\right)$$

$$= \left[\frac{e^{x_1 - \max}}{S}, \frac{e^{x_2 - \max}}{S}, \frac{e^{x_3 - \max}}{S}, \frac{e^{x_4 - \max}}{S}, \frac{e^{x_5 - \max}}{S}\right]$$

$$\text{where}: S = \left(e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4} + e^{x_5}\right) e^{-\max}$$

for numerical stability, add $-\max(x_1, x_2, x_3, x_4, x_5)$

---

Online softmax: compute softmax incrementally

$$\text{softmax}\left(\left[x_1, x_2, x_3\right]\right) \longrightarrow \text{softmax}\left(\left[x_1, x_2, x_3, \boxed{x_4, x_5}\right]\right)$$

new!

Goal: We want to compute softmax block by block

$$\text{softmax}\left(\left[x_1, x_2, x_3\right]\right)$$

$$= \left[\frac{e^{x_1 - m_1}}{S_1}, \frac{e^{x_2 - m_1}}{S_1}, \frac{e^{x_3 - m_1}}{S_1}\right]$$

Where:
$$m_1 = \max(x_1, x_2, x_3)$$
$$S_1 = e^{x_1 - m_1} + e^{x_2 - m_1} + e^{x_3 - m_1}$$

Then: How to compute $\text{softmax}([x_1, x_2, x_3, x_4, x_5])$
by reusing $\text{softmax}([x_1 \ x_2 \ x_3])$ ?

Observe:

The end goal is:

old        new
$$m_2 = \max(\underbrace{x_1 \ x_2 \ x_3}_{} \ \underbrace{x_4 \ x_5}_{})$$
$$S_2 = e^{x_1 - m_2} + e^{x_2 - m_2} + e^{x_3 - m_2} + e^{x_4 - m_2} + e^{x_5 - m_2}$$

$$\left[ \frac{e^{x_1 - m_2}}{S_2} \quad \frac{e^{x_2 - m_2}}{S_2} \quad \frac{e^{x_3 - m_2}}{S_2} \quad \frac{e^{x_4 - m_2}}{S_2} \quad \frac{e^{x_5 - m_2}}{S_2} \right]$$

and we already have

$$\left[ \frac{e^{x_1 - m_1}}{S_1} , \frac{e^{x_2 - m_1}}{S_1} , \frac{e^{x_3 - m_1}}{S_1} \right]$$

For the first 3 elements
we can rescale them by $\quad \cdot S_1 \cdot S_2^{-1} \cdot e^{m_1} \cdot e^{-m_2}$

$$= \frac{S_1}{S_2} \cdot \frac{m_1}{m_2}$$

this scaling factor can make them

$$\left[\frac{e^{x_1-m_2}}{S_2} \quad \frac{e^{x_2-m_2}}{S_2} \quad \frac{e^{x_3-m_2}}{S_2}\right]$$

---

In the original paper,
it stores softmax result into 3 components

$$\text{softmax}([x_1, x_2)] = \begin{cases} ① & [e^{x_1-m_1} \quad e^{x_2-m_1}] \\ ② & m_1 = \max(x_1, x_2) \\ ③ & S_1 = e^{x_1-m_1} + e^{x_2-m_1} \end{cases}$$

the benefit of this

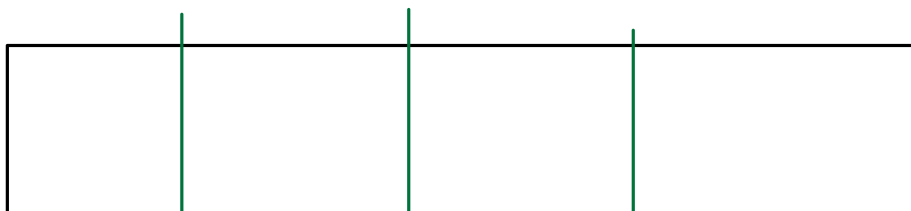$$\text{softmax}([x_1, x_2]) \quad , \quad \text{softmax}([x_3 \, x_4 \, x_5])$$

can express softmax$([x_1 \, x_2 \, x_3 \, x_4 \, x_5])$

$$\text{by} \begin{cases} ① & \frac{e^{m_1}}{e^{m'}}[e^{x_1-m_1} \quad e^{x_2-m_1}] : \frac{e^{m_2}}{e^{m'}}[e^{x_3-m_2} \quad e^{x_4-m_2} \quad e^{x_5-m_2}] \\ ② & m' = \max(m_1, m_2) \\ ③ & S' = \frac{e^{m_1}}{e^{m'}} S_1 + \frac{e^{m_2}}{e^{m'}} S_2 \end{cases}$$

Where $m_1 = \max(x_1, x_2)$  $m_2 = \max(x_3, x_4, x_5)$
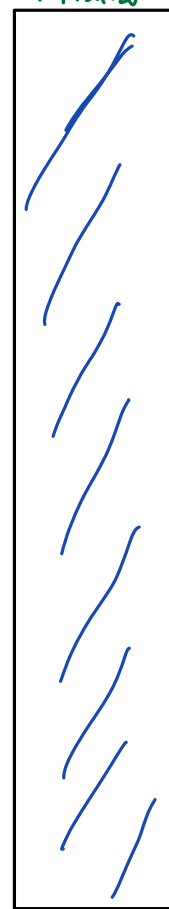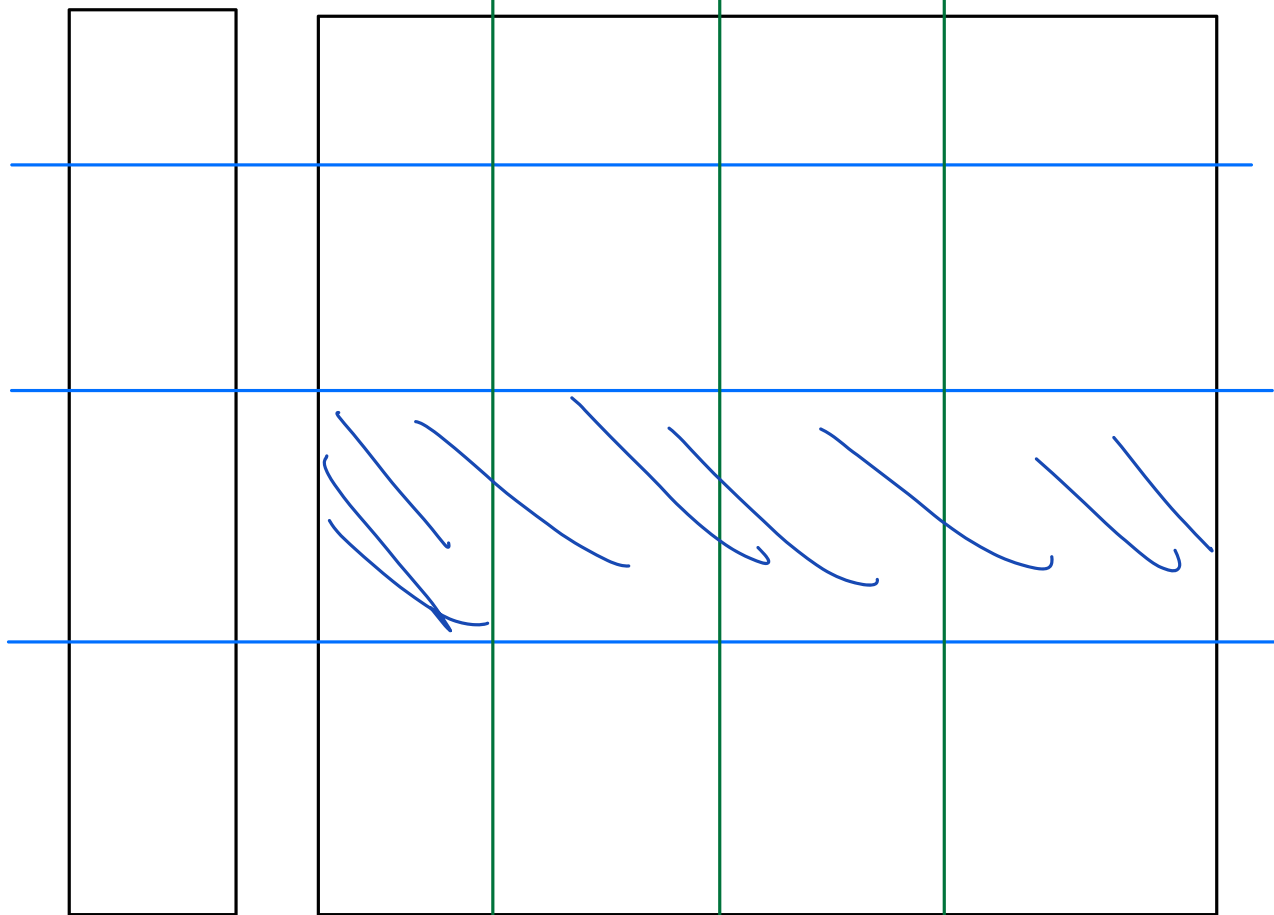
$S_1 = e^{x_1 - m_1} + e^{x_2 - m_1}$  $S_2 = e^{x_3 - m_2} + e^{x_4 - m_2} + e^{x_5 - m_2}$
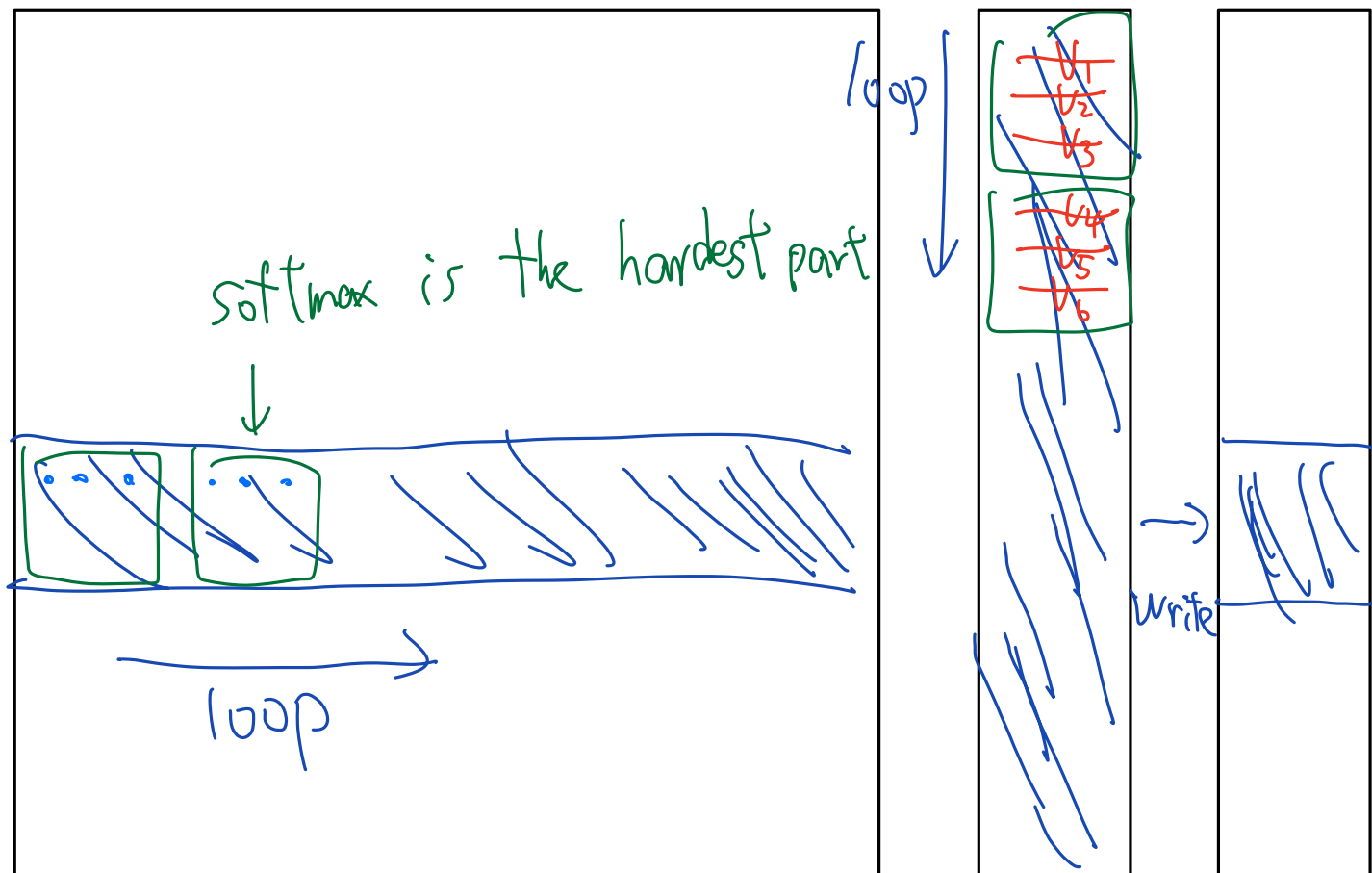
$K^T: d \times N$

$Q: N \times d$

$QK^T \in N \times N$ ($\leftarrow$ don't want to materialize)

$V: N \times d$

softmax$(QK^T)$ 

$V$



softmax is the hardest part

loop

loop

write

Can we express softmax$([x_1 x_2 x_3 x_4 x_5 x_6])$

$[v_1 v_2 v_3 v_4 v_5 v_6]$

softmax$([x_1 x_2 x_3]) \cdot [v_1 v_2 v_3]$

by

softmax$([x_4 x_5 x_6]) \cdot [v_4 v_5 v_6]$  ?

Recap:

softmax$([\ x_1, x_2\ x_3\ x_4\ x_5\ x_6\ ])$

$$\begin{cases} \textcircled{1} & \frac{e^{m_1}}{e^{m'}}[e^{x_1-m_1}\ e^{x_2-m_1}\ e^{x_3-m_1}] : \frac{e^{m_2}}{e^{m'}}[e^{x_4-m_2}\ e^{x_5-m_2}\ e^{x_6-m_2}] \\ \textcircled{2} & m' = \max(m_1, m_2) \\ \textcircled{3} & S' = \frac{e^{m_1}}{e^{m'}}S_1 + \frac{e^{m_2}}{e^{m'}}S_2 \end{cases}$$

Computed independently

So softmax$([x_1\ x_2\ x_3\ x_4\ x_5\ x_6]) \cdot [v_1\ v_2\ v_3\ v_4\ v_5\ v_6]$

will be represented as:

$$\begin{cases} \textcircled{1} & \frac{e^{m_1}}{e^{m'}}\left[v_1 e^{x_1-m_1} + v_2 e^{x_2-m_1} + v_3 e^{x_3-m_1}\right) \\ & + \frac{e^{m_2}}{e^{m'}}\left[v_4 e^{x_4-m_2} + v_5 e^{x_5-m_2} + v_6 e^{x_6-m_2}\right) \\ \textcircled{2} & m' = \max(m_1, m_2) \\ \textcircled{3} & S' = \frac{e^{m_1}}{e^{m'}}S_1 + \frac{e^{m_2}}{e^{m'}}S_2 \end{cases}$$

Sanity check: $\frac{\textcircled{1}}{\textcircled{3}}$ will be the actual softmax result