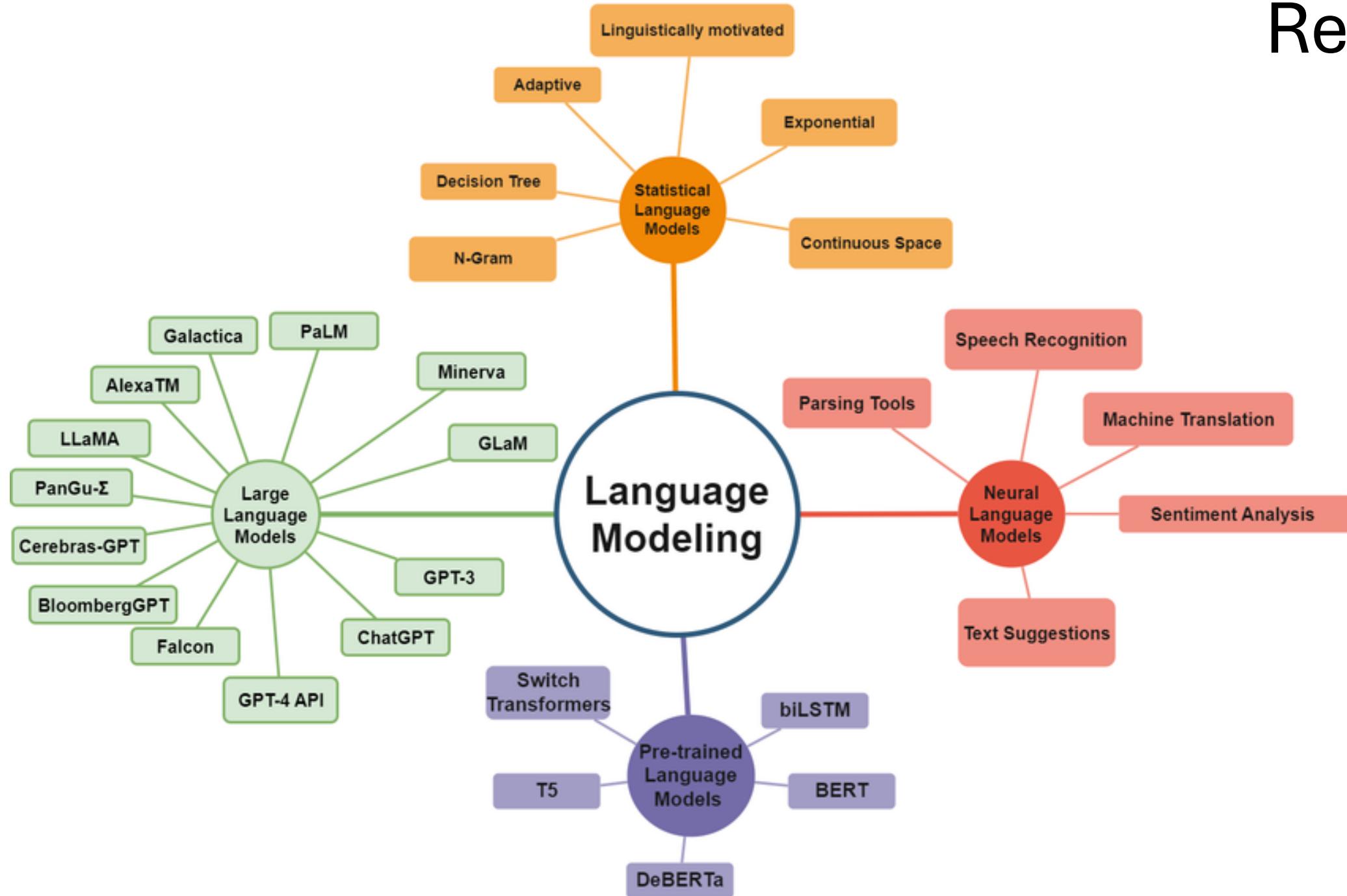


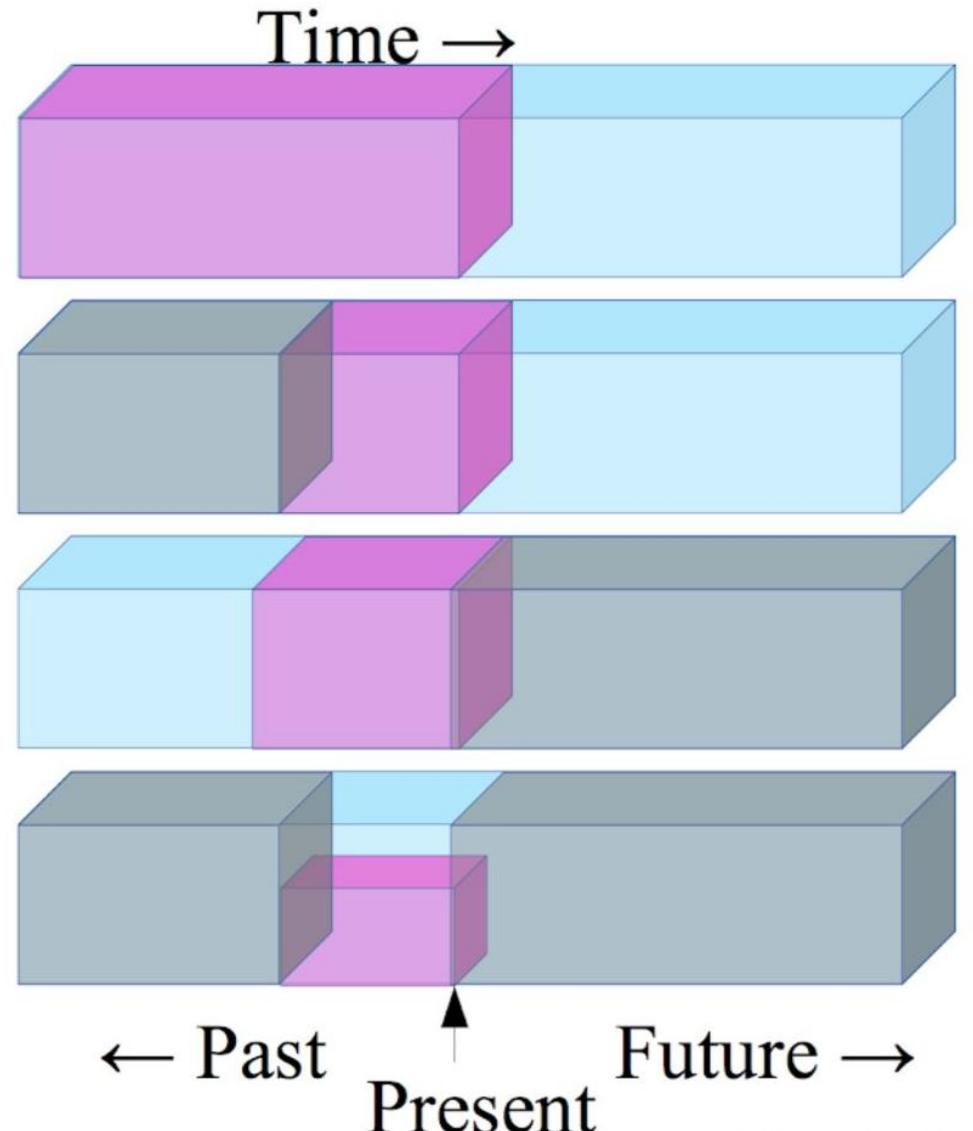
IR with LLMs and RAGs

- Contextual LLMs – Recap
- Non-Textual LLMs
- IR with LLMs
- IR with LLMs – Problems
- Solutions - Fine Tuning
- Solutions - RAGs
- Solutions – Prompting
- Agentic Rags

Recap



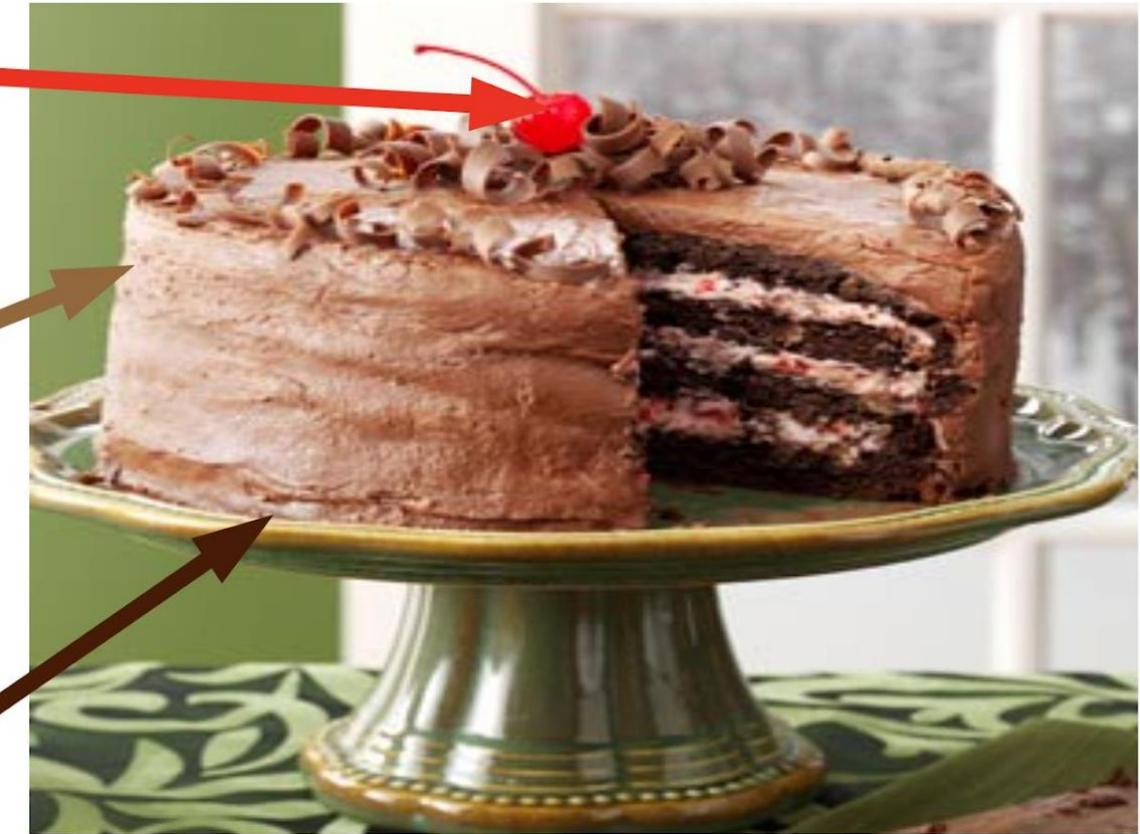
- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ Pretend there is a part of the input you don't know and predict that.



Slide: LeCun

Recap

- ▶ “Pure” Reinforcement Learning (**cherry**)
 - ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**
- ▶ Supervised Learning (**icing**)
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10→10,000 bits per sample**
- ▶ Self-Supervised Learning (**cake génoise**)
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**



Predict missing pieces

Recap

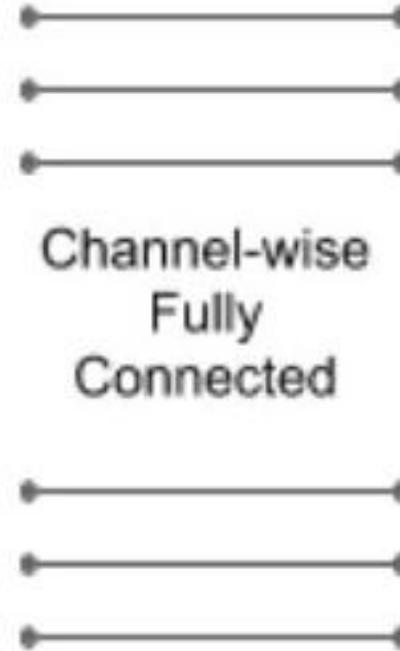


Context Encoders



Encoder

Encoder Features



Channel-wise
Fully
Connected

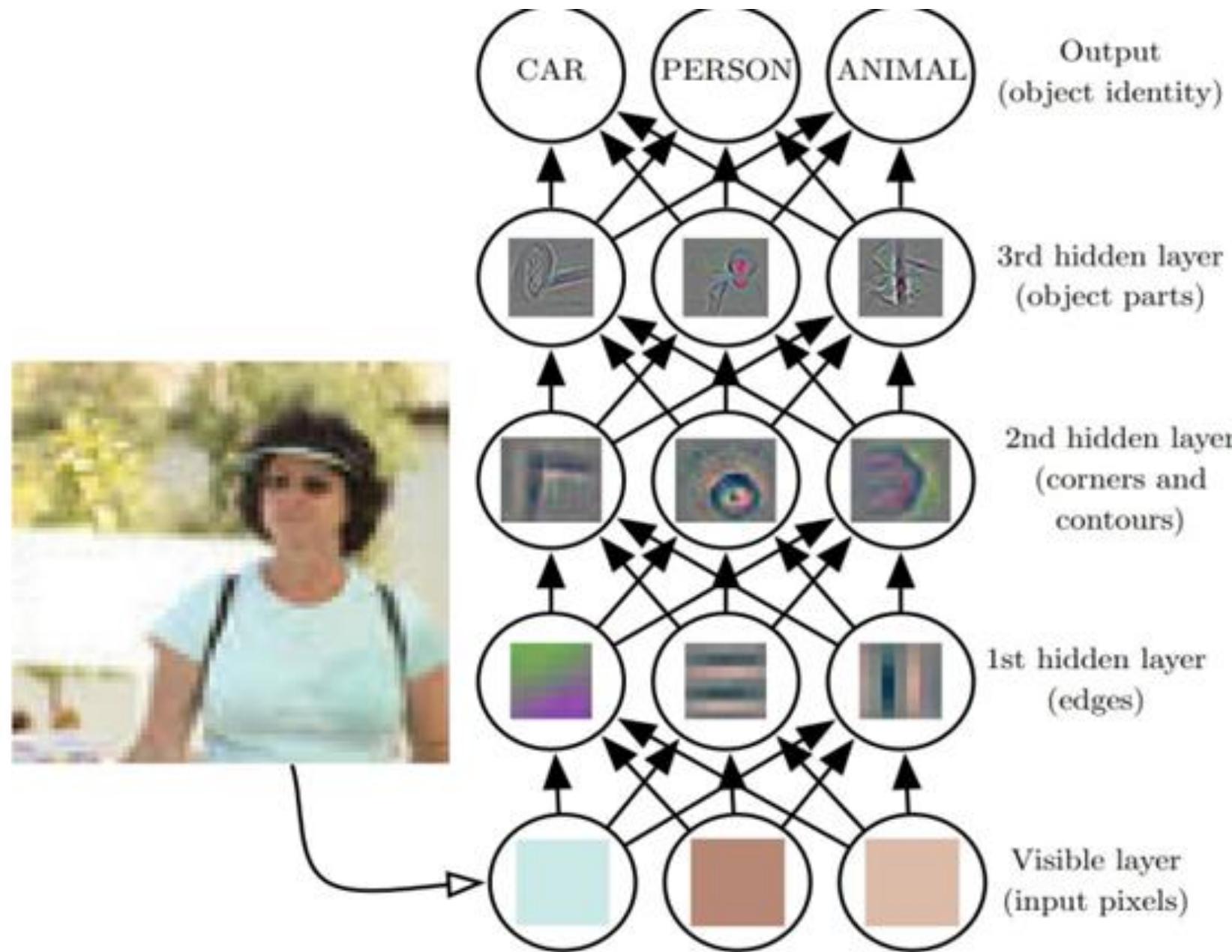
Decoder Features

Decoder

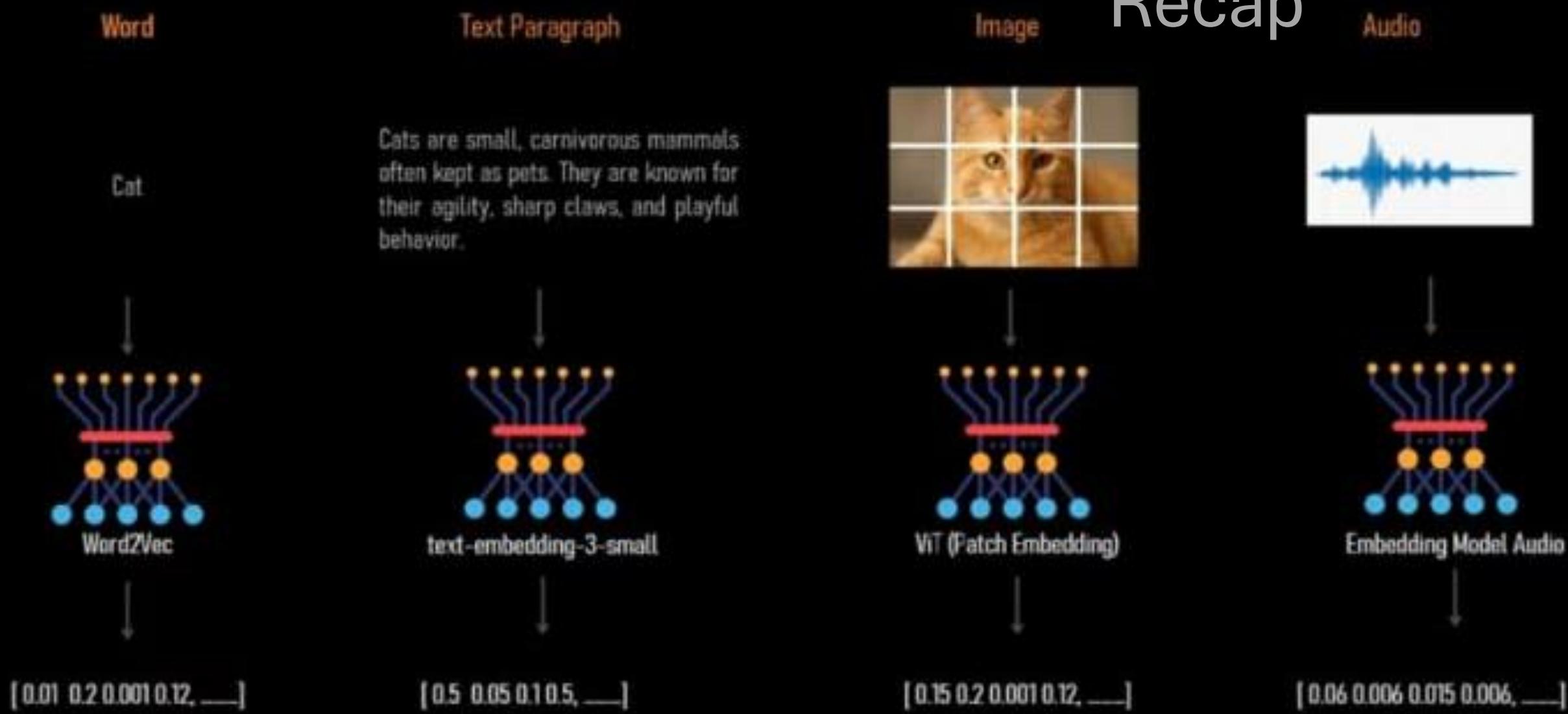
\mathcal{L}



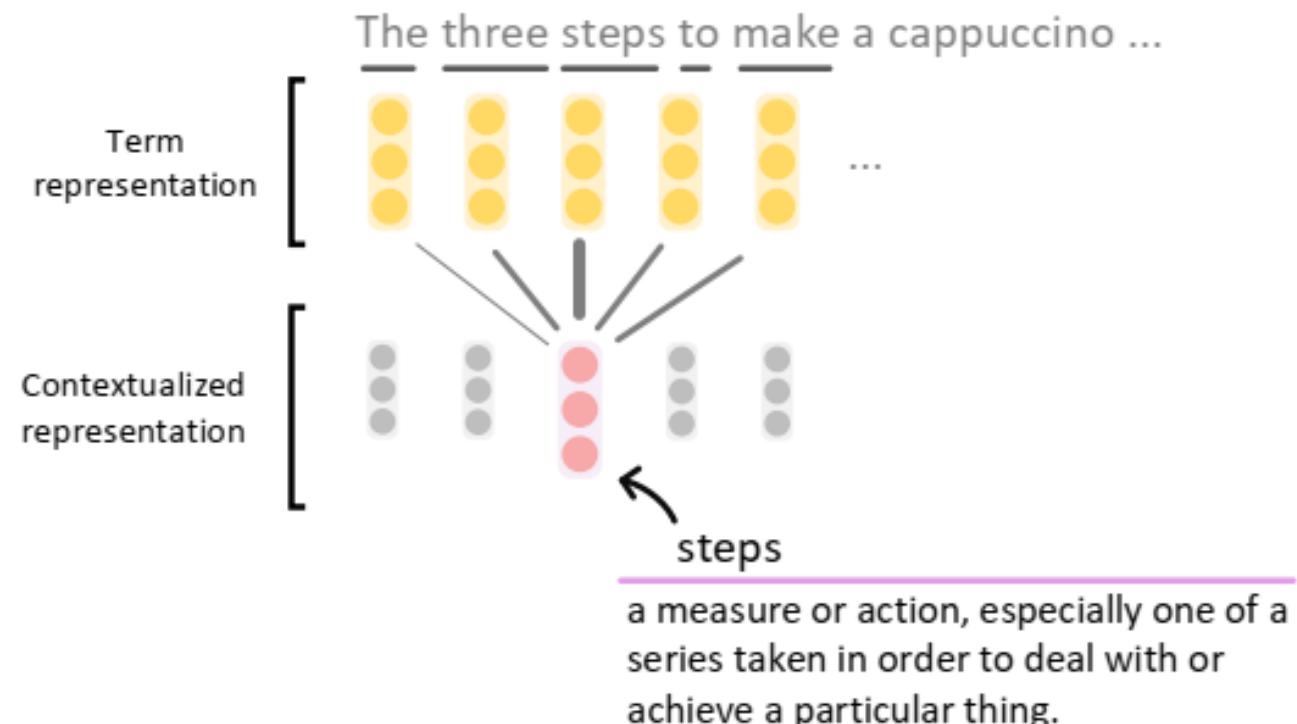
Recap



Recap

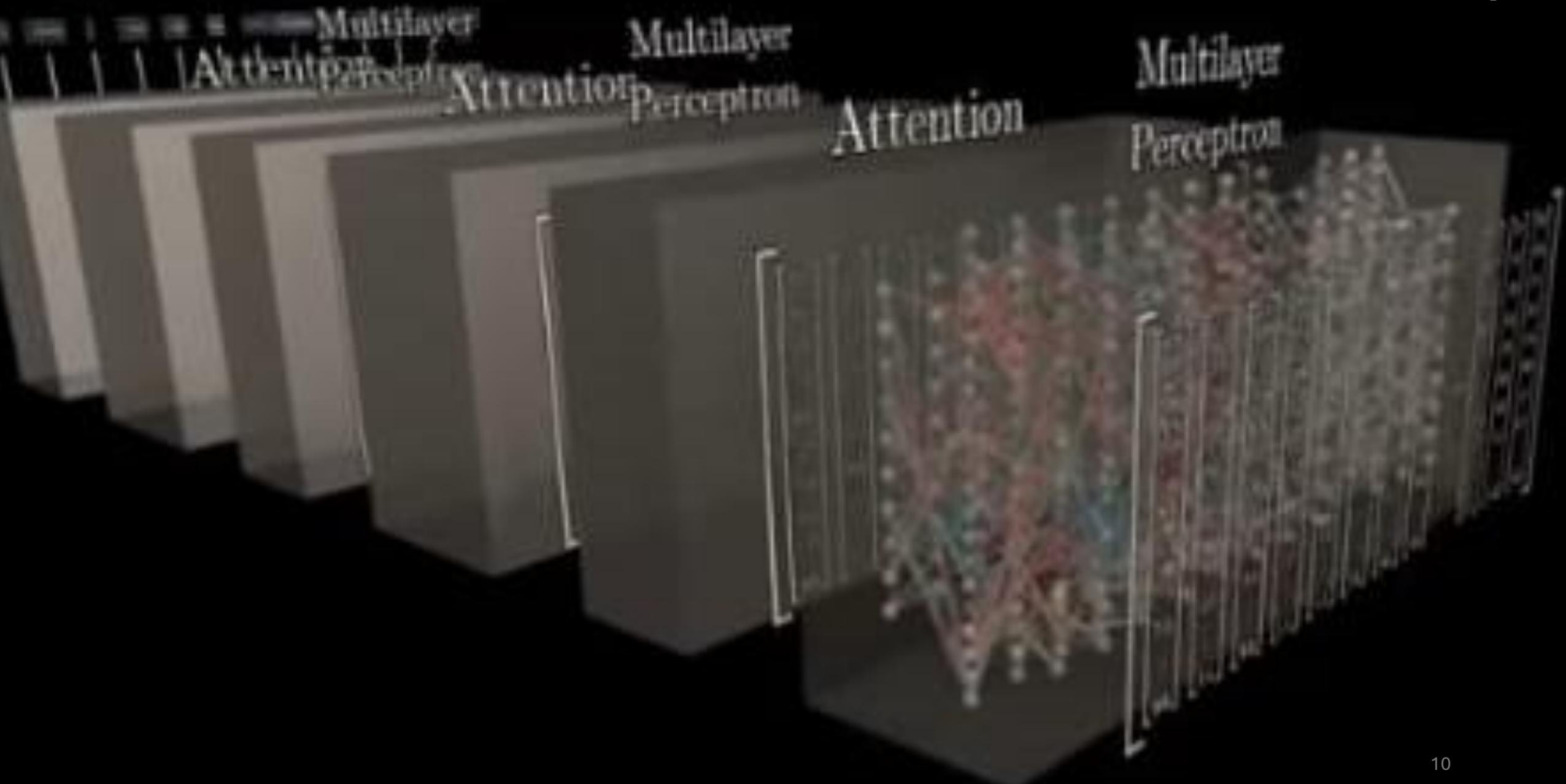


Contextualization via Self-Attention



- Learn meaning based on surrounding context for every word occurrence
- This *contextualization* combines representations
- Context here is local to the sequence (not necessarily a fixed window)
- Is computationally intensive $O(n^2)$
 - Every token attends to every other token

Recap



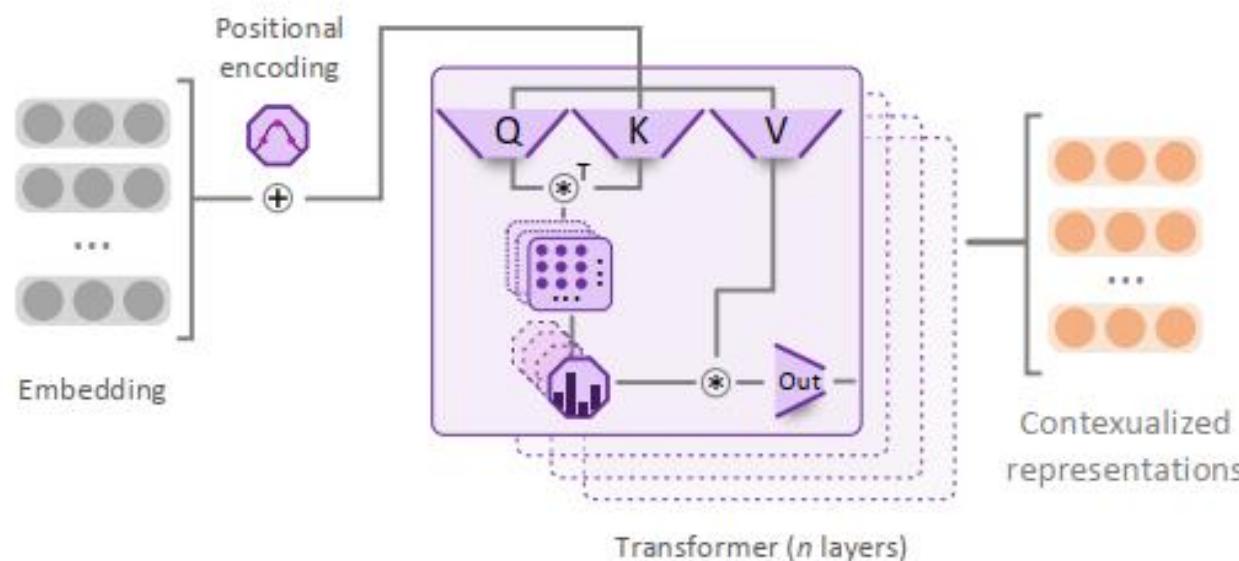
Transformer

- Transformers contextualize with multi-head self-attention
 - Every token attends to every other token $O(n^2)$ complexity
- Commonly Transformers stack many layers
- Can be utilized as encoder-only or encoder-decoder combination
- Do not require any recurrence
 - The attention breaks down to a series of matrix multiplications over the sequence
- Initially proposed in translation
 - Now the backbone of virtually every NLP advancement in the last years

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, et al.

Attention is all you need. In NeurIPS. 2017.

Transformer – Architecture



- We embed (subword) tokens
- We add a positional encoding
- In each Transformer-Layer:
 - Project each vector with 3 linear layers to **Query**, **Key**, **Value**
 - Transform projections to another multi-head dimension
 - Matrix-multiply Query & Key
 - Get Q-K attention via softmax
 - Multiply attention with Values and project back to output

Nice detailed walkthrough code + paper:

<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

Self-Attention Definition

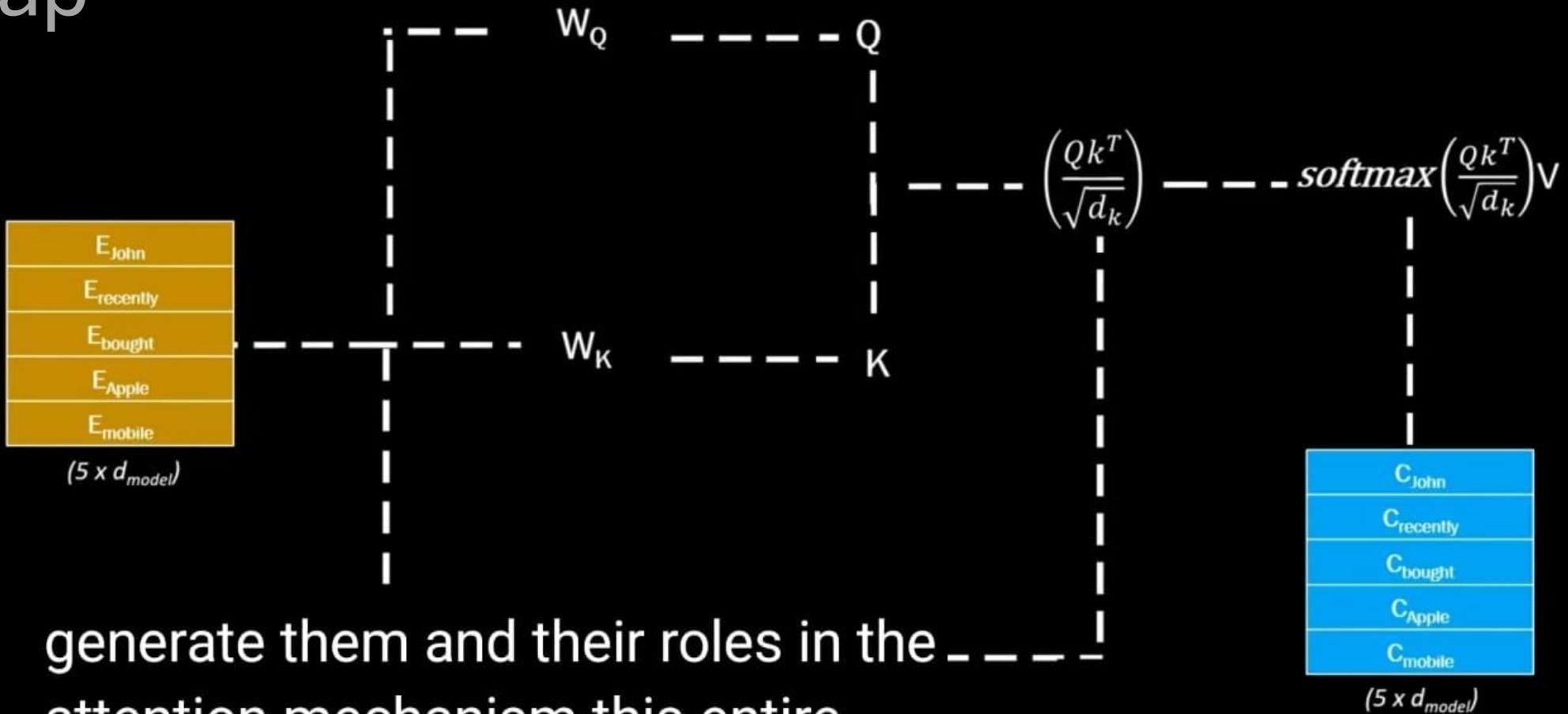
- The Transformer Self-Attention is defined as:

$$\text{SelfAttention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) * V$$

- Q, K, V are projections of the **same input** sequence
- This definition hides quite a bit of complexity, visible in the code

Q	Attention “Query”
K	Attention “Key”
V	Attention “Value”
d_k	Dimension of key embeddings

Recap



Recap

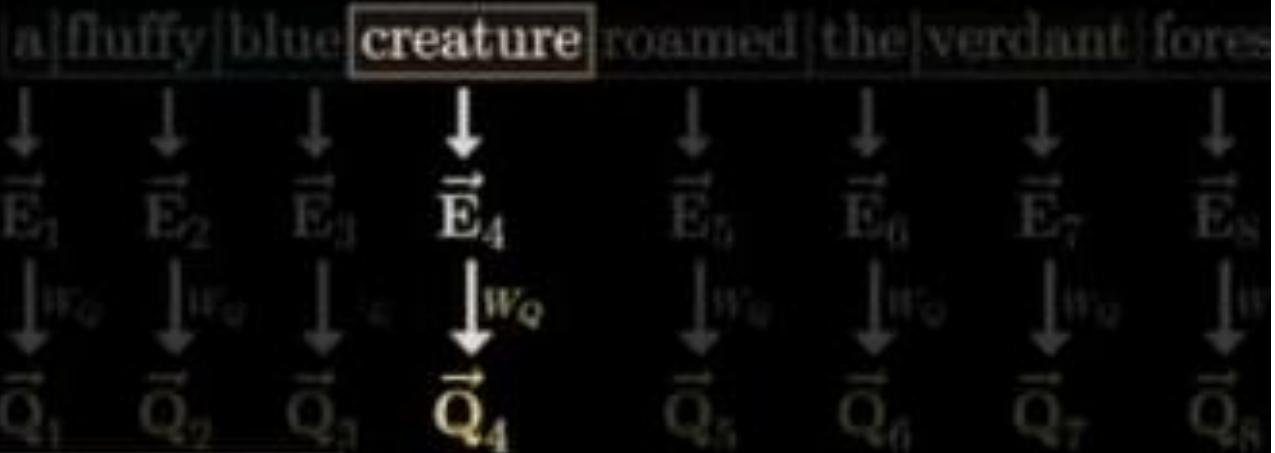


a|fluffy|blue|creature|roamed|the|verdant|forest



Recap

- S



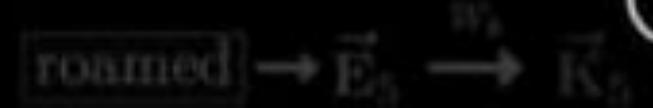
I'm an adjective!
I'm there!



Any adjectives
in front of me?

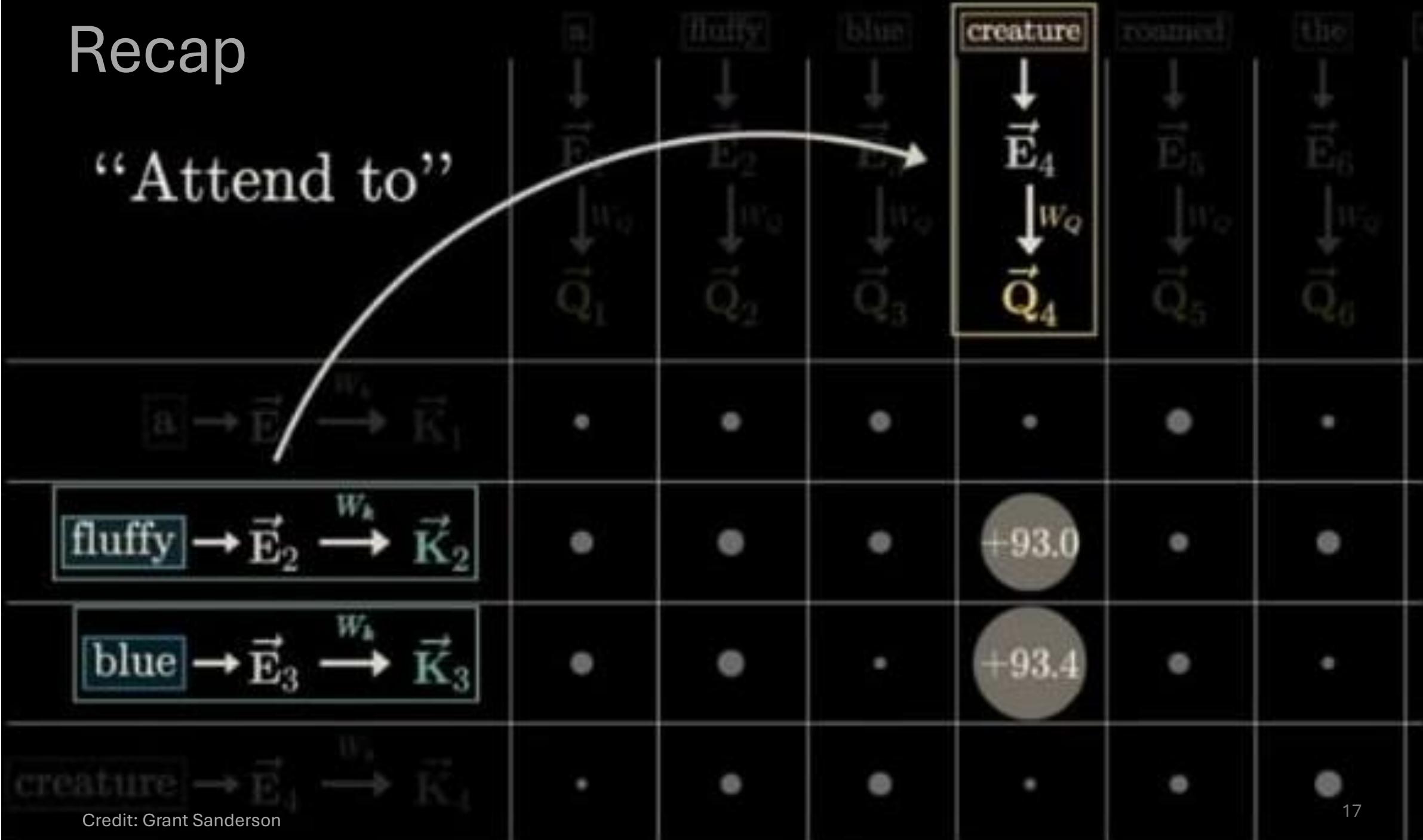


I'm an adjective!
I'm there!



Recap

“Attend to”



Recap

	\downarrow	fluffy	blue	creature	roamed	the	verdant	forest
	\vec{E}_1	\vec{E}_2	\vec{E}_3	\vec{E}_4	\vec{E}_5	\vec{E}_6	\vec{E}_7	\vec{E}_8
	$\downarrow w_0$							
	\vec{Q}_1	\vec{Q}_2	\vec{Q}_3	\vec{Q}_4	\vec{Q}_5	\vec{Q}_6	\vec{Q}_7	\vec{Q}_8
$\blacksquare \rightarrow \vec{E}_1 \xrightarrow{w_1} \vec{K}_1$	$\vec{E}_1 \cdot \vec{Q}_1$	$\vec{E}_1 \cdot \vec{Q}_2$	$\vec{E}_1 \cdot \vec{Q}_3$	$\vec{E}_1 \cdot \vec{Q}_4$	$\vec{E}_1 \cdot \vec{Q}_5$	$\vec{E}_1 \cdot \vec{Q}_6$	$\vec{E}_1 \cdot \vec{Q}_7$	$\vec{E}_1 \cdot \vec{Q}_8$
$\text{fluffy} \rightarrow \vec{E}_2 \xrightarrow{w_1} \vec{K}_2$	$\vec{E}_2 \cdot \vec{Q}_1$	$\vec{E}_2 \cdot \vec{Q}_2$	$\vec{E}_2 \cdot \vec{Q}_3$	$\vec{E}_2 \cdot \vec{Q}_4$	$\vec{E}_2 \cdot \vec{Q}_5$	$\vec{E}_2 \cdot \vec{Q}_6$	$\vec{E}_2 \cdot \vec{Q}_7$	$\vec{E}_2 \cdot \vec{Q}_8$
$\text{blue} \rightarrow \vec{E}_3 \xrightarrow{w_1} \vec{K}_3$	$\vec{E}_3 \cdot \vec{Q}_1$	$\vec{E}_3 \cdot \vec{Q}_2$	$\vec{E}_3 \cdot \vec{Q}_3$	$\vec{E}_3 \cdot \vec{Q}_4$	$\vec{E}_3 \cdot \vec{Q}_5$	$\vec{E}_3 \cdot \vec{Q}_6$	$\vec{E}_3 \cdot \vec{Q}_7$	$\vec{E}_3 \cdot \vec{Q}_8$
$\text{creature} \rightarrow \vec{E}_4 \xrightarrow{w_1} \vec{K}_4$	$\vec{E}_4 \cdot \vec{Q}_1$	$\vec{E}_4 \cdot \vec{Q}_2$	$\vec{E}_4 \cdot \vec{Q}_3$	$\vec{E}_4 \cdot \vec{Q}_4$	$\vec{E}_4 \cdot \vec{Q}_5$	$\vec{E}_4 \cdot \vec{Q}_6$	$\vec{E}_4 \cdot \vec{Q}_7$	$\vec{E}_4 \cdot \vec{Q}_8$
$\text{roamed} \rightarrow \vec{E}_5 \xrightarrow{w_1} \vec{K}_5$	$\vec{E}_5 \cdot \vec{Q}_1$	$\vec{E}_5 \cdot \vec{Q}_2$	$\vec{E}_5 \cdot \vec{Q}_3$	$\vec{E}_5 \cdot \vec{Q}_4$	$\vec{E}_5 \cdot \vec{Q}_5$	$\vec{E}_5 \cdot \vec{Q}_6$	$\vec{E}_5 \cdot \vec{Q}_7$	$\vec{E}_5 \cdot \vec{Q}_8$
$\text{the} \rightarrow \vec{E}_6 \xrightarrow{w_1} \vec{K}_6$	$\vec{E}_6 \cdot \vec{Q}_1$	$\vec{E}_6 \cdot \vec{Q}_2$	$\vec{E}_6 \cdot \vec{Q}_3$	$\vec{E}_6 \cdot \vec{Q}_4$	$\vec{E}_6 \cdot \vec{Q}_5$	$\vec{E}_6 \cdot \vec{Q}_6$	$\vec{E}_6 \cdot \vec{Q}_7$	$\vec{E}_6 \cdot \vec{Q}_8$
$\text{verdant} \rightarrow \vec{E}_7 \xrightarrow{w_1} \vec{K}_7$	$\vec{E}_7 \cdot \vec{Q}_1$	$\vec{E}_7 \cdot \vec{Q}_2$	$\vec{E}_7 \cdot \vec{Q}_3$	$\vec{E}_7 \cdot \vec{Q}_4$	$\vec{E}_7 \cdot \vec{Q}_5$	$\vec{E}_7 \cdot \vec{Q}_6$	$\vec{E}_7 \cdot \vec{Q}_7$	$\vec{E}_7 \cdot \vec{Q}_8$
$\text{forest} \rightarrow \vec{E}_8 \xrightarrow{w_1} \vec{K}_8$	$\vec{E}_8 \cdot \vec{Q}_1$	$\vec{E}_8 \cdot \vec{Q}_2$	$\vec{E}_8 \cdot \vec{Q}_3$	$\vec{E}_8 \cdot \vec{Q}_4$	$\vec{E}_8 \cdot \vec{Q}_5$	$\vec{E}_8 \cdot \vec{Q}_6$	$\vec{E}_8 \cdot \vec{Q}_7$	$\vec{E}_8 \cdot \vec{Q}_8$

Recap

	a	fluffy	blue	creature	roamed	the	verdant	forest
	\vec{E}_1	\vec{E}_2	\vec{E}_3	\vec{E}_4	\vec{E}_5	\vec{E}_6	\vec{E}_7	\vec{E}_8
	w_0	w_0	w_0	w_0	w_0	w_0	w_0	w_0
a $\rightarrow \vec{E}_1 \xrightarrow{w_0} \vec{K}_1$	-0.7	83.7	-26.7	-27.5	5.2	-89.3	-15.2	36.1
fluffy $\rightarrow \vec{E}_2 \xrightarrow{w_0} \vec{K}_2$	73.4	-12.9	5.4	+93.0				
blue $\rightarrow \vec{E}_3 \xrightarrow{w_0} \vec{K}_3$	62.4	8.2	-1.8	+93.4				
creature $\rightarrow \vec{E}_4 \xrightarrow{w_0} \vec{K}_4$	21.5	20.7	56.1	+4.9	32.4	92.3	9.1	-26.1
roamed $\rightarrow \vec{E}_5 \xrightarrow{w_0} \vec{K}_5$	20.1	-81.9	-87.8	-55.4	-0.6	64.7	96.2	-15.9
the $\rightarrow \vec{E}_6 \xrightarrow{w_0} \vec{K}_6$	-67.9	-33.3	-22.8	-31.4	-5.5	-0.6	-4.6	-96.8
verdant $\rightarrow \vec{E}_7 \xrightarrow{w_0} \vec{K}_7$	41.2	50.5	-42.3	-59.8	-70.0	97.9	-0.7	+93.8
forest $\rightarrow \vec{E}_8 \xrightarrow{w_0} \vec{K}_8$	38.9	75.5	91.1	-90.6	-75.6	-93.0	-70.8	-1.7

We want these to act like weights

Recap



Recap

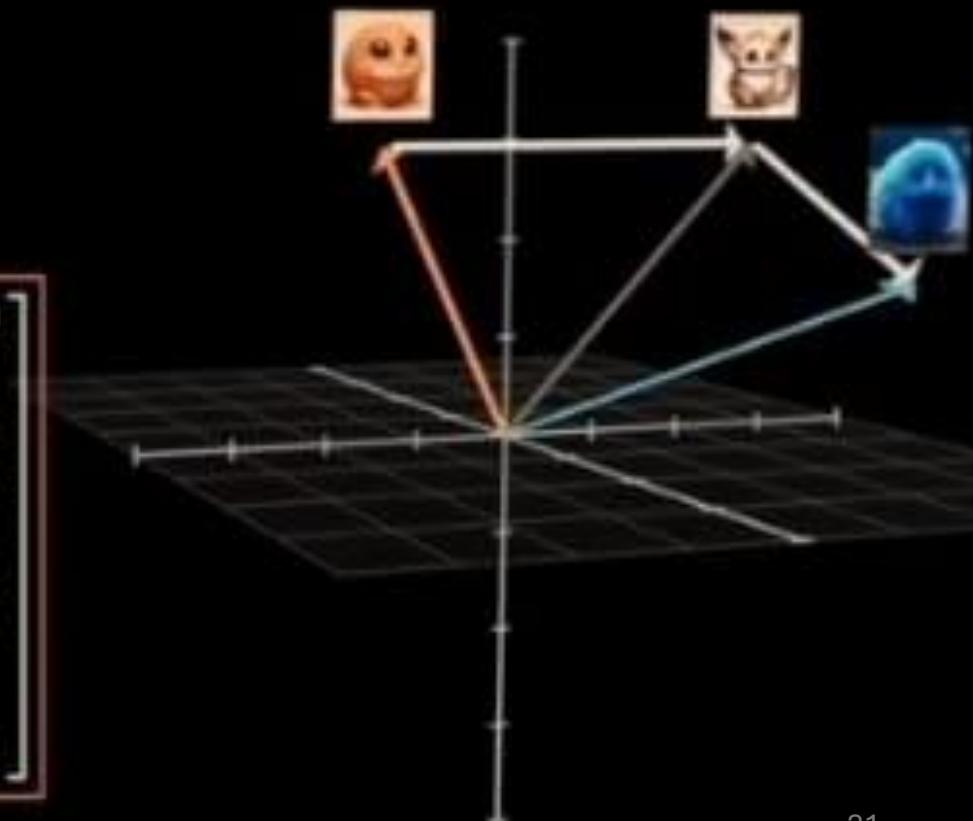
blue
 fluffy
 creature

\downarrow \downarrow \downarrow

$$\begin{bmatrix} +1.0 \\ +4.3 \\ +2.9 \\ +6.9 \\ -1.9 \\ +9.9 \\ \vdots \\ +7.9 \end{bmatrix} \quad
 \begin{bmatrix} +6.2 \\ -3.2 \\ +1.8 \\ +6.6 \\ +1.8 \\ +6.4 \\ \vdots \\ -6.2 \end{bmatrix} \quad
 \begin{bmatrix} -7.6 \\ +1.9 \\ -7.2 \\ +6.8 \\ +1.4 \\ -6.7 \\ \vdots \\ -4.7 \end{bmatrix}$$

W_V

$$\underbrace{\begin{bmatrix} -3.6 & -1.7 & -8.6 & +3.5 & +1.3 & -4.6 & \cdots & -8.0 \\ +1.5 & +8.5 & -3.6 & +3.3 & -7.3 & +4.3 & \cdots & -6.3 \\ +1.7 & -9.5 & +6.5 & -9.8 & +3.5 & -4.6 & \cdots & +9.2 \\ -5.0 & +1.5 & +1.8 & +1.4 & -5.5 & +9.0 & \cdots & +6.9 \\ +3.9 & -4.0 & +6.2 & -2.0 & +7.5 & +1.6 & \cdots & +3.8 \\ +4.5 & +0.0 & +9.0 & +2.9 & -1.5 & +2.1 & \cdots & -3.9 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ +1.5 & +3.0 & +3.0 & -1.4 & +7.9 & -2.6 & \cdots & +7.8 \end{bmatrix}}_{W_V} \times \begin{bmatrix} +1.0 \\ +4.3 \\ +2.9 \\ +6.9 \\ -1.9 \\ +9.9 \\ \vdots \\ +7.9 \end{bmatrix} = \begin{bmatrix} -103.0 \\ +13.1 \\ +12.6 \\ +95.7 \\ +11.2 \\ +14.4 \\ \vdots \\ +62.1 \end{bmatrix}$$



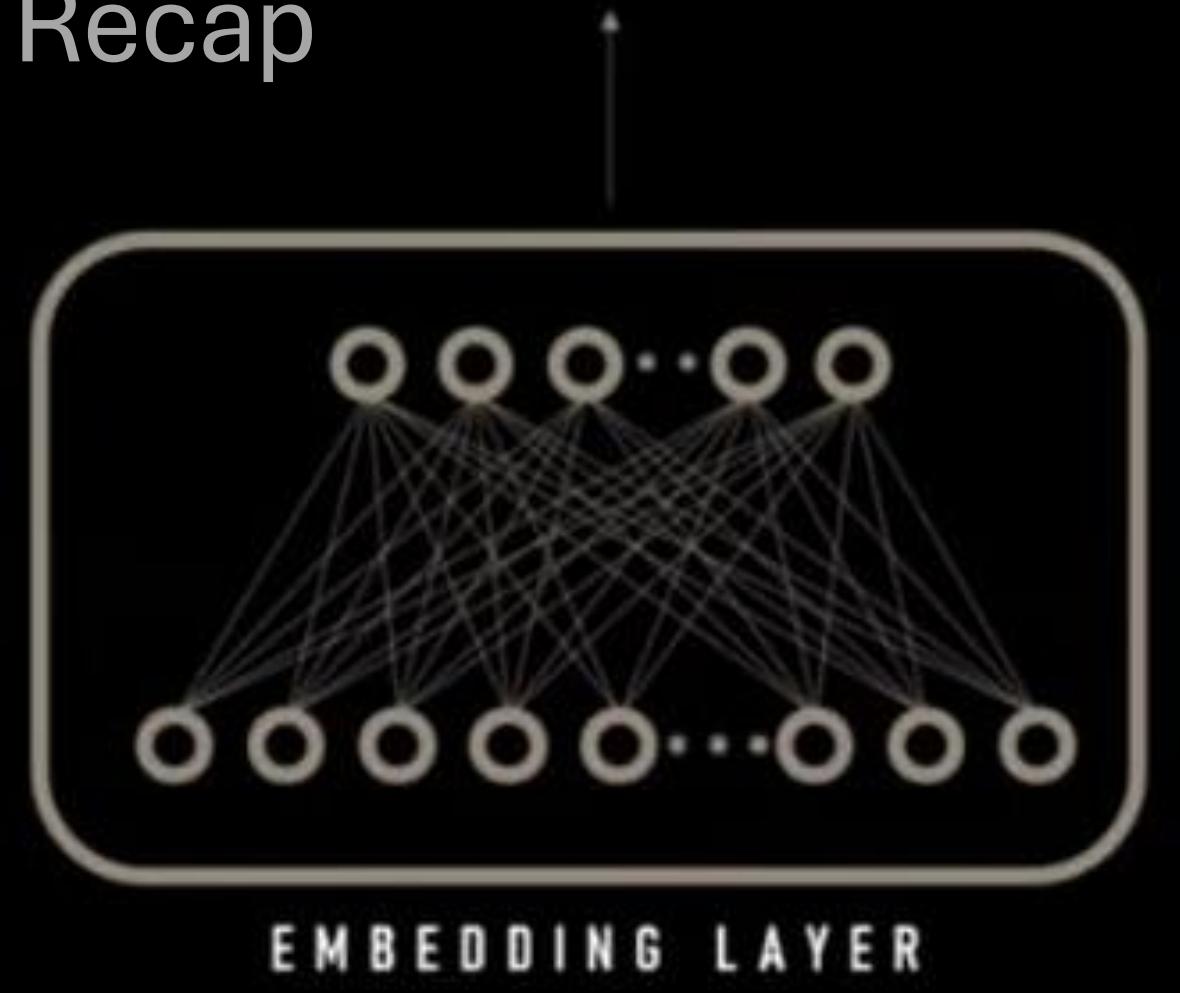
[<SOS>] [He] [deposited] [money] [in] [the] [bank] [] [<EOS>]

Recap

TOKENIZATION

<SOS> He deposited money in the bank. <EOS>

N-Dim Embedding Vector (Output) Recap



<SOS> He deposited money in the bank. <EOS>
<SOS> money in the bank. <EOS> <PAD> <PAD>
<SOS> money. <EOS> <PAD> <PAD> <PAD> <PAD>



[1, 80, 210, 410, 30, 60, 110, 90, 2]
[1, 410, 30, 60, 110, 90, 2, 0, 0]
[1, 410, 90, 2, 0, 0, 0, 0, 0]

(input tokens)

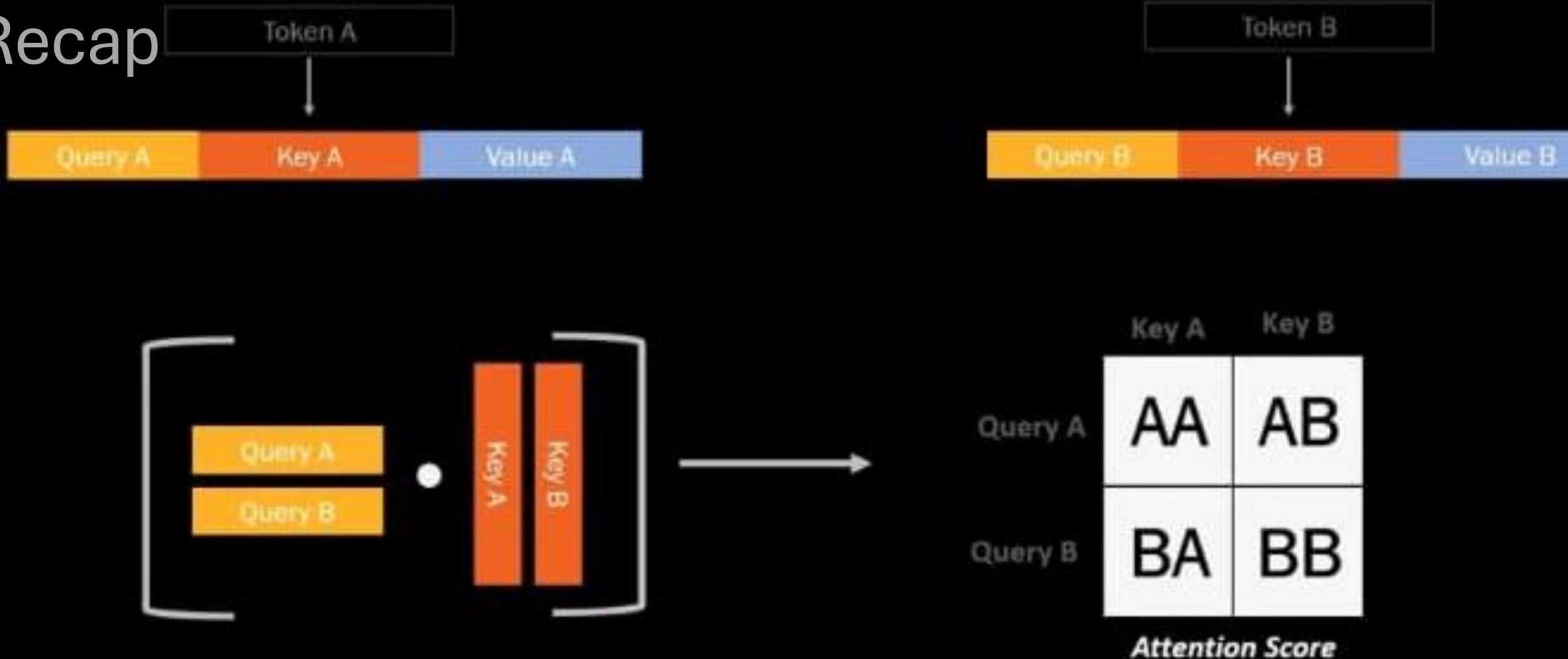
Recap



7 x 7 Attention Scores

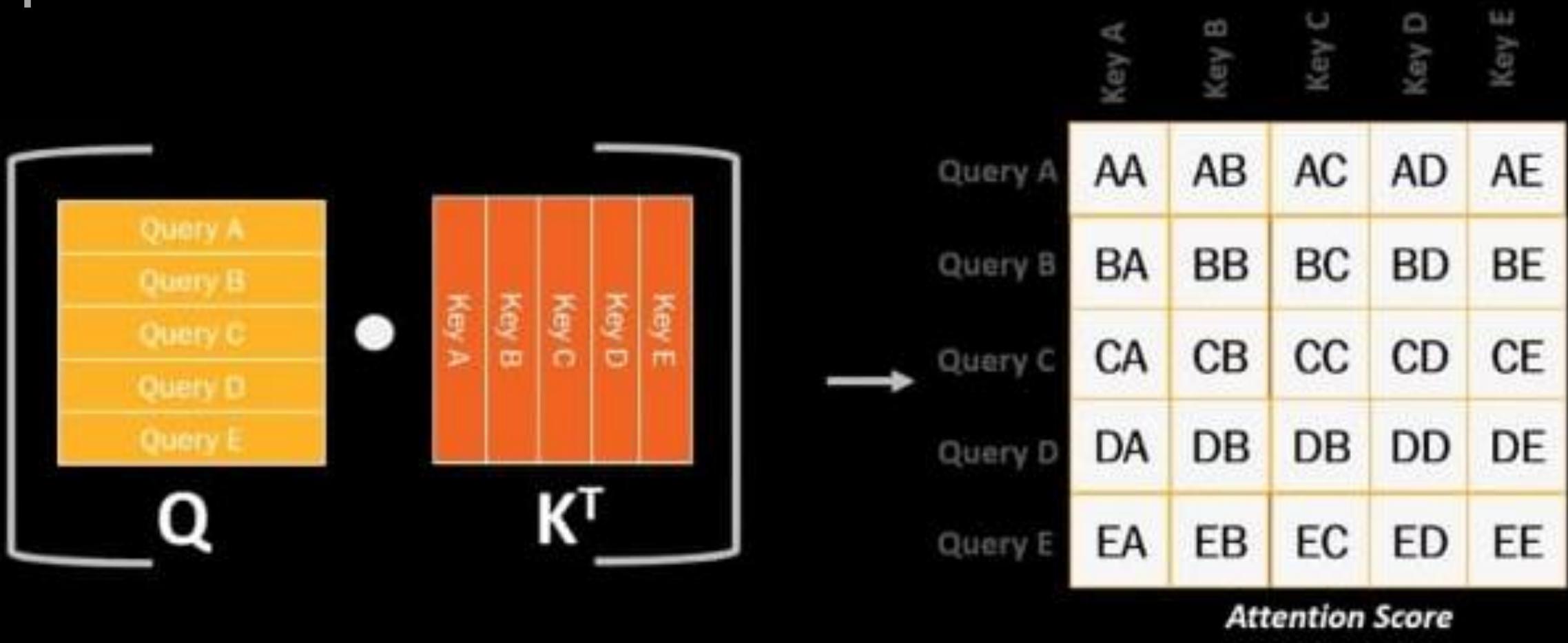
Decoder Masked Multihead Attention

Recap



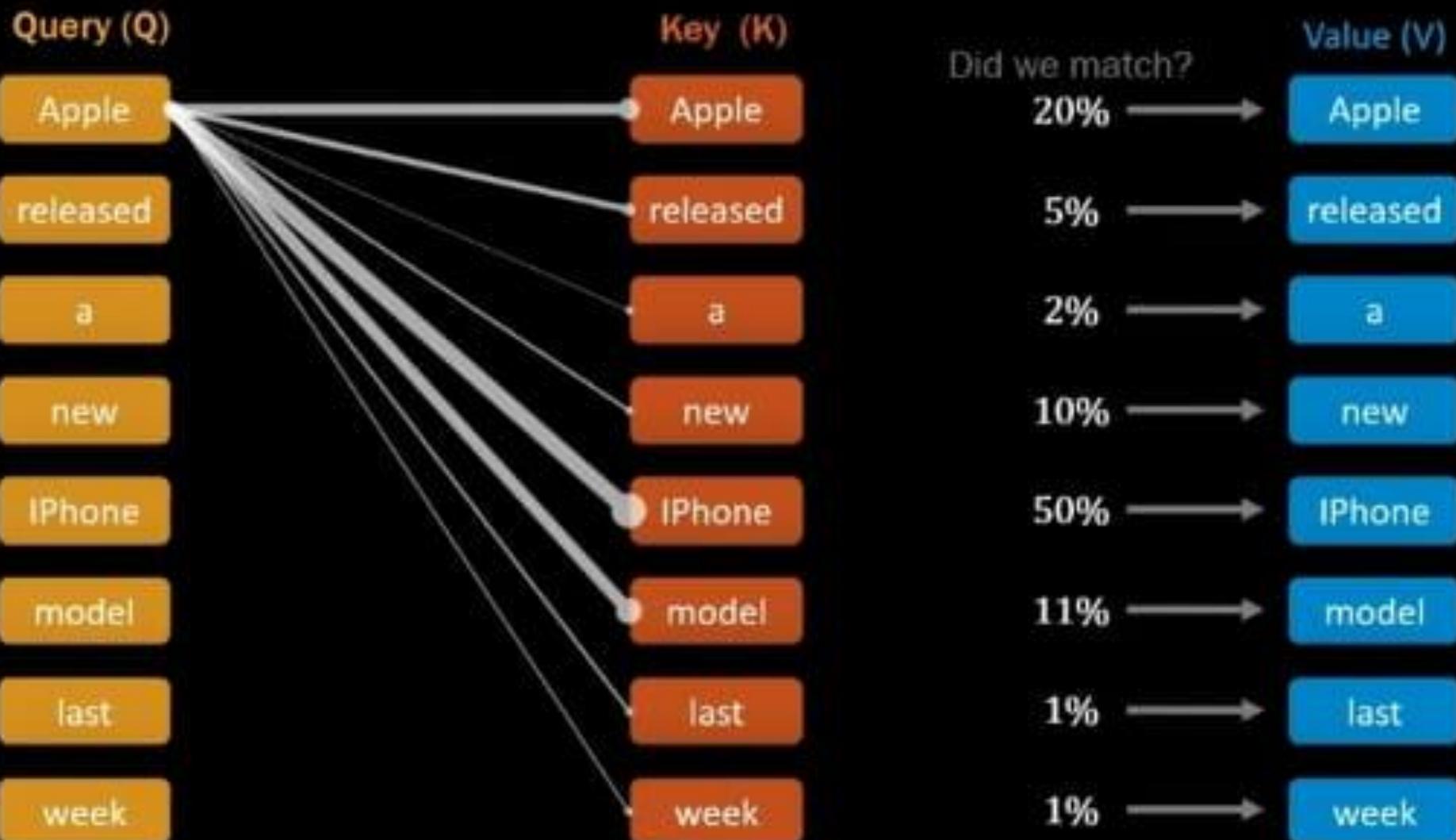


Recap



Recap

“Apple released a new iPhone model last week.”



Hey do you have anything useful
for me?

Here's what I can offer

Here's the actual information if we match

Recap

$$W_Q$$

W_K

↑W_V

↓WV

$$\begin{bmatrix} -0.9 & -0.8 & -0.7 & -0.7 & -0.7 & -0.7 & \dots & -0.7 \\ -0.2 & -0.1 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \\ 1 & 1 & 1 & 1 & 1 & 1 & \dots & 1 \\ -0.2 & -0.1 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \end{bmatrix}$$

$$\begin{bmatrix} -0.1 & -0.8 & \dots & -0.8 \\ -0.8 & -0.1 & \dots & -0.1 \\ \vdots & \vdots & \ddots & \vdots \\ -0.8 & -0.1 & \dots & -0.8 \\ -0.1 & -0.8 & \dots & -0.8 \\ 0.0 & 0.0 & \dots & 0.0 \\ 0.0 & 0.0 & \dots & 0.0 \\ 0.0 & 0.0 & \dots & 0.0 \end{bmatrix} = \begin{bmatrix} -0.7 & -0.3 & -0.8 & -0.7 & -0.3 & -0.4 & \dots & -0.8 \\ -0.3 & -0.9 & -0.7 & -0.8 & -0.8 & -0.3 & \dots & -0.8 \\ -0.8 & -0.7 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \\ -0.7 & -0.8 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \\ -0.3 & -0.8 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \\ -0.4 & -0.3 & -0.1 & -0.1 & -0.1 & -0.1 & \dots & -0.1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ -0.8 & -0.3 & -0.8 & -0.7 & -0.7 & -0.7 & \dots & -0.8 \end{bmatrix}$$

John hit the brakes sharply, they screeched loud

$$\begin{bmatrix} -0.1 & -0.2 & -0.3 & -0.4 & -0.5 & -0.6 & \cdots & -0.9 \\ -0.4 & -0.7 & -0.3 & -0.7 & -0.8 & -0.6 & \cdots & -0.2 \\ 1 & 1 & 1 & 1 & 1 & 1 & \cdots & 1 \\ -0.9 & -0.6 & -0.1 & -0.1 & -0.5 & -0.8 & \cdots & -0.9 \end{bmatrix}$$

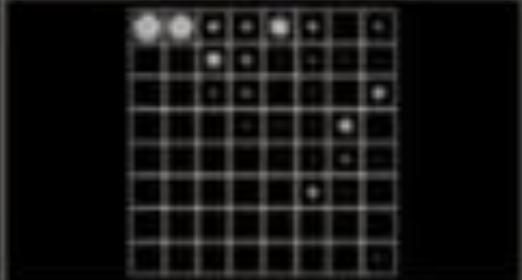
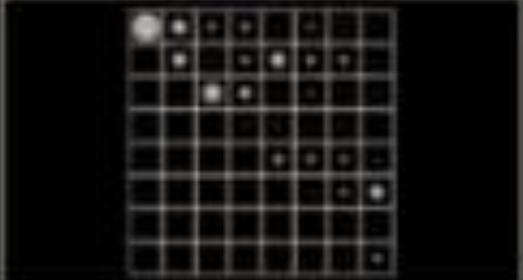
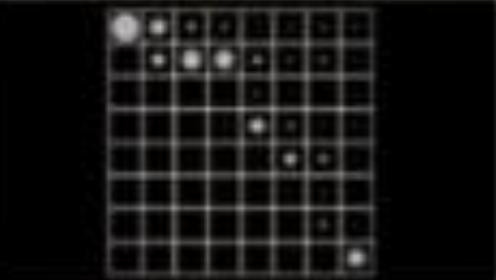
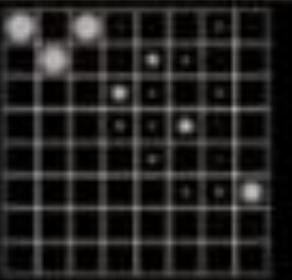
-0.1	-0.2	-0.3	-0.4	-0.5	-0.6	-0.7	-0.8
-0.7	-0.8	-0.9	-0.8	-0.7	-0.6	-0.5	-0.4
-0.1	-0.2	-0.3	-0.4	-0.5	-0.6	-0.7	-0.8
-0.9	-0.8	-0.7	-0.6	-0.5	-0.4	-0.3	-0.2
-0.1	-0.2	-0.3	-0.4	-0.5	-0.6	-0.7	-0.8

$$\begin{bmatrix} -0.1 & -0.8 & -0.9 \\ 0.8 & -0.1 & -0.3 \\ -0.2 & -0.4 & -0.1 \end{bmatrix} = \begin{bmatrix} 0.4 & -0.1 & -0.8 & -0.6 & -0.1 & -0.4 & -0.8 \\ -0.7 & -0.1 & -0.2 & -0.8 & -0.4 & -0.8 & -0.8 \\ -0.1 & -0.1 & -0.1 & -0.2 & -0.1 & -0.2 & -0.1 \end{bmatrix}$$

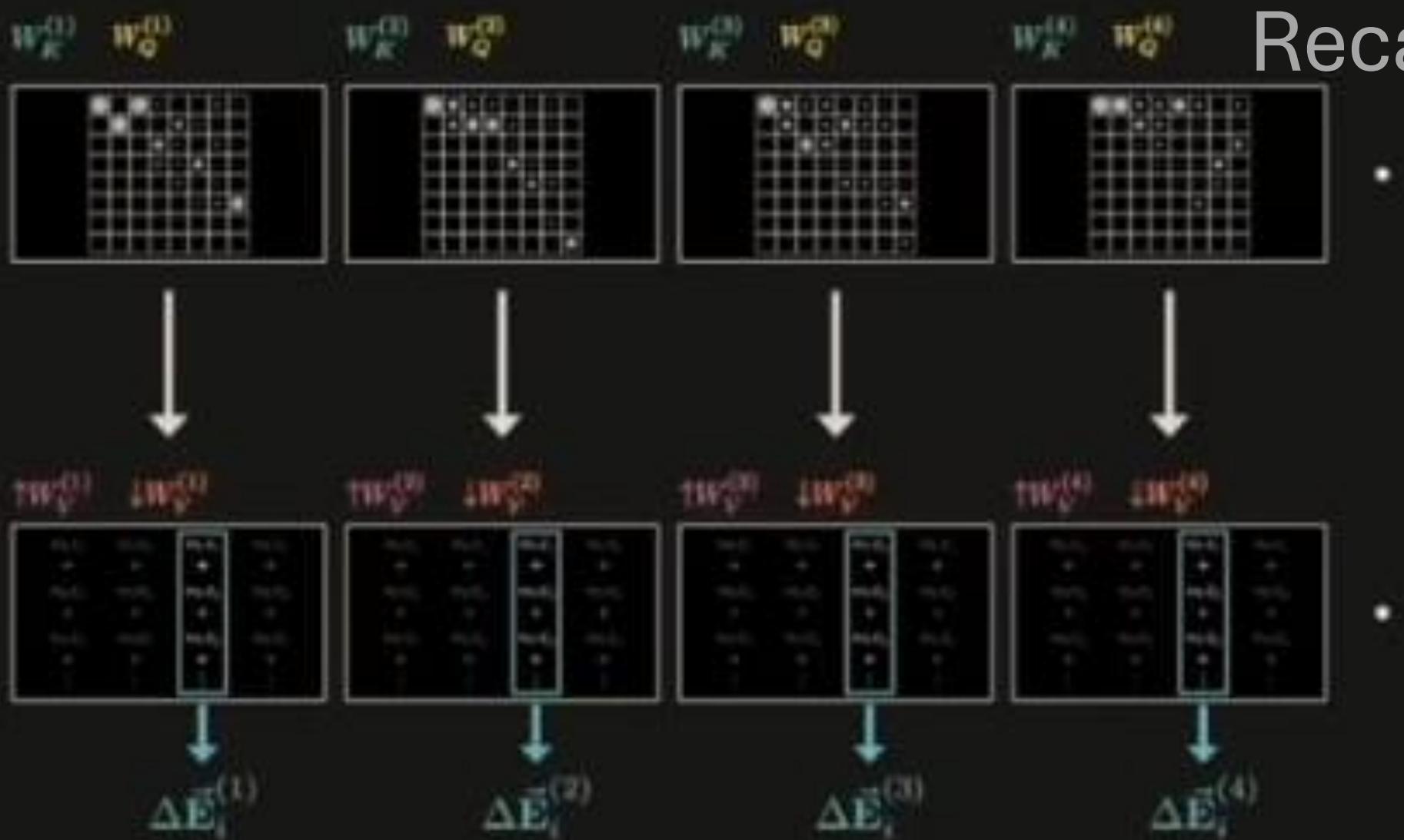
John hit the brakes sharply, they screeched loudly, and he jolted forward.

$W_K^{(1)}$ $W_Q^{(1)}$ $W_K^{(2)}$ $W_Q^{(2)}$ $W_K^{(3)}$ $W_Q^{(3)}$ $W_K^{(4)}$ $W_Q^{(4)}$

Recap

 $\dots \dots \dots$  $\uparrow W_V^{(1)} \quad \downarrow W_V^{(1)}$ $\uparrow W_V^{(2)} \quad \downarrow W_V^{(2)}$ $\uparrow W_V^{(3)} \quad \downarrow W_V^{(3)}$ $\uparrow W_V^{(4)} \quad \downarrow W_V^{(4)}$ $\dots \dots \dots$ $\vec{v}_1^{(1)} \quad \vec{v}_2^{(1)} \quad \vec{v}_3^{(1)} \quad \dots$ $\vec{v}_1^{(2)} \quad \vec{v}_2^{(2)} \quad \vec{v}_3^{(2)} \quad \dots$ $\vec{v}_1^{(3)} \quad \vec{v}_2^{(3)} \quad \vec{v}_3^{(3)} \quad \dots$ $\vec{v}_1^{(4)} \quad \vec{v}_2^{(4)} \quad \vec{v}_3^{(4)} \quad \dots$ $\dots \dots \dots$

Recap



Original
embedding

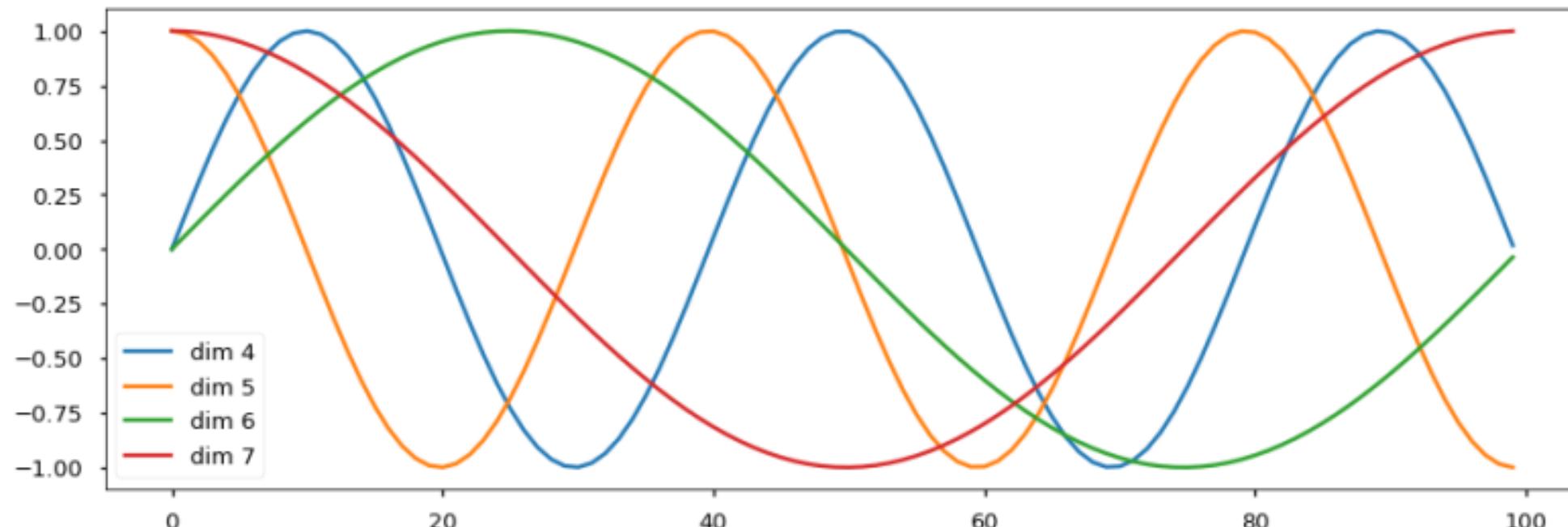
$$\vec{E}_i + \Delta \vec{E}_i^{(1)} + \Delta \vec{E}_i^{(2)} + \Delta \vec{E}_i^{(3)} + \Delta \vec{E}_i^{(4)} + \dots$$

Credit: Grant Sanderson

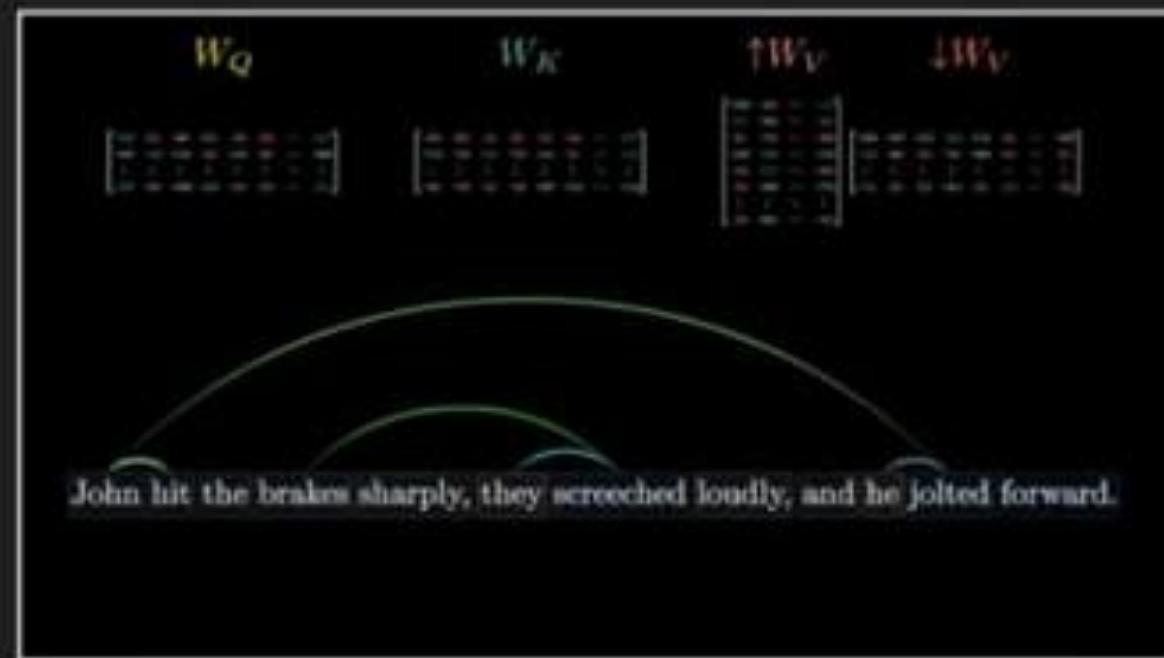
Transformer – Positional Encoding

Recap

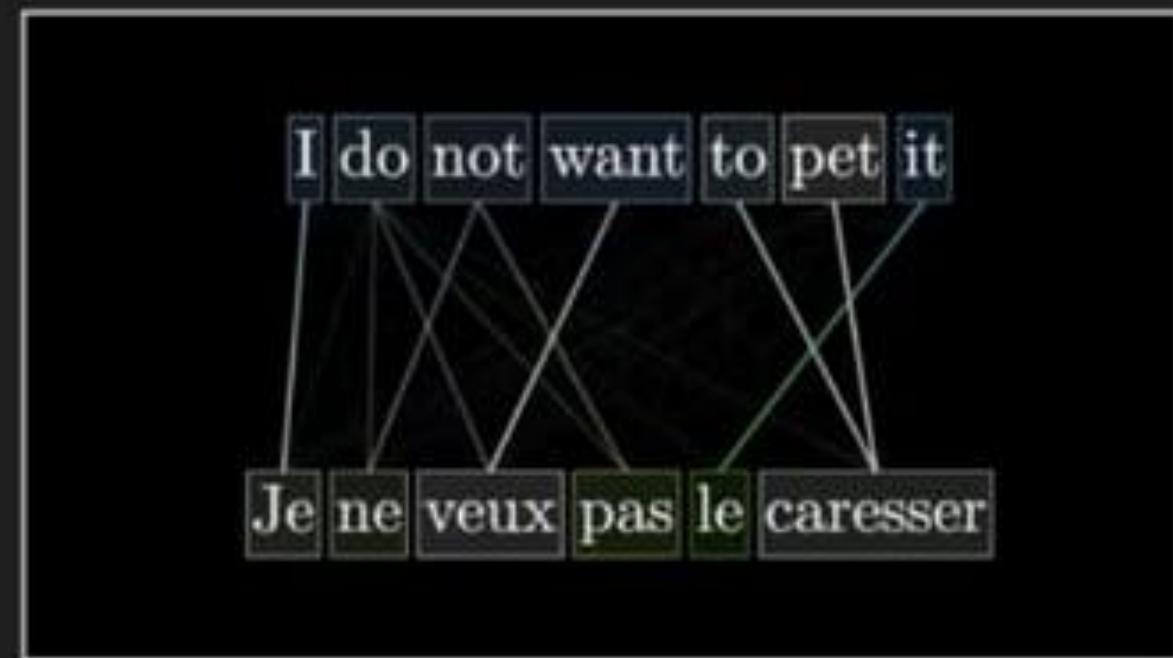
- Transformers add sinusoid curves to the input, before the attention
 - Informs about relative position inside the sequence
 - Removes need for explicit recurrence patterns



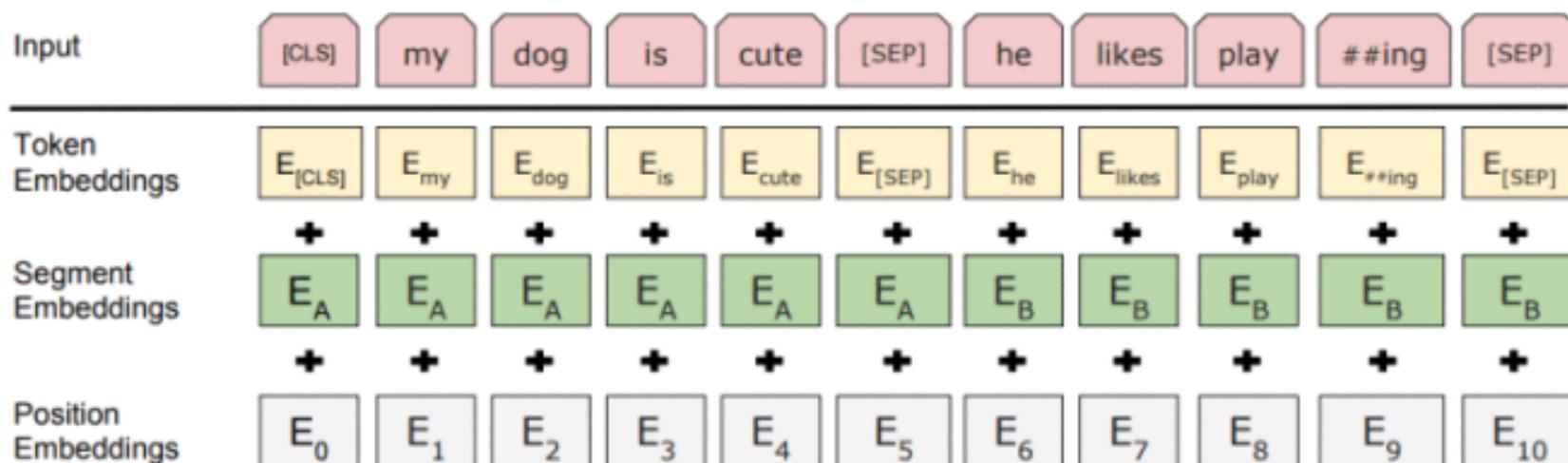
Self-attention



Cross-attention

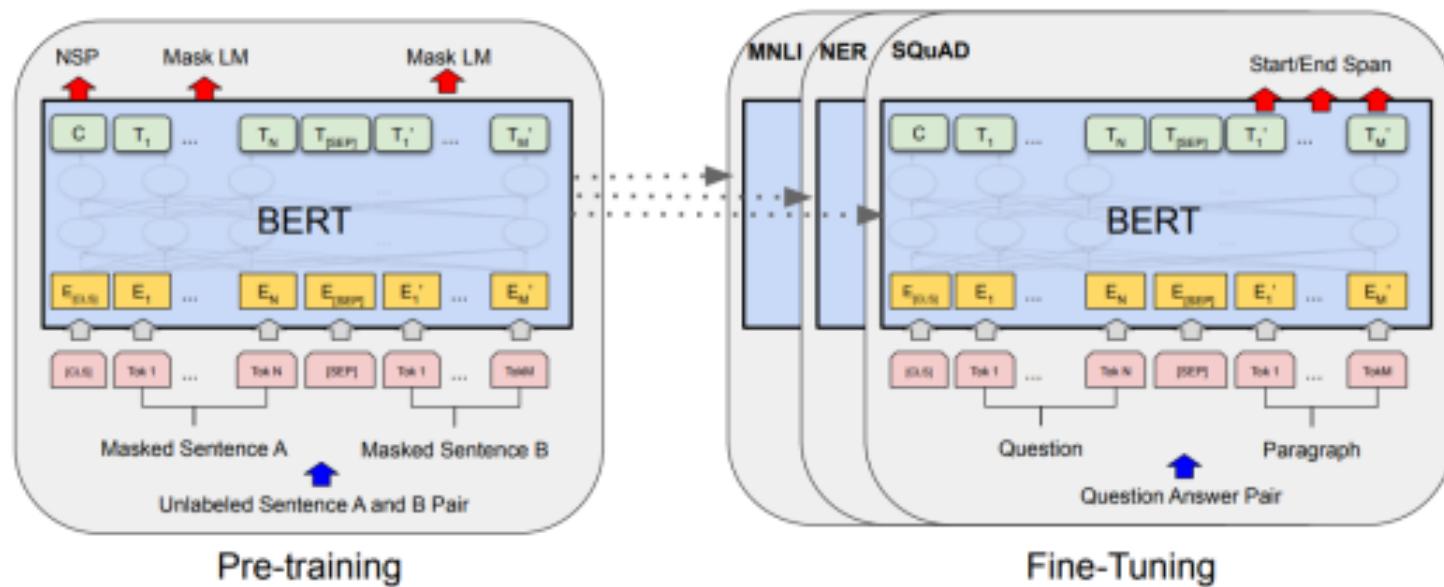


- Either one or two sentences, always prepended with [CLS]
 - BERT adds trained position embeddings & sequence embeddings



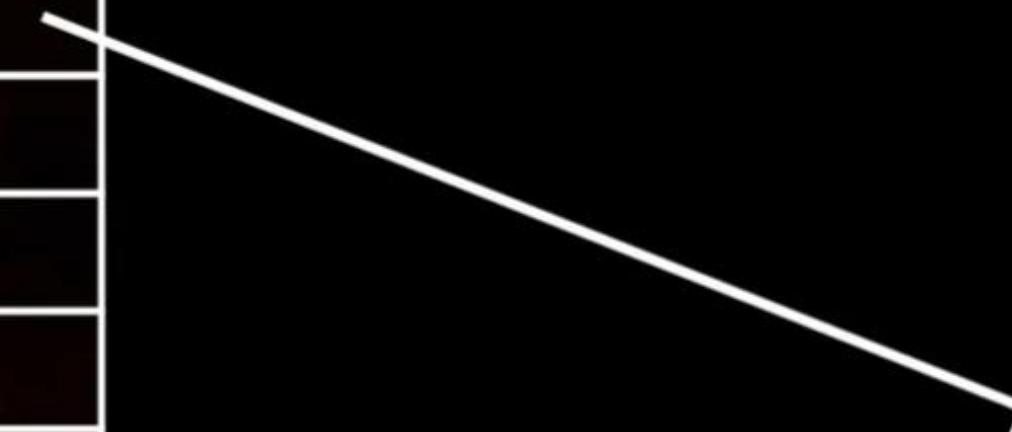
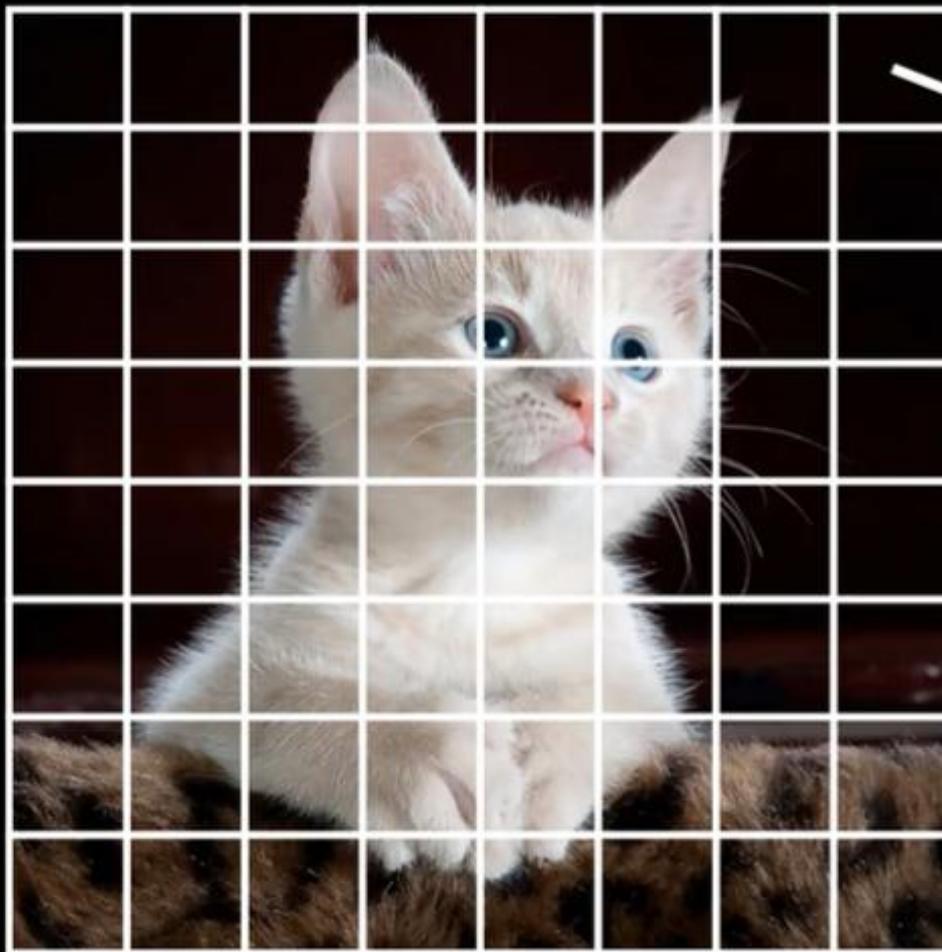
- Model itself is quite simple: n Layers of stacked Transformers
 - Using LayerNorm, GeLU activations (like ReLU, but with a grace swing under 0)
 - Task specific heads on top to pool [CLS] or individual token representations
 - Every Transformer layer receives as input the output of the previous one
- The [CLS] token itself is only special because we train it to be
 - No mechanism inside the model that differentiates it from other tokens
- Novel contributions center around pre-training & workflow

- Someone with lots of compute or time pre-trains a large model
 - BERT uses Masked Language Modelling [MASK] and Next Sentence Prediction [CLS]
- We download it and fine-tune on our task



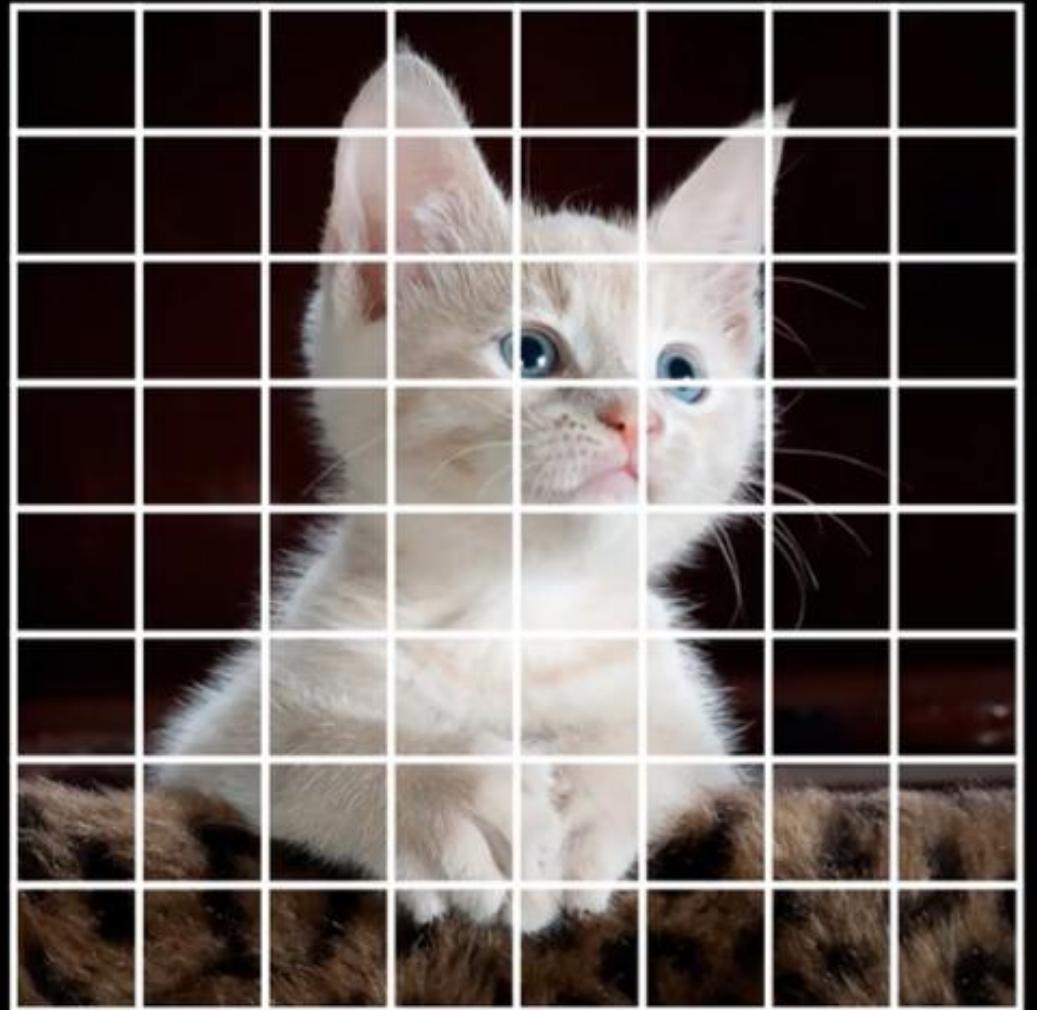
- Same as with Transformer variations, there are now many BERT variants
 - For many languages
 - Domains like biomedical publications
 - Different architectures, but similar workflow:
Roberta, Transformer-XL, XLNet, Longformer ...
- Main themes for adapted architectures:
 - Bigger
 - More efficient
 - Allowing for longer sequences (BERT is capped at 512 tokens in total)

Vision Transformer(ViT)



102	230	225	96	179	61	234	203
92	3	98	243	14	149	245	46
106	244	99	187	71	212	153	199
188	174	65	153	20	44	208	152
102	214	240	39	121	24	34	114
210	65	239	39	214	244	151	25
74	145	222	14	202	85	145	117
87	184	189	221	116	237	109	85

Vision Transformer(ViT)

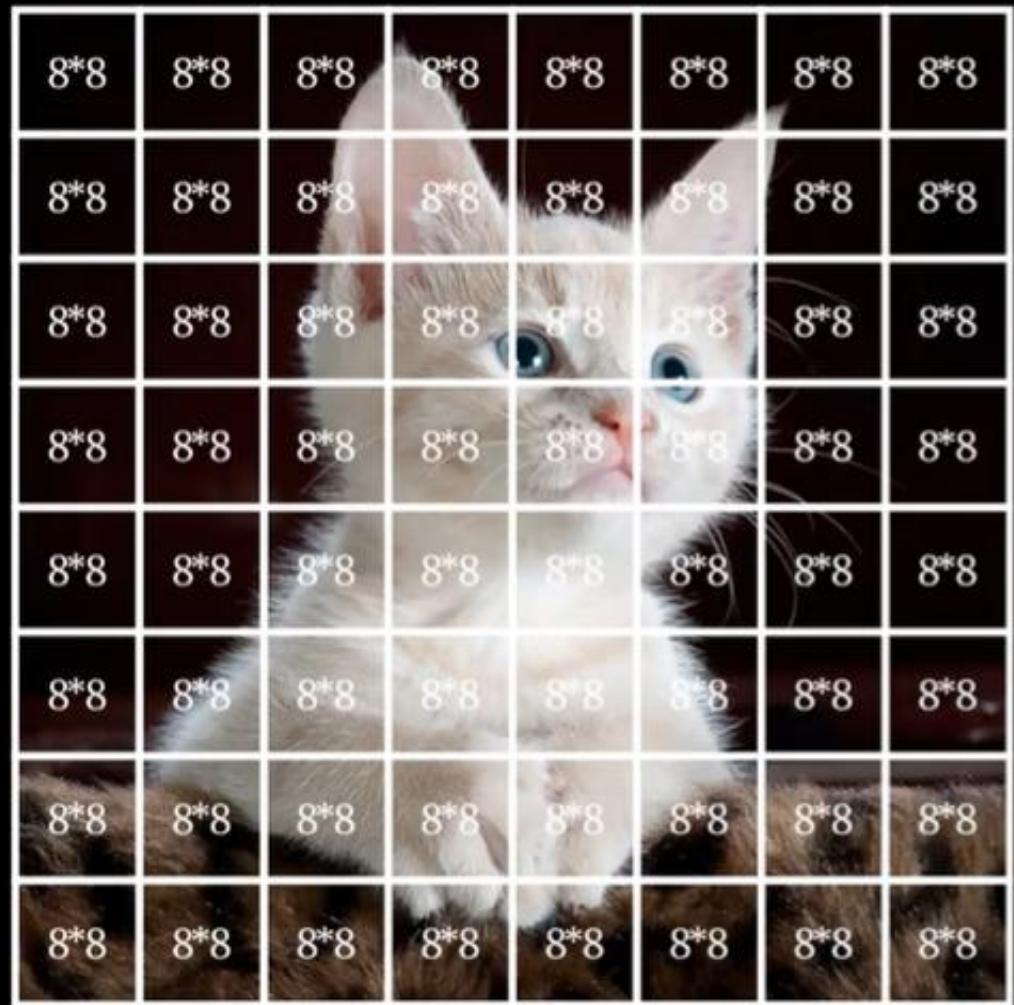


0.39	0.60	0.57	0.24	0.51	0.65	0.27	0.99
0.00	0.97	0.83	0.96	0.62	0.94	0.36	0.04
0.75	0.18	0.00	0.60	0.63	0.62	0.89	0.30
0.08	0.53	0.81	0.95	0.35	0.02	0.55	0.63
0.23	0.61	0.36	0.87	1.00	0.04	0.79	0.04
0.81	0.36	0.59	0.00	0.74	0.09	0.00	0.55
0.20	0.20	0.91	0.87	0.96	0.98	0.66	0.77
0.51	0.04	0.15	0.06	0.53	0.68	0.65	0.92

0.67	0.40	0.14	0.27	0.63	0.58	0.02	0.44
0.01	0.34	0.72	0.10	0.99	0.15	0.21	0.20
0.08	0.73	1.00	0.95	0.27	0.80	0.30	0.58
0.44	0.99	0.00	0.58	0.42	0.19	0.29	0.34
0.39	1.00	0.40	0.71	0.28	0.86	0.29	0.28
0.89	0.05	0.23	0.76	0.98	0.68	0.79	0.31
0.82	0.42	0.38	0.64	0.02	0.25	0.59	0.49
0.17	0.89	0.86	0.53	0.64	0.08	0.17	0.33

0.89	0.92	0.59	0.71	0.63	0.02	0.20	0.05
0.30	0.91	0.92	0.84	0.27	0.86	0.51	0.00
0.55	0.58	0.08	0.18	0.07	0.82	0.22	0.61
0.34	0.09	0.92	0.43	0.91	1.00	0.85	0.53
0.64	0.29	0.66	0.14	0.96	0.99	0.73	0.98
0.79	0.91	0.74	0.45	0.19	0.29	0.74	0.62
0.11	0.08	0.21	0.51	0.71	0.48	0.97	0.47
0.89	0.05	0.20	0.61	0.16	0.47	0.28	0.46

Vision Transformer(ViT)



Patch size = $8*8 = 64$

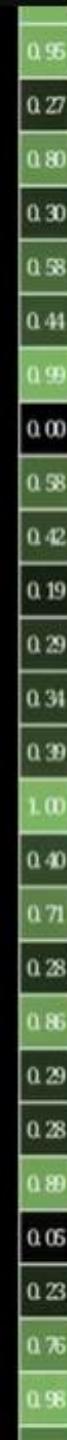
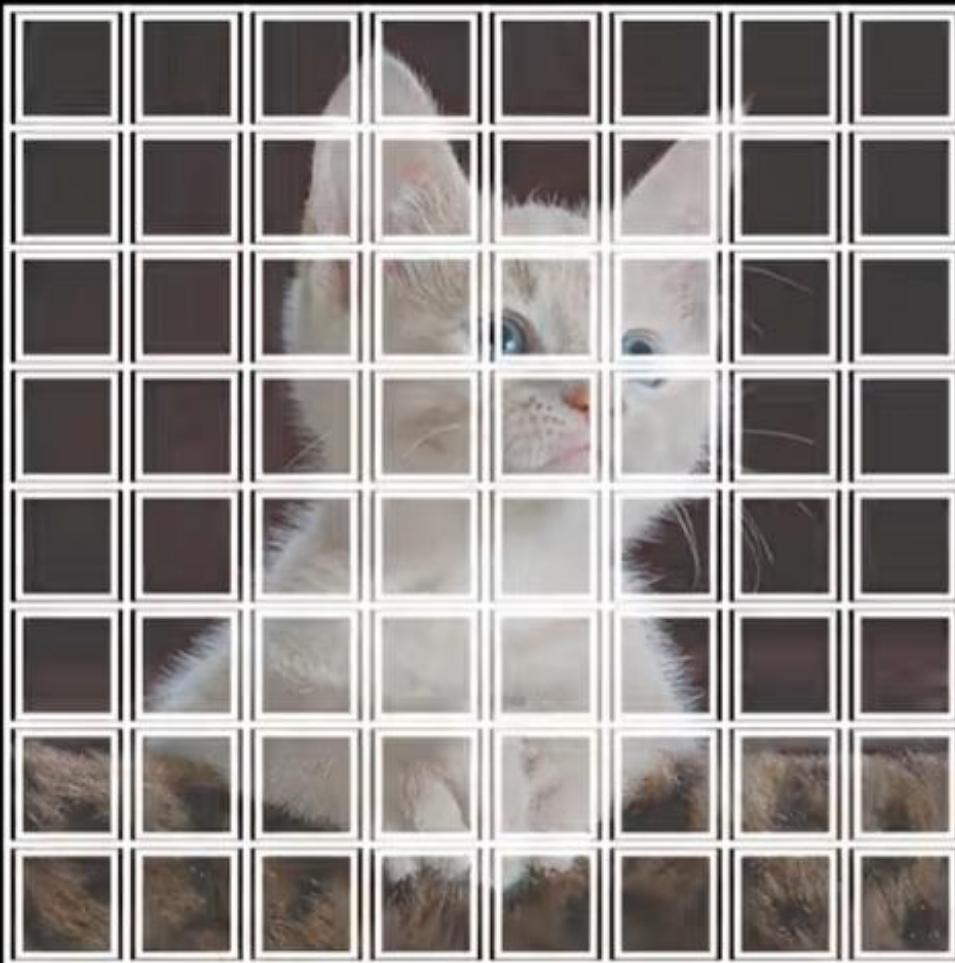
Total patches = $8*8 = 64$

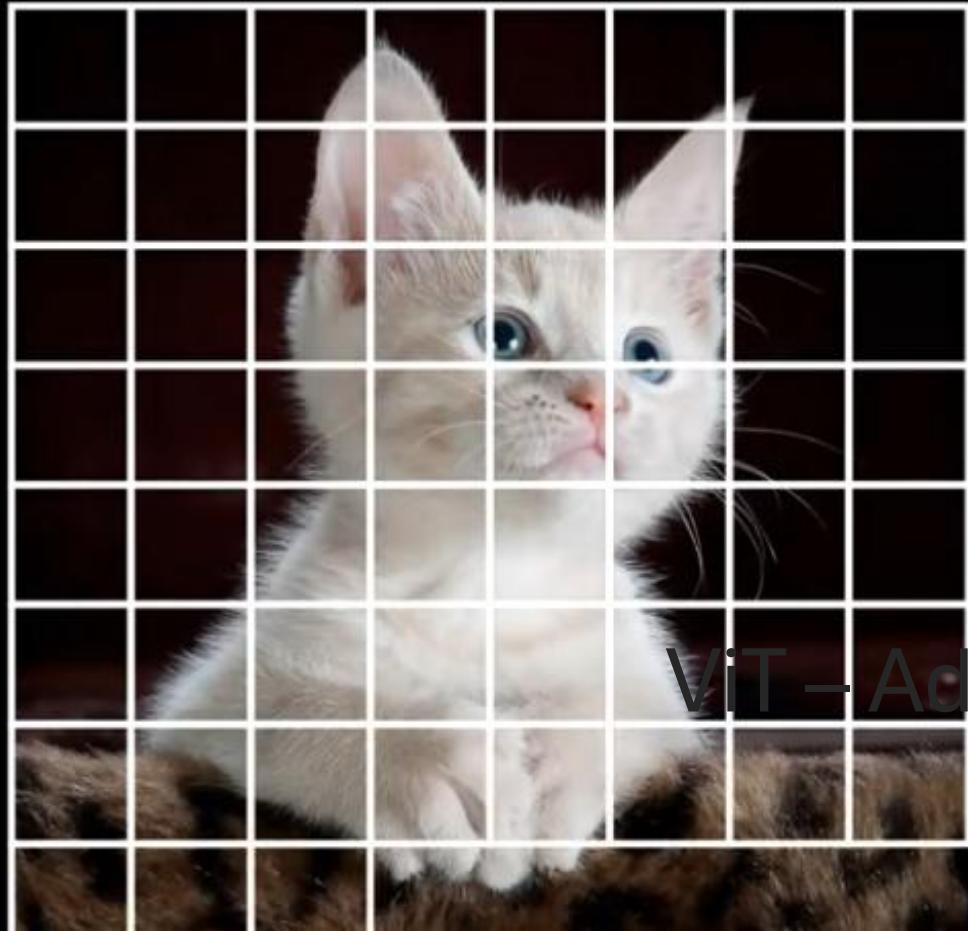
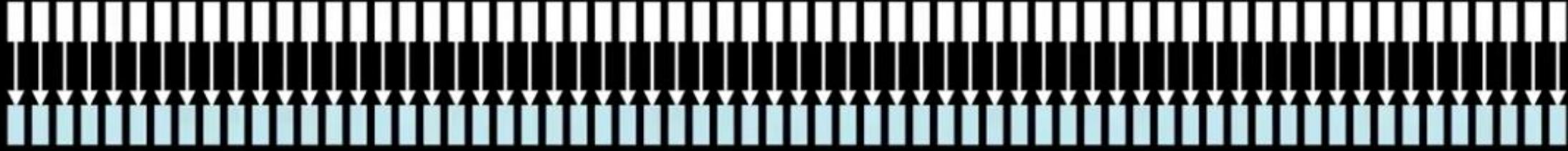
0.39	0.60	0.57	0.24	0.51	0.65	0.27	0.99
0.00	0.97	0.83	0.96	0.62	0.94	0.36	0.04
0.75	0.18	0.00	0.60	0.63	0.62	0.89	0.30
0.08	0.53	0.81	0.95	0.35	0.02	0.55	0.63
0.23	0.61	0.36	0.87	1.00	0.04	0.79	0.04
0.81	0.36	0.99	0.00	0.74	0.09	0.00	0.55
0.20	0.20	0.91	0.87	0.95	0.98	0.66	0.77
0.51	0.04	0.15	0.06	0.53	0.68	0.65	0.92

0.67	0.40	0.14	0.27	0.63	0.38	0.02	0.44
0.01	0.34	0.72	0.10	0.99	0.15	0.21	0.20
0.08	0.73	1.00	0.95	0.27	0.80	0.30	0.58
0.44	0.99	0.00	0.58	0.42	0.19	0.29	0.34
0.39	1.00	0.40	0.71	0.28	0.86	0.29	0.28
0.89	0.05	0.23	0.76	0.98	0.68	0.79	0.31
0.82	0.42	0.38	0.64	0.02	0.25	0.59	0.49
0.17	0.89	0.86	0.53	0.61	0.08	0.17	0.33

0.89	0.92	0.59	0.71	0.63	0.02	0.20	0.05
0.30	0.91	0.92	0.84	0.27	0.86	0.51	0.00
0.55	0.58	0.08	0.18	0.07	0.82	0.22	0.61
0.34	0.09	0.92	0.43	0.91	1.00	0.85	0.53
0.64	0.29	0.66	0.14	0.95	0.99	0.73	0.98
0.79	0.91	0.74	0.45	0.19	0.29	0.74	0.62
0.11	0.03	0.21	0.51	0.71	0.48	0.97	0.47
0.89	0.05	0.20	0.61	0.16	0.47	0.28	0.46

Vision Transformer(ViT)



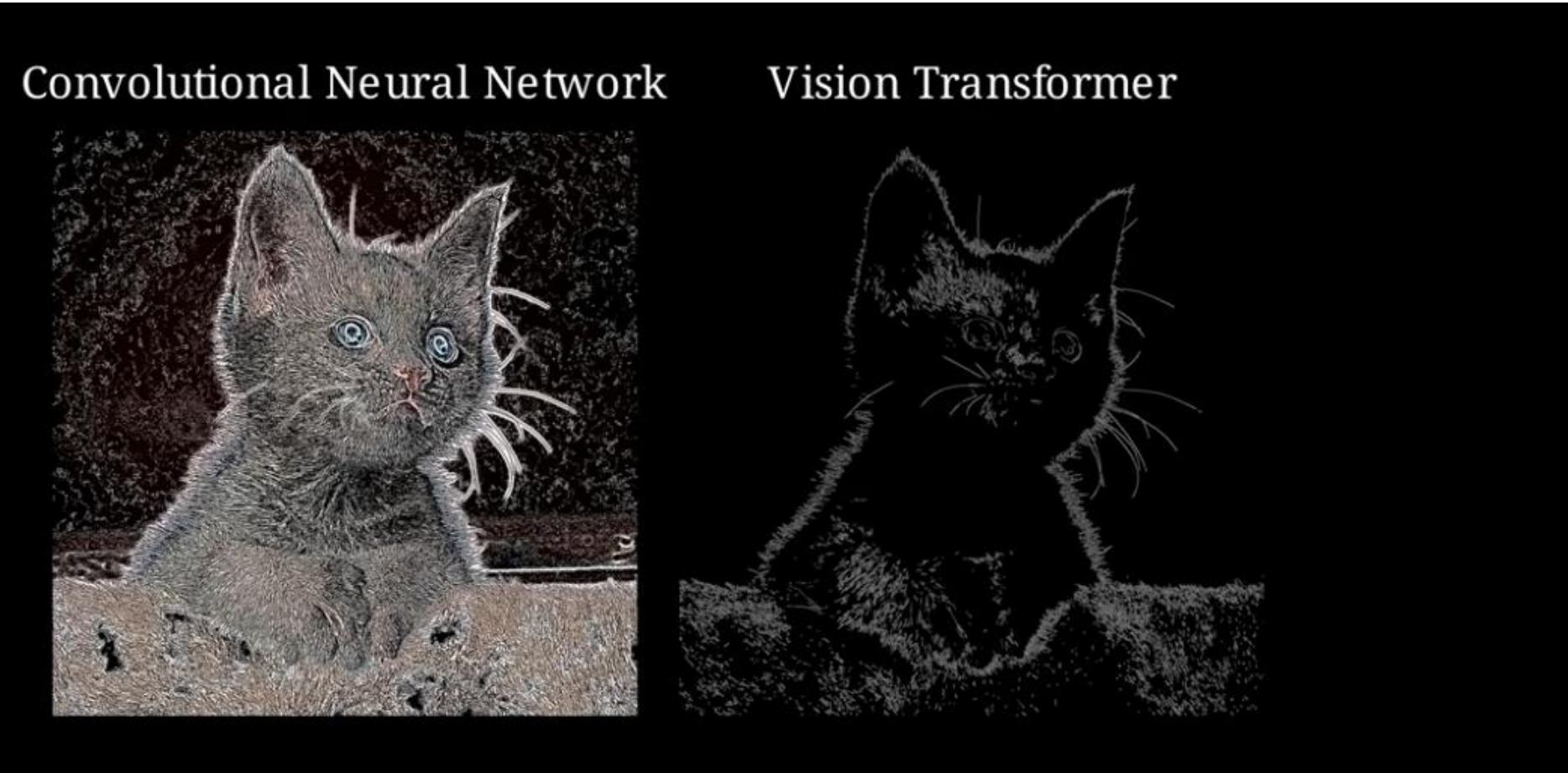


ViT – Add positional encoding

The dog chased the cat around the yard.
The cat chased the dog around the yard.

A diagram illustrating the effect of swapping words in a sentence. Two rows of words are shown, each enclosed in colored boxes: blue, red, orange, blue, green, purple, blue, orange. A dashed red 'X' is drawn through the second row, indicating that the order of words changes the meaning of the sentence.

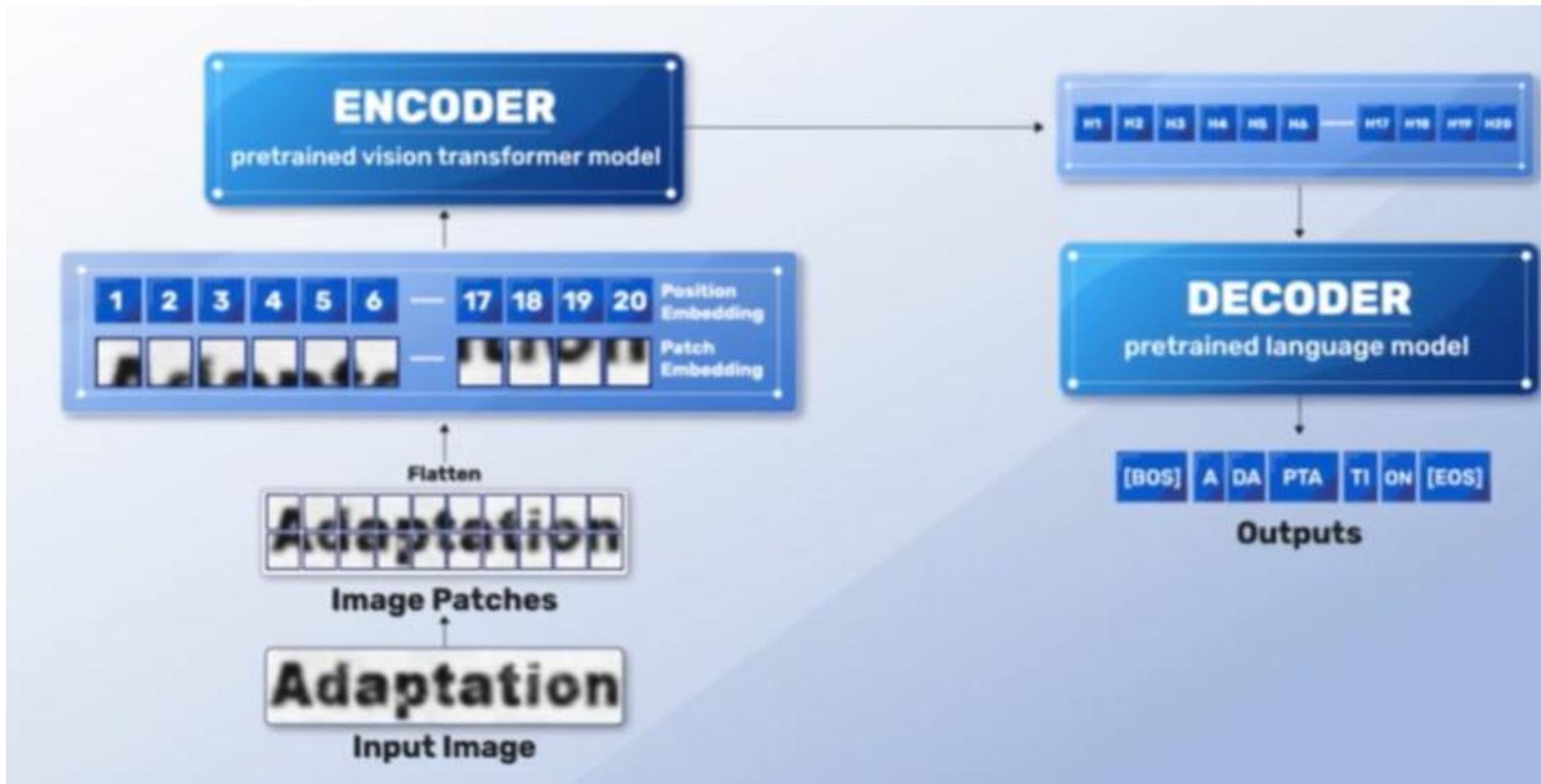
- CNN focuses on structure
- ViT focuses more on general features



OCR



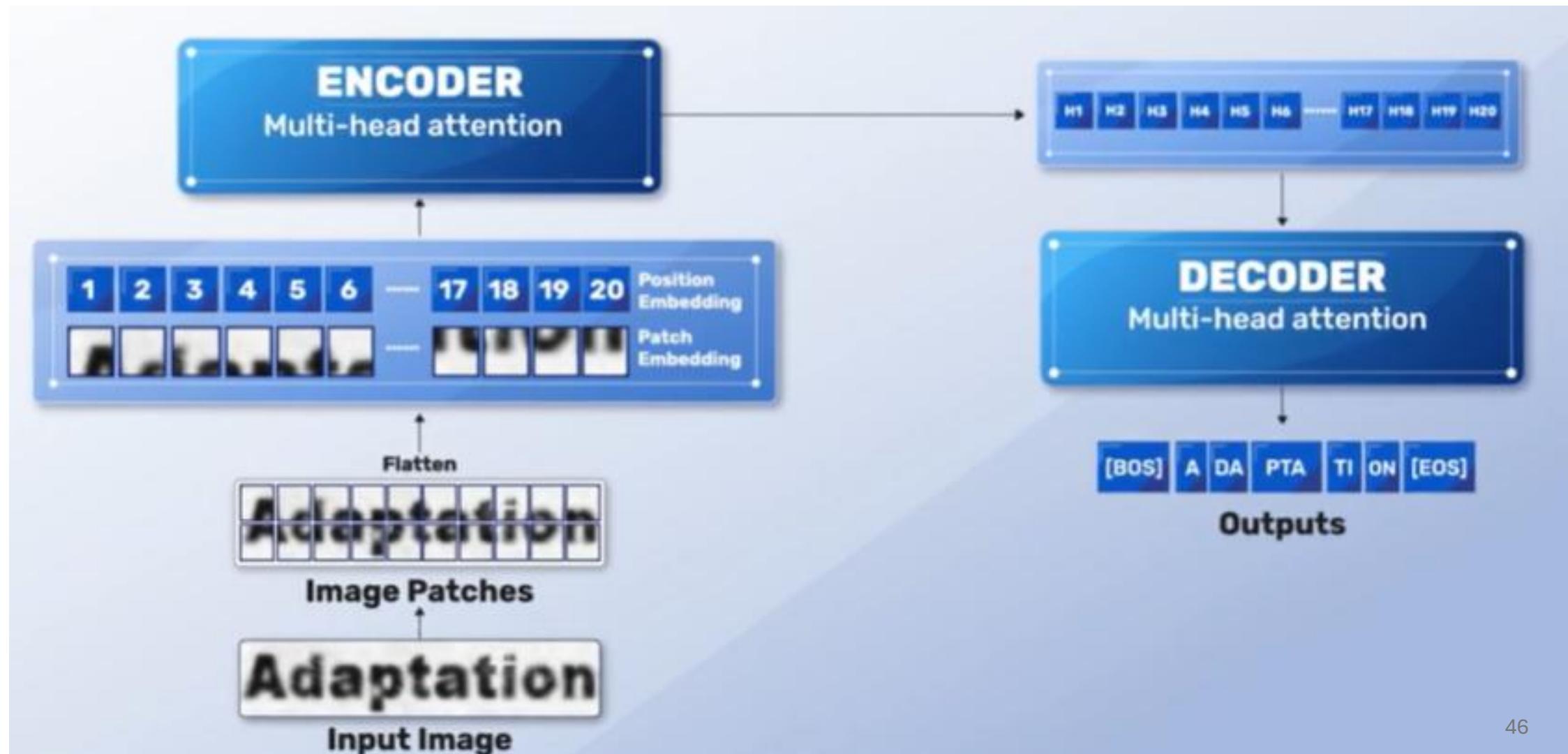
TrOCR - OCR Transformers



TrOCR - OCR Transformers

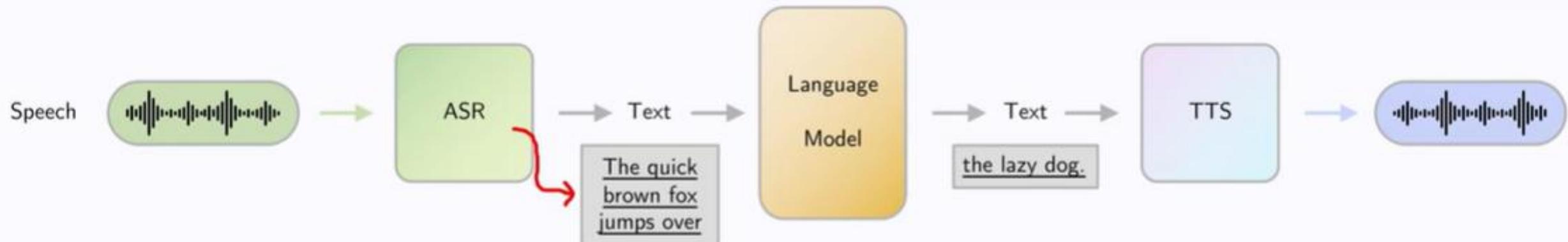


TrOCR - OCR Transformers





Cascading Models: Limitations



1 - Loss of information



"I can do that!"
"I can do that?!"

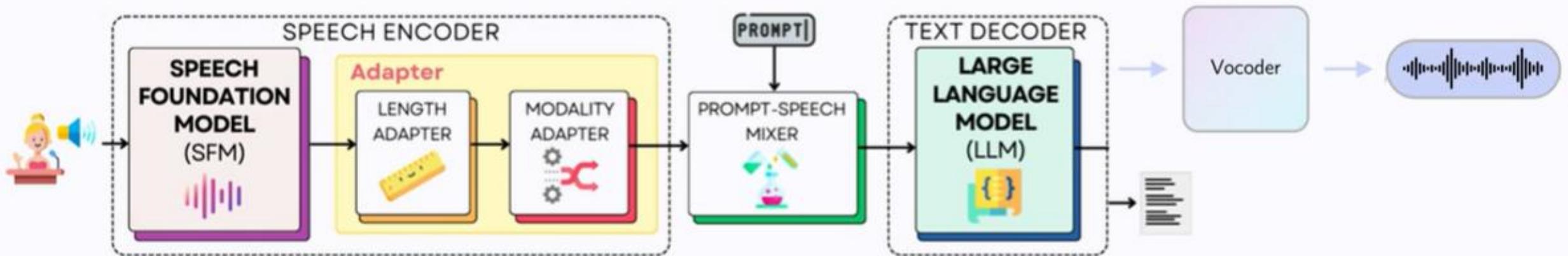
2 - Error propagation

I like beaches

I like peaches

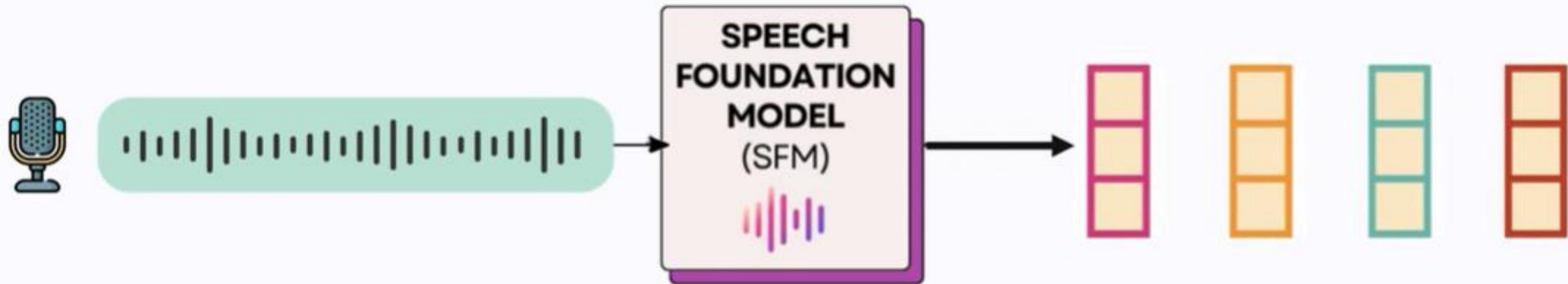
3 - High latency

Components of a Speech LLM



1. Speech Encoder (instead of speech-to-text)
2. Large language model
3. Vocoder (instead of text-to-speech)

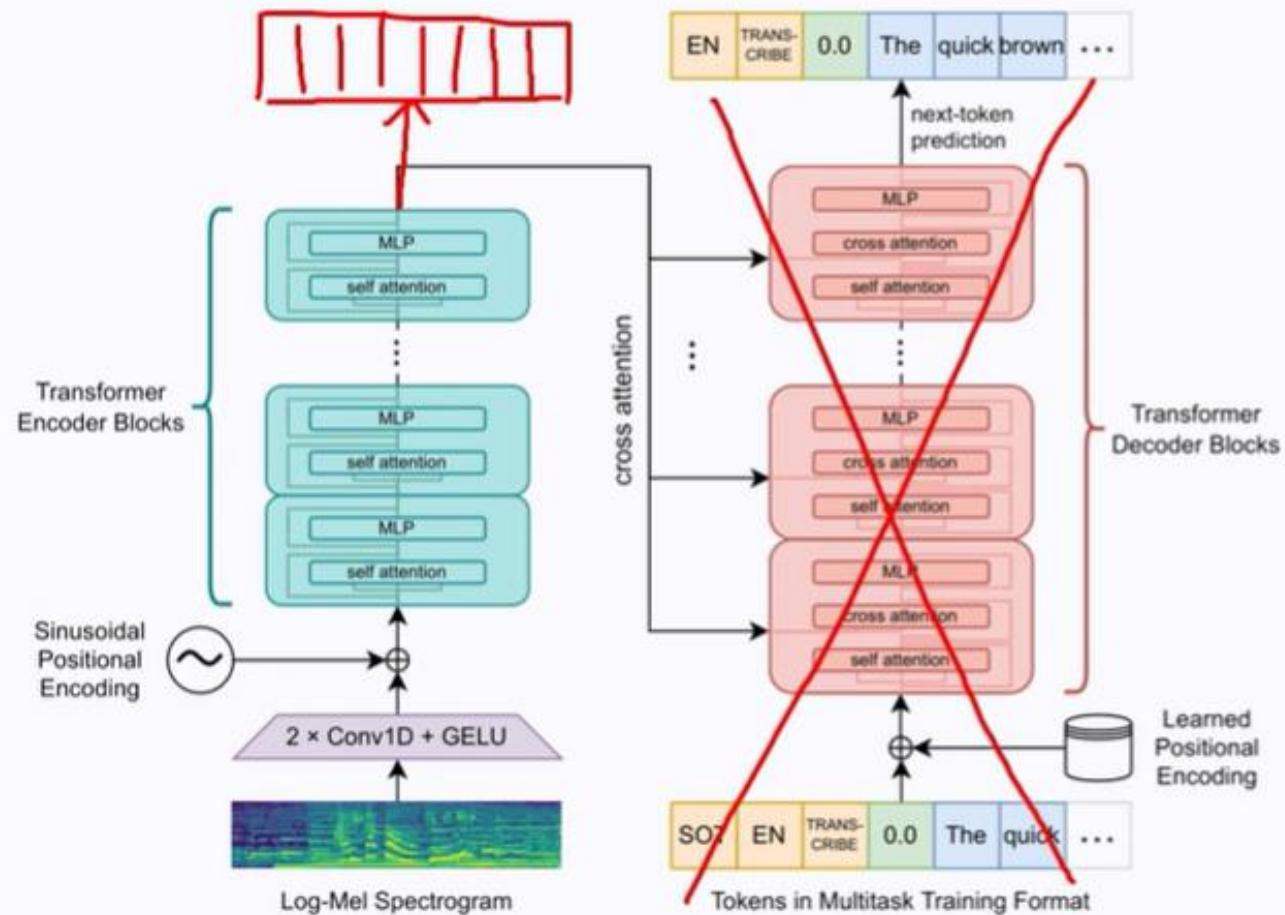
Speech Encoder



Semantic embeddings: capture meaning and content

Acoustic embeddings: capture intonation, sound quality

Speech Encoder

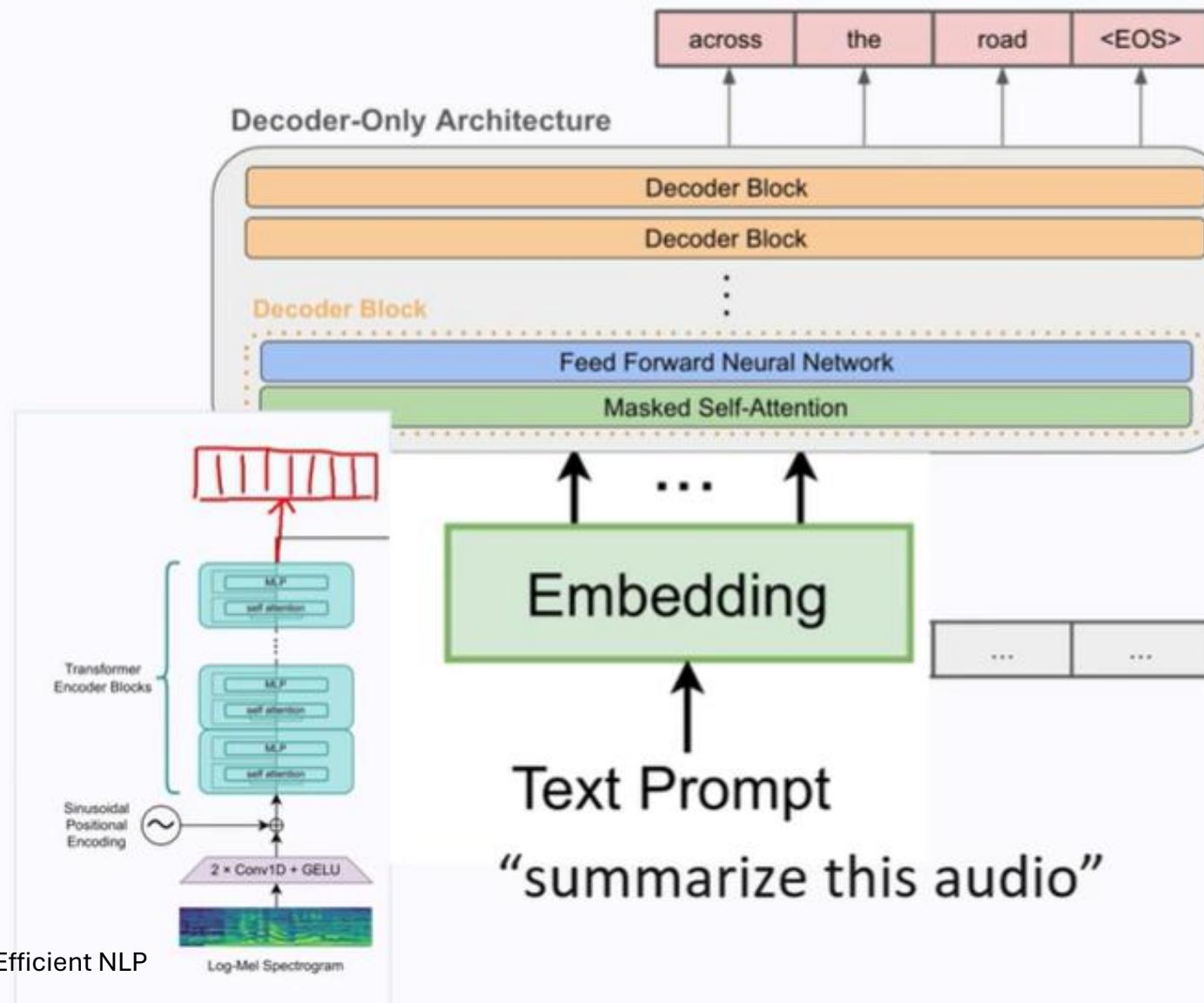


OpenAI Whisper: Encoder-decoder transformer for ASR

Other options:

- HuBERT
- Wav2vec 2.0
- W2V-BERT

Large Language Model (LLM)



LLaMA: LLMs from Meta, in various sizes (3B-70B)

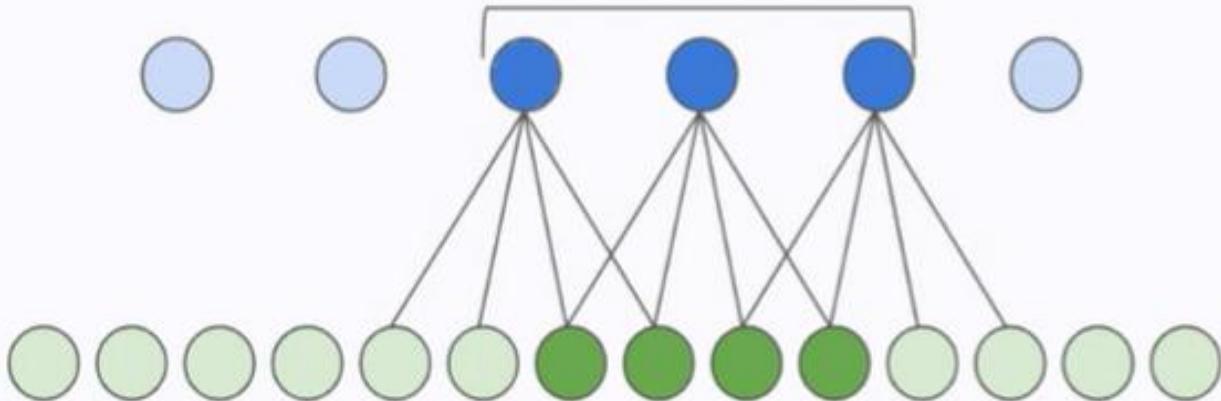
Alternatives: Mistral, Vicuna, Qwen, etc

N/A: Proprietary LLMs (OpenAI GPT-4, Anthropic Claude, Gemini, etc)

Length Adaptation

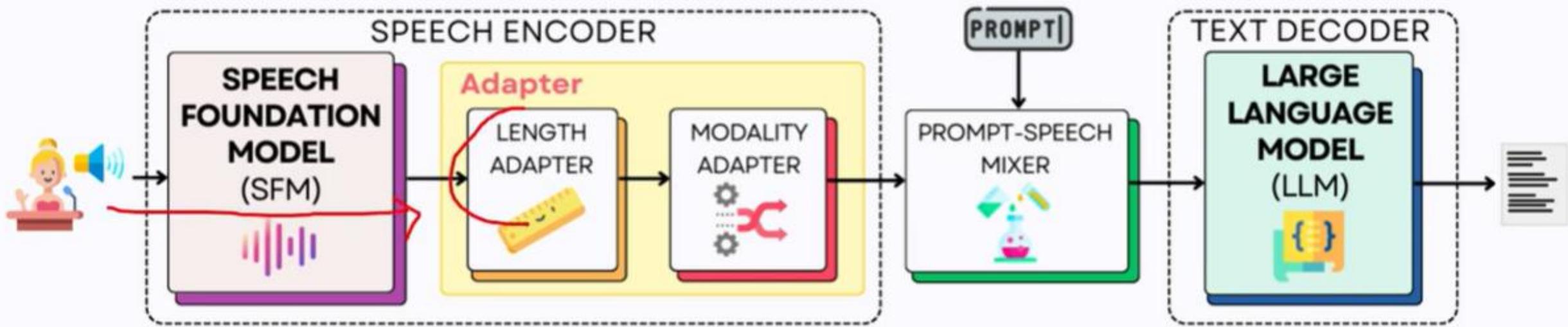
1 minute of speech

- Text: 150 words = 150 tokens = 150 input vectors
- Speech: 40 frames/sec = 2400 input vectors

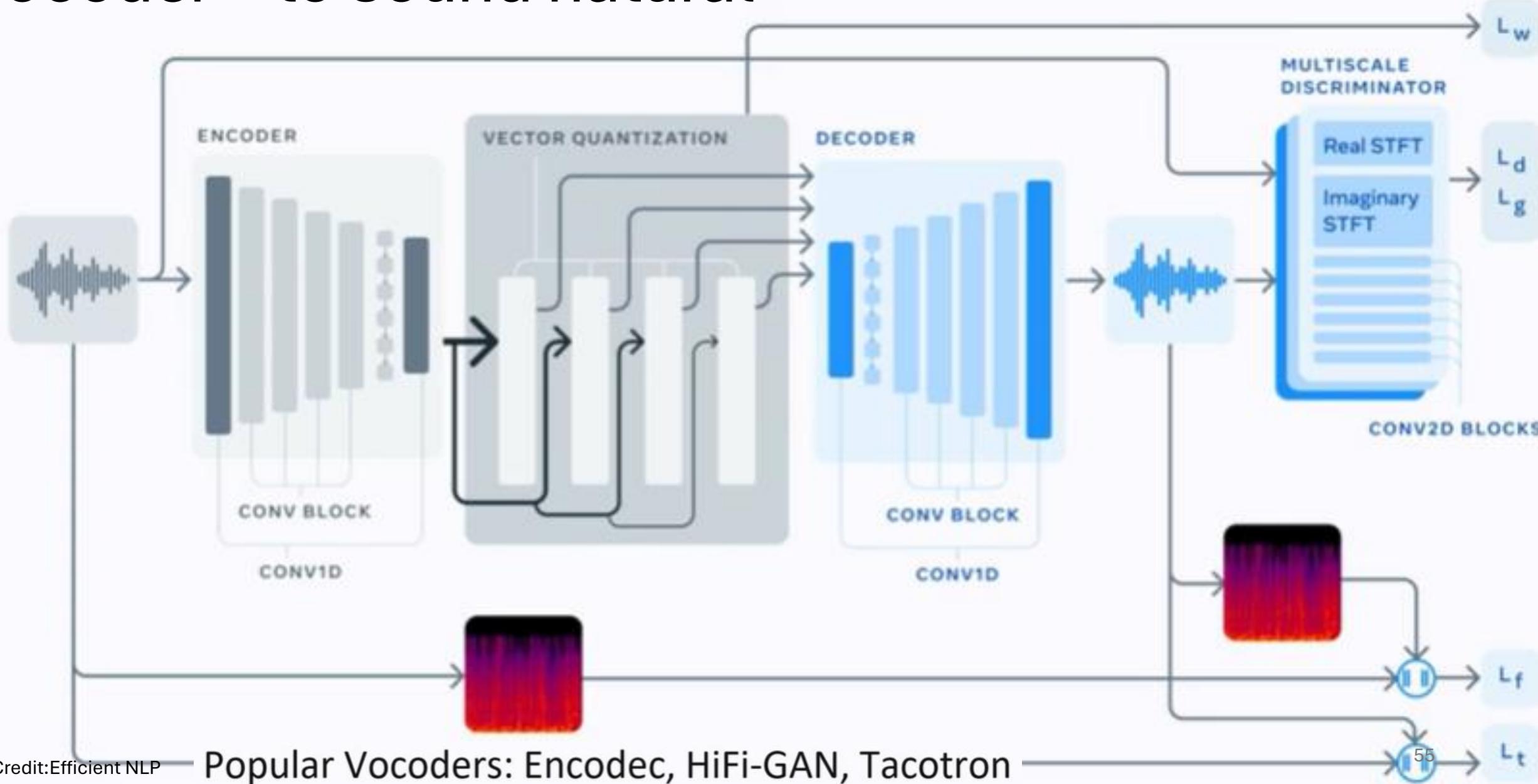


- Downsampling Conv1D layer
- CTC (Connectionist Temporal Classification)
- Averaging consecutive frames

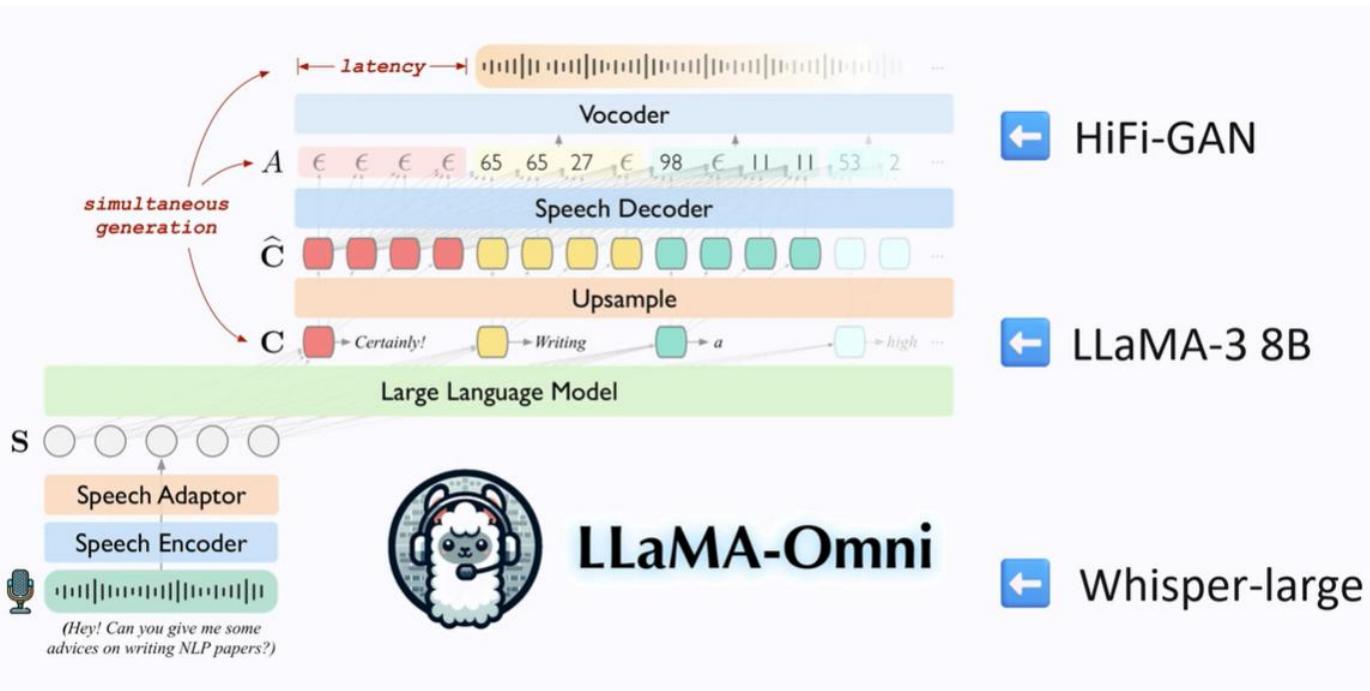
Speech LLM



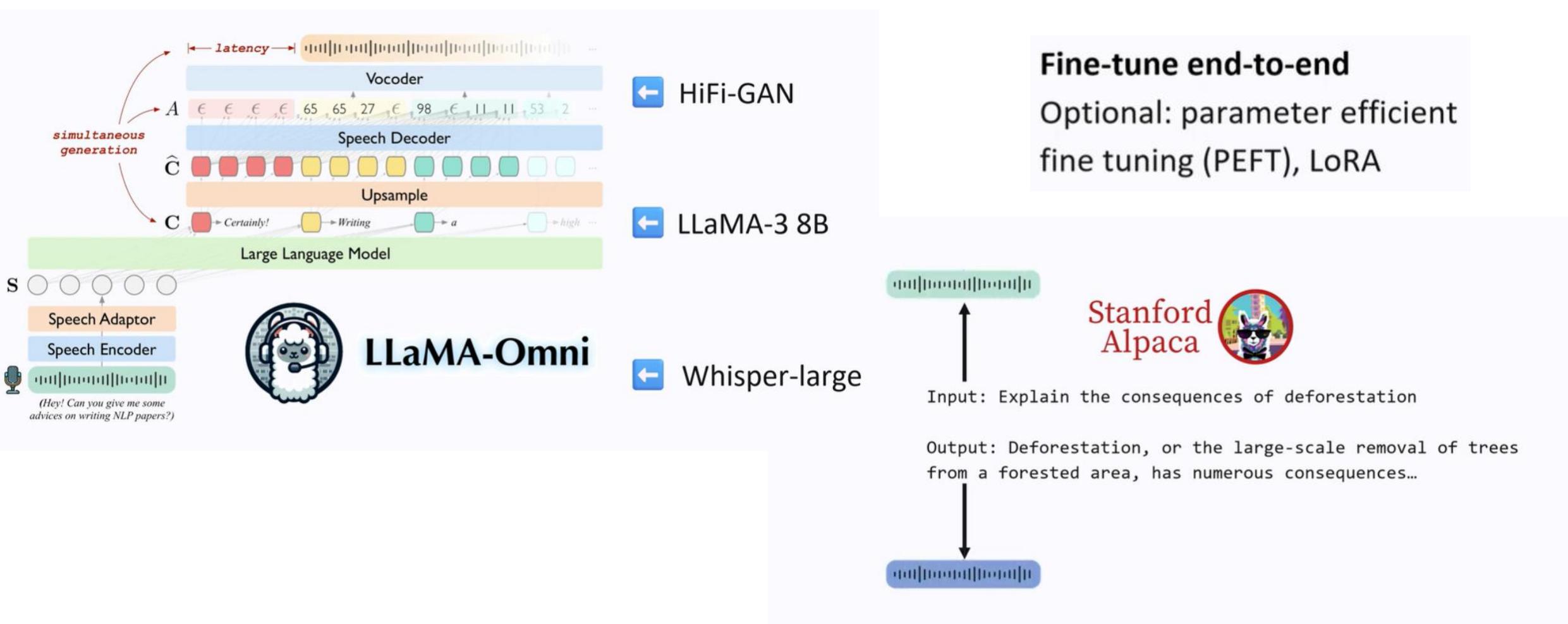
Vocoder – to sound natural



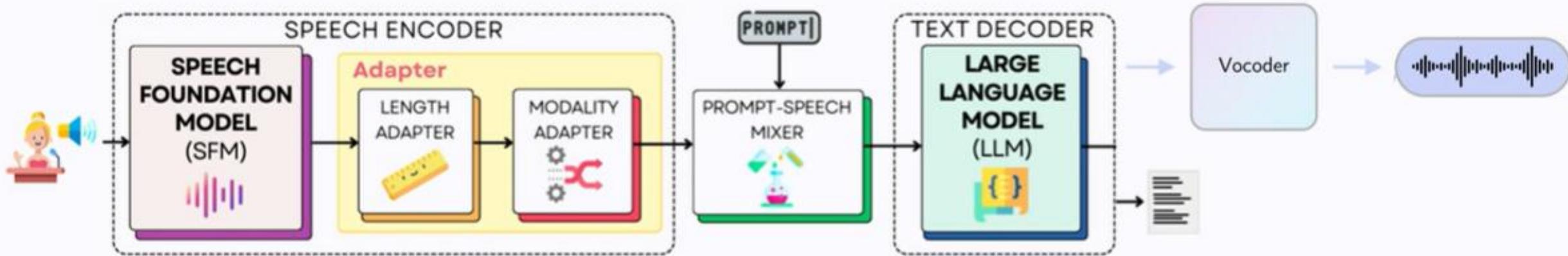
Speech LLM



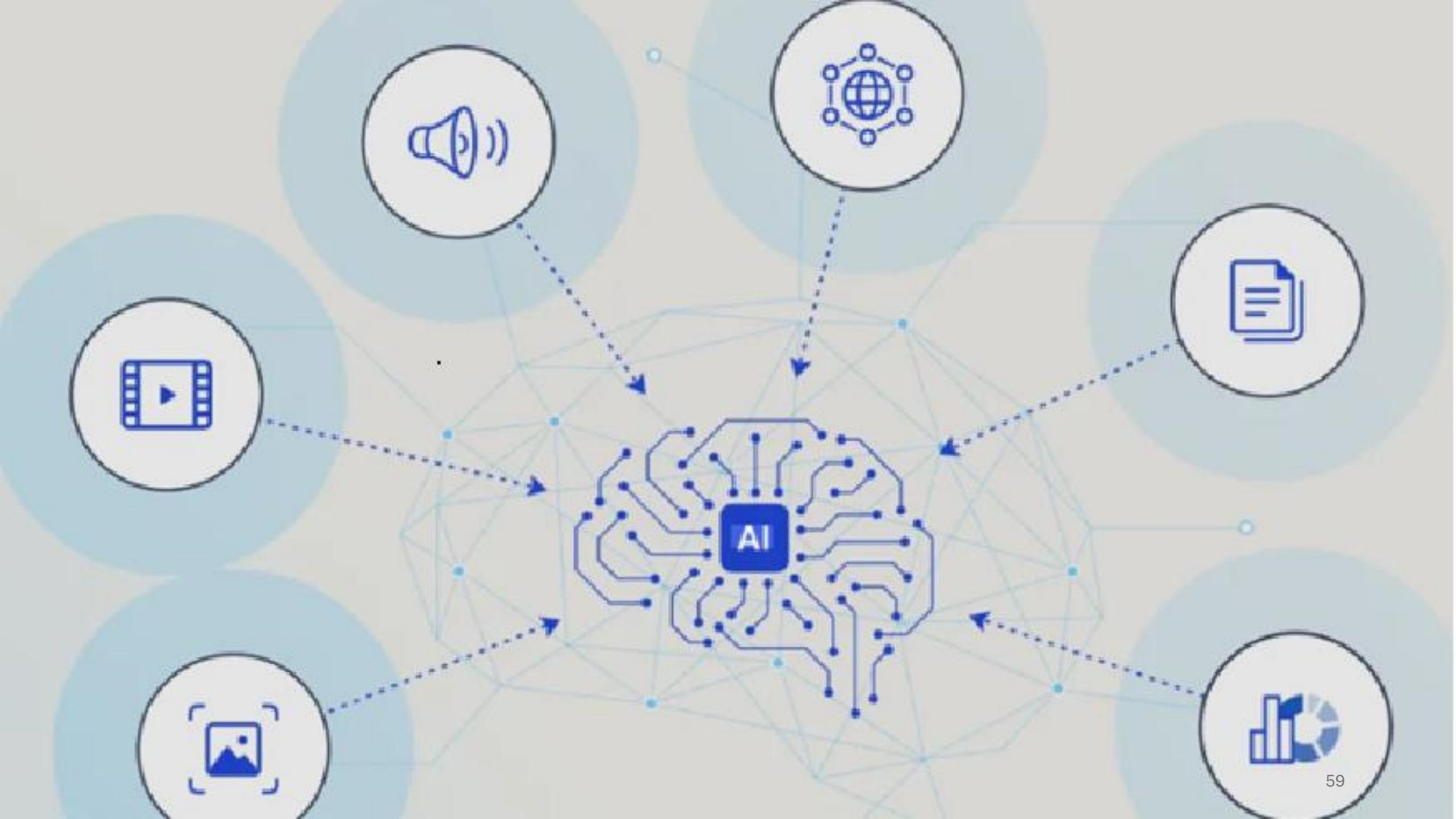
Speech LLM



Components of a Speech LLM

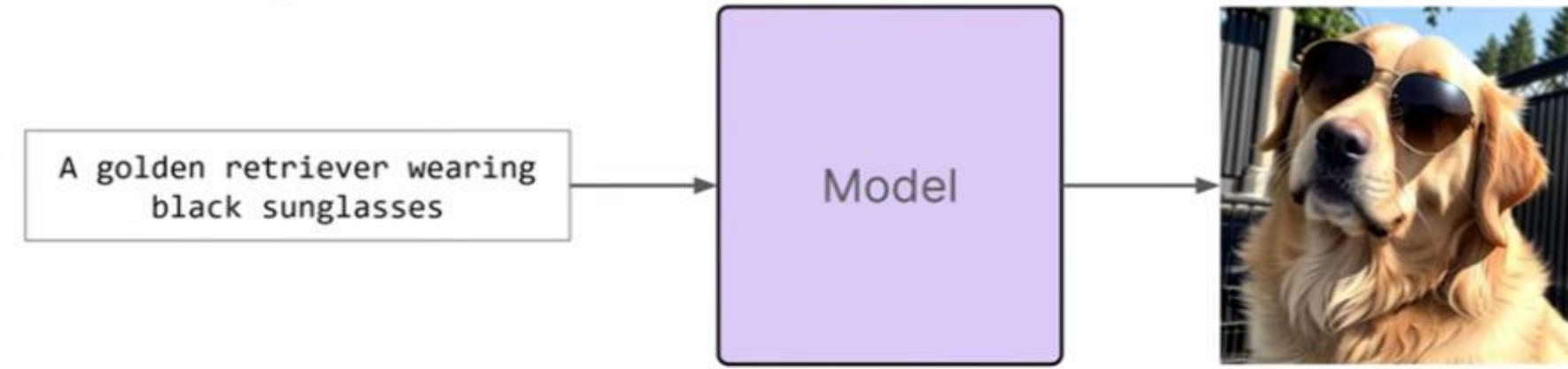


1. Speech Encoder (instead of speech-to-text)
2. Large language model
3. Vocoder (instead of text-to-speech)



Multimodal Models

Text-to-image Models



a red delicious apple

a black office chair

a corgi dog

Multimodal Models

Encode the images
and text snippets



Image Encoder

a red delicious apple

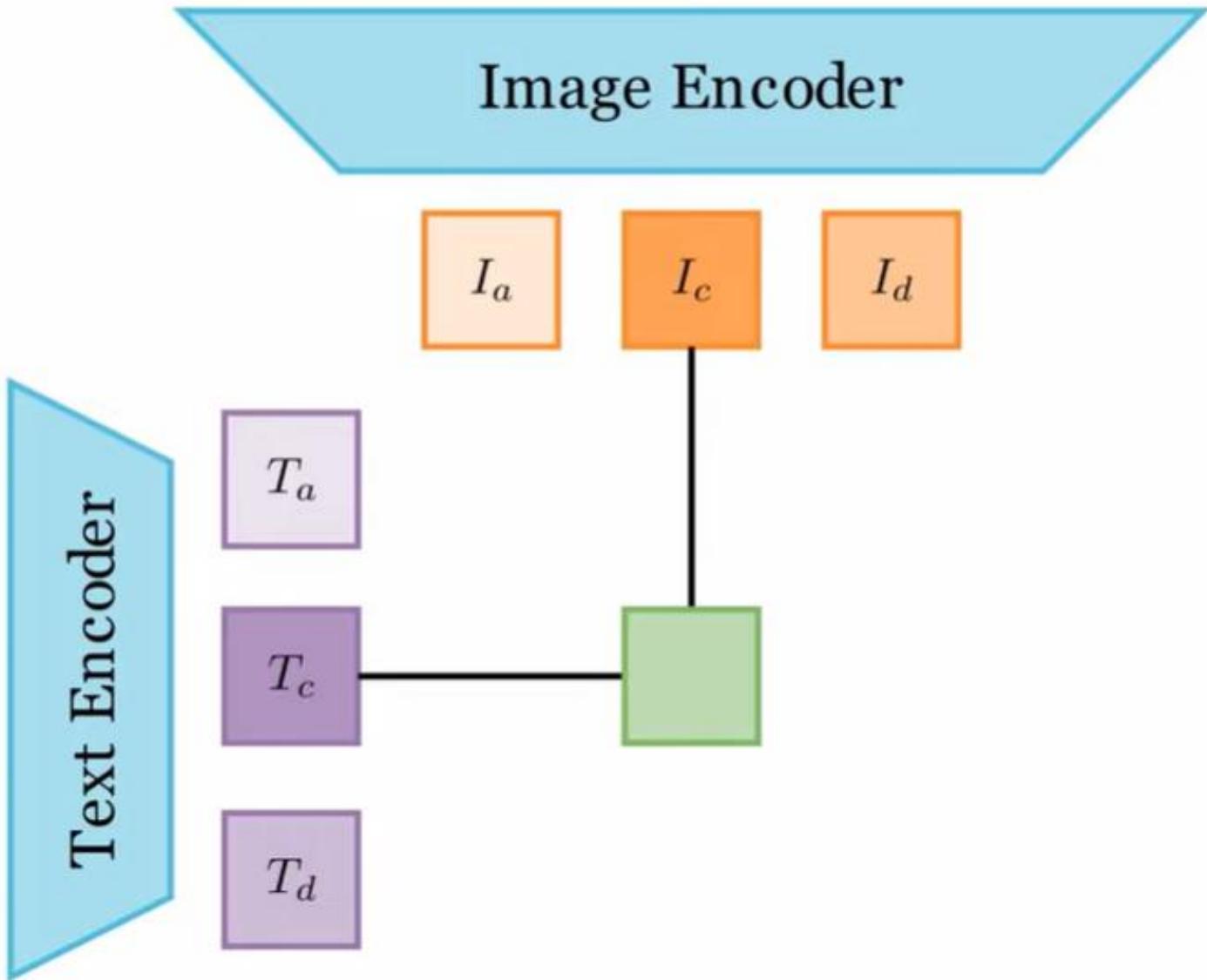
Text Encoder

a black office ch air

a corgi dog

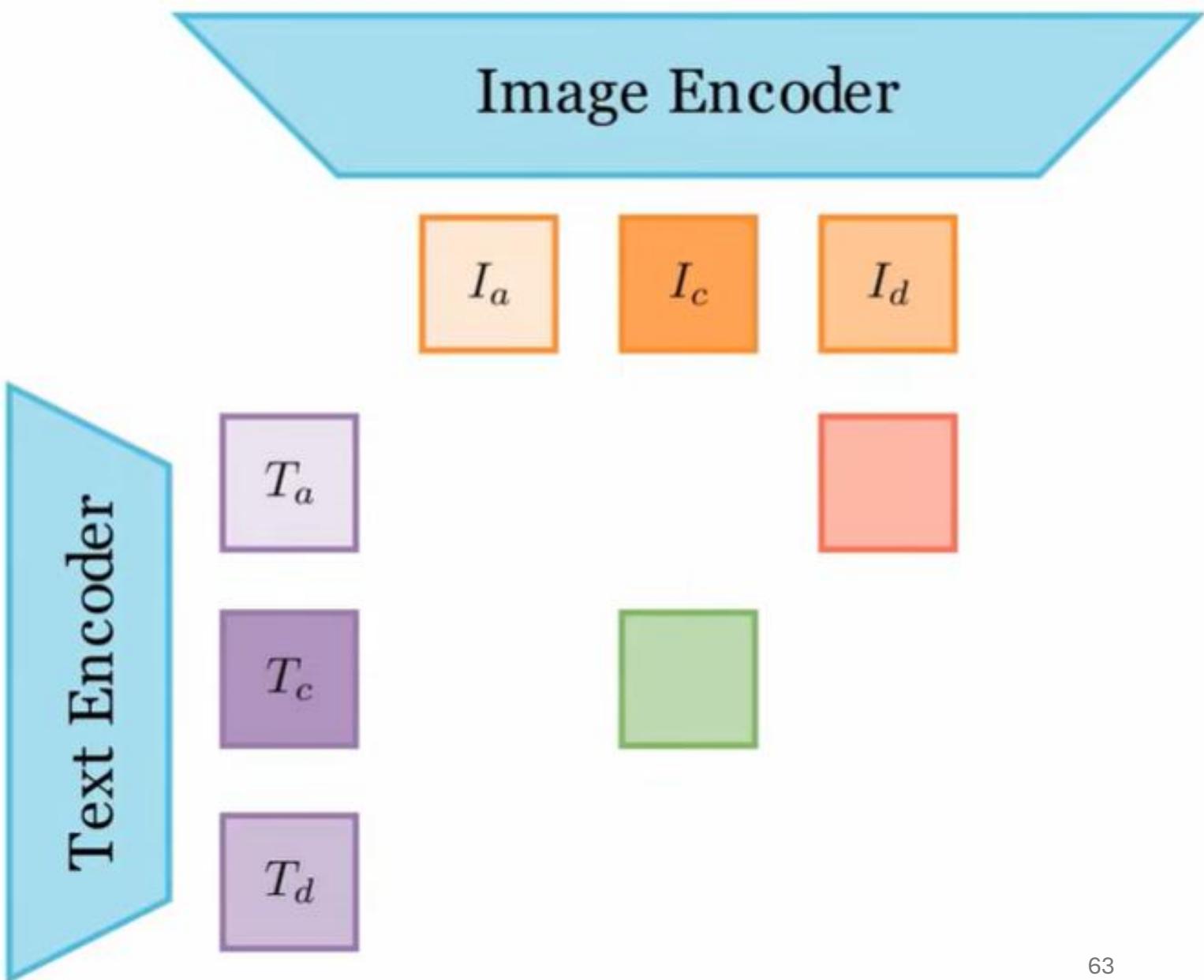
Multimodal Models

Train to increase cosine similarity
of matching image/text encoding pairs

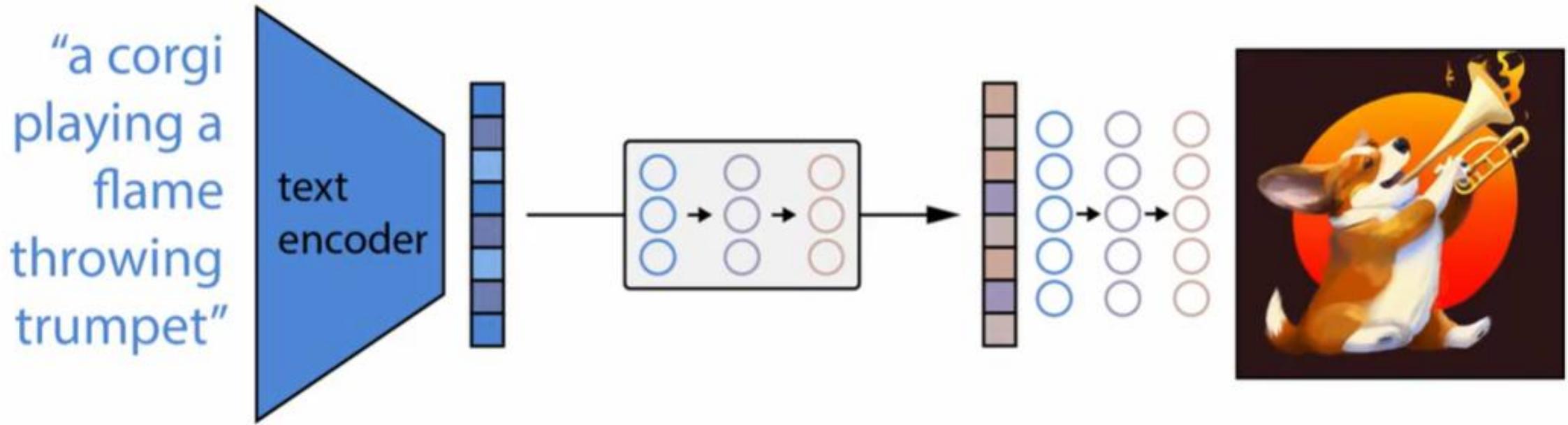


Multimodal Models

Train to decrease cosine similarity
of non-matching image/text encoding pairs



Multimodal Models



Encode the meaning of the text and decode to an image

LLM frequent tasks

Language Generation

Text generation

Code generation

Language Understanding

Text classification & sentiment analysis

Text summarization

Language translation

Intent recognition

Question-answering

Named entity recognition

Problem with classic IR

- The **vocabulary mismatch problem**
- tf-idf or BM25 cosine similarities only work if there is **exact word overlap** between query and doc!
- But query-writer can't know the exact words the doc might include!

Dense retrieval

- Instead of representing query and documents with count vectors
- Represent both with embeddings!

Hypothetical version of dense retrieval: static embeddings

Replace tf-idf vectors with, Static Word Embeddings

- Query: the mean of the embeddings of each query word
- Doc: the mean of the doc word embeddings
- Now just compute query-doc cosine as normal.

We don't do this because **contextual embeddings** work much better!

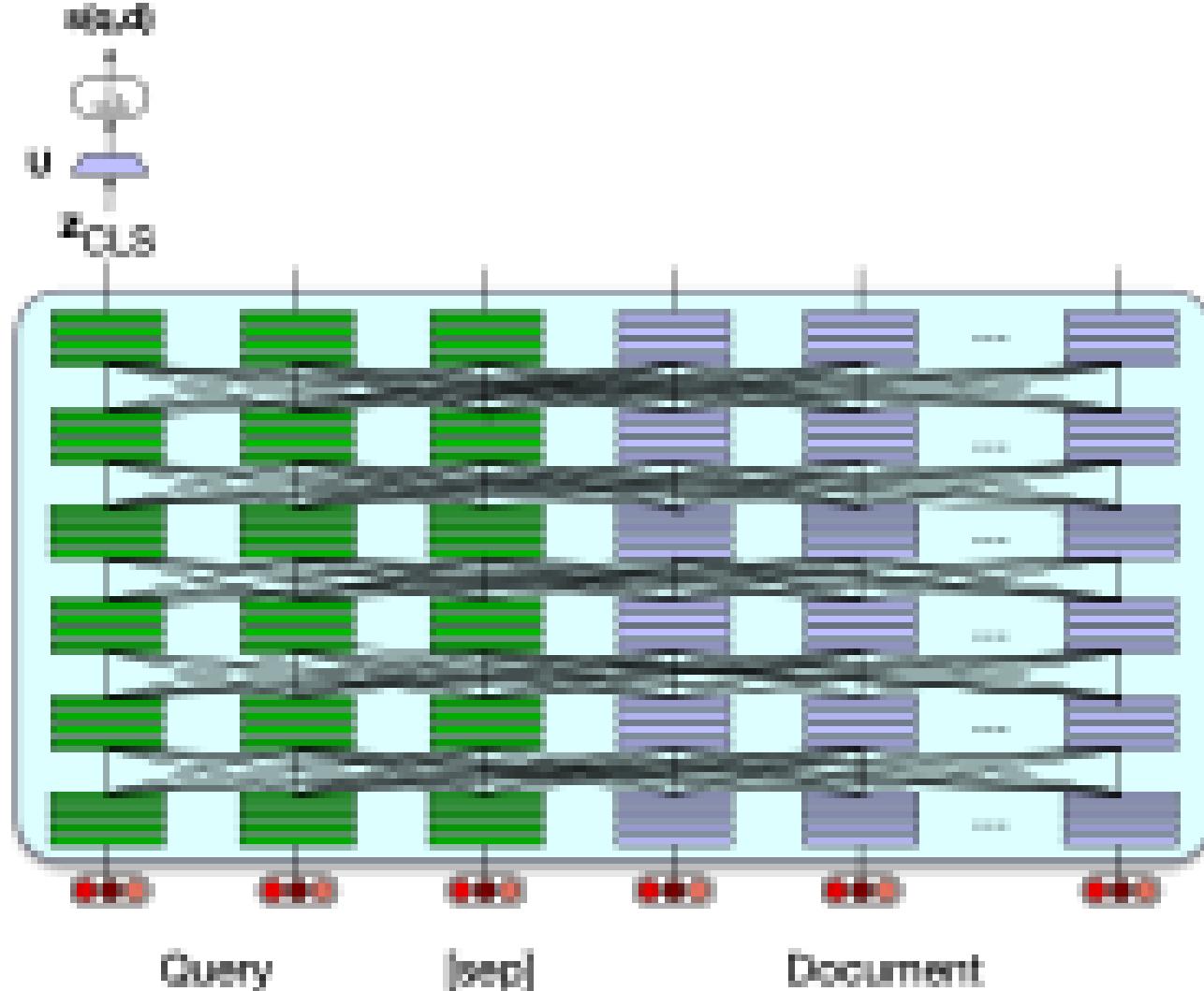
Dense IR – Structural Problems

- QA - How do we produce the answer? Could we have a conversation?
- Vectors couldn't/shouldn't represent full document
- How do we find the vector representation and similar vectors?
- How do we rank the results?

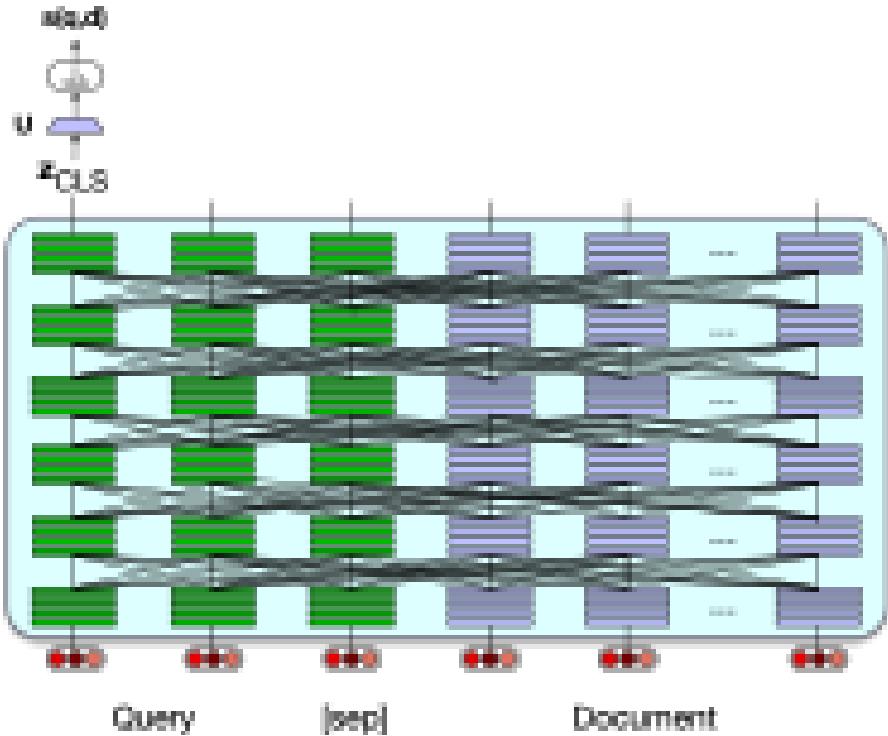
Dense IR – Structural Problems

- QA - How do we produce the answer? Could we have a conversation?
 - Solution: Use a dedicated LLM that is trained to do that
- Vectors couldn't/shouldn't represent full document
 - Solution: Chunking (will be discussed later)
- How do we find the vector representation and similar vectors?
 - Solution: VectorDB (will be discussed later)
- How do we rank the results?
 - Solution: Train the LLM to produce a score

Dense retrieval #1: Single encoder

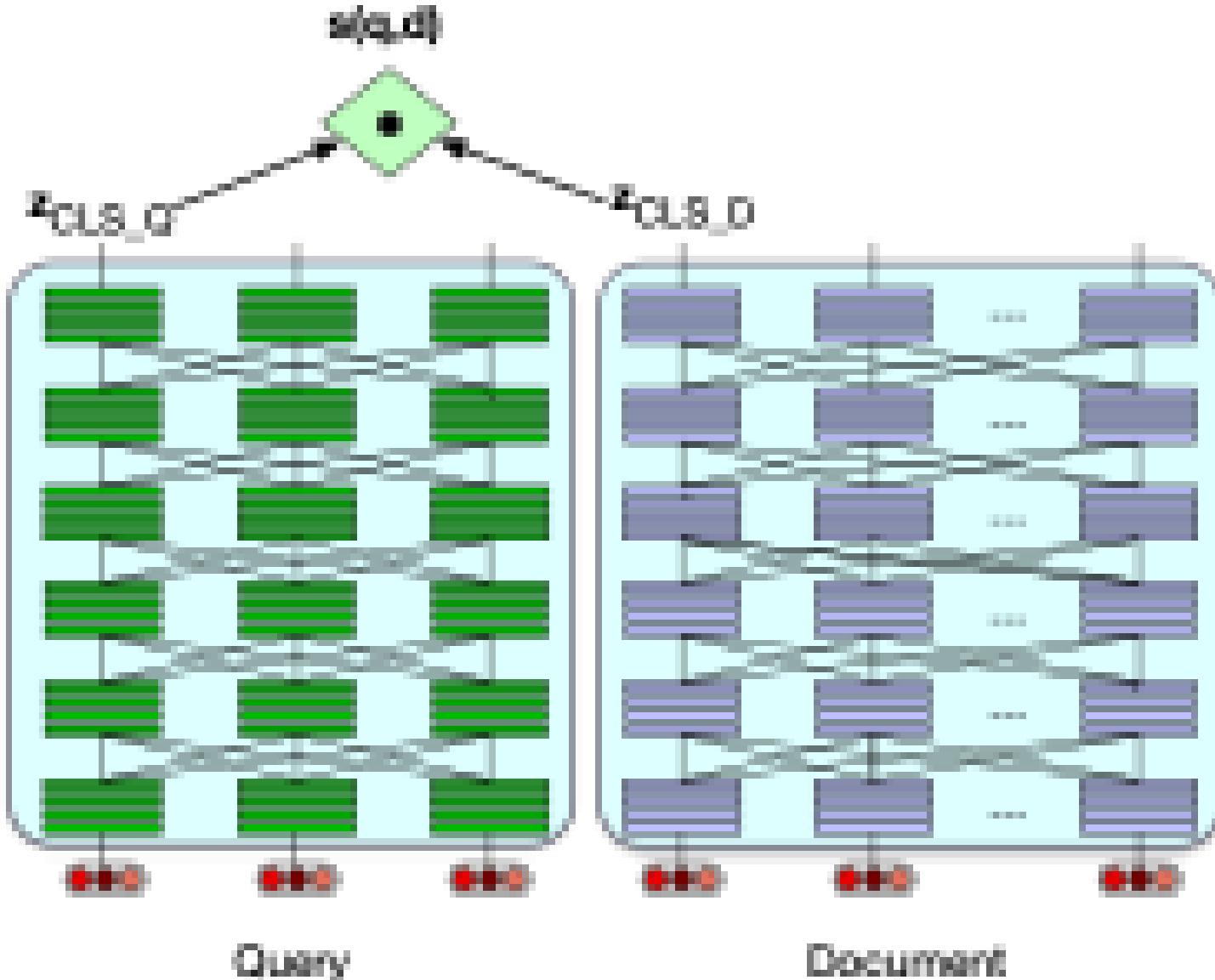


Dense retrieval #1: Single encoder



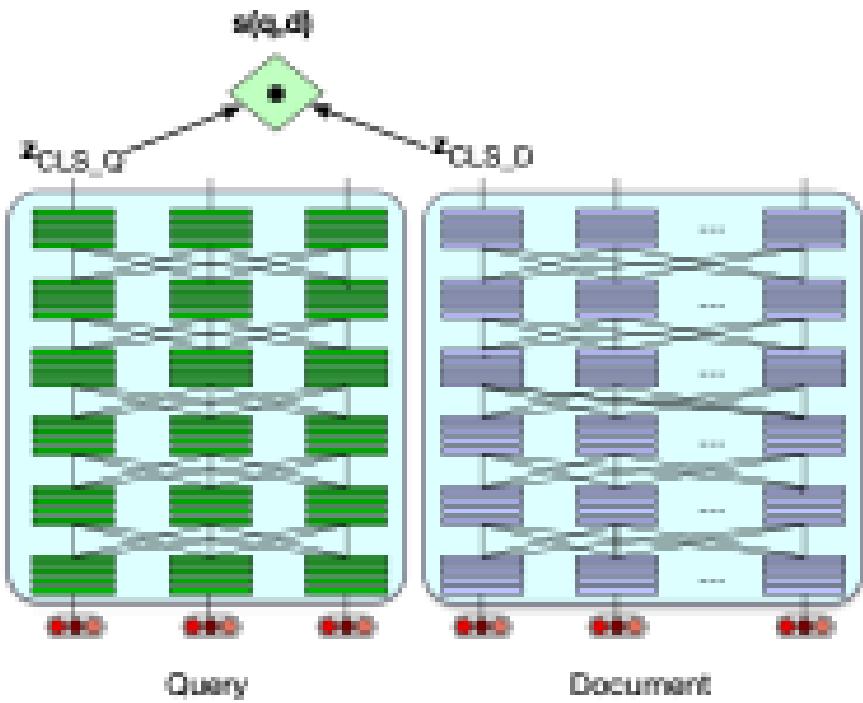
- System is run on passages (say 100 tokens) instead of whole documents
- **Training:**
 - BERT and linear layer U can then fine-tuned for relevance
 - Creating a tuning dataset of relevant and non-relevant passages.

Dense retrieval #2 (biencoder)



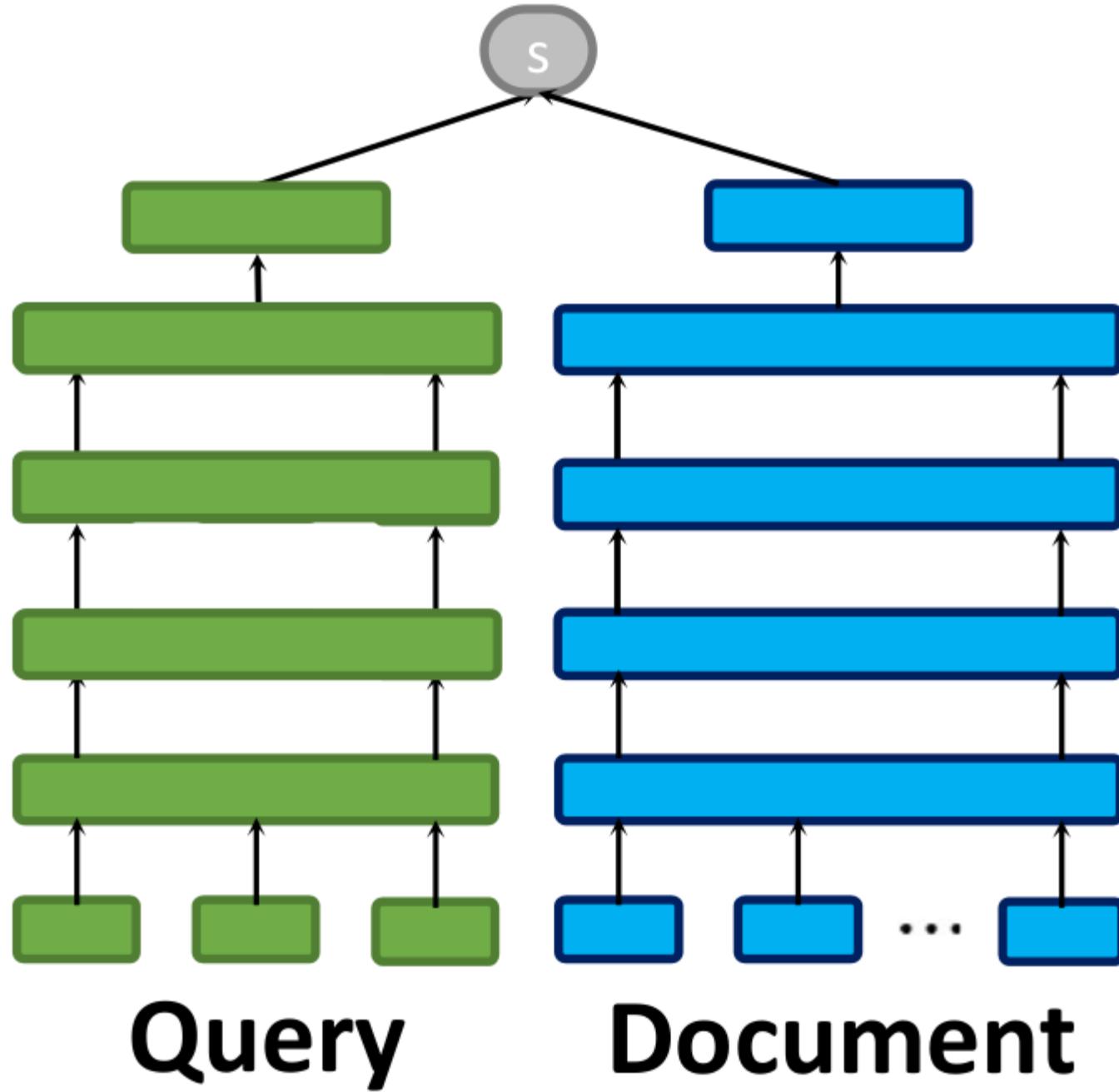
$$\begin{aligned}\mathbf{z}_q &= \text{BERT}_Q(q) [\text{CLS}] \\ \mathbf{z}_d &= \text{BERT}_D(d) [\text{CLS}] \\ \text{score}(q, d) &= \mathbf{z}_q \cdot \mathbf{z}_d\end{aligned}$$

Dense retrieval #2 (**biencoder**)

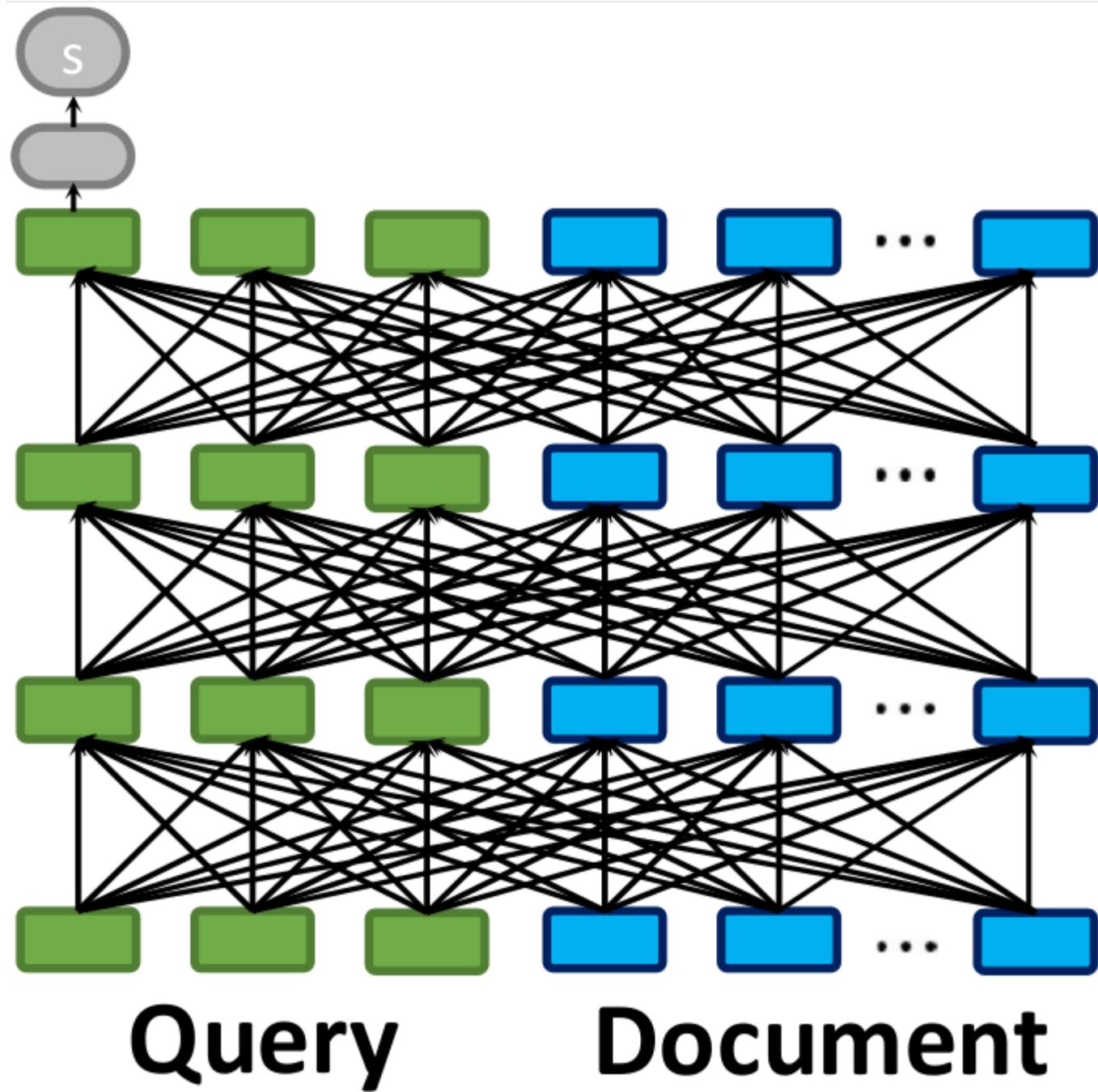


- Encode doc vectors in advance.
- Encode query when it arrives
- Score is dot product between query vector and precomputed doc vector
- Cheaper but less accurate

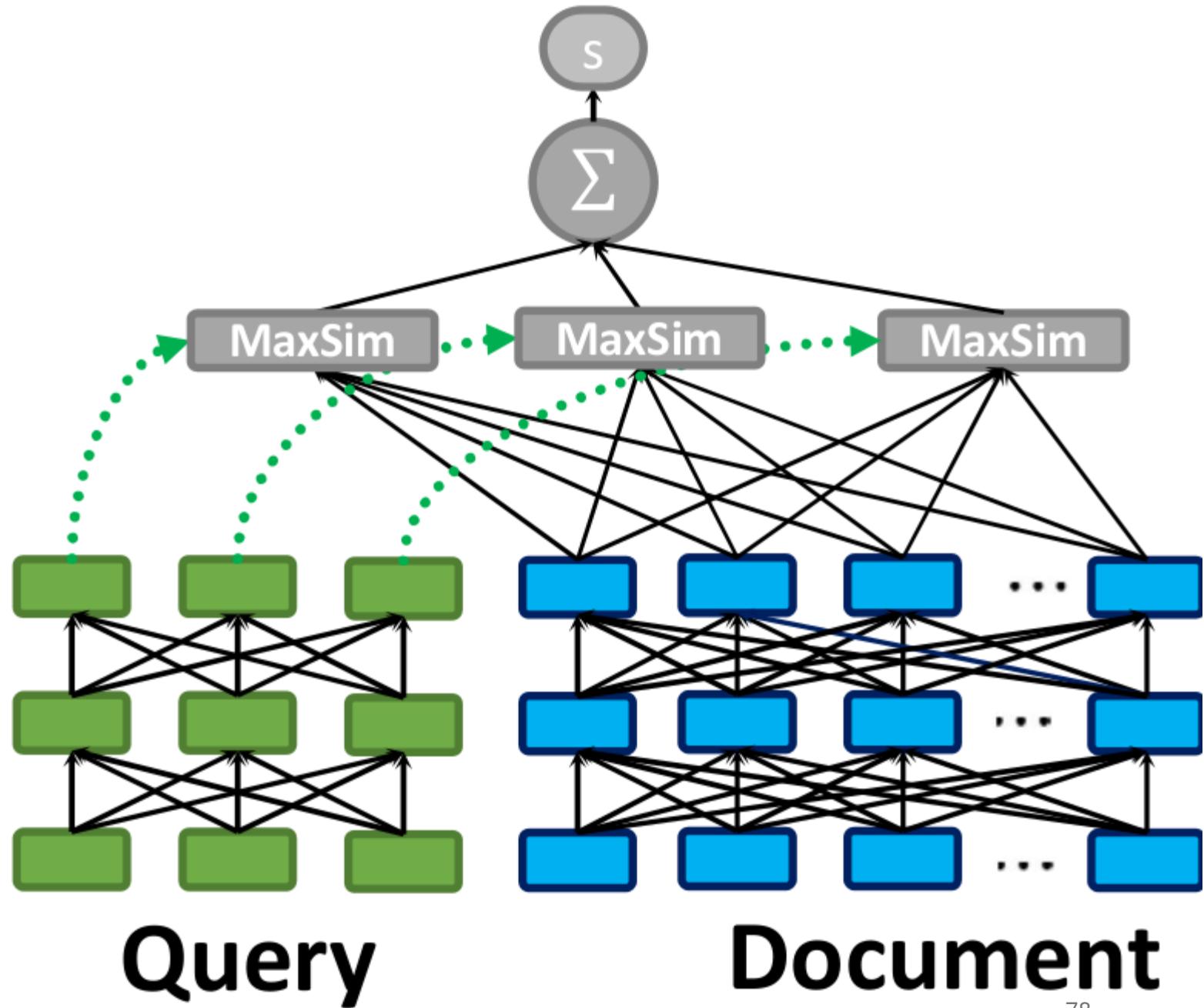
Bi-Encoders

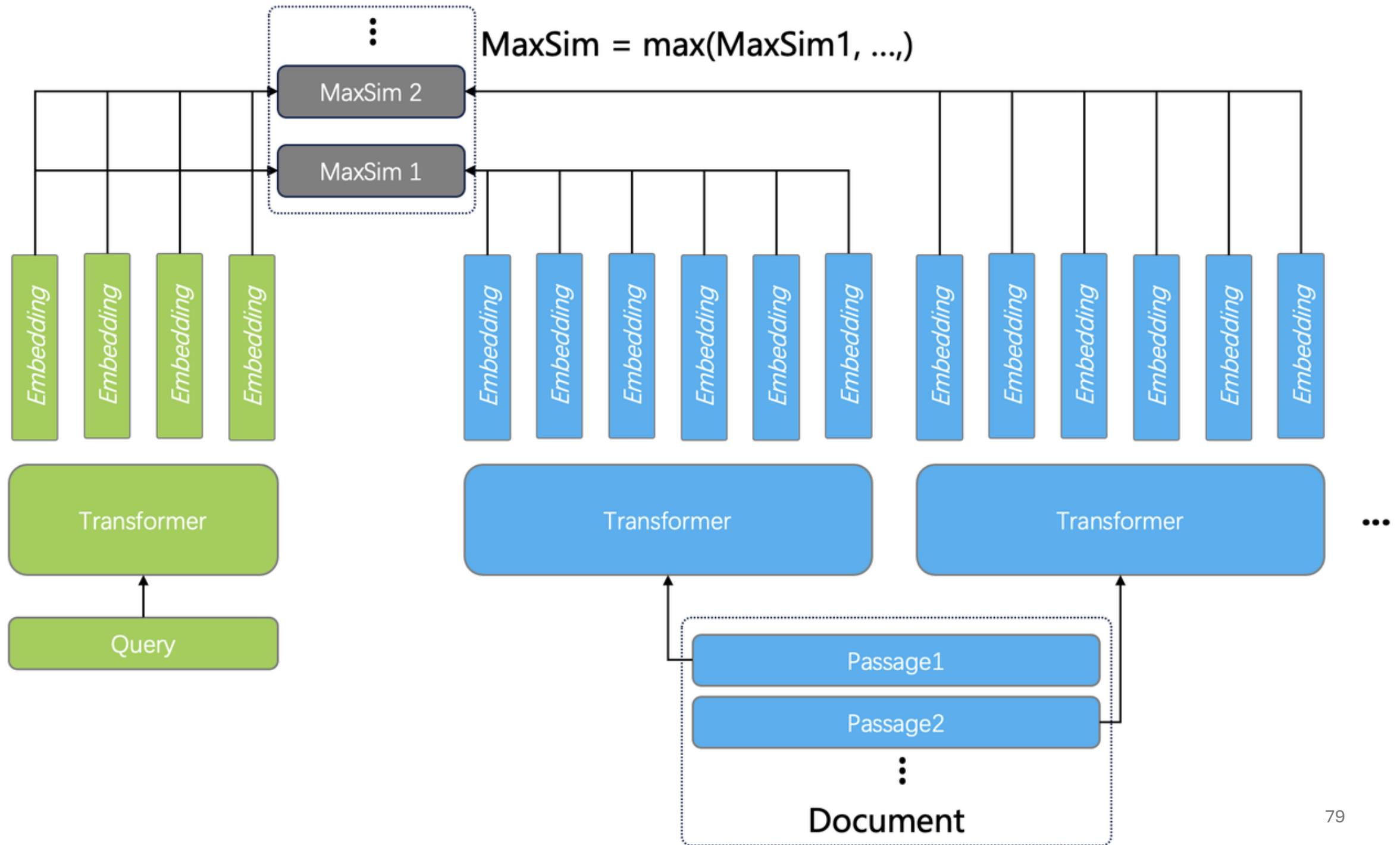


Cross-Encoder

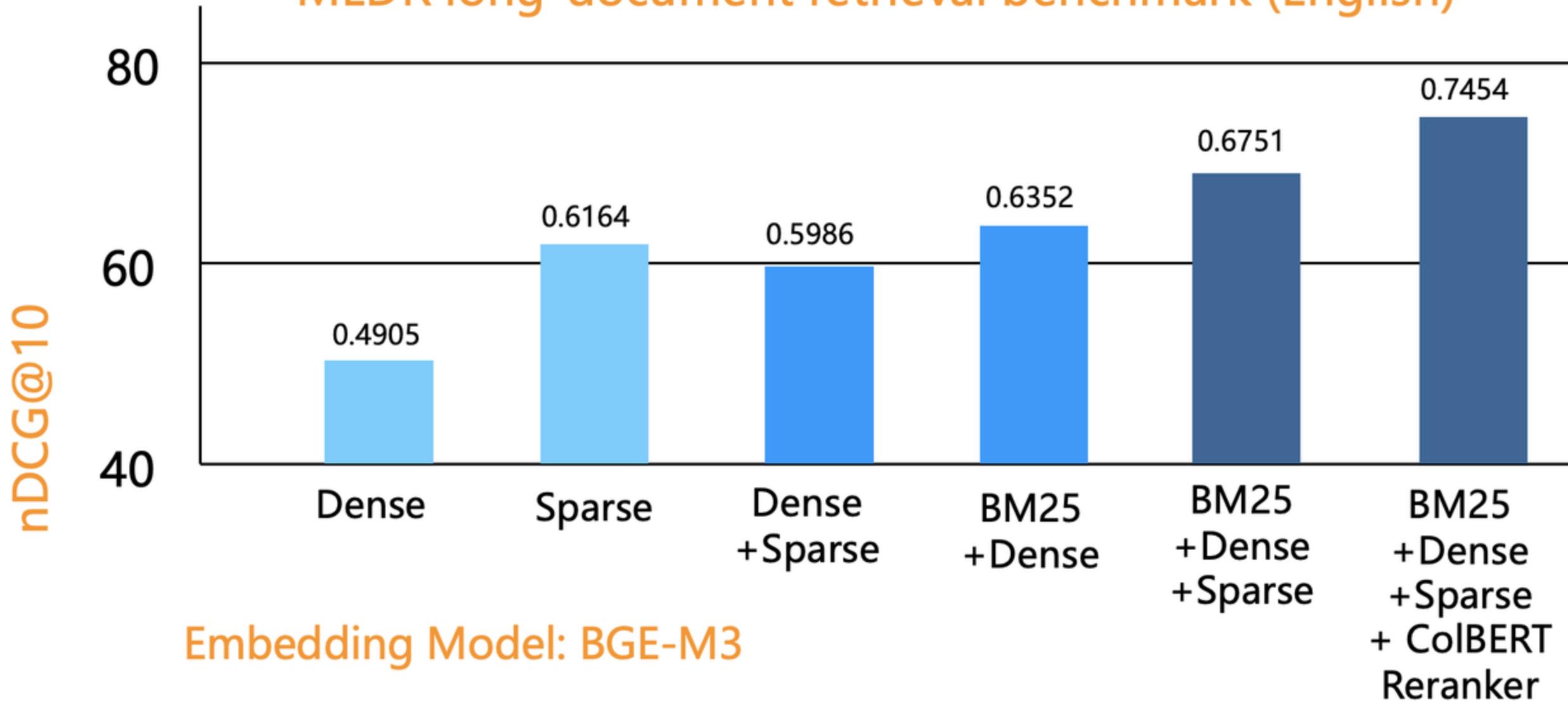


ColBERT (Contextualized Late Interaction)





MLDR long-document retrieval benchmark (English)



Embedding Model: BGE-M3

CoBERT

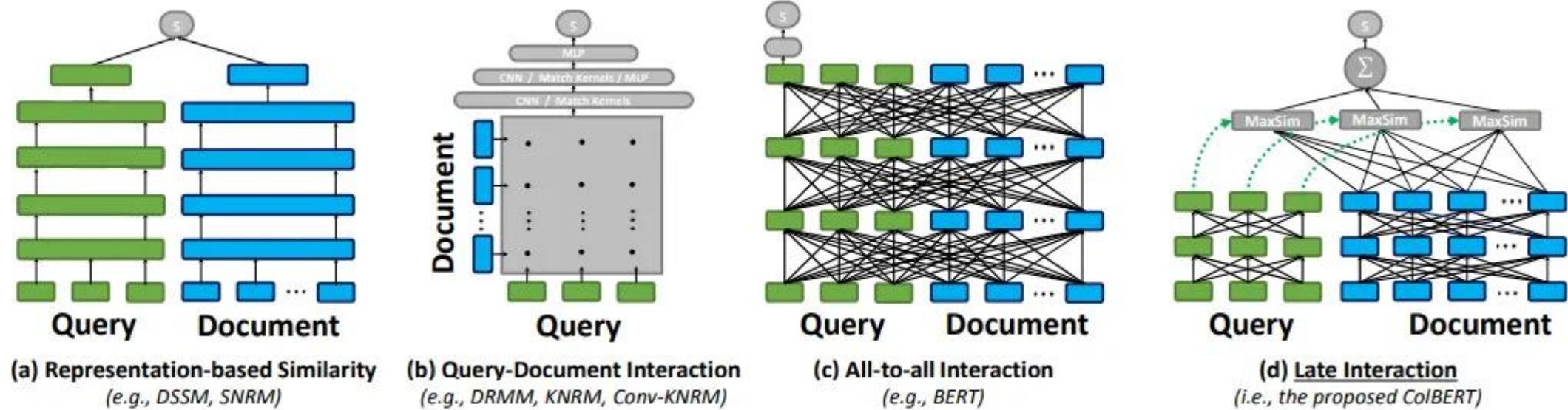
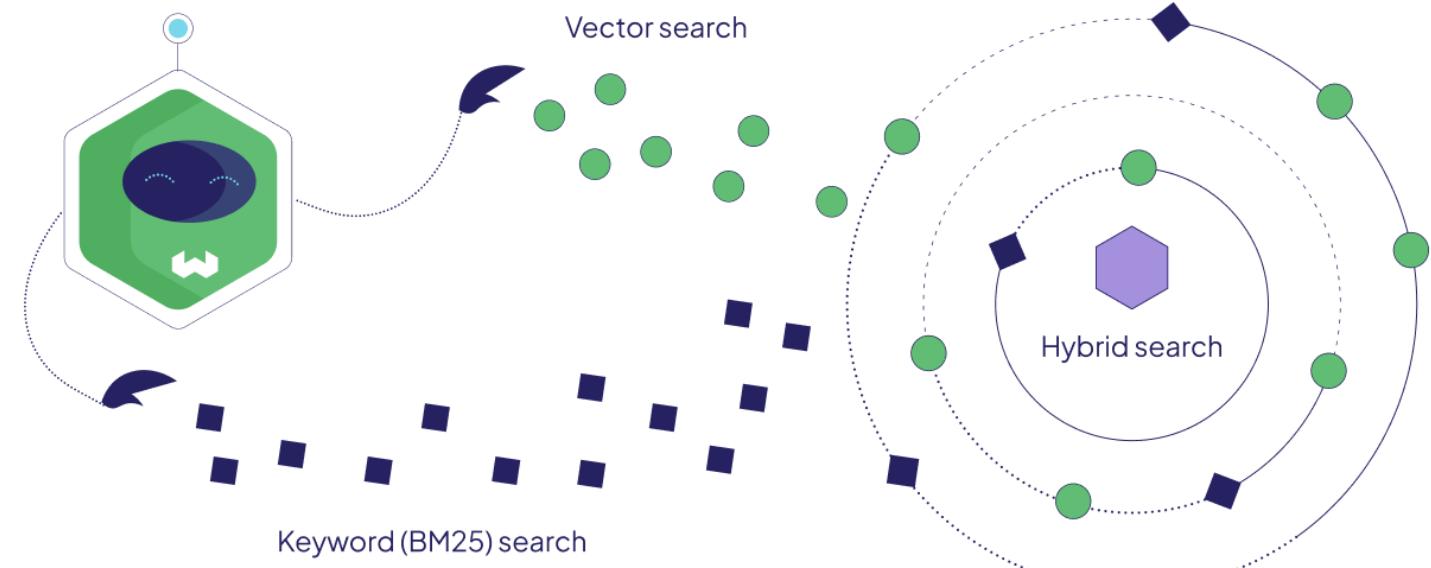


Figure 2: Schematic diagrams illustrating query-document matching paradigms in neural IR. The figure contrasts existing approaches (sub-figures (a), (b), and (c)) with the proposed late interaction paradigm (sub-figure (d)).

Hybrid Search

Advantages

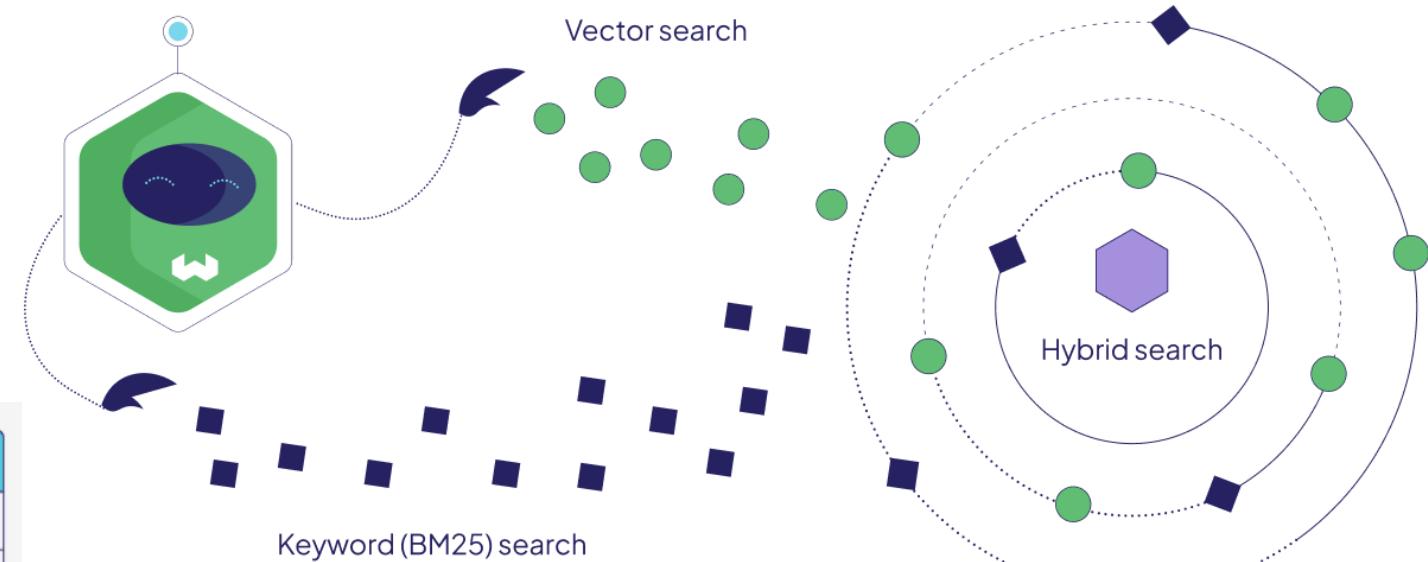
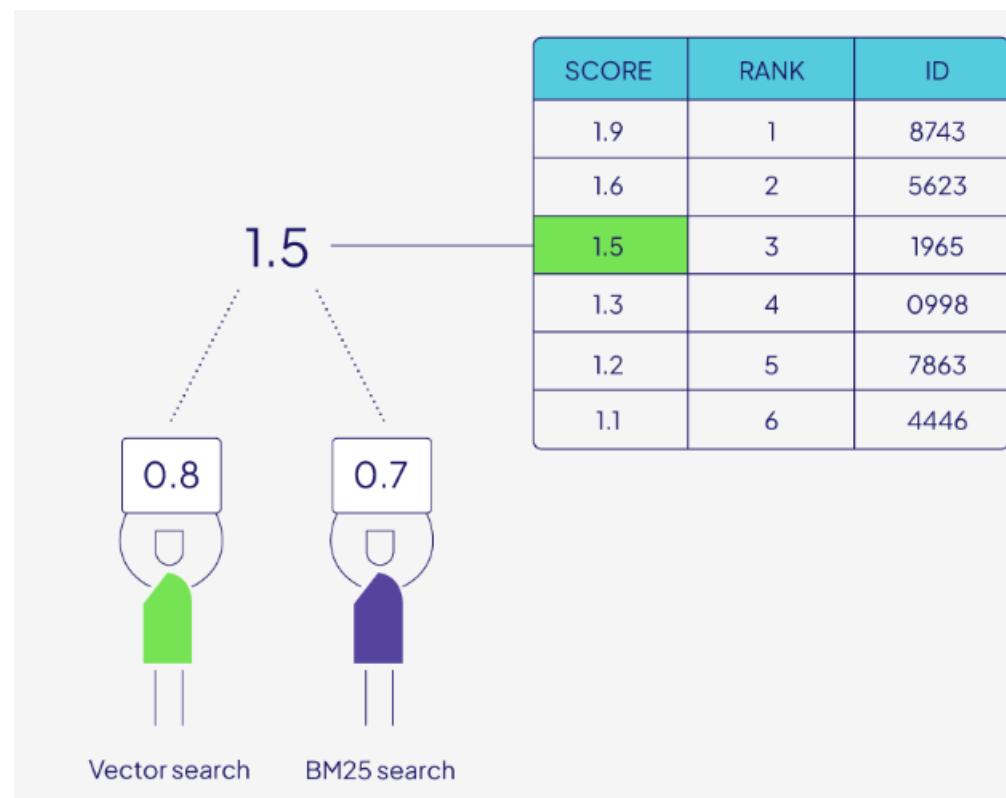
- keywords still needed:
 - ambiguous phrases,
 - domain-specific terms
 - misspellings
- Understanding the semantics
- improves user experience – non perfectly phrased queries



$$score(D, Q) = \sum_{i=1}^n IDF(q_i) \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgd_1})}$$

Source: [Wikipedia page on Okapi BM25](#)

Hybrid Search – Simple Fusion



$$score(D, Q) = \sum_{i=1}^n IDF(q_i) \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgd_1})}$$

Source: Wikipedia page on Okapi BM25

LLMs for IR - problems

Problems traditionally found in LLMs:

- **Data obsolescence (outdated)** - RAG can retrieve up-to-date information
- **Re-training costs:** don't need frequent and costly re-training
- **Hallucinations:** grounding responses in factual information retrieved from real documents

LLMs Hallucinate

- Hallucination: a response that is not faithful to the facts of the world.
- In the legal domain LLMs were shown to hallucinate up to 88% of the time!

Dahl et al. (2024)

Hallucination - Factual

Category	Type	Example Dialog	Explanation
Factuality Hallucination	Factual	Instruction: What are Thomas Edison's main contributions to science and technology?	The response is factually incorrect. In reality, Edison improved the light bulb, building on earlier designs by others, and Alexander Graham Bell invented the telephone.
	Contradiction	Response: Thomas Edison developed the first practical telephone and invented the light bulb.	
Factuality Hallucination	Factual Fabrication	Instruction: What are the major environmental impacts of the construction of the Eiffel Tower?	The response erroneously states that the construction of the Eiffel Tower led to the extinction of the 'Parisian tiger' and sparked the global green architecture movement. In fact, there is no evidence of a 'Parisian tiger' ever existing, making this a fabricated claim. Moreover, attributing the origins of green architecture to the Eiffel Tower is an exaggeration, as this movement has diverse roots and cannot be traced back to a single event.
		Response: The construction of the Eiffel Tower in 1889 led to the extinction of the Parisian tiger , a species that played a crucial role in the region's ecosystem. Additionally, it is widely recognized as the event that sparked the global green architecture movement .	

Hallucination - Faithfulness

Category	Type	Example Dialog	Explanation
	Instruction Inconsistency	<p>Instruction: Translate the English question into Spanish: "What is the capital of France?"</p> <p>Response: The capital of France is Paris.</p>	The user explicitly provided an instruction to translate the English question into Spanish. Instead of performing the translation as instructed, the LLM chose to answer the question.
Faithfulness		<p>Retrieved context: The Nile is a major north-flowing river in northeastern Africa. It's traditionally considered the longest river in the world. Originating in the Great Lakes region of central Africa, it flows through several countries before emptying into the Mediterranean Sea.</p>	
Hallucination	Context Inconsistency	<p>Instruction: Summarize the retrieval document above.</p> <p>Response: The Nile originates from the mountain ranges of central Africa and flows into the Mediterranean.</p>	The summary incorrectly states that the Nile originates from the "mountain ranges of central Africa", while the paragraph mentions that the Nile originates in the Great Lakes region of central Africa.
Logical	Inconsistency	<p>Instruction: Solve the equation $2x + 3 = 11$ step by step.</p> <p>Response: Step 1: Subtract 3 from both sides to get: $2x = 8$. Step 2: Divide both sides by 2 to get: x = 3. So the result of the equation is x=3.</p>	The first step correctly isolates the term with x to get $2x = 8$. However, the next step inaccurately divides 8 by 2 to yield a result of x = 3, which is inconsistent with the earlier reasoning.

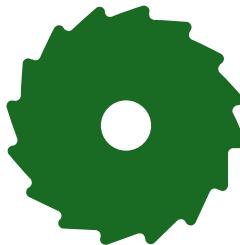
Hallucination - Misinformation and Biases

Type	Premise	Input	Model Output	Explanation
<i>Imitative Falsehood</i>	Within the LLM's training data, a frequently referenced but incorrect fact appears multiple times: <i>Thomas Edison invented the light bulb.</i>	Who is credited with the invention of the light bulb?	Thomas Edison is credited with the invention of the light bulb.	While Edison improved the light bulb's design, he was not its sole inventor . The model's answer reflects a misconception prevalent in its training data, exemplifying an imitative falsehood.
<i>Societal Biases</i>	In the model's training data, certain names may frequently be associated with specific nationalities.	Provide a summary of the following news: Dr. Kim recently published a breakthrough paper on quantum computing. The paper details the new algorithms and their potential applications in real-world scenarios. It has been hailed as revolutionary by experts.	Dr. Kim from South Korea recently published a paper on quantum computing, detailing new algorithms and their applications.	The model erroneously added "South Korea" based on the name "Kim" revealing its bias in linking certain names to specific nationalities, which can cause hallucinations in the summary.

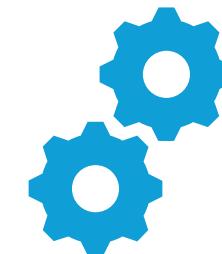
Possible Solutions



Fine tuning



RAG

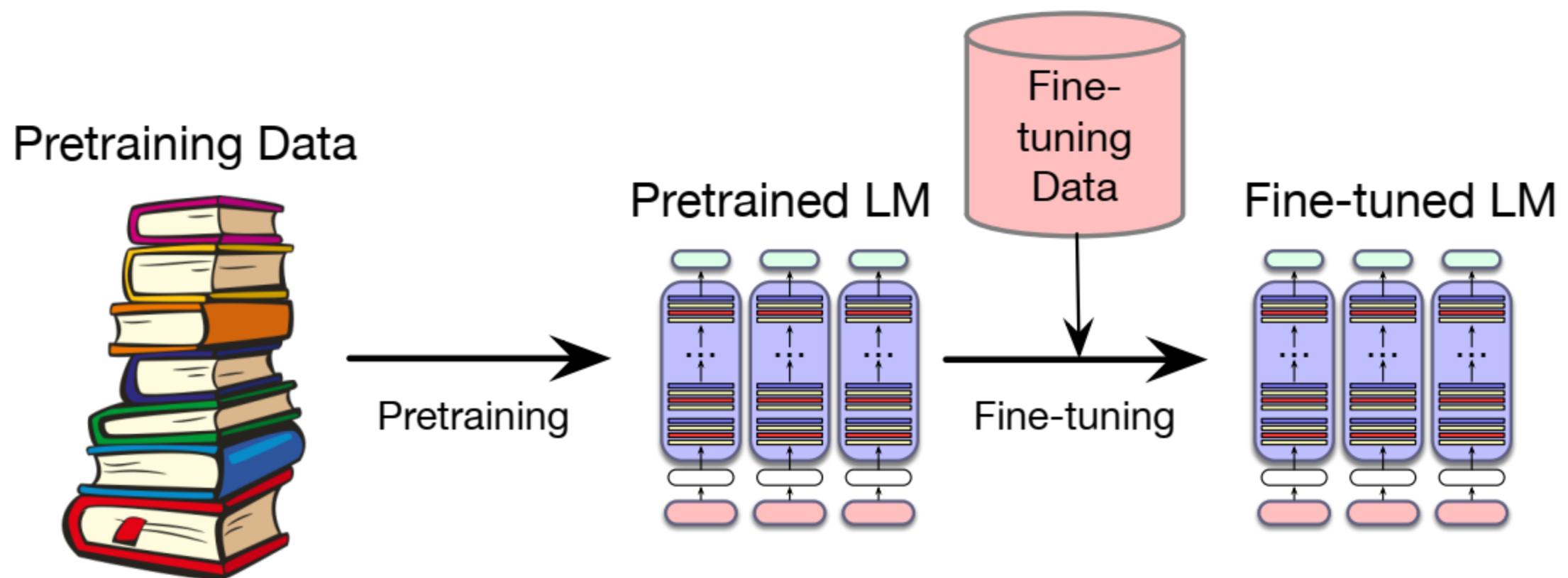


Prompt Engineering

What is LLM Fine-Tuning

- **fine-tuning** an existing (and previously trained) LLM consists of **adjusting its model parameters** using **additional, specialized training data** to enhance its performance in a specific use case or application domain.
- Fine-tuning is an important part of LLM development, maintenance, and reuse, for two reasons:
 - It allows the model to adapt to a more domain-specific, often smaller, dataset, improving its accuracy and relevance in specialized areas such as legal, medical, or technical fields. See the example in the image below.
 - It ensures the LLM stays up-to-date on evolving knowledge and language patterns, avoiding issues like outdated information, hallucinations, or misalignment with current facts and best practices

Fine tuning process



Finetuning as "continued pretraining" on new data

- Further train all the parameters of model on new data
- using the same method (word prediction) and loss function
- (cross-entropy loss) as for pretraining.
- as if the new data were at the tail end of the pretraining data
- Hence sometimes called **continued pretraining**

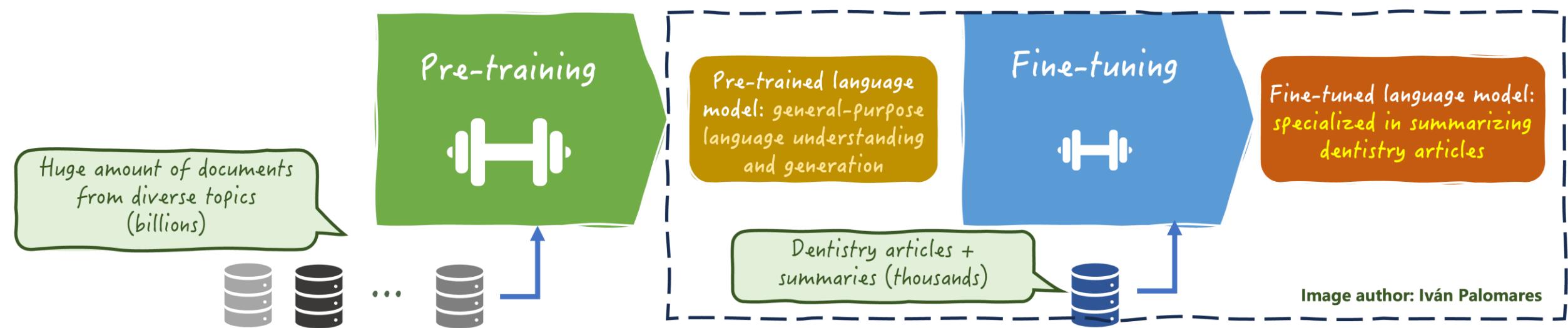
Fine tuning for specialized domain

Fine-tuning involves training the LLM further using your proprietary data. It essentially “teaches” the model to understand your specific domain better.

How It Works

1. Take an existing LLM (like GPT-4).
2. Feed it your domain-specific data (e.g., medical records, customer support logs).
3. Update the model's parameters so it “learns” to generate more accurate, relevant responses for your use case.

Fine Tuning for specialized task

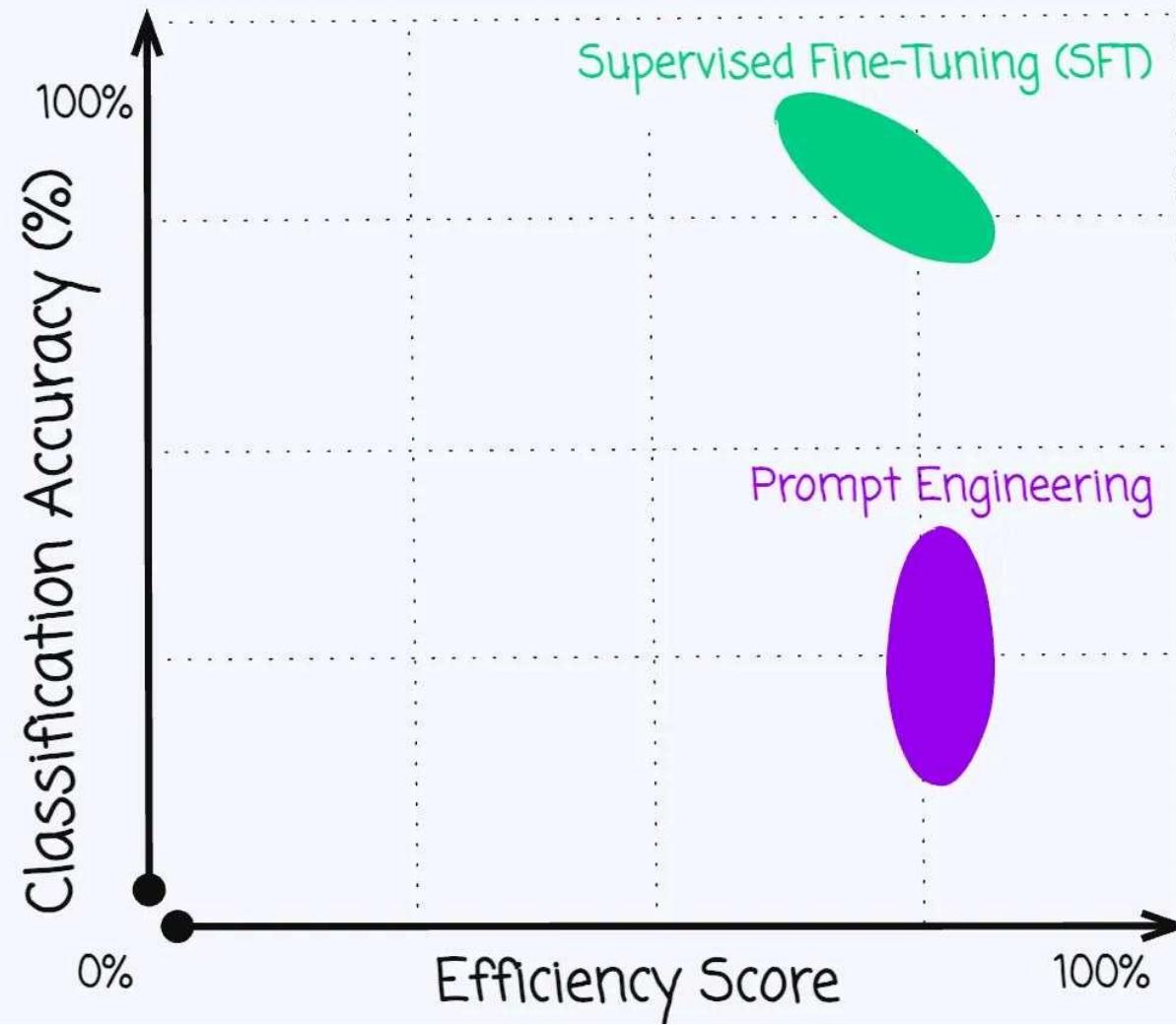


Fine Tuning

- **Domain-Adaptive Pre-training (DAP)**
- Despite its name, DAP can be used as an intermediate strategy between general model pretraining and task-specific fine-tuning of base LLMs inside RAG. It consists in utilizing a domain-specific corpus to have the model gain a better understanding of a certain domain, including jargon, writing styles, etc. Unlike conventional fine-tuning, it may still use a relatively large dataset, and it often is done before integrating the LLM with the rest of the RAG system, after which more focused and task-specific fine-tuning on smaller datasets would take place instead.
- **Retrieval Augmented Fine-Tuning**
- This is an interested and more RAG-specific fine-tuning strategy, whereby the LLM is specifically retrained on examples that incorporate both the retrieved context — augmented LLM input — and the desired response. This makes the LLM more skilled at leveraging and optimally utilizing retrieved knowledge, generating responses that will better integrate that knowledge. In other words, through this strategy, the LLM gets more skilled in properly using the RAG architecture it sits on.

Small Language Models for Text Classification

Source: <https://arxiv.org/abs/2505.16078>



Best Choice ★

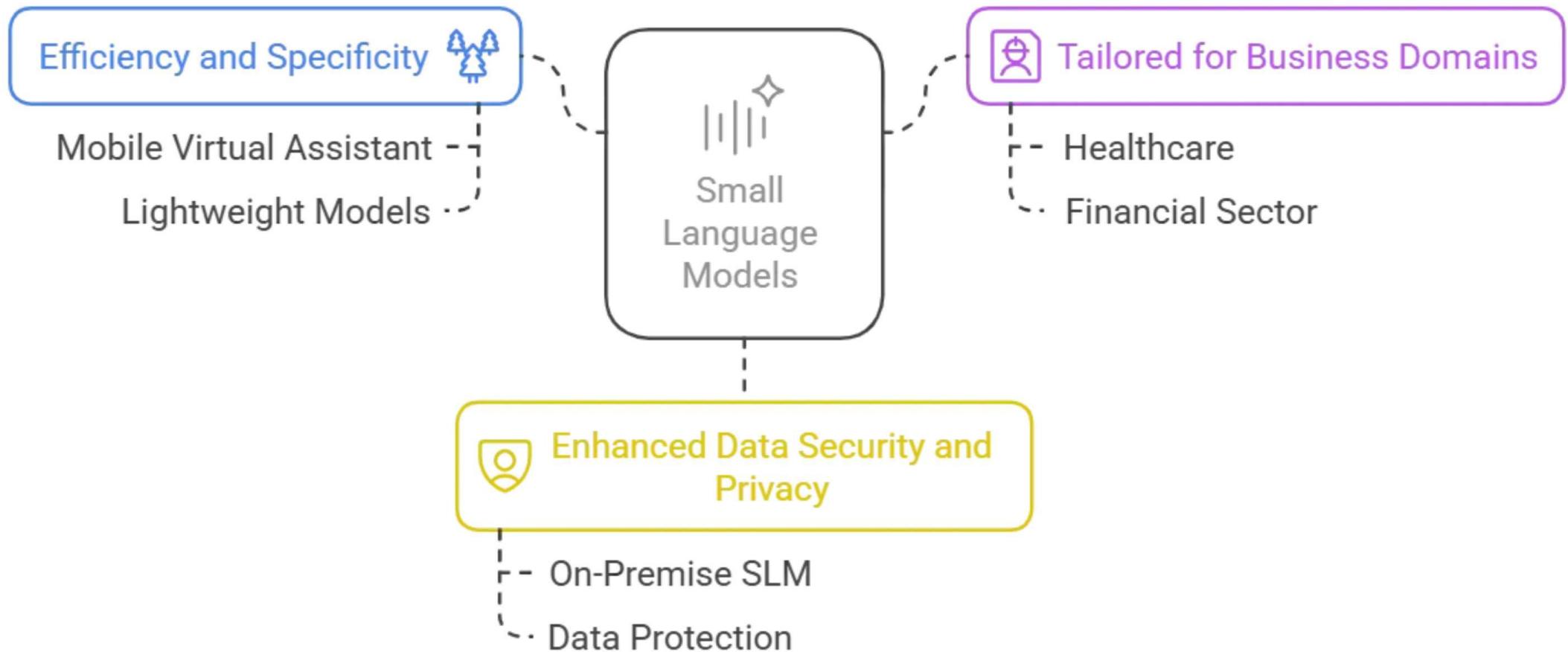
Supervised Fine-Tuning on 1B models achieves 89% accuracy with optimal efficiency.

Most Efficient 🚀

Prompt Engineering (few-shot) is fast and lightweight but limited performance.

💡 Key Insight: 1B models with supervised fine-tuning for near-perfect accuracy while using 5-10x less memory than larger alternatives.

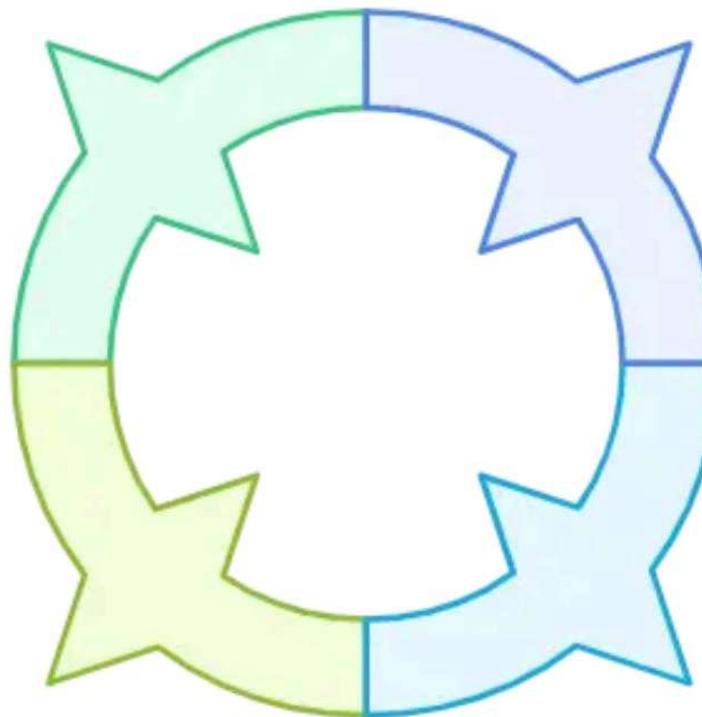
SLMs (Small Language Models)



BioBERT in Clinical Decision-Making

BioBERT enhances medical research with precise language processing.

Healthcare



Domain-Specific Models

Micro Models

Customer Support

E-commerce Chatbot

E-commerce chatbot improves customer satisfaction with quick responses.

SLM Creation – via distillation

The larger model (teacher) is trained
on a dataset.



The smaller model (student) is
trained to mimic the outputs of the
teacher model.



A specialized loss function measures
the difference between the outputs
of the two models, guiding the
student's learning.

SLM Creation – via distillation

Response-Based

The student model learns to replicate the final output layer of the teacher model, often using "soft targets" for more nuanced information.

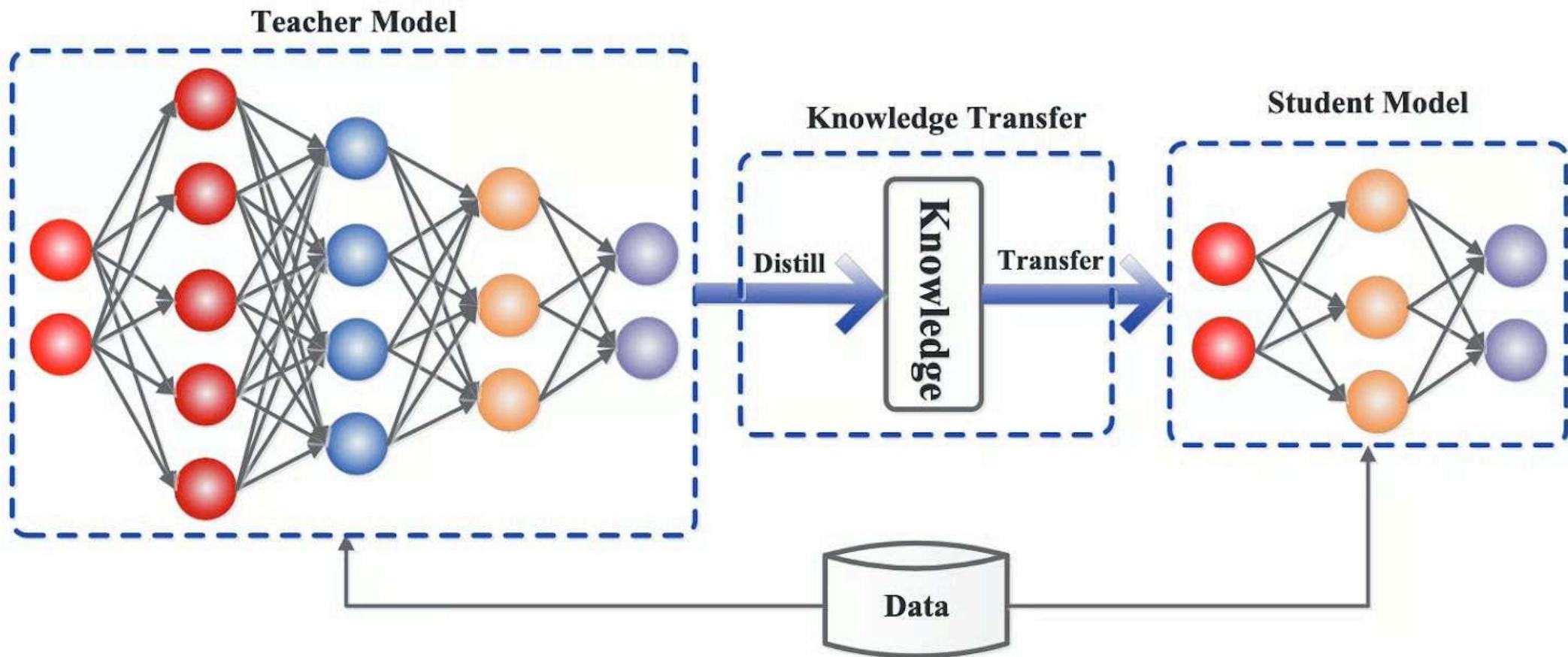
Feature-Based

Focuses on replicating the intermediate layers of the teacher model, helping the student extract similar patterns from data.

Relation-Based

Trains the student to understand relationships between different parts of the teacher model, emulating complex reasoning processes.

SLM Creation – via distillation



SLM Creation – via pruning

A large model is fully trained.



Parameters that contribute least to
the model's performance are
identified.



These less important parameters are
removed, "pruning" the model to a
smaller size.



The pruned model is often fine-tuned
to regain any lost performance.

SLM Creation – via quantization

The original model uses high-precision floating-point numbers for its parameters.



These parameters are converted to lower-precision representations.



This reduction in precision leads to smaller model sizes and faster computation.

ASPECT

LARGE LANGUAGE MODELS

SPECIALIZED LANGUAGE MODELS

Model Size

Very large, often exceeding billions of parameters

Smaller, typically ranging from millions to hundreds of millions of parameters

Training

Trained on vast datasets using deep learning techniques

Can be trained on smaller, more specific datasets

Data Specificity

General-purpose, designed to handle a wide range of topics and contexts

Highly specialized, optimized for specific domains or tasks

Resource Requirements

High computational power, memory, and infrastructure needed for training and use

More efficient, requiring less computational power and memory

Performance

Excels in generating complex, nuanced text across diverse topics

Performs well on specific tasks but may struggle with broader generalization

Cost of Implementation

Expensive to build and maintain due to extensive resources and infrastructure

More cost-effective, requiring fewer resources for training and deployment

Risk of Hallucinations

Higher chance of generating incorrect or irrelevant information (hallucinations)

Lower risk of hallucinations due to domain-specific focus

Enterprise Usability

May face issues accessing data behind firewalls, especially in enterprise contexts

More easily integrated into enterprise environments, especially for specific tasks

Examples

OpenAI's GPT-4o, Google's Bard

Microsoft's phi-3-mini, OpenAI's GPT-4o mini

LLMs for IR - problems

Without RAG

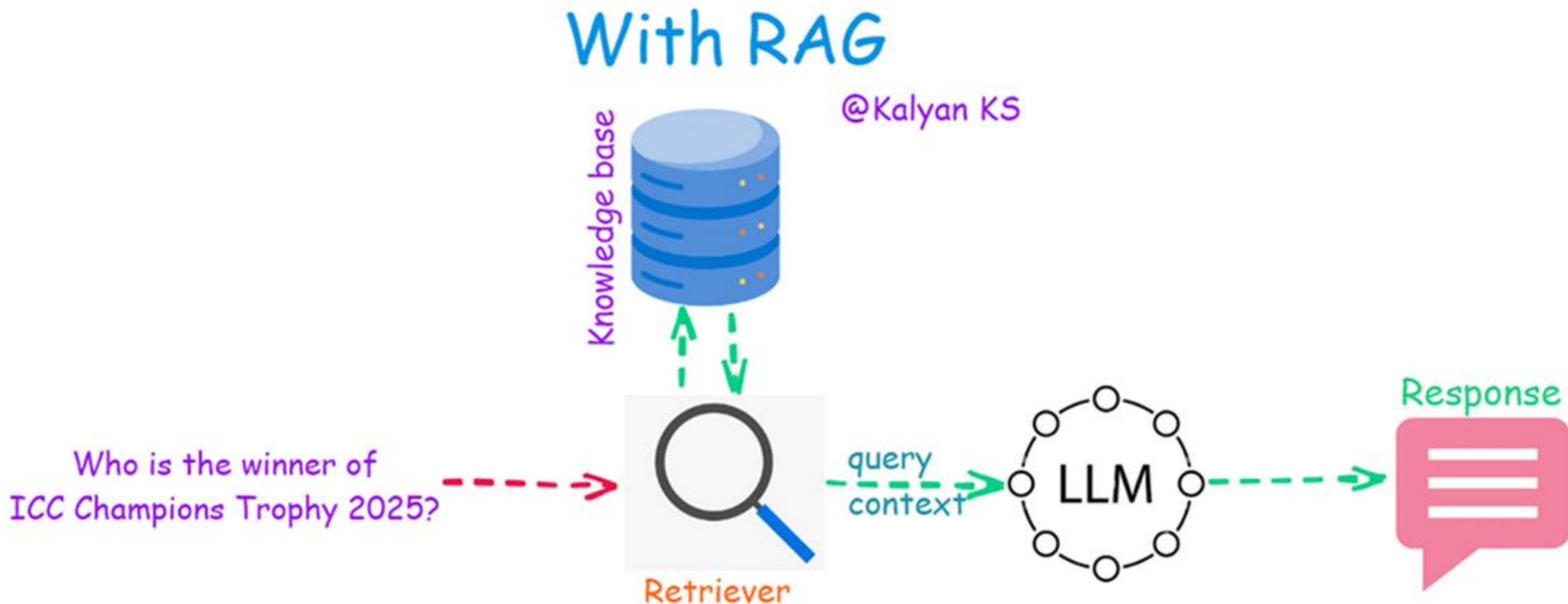
@Kalyan KS

Who is the winner of
ICC Champions Trophy 2025?
(query related to a latest event beyond
LLM knowledge cut-off date)

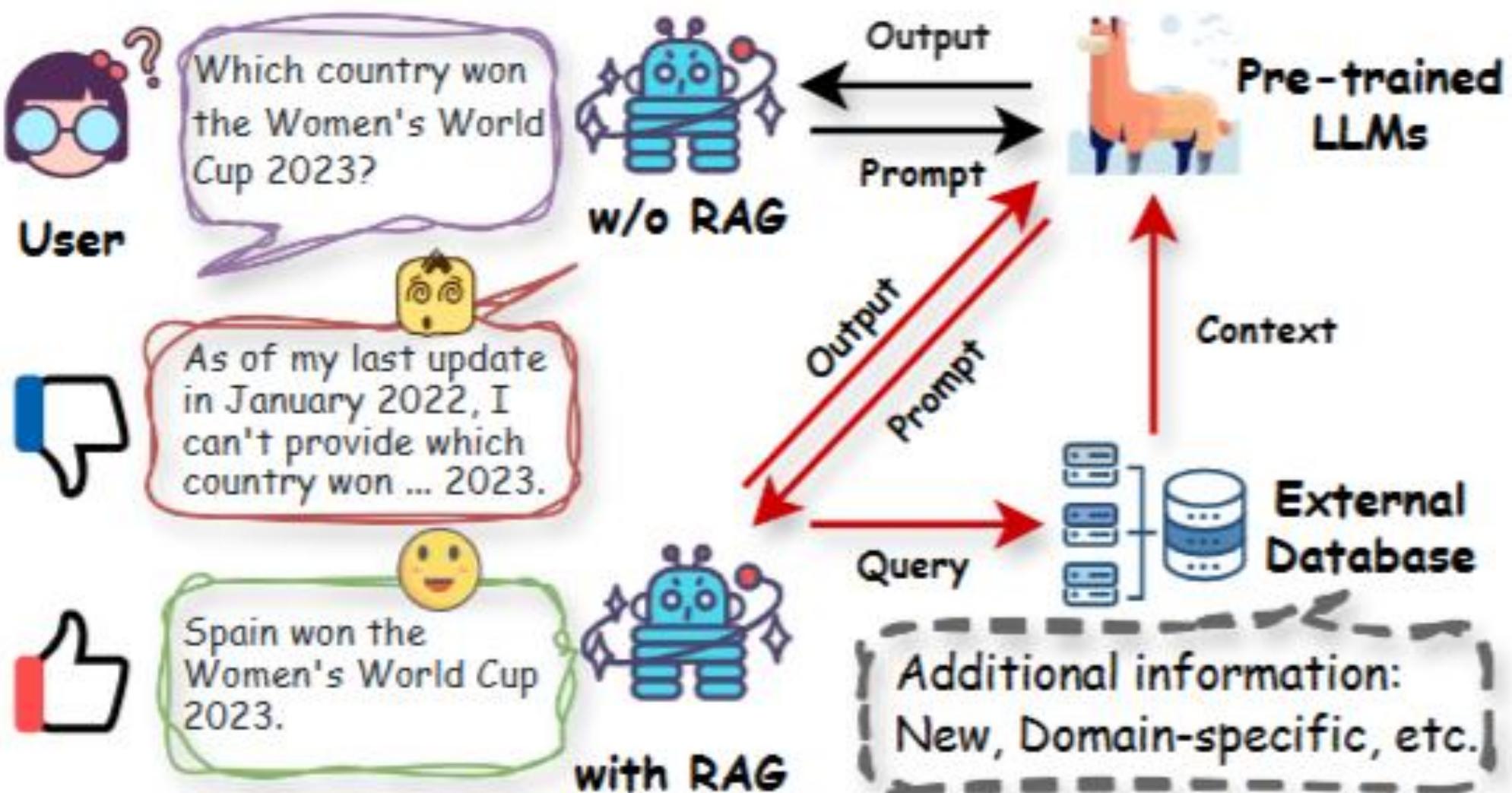


No Answer
or
Wrong Answer (Hallucination)

RAG - Motivation



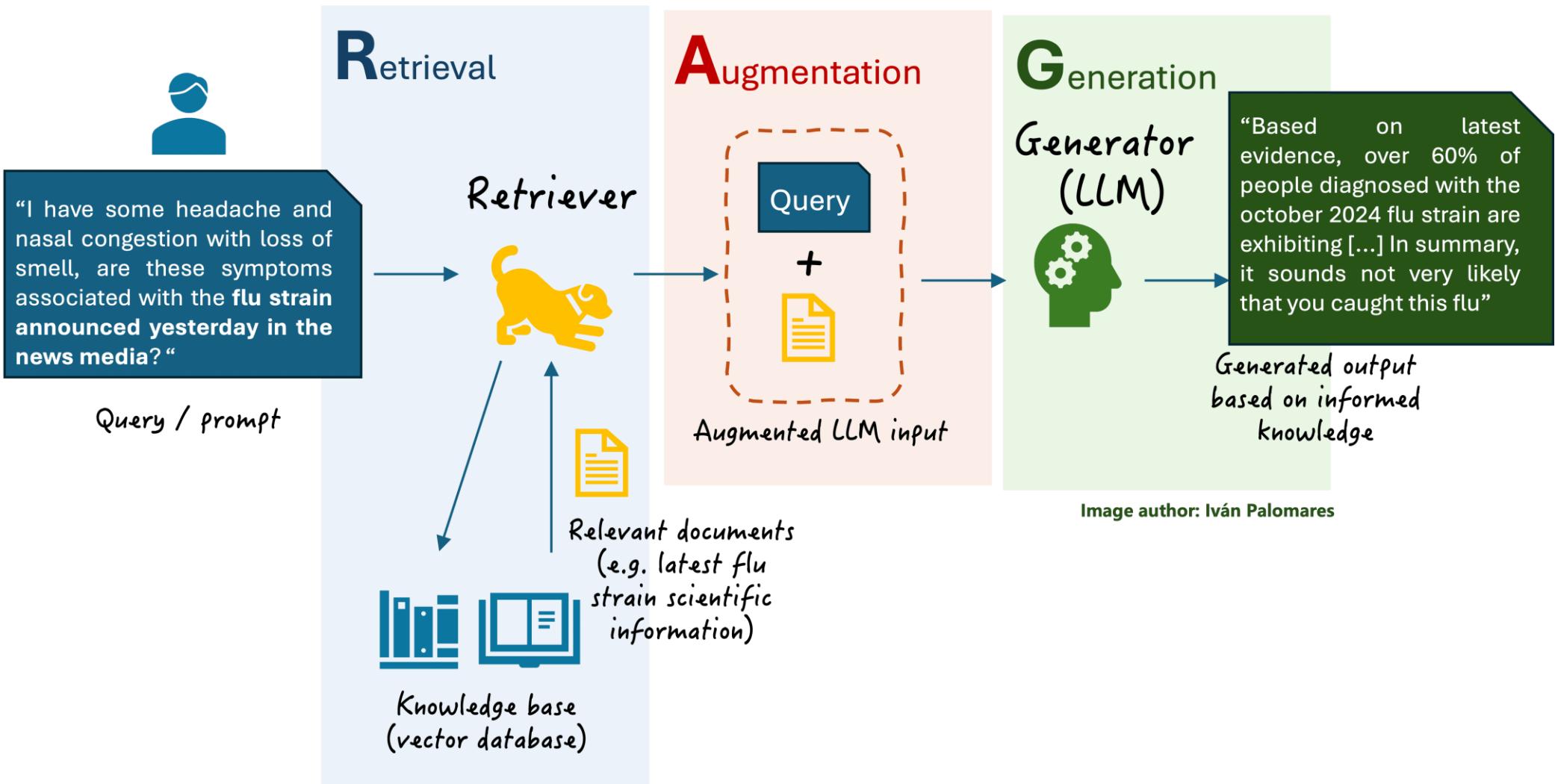
Adding a RAG



RAG – typical key data-related components

- **Retrieval:** a component called **retriever** accesses the vector database to find and retrieve relevant documents to the user query.
- **Augmentation:** the original user query is augmented by incorporating contextual knowledge from the retrieved documents.
- **Generation:** the LLM -also commonly referred to as **generator** from an RAG viewpoint- receives the user query augmented with relevant contextual information and generates a more precise and truthful text response.

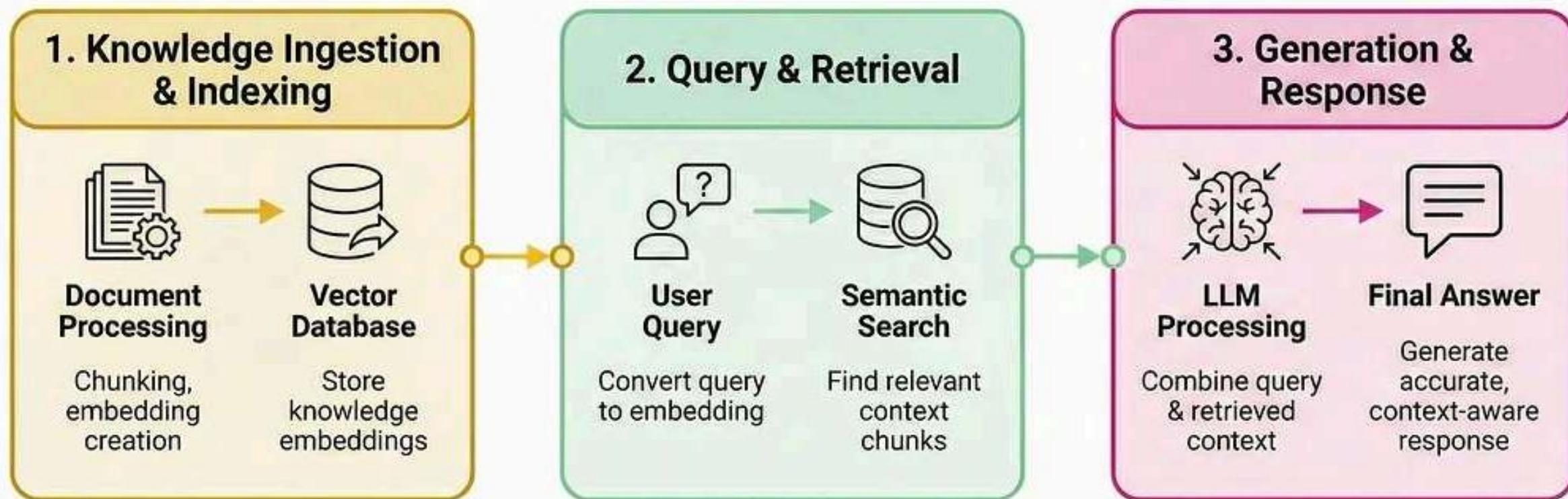
Typical RAG system



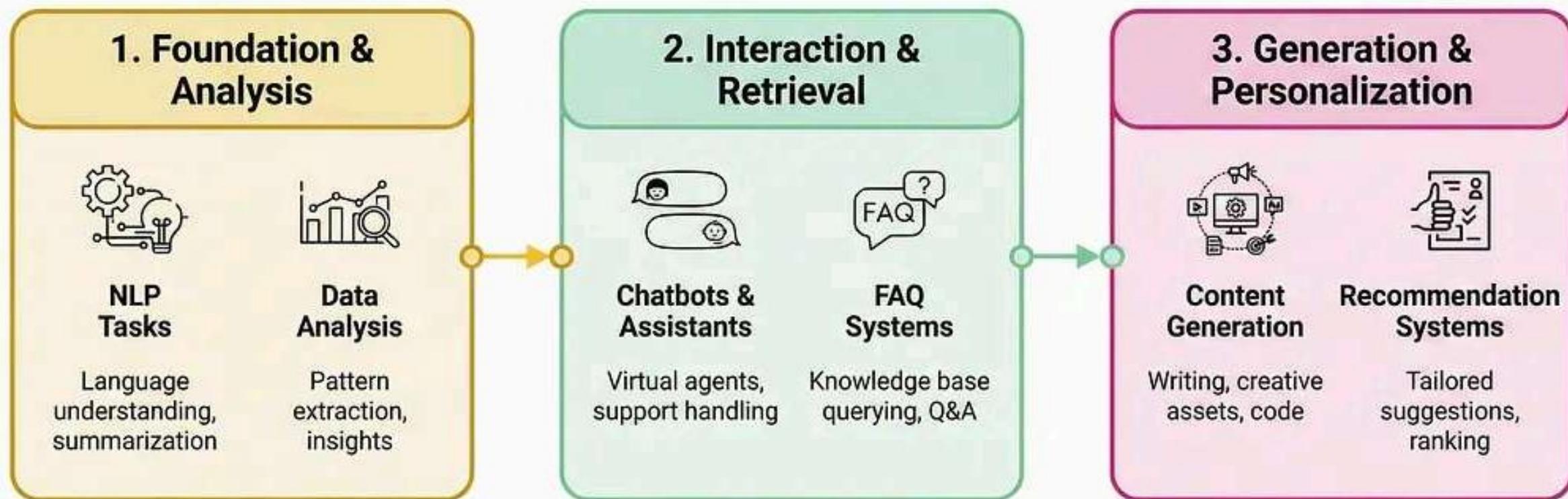
Prompting

- Prompting in LLMs refers to the technique of providing a specific instruction or context to guide the model's generation of text. The prompt serves as a conditioning signal and influences the language generation process of the model. Existing prompting strategies can be roughly categorized into three groups: zero-shot prompting, few-shot prompting, and chain-of-thought (CoT) prompting

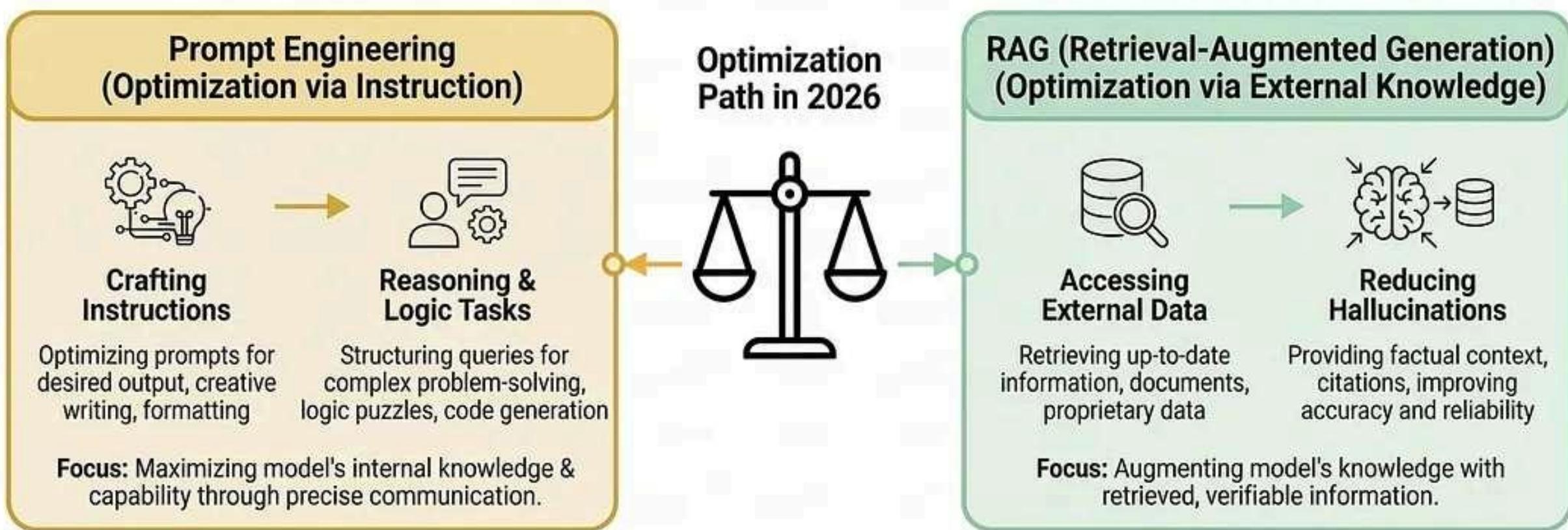
How RAG (Retrieval-Augmented Generation) Works



Where to use Prompt Engineering?

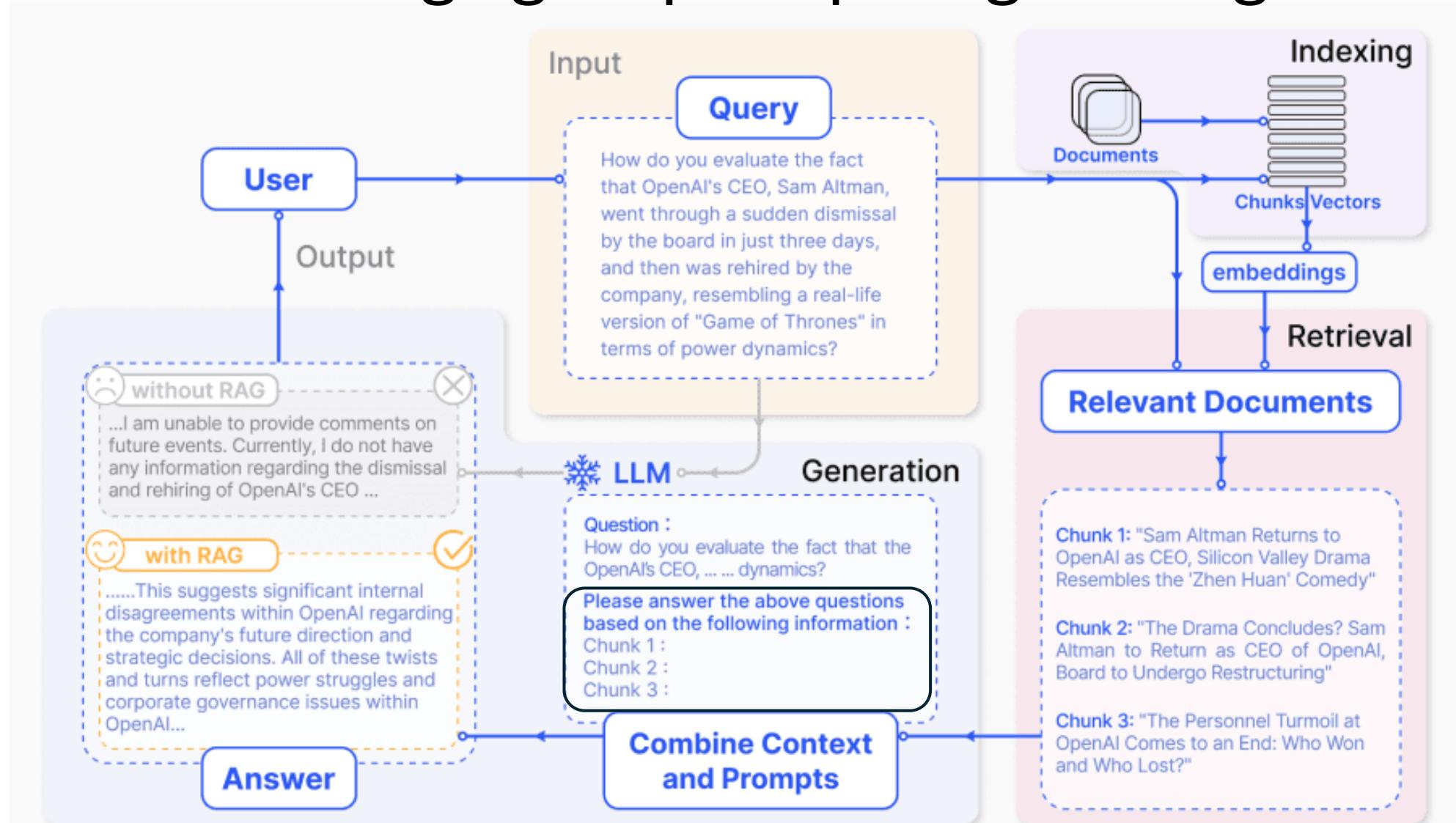


Prompt Engineering vs RAG

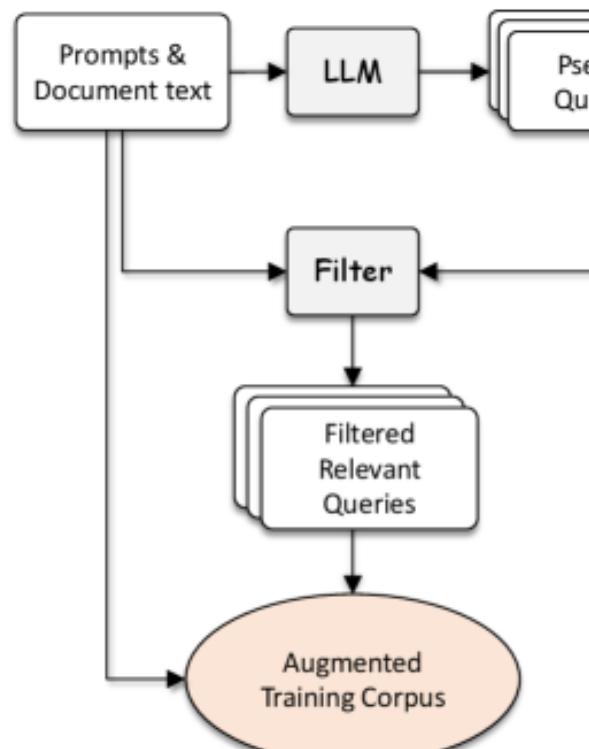
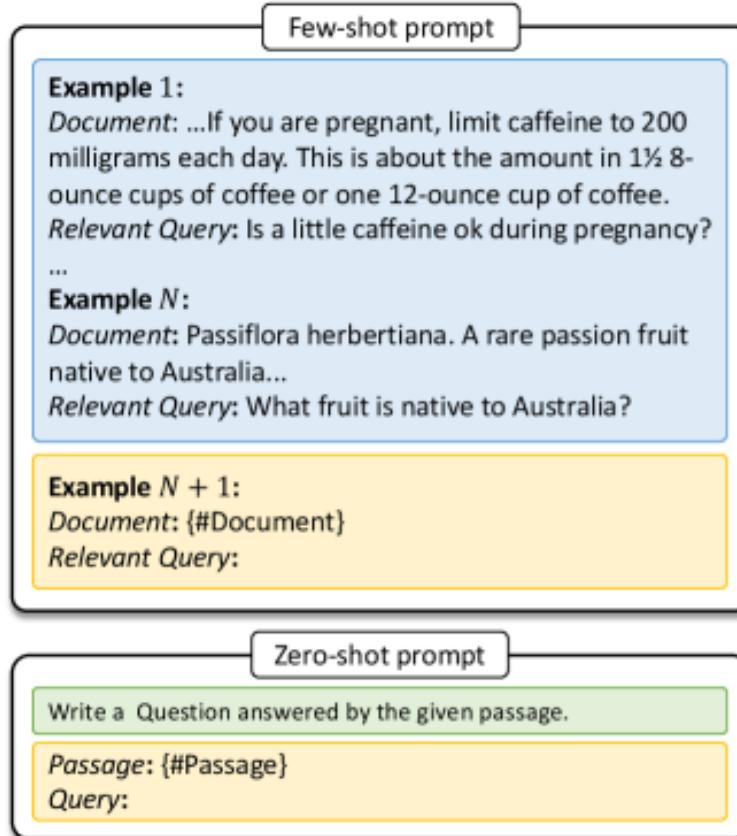


Choose Your Path Based on Data Needs, Accuracy Requirements, and Task Complexity.

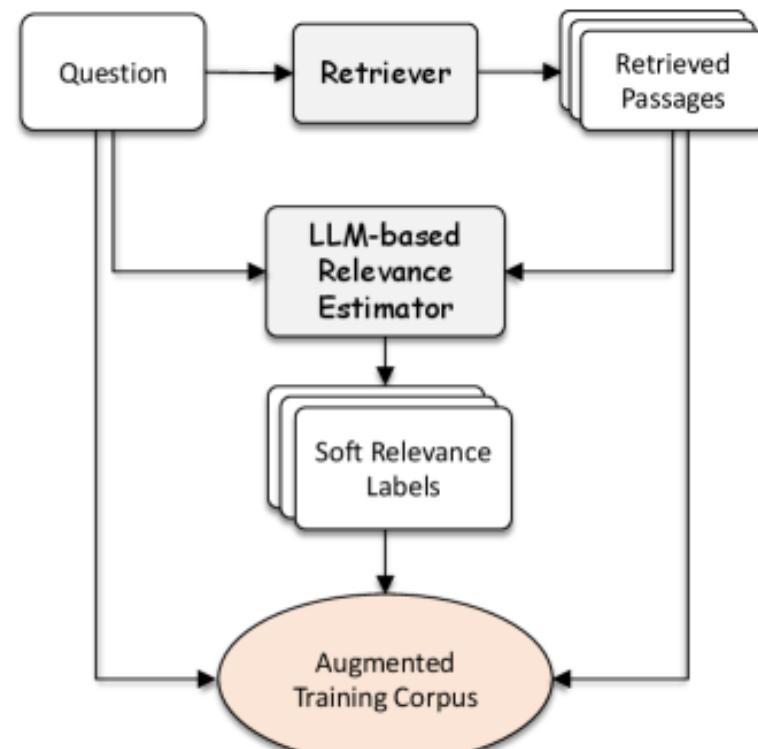
AI Chat Bot (RAG system with large close LLM) – Leveraging for prompt engineering



Leveraging LLMs to Generate Search Data on an AI Chat Bot



Framework of pseudo query generation

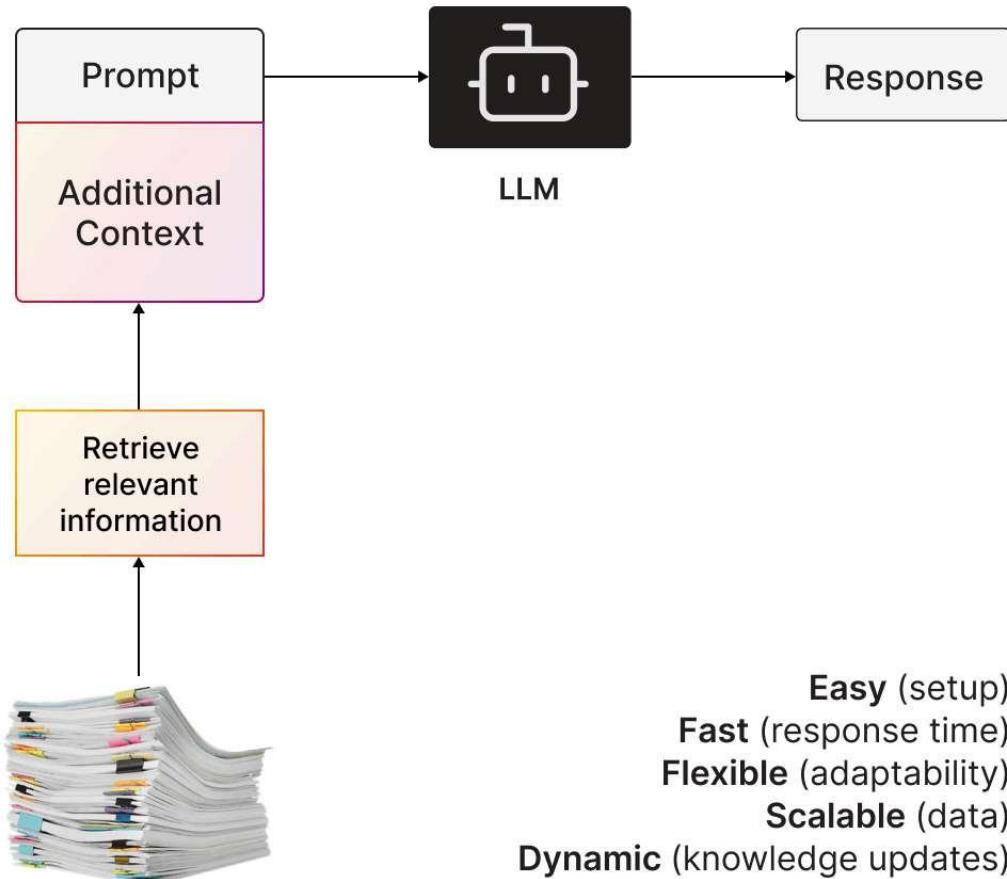


Framework of relevance label generation

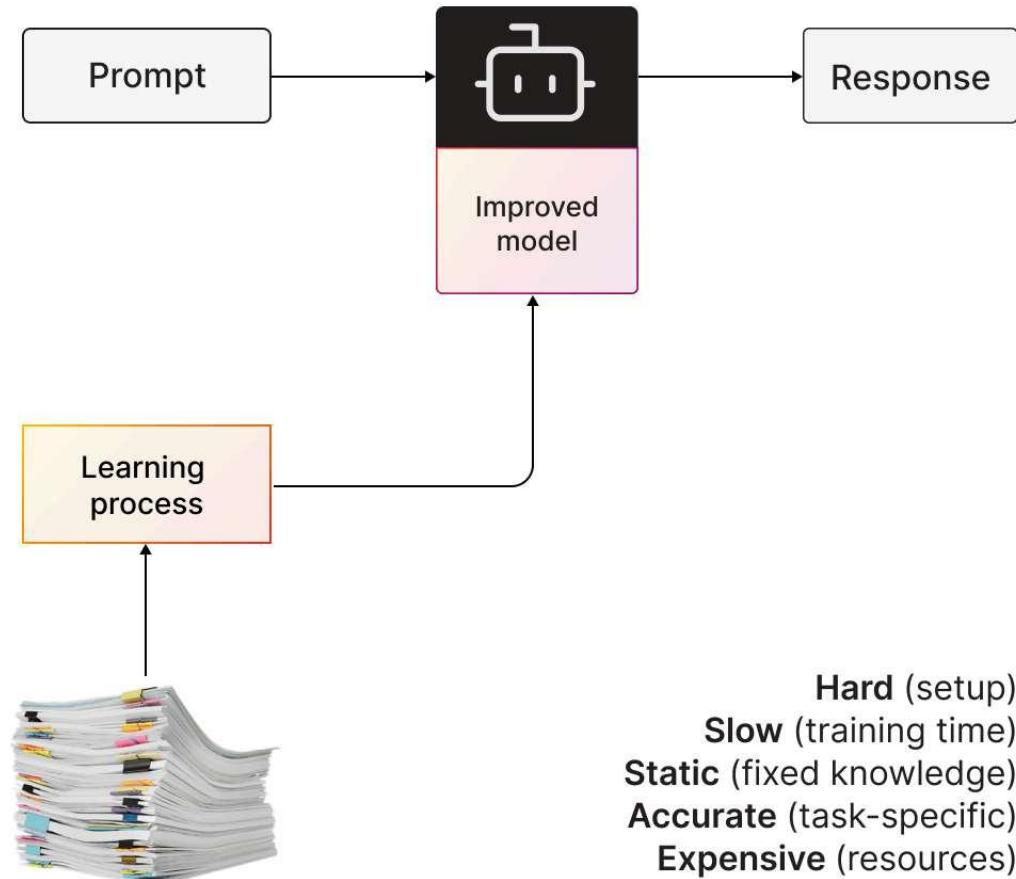
Fine-Tuning vs. RAG

Feature	Fine-Tuning	RAG
Cost	High: Requires GPUs and training time.	Low: Minimal hardware requirements.
Setup Time	Long: Training takes weeks or months.	Short: RAG can be implemented quickly.
Maintenance	High: Frequent retraining required.	Low: Automatically updates with new data.
Flexibility	Low: Fixed to the trained model.	High: Can adapt to any LLM or dataset.
Use Cases	Repeated, predictable queries.	Dynamic, real-time information needs.

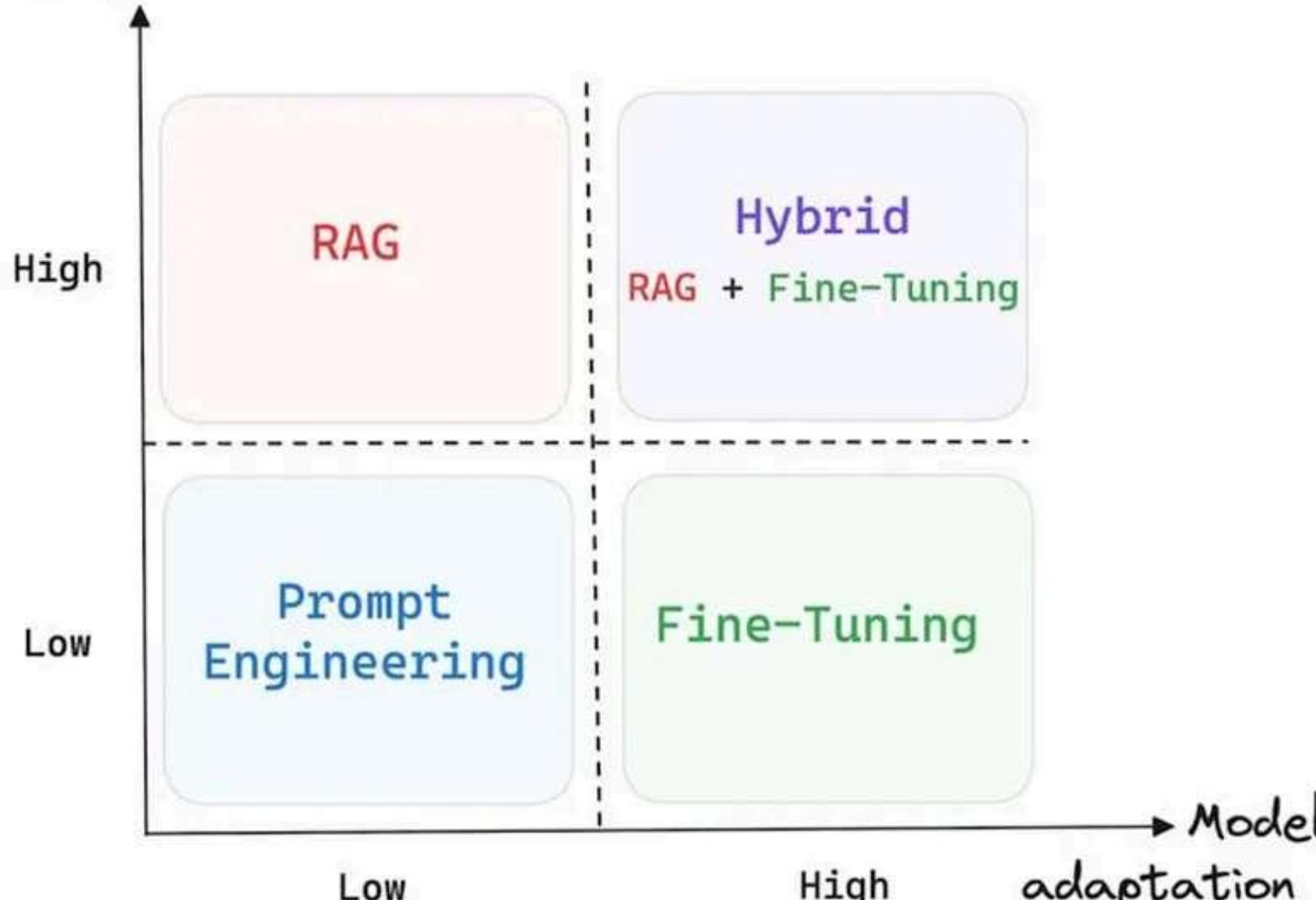
Retrieval-augmented generation (RAG)



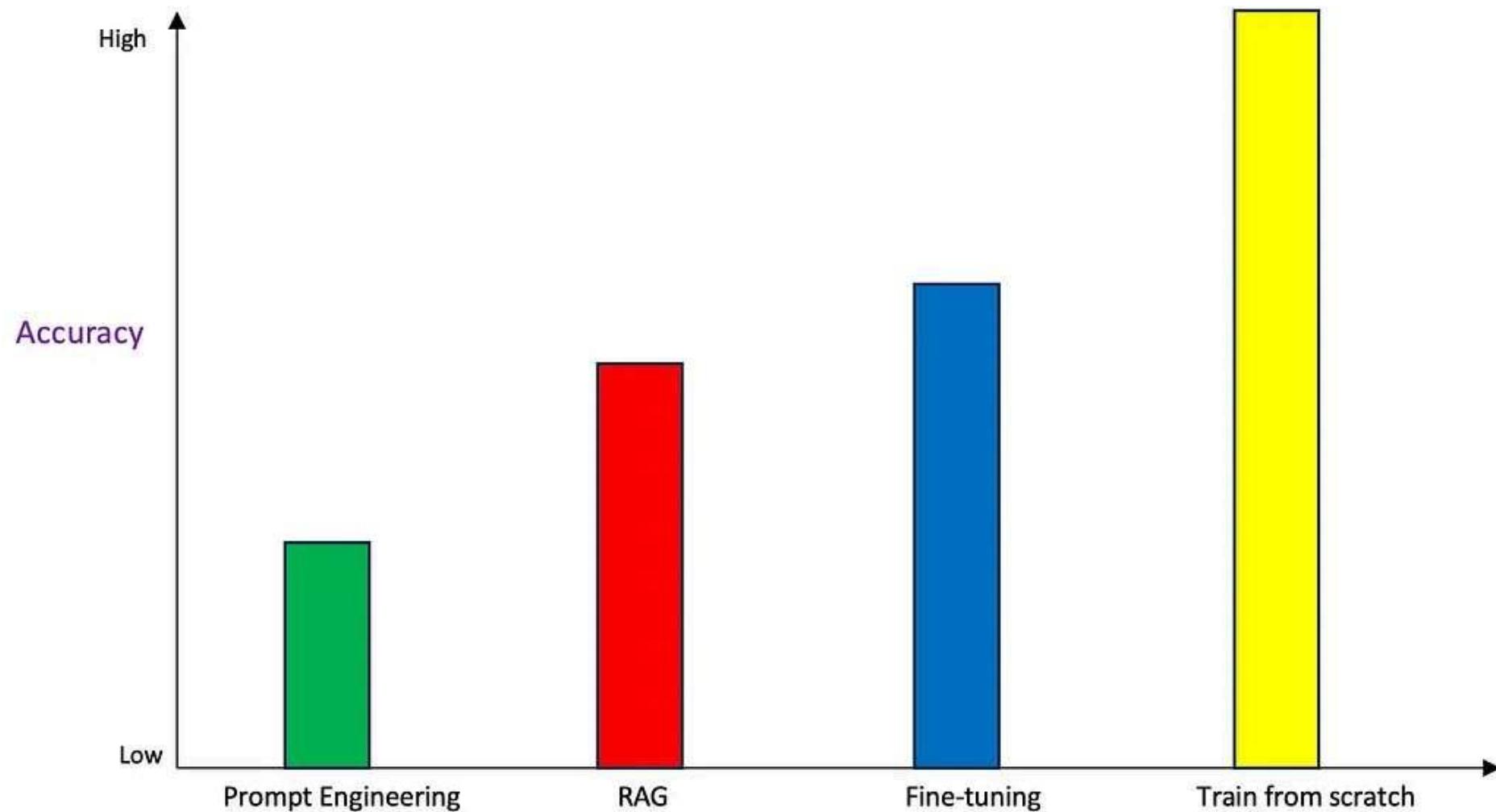
Fine-tuning the Model

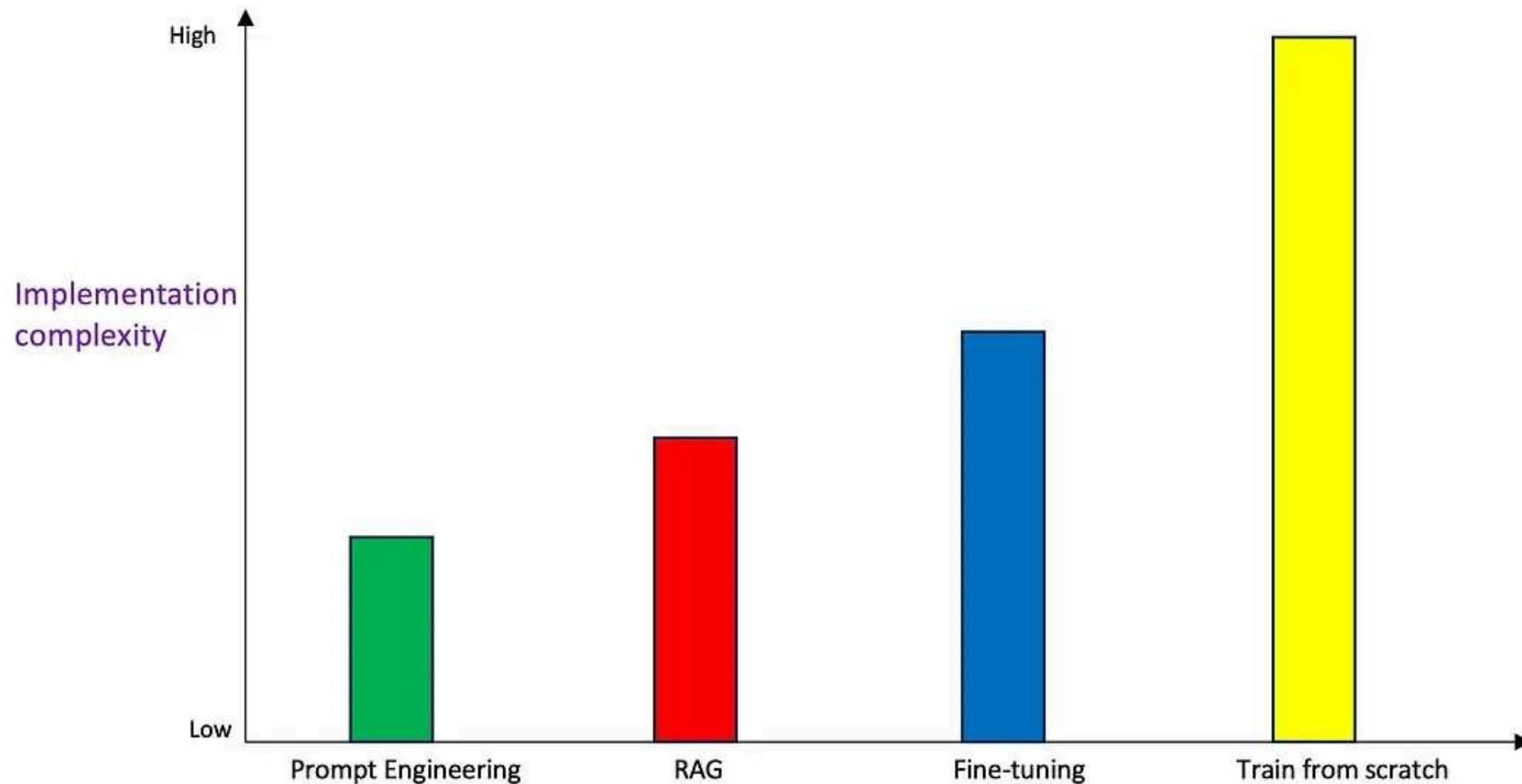


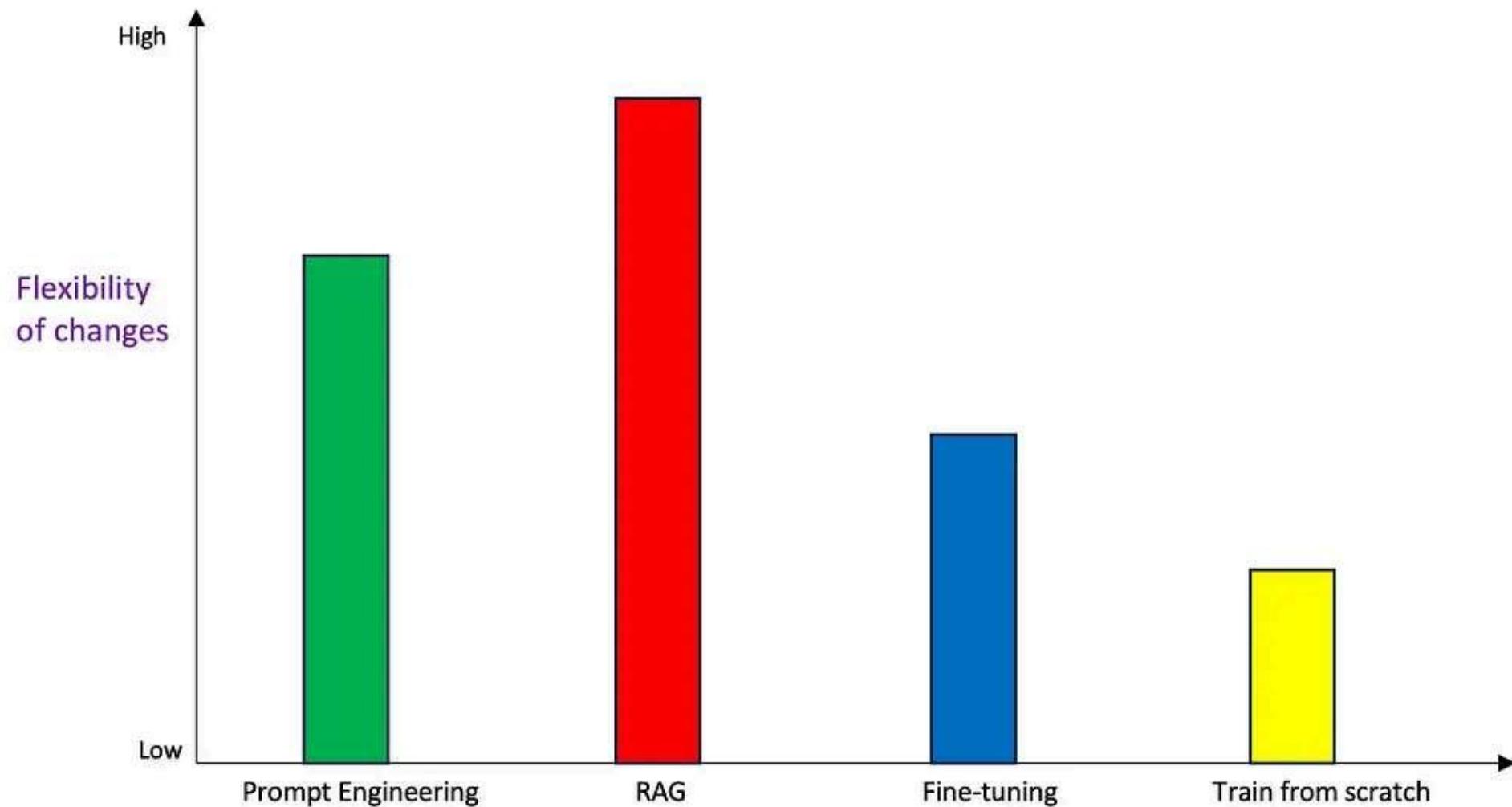
External Knowledge Required



adaptation required
(eg. behaviour, writing style,
vocabulary)







The Classic RAG Workflow

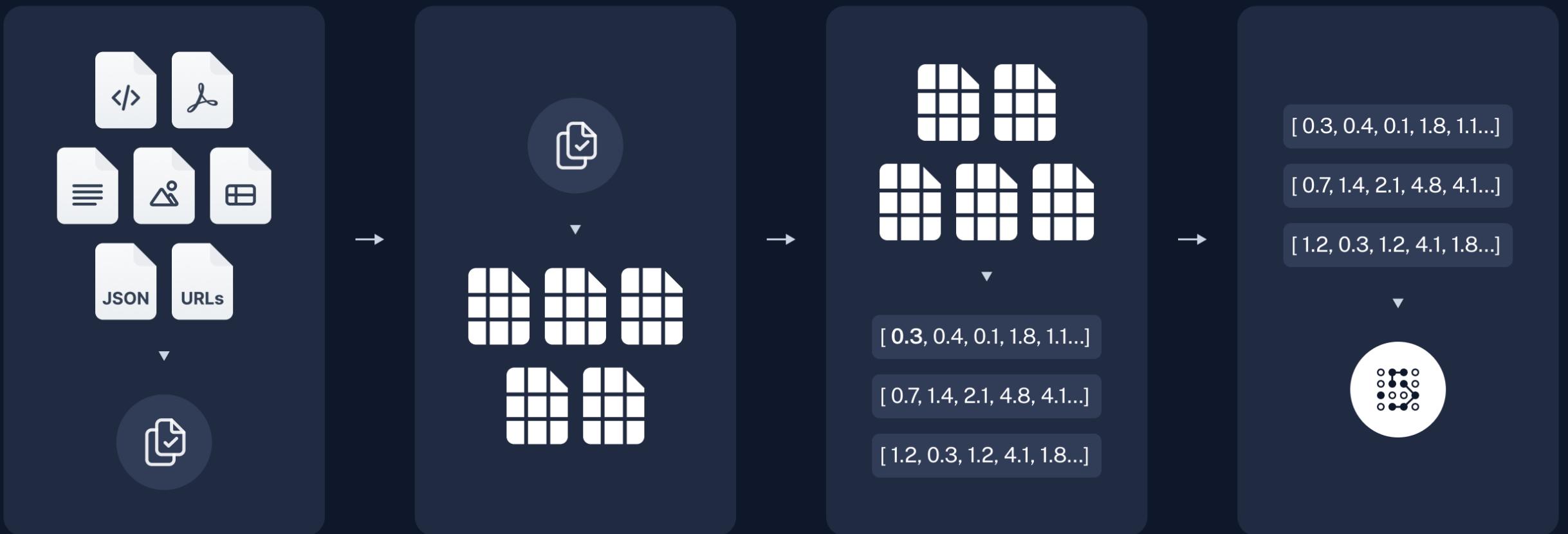
- A typical RAG system (depicted in the diagram below) handles three key data-related components:
- An **LLM** that has acquired knowledge from the data it has been trained with, typically millions to billions of text documents.
- A **vector database**, also called **knowledge base** storing text documents.
 - Vectors represent words, sentences, or entire documents, maintaining key properties of the original text such that two similar vectors are associated with words, sentences, or pieces of text with similar semantics.
 - Storing text as numerical vectors enhances the system's efficiency, such that relevant documents are quickly found and retrieved.
- **Queries or prompts** (for AI chat bots) formulated by the user in natural language.

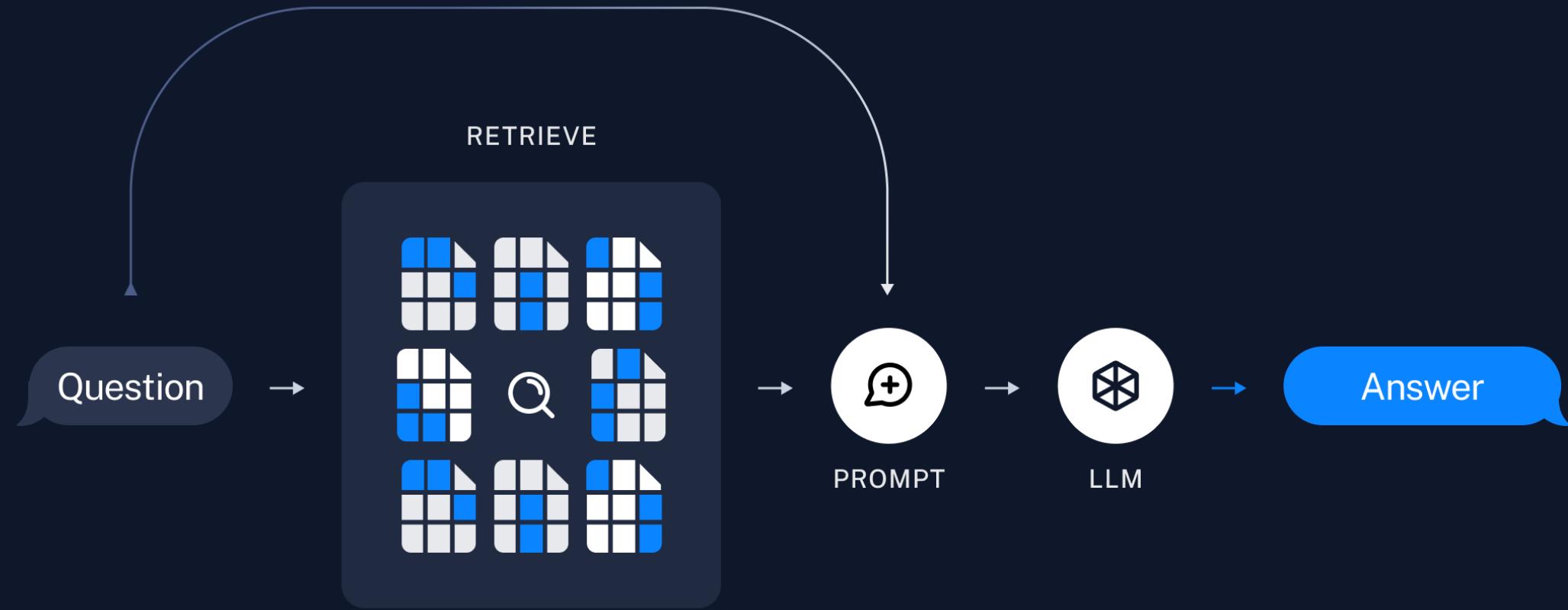
RAG - Ingestion

- In simple traditional RAG, you'll retrieve data from a vector database, using semantic search to find the true meaning of the user's query and retrieve relevant information instead of simply matching keywords in the query.

But before we can retrieve the data, you have to ingest the data.
Here are steps to get data into your database:

- Chunk the data
- Create vector embeddings
- Load data into a vector database





vs. Information Retrieval (IR)

- Document Search
- Ranking Only

vs. Question Answering (QA)

- Fixed Dataset
- Span Extraction

RAG System

RETRIEVAL

- Knowledge Base
- Embeddings
- Retriever
- Chunking

User Query

Retrieved Documents

GENERATION

- LLM
- Prompt Constructor
- Post-Processing

vs. Recommendation Systems

- Item Suggestions
- User Profiles

vs. LLMs

- Parametric Memory
- Hallucinations

RAG – Ingestion – Chunking - considerations

- **Vector count:** How many chunks you will create directly affects the vector count, which in turn affects search time.
- **Vector quality:** How you chunk your data also determines the semantic coherence and search precision of each vector.
- **Index structure:** Chunk characteristics affect clustering efficiency and index transversal speed.

RAG – Ingestion - Chunking

- Several variables play a role in determining the best chunking strategy, and these variables vary depending on the use case. Here are some key aspects to keep in mind:
- **What kind of data is being chunked?** Are you working with long documents, such as articles or books, or shorter content, like tweets, product descriptions, or chat messages? Small documents may not need to be chunked at all, while larger ones may exhibit certain structure that will inform chunking strategy, such as sub-headers or chapters.
- **Which embedding model are you using?** Different embedding models have differing capacities for information, especially on specialized domains like code, finance, medical, or legal information. And, the way these models are trained can strongly affect how they perform in practice. After choosing an appropriate model for your domain, be sure to adapt your chunking strategy to align with expected document types the model has been trained on.
- **What are your expectations for the length and complexity of user queries?** Will they be short and specific or long and complex? This may inform the way you choose to chunk your content as well so that there's a closer correlation between the embedded query and embedded chunks.
- **How will the retrieved results be utilized within your specific application?** For example, will they be used for semantic search, question answering, retrieval augmented generation, or even an agentic workflow? For example, the amount of information a human may review from a search result may be smaller or larger than what an LLM may need to generate a response. These users determine how your data should be represented within the vector database.

RAG – Ingestion – Chunking – Fixed-size chunking

- This is the most common and straightforward approach to chunking: we simply decide the number of tokens in our chunk, and use this number to break up our documents into fixed size chunks.
- Usually, this number is the max context window size of the embedding model (such as 1024 for llama-text-embed-v2, or 8196 for text-embedding-3-small).
- Keep in mind that different embedding models may tokenize text differently, so you will need to estimate token counts accurately.
- Fixed-sized chunking will be the best path in most cases, and we recommend starting here and iterating only after determining it insufficient.

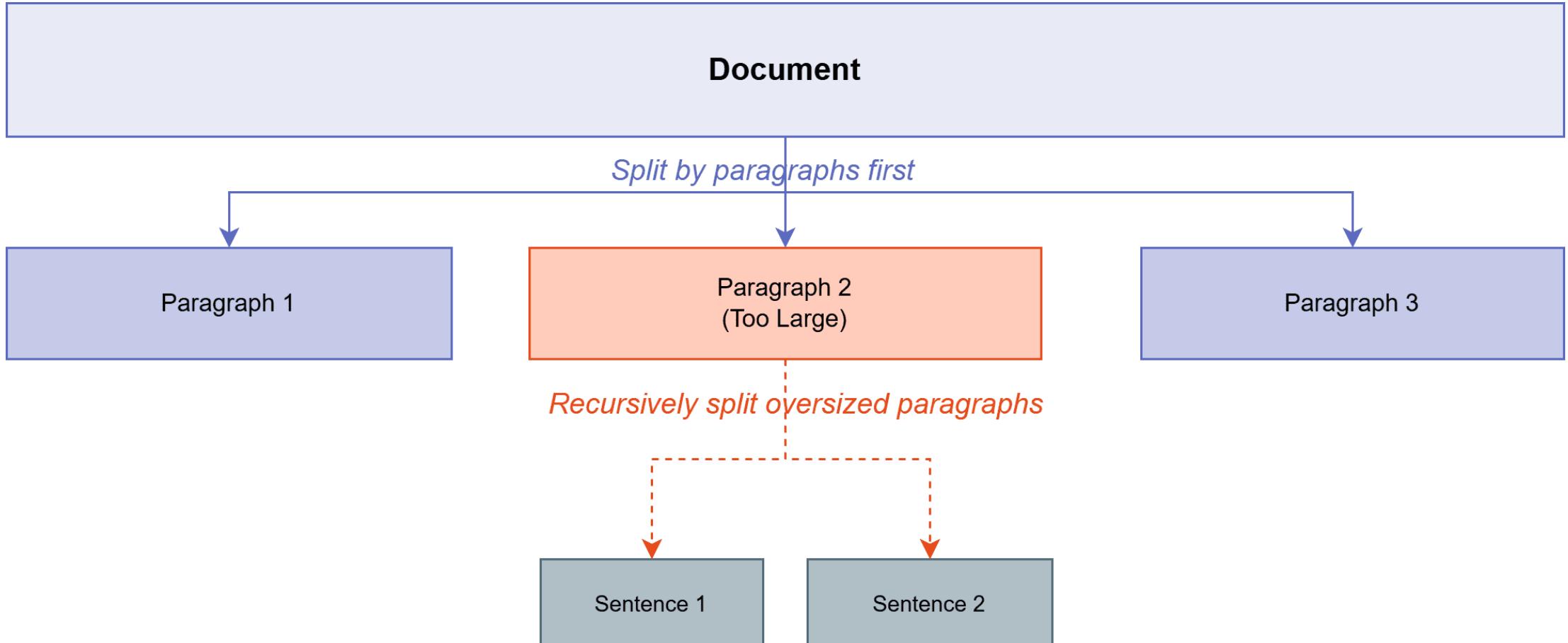
RAG – Ingestion – Chunking – Simple Sentence and Paragraph splitting

- As we mentioned before, some embedding models are optimized for embedding sentence-level content. But sometimes, sentences need to be mined from larger text datasets that aren't preprocessed. In these cases, it's necessary to use sentence chunking, and there are several approaches and tools available to do this:
- **Naive splitting:** The most naive approach would be to split sentences by periods ("."), new lines, or white space.
- **NLTK:** The Natural Language Toolkit (NLTK) is a popular Python library for working with human language data. It provides a trained sentence tokenizer that can split the text into sentences, helping to create more meaningful chunks.
- **spaCy:** spaCy is another powerful Python library for NLP tasks. It offers a sophisticated sentence segmentation feature that can efficiently divide the text into separate sentences, enabling better context preservation in the resulting chunks.

RAG – Ingestion – Chunking – Recursive Character Level Chunking

- LangChain implements a RecursiveCharacterTextSplitter that tries to split text using separators in a given order.
- The default behavior of the splitter uses the ["\n\n", "\n", " ", ""] separators to break paragraphs, sentences and words depending on a given chunk size.
- This is a great middle ground between always splitting on a specific character and using a more semantic splitter, while also ensuring fixed chunk sizes when possible.

RAG – Ingestion – Chunking – Recursive Character Level Chunking

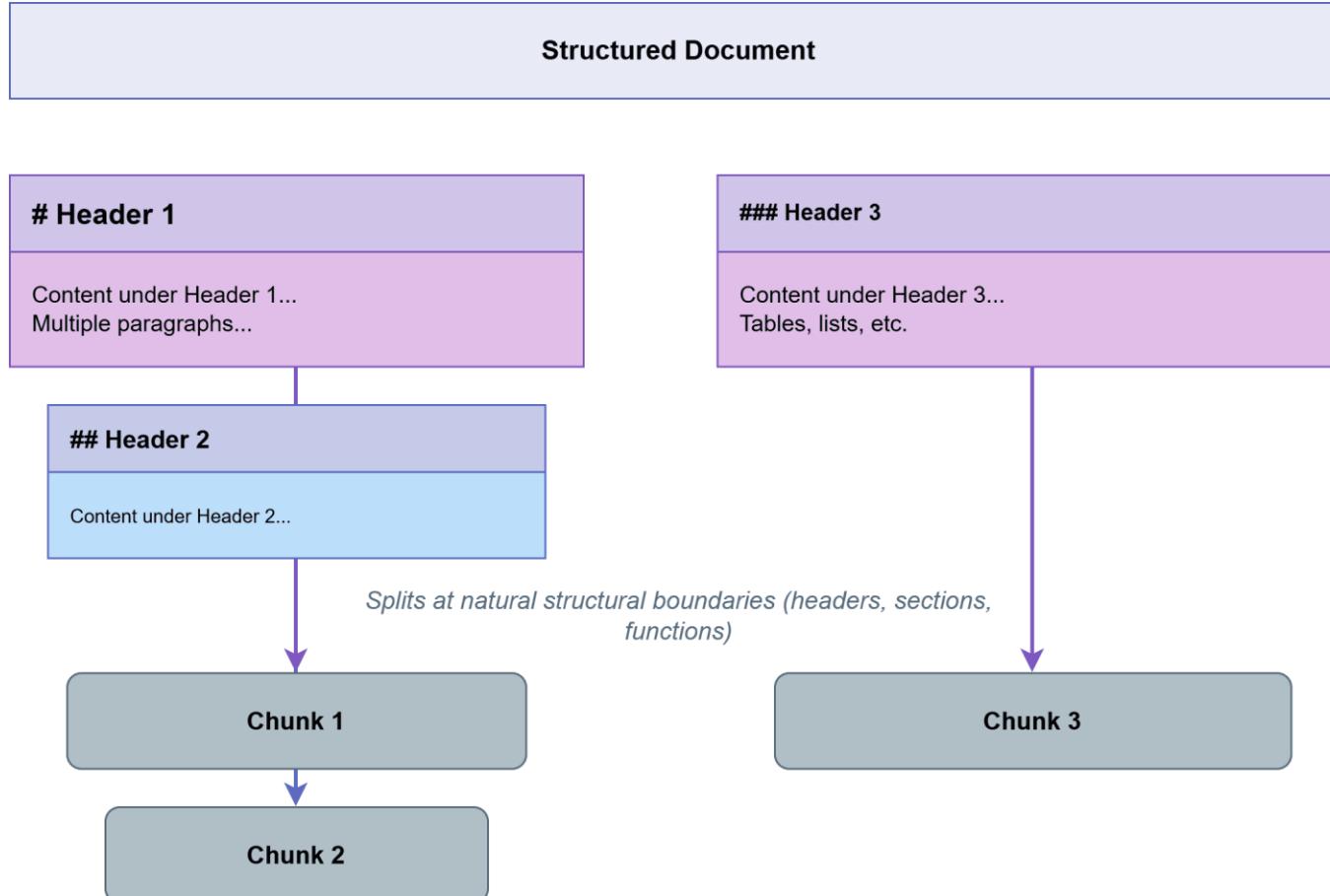


RAG – Ingestion – Chunking –

Document structure-based chunking

- When chunking large documents such as PDFs, DOCX, HTML, code snippets, Markdown files and LaTex, specialized chunking methods can help preserve the original structure of the content during chunk creation.
- PDF documents contain loads of headers, text, tables, and other bits and pieces that require preprocessing to chunk. LangChain has some handy utilities to help process these documents, while Pinecone Assistant can chunk and processes these for you
- HTML, from scraped web pages can contain tags (`<p>` for paragraphs, or `<title>` for titles) that can inform text to be broken up or identified, like on product pages or blog posts. Roll your own parser, or use LangChain splitters here to process these for chunking
- Markdown: Markdown is a lightweight markup language commonly used for formatting text. By recognizing the Markdown syntax (e.g., headings, lists, and code blocks), you can intelligently divide the content based on its structure and hierarchy, resulting in more semantically coherent chunks.
- LaTex: LaTeX is a document preparation system and markup language often used for academic papers and technical documents. By parsing the LaTeX commands and environments, you can create chunks that respect the logical organization of the content (e.g., sections, subsections, and equations), leading to more accurate and contextually relevant results.

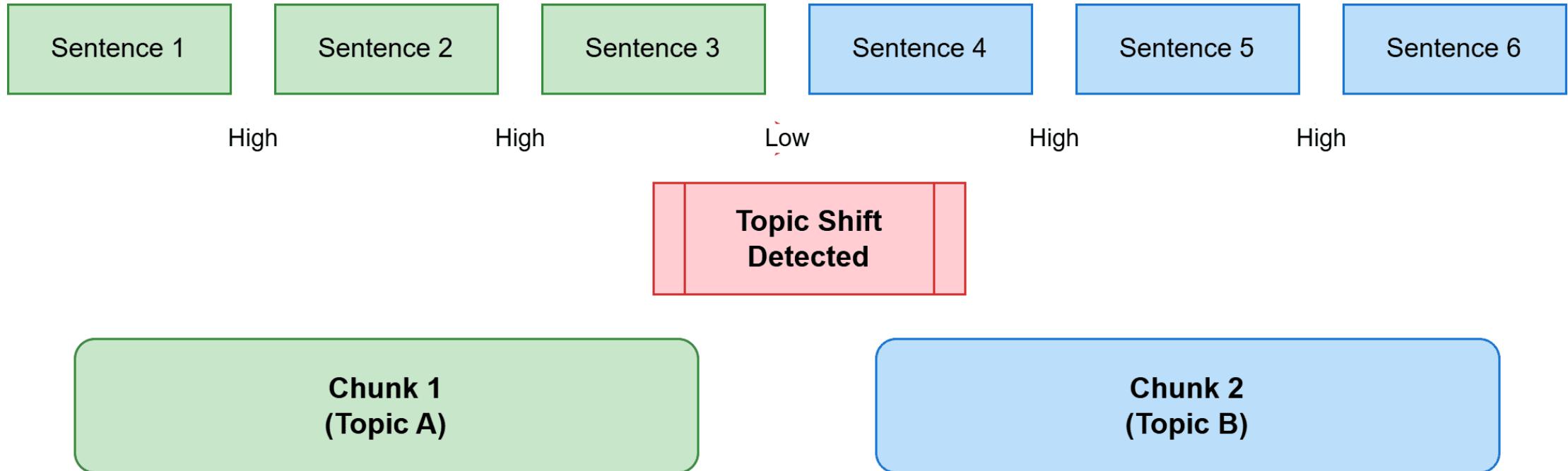
RAG – Ingestion – Chunking – Document structure-based chunking



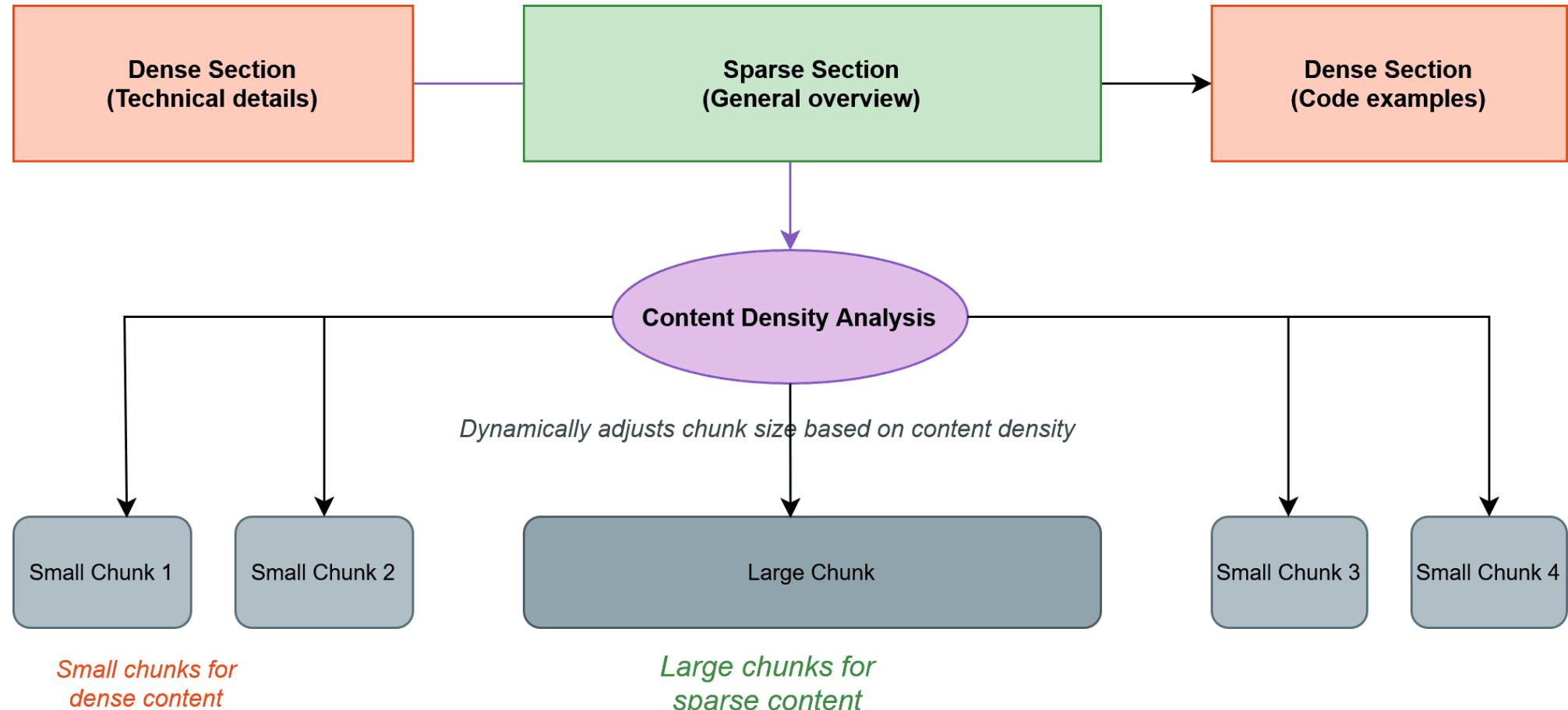
RAG – Ingestion – Chunking – Semantic Chunking

- A new experimental technique for approaching chunking was first introduced by Greg Kamradt. In his notebook, Kamradt rightfully points to the fact that a global chunking size may be too trivial of a mechanism to consider the meaning of segments within the document. If we use this type of mechanism, we can't know if we're combining segments that have anything to do with one another.
- Luckily, if you're building an application with LLMs, you most likely already can create embeddings - and embeddings can be used to extract the semantic meaning present in your data. This semantic analysis can be used to create chunks that are made up of sentences that talk about the same theme or topic.
- Semantic chunking involves breaking a document into sentences, grouping each sentence with its surrounding sentences, and generating embeddings for these groups. By comparing the semantic distance between each group and its predecessor, you can identify where the topic or theme shifts, which defines the chunk boundaries. You can learn more about applying semantic chunking with Pinecone [here](#).

RAG – Ingestion – Chunking – Semantic Chunking



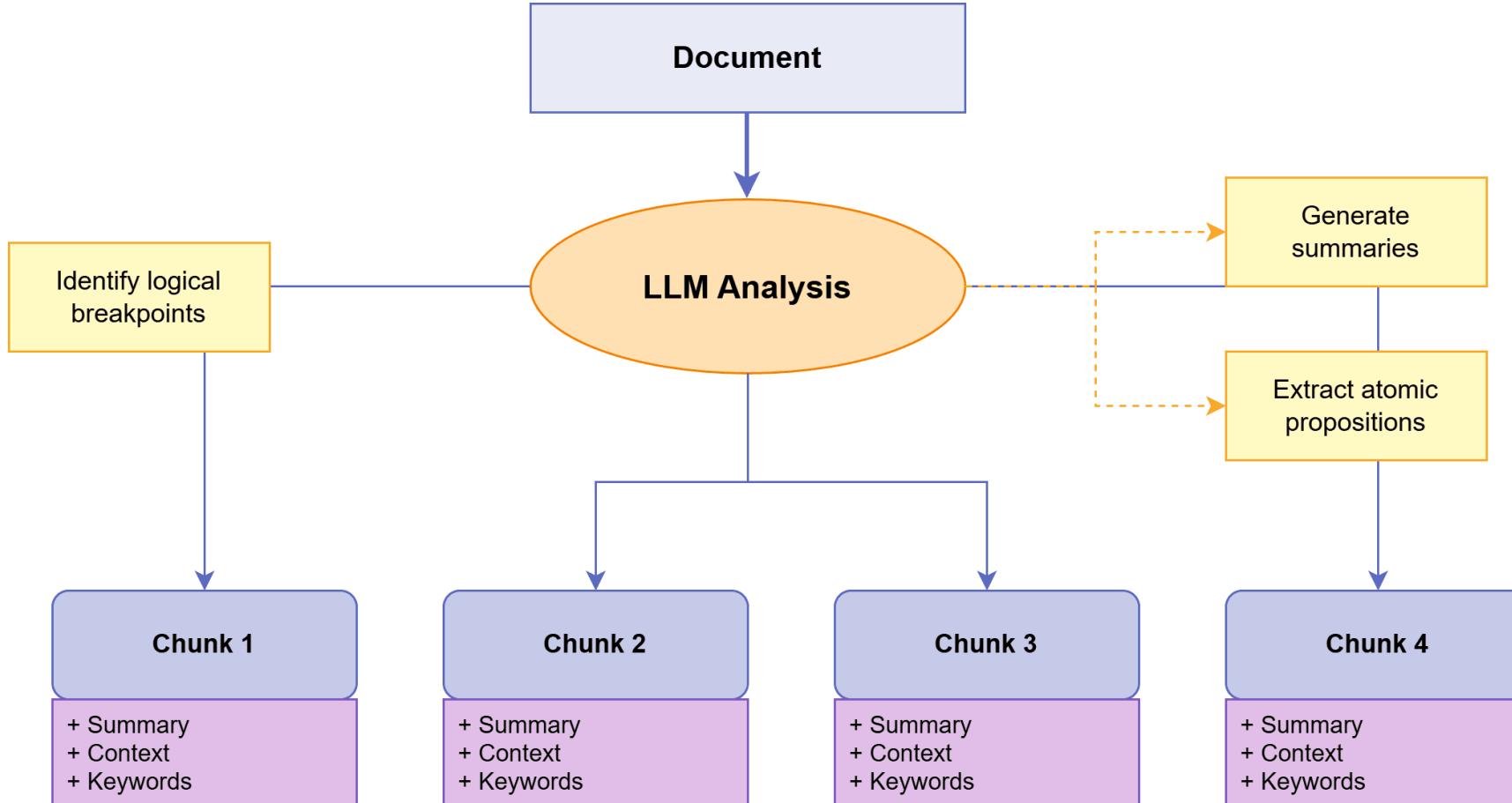
RAG – Ingestion – Chunking – Adaptive Chunking



RAG – Ingestion – Chunking – Contextual Chunking with LLMs

- Sometimes, it's not possible to chunk information from a larger complex document without losing the context entirely. This can happen when the documents are many hundreds of pages, change topics frequently, or require understanding from many related portions of the document. Anthropic introduced contextual retrieval in 2024 to help address this problem.
- Anthropic prompted a Claude instance with an entire document and its chunk, in order to generate a contextualized description, which is appended to the chunk and then embedded. The description helps retain the high-level summary meaning of the document to the chunk, which exposes this information to incoming queries. To avoid processing the document each time, it's cached within the prompt for all necessary chunks. You can learn more about contextual retrieval in our video [here](#) and our code example [here](#).

RAG – Ingestion – Chunking – Contextual Chunking with LLMs



RAG – Ingestion – Create vector embeddings

- Then, using an embedding model, you'll embed each chunk and load it into the vector database. The embedding model is a special type of LLM that converts the data chunk into a vector embedding, a numerical representation of the data's meaning. This allows computers to search for similar items based on the vector representation of the stored data.

RAG - Vector Databases

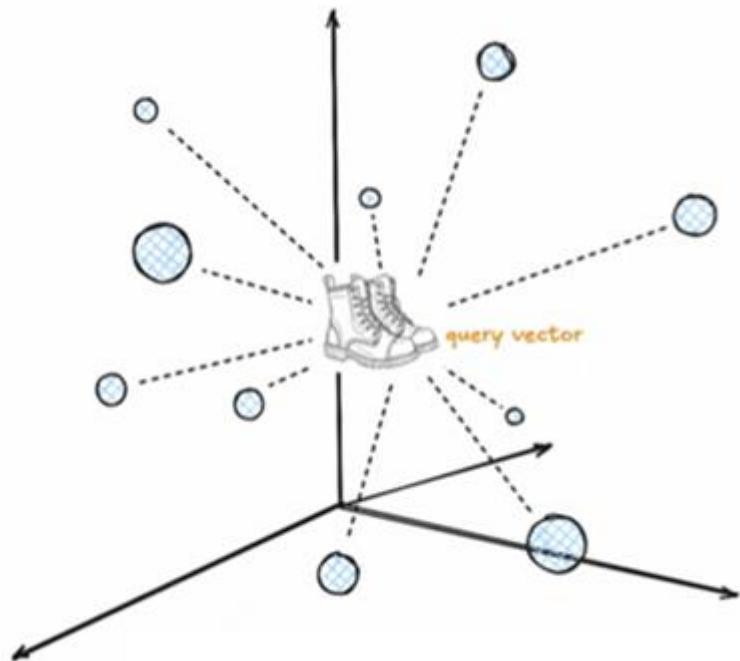
- **Vector DB** - a specialized for the high-dimensional vectors
- Enable efficient similarity-based searches
 - For instance, vectors associated with two Mediterranean restaurant reviews would be relatively close
- Vector DB in IR - documents relevant to user's query will be retrieved efficiently

RAG – Ingestion – Load data into a vector database

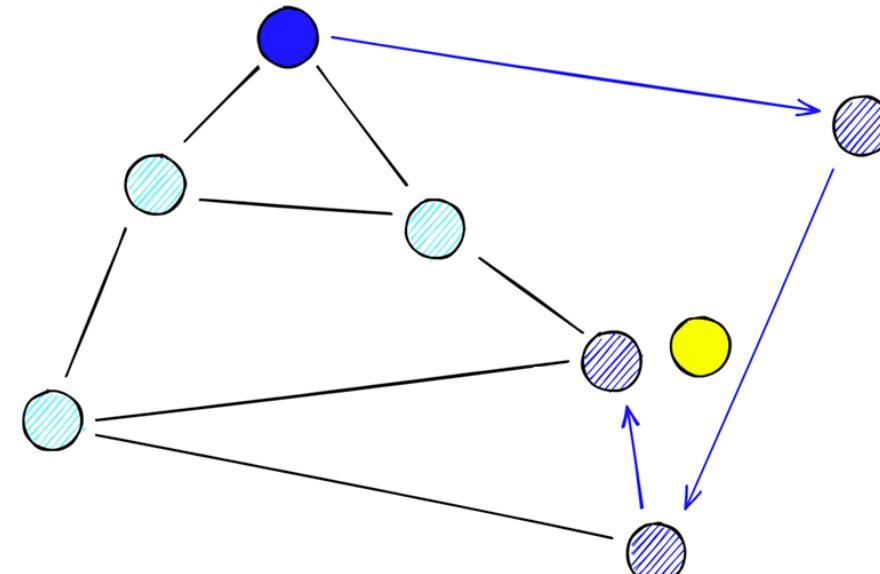
- Then, using an embedding model, you'll embed each chunk and load it into the vector database. The embedding model is a special type of LLM that converts the data chunk into a vector embedding, a numerical representation of the data's meaning. This allows computers to search for similar items based on the vector representation of the stored data.

RAG - Vector DB Indexing Strategies

Flat indexing

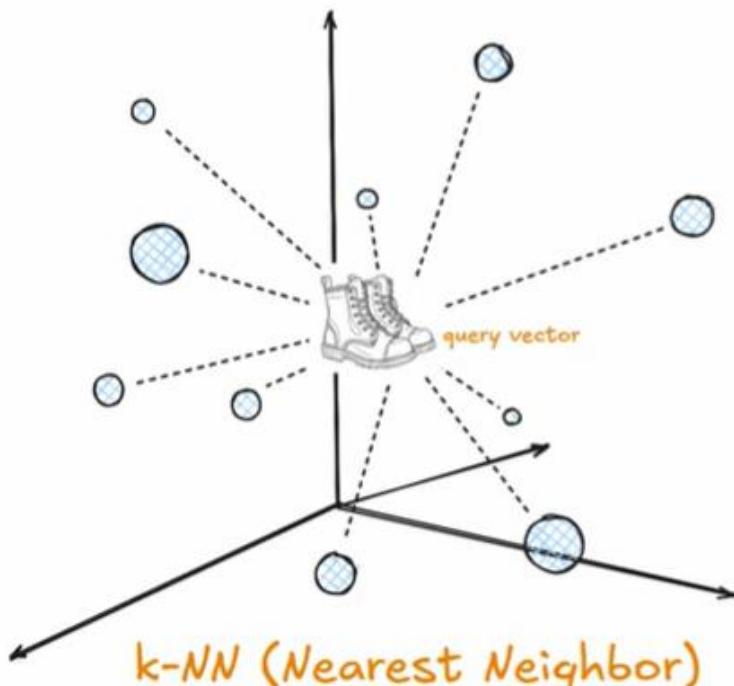


Approximating



RAG - Indexing Strategies – Flat Indexing

Flat Indexing



query vector: Product for which we want to find the similar items

- ▷ Simplest form of indexing
- ▷ Finds the nearest neighbors between the query embedding and every other embedding.
- ▷ Accurate but very slow for large datasets (billions or trillions)

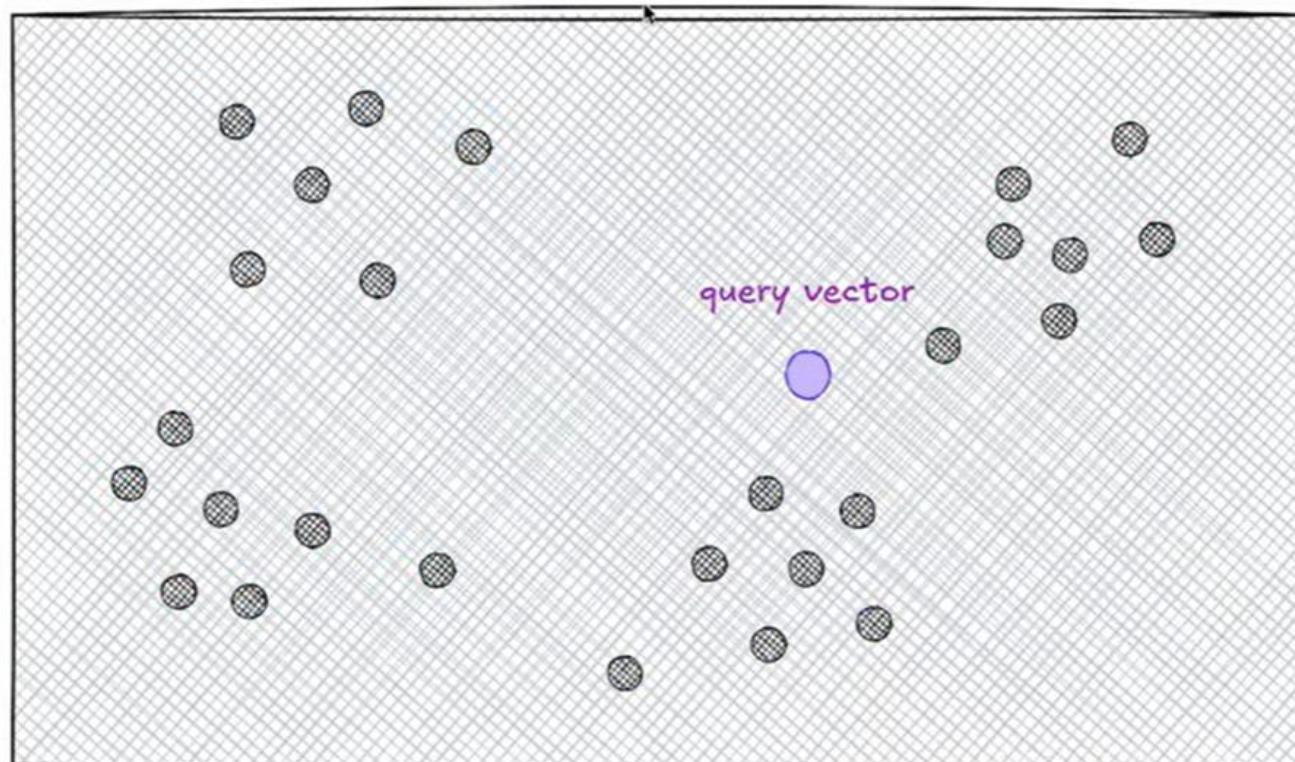
RAG - Indexing Strategies - ANN

ANN !!!

Approximate Nearest Neighbor

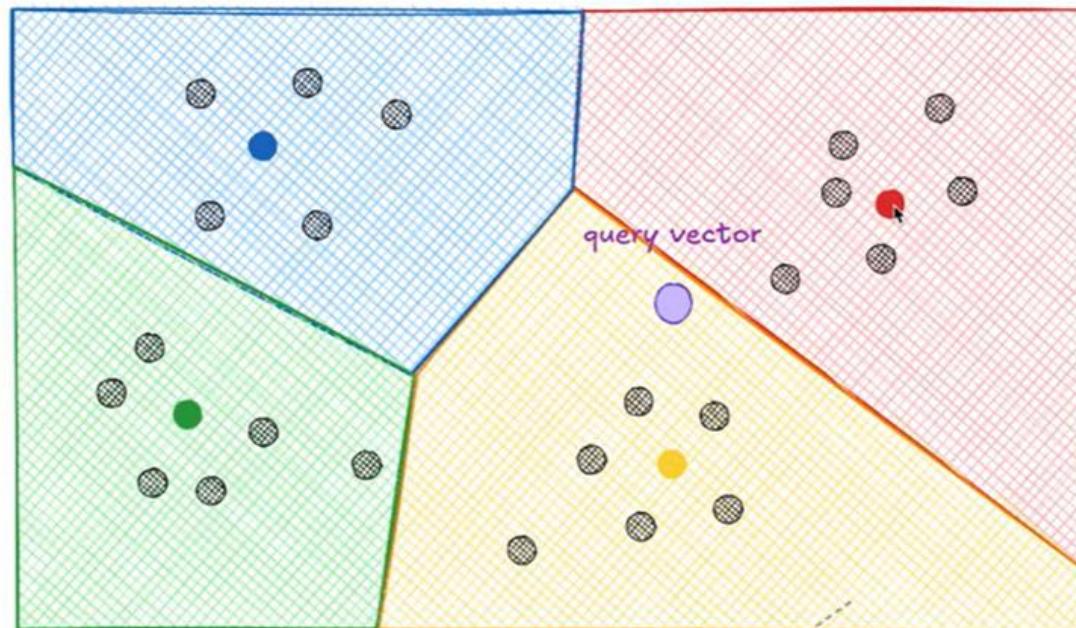
RAG - Indexing Strategies - IVF

IVF (Inverted File Indexing)



RAG - Indexing Strategies - IVF

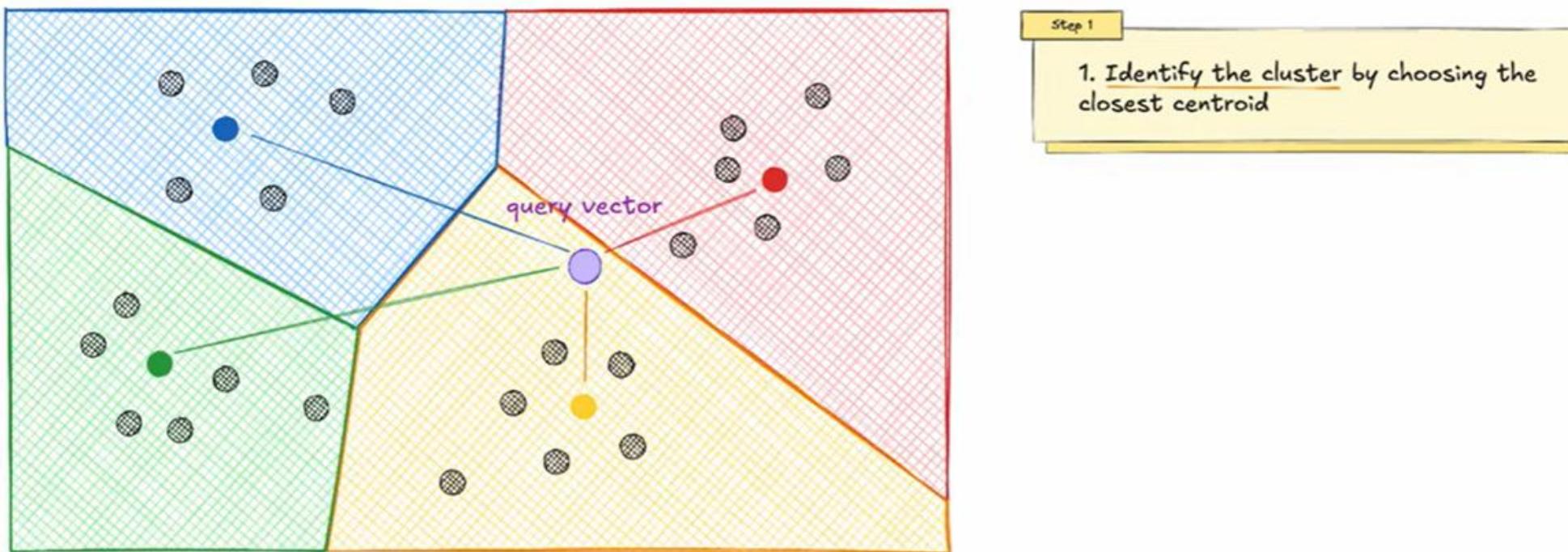
IVF (Inverted File Indexing)



voronoi cells (clusters)

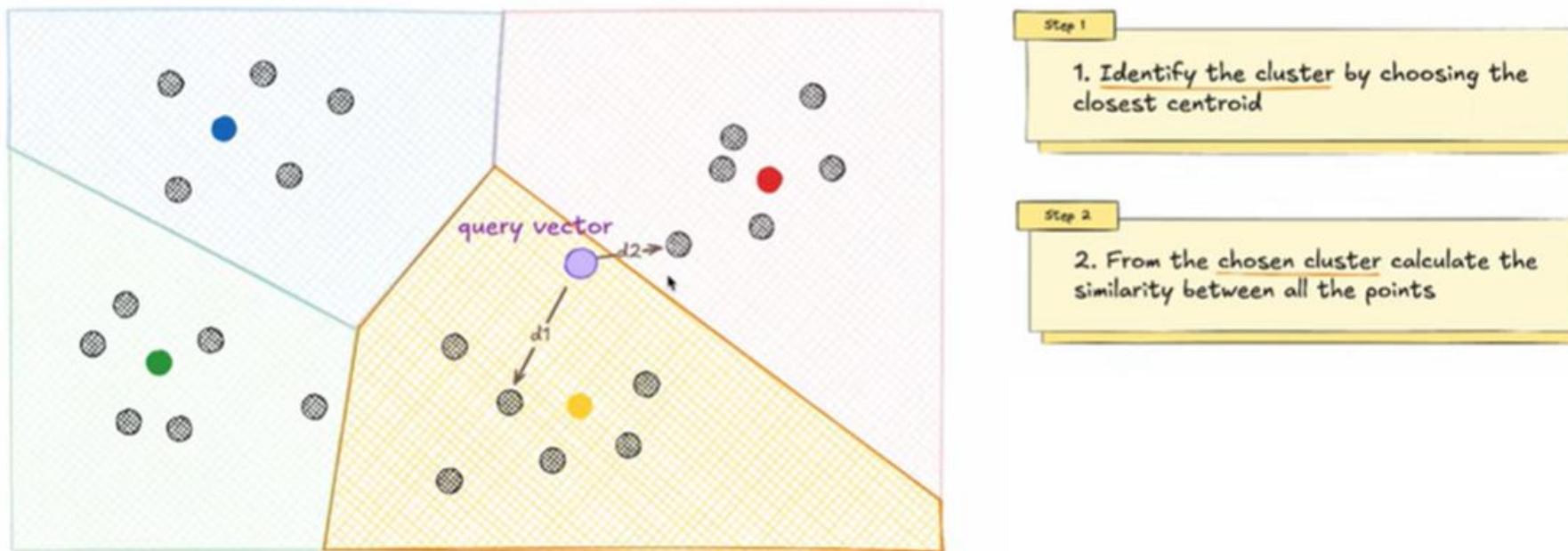
RAG - Indexing Strategies - IVF

IVF (Inverted File Indexing)



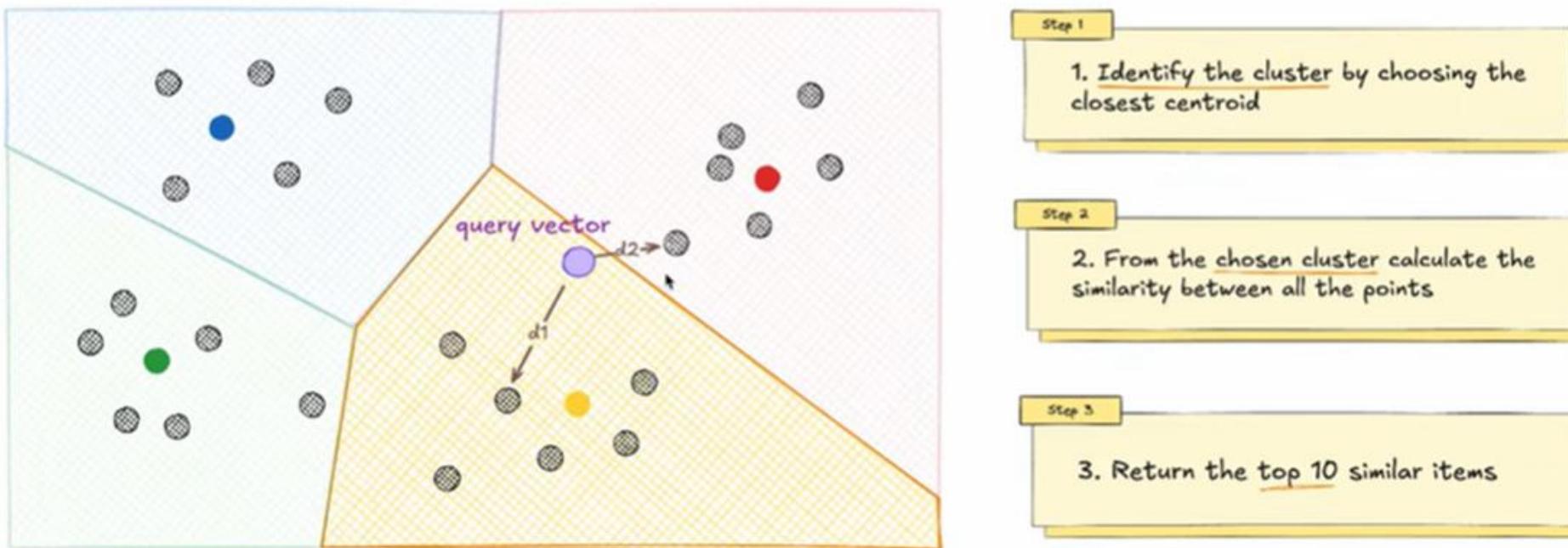
RAG - Indexing Strategies - IVF

IVF (Inverted File Indexing)



RAG - Vector DB Indexing Strategies

IVF (Inverted File Indexing)



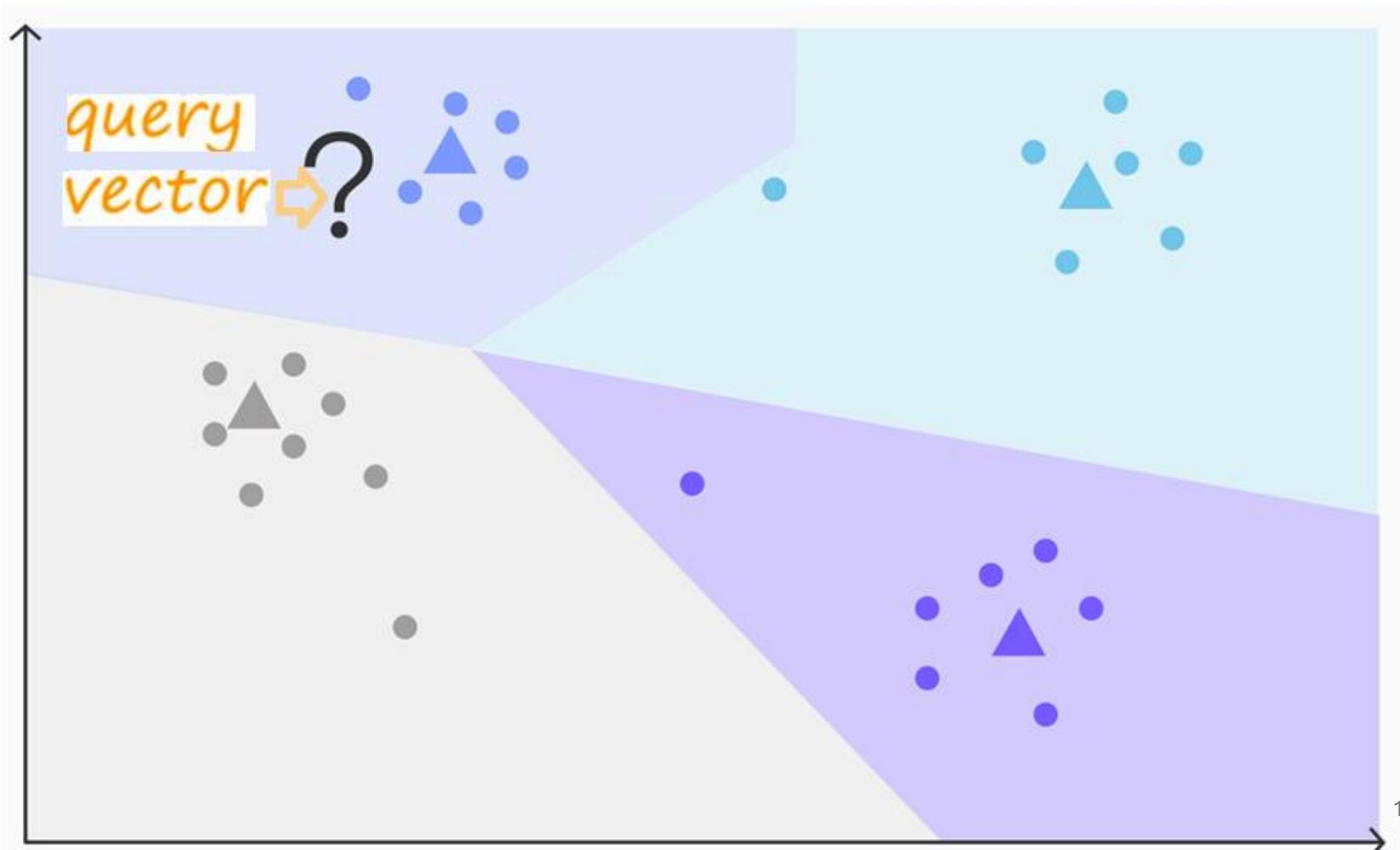
RAG - Indexing Strategies - IVF – The nList Parameter

IVF (Inverted File Index)

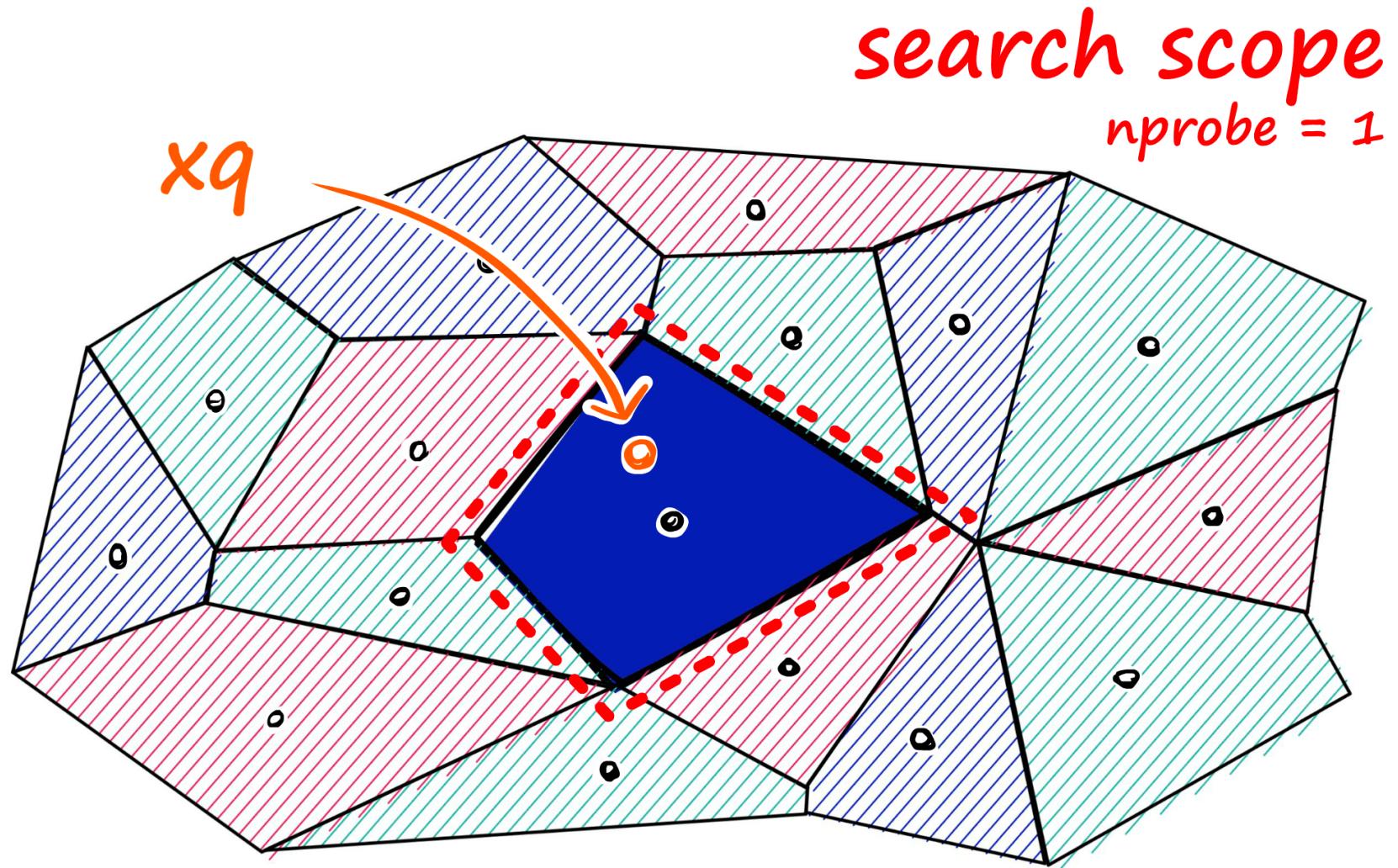
It partitions your vector space into clusters, searching only the relevant clusters. Great for large datasets, with a little less accuracy.

- ❑ nList parameter is the number of clusters
- ❑ Best practice: square root of dataset size
- ❑ Ideal for billion-scale datasets

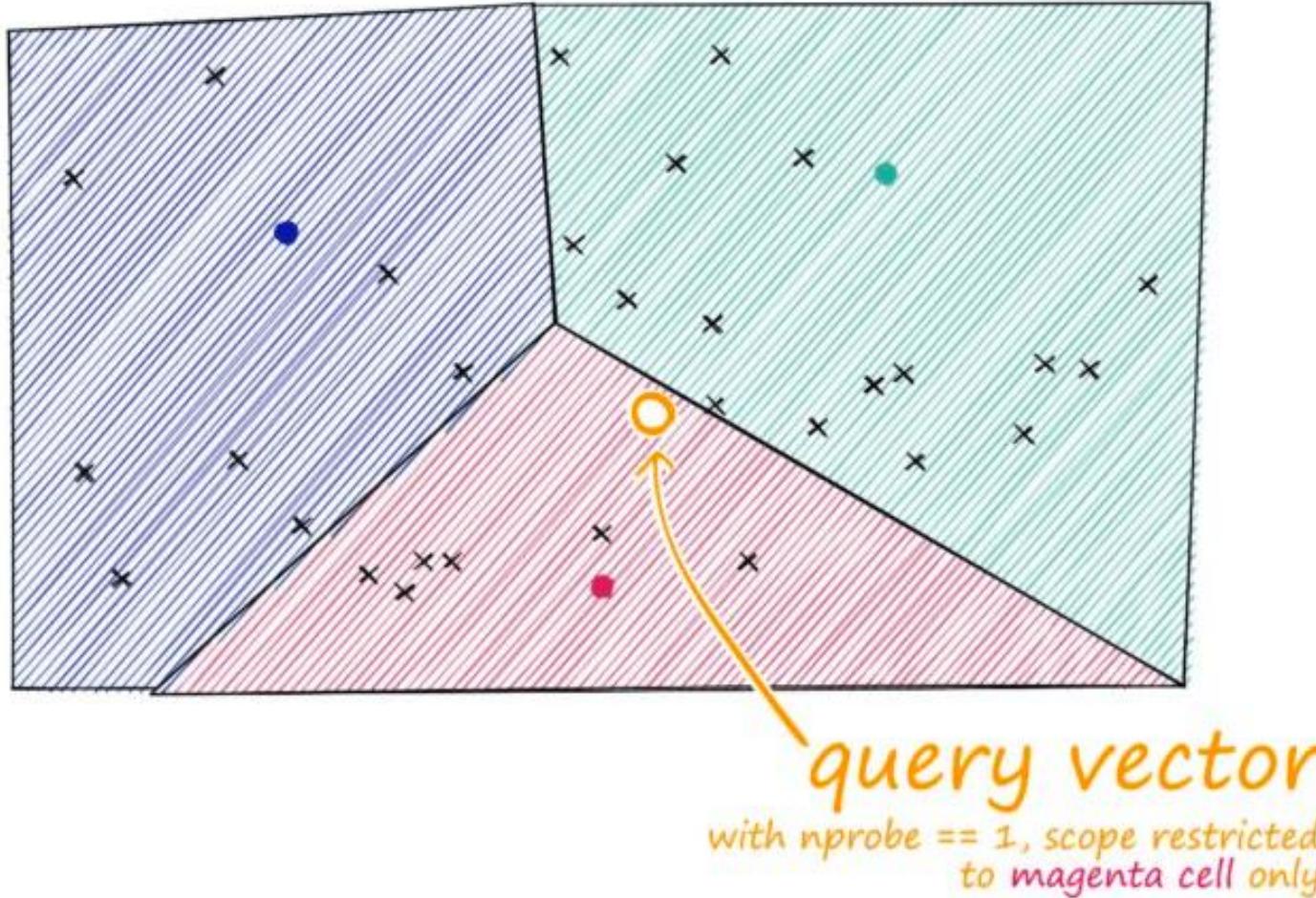
RAG - Indexing Strategies - IVF



RAG - Indexing Strategies - IVF – search in closest cluster

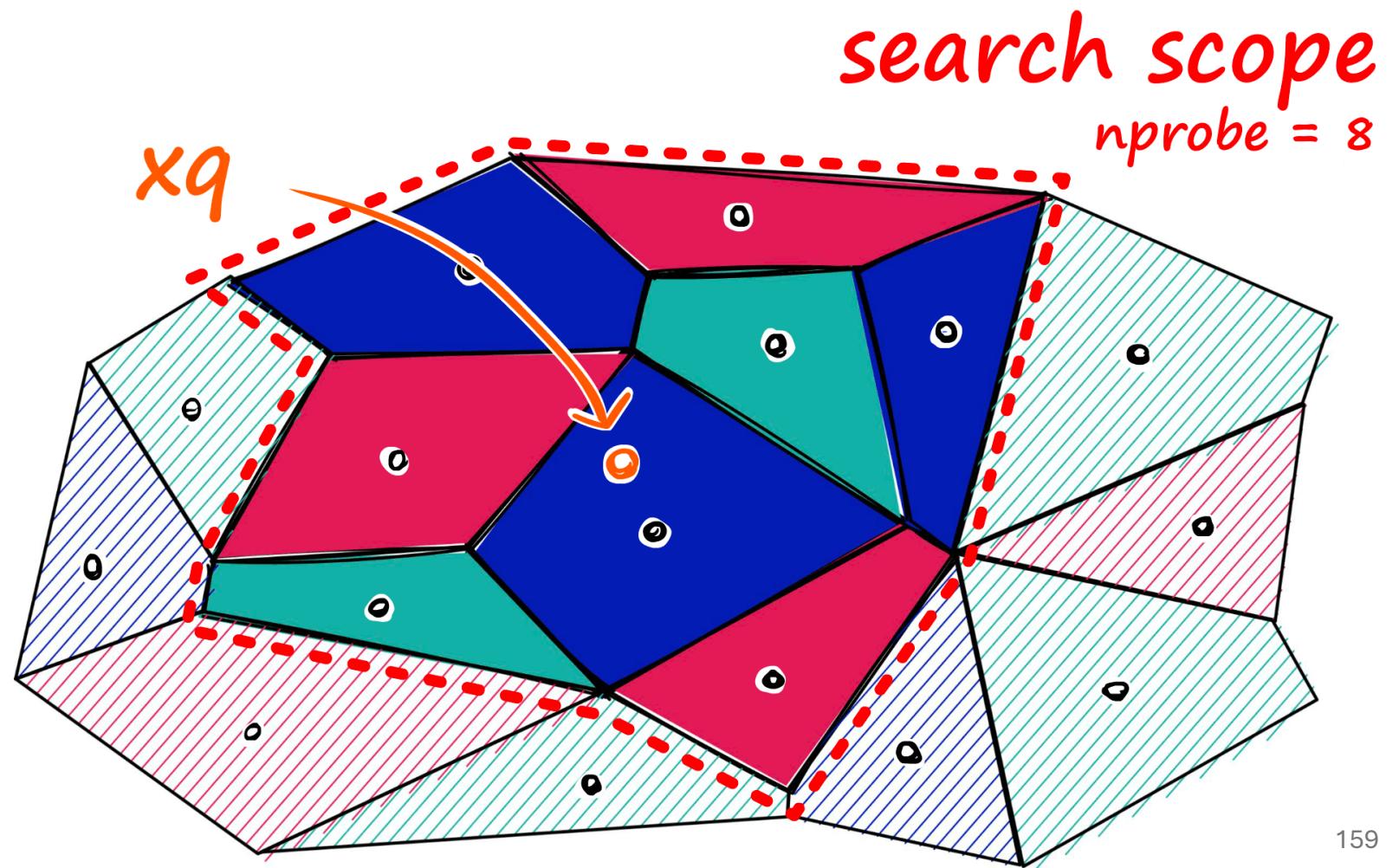


RAG - Indexing Strategies - IVF – search in closest cluster – insufficient



query vector
with $nprobe == 1$, scope restricted
to magenta cell only

RAG - Indexing Strategies - IVF – Increase nprobe - Search in neighboring clusters



RAG - Indexing Strategies - IVF – Increase nprobe - Search in neighboring clusters

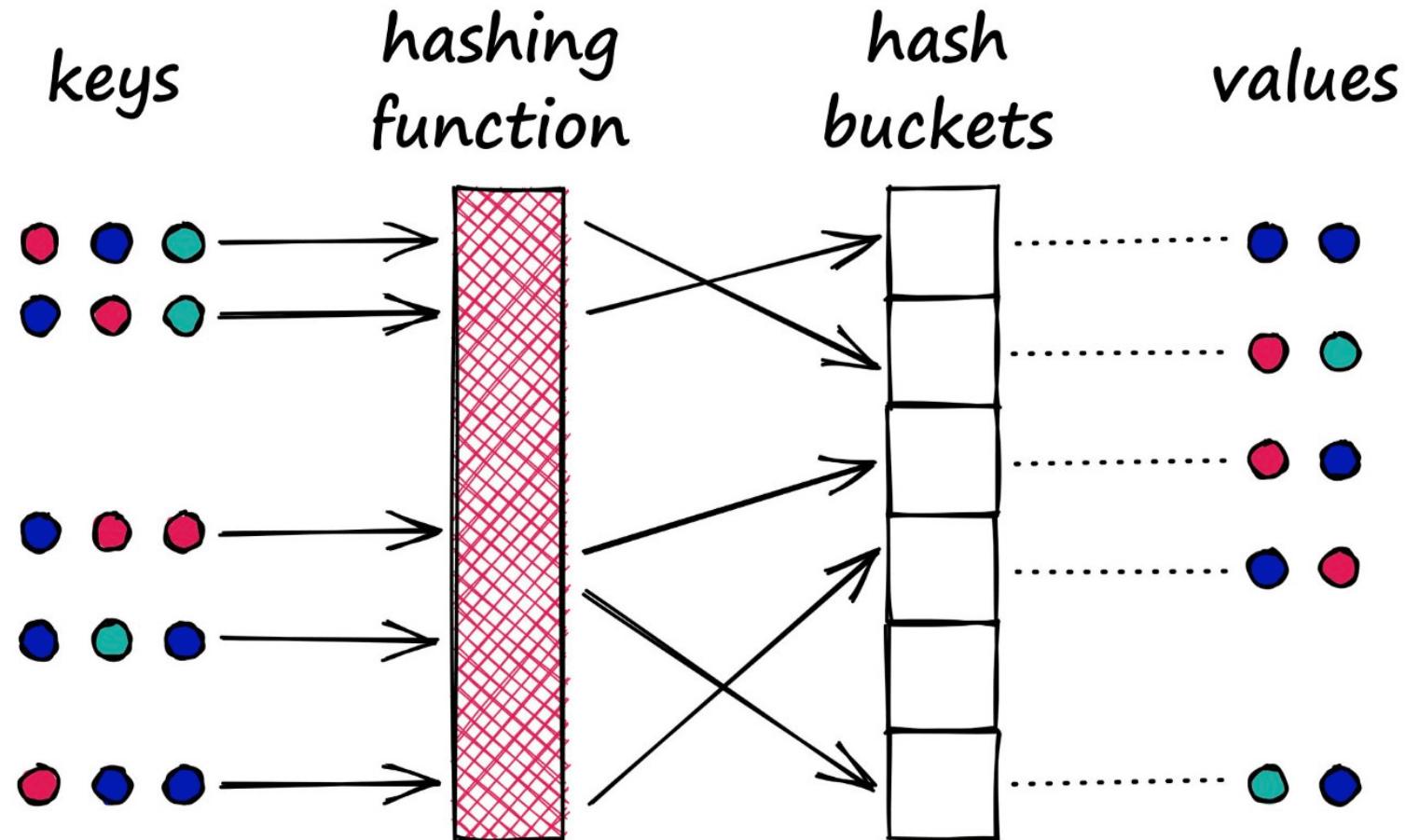
IVF (Inverted File Indexing)

<https://sfotini.github.io/blog/use-ivfflat-index-on-azure-cosmos-db-for-postgresql-for-similarity-search/>

- Step 1
1. Identify the cluster by choosing the closest centroid
- Step 2
2. From the chosen cluster calculate the similarity between all the points
- Step 3
3. Return the top 10 similar items

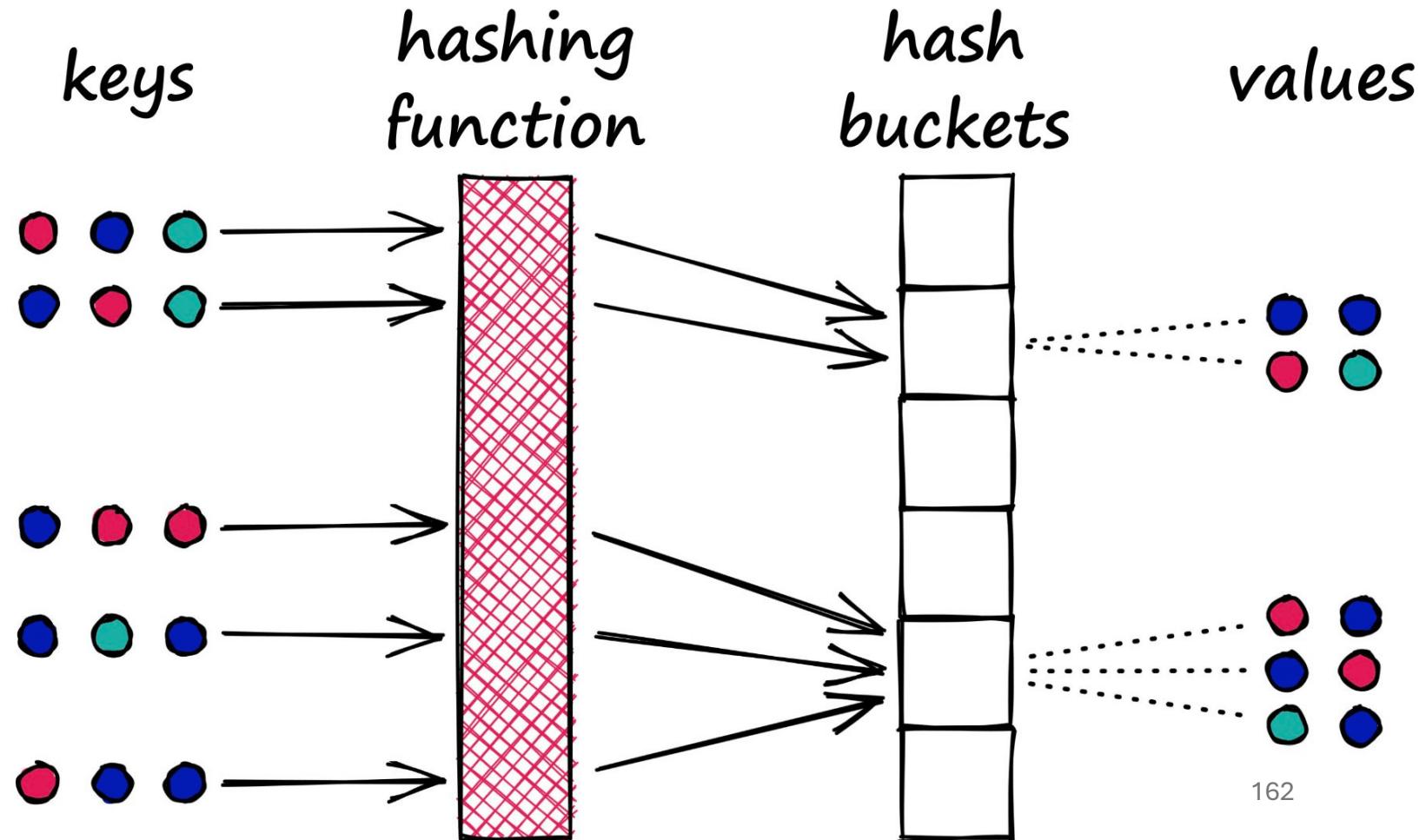
RAG - Indexing Strategies – LSH –

A typical hash function



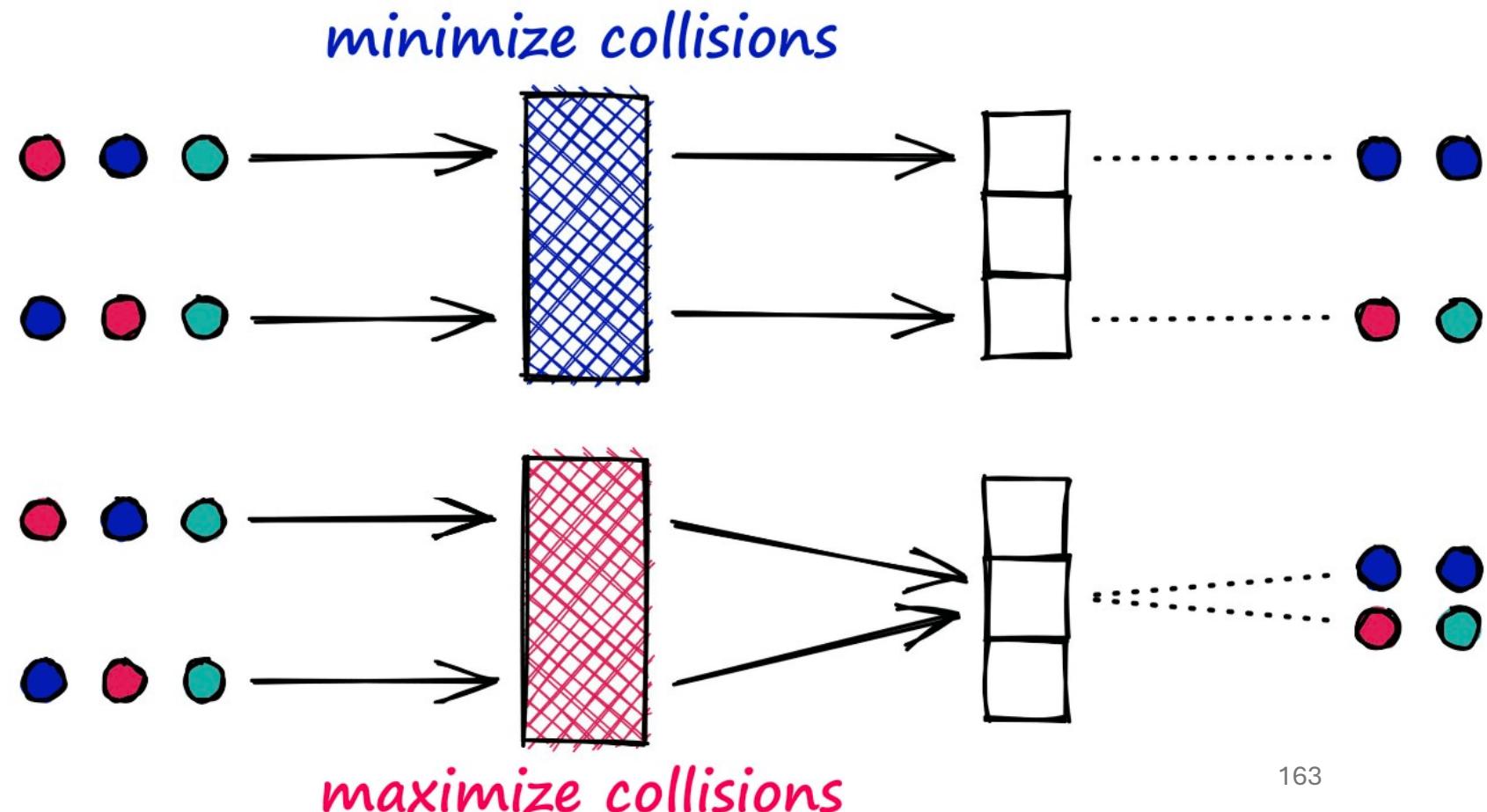
RAG - Indexing Strategies – LSH – Locality sensitive hashing (LSH)

An LSH function aims to place similar values into the same buckets.



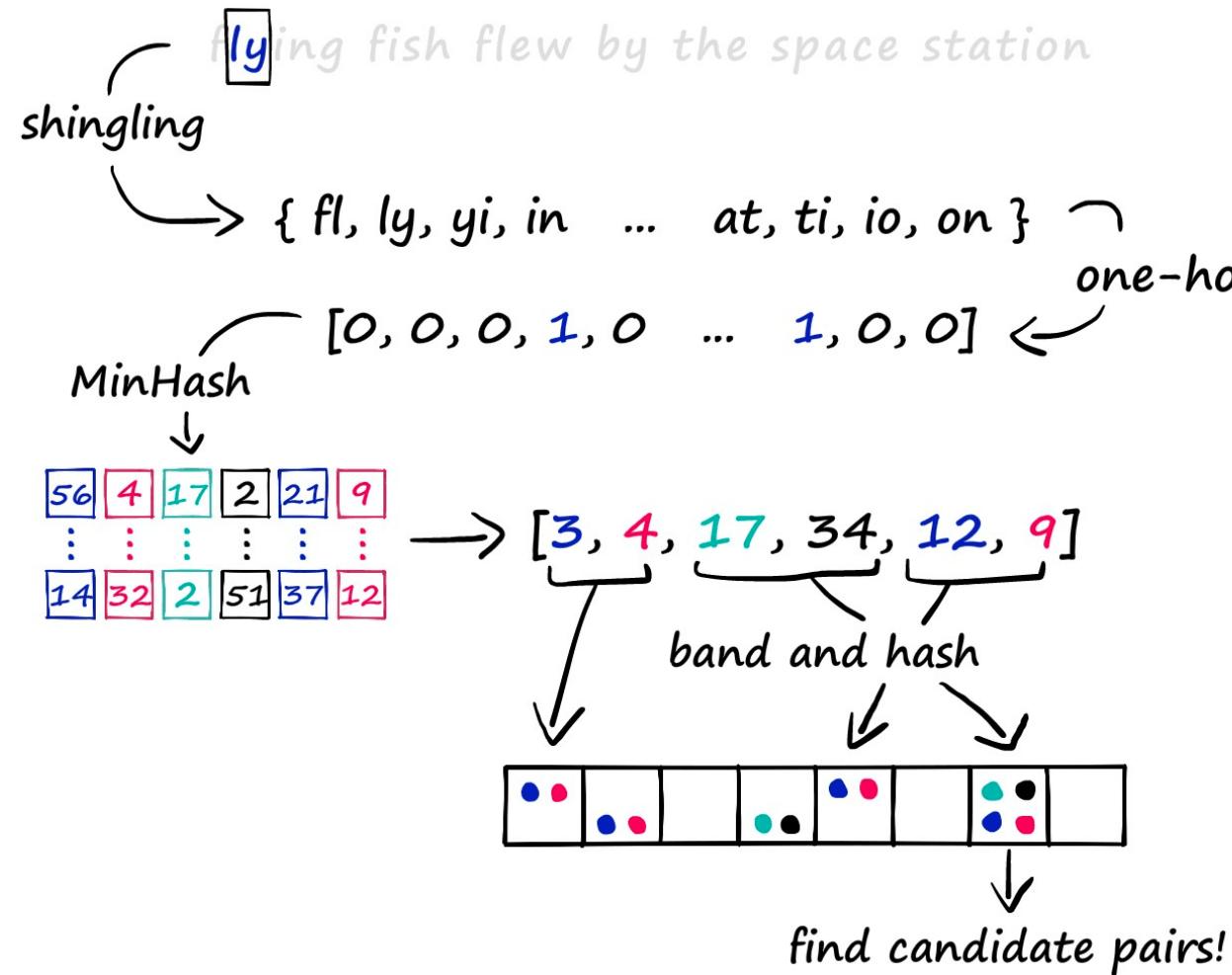
RAG - Indexing Strategies – LSH – Locality sensitive hashing (LSH)

- Typical Hashing - the top (blue) minimizes hashing collisions.
- LSH - The bottom (magenta) maximizes hashing collisions
 - aims to maximize collisions between similar items.



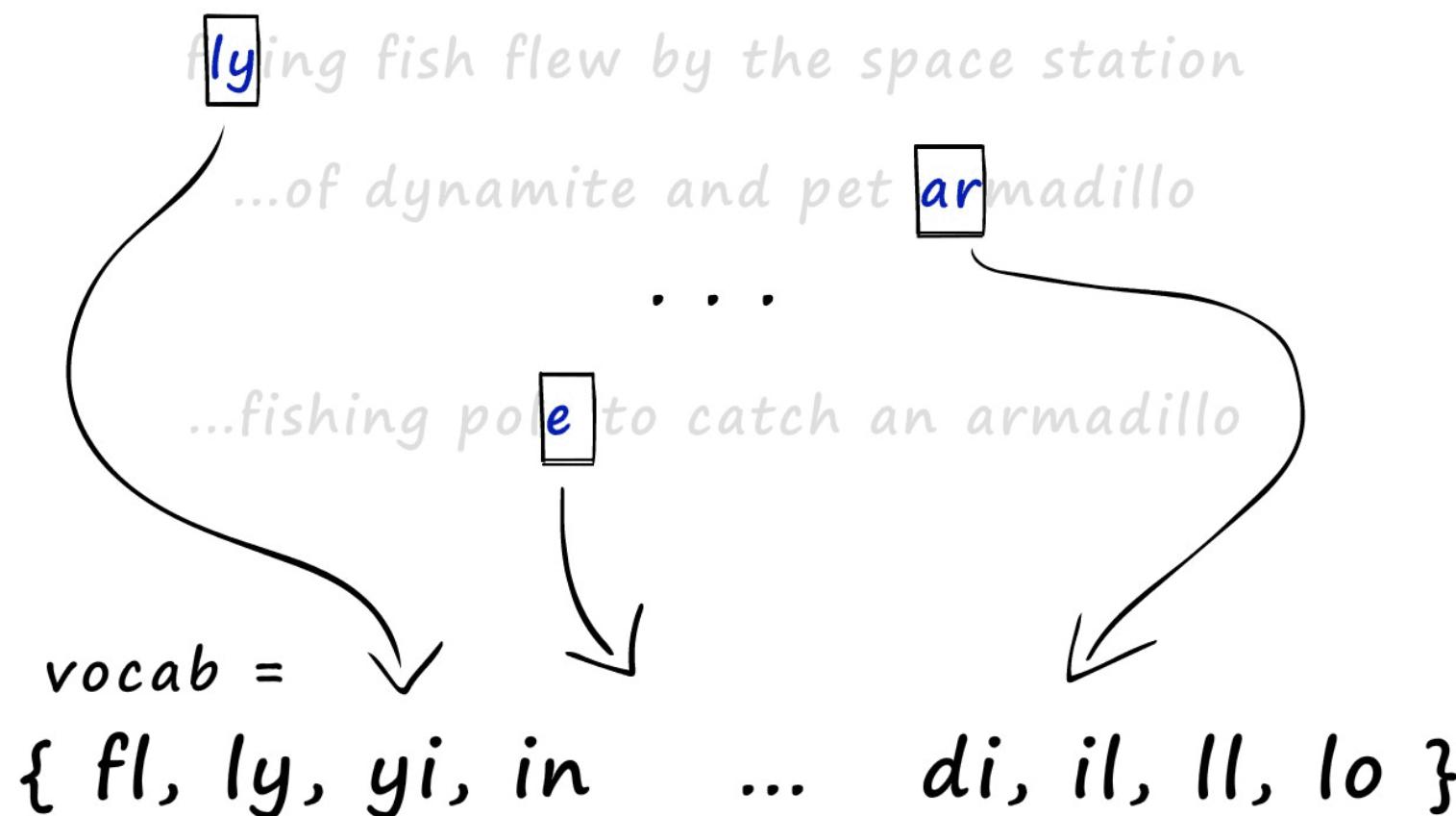
RAG - Indexing Strategies – LSH – LHS Process: Shingling, MinHashing, and LSH

A high-level view of the LSH process



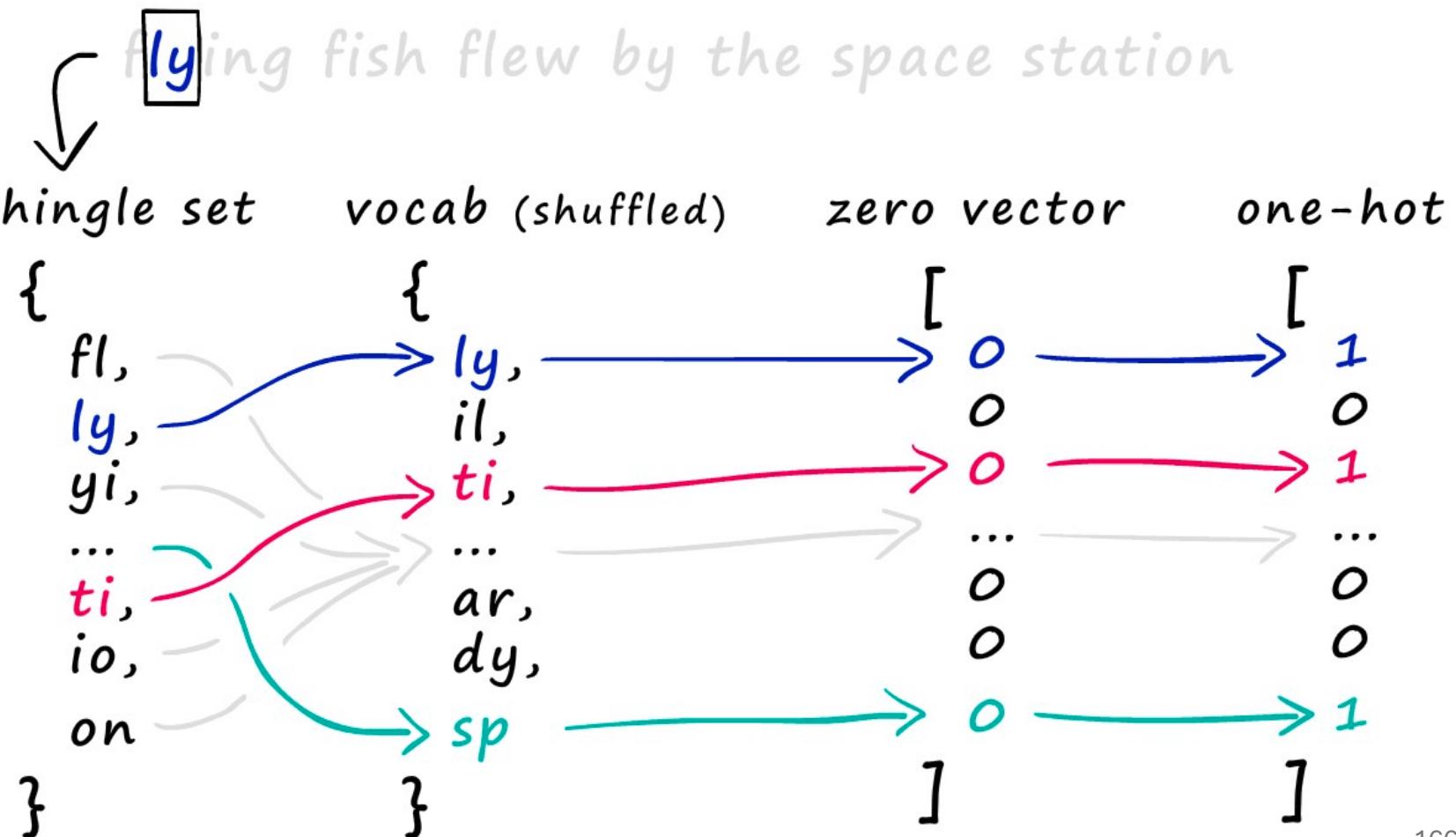
RAG - Indexing Strategies – LSH – LHS Process: (k-)Shingling

- **k-Shingling** -
the process of
converting a
string of text
into a set of
'shingles'.
- Similar to a
moving window
of length k.



RAG - Indexing Strategies – LSH – LHS Process: (k-)Shingling

- Create one hot encodings
(determined by shingle set vocab (shuffled) zero vector one-hot vocabulary)



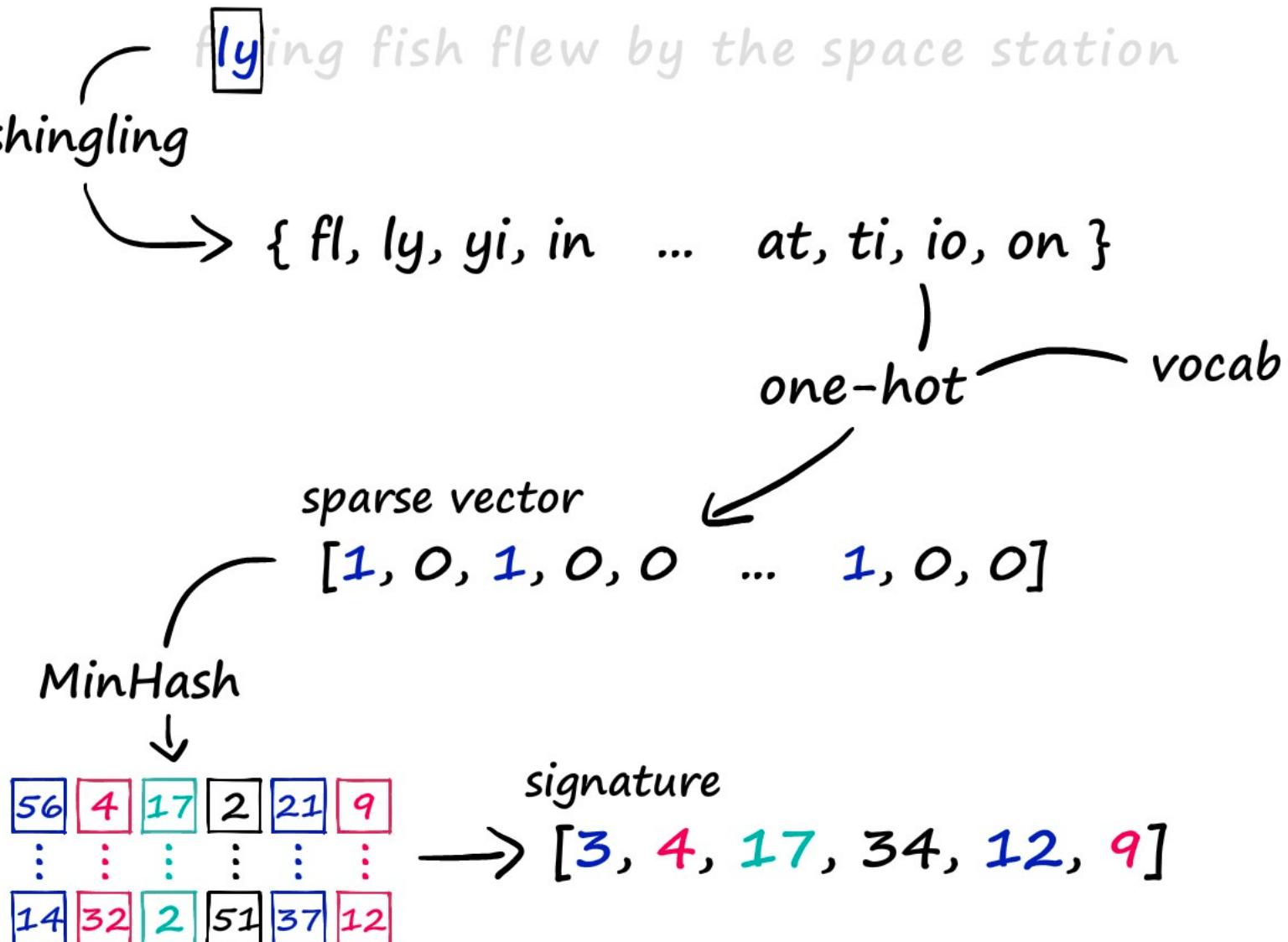
RAG - Indexing Strategies – LSH – LHS Process: Minhashing

- Generate (randomly) one Minhash function for every position in our signature
- For dense vectors - if we wanted to create a dense vector/signature of 20 numbers - we would use 20 minhash functions
- MinHash functions are simply a randomized order of numbers - and we count from 1 to the final number (which is $\text{len}(\text{vocab})$).

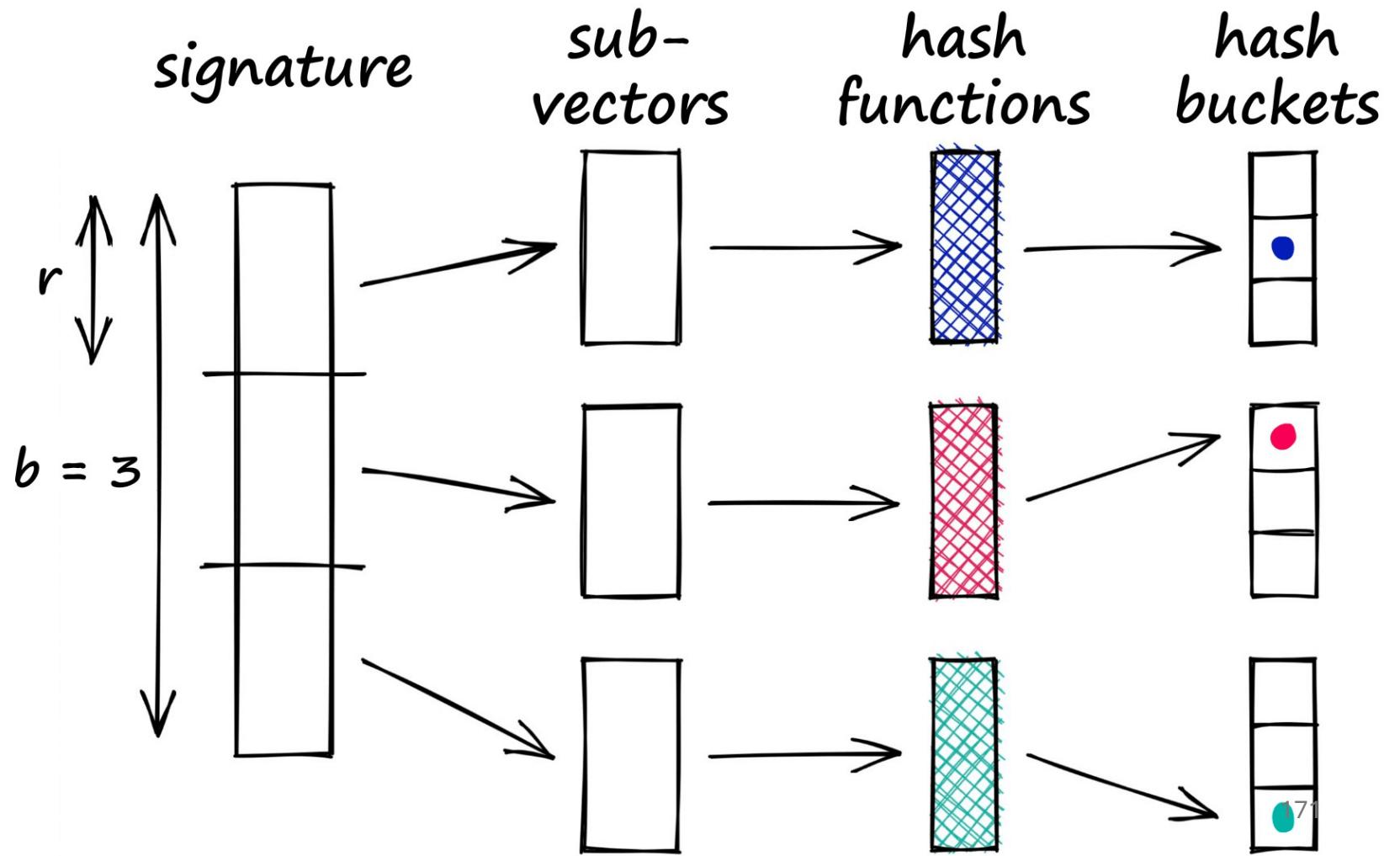
RAG - Indexing Strategies – LSH – LHS Process: Minhashing

- We look at our sparse vector and say, “did this shingle at vocab[1] exist in our set?”.
 - If it did — the sparse vector value will be 1 (0 otherwise).
 - We return the sum of resulting 1s.
-
- We need to produce many (20 or more) of these values. So, we assign a different minhash function to each signature position

- A high-level view of text, build a signature and process it through

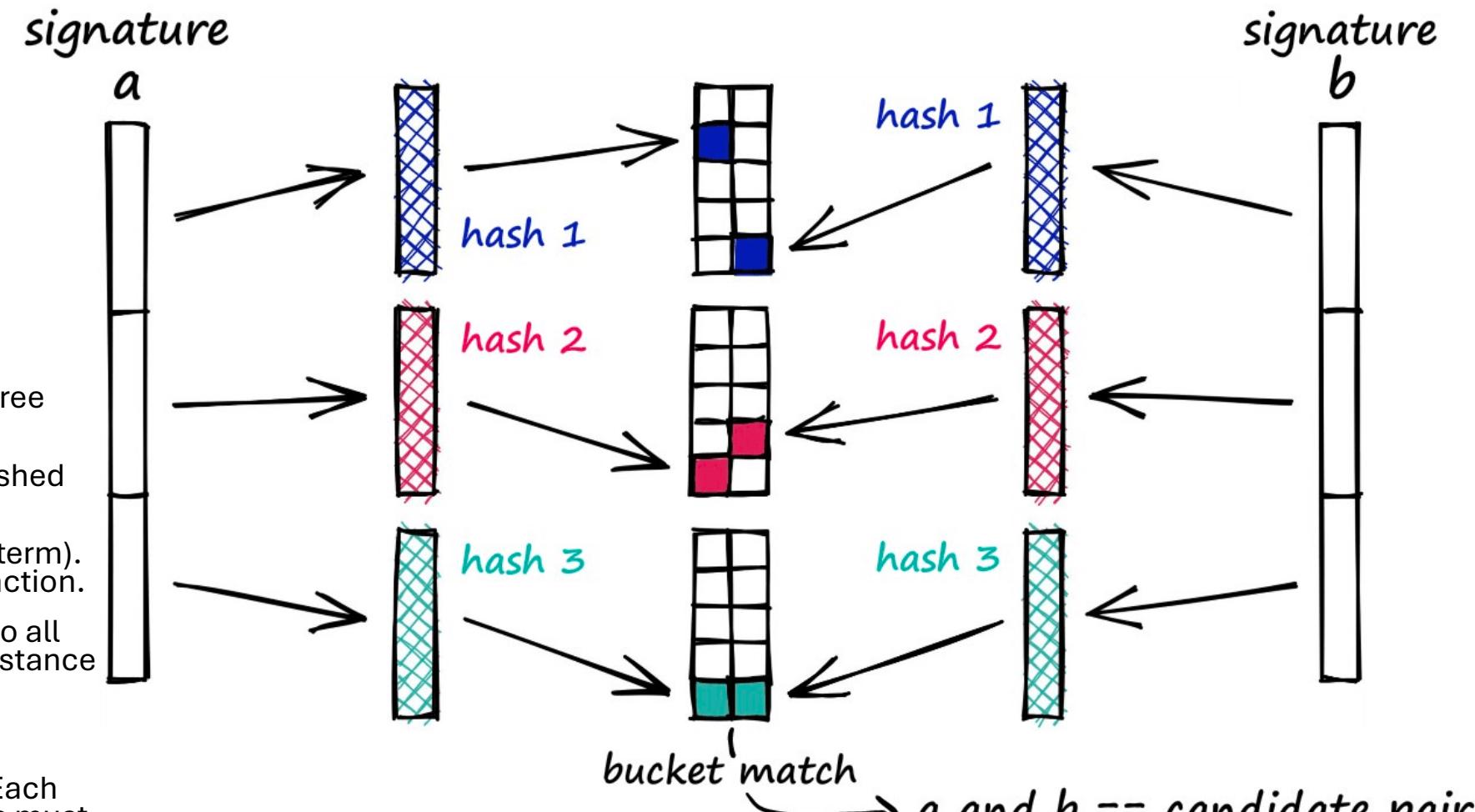


- We split our sig through a hash hash functions



Performing search with LSH consists of three steps:

1. Index all of our vectors into their hashed vectors.
2. Introduce our query vector (search term). It is hashed using the same LSH function.
3. Compare our hashed query vector to all other hash buckets via Hamming distance — identifying the nearest



We split the signatures into subvectors. Each equivalent subvector across all signatures must be processed through the same hash function. However, it is not necessary to use different hash functions for each subvector (we can use just one hash function for them all).

a and b == candidate pair

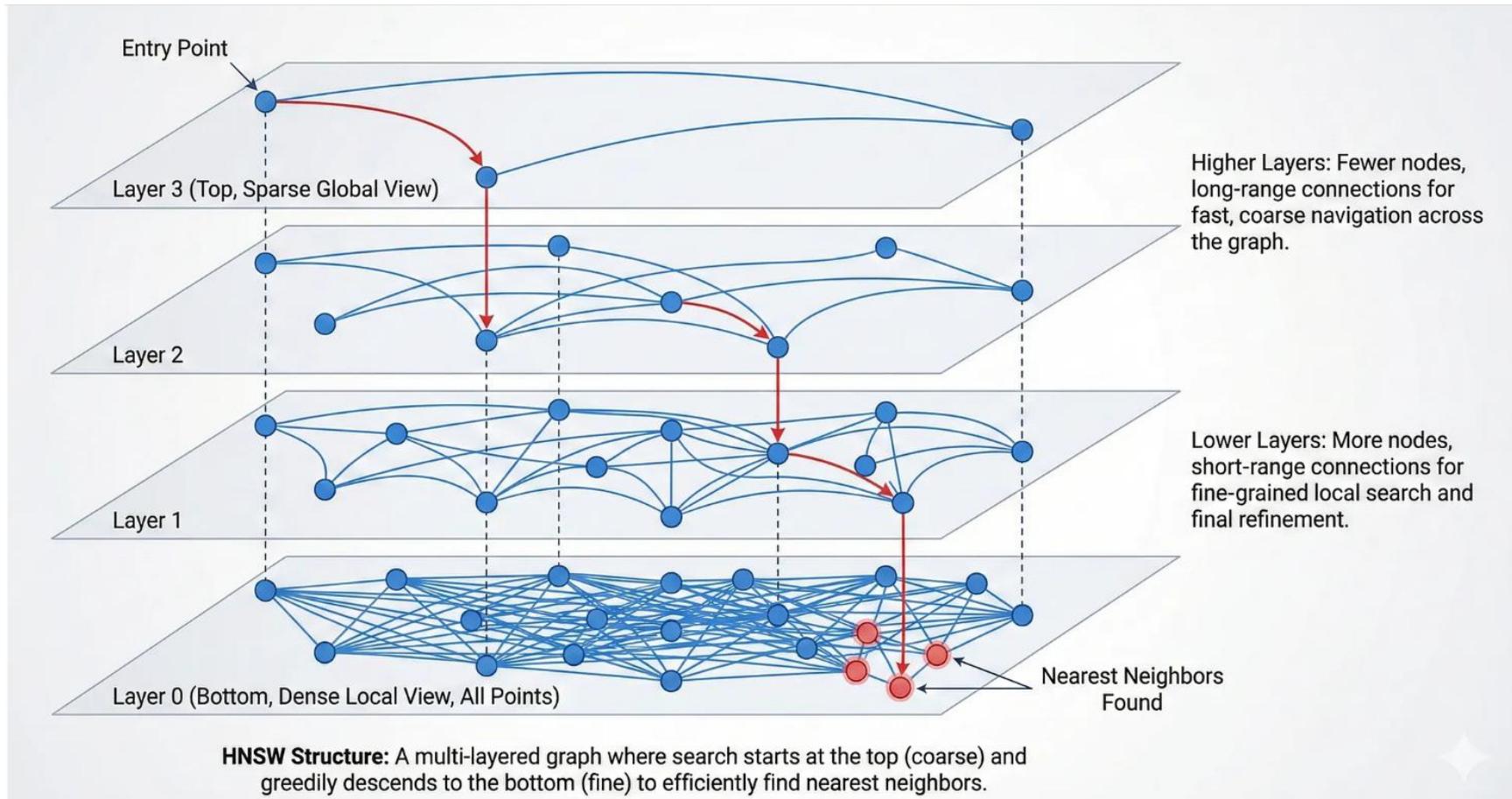
RAG – Ingestion – Chunking

HNSW (Hierarchical Navigable Small World)

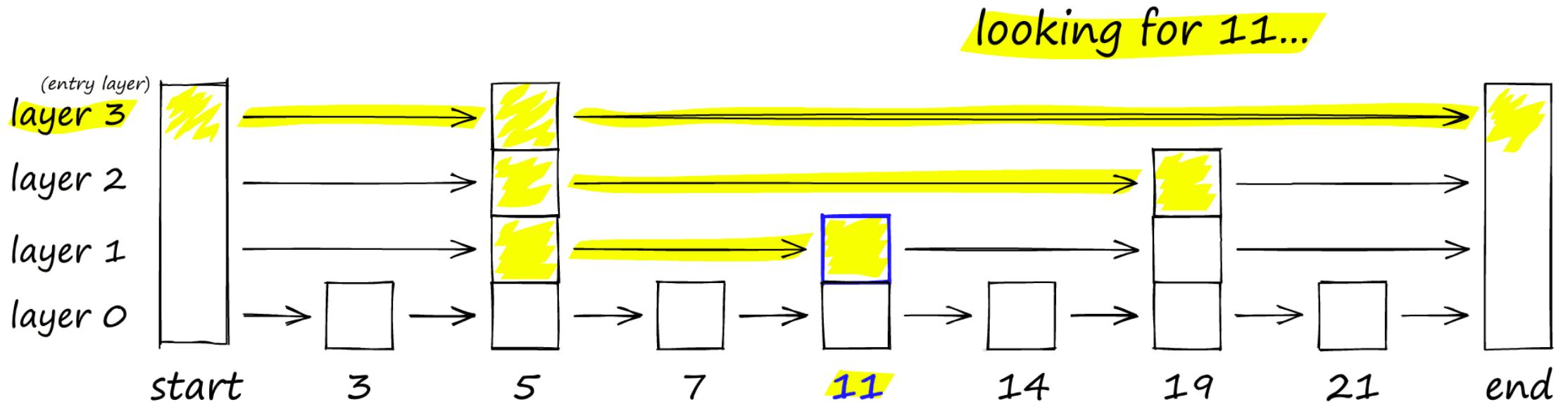
Multi-layered graph structure with great query speed, but comes with significant memory consumption.

- ❑ Fast query with 95%+ recall accuracy
- ❑ Memory: 50-100 bytes per vector
- ❑ Query time: 1-5ms for +1M datasets

Hierarchical Navigable Small World (HNSW)

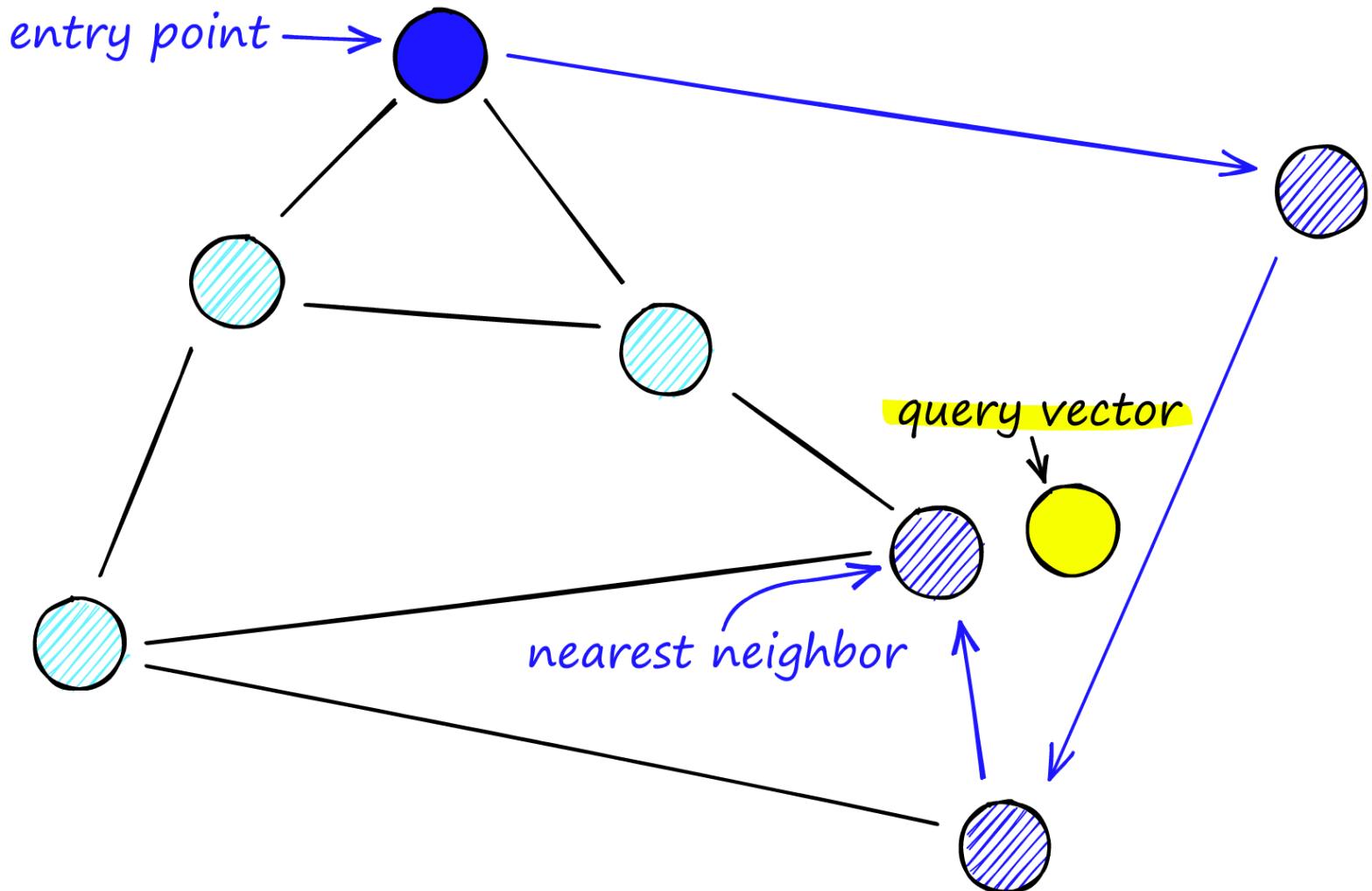


Inspiration: Probabilistic Skip List



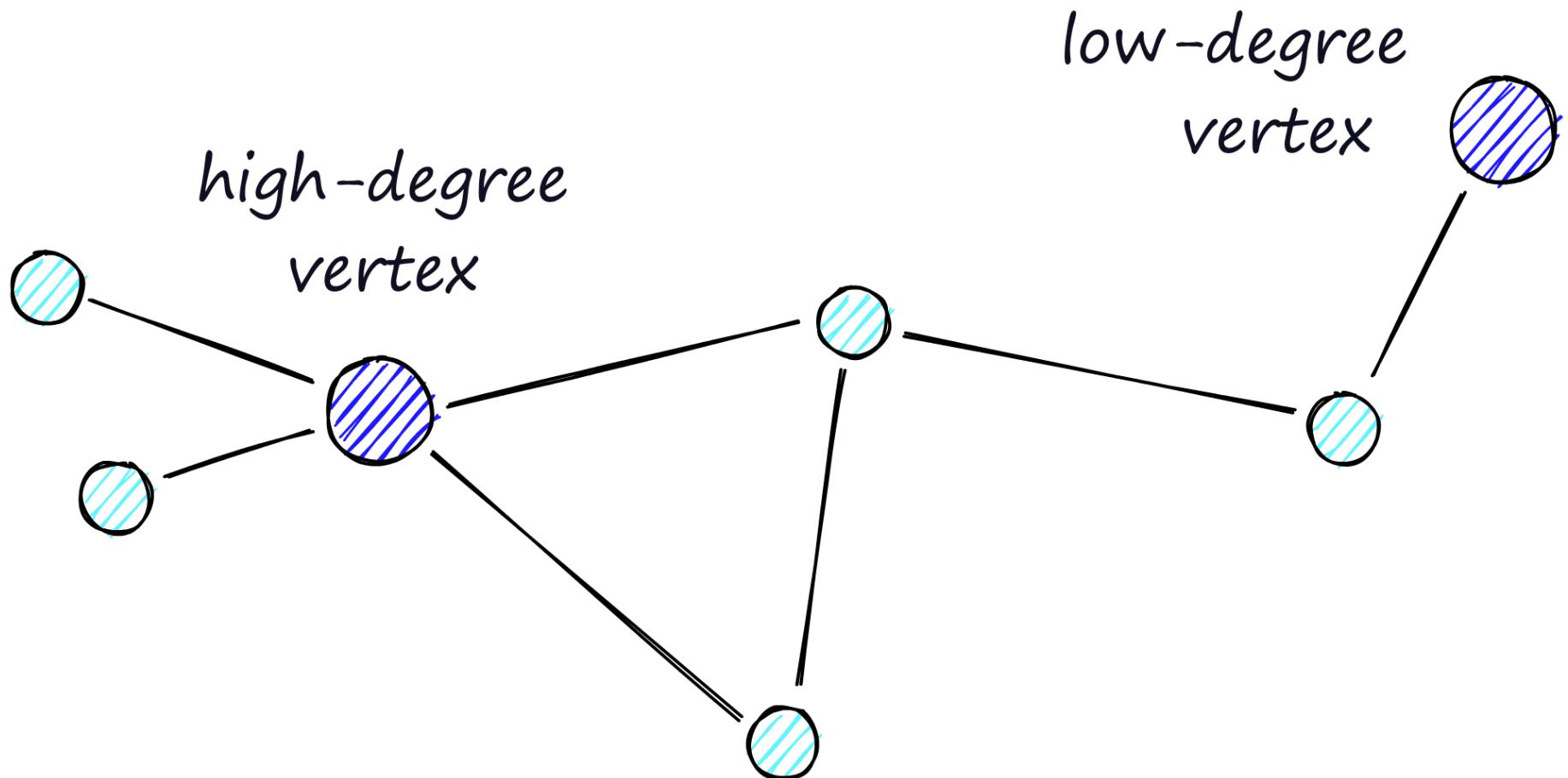
Navigable Small World Graphs

The search process through a NSW graph. Starting at a pre-defined entry point, the algorithm greedily traverses to connected vertices that are nearer to the query vector.

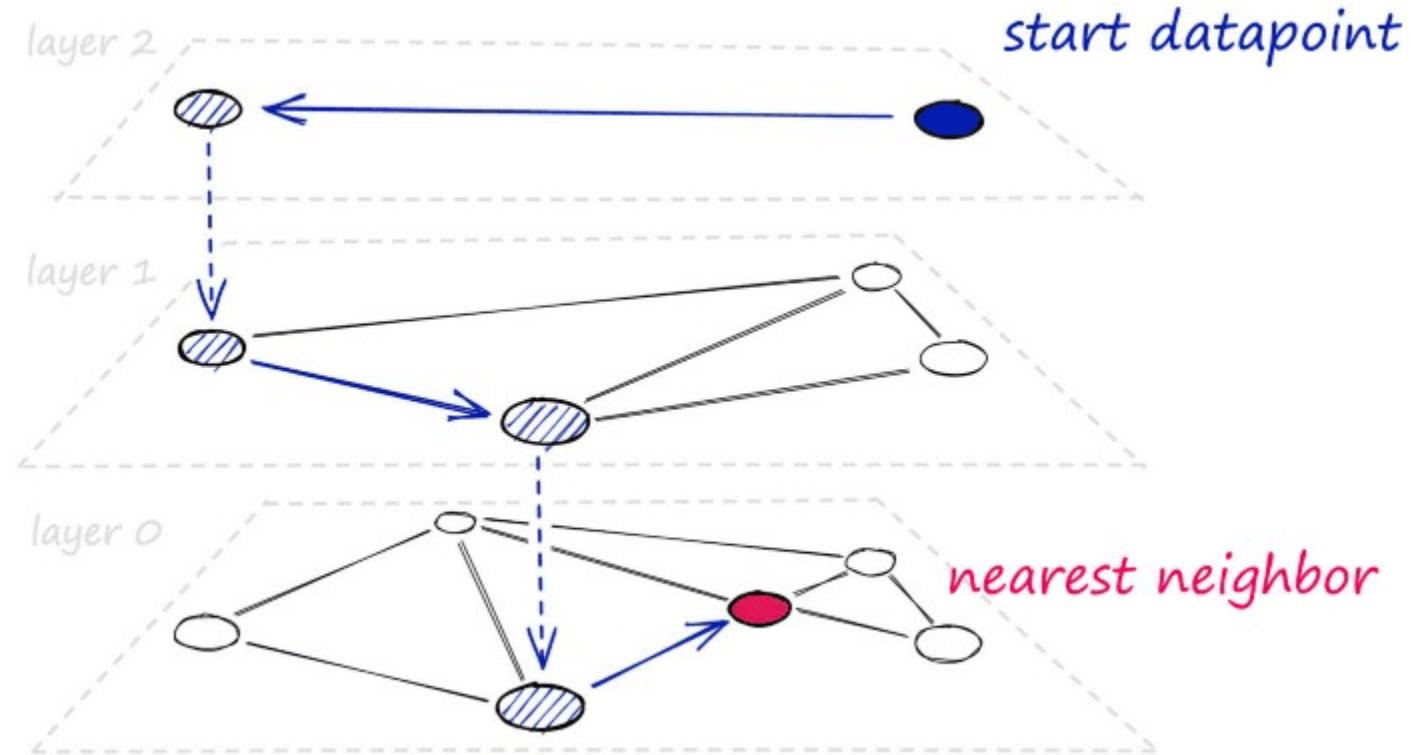


The routing

High-degree vertices have many links, whereas low-degree vertices have very few links.

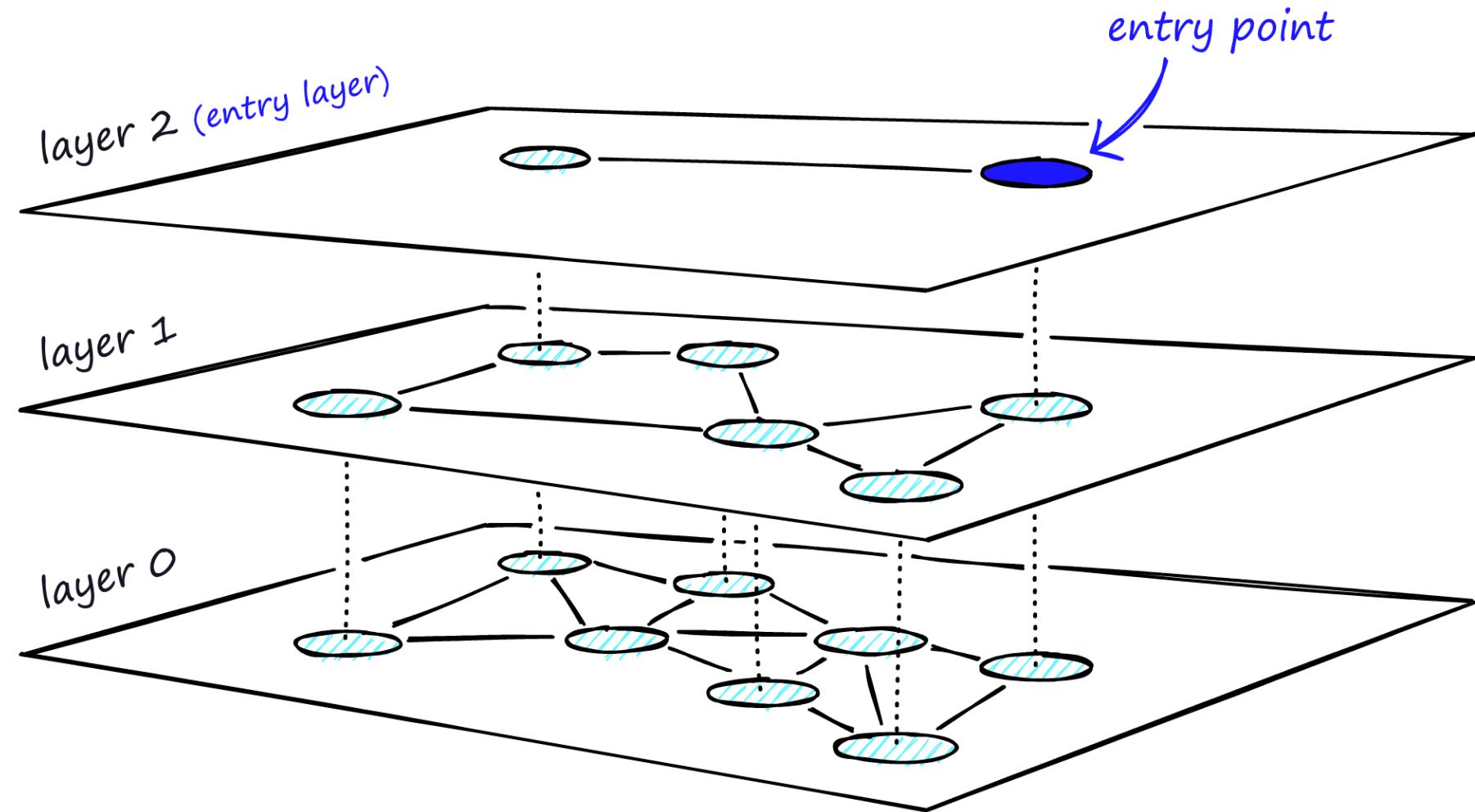


HNSW - Searching Nearest Neighbors



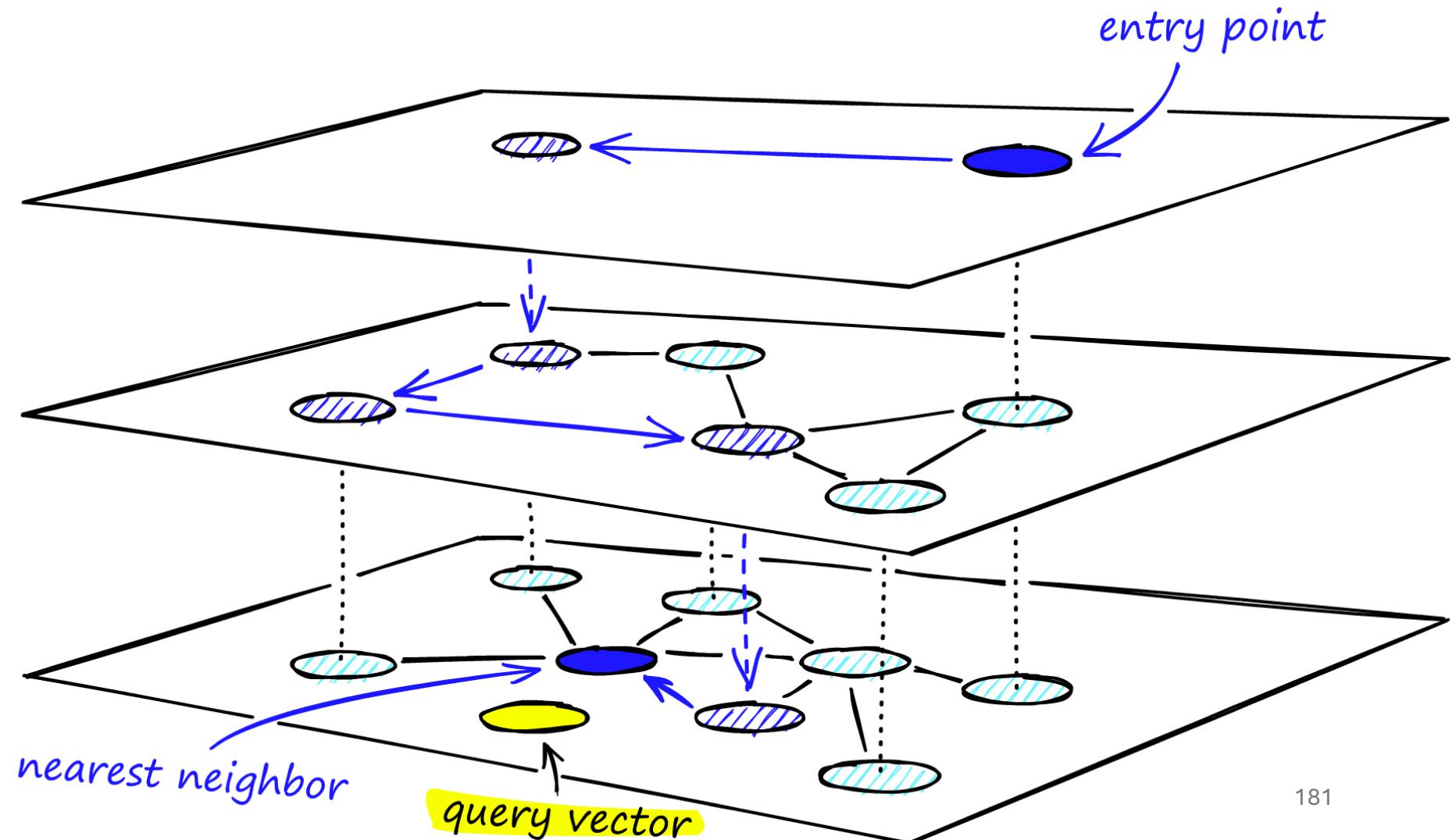
Creating HNSW

Layered graph of HNSW, the top layer is our entry point and contains only the longest links, as we move down the layers, the link lengths become shorter and more numerous.



Creating HNSW

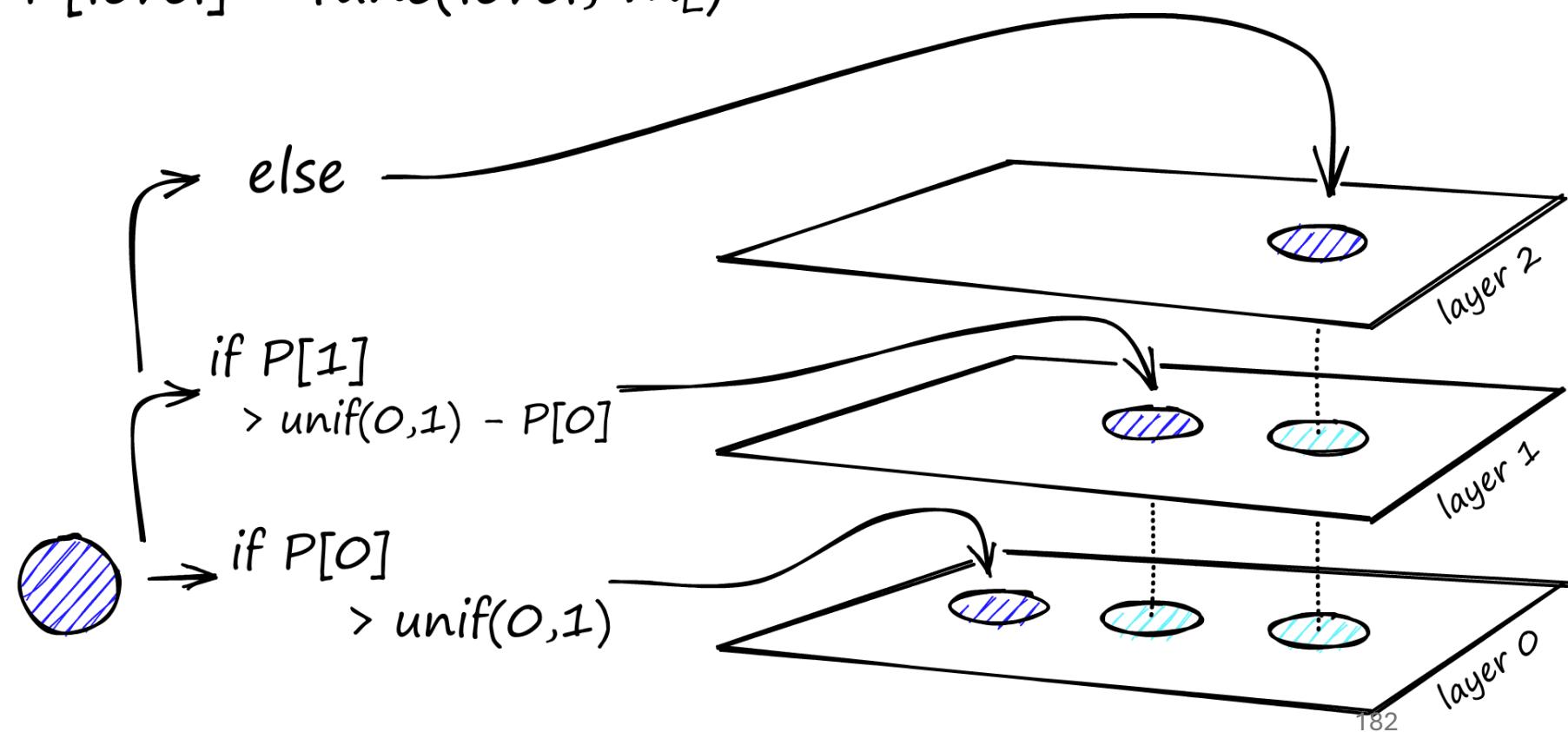
The search process through the multi-layer structure of an HNSW graph.



Graph Construction

The probability function is repeated for each layer (other than layer 0). The vector is added to its insertion layer and every layer below it.

$$P[\text{level}] = \text{func}(\text{level}, m_L)$$



Graph Construction

Explanation
of the
number of
links
assigned to
each vertex
and the
effect of M,
 M_{\max} , and
 $M_{\max 0}$.

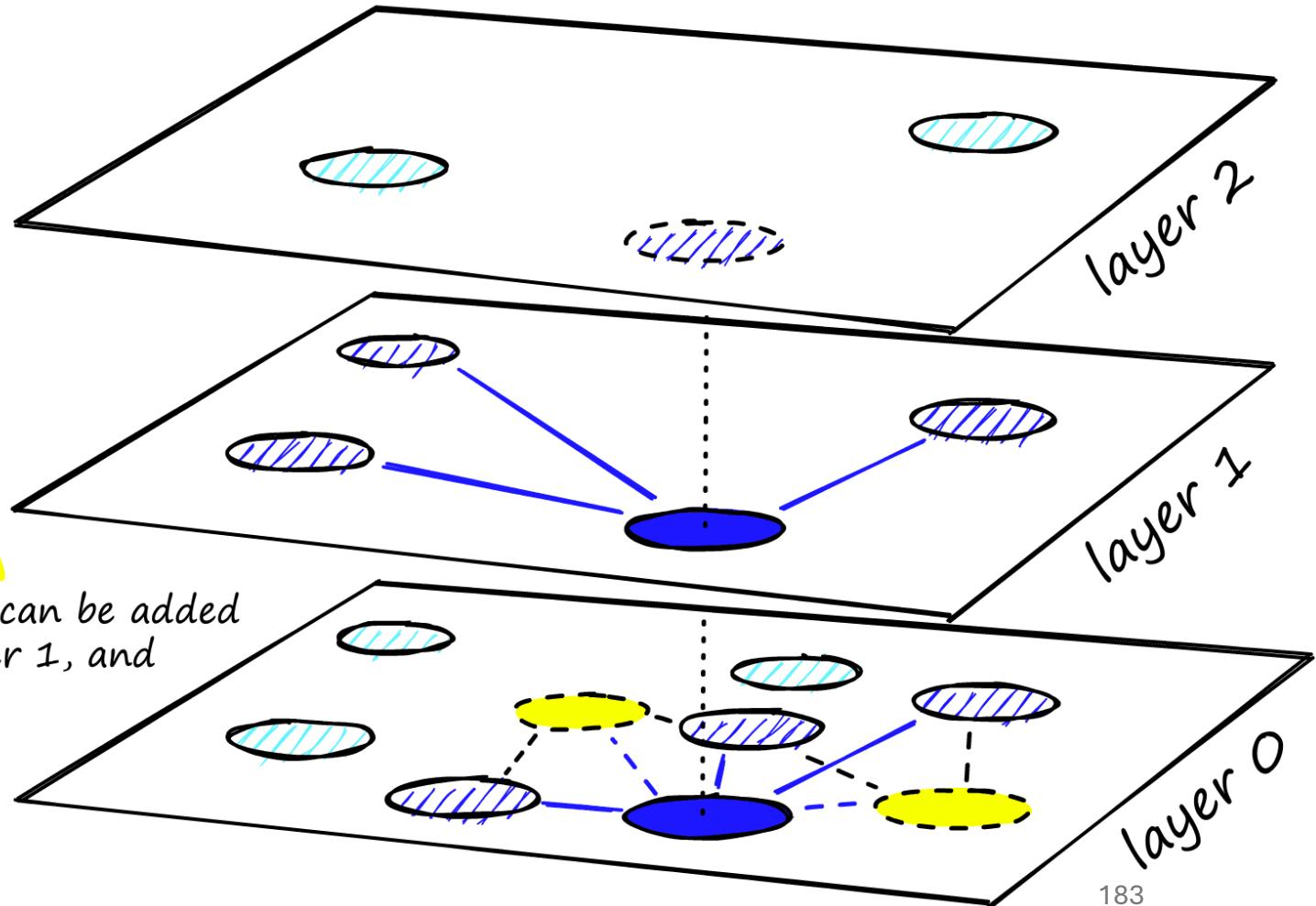
*insert vector
at layer 1*

with $M = 3$
layer 1 and 0
find 3 links

as **more vertices are inserted**, more links can be added
- up to M_{\max} for layer 1, and
 $M_{\max 0}$ for layer 0

$$M_{\max} = 3$$

$$M_{\max 0} = 5$$



RAG – Augmentation

- create an augmented prompt with both the search results and the user's query to send to the LLM. This is where the magic happens.
- An augmented prompt might look like this:

QUESTION:

<the user's question>

CONTEXT:

<the search results to use as context>

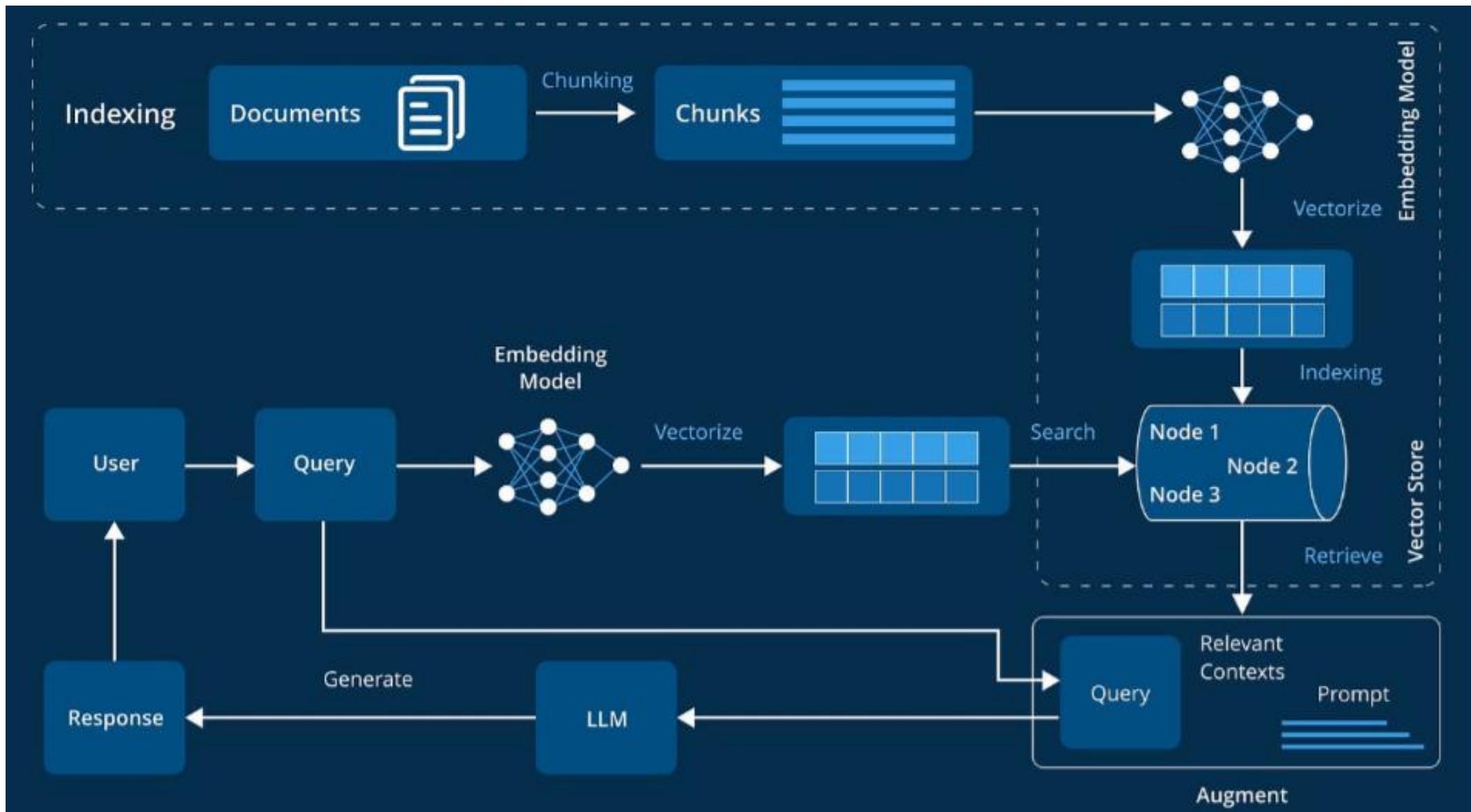
Using the CONTEXT provided, answer the QUESTION. Keep your answer grounded in the facts of the CONTEXT. If the CONTEXT doesn't contain the answer to the QUESTION, say you don't know.

- By sending both the search results and the user's question as context to the LLM, you are encouraging it to use the more accurate and relevant info from the search results during the next generation step.

RAG – Generation

- Using the augmented prompt, the LLM now has access to the most pertinent and grounding facts from your vector database so your application can provide an accurate answer for your user, reducing the likelihood of hallucination.

Typical RAG system



Without Query Expansion

Query: What is a **bridge** for **teeth**?

"What is a **bridge**?" Are you missing a tooth? A **bridge** is one of the dental treatments available to replace missing **teeth**. Bridges are fixed into the mouth by attaching them to neighboring **teeth**. Therefore, unlike dentures (which are removable), these stay in your mouth. The most common type is a conventional **bridge** where the **teeth** on either-side of the gap are crowned, with a "false" tooth stuck to those crowns, in order to fill the gap.

With Query Expansion

Query: What is a bridge for teeth?

[with, bridges, gap]

[dental, tooth, mouth]

"What is a **bridge**?" Are you missing a **tooth**? A **bridge** is one of the dental treatments available to replace missing **teeth**. **Bridges** are fixed into the **mouth** by attaching them to neighboring **teeth**. Therefore, unlike dentures (which are removable), these stay in your **mouth**. The most common type is a conventional **bridge** where the **teeth** on either-side of the **gap** are crowned, **with** a "false" **tooth** stuck to those crowns, in order to fill the **gap**.

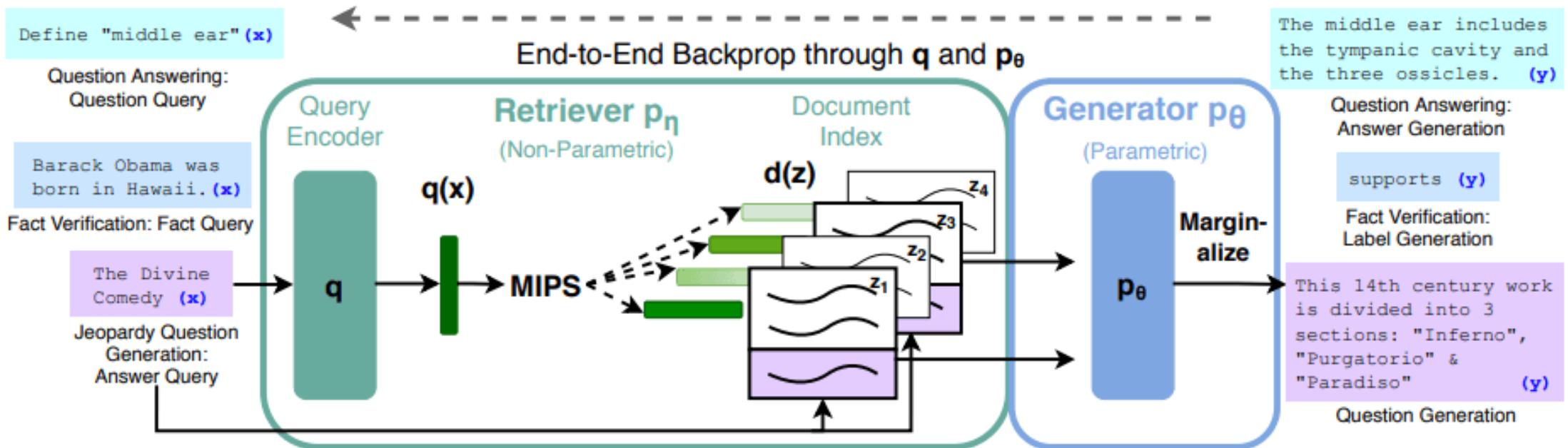


Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query x , we use Maximum Inner Product Search (MIPS) to find the top-K documents z_i . For final prediction y , we treat z as a latent variable and marginalize over seq2seq predictions given different documents.

RAG Components – Developer's Stack

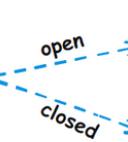


Kalyan KS

RAG Developer's Stack

www.linkedin.com/in/kalyanksnlp

LLMs



Llama 3.3



Phi-4



Gemma 3



Qwen 2.5



Mistral



DeepSeek

Frameworks



LangChain



Llama Index



Haystack



txtai

Vector Databases



Chroma



Pincone



Qdrant



Weaviate



Milvus

Data Extraction



Crawl4AI



FireCrawl



ScrapeGraphAI



MegaParser



Docling



Llama Parse



ExtractThinker

Open LLMs Access



Hugging Face



Ollama

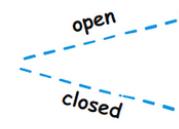


Groq



Together AI

Text Embeddings



NOMIC



SBERT



BGE



Ollama



OpenAI



Voyage AI



Google



Cohere

Evaluation



Giskard

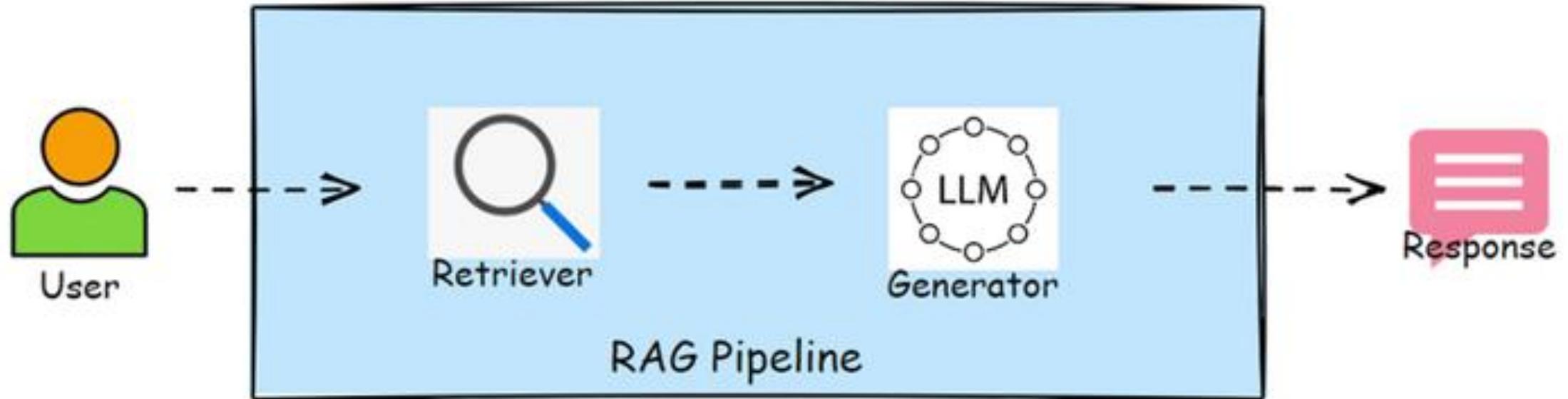


ragas



trulens

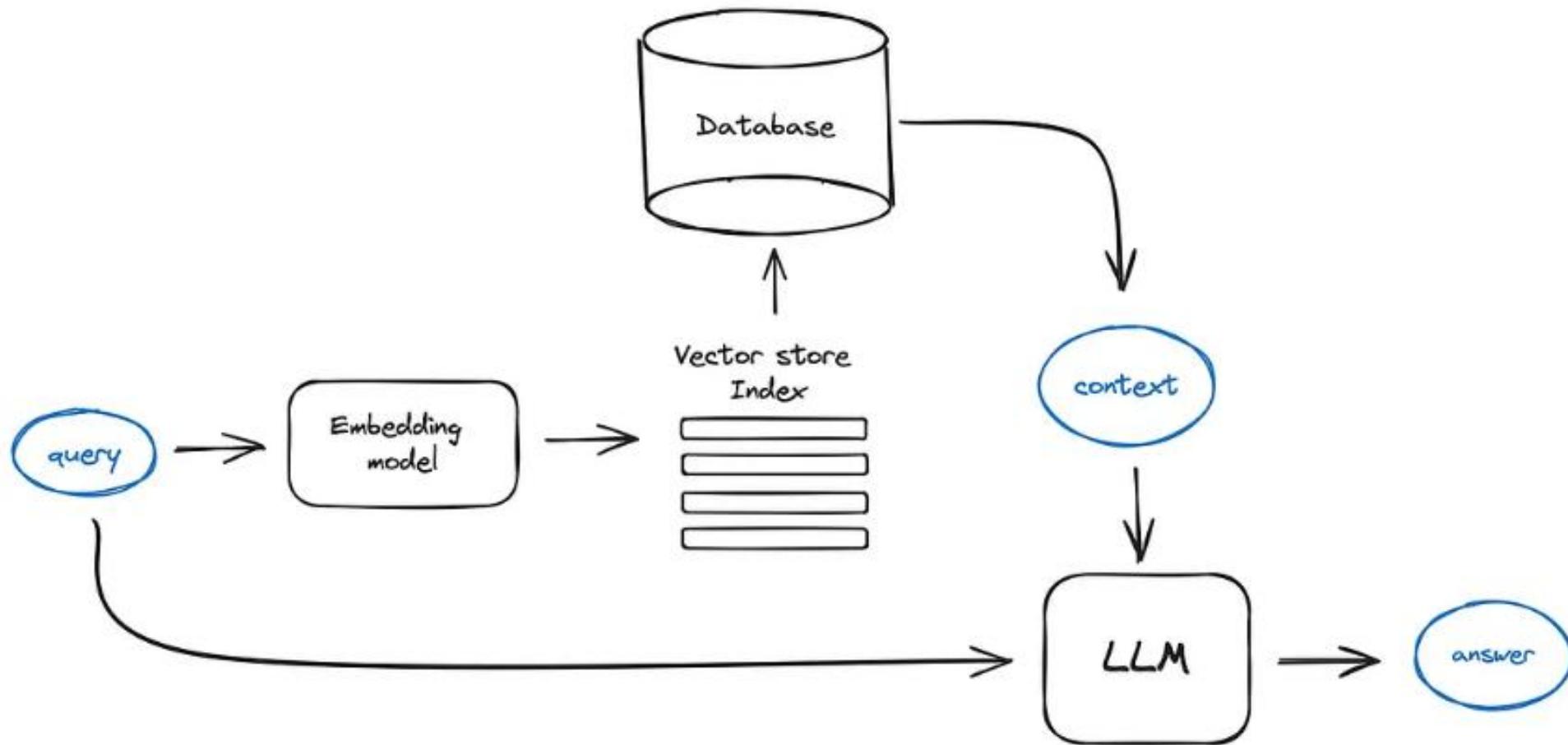
Evaluation Metrics



Agentic RAG

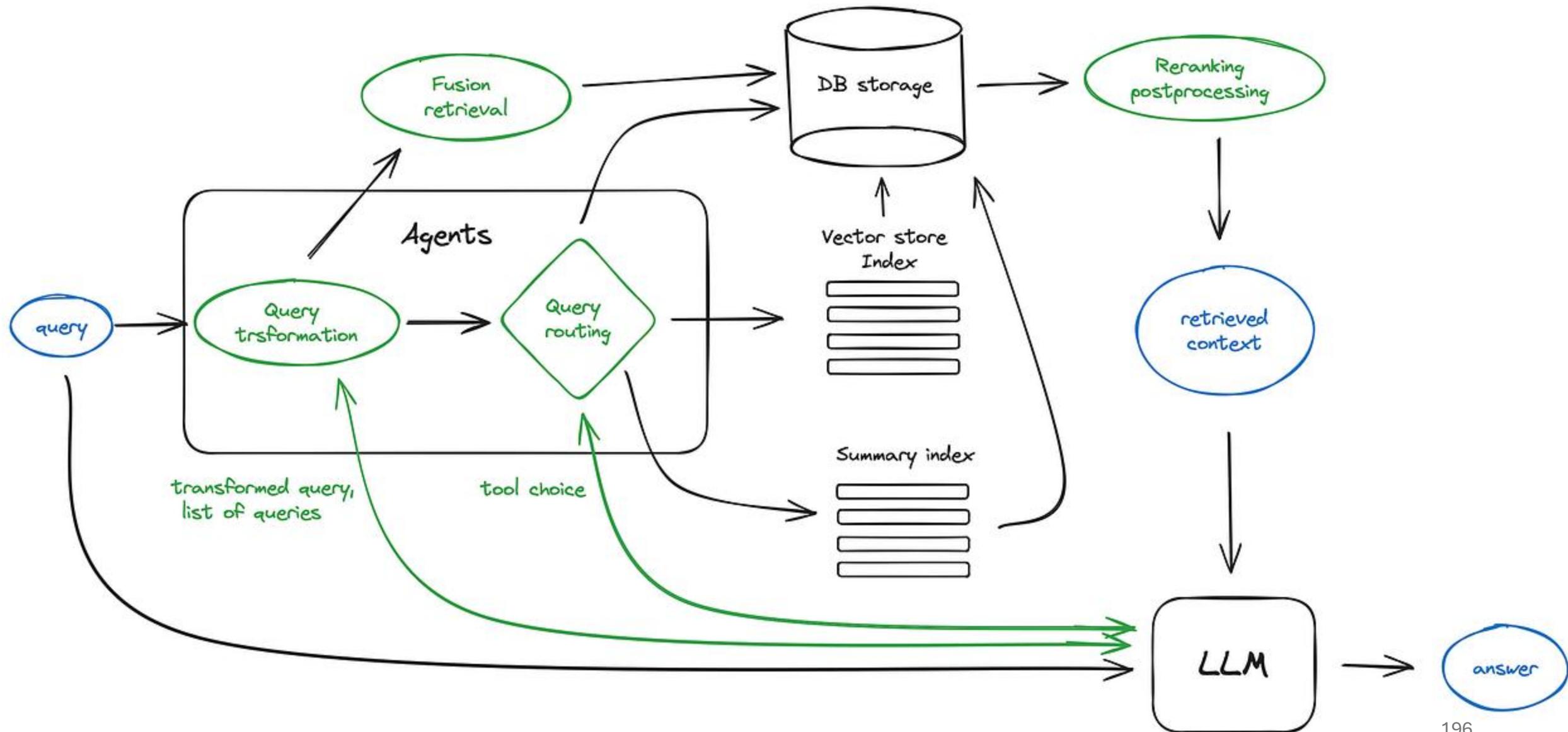
- **Differentiator:** Uses autonomous agents to plan, choose tools (vector DBs, web/APIs, SQL), and iteratively refine retrieval steps, turning a static pipeline into an adaptive reasoning loop.
- **When to use:** For queries that may need tool choice and escalation, such as “Summarize current pricing for Vendor Z and verify with their API if the document set lacks 2025 data.”
- **Costs/trade-offs:** Multi-step planning/tool calls add latency and token spend, and orchestration complexity raises observability and failure-handling burdens.

Naive RAG

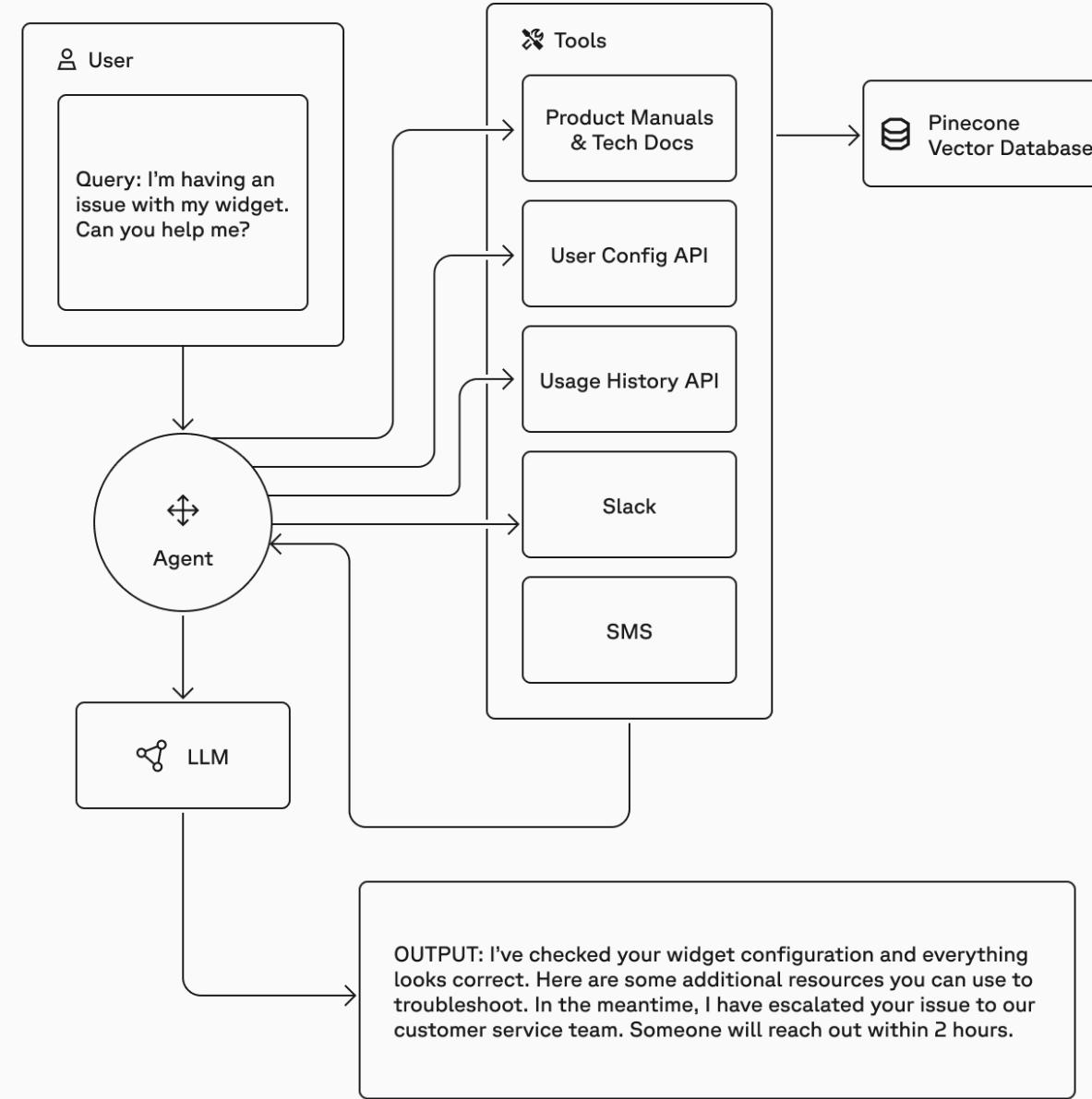


A scheme by author, as well all the schemes further in the text

Advanced RAG



Agentic RAG



Feature	Traditional RAG	Agentic RAG
Data Processing	Linear, retrieves data from specific sources	Multi-source, refines and adapts data in real-time
Decision-Making	Passive, follows user-guided queries	Autonomous, adjusts responses based on real-time analysis
Complexity Handling	Suitable for simple or direct queries	Handles multi-step reasoning and dynamic queries
Application Examples	Customer support, FAQ bots	Healthcare diagnostics, legal research, financial analytics

Until the next time 😊

