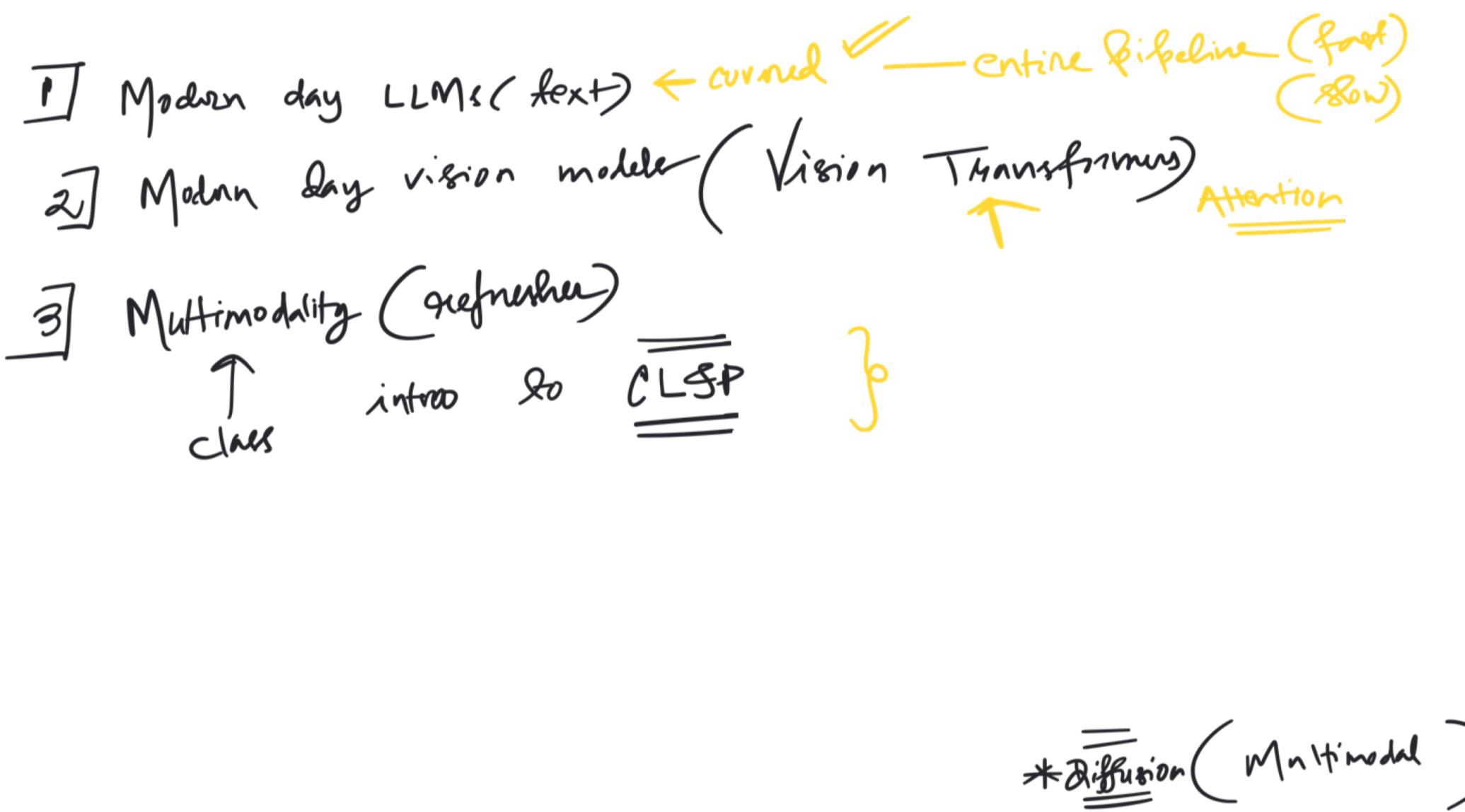


Refreshers Modern Text & Image AI Model Architectures

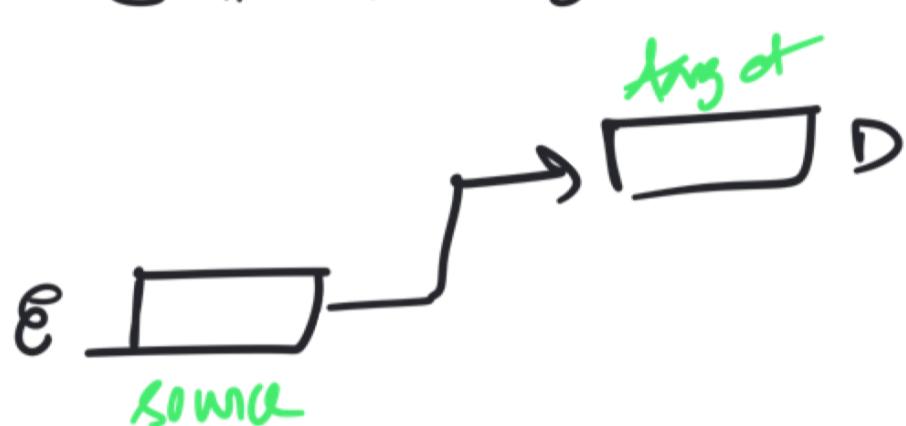


1] Encoder only (BERT)

Chatting

Sentiment Analysis

2] Encoder- Decoder [T5]



translation - a data diff
source
target

3] Decode only GPT-2

No need of separate encoder

Training Dataset

Use a
special token

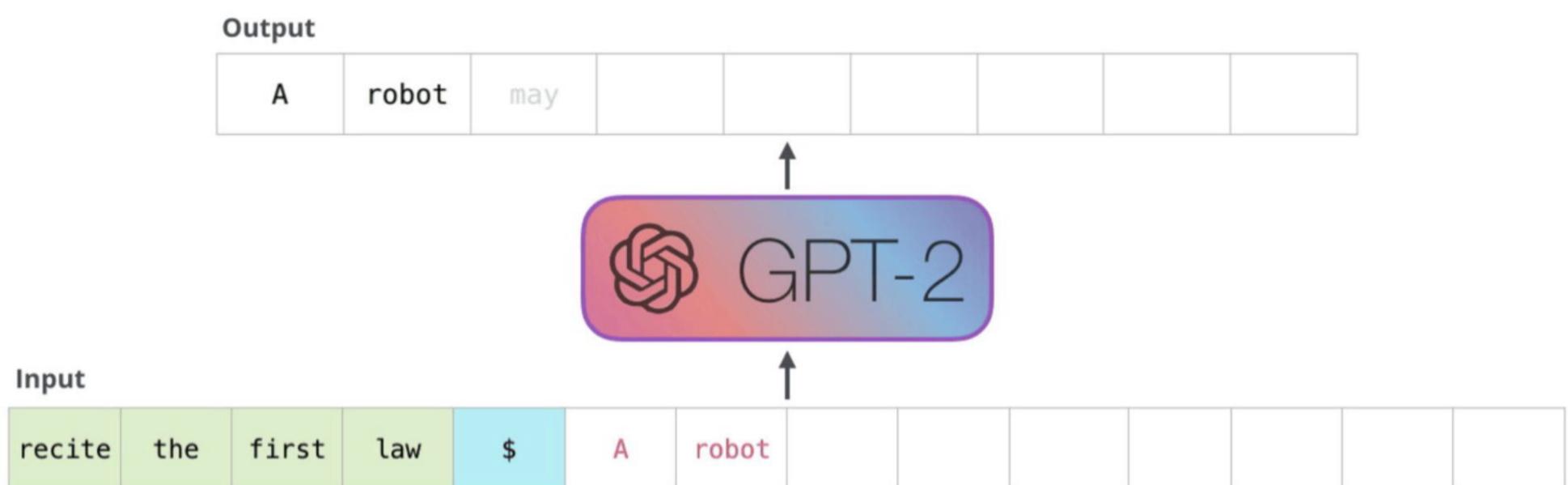
I	am	a	student	<to-fr>	je	suis	étudiant
let	them	eat	cake	<to-fr>	Qu'ils	mangent	de
good	morning	<to-fr>	Bonjour				

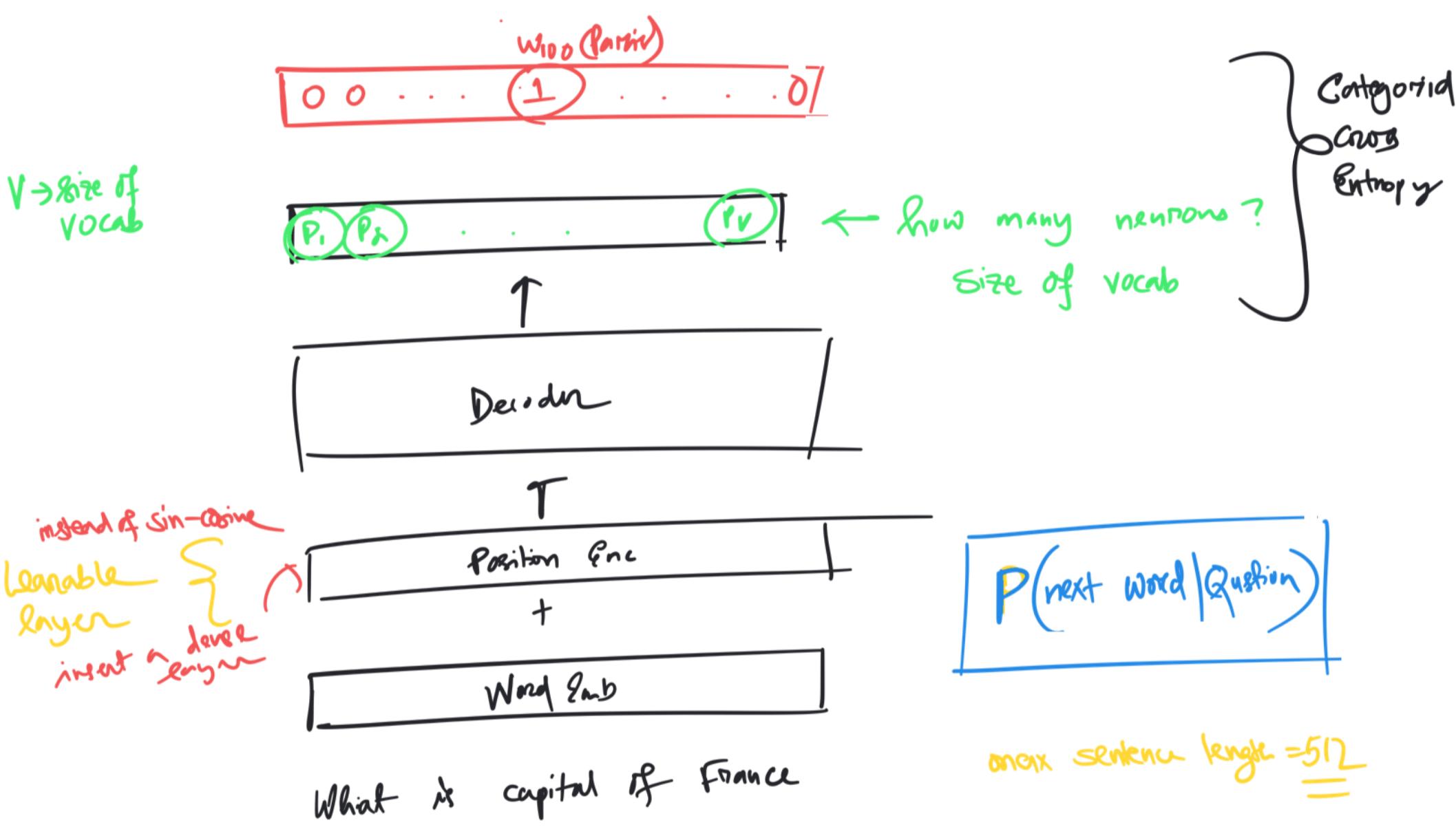
Training Dataset

*Special tokens
for summarization*

Article #1 tokens	<summarize>	Article #1 Summary
Article #2 tokens	<summarize>	Article #2 Summary padding
Article #3 tokens	<summarize>	Article #3 Summary

Next token prediction models



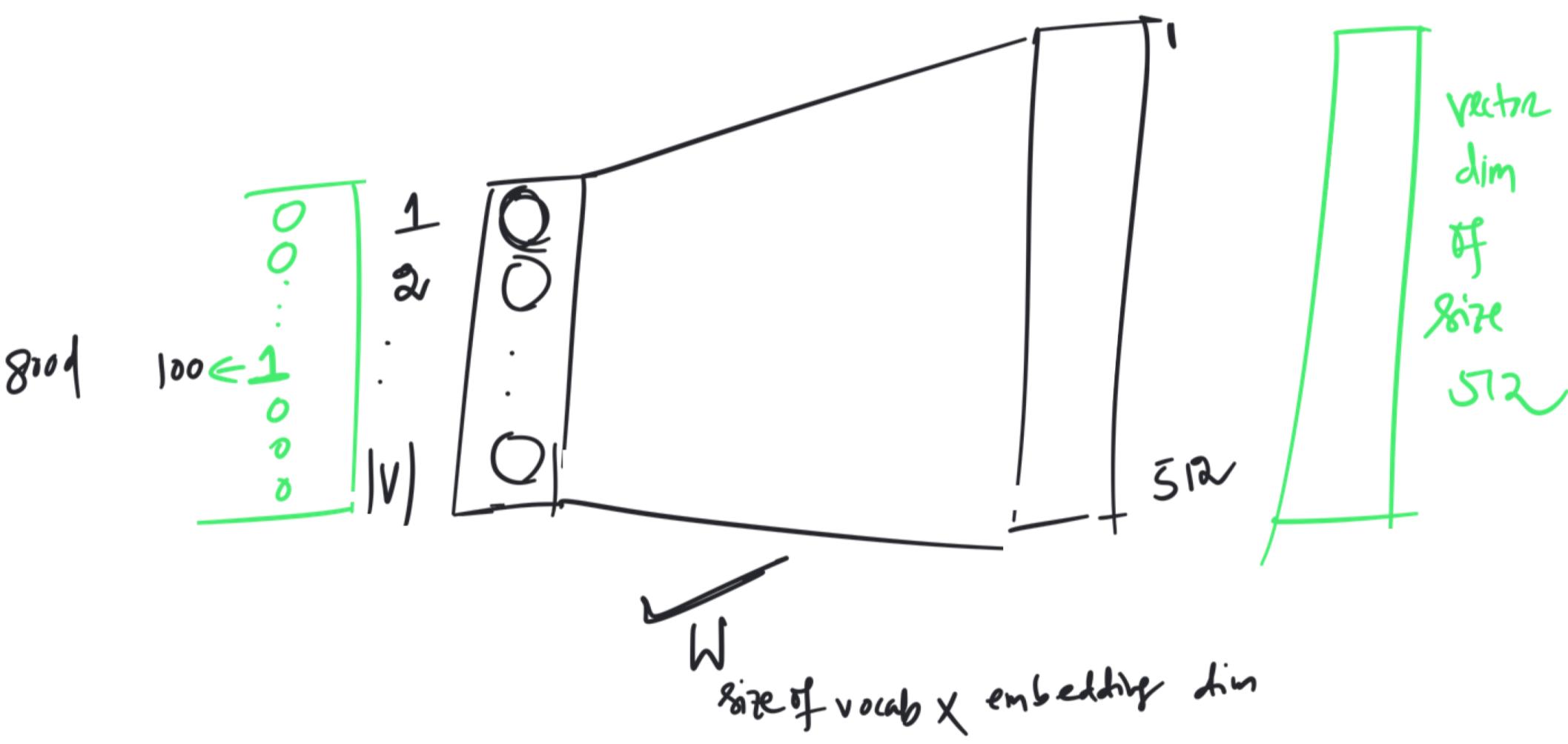


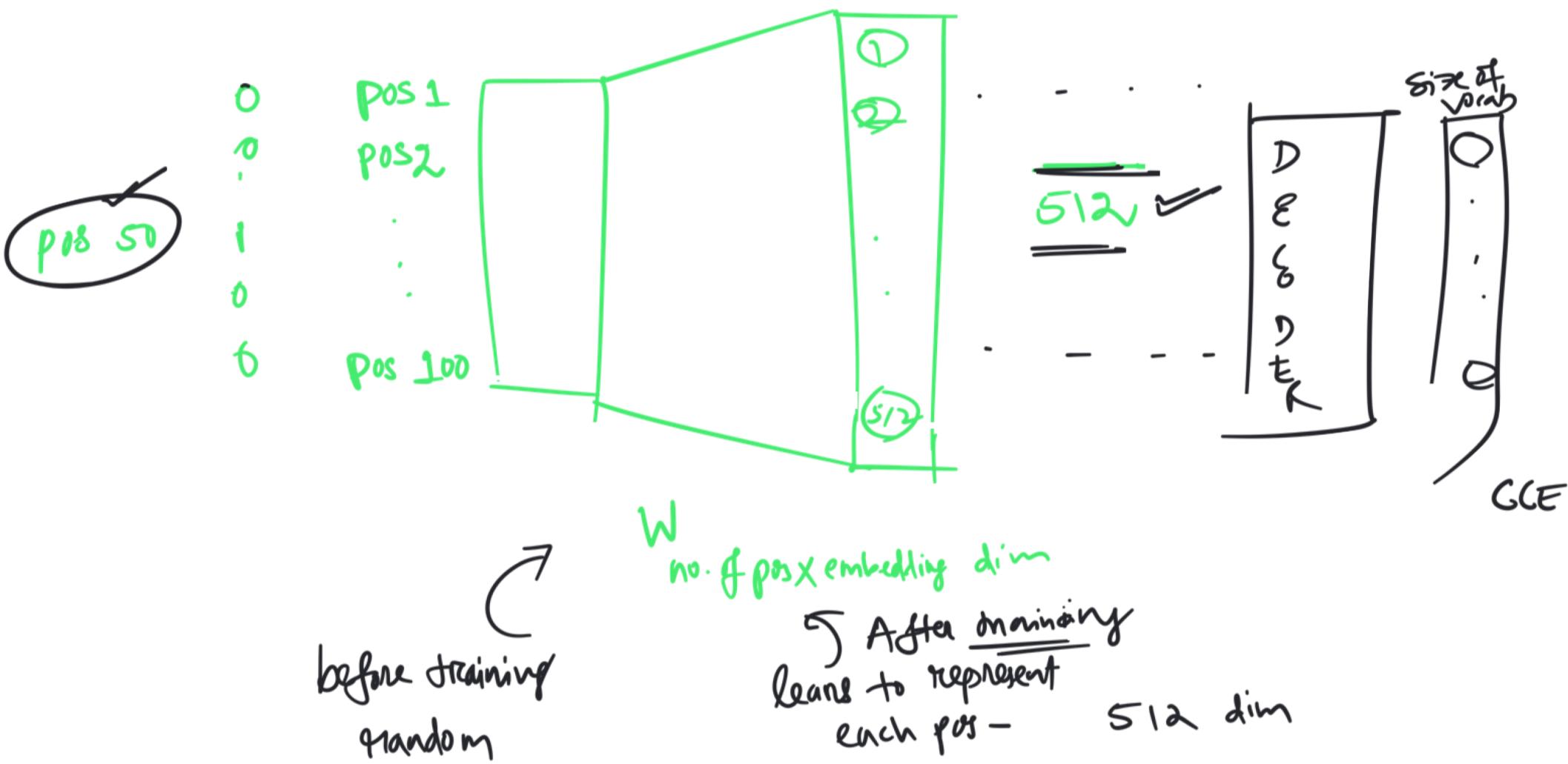


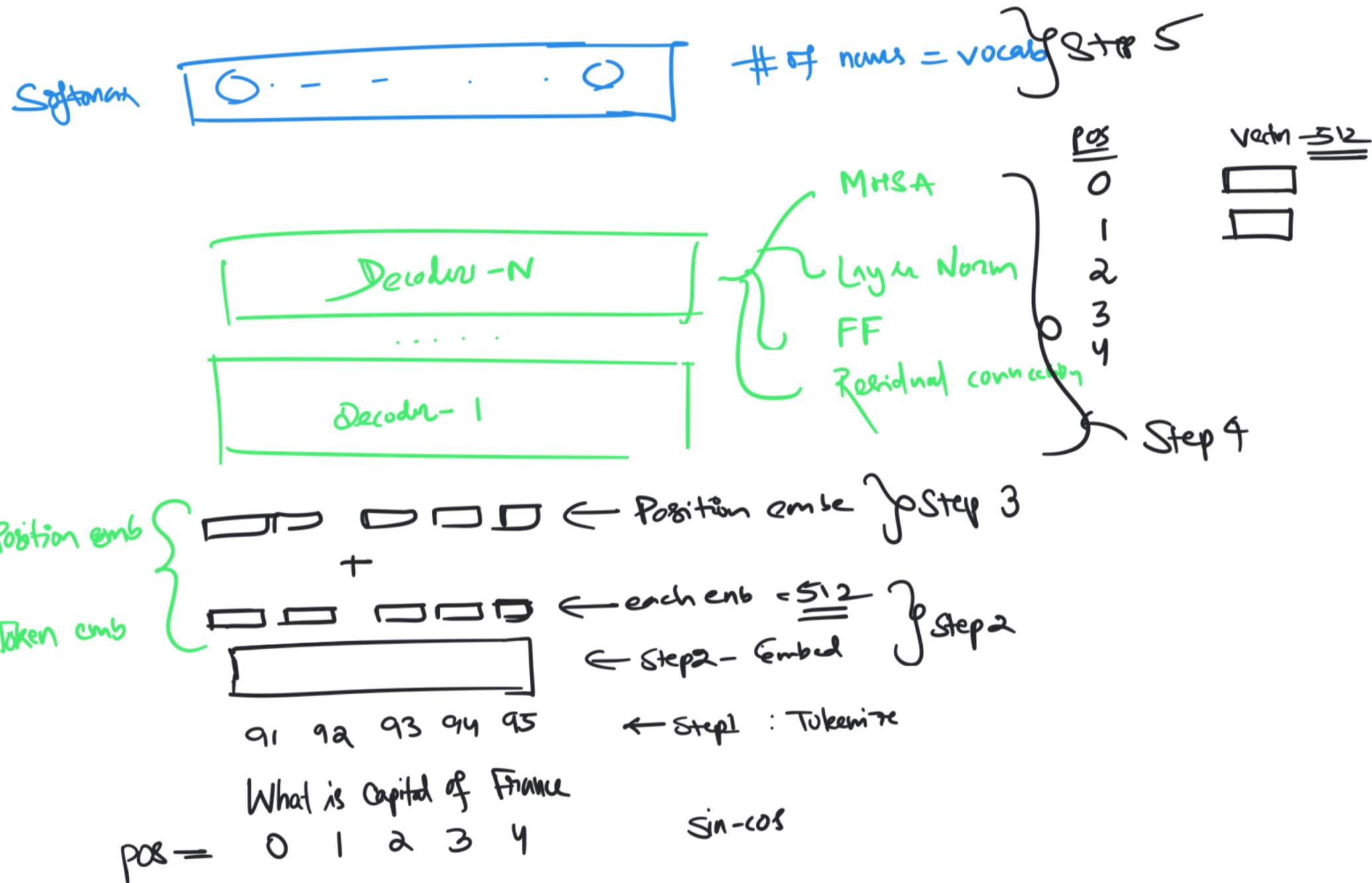
causal mask

Matrix -mult \rightarrow Batching \rightarrow Training faster

$\langle \text{sos} \rangle$	g	am	good	ginger	
$\langle \text{sos} \rangle$	0	0	0	0	$\begin{matrix} 1 \\ g \end{matrix}$
$\langle \text{sos} \rangle$ g		0	0	0	am
$\langle \text{sos} \rangle$ f am			0	0	good







Base GPT

Base Klara

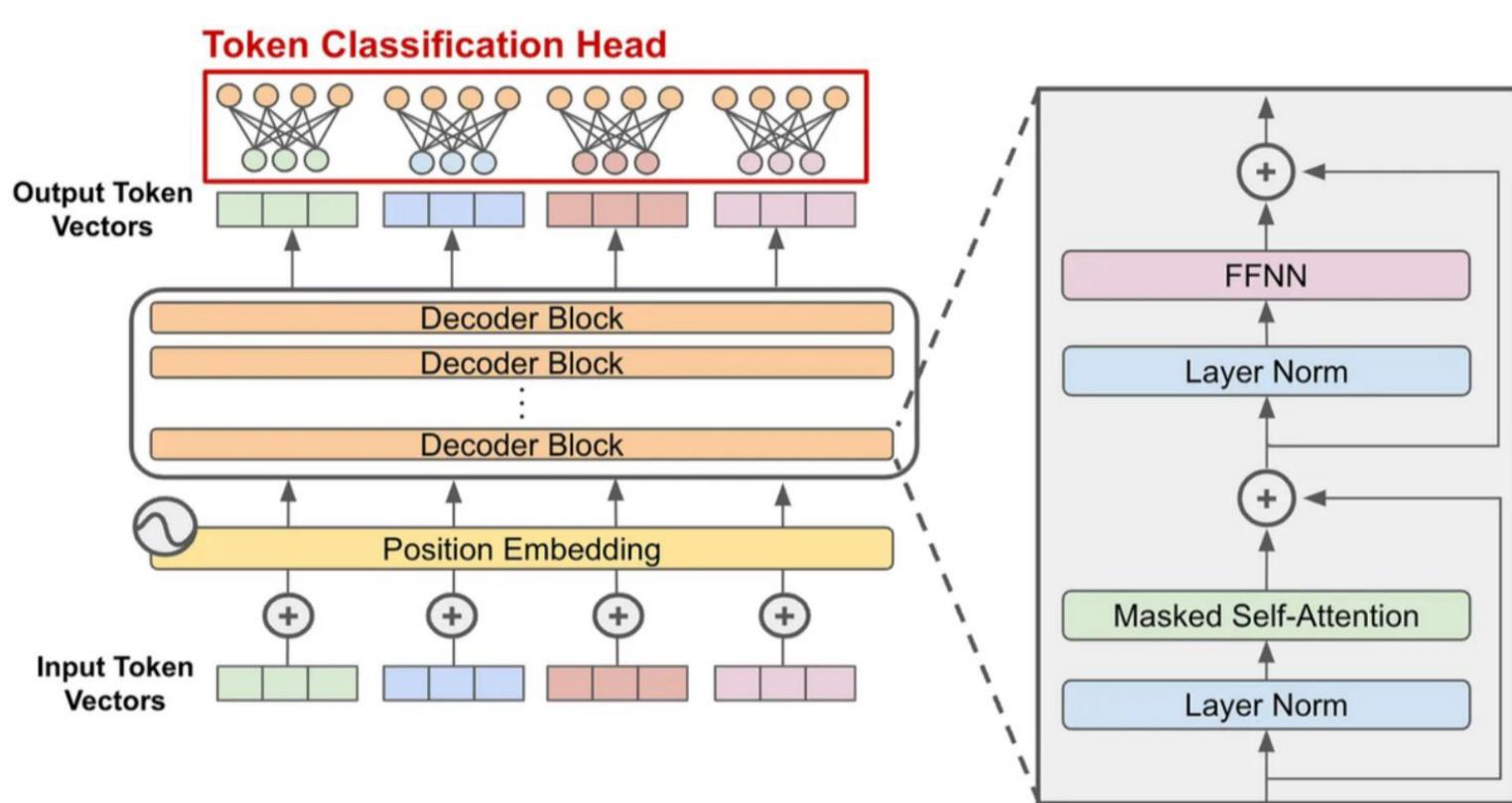
BA Pass

They understand that language

(summarize)

< >

Expected output



input

Rosain Max positions = 512

inference = 1024 γ 512

Absolute position encodings — Context length

Absolute

Relative

0 1 2 3 4
The dog chased another dog

-2 -1 0 +1 +2
the dog chased another dog

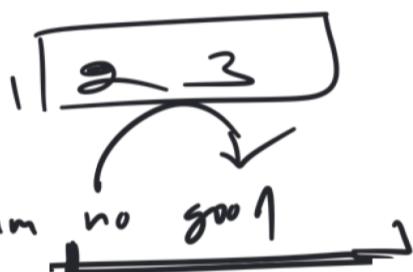
no good



phrase has same meaning

pos token no.	10	"
posn token	110	111

0 1 2 3



I am no good

Evegen believes that I am no good

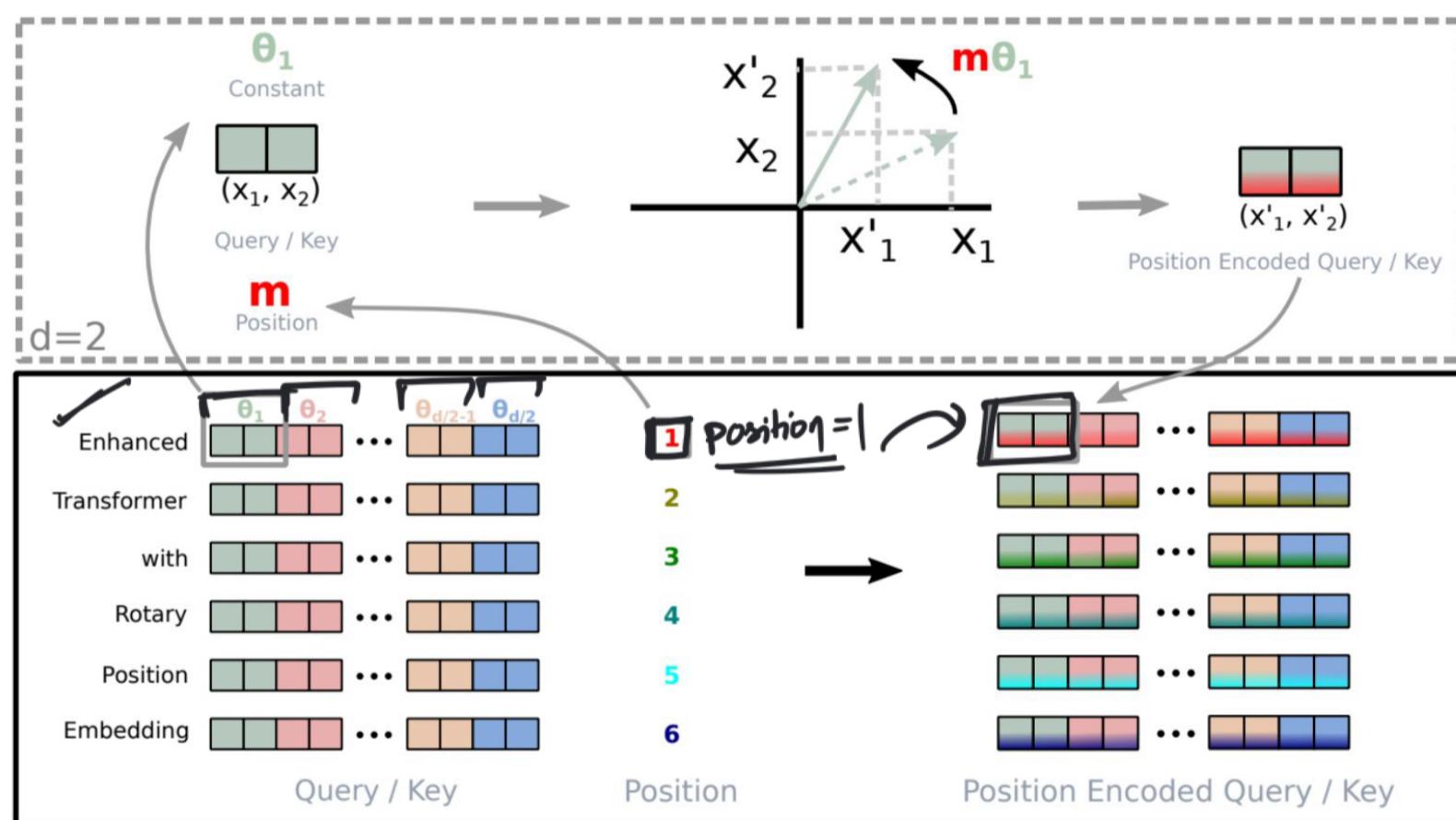


0 1 2 3 4 5 6

RoPE → encode the position as a rotation in vector space

instead of adding the pos info, rotate the Q, K vectors
↑
proportional to the position

rotate by pairs of the vector embedding



$Q_P \cdot K_Q$

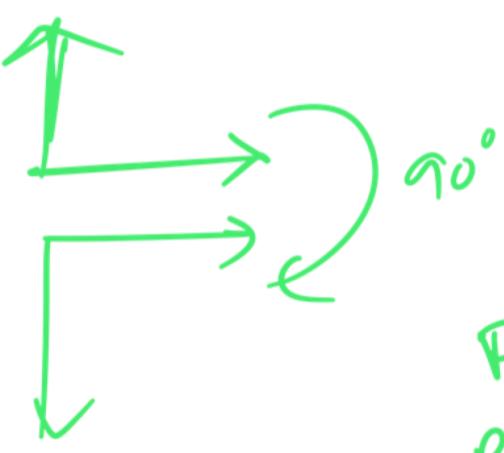
(Q rotated by P) . K rotated by q

1024

2028 ✓

info is not \downarrow P and q } no dependency on
info \Downarrow $(P - q)$ } absolute pos
} dependence on
diff of two positions

context length = 1M



original emb

Rotated

$x = \begin{bmatrix} \text{Pair 1} & \text{Pair 2} \\ [1, 0] & [0, 1] \end{bmatrix} \xrightarrow{\text{L}} [0, -1, 0, 1]$

Pair 1 = 90°
 Pair 2 = 45° $(x, y) = 1, 0$

$\boxed{(x', y') = (x \cos \theta - y \sin \theta, x \sin \theta + y \cos \theta)}$
 $\boxed{(1, 0) \xrightarrow{90^\circ} (0, 1)}$

$(0, 1) \xrightarrow{45^\circ} (-0.70, 0.70)$

$\text{Word}_1 \text{ Emb}$ $\text{Word}_2 \text{-emb}$

$$R_{pt}(\text{Word}_1, \theta_{251}) \cdot R_{pt}(\text{Word}_2, \theta_{252})$$

Q K

$Q \cdot K = \text{Word}_1 \cdot \text{Word}_2 \left(\theta_{252} - \theta_{251} \right)$ } vector

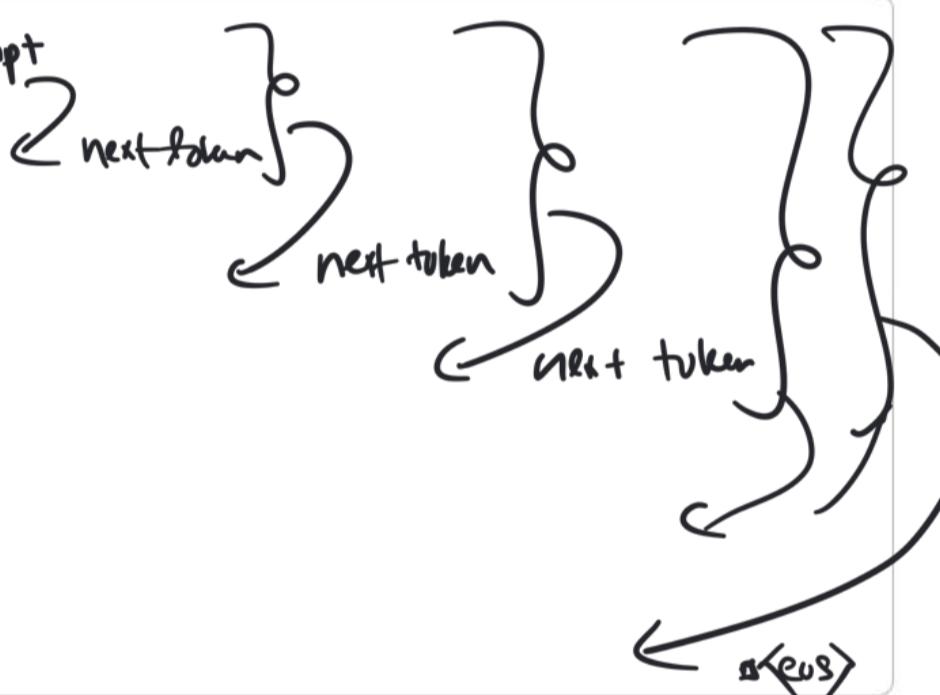
Deep dives RoPE }
}

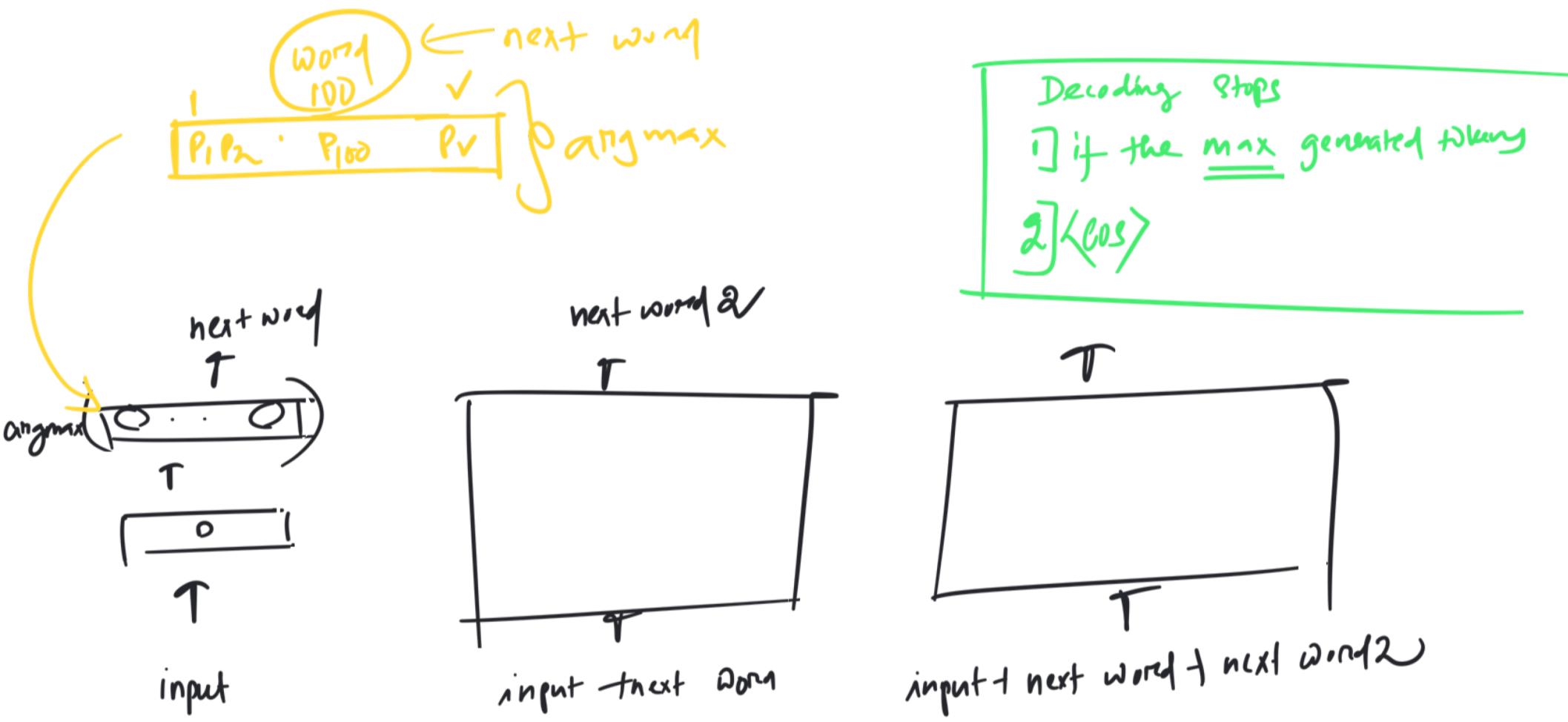
θ_k = position of token, rotation pair index

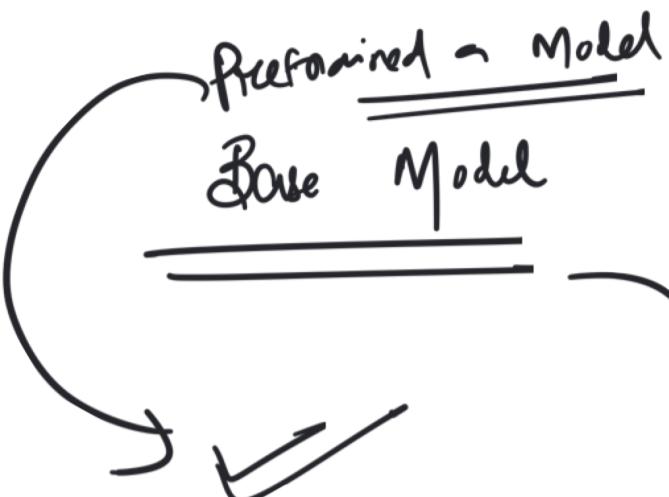


I have a dream (100%)
|
of (17.00%)
|
being (9.68%)
|
a (8.92%)
|
doctor (2.86%)
|
. (5.45%)

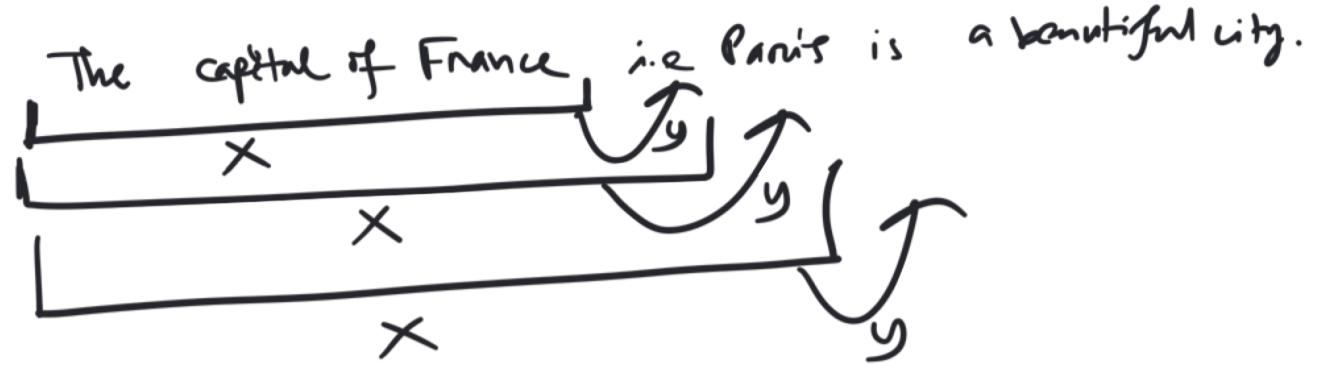
Prompt







Instruction Model



take the base

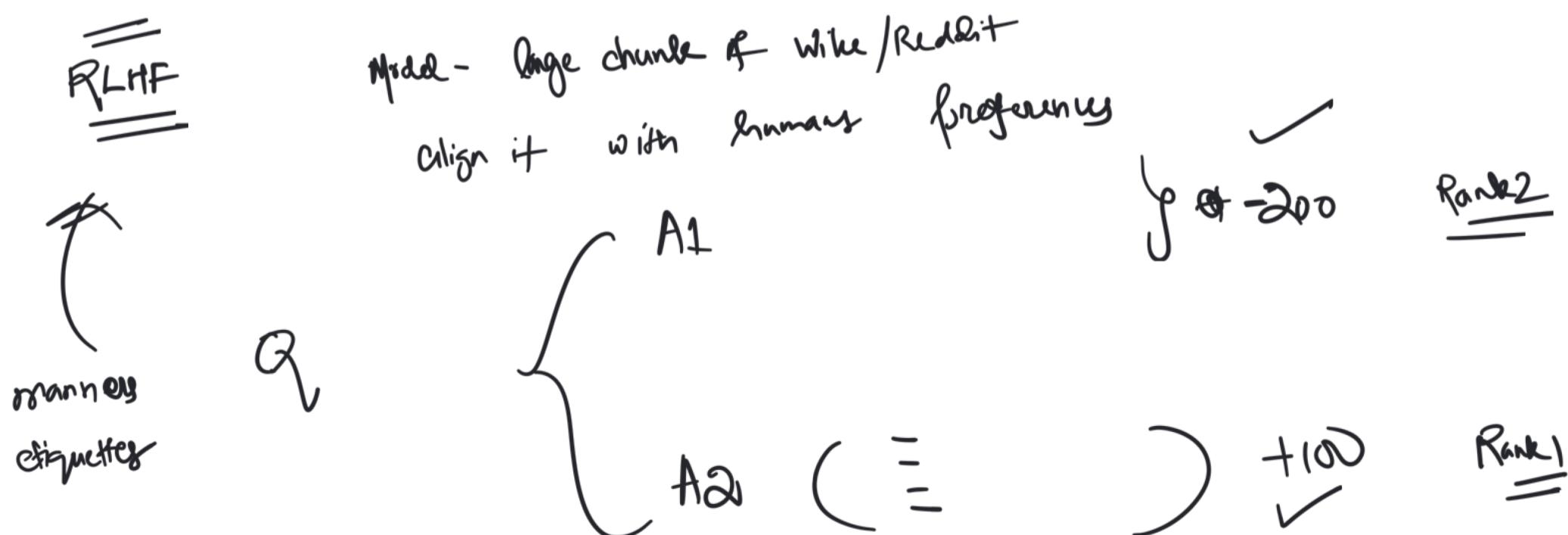
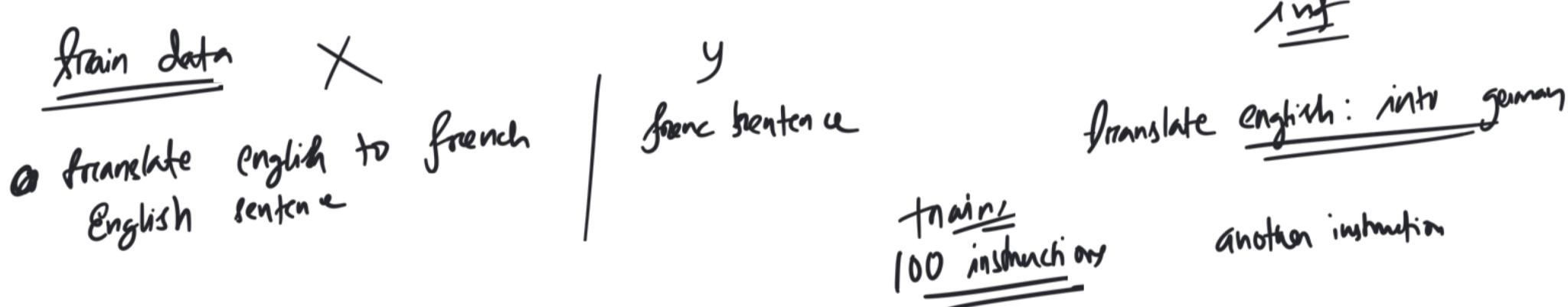
base What is capital of france Paris

rich data

translate English to French:
I am good

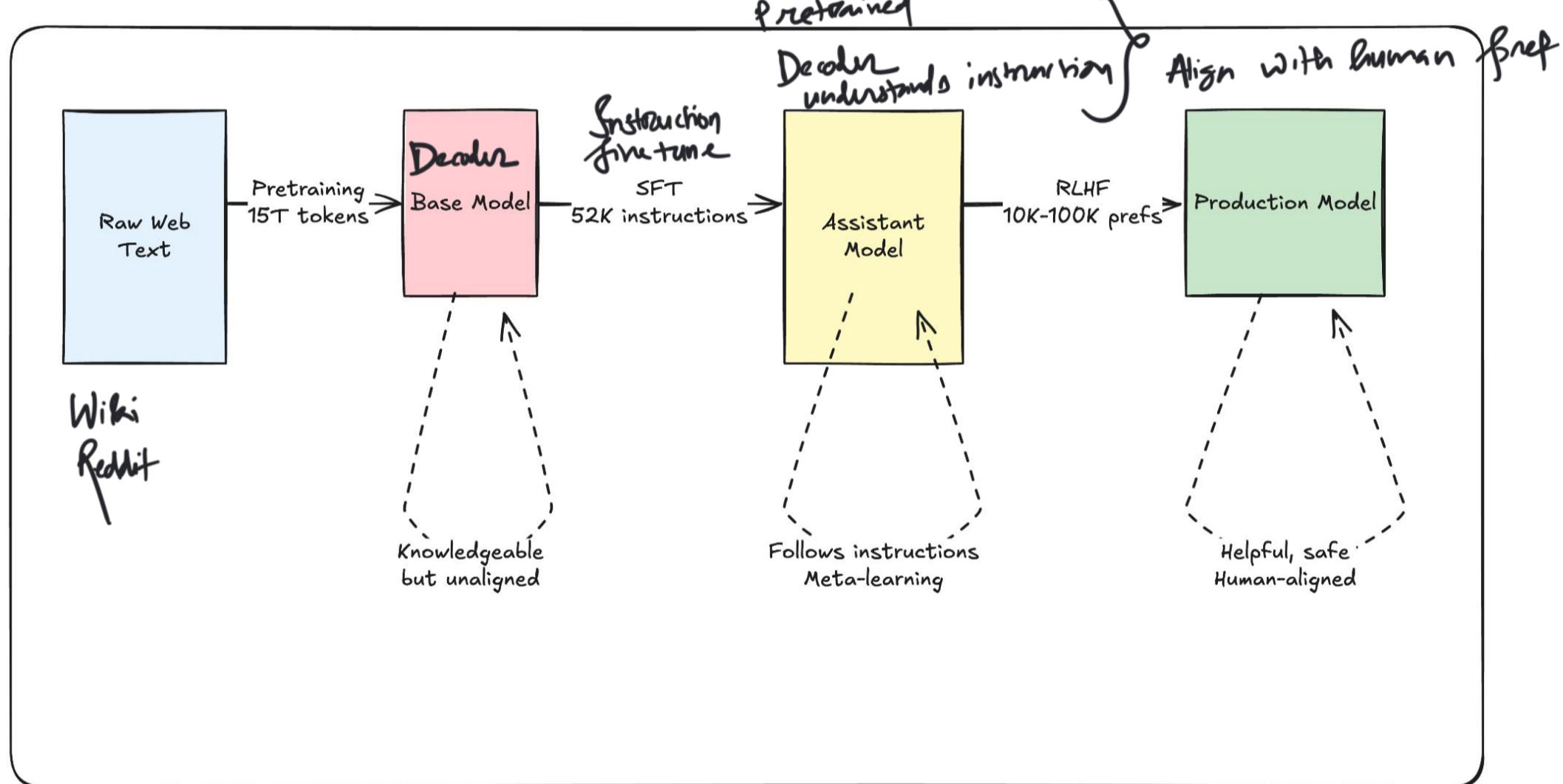


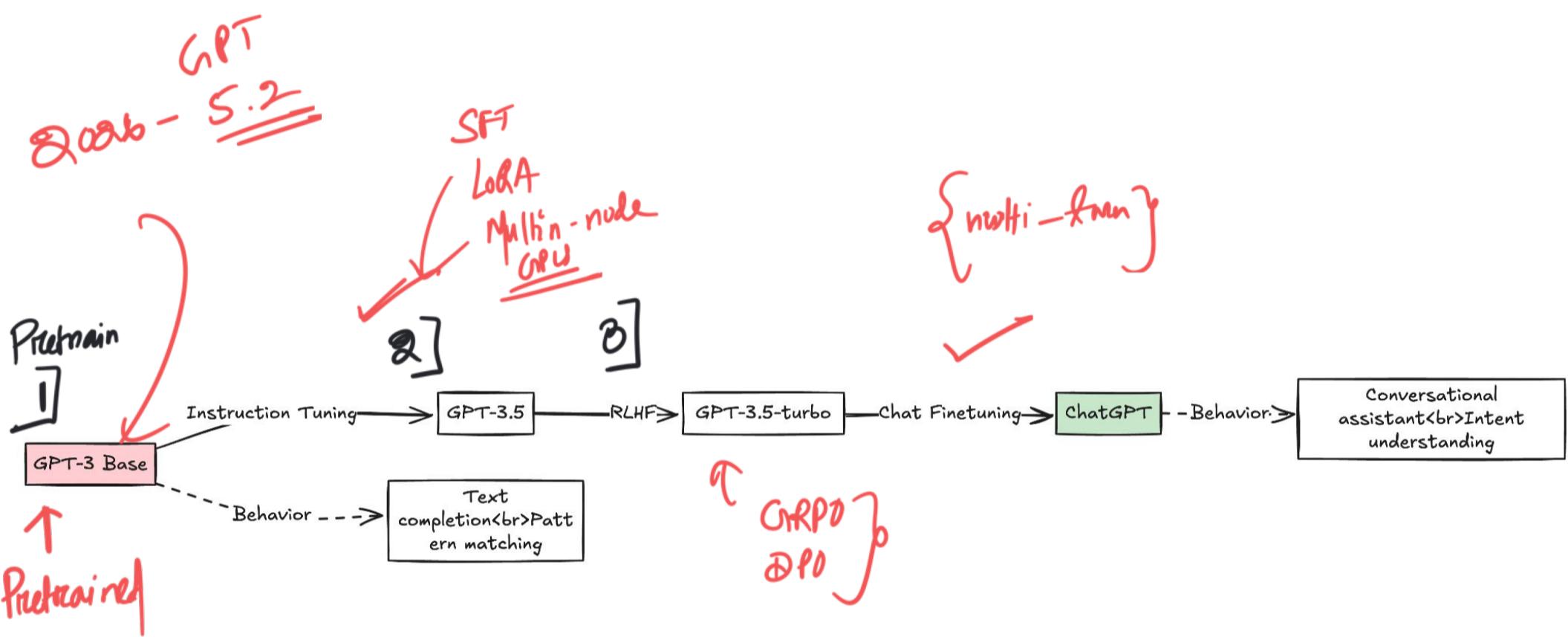
* Instruction finetuning: Meta-learning



```
{  
  "instruction": "Explain photosynthesis to a 10-year-old",  
  "response": "Photosynthesis is how plants make food from sunlight! The green  
leaves act like tiny solar panels. They take in sunlight, water from the  
ground, and carbon dioxide from the air. Then, like a magical kitchen, they  
mix it all together to make sugar (food for the plant) and release oxygen  
that we breathe. That's why plants are so important—they make their own food  
and give us fresh air!"  
}  
  
{  
  "instruction": "Translate to French: Hello, how are you?",  
  "response": "Bonjour, comment allez-vous?"  
}  
  
{  
  "instruction": "Write a Python function to reverse a string",  
  "response": "Here's a simple function:\n\n```python\ndef reverse_string(s):\n    return s[::-1]\n\n# Example usage\nprint(reverse_string('hello'))\n# Output: 'olleh'\n```\n"}  

```





1] Pretraining → Decoder

Supervised

Self-supervised

Unsupervised



2] Instruction FT

✓
human annotation

✓
human preference

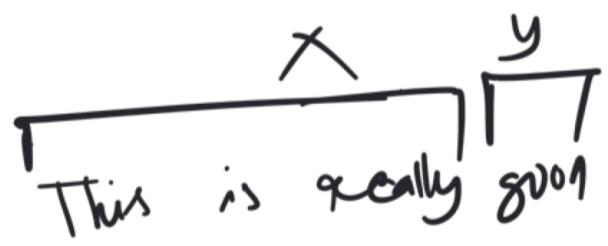
3] RL

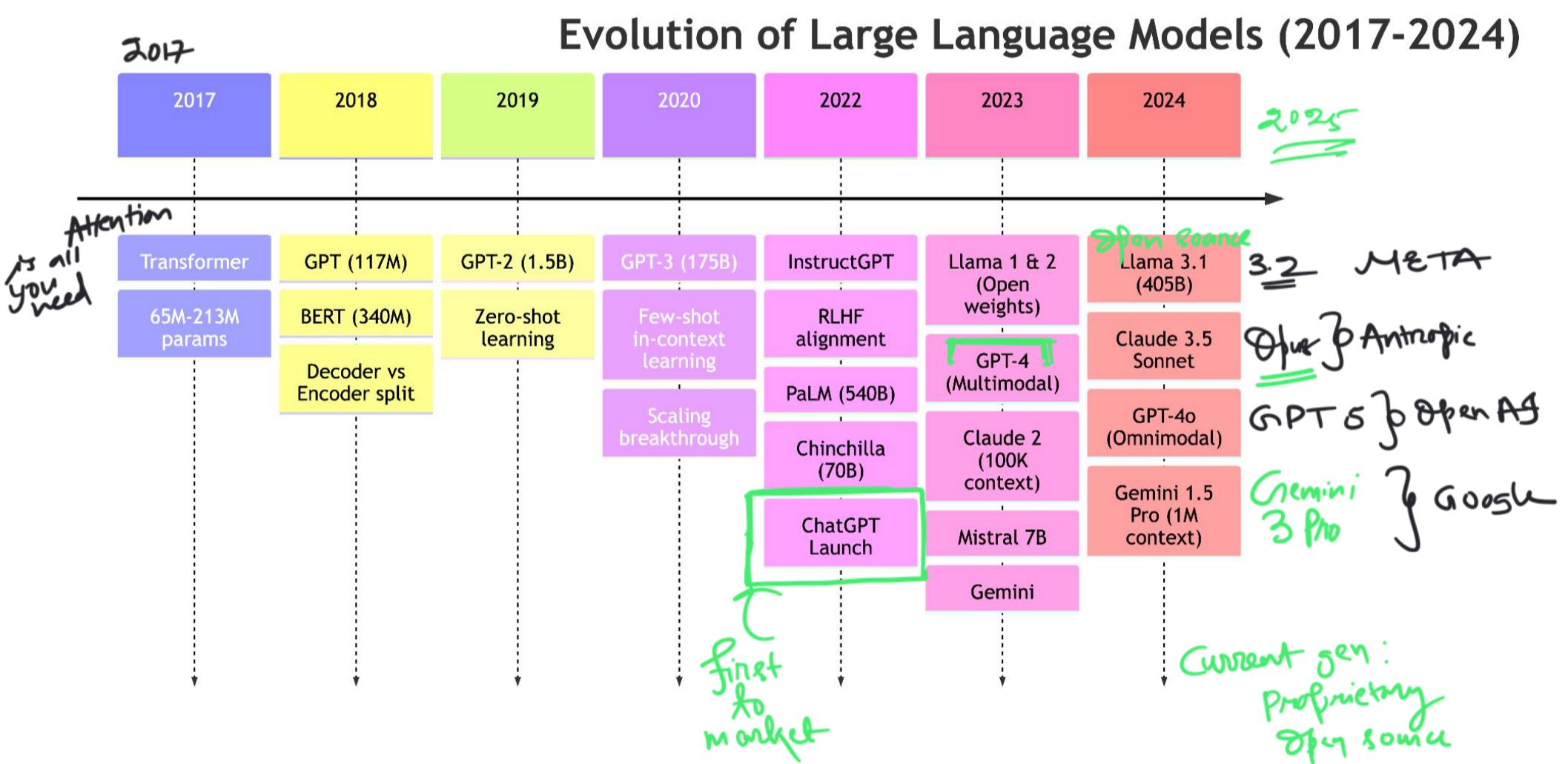
① Decoder only

② Modern LLM (long context)

③ Only pretraining — SFT

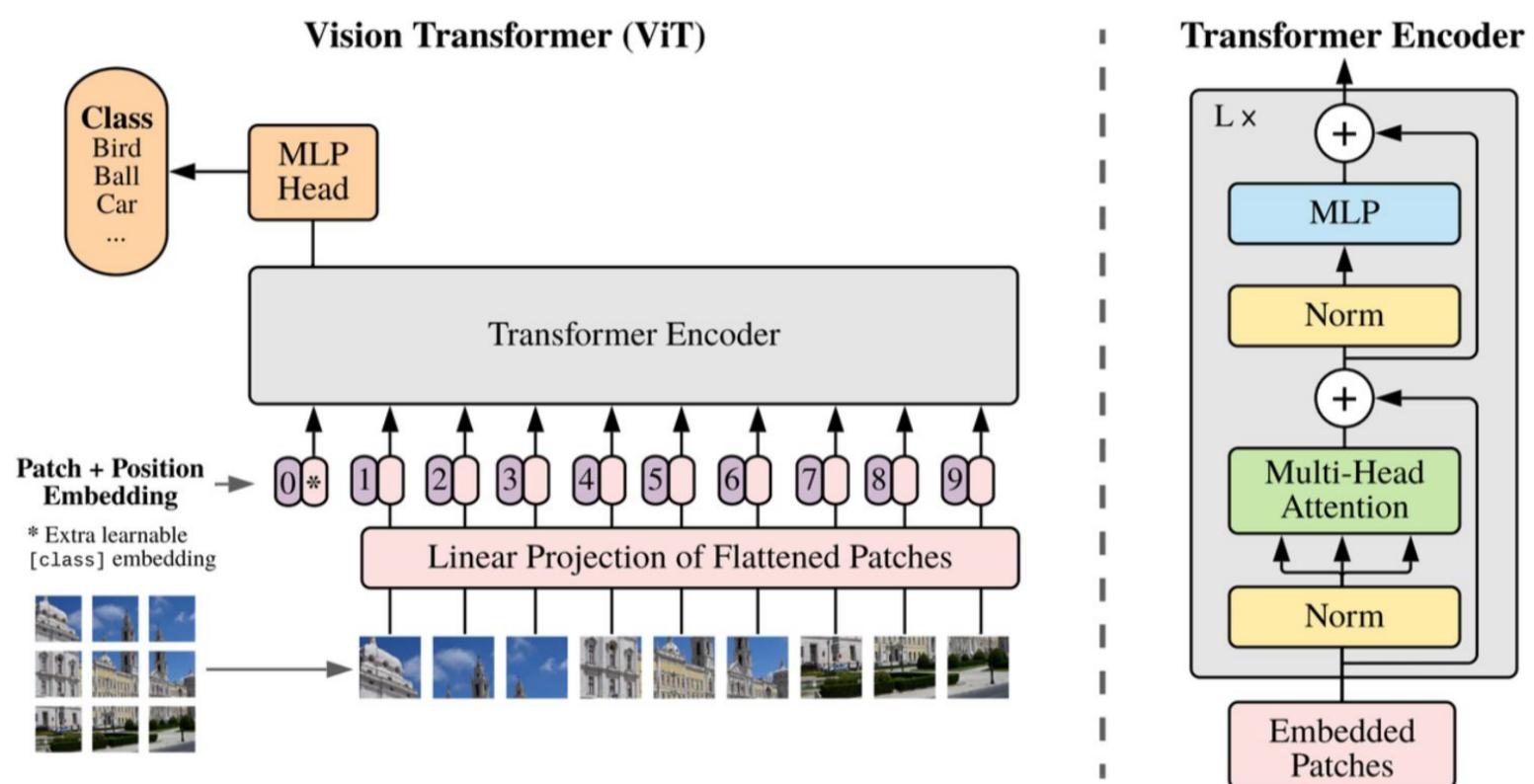
④ Align — RL

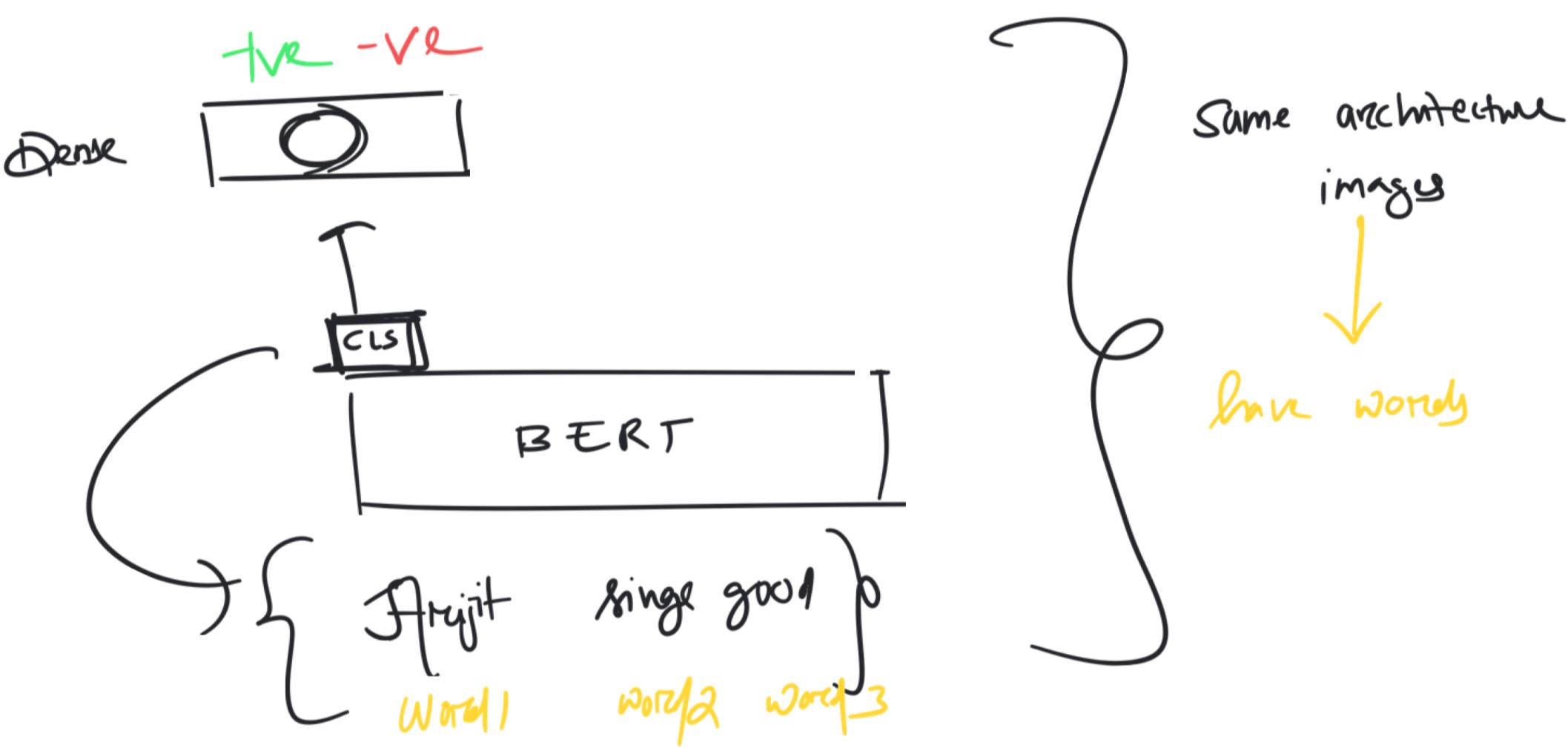


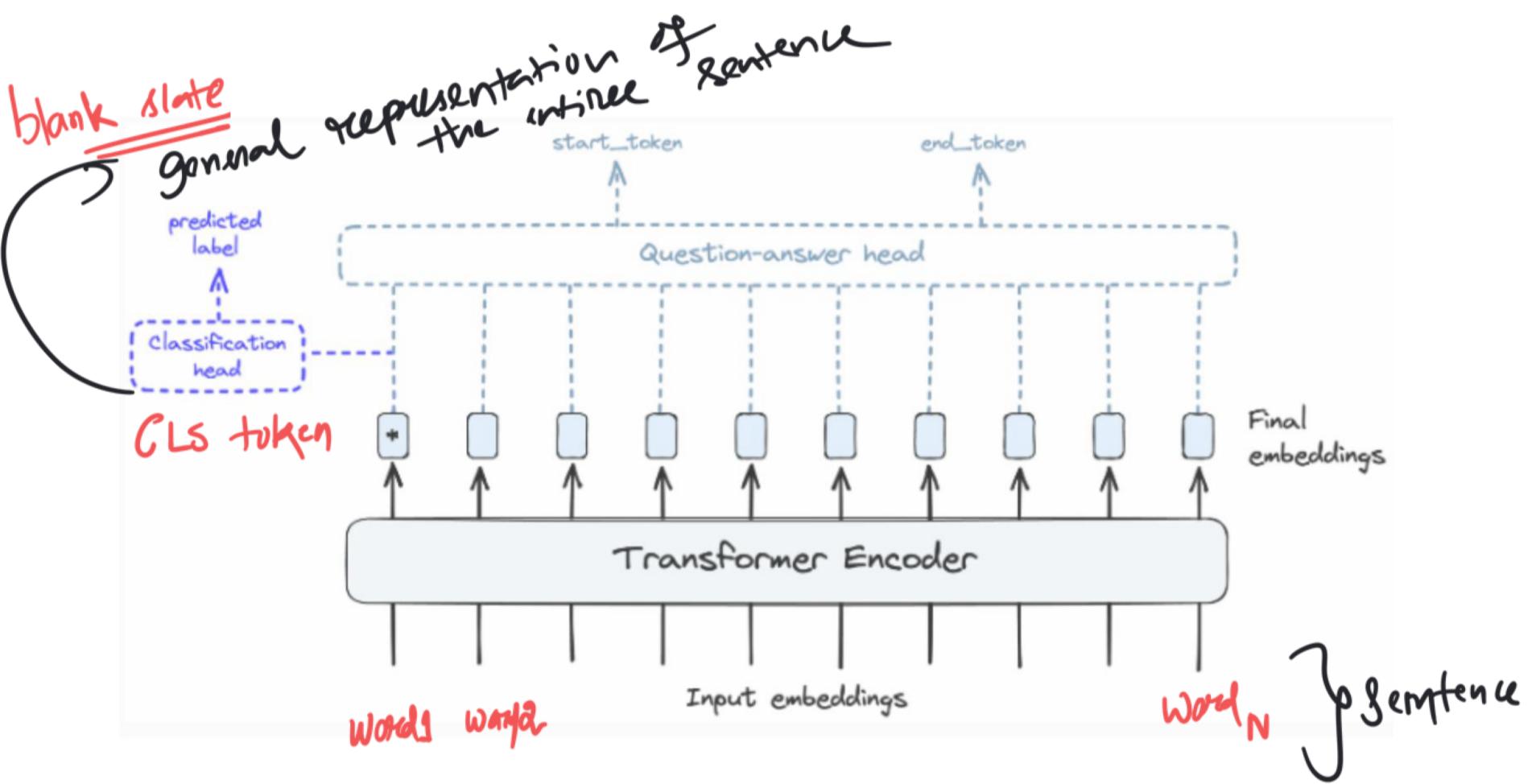


BERT

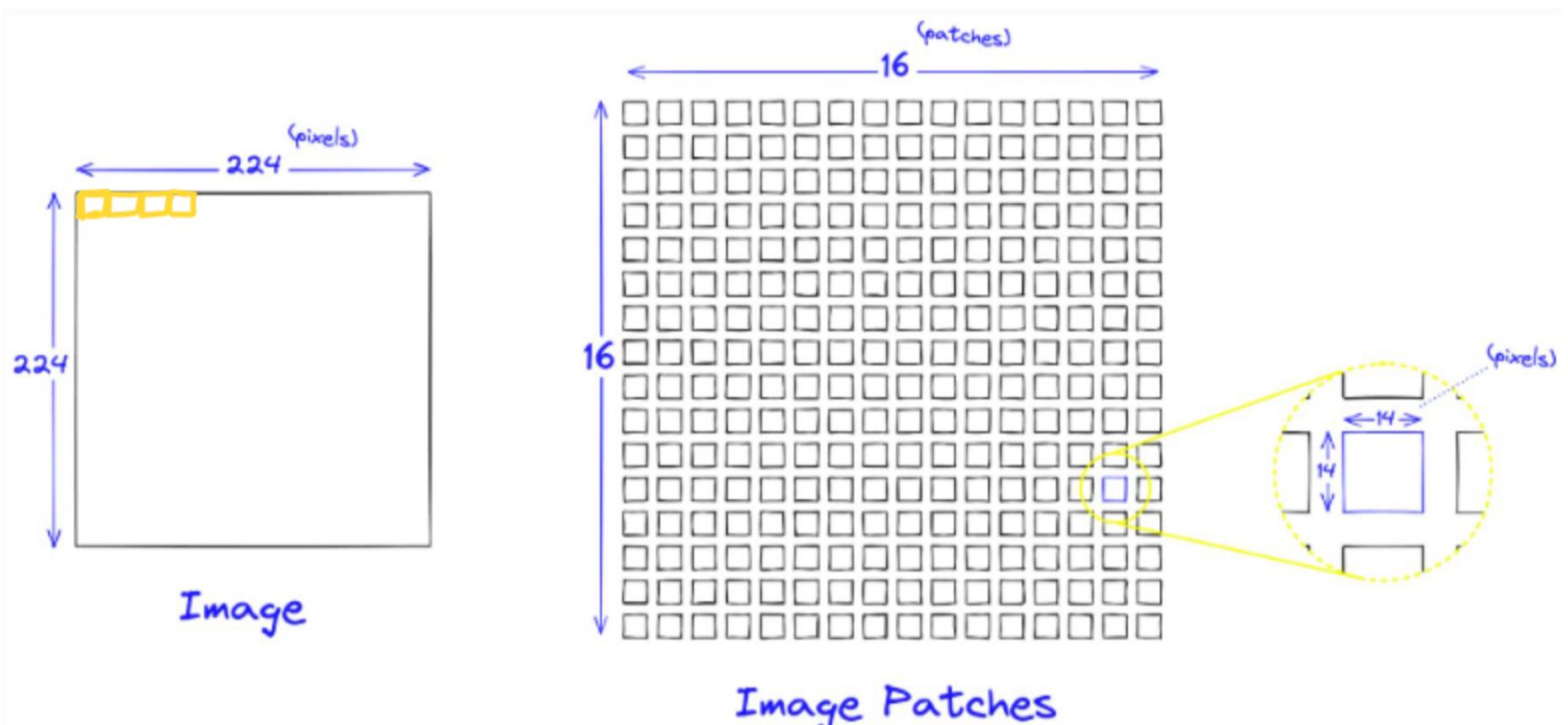
Vision Transformers



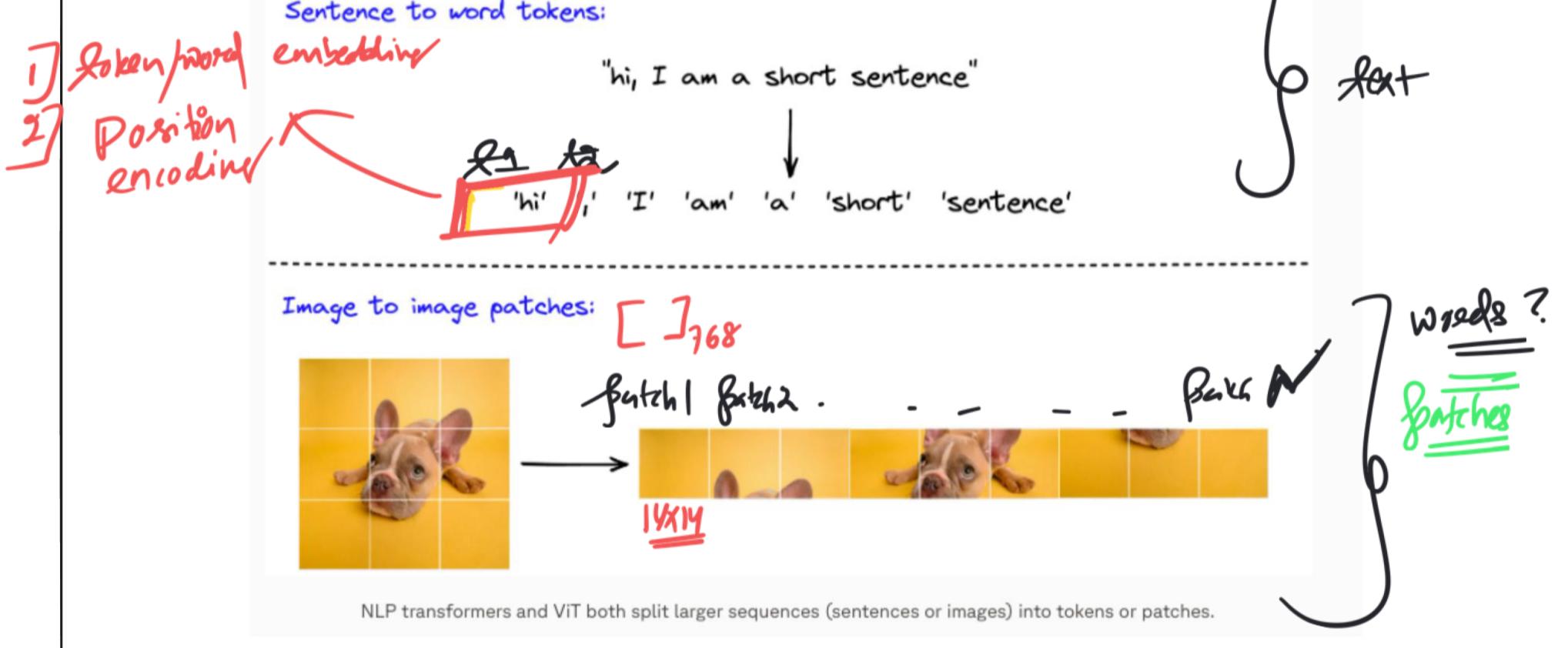


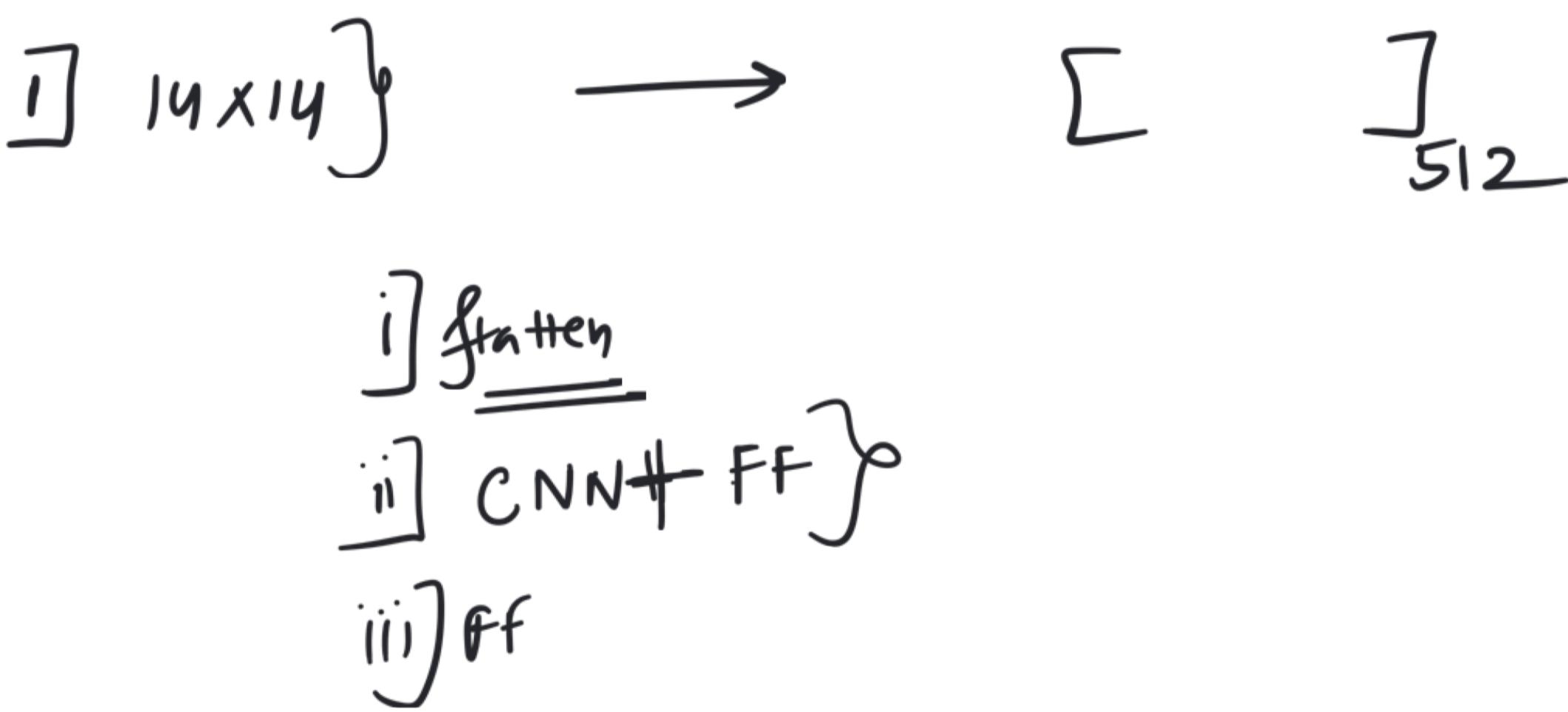


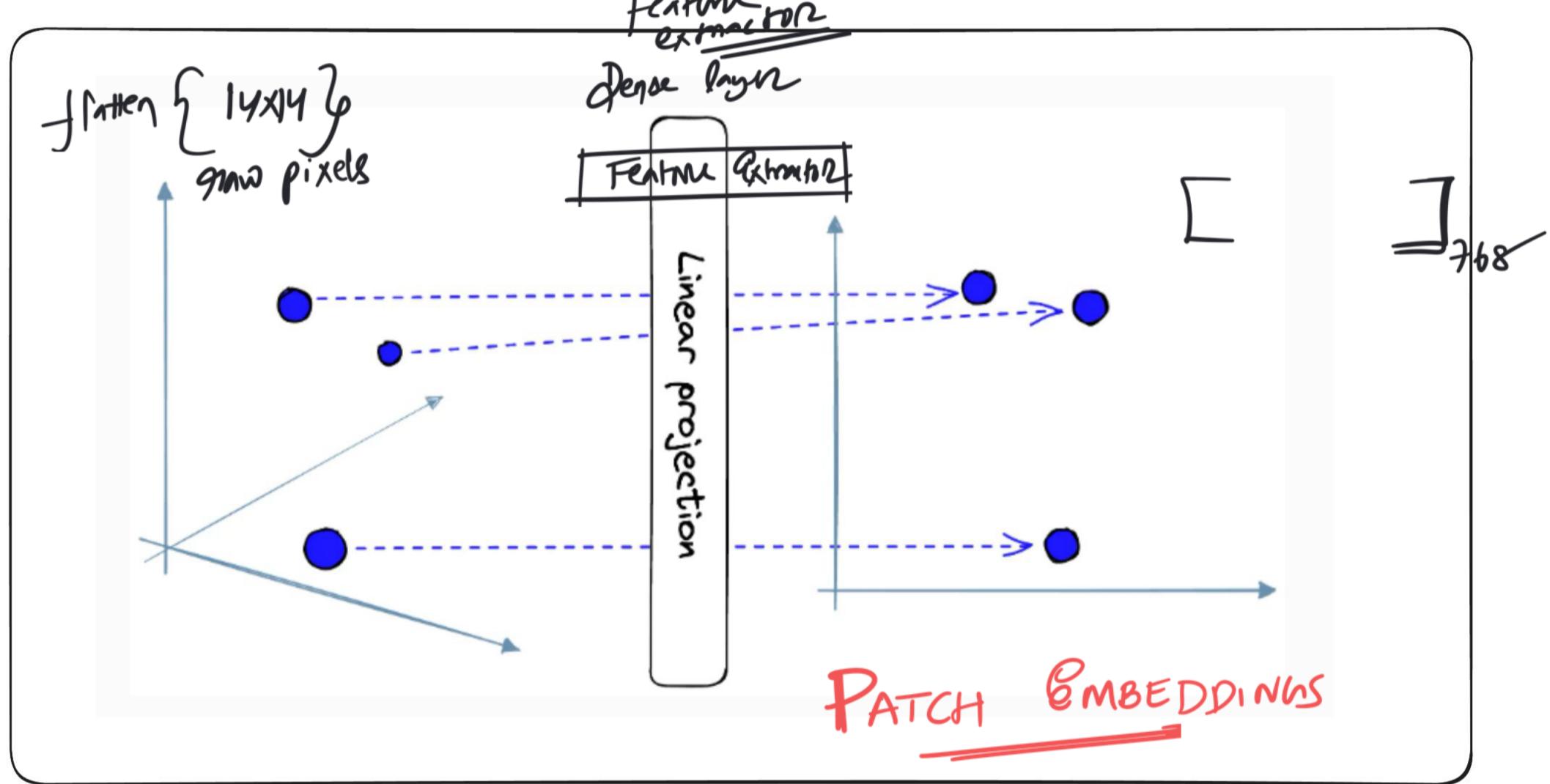
individual pixel values = Words 224×224 ✓

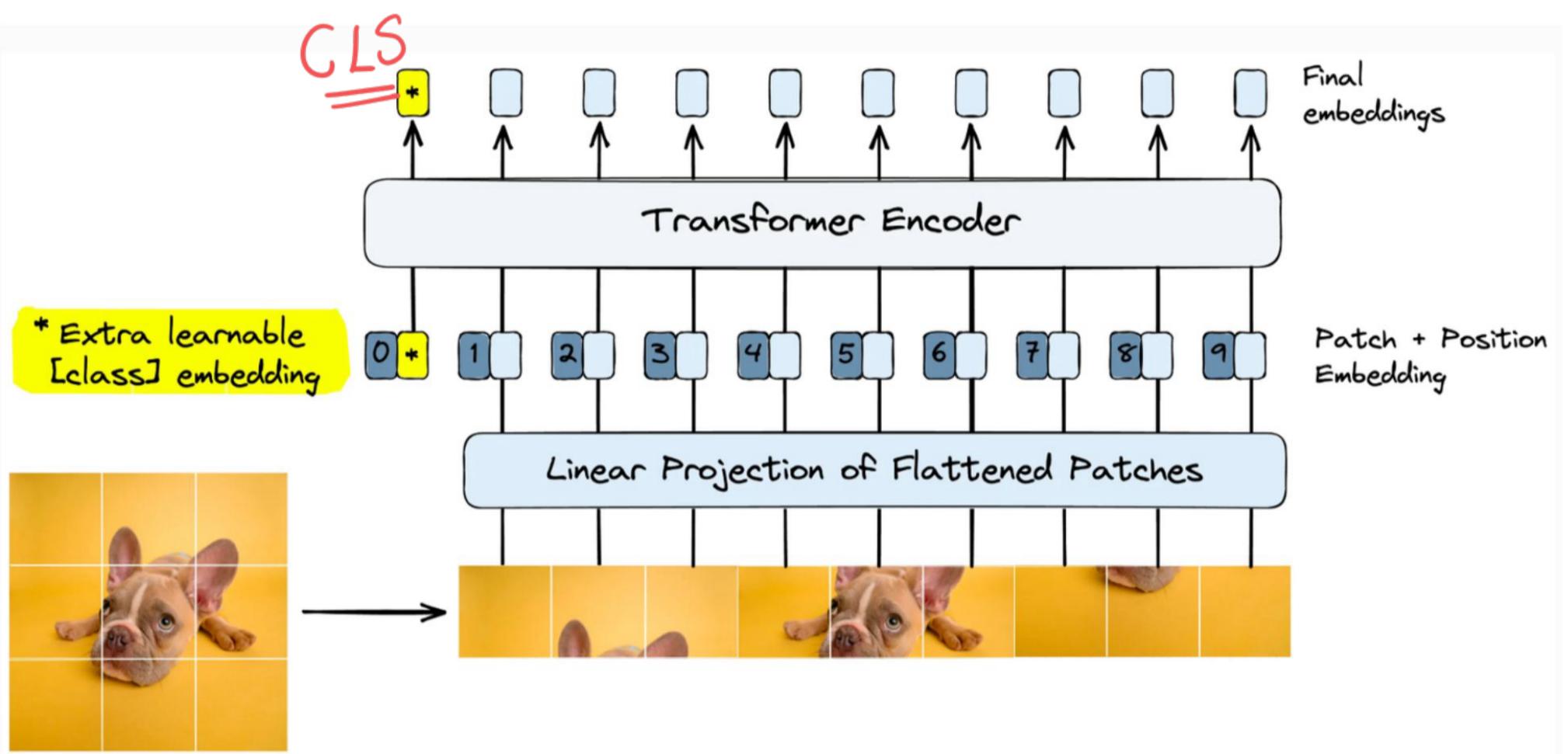


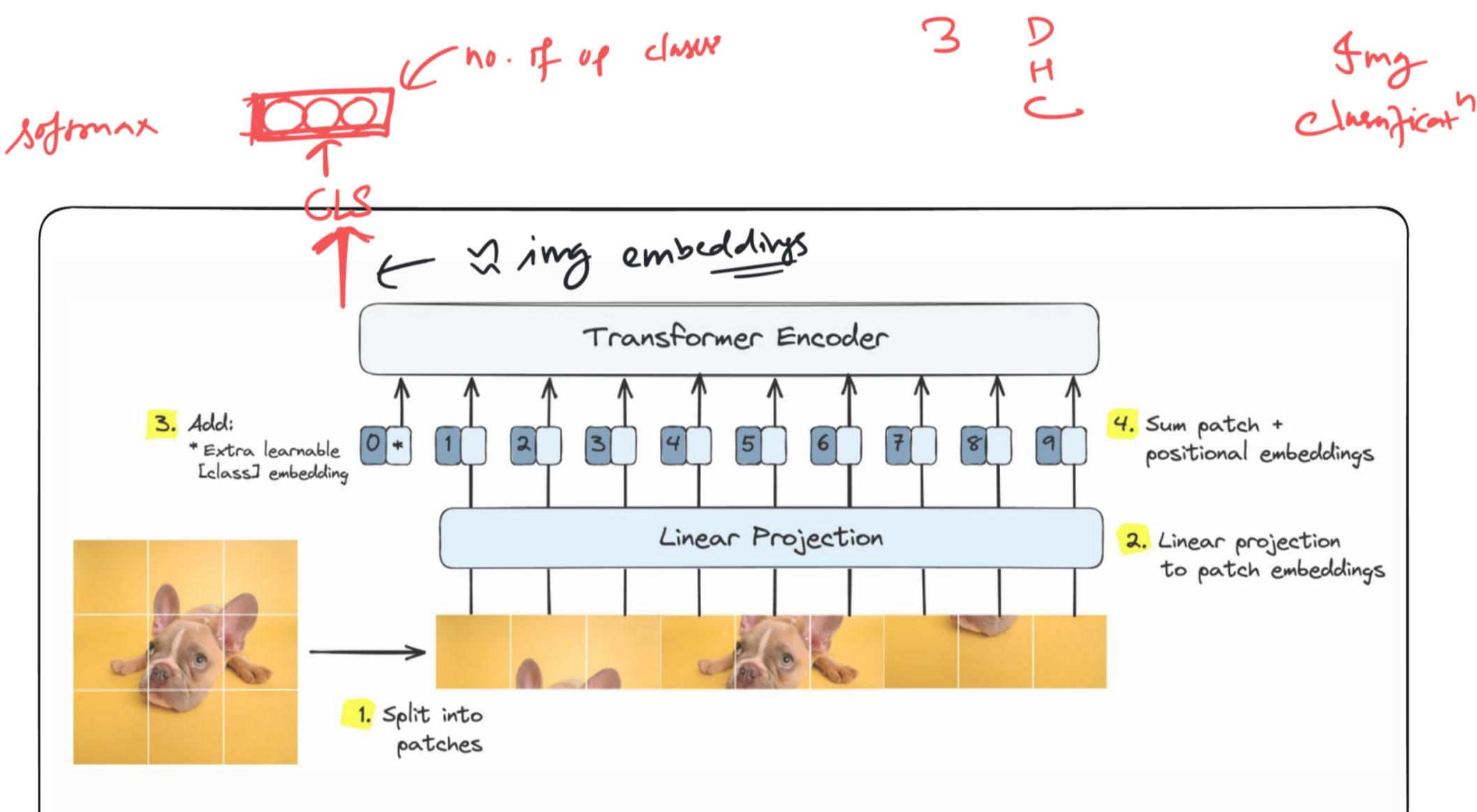
Conversion of 224x224 pixel image into 256 14x14 pixel image patches. ✓











Step 0 $224 \times 224 \times 3$

~~P1~~ P2

Step 1 Patches - 196 Patch

Step 2 Conv/Dense - $196 \times 16 \times 16 \times 3 \rightarrow 196 \times 768$

Step 3 Position encoding -

$196 \times 768 \}$ Position
on one

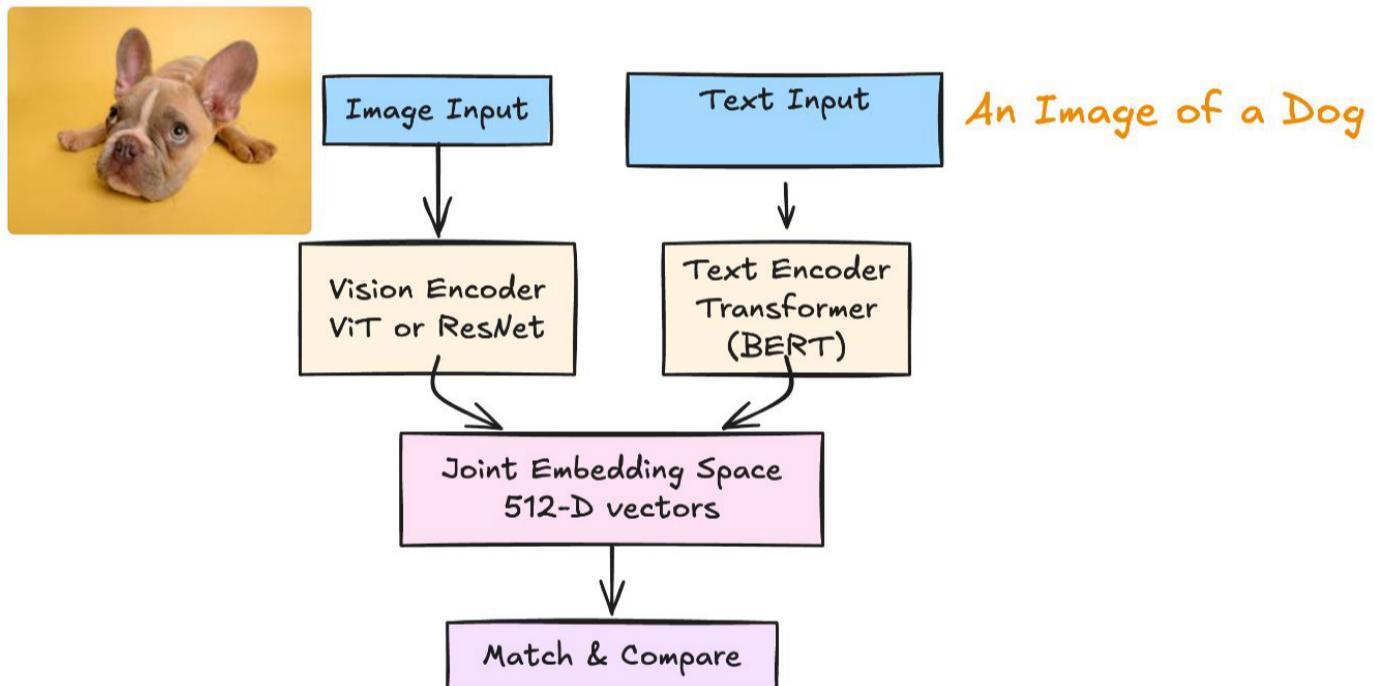
Step 4 CLS $(196 + 1) \times 768$

CLS

Step 5 Transformer

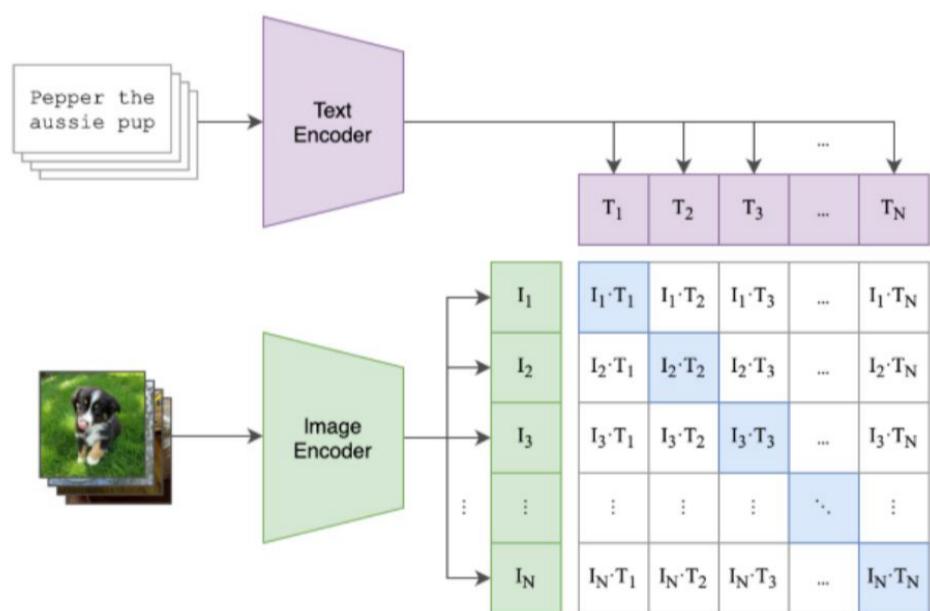
Step 6 CLS emb

Step 7 - Dense (768 , no. of class)

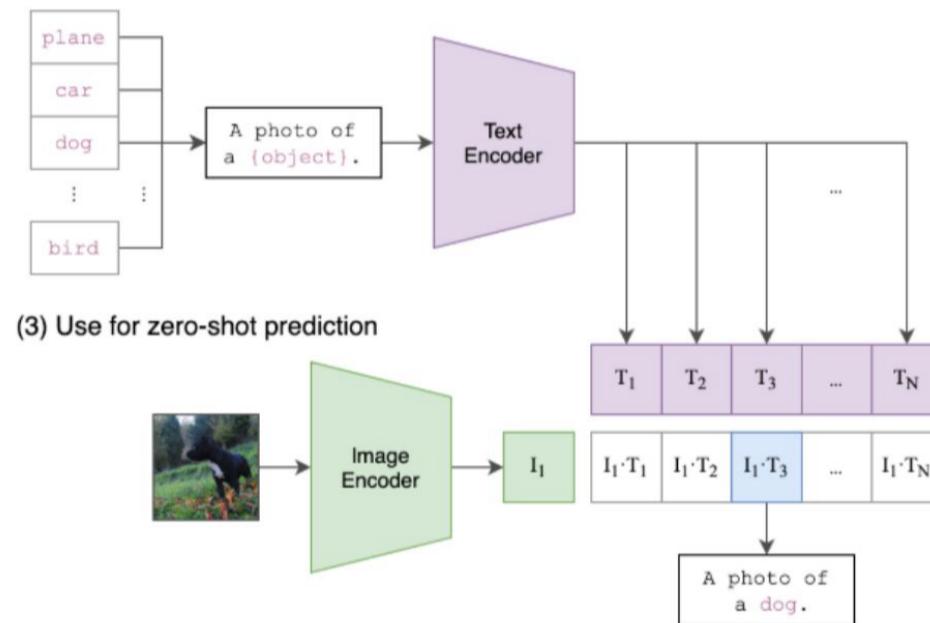


Approach

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

