# Transfer learning

pretrained models

not in
Raw

use + specific data
for finetuning
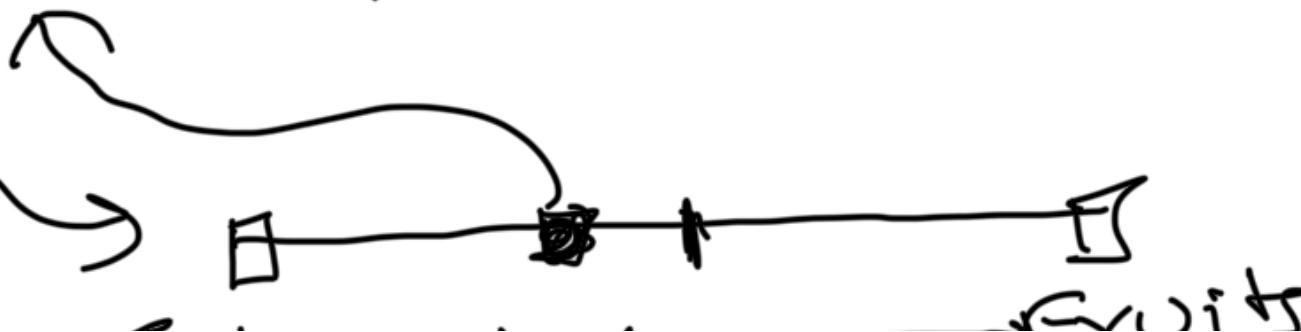
## ⊛ Word Embeddings

Love apple phones.

$$Apple = 0.1 \times E_{love} + 0.5 \times E_{Apple} + 0.4 \times E_{phones}$$

Word Emb.
Space

phony d. Apple.

→ juice
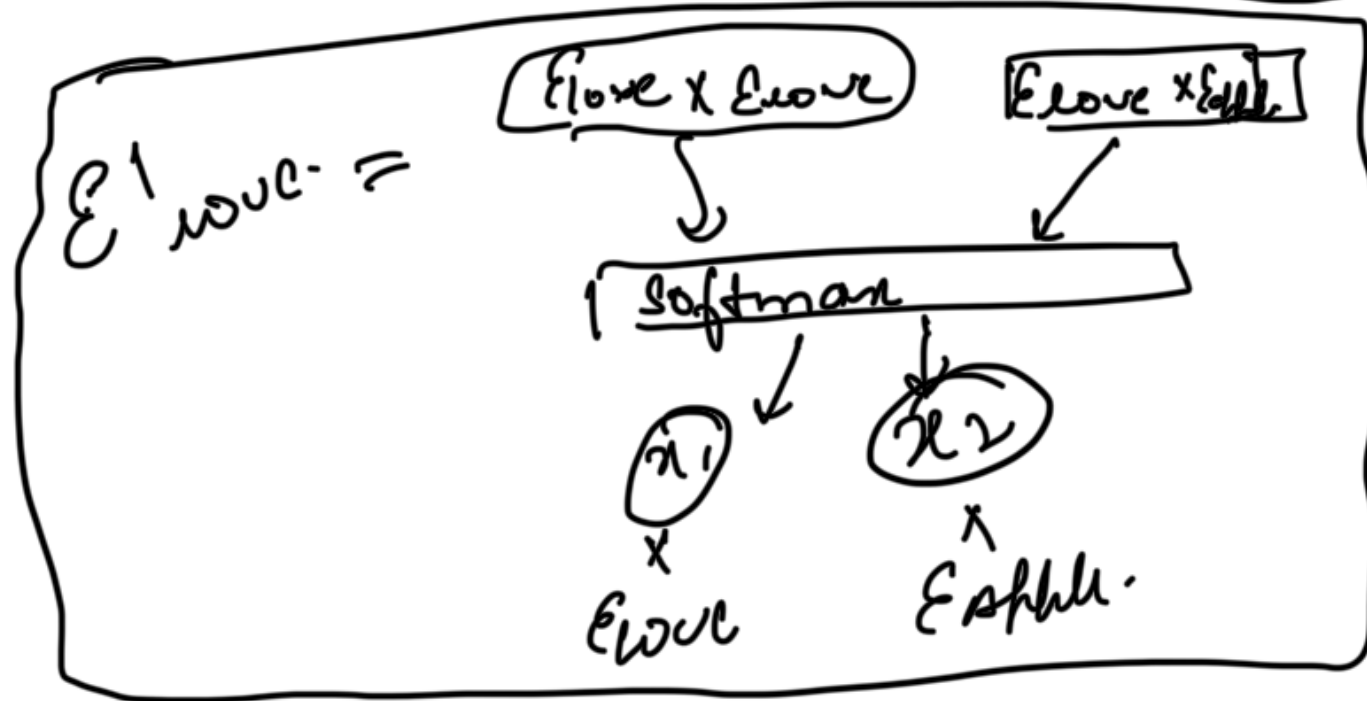
Tech $\leftarrow x \rightarrow$ $\leftarrow y \rightarrow$

$y > x \rightarrow$ more similar to tech.

$0.1 \rightarrow$ rel$^n$ b/w apple & love

$\overset{\curvearrowright}{E_{Apple}} \cdot \overset{\curvearrowleft}{E_{love}^T} \rightarrow$ # passed to Softmax $f \cong$

$\downarrow$

to get the prob.

$E'_{love} =$

$\boxed{E_{love} \times E_{love}}$  $\boxed{E_{love} \times E_{apple}}$

Softmax

$\overset{x}{\boxed{x_1}}$  $\overset{\wedge}{\boxed{x_2}}$

$\overset{x}{E_{love}}$  $\overset{\wedge}{E_{apple}}$

Similarity $\dfrac{0.1}{0.1}$

Problems

No learnable Params.

Use of Same vector.

To solve

$E_{apple} \cdot \boxed{w q}$

$E_{love} \boxed{w_k^T}$

$E_{apple}$

$W_K$ $W_{as}$ $W_v$.

$K_{apple}$ $Q_{apple}$ $V_{ehicle}$

Software

$x_1$

$E_{core}^K$ ...

Eg →

$Q_{Apple}$ $Q_{Apple}$ $Q_{Apple}$

$V_{core}$ $K_{apple}$ $K_{phone}$

Singular O/P

Softmax

$x_1$ $x_2$ $x_3$

$E_{apple}' =$ * * *

$V_{apple}$ $V_{core}$ $V_{phone}$.

$$E_{Apple} = x_1 \cdot v_{apple} + x_2 \cdot v_{love} + x_3 \cdot v_{phone}$$

no dependence of $E$ of love & phones.

So, they eliminate sequential processing & allows parallization.

Long range dependencies

Softmax $(Q \cdot K^T) \cdot V$ → done to add learnable params in the matrices

(*) Linear Transformation → changes dimension → extracting such factors

.Juice

. ice          .O beer
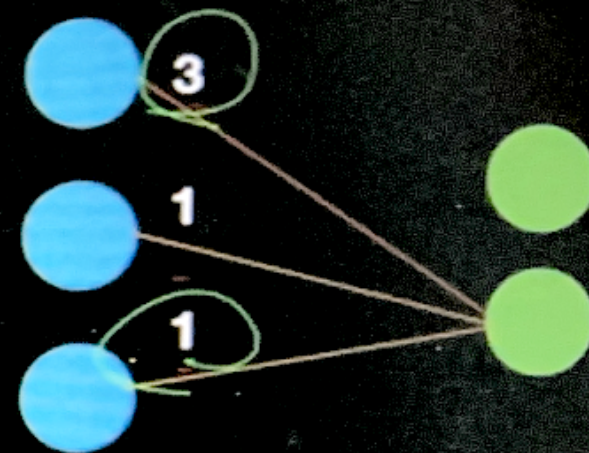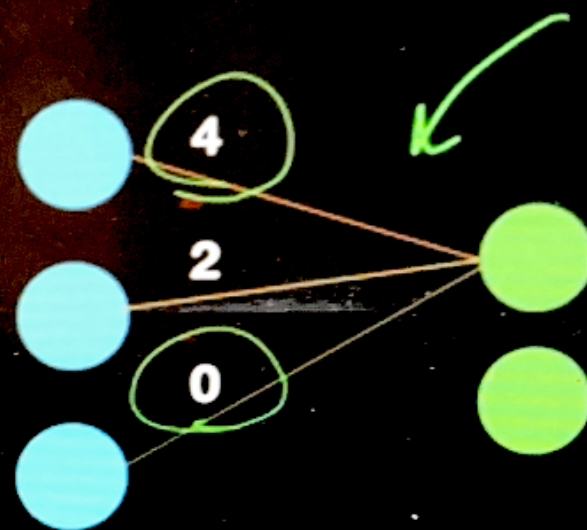
-Apple                                    . Watch

very abart          time

• Tech

in space so
not much relation can be captured.

$$
\begin{cases}
\underleftarrow{\text{Change}} \ \overrightarrow{\text{dimension}} \\
y = w \cdot x \quad \sim g \ just \ like \\
\qquad\qquad\qquad\qquad ANN.
\end{cases}
$$

|        | Apple | Orange | Watch |
|--------|-------|--------|-------|
| isTech | 2     | 0      | 1     |
| isWatch| 0     | 0      | 3     |
| isFruit| 2     | 3      | 0     |



$$[2, 0, 2], \begin{bmatrix} 4 & 3 \\ 2 & 1 \\ 0 & 1 \end{bmatrix} = \boxed{8 \quad 8}$$

↑
Apple

$2*4 + 0*2 + 2*0$

$2*3 + 0*1 + 2*1$

$$[0, 0, 3], \begin{bmatrix} 4 & 3 \\ 2 & 1 \\ 0 & 1 \end{bmatrix}$$

$$\boxed{[0, 3]}$$

↑
Orange

$$[1, 3, 0], \begin{bmatrix} 4 & 3 \\ 2 & 1 \\ 0 & 1 \end{bmatrix}$$

$$[10, 6]$$

↑
Watch

(*) Query, key, value
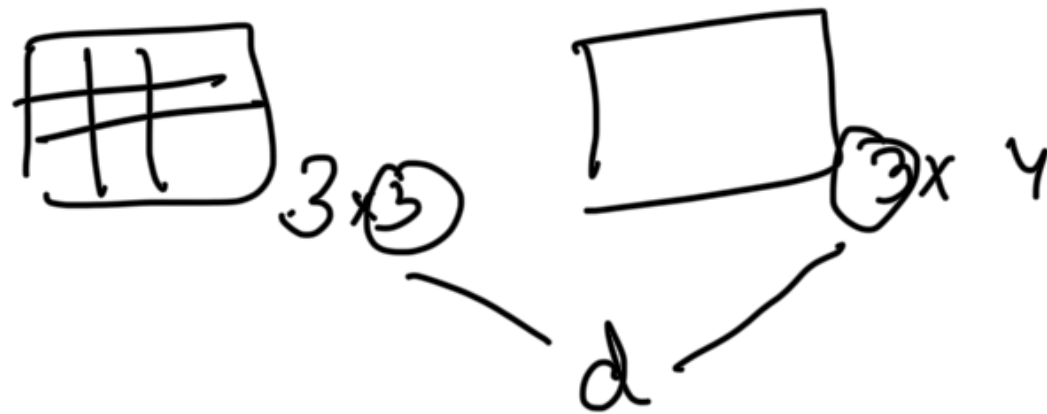Q, K, V.

Sound of [Dog] = Bau
↓          ↓        ↓
query    key     value.

Query. ← Qapple

[ ]
K love → key

↓
[ ] → (V eou E) → value.

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{}}\right)V$$

⊛ why $\sqrt{d_k}$ → dim. of $\mathcal{Q}, K, \mathcal{V}$ - hyperparameter.

various $(A \cdot B) \propto n$, $n$ = common dim. of A & B.



3×3

3×Y

$d$

High var → Vanishing grad
$\downarrow$
$\downarrow$
Hinders training
$\downarrow$
so we use $\sqrt{d_k}$

$C = A \cdot B$

$$Var(c) = n \cdot Var(A) \cdot Var(B)$$

$\downarrow$
→ zero mean

→ Indep random var.

→ Entries of A & B have identical values.

$$X = Q \cdot K^T$$

To reduce ~~var~~
we use $\sqrt{d_k}$ $\rightarrow$ keep
var$(\cancel{X})$ same.

$var(x) \propto d_k$

$X \Rightarrow var(x)$

$a \cdot x \rightarrow a^2 \, var(x)$

Ⓧ **Multi_Head_Attention**

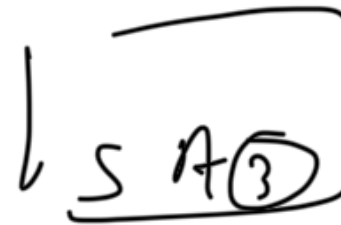As complexity of model ⓝ It is tough to capture the meaning of different words & other things
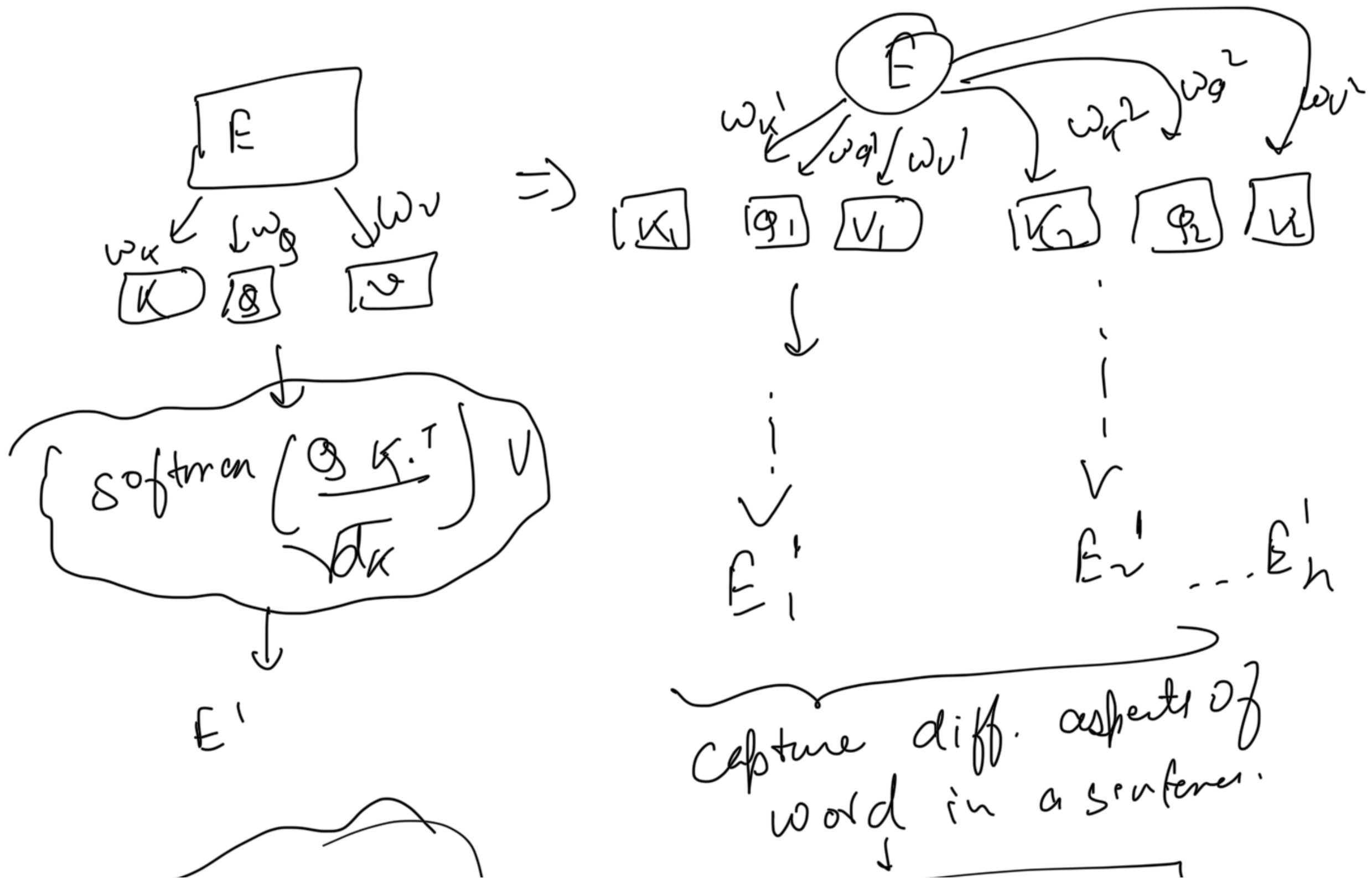
Spatial Rel$^{s}$

Sub-verb
Obj

Time $\rightsquigarrow$

[self-att.①]

[self-att②]

$\downarrow$ S A③

to provide diverse content
just like cnn

E

$W_k \downarrow \quad \downarrow W_q \quad \downarrow W_v$

K    Q    V

$$\text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V$$

E'

$\Rightarrow$

E

$W_k^1 \swarrow \swarrow W_q^1 / W_v^1 \qquad W_k^2 \searrow W_q^2 \searrow W_v^2$

$K_1$   $Q_1$   $V_1$      $K_2$   $Q_2$   $V_2$

$\downarrow$                          $\downarrow$

$\vdots$                              $\vdots$

$\downarrow$                          $\downarrow$

$E_1'$                          $E_2' \; \ldots E_n'$

capture diff. aspects of
word in a sentence.

In GPT₂ 96 Attention Heads

96 diff perspective

Keeps only the relevant features

$$\boxed{E_1' \mid E_2' \mid \cdots \cdots \mid E_n'}$$

$\downarrow * \textcircled{w} \rightarrow$ learn parameter

$\downarrow$

(linear transformation)

$\textcircled{*}$ Positional Encoding

$\rightarrow$ due to lack of sequential understanding.

If Individually $\sin/\cos$ is used then $\rightarrow$ due to periodicity 2 words

may be identical
↓
use $\sin(x) \cos(x)$
together

In og transformer

512 dim → 264 pair s'n cosin

| $\sin(u)$ | $\cos(u)$ | $\sin\left(\frac{u}{2}\right)$ | $\cos\left(\frac{u}{2}\right)$ | — — — — — |

reduced freq

$$PE_{(pos, 2i)} = \sin(pos) \, 10000^{\frac{2i}{d_{model}}}$$

$$PE_{pos}(2i+1) = \cos\left(pos / 10000^{\frac{2i}{d_{model}}}\right)$$

If we know

$PE_{POS} \rightarrow PE_{POS+(k)} \rightarrow$ offset

con be calculated

as there is a Transformation

matrix $T_k$.

$$PE_{POS+k} = T_k \times PE_{POS}$$
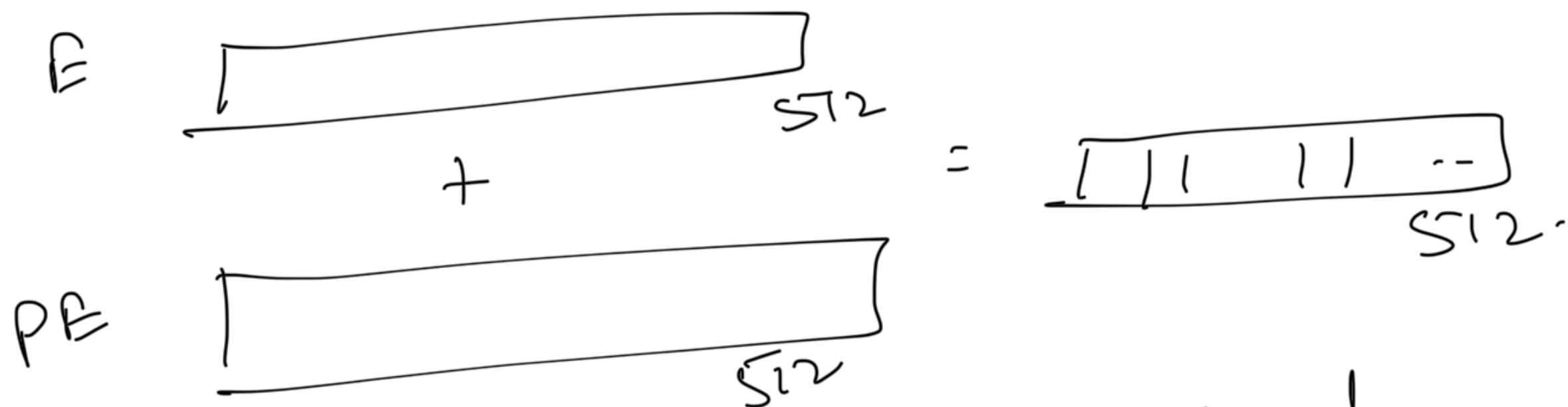


word
Embedding
$(\mathcal{E})$

Positional Encoder
(PE)

out $\rightarrow (E + PE) * \textcircled{W}$  $W_{K, V, Q}$

So, dim of $W$
increases

$\textcircled{F}$ additional

Hence we use ① ... operation

E $\boxed{\phantom{xxxxxxxxxxxxxxxxxxx}}$
                512

$+$                                  $=$  $\boxed{[ \ |\ |\ \ \ \ |\ |\ \ \ - -}$
                                                           512.

PE $\boxed{\phantom{xxxxxxxxxxxxxxxxxxx}}$
                512

①→ reduces   computational
           overhead