

Sesión 8.0

Large language models

Training, quantization



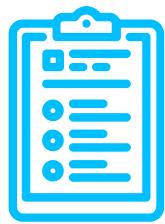
Simplest
Diffusion paper



Hardest LLM paper



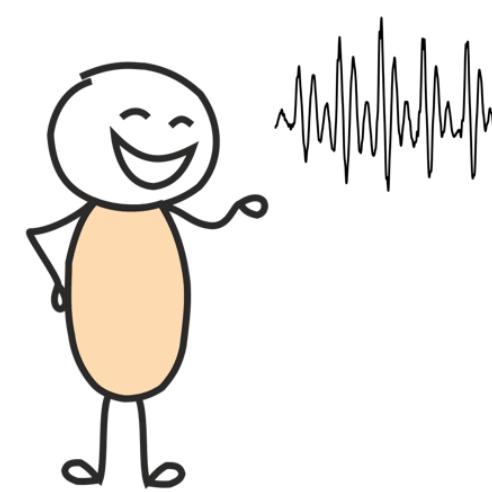
1.



Natural language *processing*

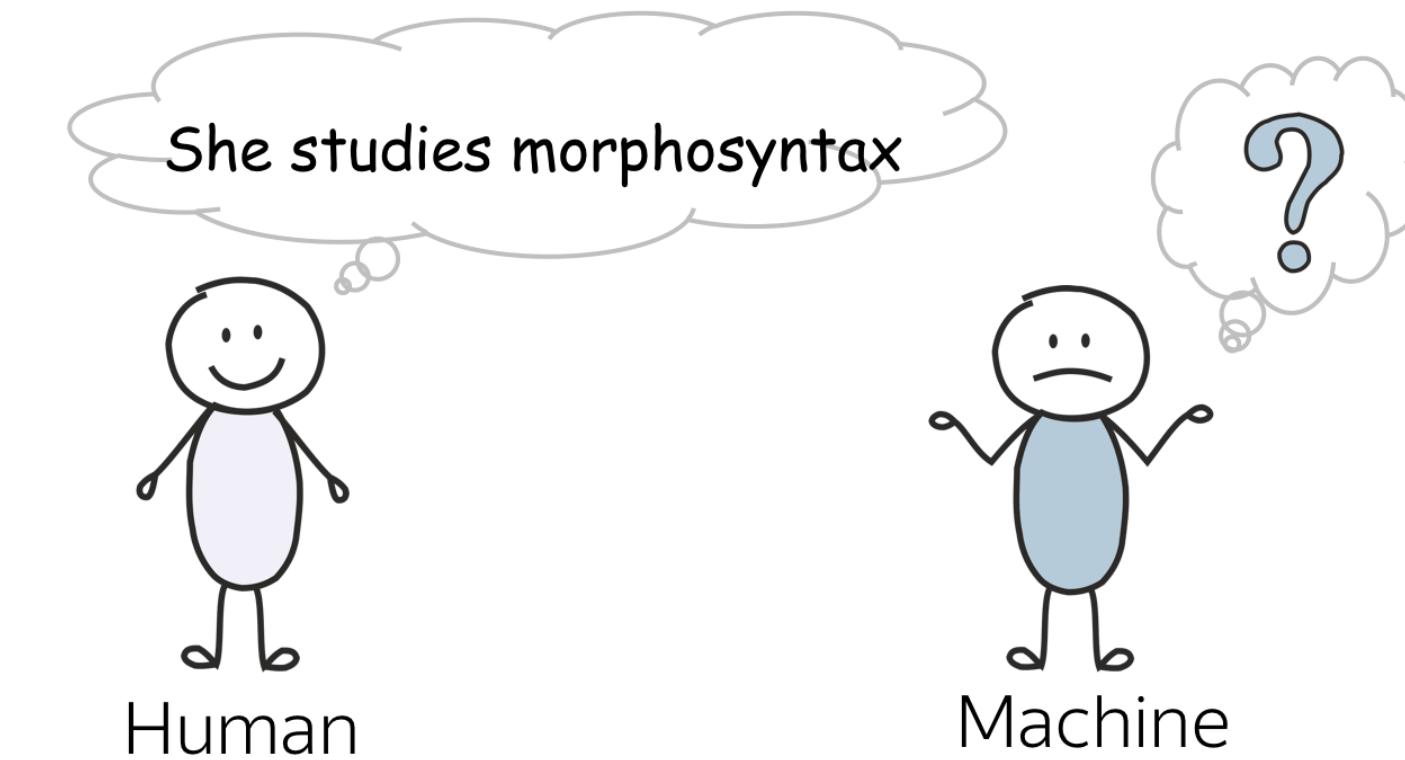


Natural language processing



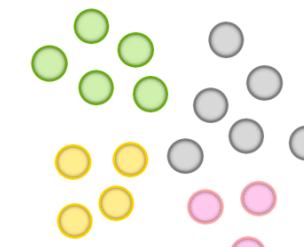
Similarly sounding options

She studies morphosyntax
She studies more faux syntax
She studies morph or syntax
....



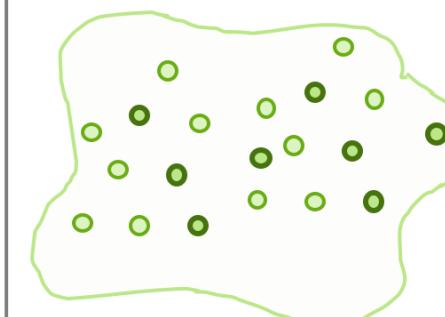
Natural language processing

What is the probability to pick a green ball?



$$\frac{5}{5 + 6 + 4 + 3} = \frac{5}{18}$$

Can we do the same for sentences?



Text corpus

$$P(\text{the mut is tinning the tebn}) = \frac{0}{|\text{corpus}|} = 0$$

$$P(\text{mut the tinning tebn is the}) = \frac{0}{|\text{corpus}|} = 0$$

With this approach, sentences that never occurred in the corpus will receive zero probability

But the first sentence is “more likely” than the second!

This method is not good!



Natural language processing

P(**I**) =

P(**I**)



Probability of **I**

$$P(y_1, y_2, \dots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \dots \cdot P(y_n|y_1, \dots, y_{n-1}) = \prod_{t=1}^n P(y_t|y_{<t}).$$



N-gram

$$P(y_1, y_2, \dots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \dots \cdot P(y_n|y_1, \dots, y_{n-1}) = \prod_{t=1}^n P(y_t|y_{<t}).$$

La probabilidad de una palabra solo depende de un número fijo de palabras anteriores.

n=3 (trigram model) $P(y_t|y_1, \dots, y_{t-1}) = P(y_t|y_{t-2}, y_{t-1})$

n=4 (bigram model) $P(y_t|y_1, \dots, y_{t-1}) = P(y_t|y_{t-1})$

n=1 (unigram model) $P(y_t|y_1, \dots, y_{t-1}) = P(y_t)$



N-gram

Before

$$P(I \text{ saw a cat on a mat}) =$$

- $P(I)$
- $P(\text{saw} | I)$
- $P(a | I \text{ saw})$
- $P(\text{cat} | I \text{ saw a})$
- $P(\text{on} | I \text{ saw a cat})$
- $P(a | I \text{ saw a cat on})$
- $P(\text{mat} | I \text{ saw a cat on a})$

After (3-gram)

$$P(I \text{ saw a cat on a mat}) =$$



- $P(I)$ → $P(I)$
 - $P(\text{saw} | I)$ → • $P(\text{saw} | I)$
 - $P(a | I \text{ saw})$ → • $P(a | I \text{ saw})$
 - $P(\text{cat} | I \text{ saw a})$ → • $P(\text{cat} | \text{saw a})$
 - $P(\text{on} | I \text{ saw a cat})$ → • $P(\text{on} | a \text{ cat})$
 - $P(a | I \text{ saw a cat on})$ → • $P(a | \text{cat on})$
 - $P(\text{mat} | I \text{ saw a cat on a})$ → • $P(\text{mat} | \text{on a})$
- ignore use

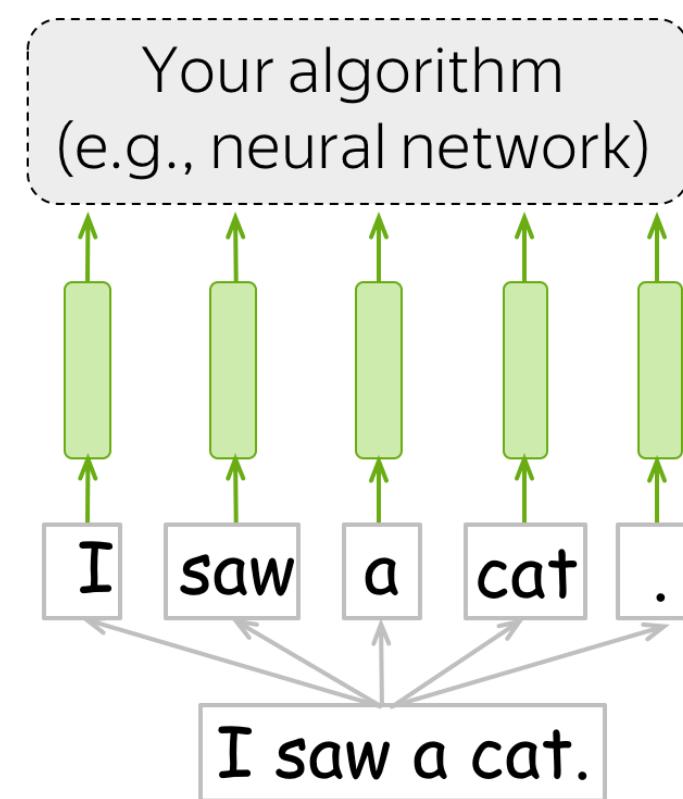


Modelo de *lenguaje*

I _____



Word *Embeddings*



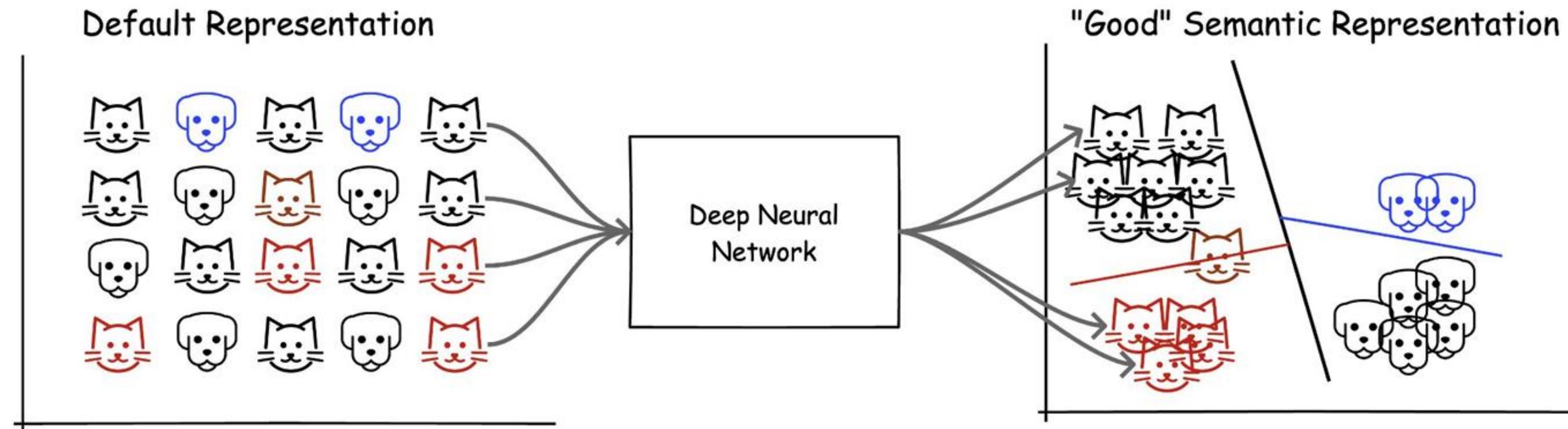
Any algorithm for solving a task

Word representation - vector
(input for your model/algorithm)

Sequence of tokens

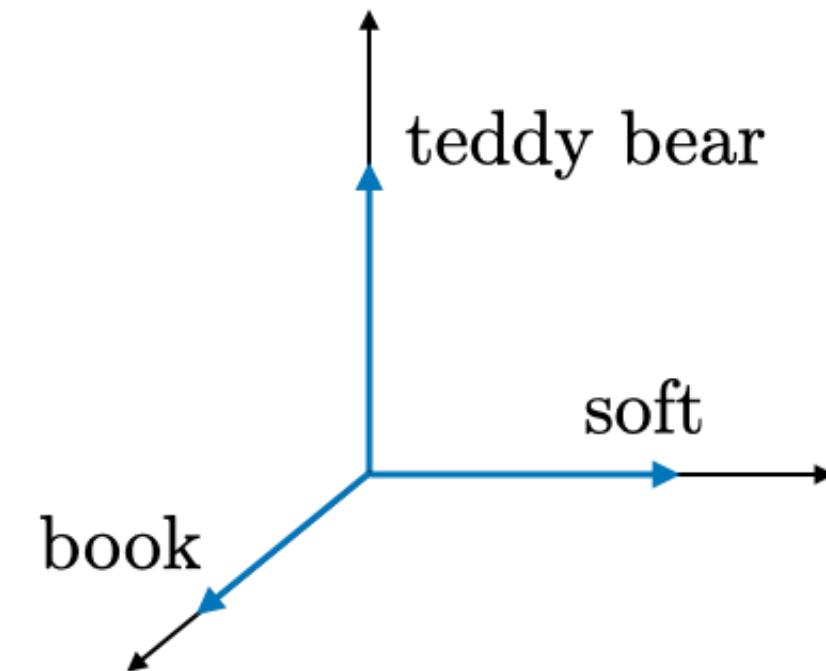
Text (your input)

Word *Embeddings*



One-hot Encoding

Rome = [1, 0, 0, 0, 0, 0, ..., 0]
 Rome = [1, 0, 0, 0, 0, 0, ..., 0] word V
 Paris = [0, 1, 0, 0, 0, 0, ..., 0]
 Italy = [0, 0, 1, 0, 0, 0, ..., 0]
 France = [0, 0, 0, 1, 0, 0, ..., 0]



Word2Vec

$$P(w_{t-2}|w_t) \ P(w_{t-1}|w_t) \ P(w_{t+1}|w_t) \ P(w_{t+2}|w_t)$$

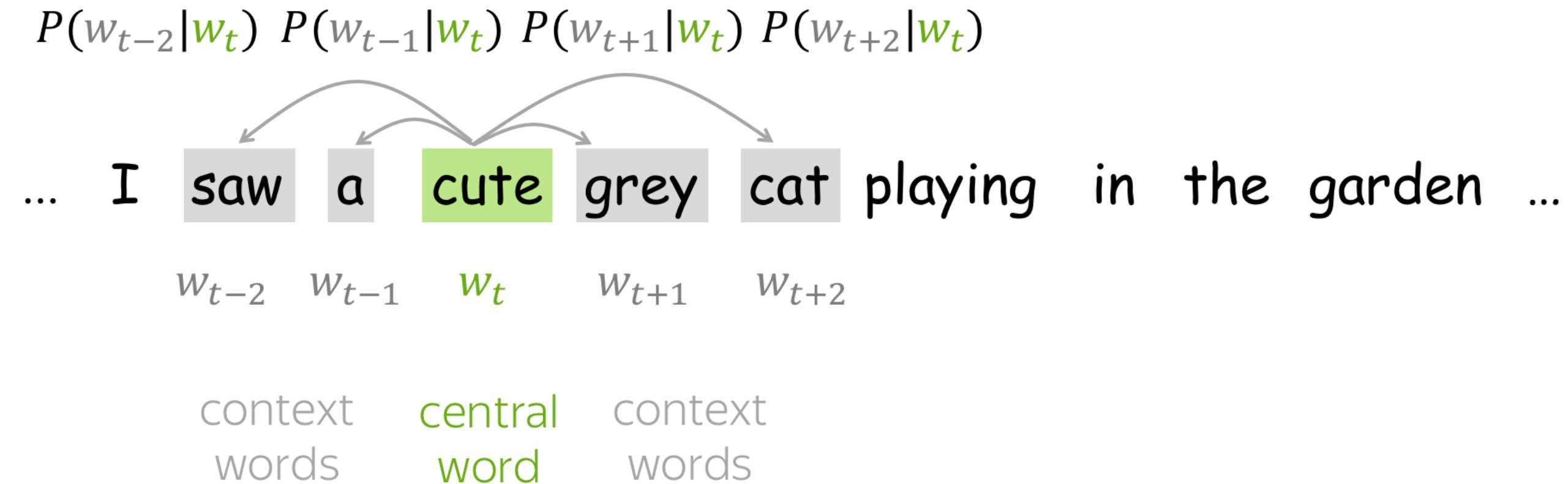


$w_{t-2} \quad w_{t-1} \quad w_t \quad w_{t+1} \quad w_{t+2}$

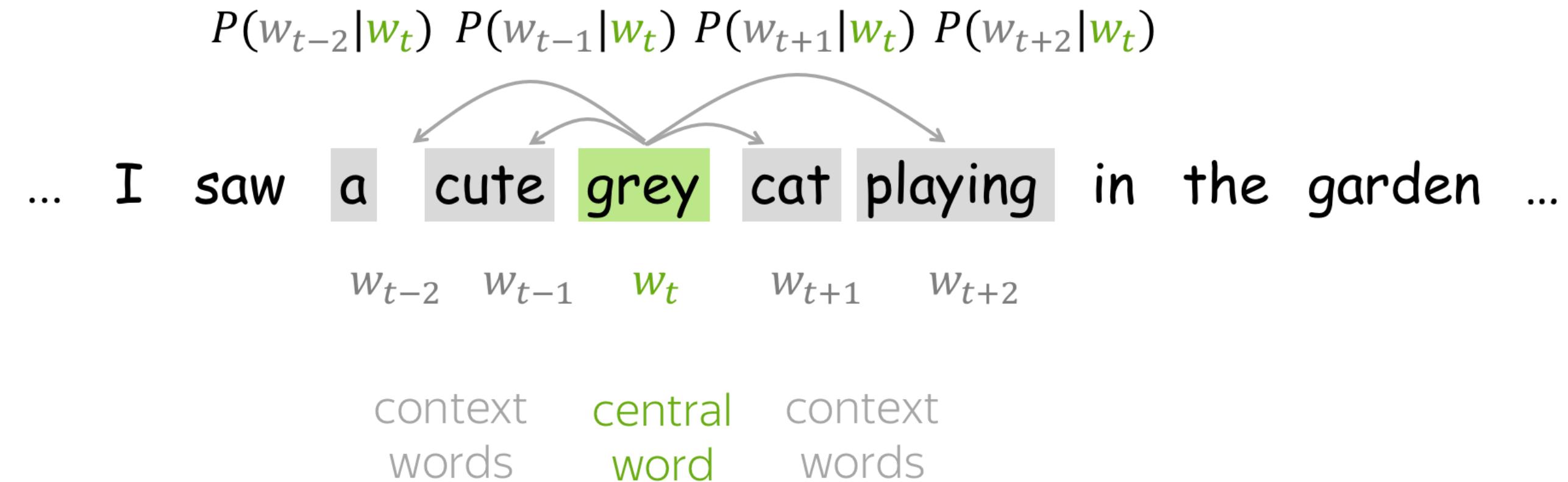
context central context
words word words



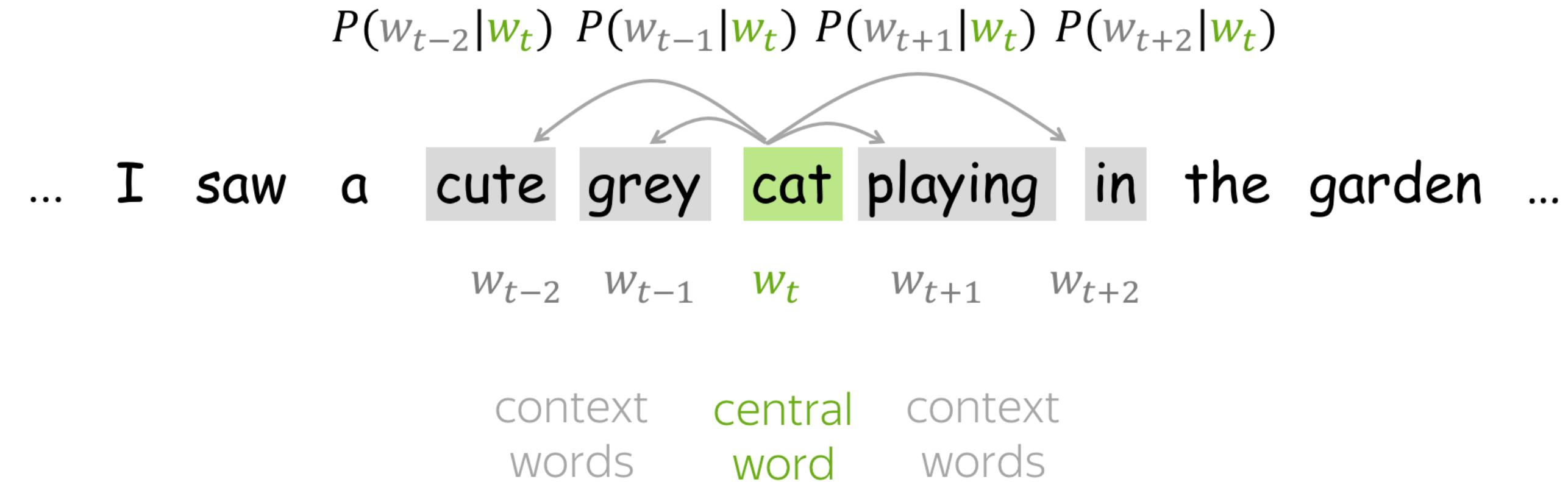
Word2Vec



Word2Vec

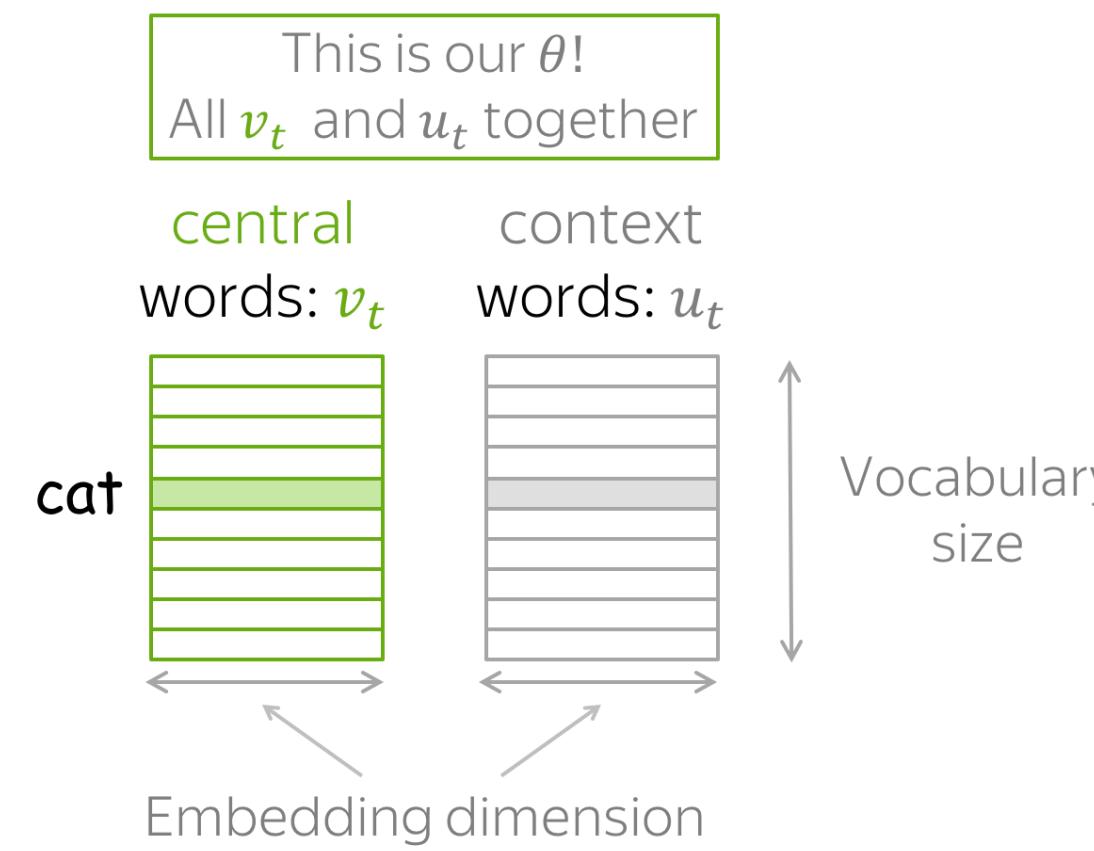


Word2Vec



Word2Vec

$$\log P(w_{t+j} | \mathbf{w}_t, \theta)$$



$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

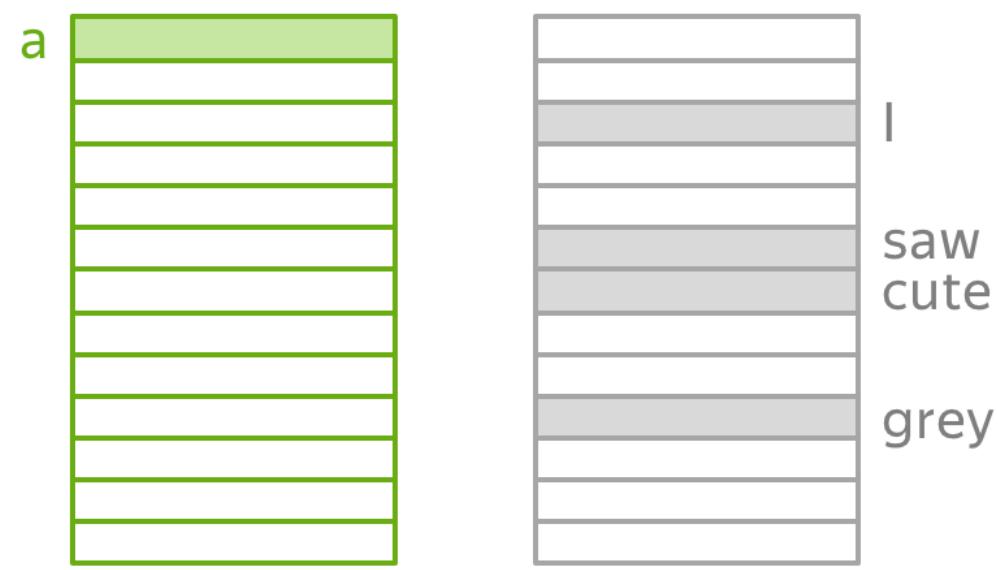
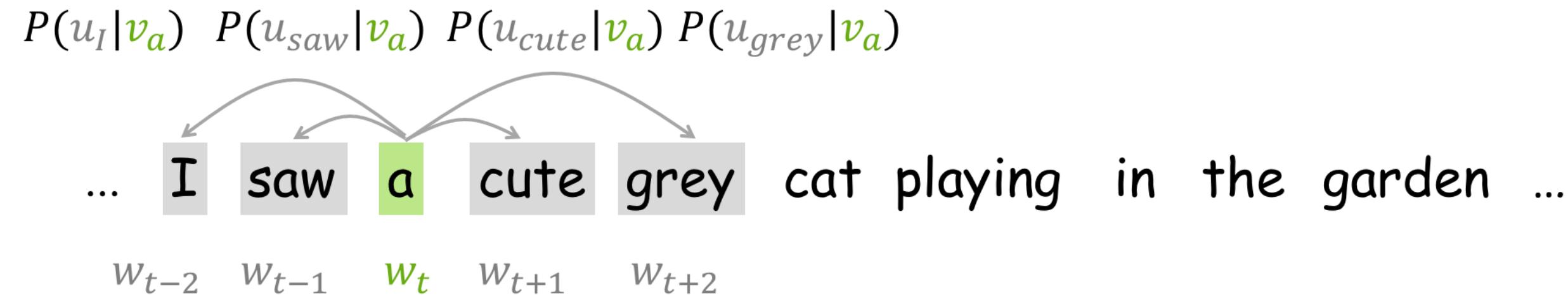
Dot product: measures similarity of o and c
Larger dot product = larger probability

Normalize over entire vocabulary
to get probability distribution

v es una palabra central.
 u es una palabra contexto.



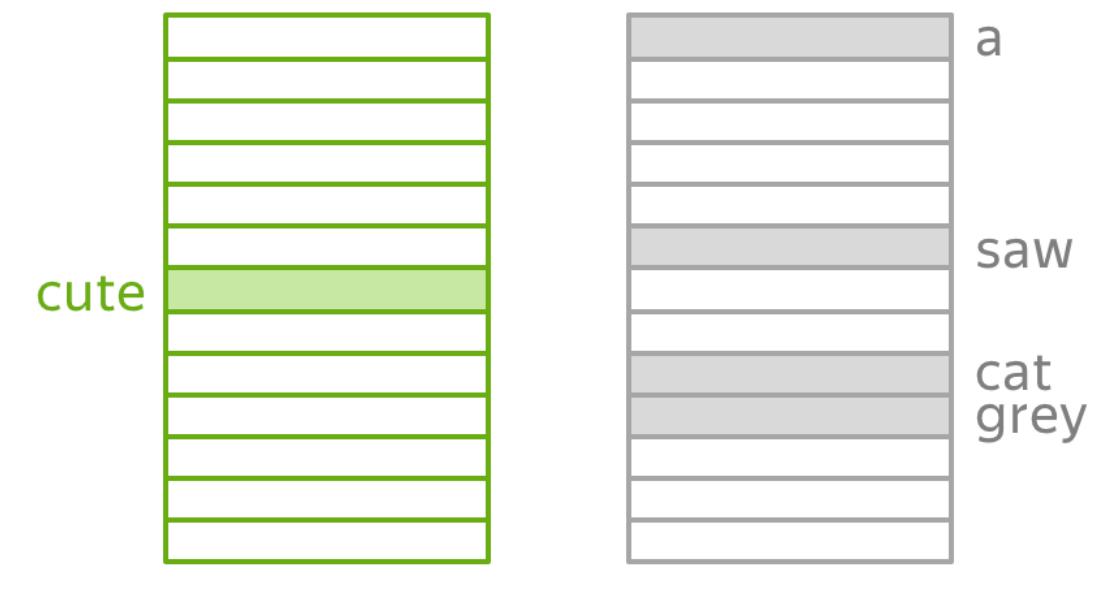
Word2Vec



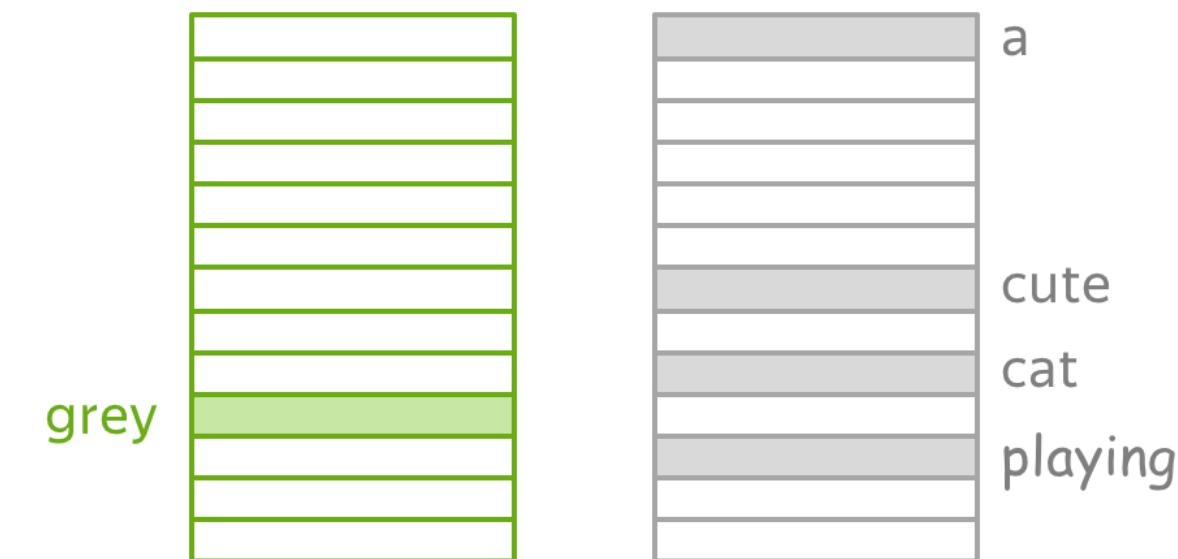
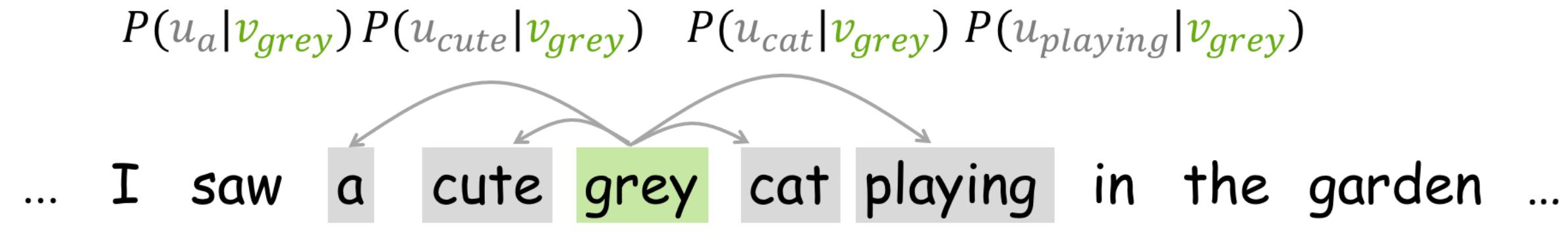
Word2Vec

$P(u_{saw}|\nu_{cute}) P(u_a|\nu_{cute}) P(u_{grey}|\nu_{cute}) P(u_{cat}|\nu_{cute})$
 ... I saw a cute grey cat playing in the garden ...

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}



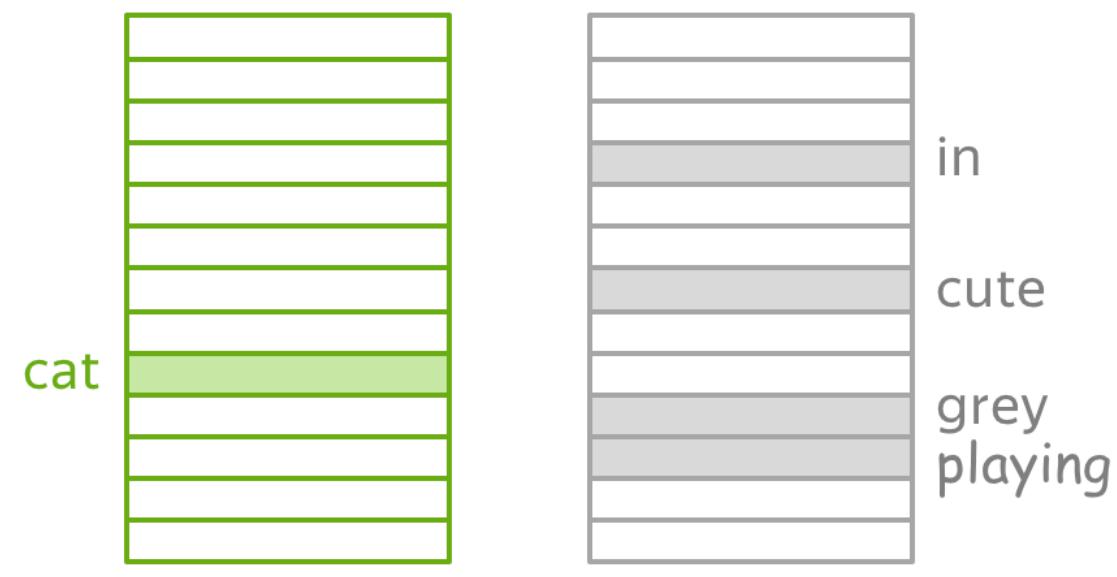
Word2Vec



Word2Vec

$P(u_{cute}|\nu_{cat})$ $P(u_{grey}|\nu_{cat})$ $P(u_{playing}|\nu_{cat})$ $P(u_{in}|\nu_{cat})$
 ... I saw a **cute** **grey** **cat** **playing** **in** the garden ...

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}



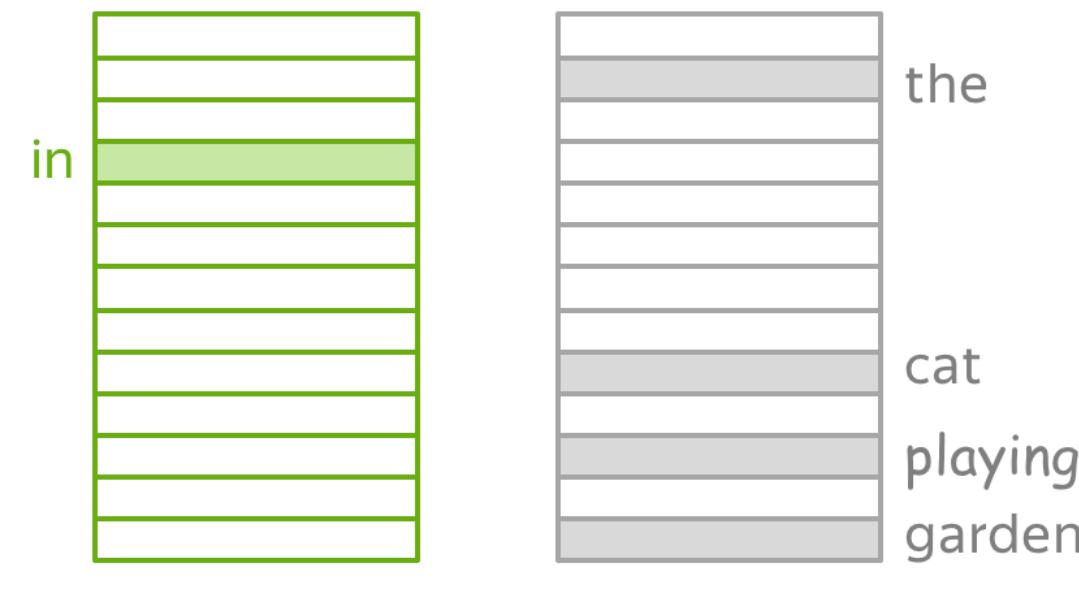
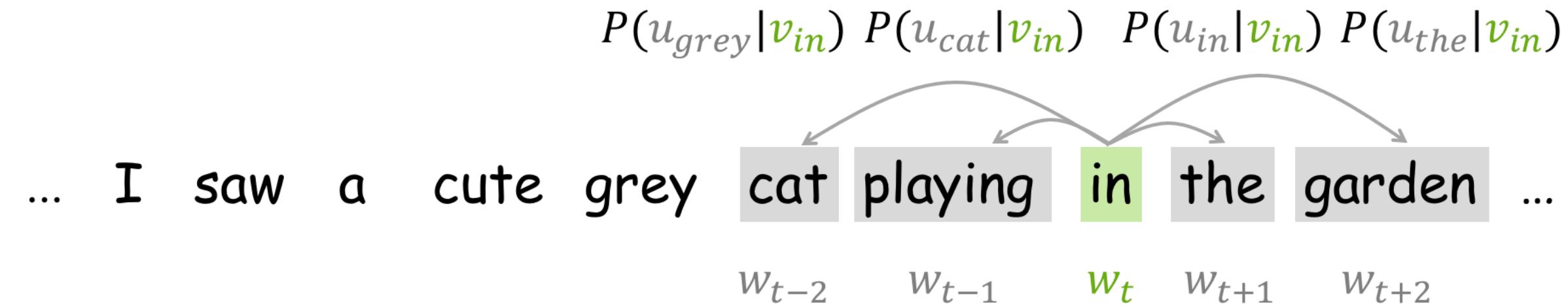
Word2Vec

$P(u_{grey}|v_{playing})$ $P(u_{cat}|v_{playing})$ $P(u_{in}|v_{playing})$ $P(u_{the}|v_{playing})$
 ... I saw a cute grey cat playing in the garden ...

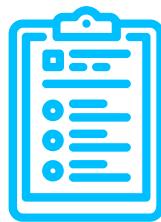
w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}



Word2Vec



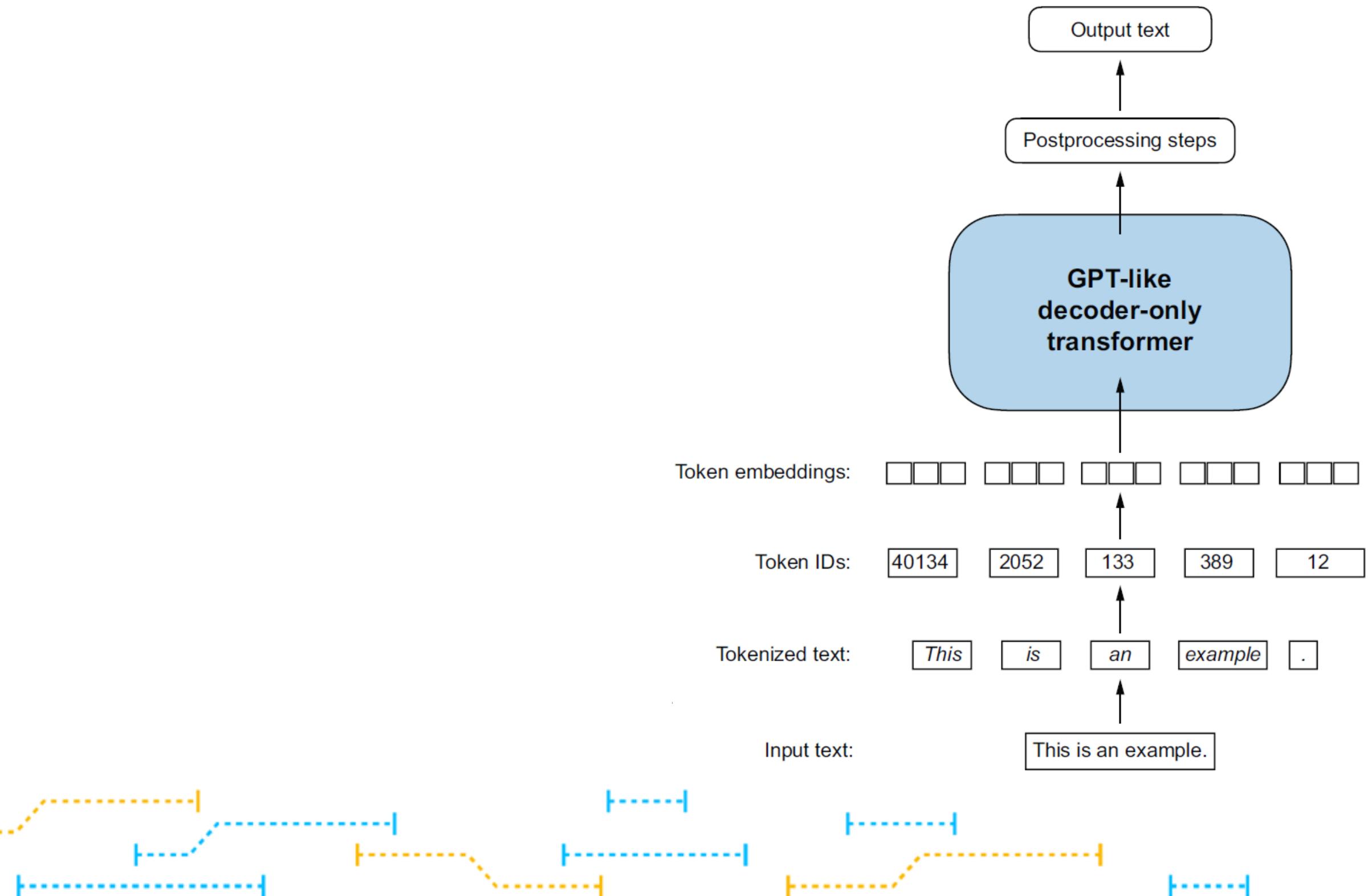
2.



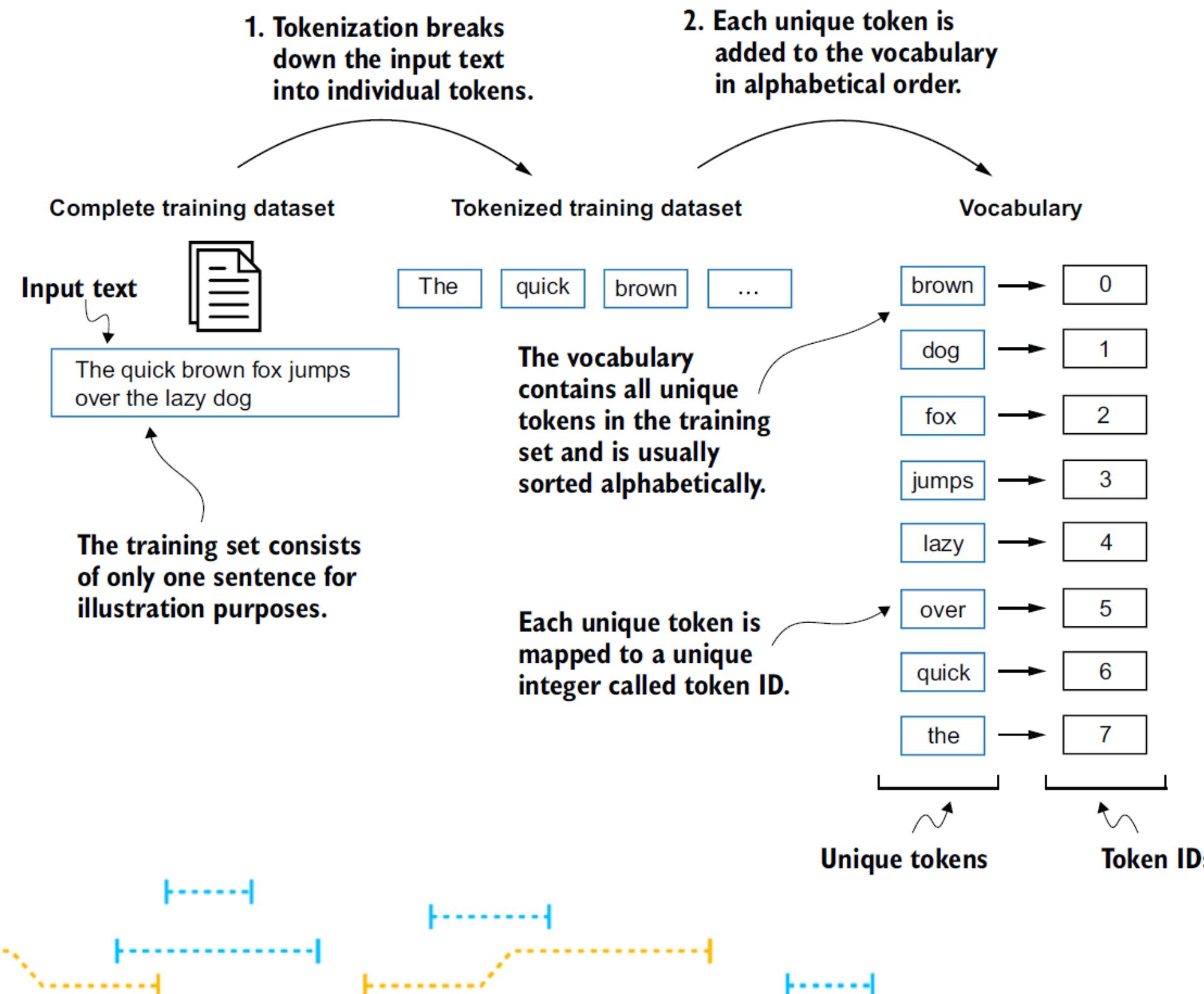
Large *Language Model*



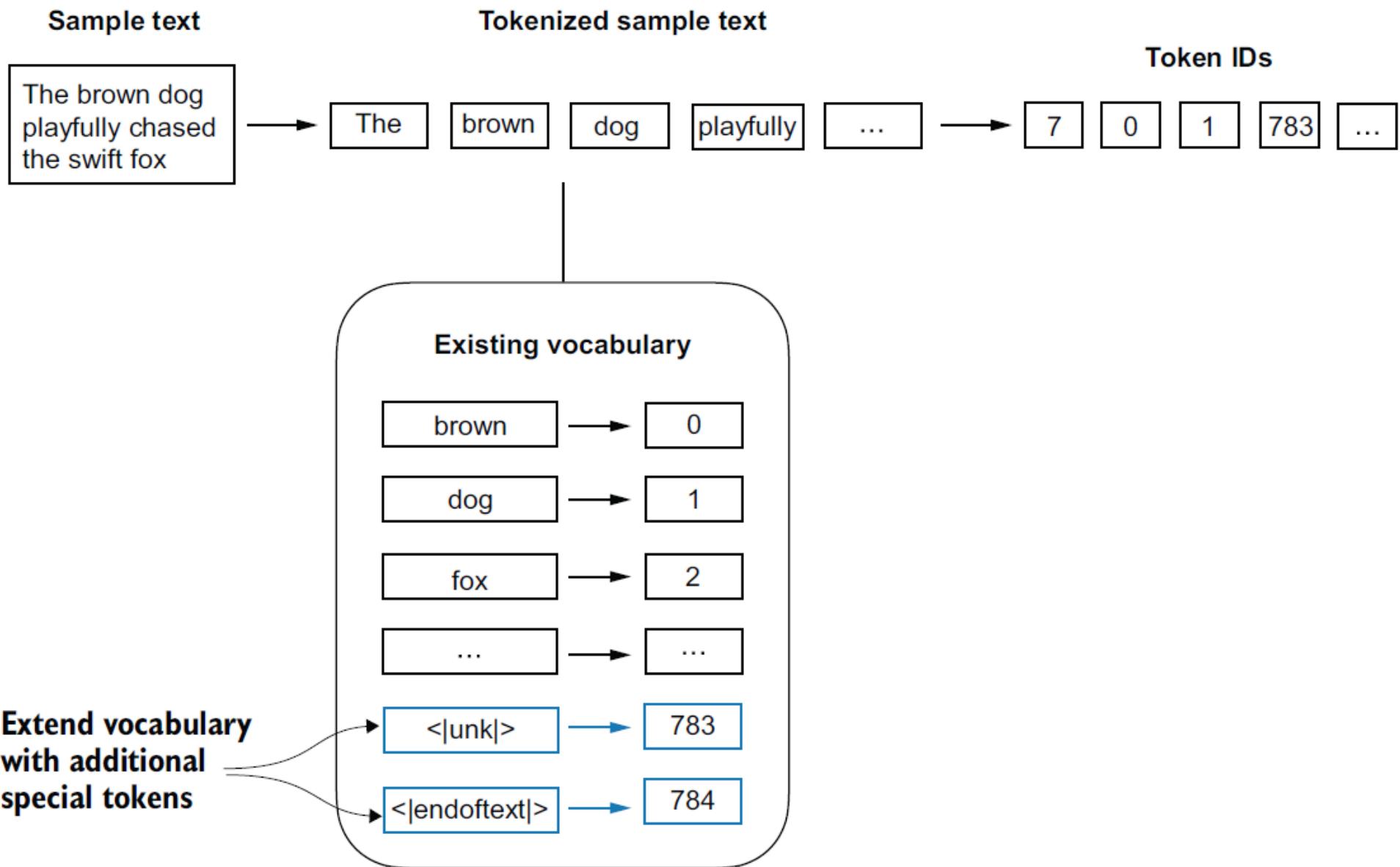
Natural language processing



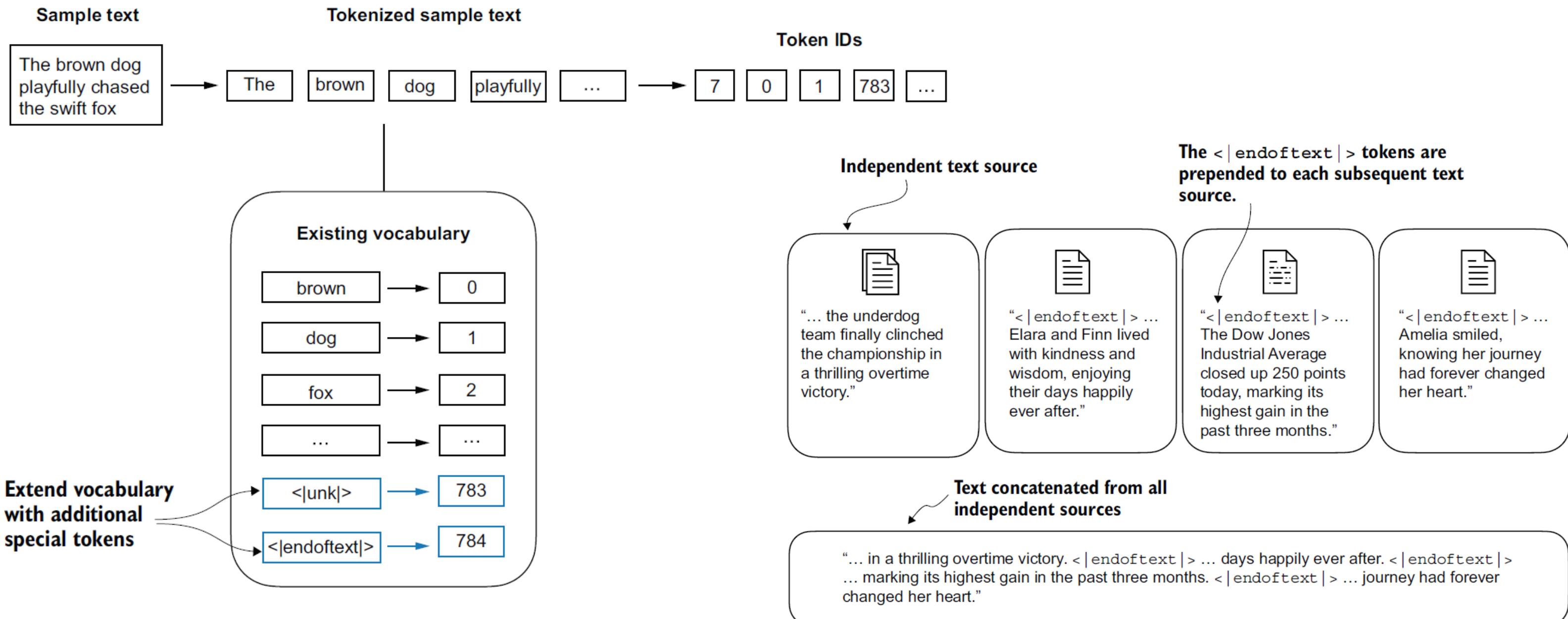
Natural language processing



Special tokens



Special tokens



Byte-Pair *Encoding*

this is the hugging face course . this chapter is about
tokenization . this section shows several tokenizer
algorithms .



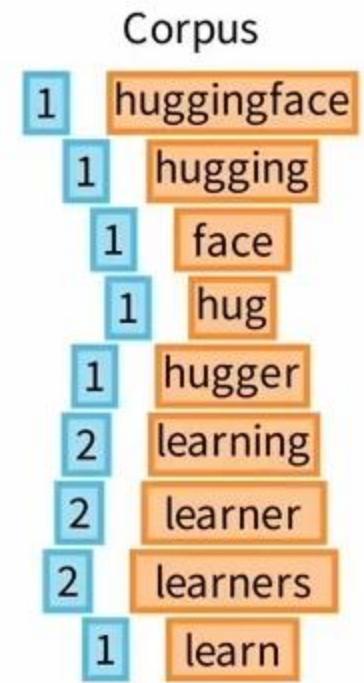
Byte-Pair Encoding

this is the hugging face course . this chapter is about
tokenization . this section shows several tokenizer
algorithms .

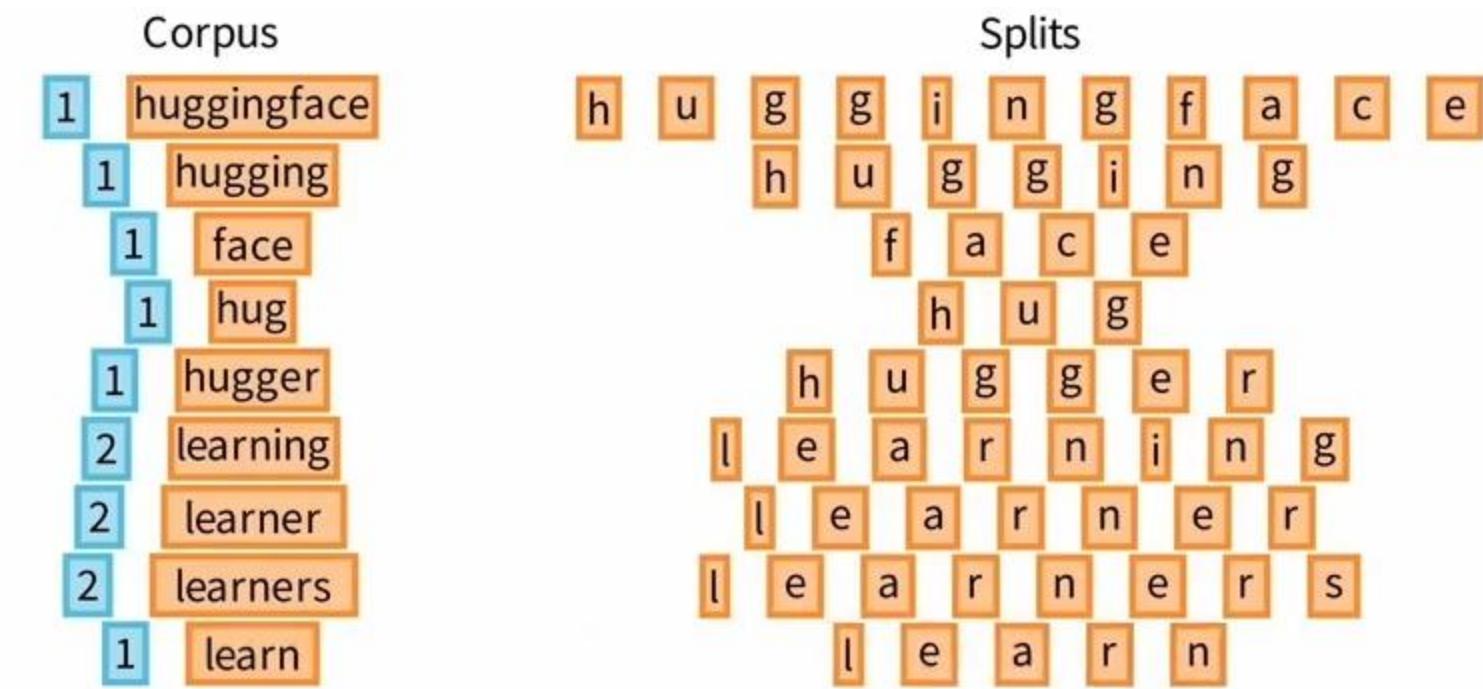
3x this 2x is 1x the
1x hugging 1x face
1x course 3x . 1x chapter
1x about 1x tokenization
1x section 1x shows
1x several 1x tokenizer
1x algorithms



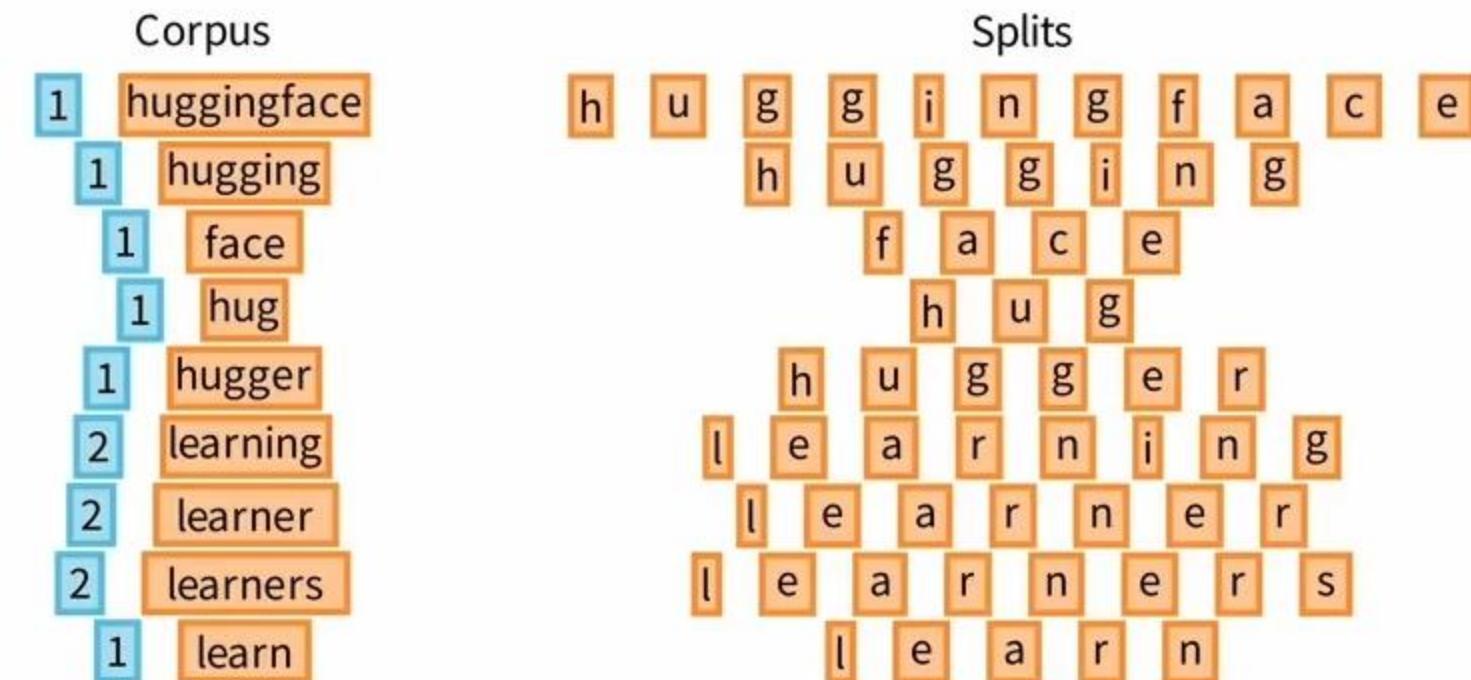
Byte-Pair Encoding



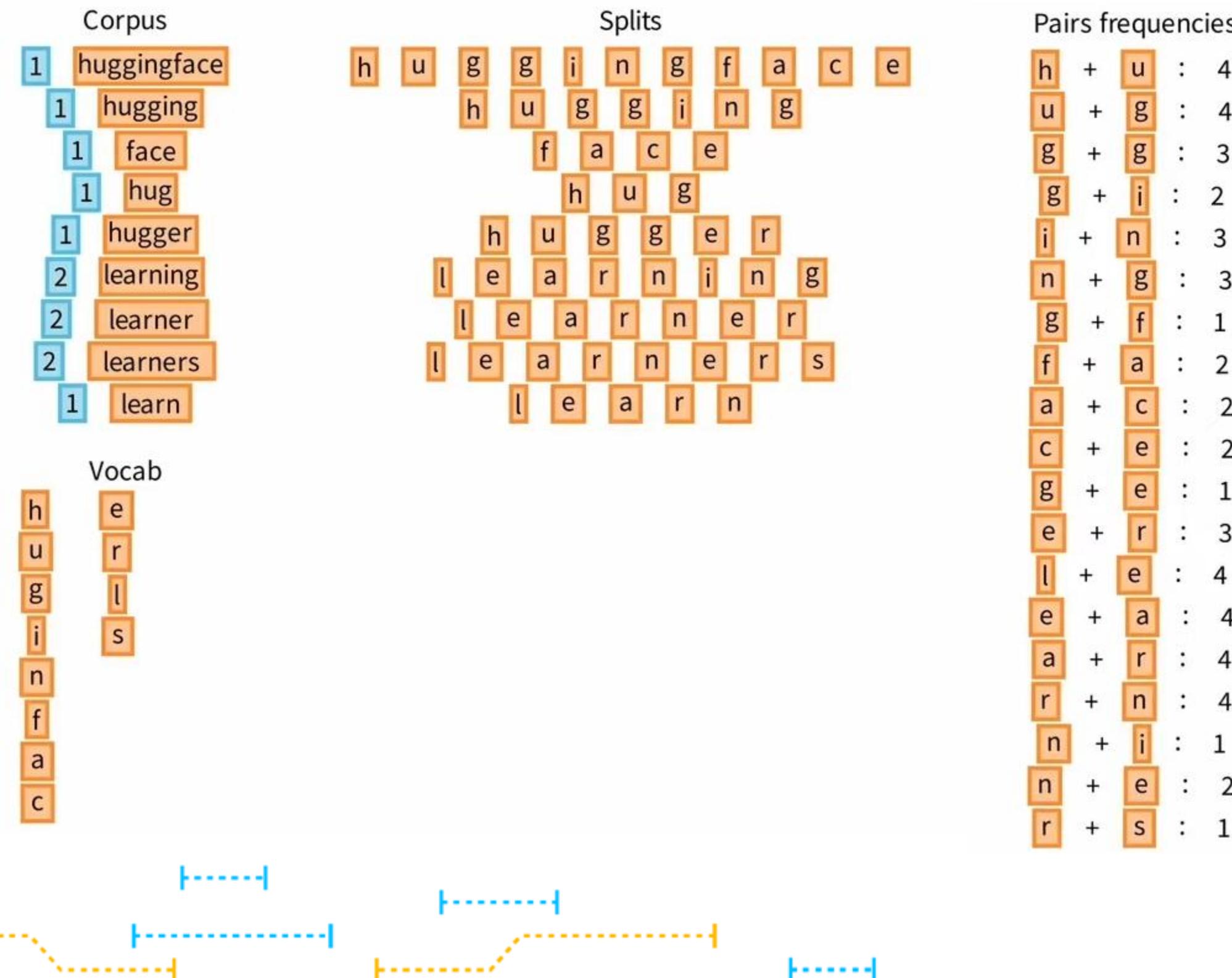
Byte-Pair Encoding



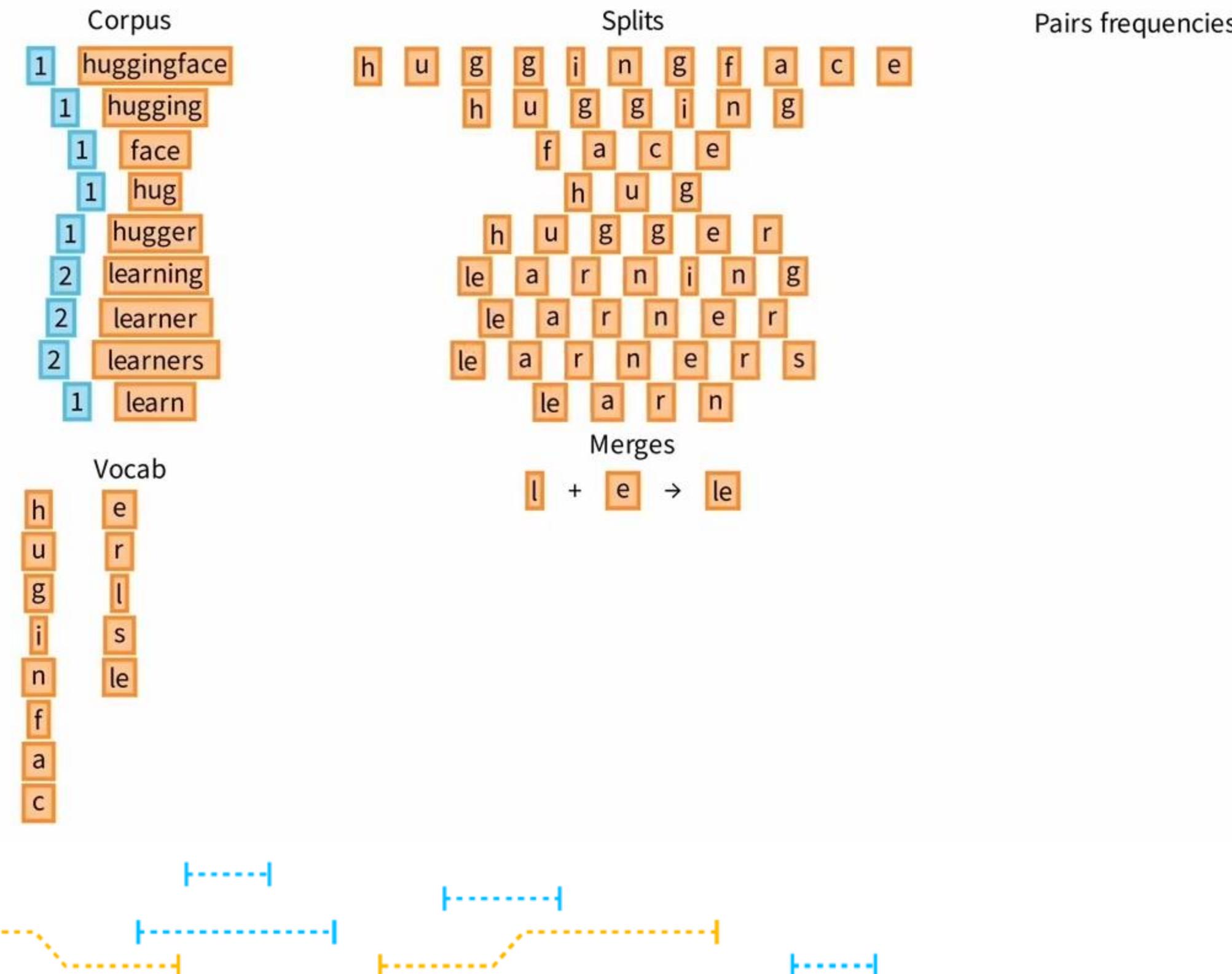
Byte-Pair Encoding



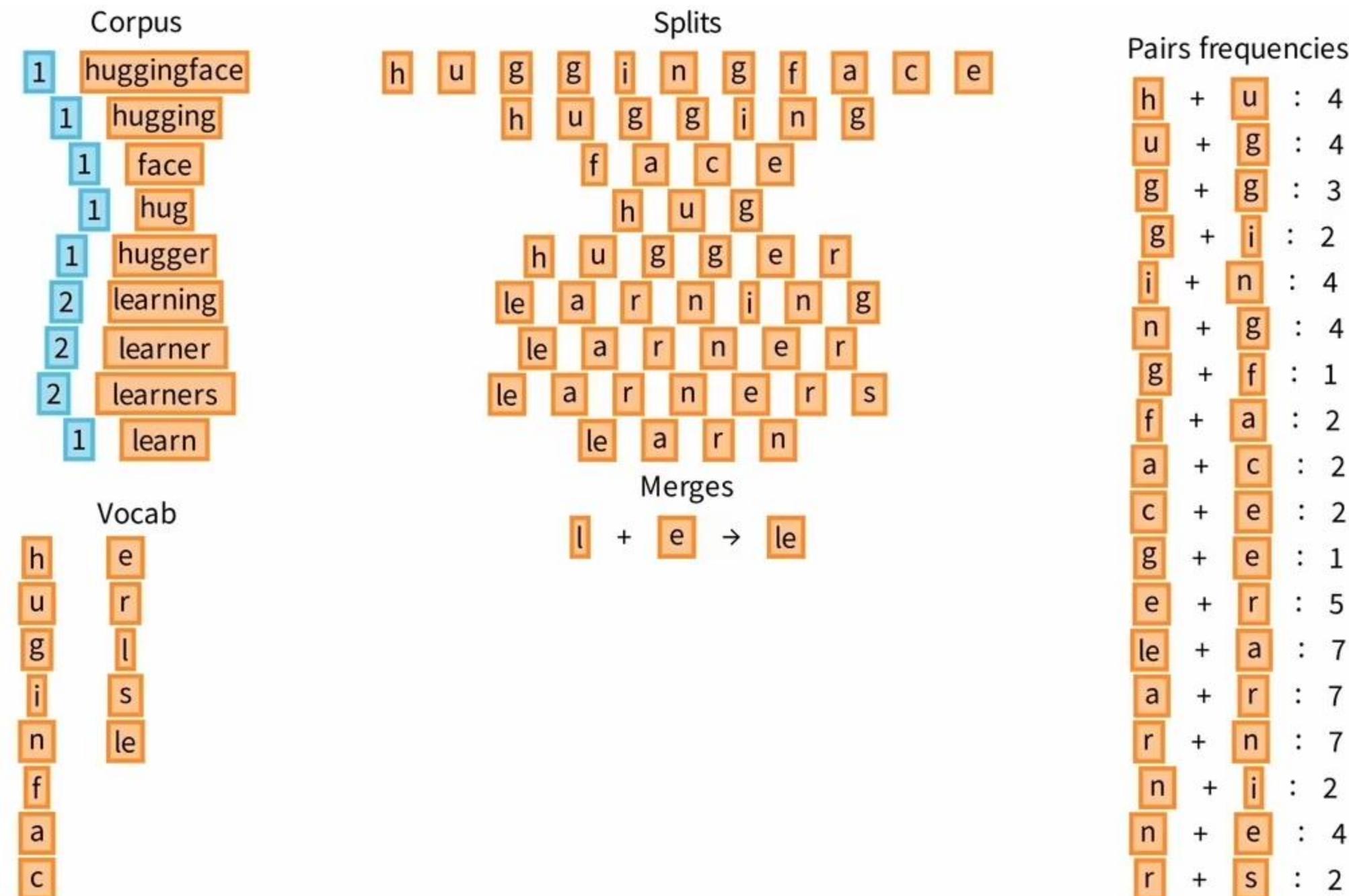
Byte-Pair Encoding



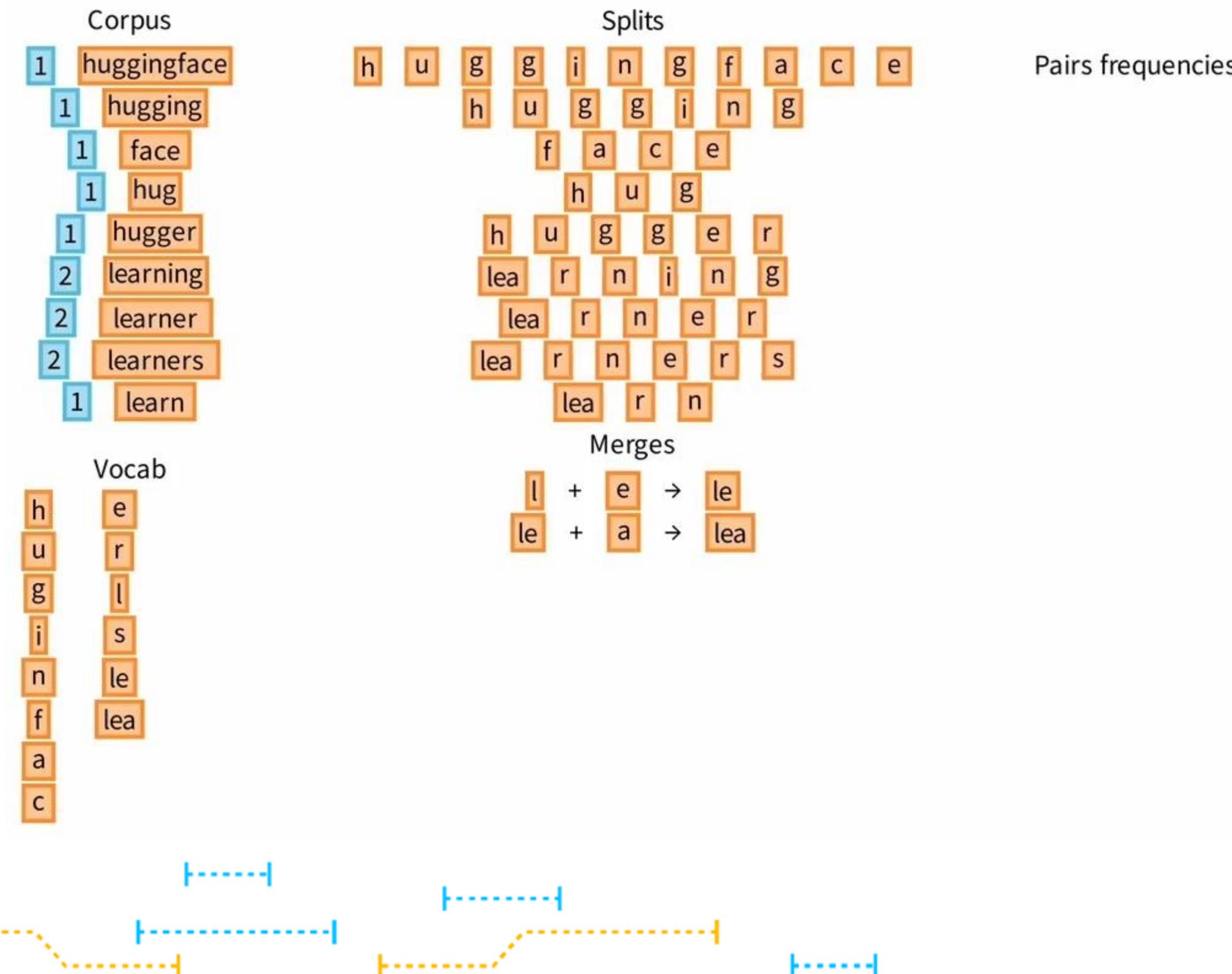
Byte-Pair Encoding



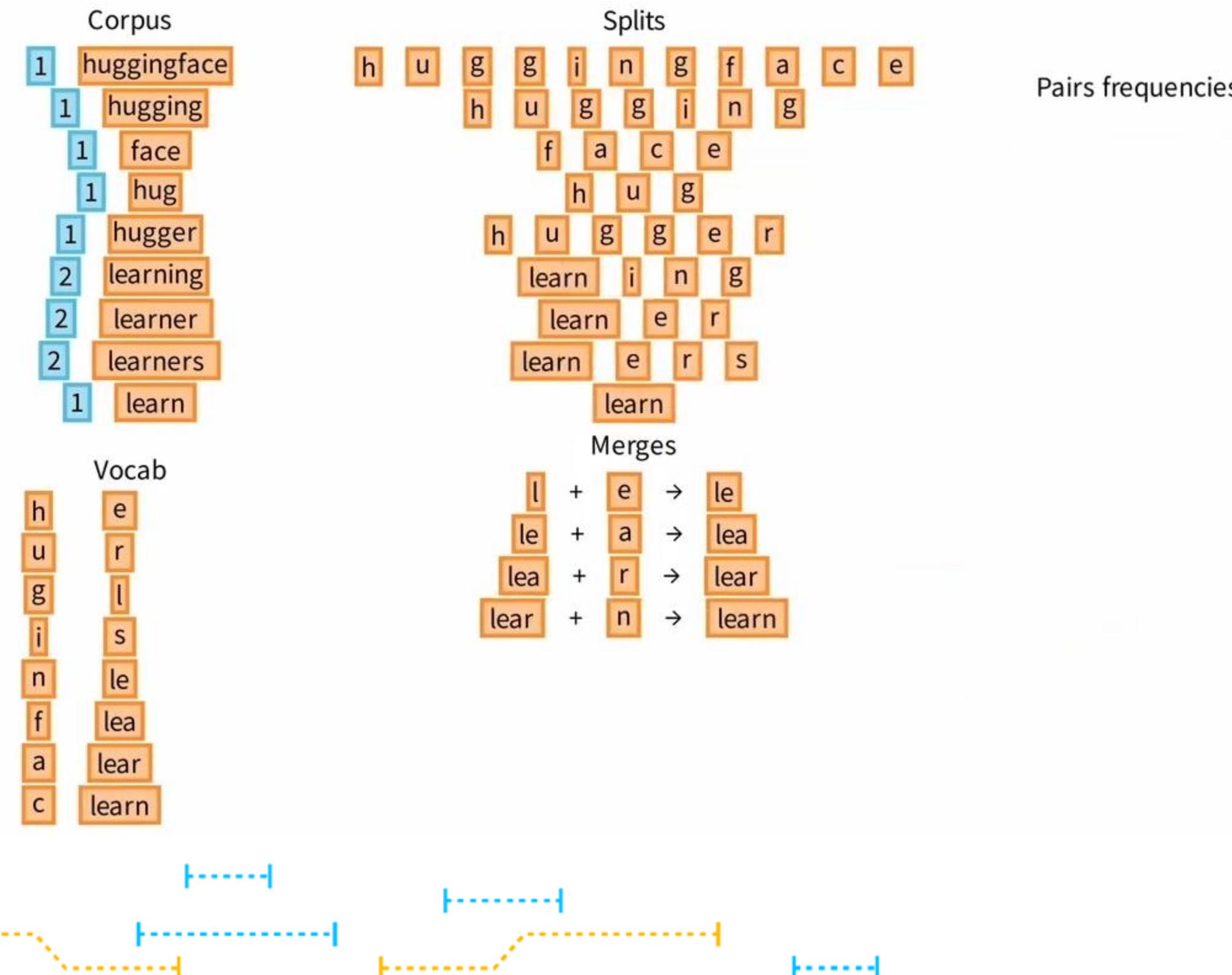
Byte-Pair Encoding



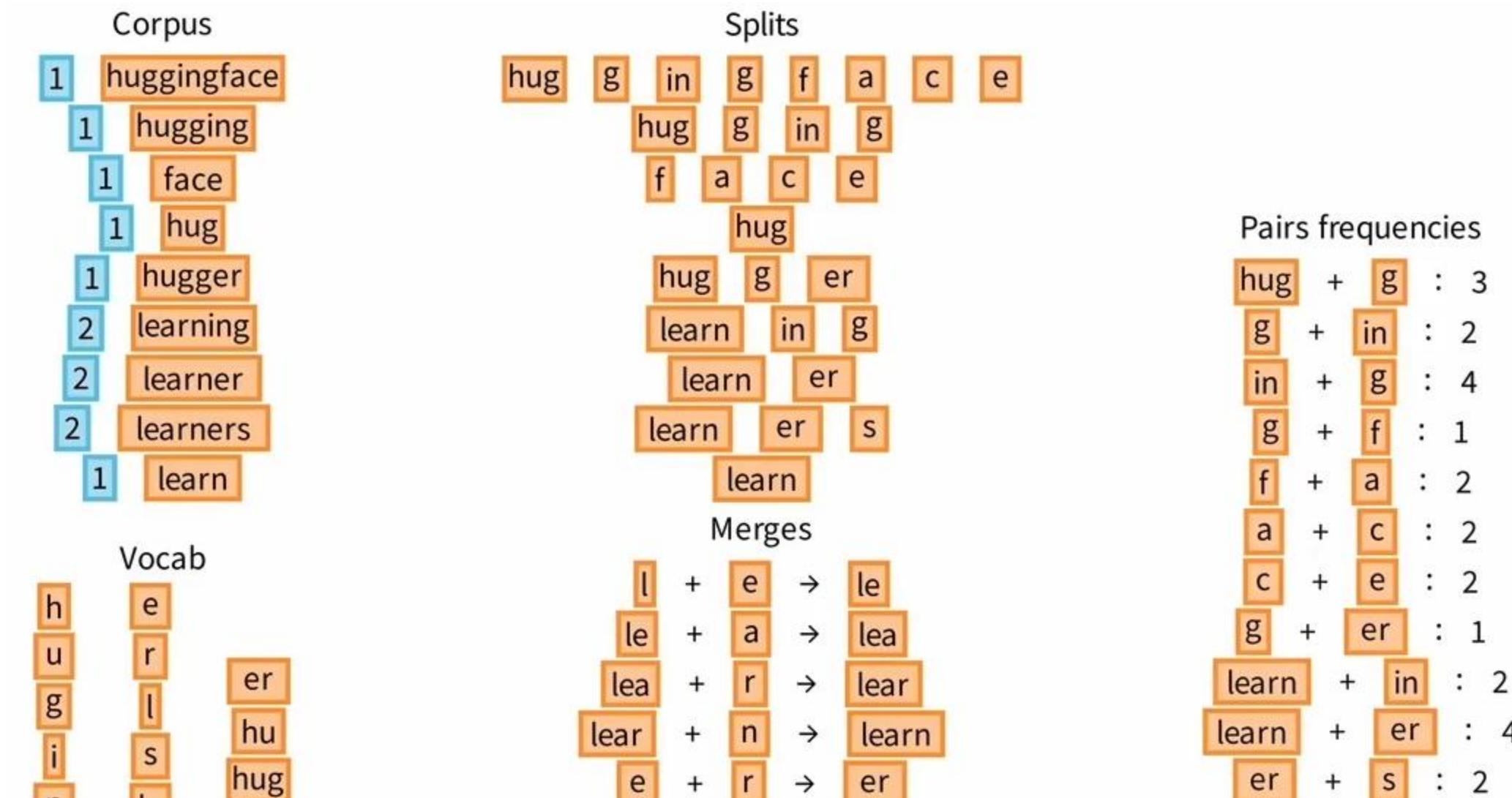
Byte-Pair Encoding



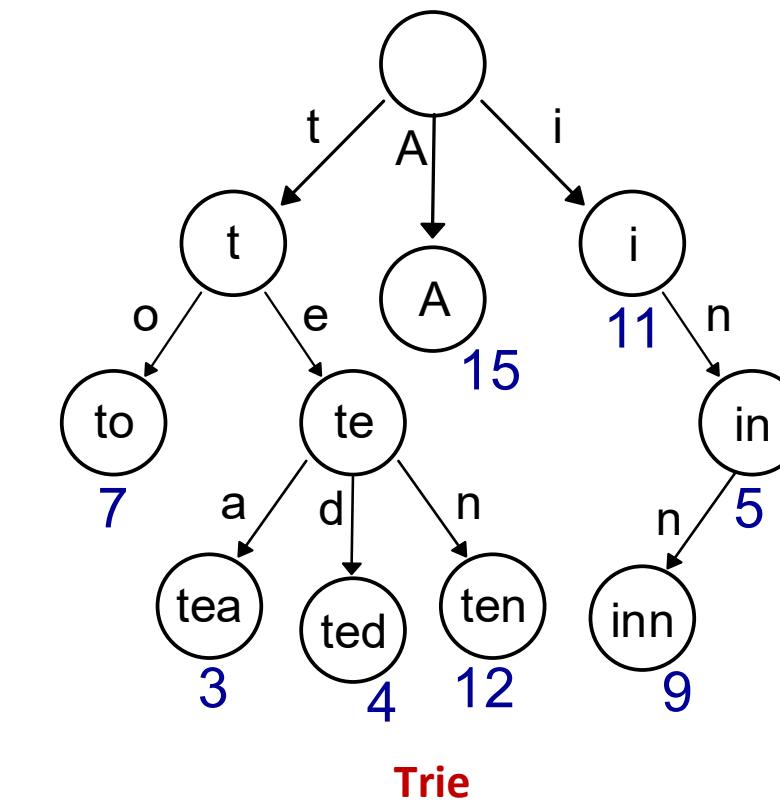
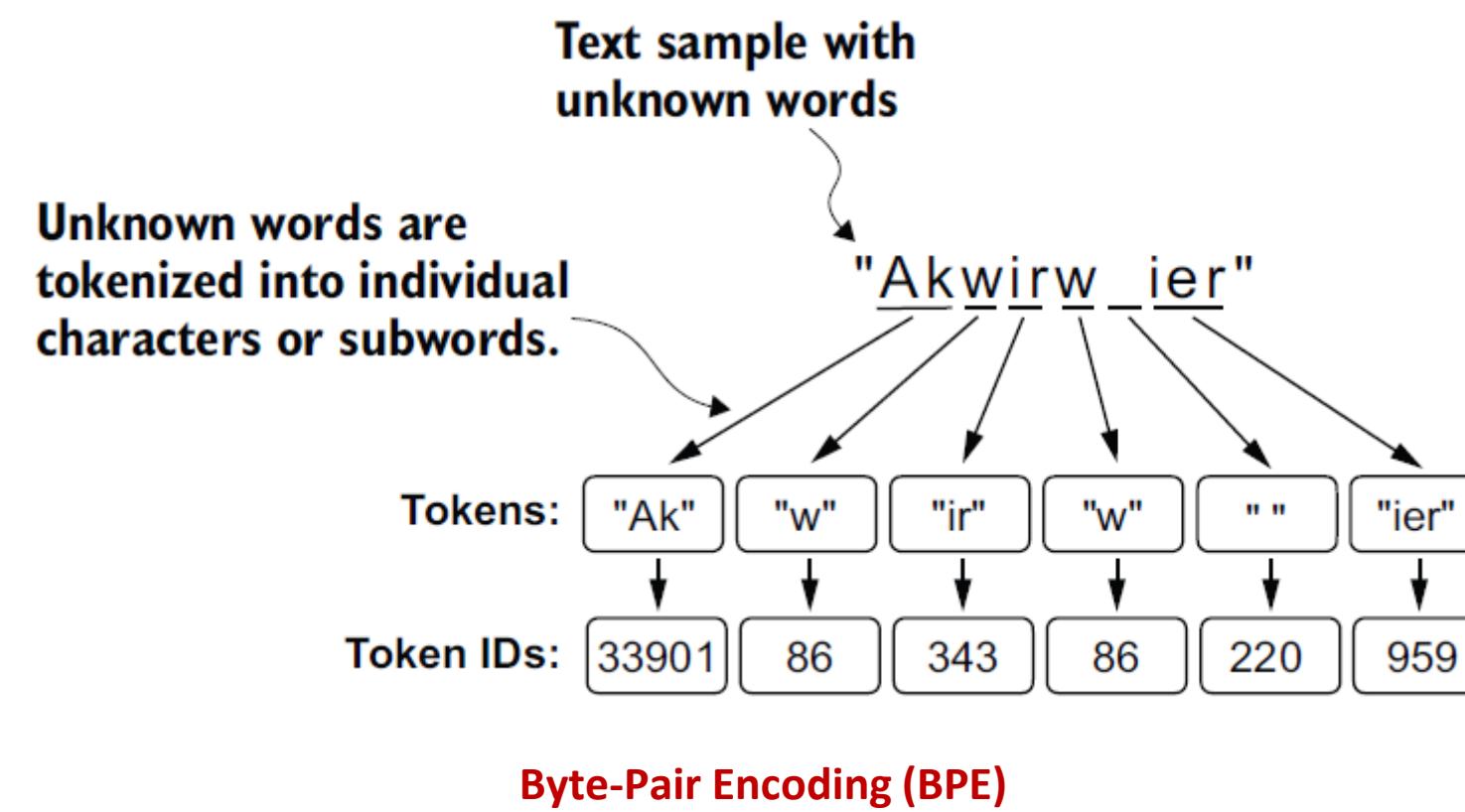
Byte-Pair Encoding



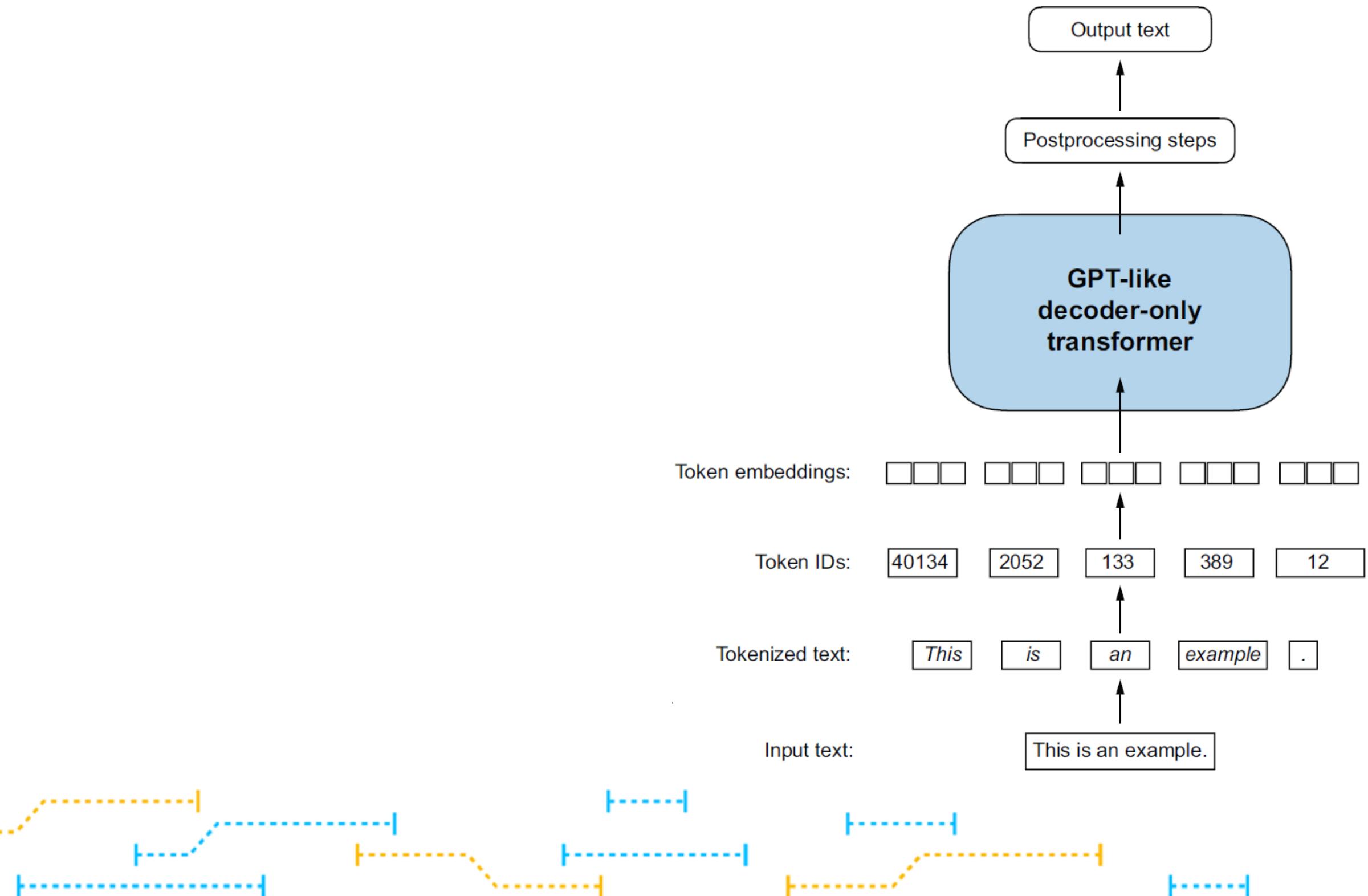
Byte-Pair Encoding



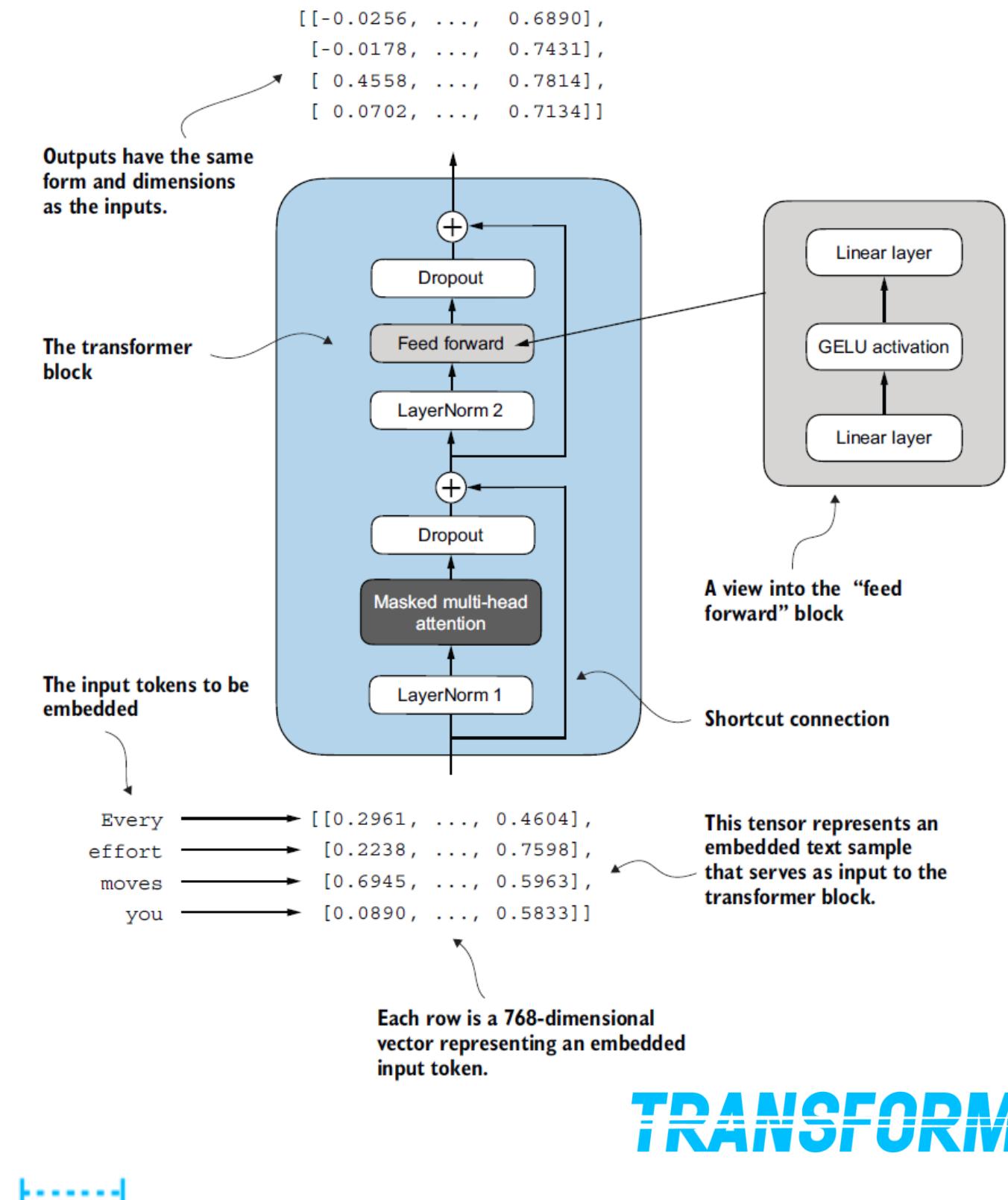
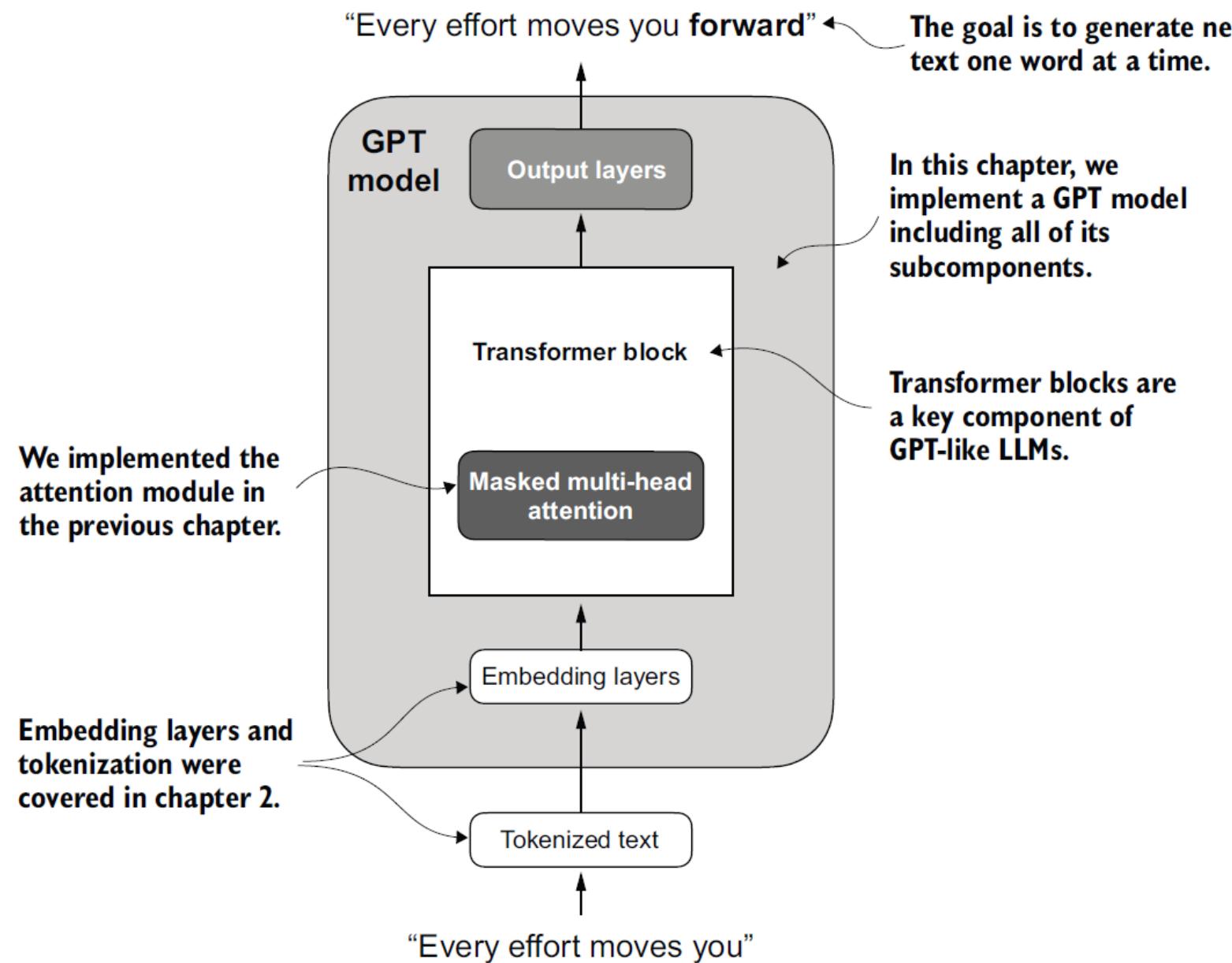
TikToken



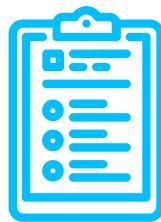
Natural language processing



GPT-model



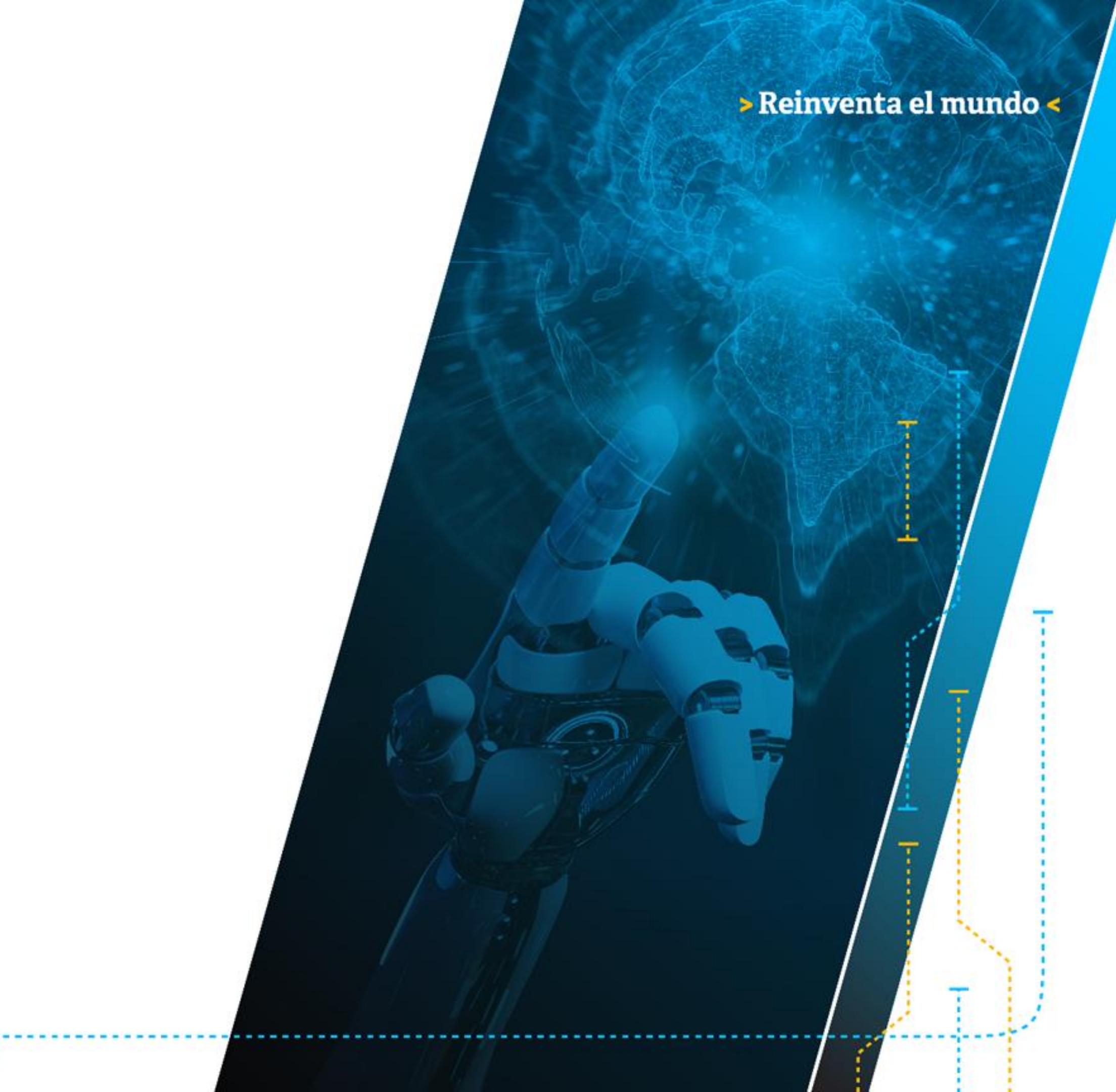
2.



Llama *Models*

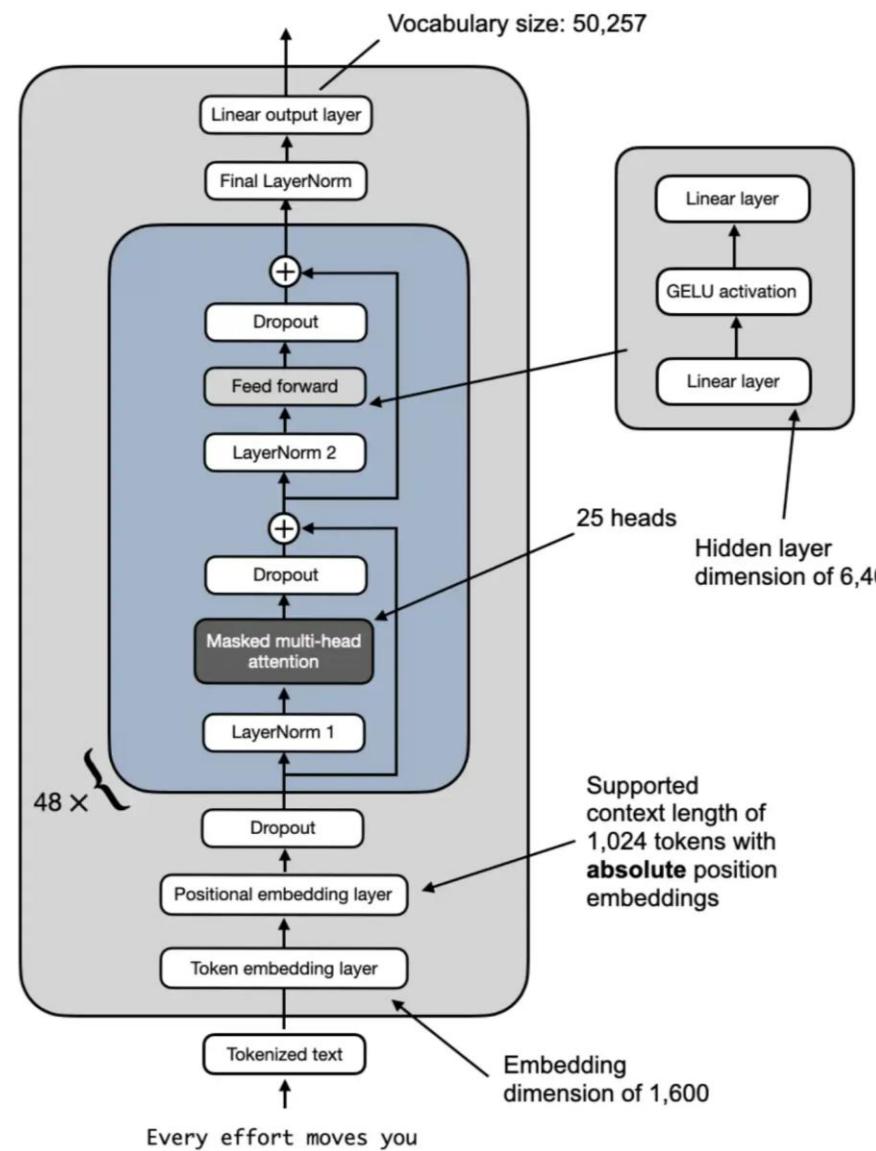
TRANSFORMATEC

> Reinventa el mundo <

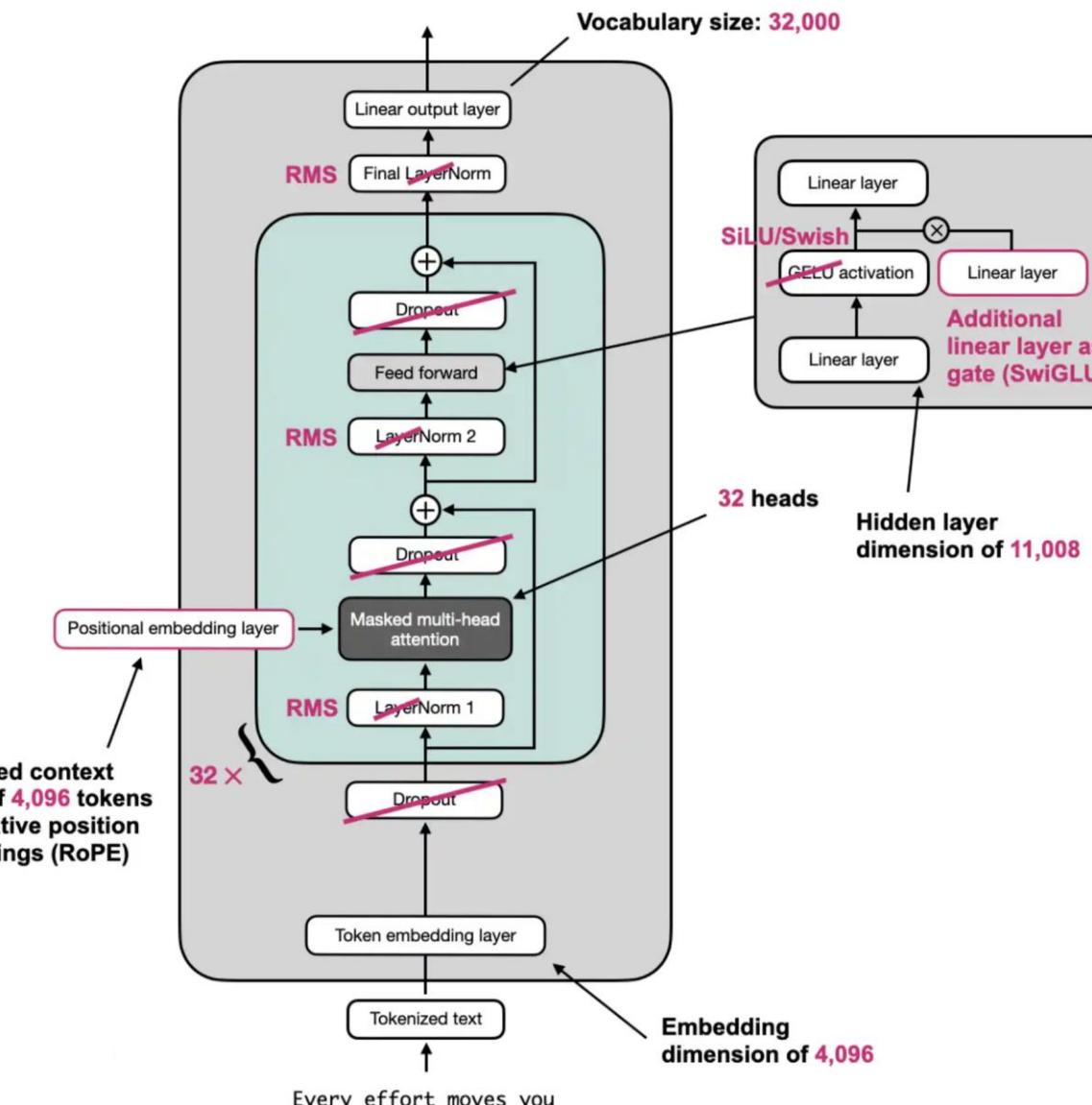


Llama 2 7B

GPT-2 XL 1.5B



Llama 2 7B



Llama 3

RMSNorm

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\mathbf{a})} g_i, \quad \text{where } \text{RMS}(\mathbf{a}) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}.$$

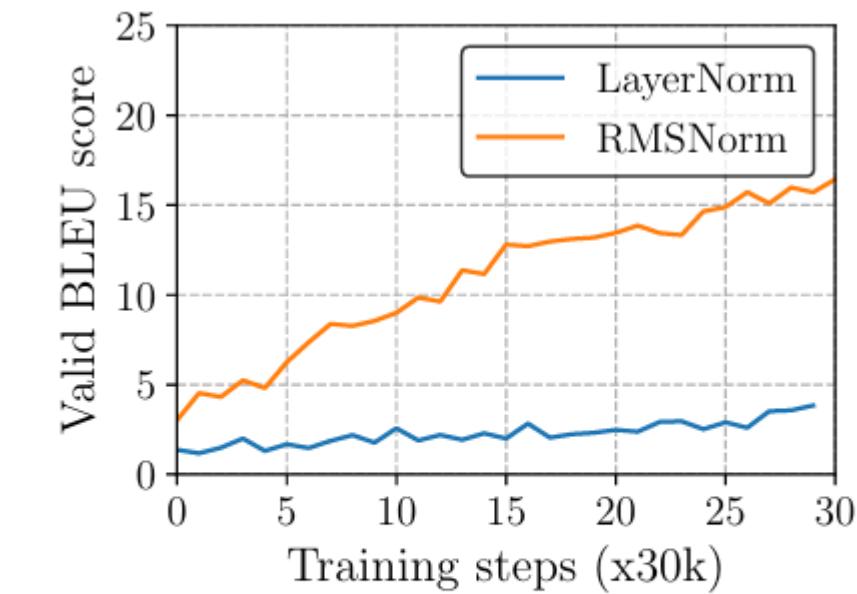
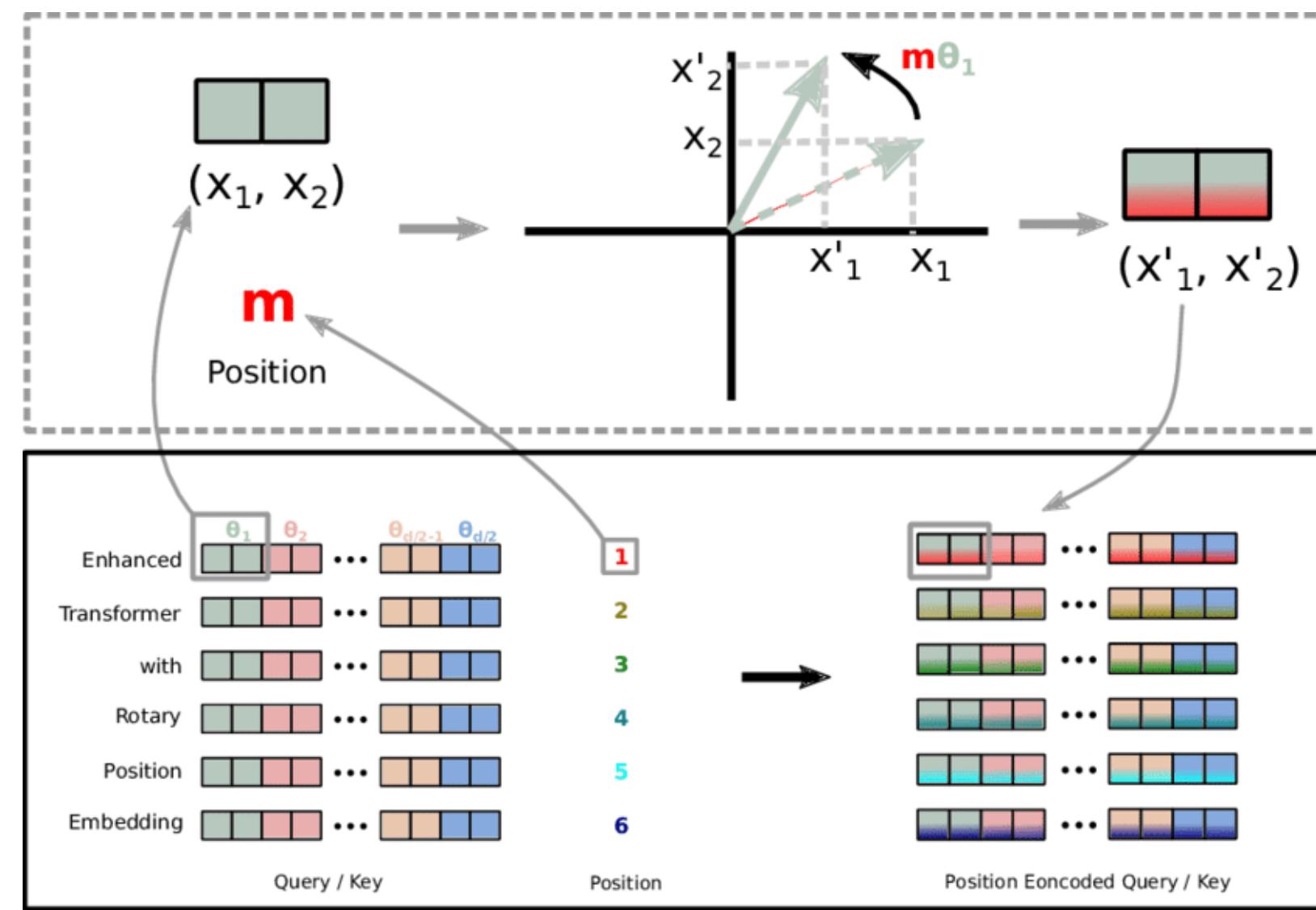


Figure 4: SacreBLEU score curve of LayerNorm and RMSNorm on newstest2013 (de-vset) when the initialization center is 0.2.



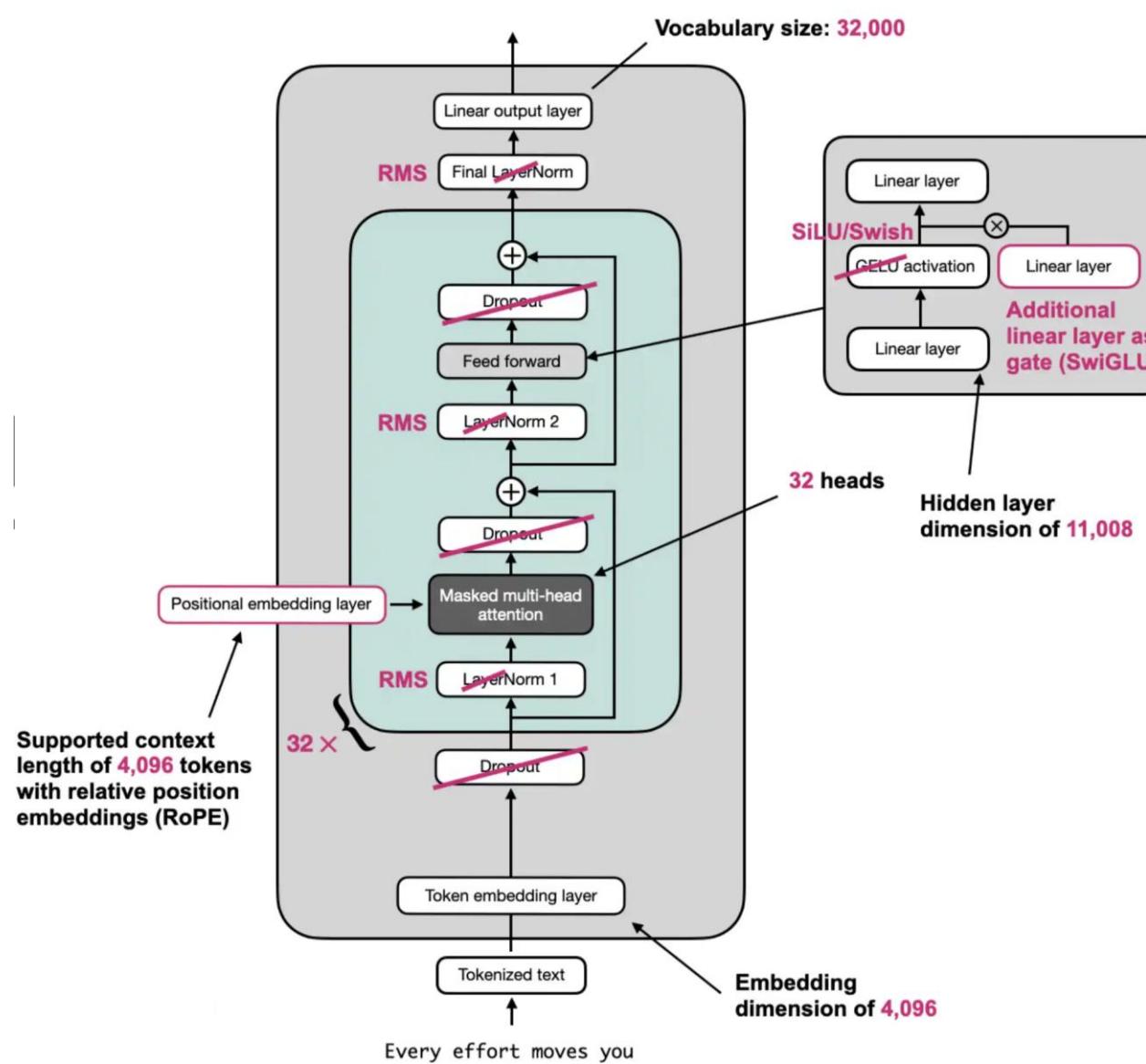
Llama 3

Rotary Positional Encoding (RoPE)

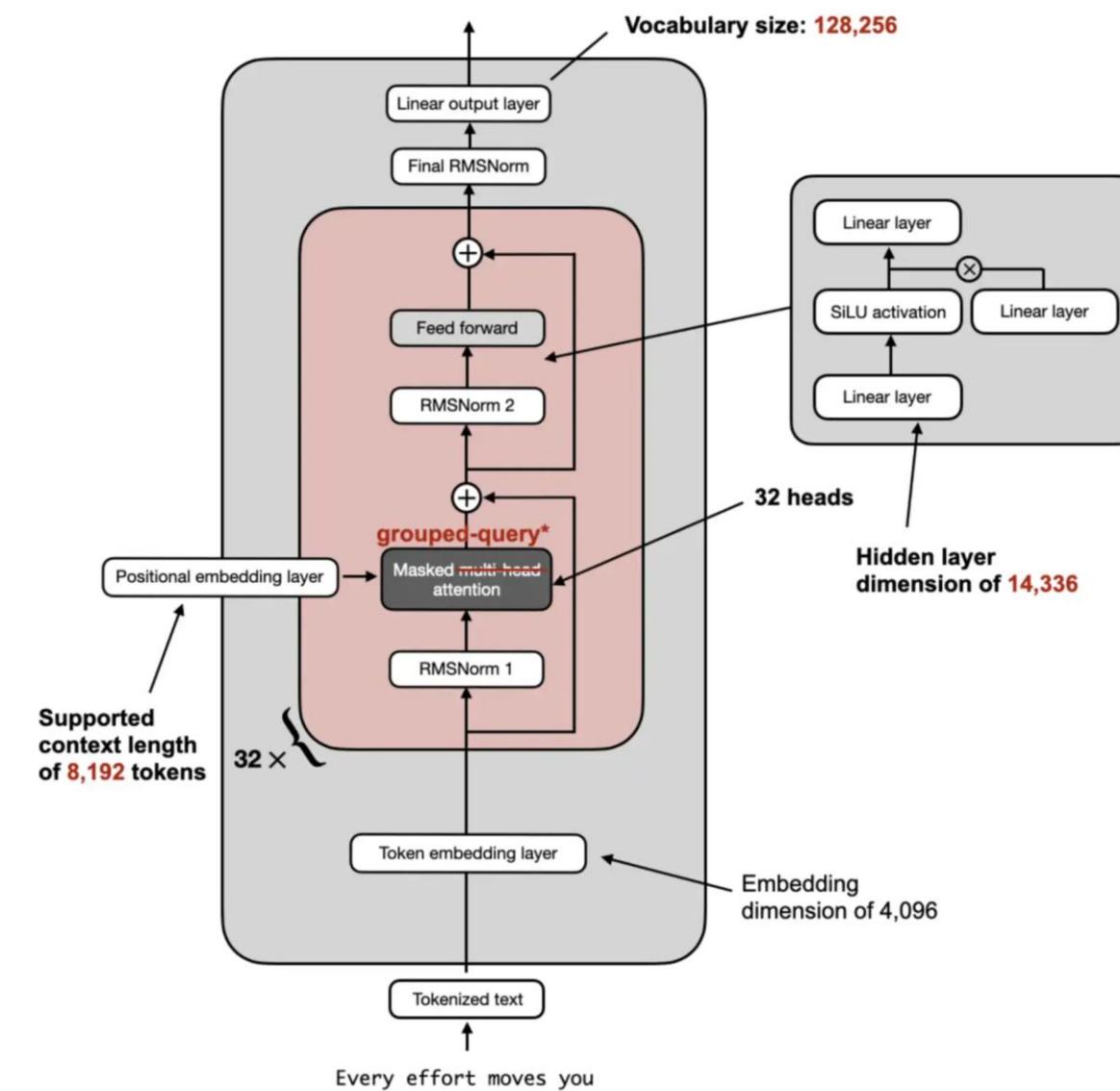


Llama 3 8B

Llama 2 7B



Llama 3 8B

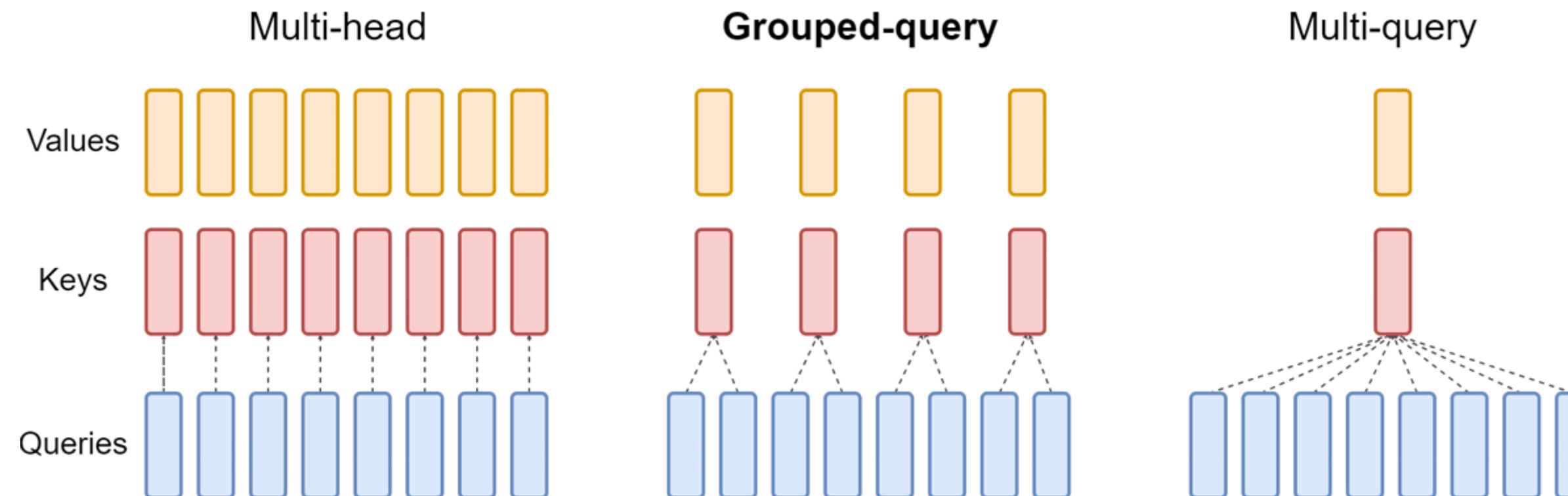


* The larger Llama 2 34B and 70B also used grouped-query attention

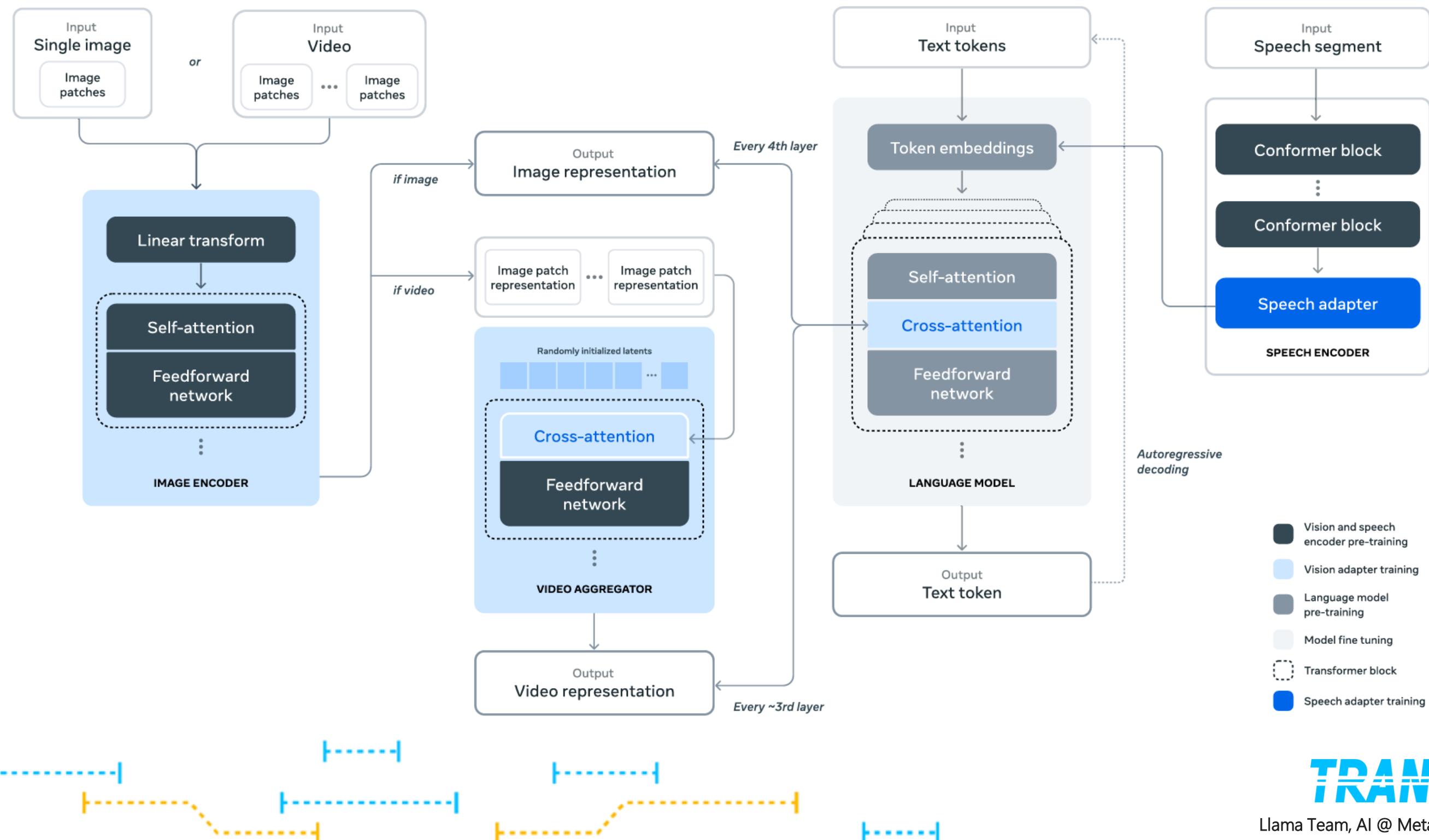


Llama 3

Grouped Query Attention (GQA)



Llama 3



TRANSFORMATEC

Llama Team, AI @ Meta (2024) "The Llama 3 Herd of Models".
arXiv preprint arXiv:2407.21783.

Llama 3

Preentrenamiento

Entrenamiento Inicial

- Se utiliza un tamaño de secuencia corto (4K tokens) para aprender la estructura del lenguaje y evitar el costo cuadrático en atención.
- Se comienza con un batch size bajo, incrementándolo gradualmente a medida que el modelo se estabiliza.

Entrenamiento de Contexto Largo

- Luego de entrenar en secuencias cortas, se entrena con contextos largos de hasta 128K tokens.
- Esta etapa se realiza de manera gradual para asegurar que el modelo se adapte sin perder rendimiento en tareas más cortas.
- Se utiliza Context Parallelism (CP) para dividir la secuencia entre GPUs, y se ajuste de RoPE para soportar un contexto de 500000.

Annealing Final

- Durante los últimos 40M tokens, se reduce la tasa de aprendizaje de manera lineal a cero.
- Se hace un ajuste fino de la mezcla de datos, aumentando la frecuencia se sampleo a datos de alta calidad (por ejemplo, matemáticos y código) para refinar el desempeño del modelo.
- Se realiza un promedio de los checkpoints para estabilizar la versión final del modelo.



Llama 3

Postentrenamiento

Supervised Finetuning (SFT)

- Se entrena el modelo usando un conjunto de datos etiquetado (instrucciones y respuestas).
- Se realiza sobre una mezcla de datos humanos y generados, ajustando el modelo para mejorar su capacidad de seguir instrucciones.

Direct Preference Optimization (DPO)

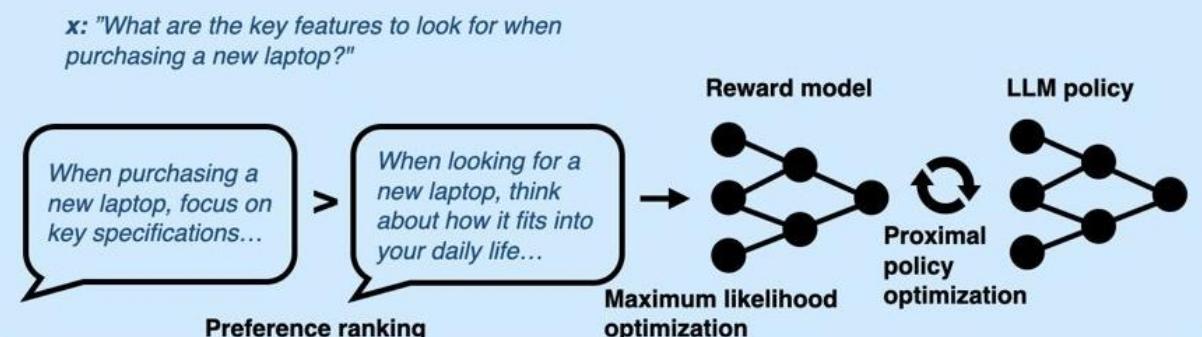
- Posteriormente, se alinea el modelo con las preferencias humanas directamente usando técnicas contrastivas.
- Se incluyen modificaciones como enmascarar tokens de formato para evitar comportamientos indeseados, y se incorpora un término de Negative Log-Likelihood (NLL) como regularización.

Model Averaging

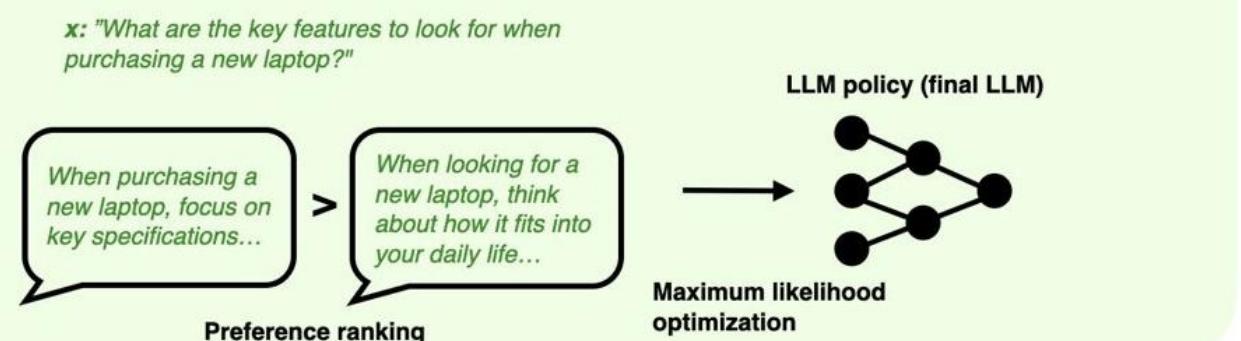
- Se realizan múltiples rondas de entrenamiento, luego promediamos checkpoints para obtener una versión robusta y estable.



Reinforcement Learning with Human Feedback (RLHF)

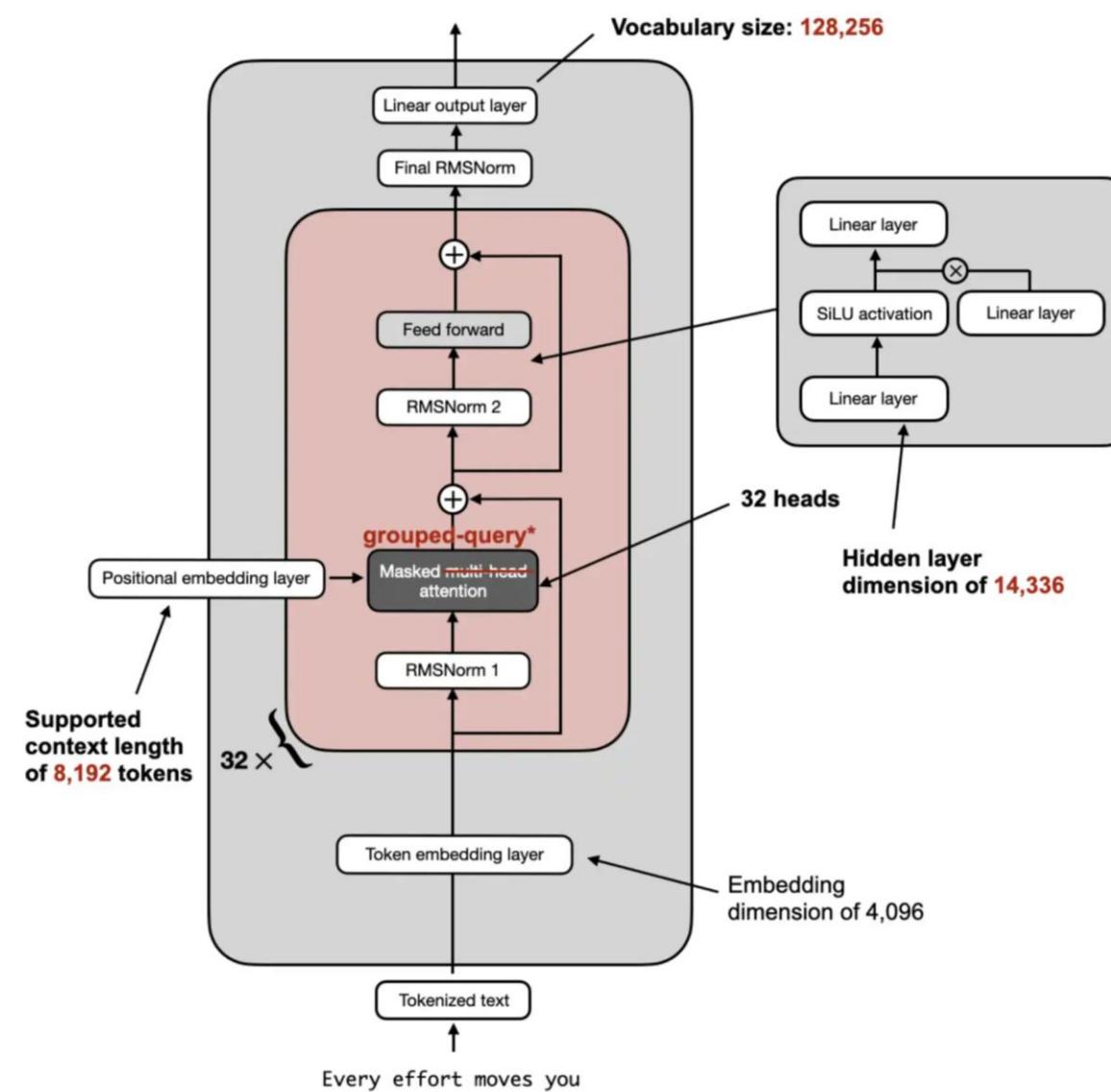


Direct Preference Optimization (DPO)



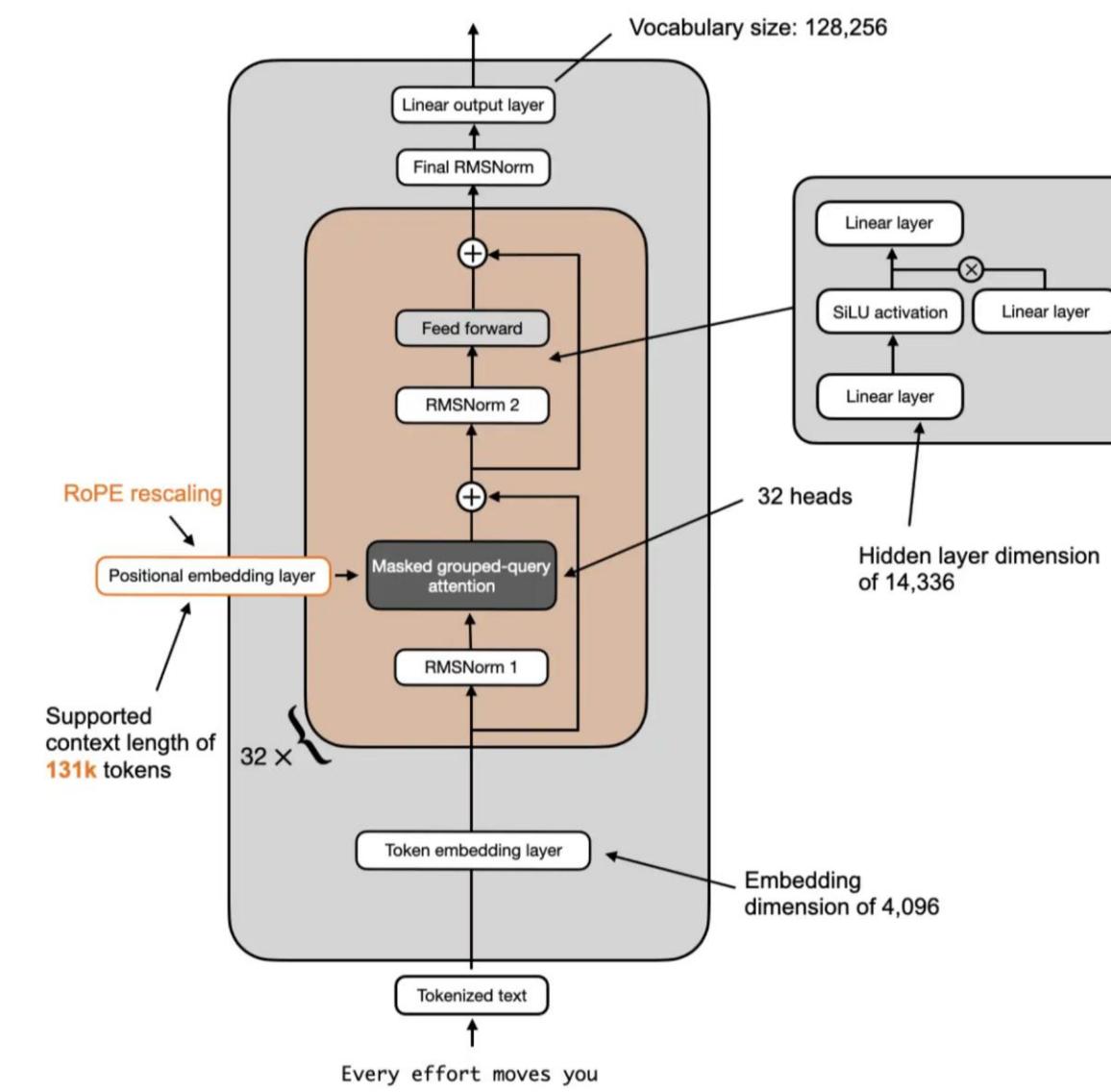
Llama 3.1 8B

Llama 3 8B



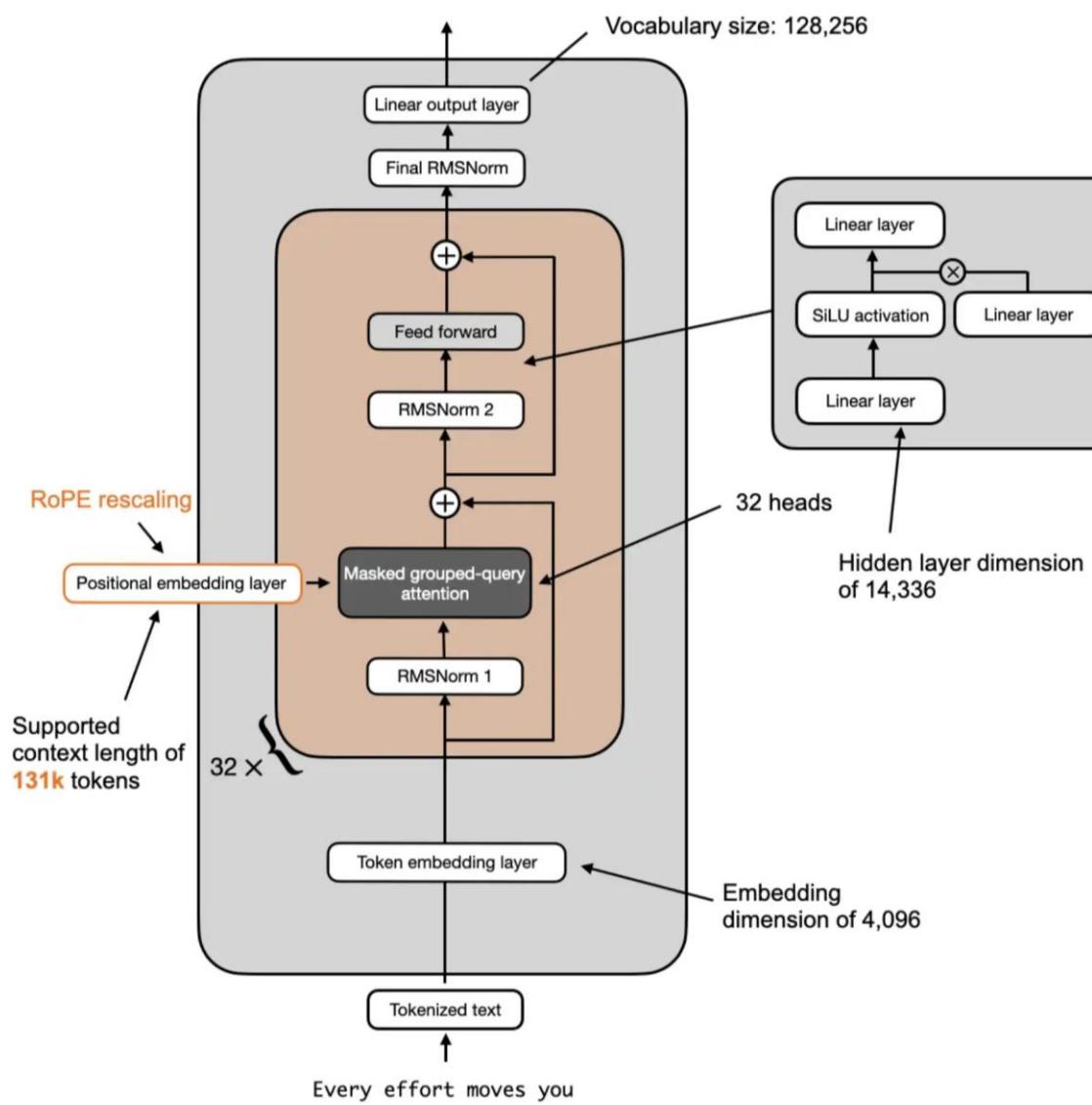
* The larger Llama 2 34B and 70B also used grouped-query attention

Llama 3.1 8B

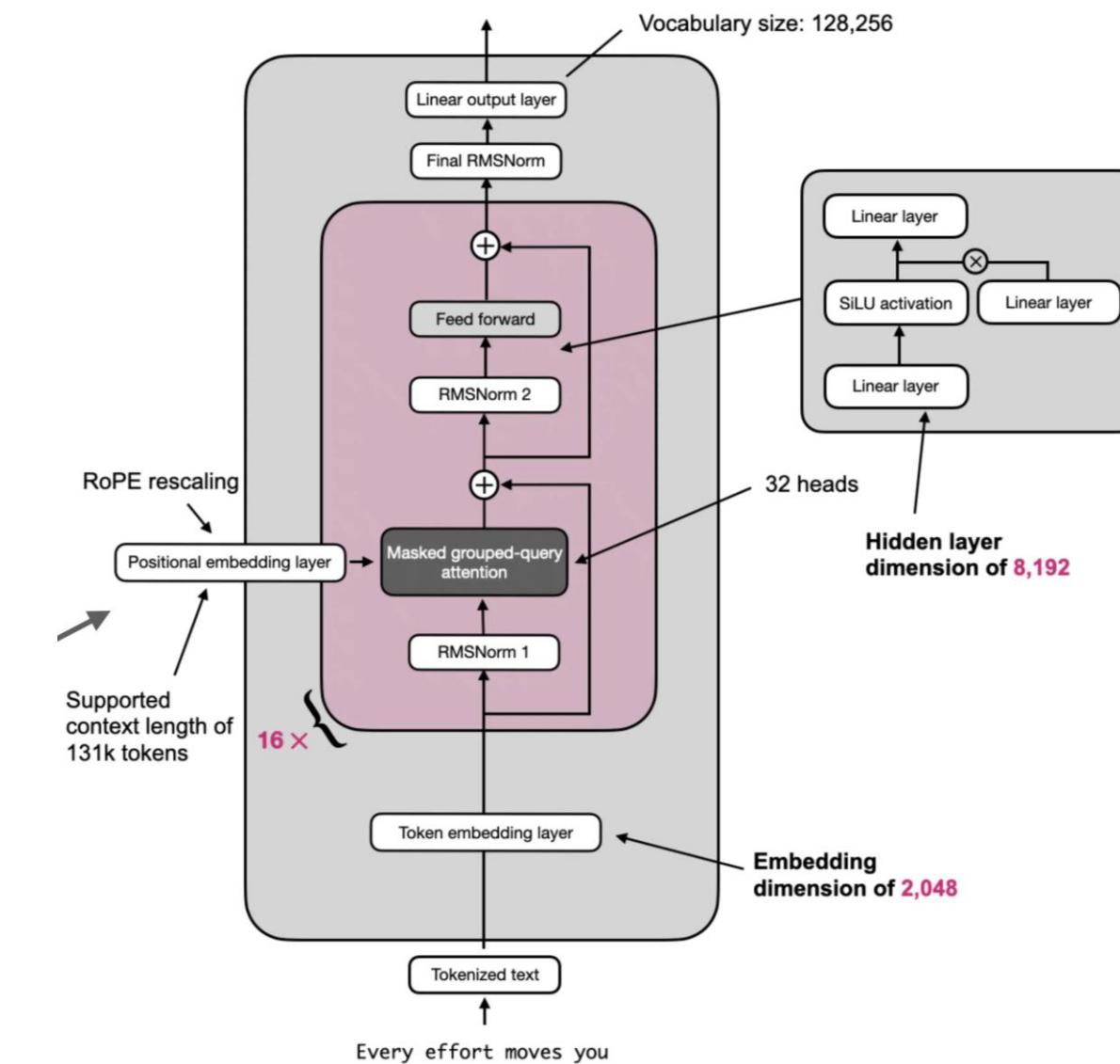


Llama 3.2 1B

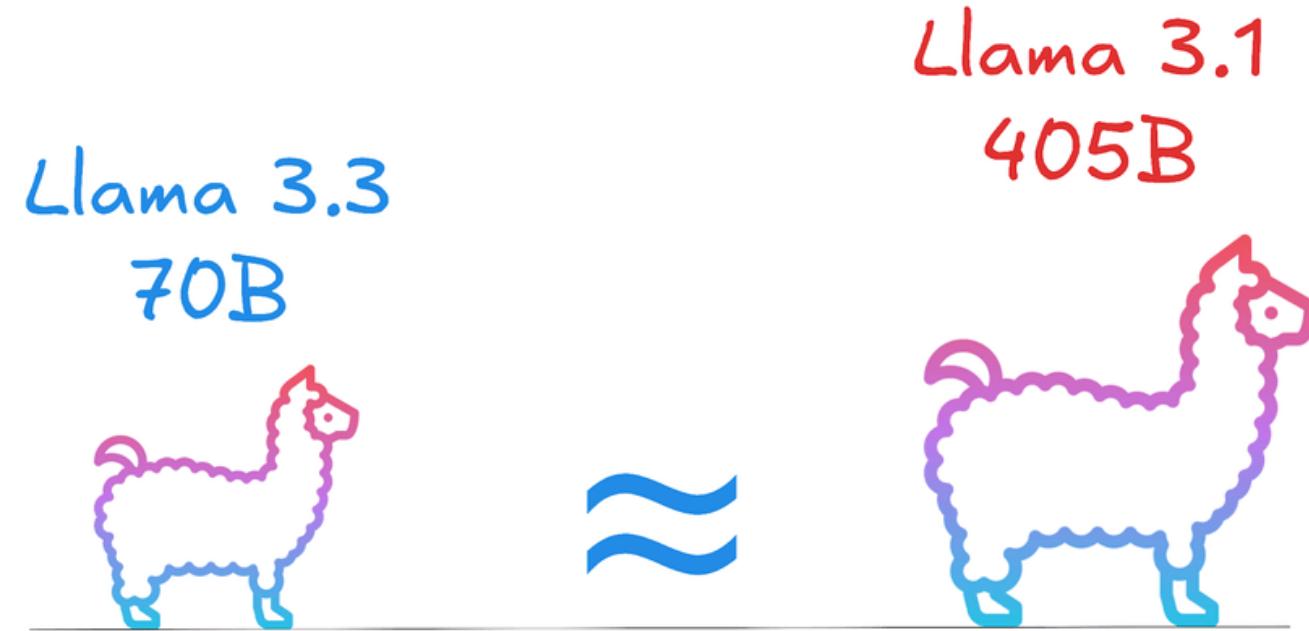
Llama 3.1 8B



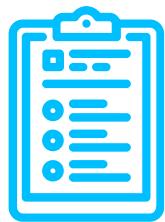
Llama 3.2 1B



Llama 3.3



3.



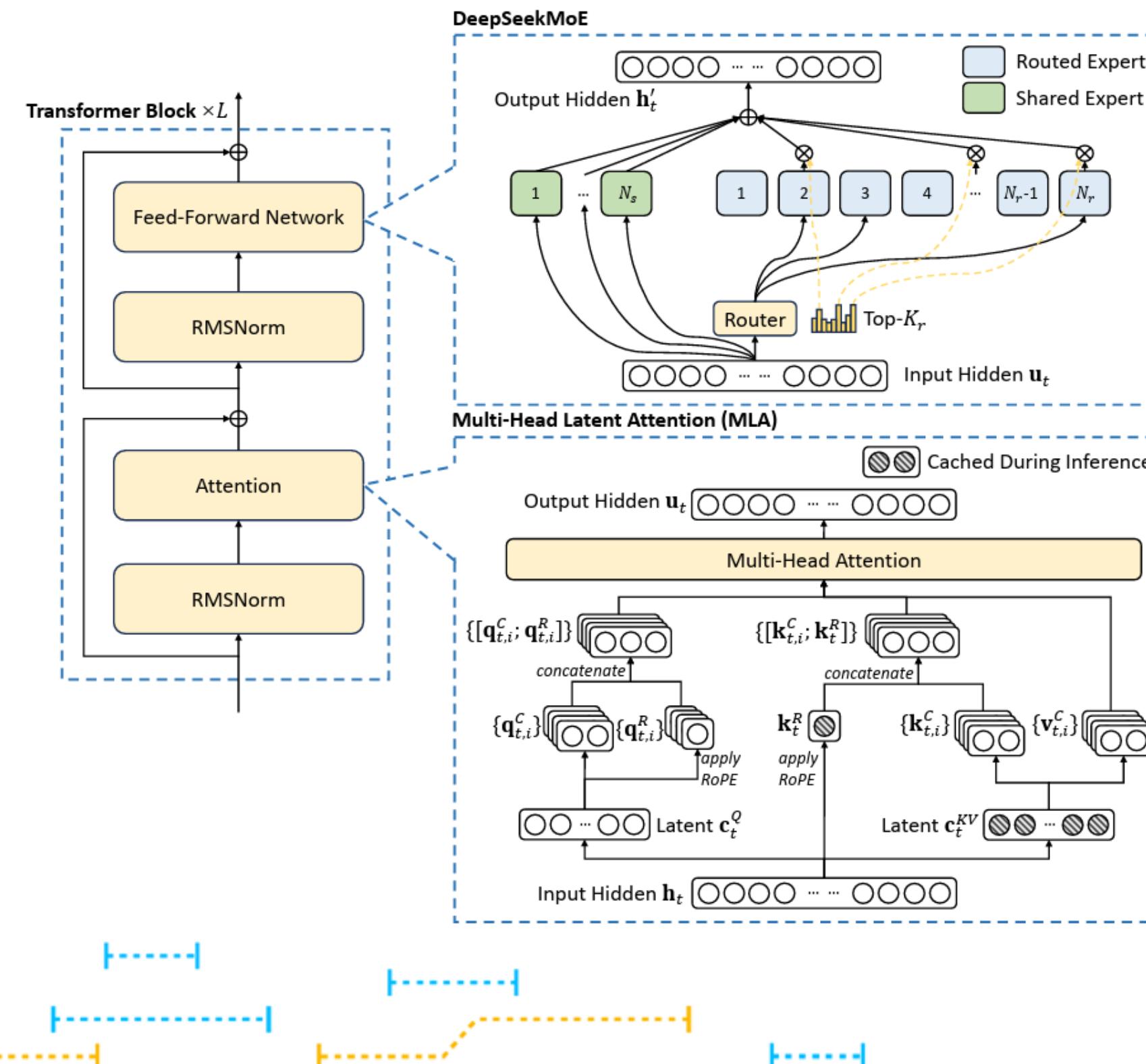
DeepSeek

TRANSFORMATEC

> Reinventa el mundo <

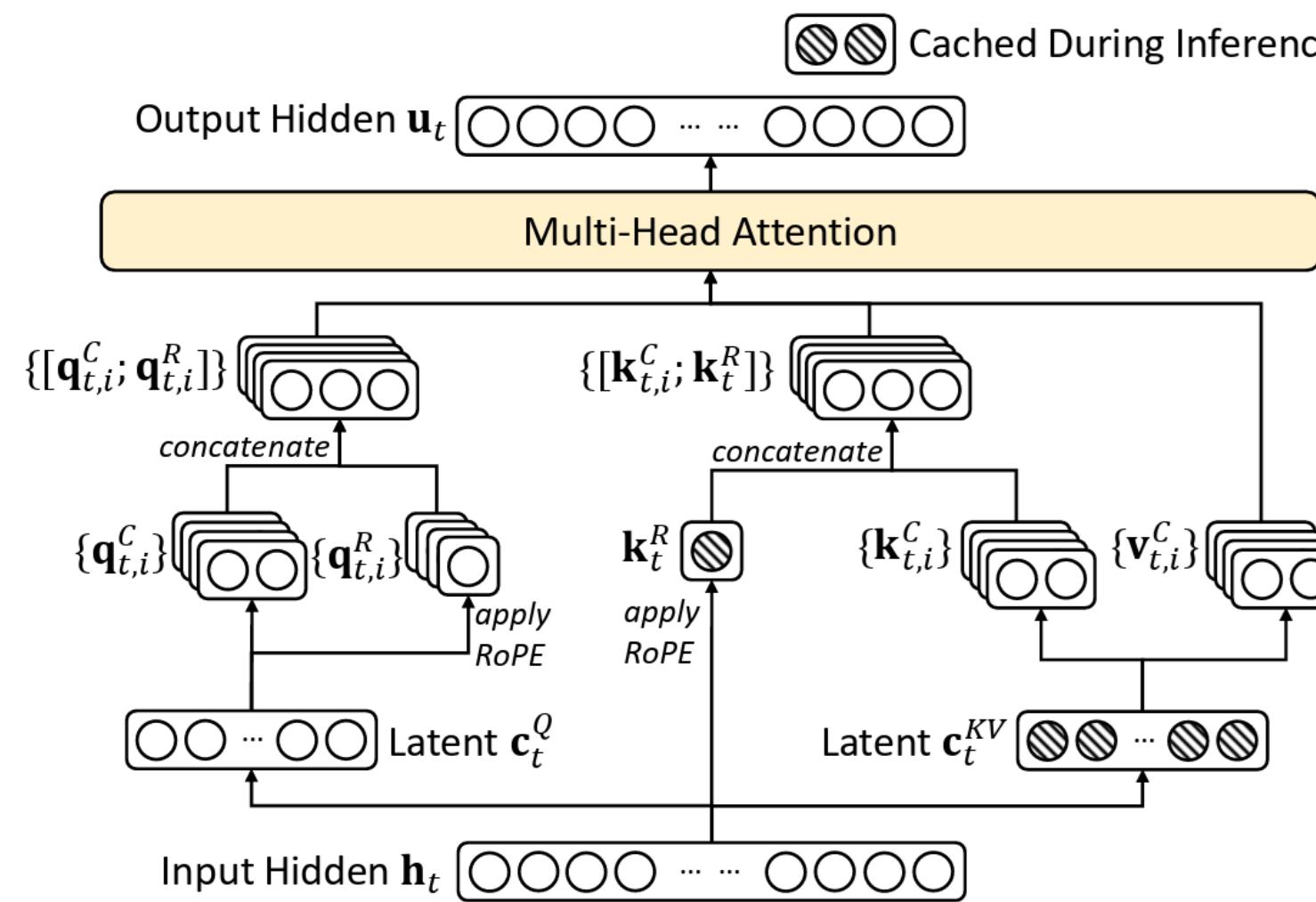


DeepSeek-V3



DeepSeek-V3

Multi-Head Latent Attention (MLA)



Key

$$\begin{aligned}\mathbf{c}_t^{KV} &= W^{DKV} \mathbf{h}_t, \\ [\mathbf{k}_{t,1}^C; \mathbf{k}_{t,2}^C; \dots; \mathbf{k}_{t,n_h}^C] &= \mathbf{k}_t^C = W^{UK} \mathbf{c}_t^{KV}, \\ \mathbf{k}_t^R &= \text{RoPE}(W^{KR} \mathbf{h}_t), \\ \mathbf{k}_{t,i} &= [\mathbf{k}_{t,i}^C; \mathbf{k}_t^R], \\ [\mathbf{v}_{t,1}^C; \mathbf{v}_{t,2}^C; \dots; \mathbf{v}_{t,n_h}^C] &= \mathbf{v}_t^C = W^{UV} \mathbf{c}_t^{KV},\end{aligned}$$

Query

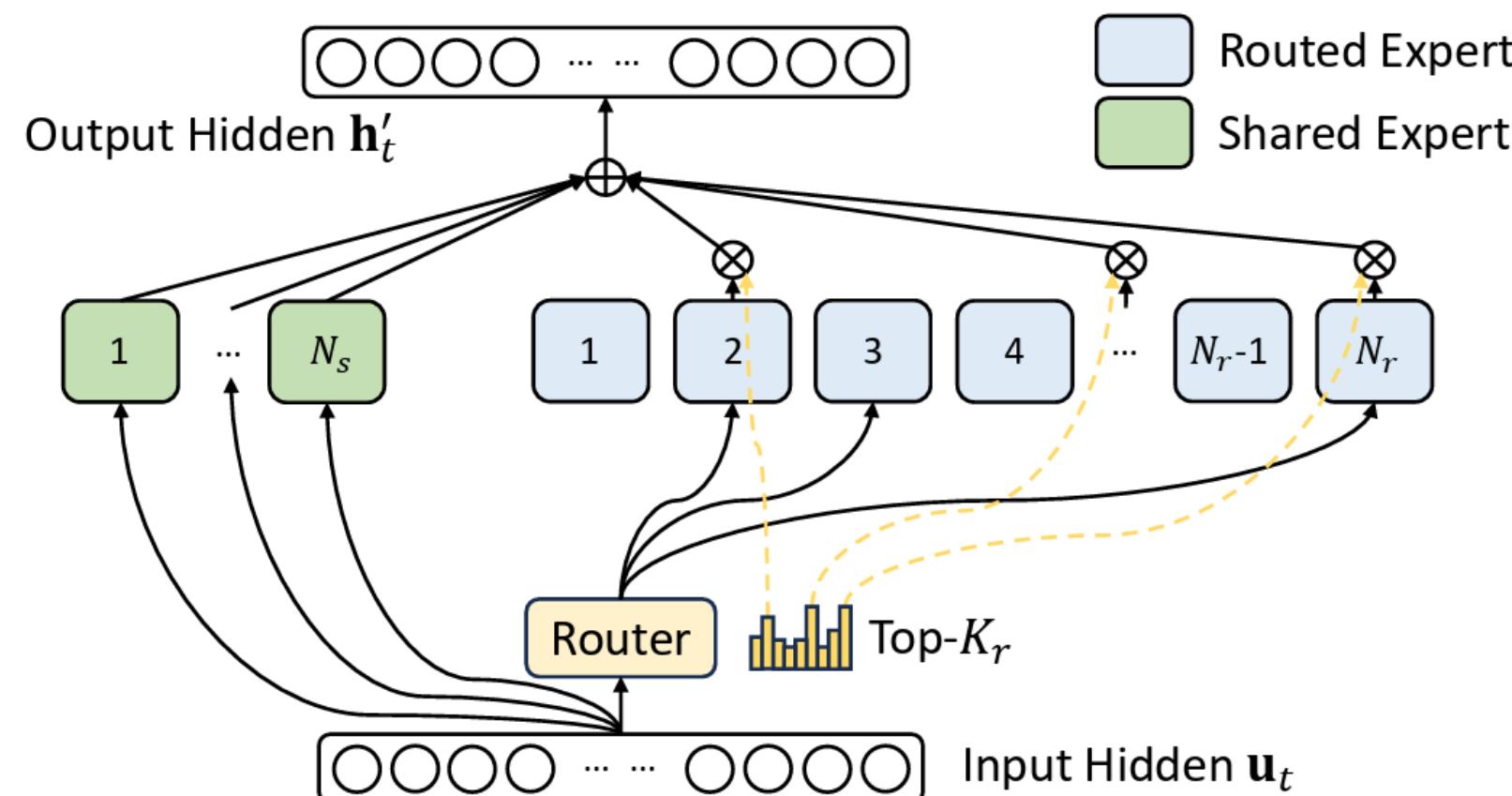
$$\begin{aligned}\mathbf{c}_t^Q &= W^{DQ} \mathbf{h}_t, \\ [\mathbf{q}_{t,1}^C; \mathbf{q}_{t,2}^C; \dots; \mathbf{q}_{t,n_h}^C] &= \mathbf{q}_t^C = W^{UQ} \mathbf{c}_t^Q, \\ [\mathbf{q}_{t,1}^R; \mathbf{q}_{t,2}^R; \dots; \mathbf{q}_{t,n_h}^R] &= \mathbf{q}_t^R = \text{RoPE}(W^{QR} \mathbf{c}_t^Q), \\ \mathbf{q}_{t,i} &= [\mathbf{q}_{t,i}^C; \mathbf{q}_{t,i}^R],\end{aligned}$$

Attention

$$\begin{aligned}\mathbf{o}_{t,i} &= \sum_{j=1}^t \text{Softmax}_j \left(\frac{\mathbf{q}_{t,i}^T \mathbf{k}_{j,i}}{\sqrt{d_h + d_h^R}} \right) \mathbf{v}_{j,i}^C, \\ \mathbf{u}_t &= W^O [\mathbf{o}_{t,1}; \mathbf{o}_{t,2}; \dots; \mathbf{o}_{t,n_h}],\end{aligned}$$

DeepSeek-V3

DeepSeekMoE



$$\mathbf{h}'_t = \mathbf{u}_t + \sum_{i=1}^{N_s} \text{FFN}_i^{(s)} (\mathbf{u}_t) + \sum_{i=1}^{N_r} g_{i,t} \text{FFN}_i^{(r)} (\mathbf{u}_t),$$

$$g_{i,t} = \frac{g'_{i,t}}{\sum_{j=1}^{N_r} g'_{j,t}},$$

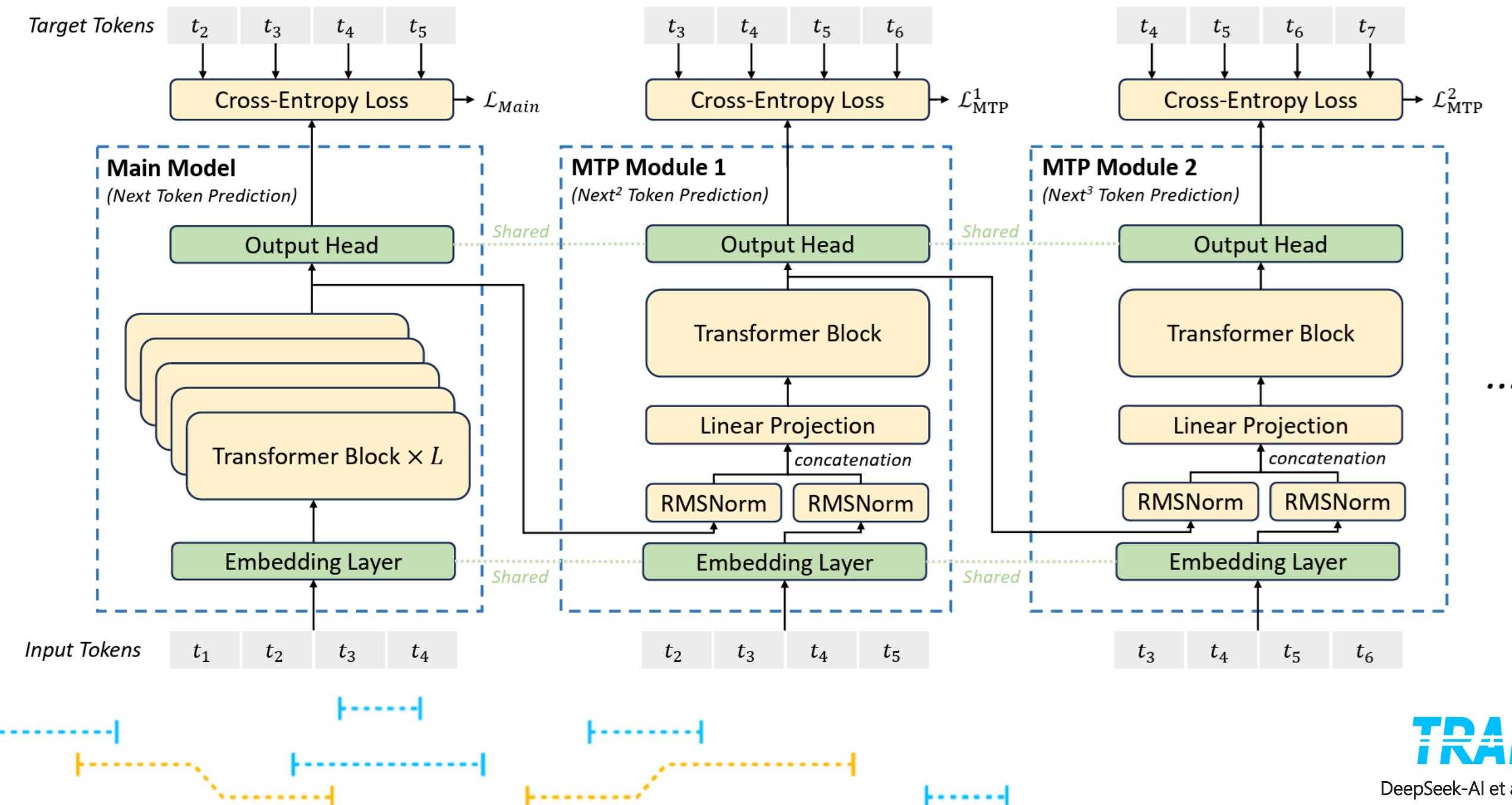
$$g'_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} \in \text{Topk}(\{s_{j,t} | 1 \leq j \leq N_r\}, K_r), \\ 0, & \text{otherwise,} \end{cases}$$

$$s_{i,t} = \text{Sigmoid} (\mathbf{u}_t^T \mathbf{e}_i),$$



DeepSeek-V3

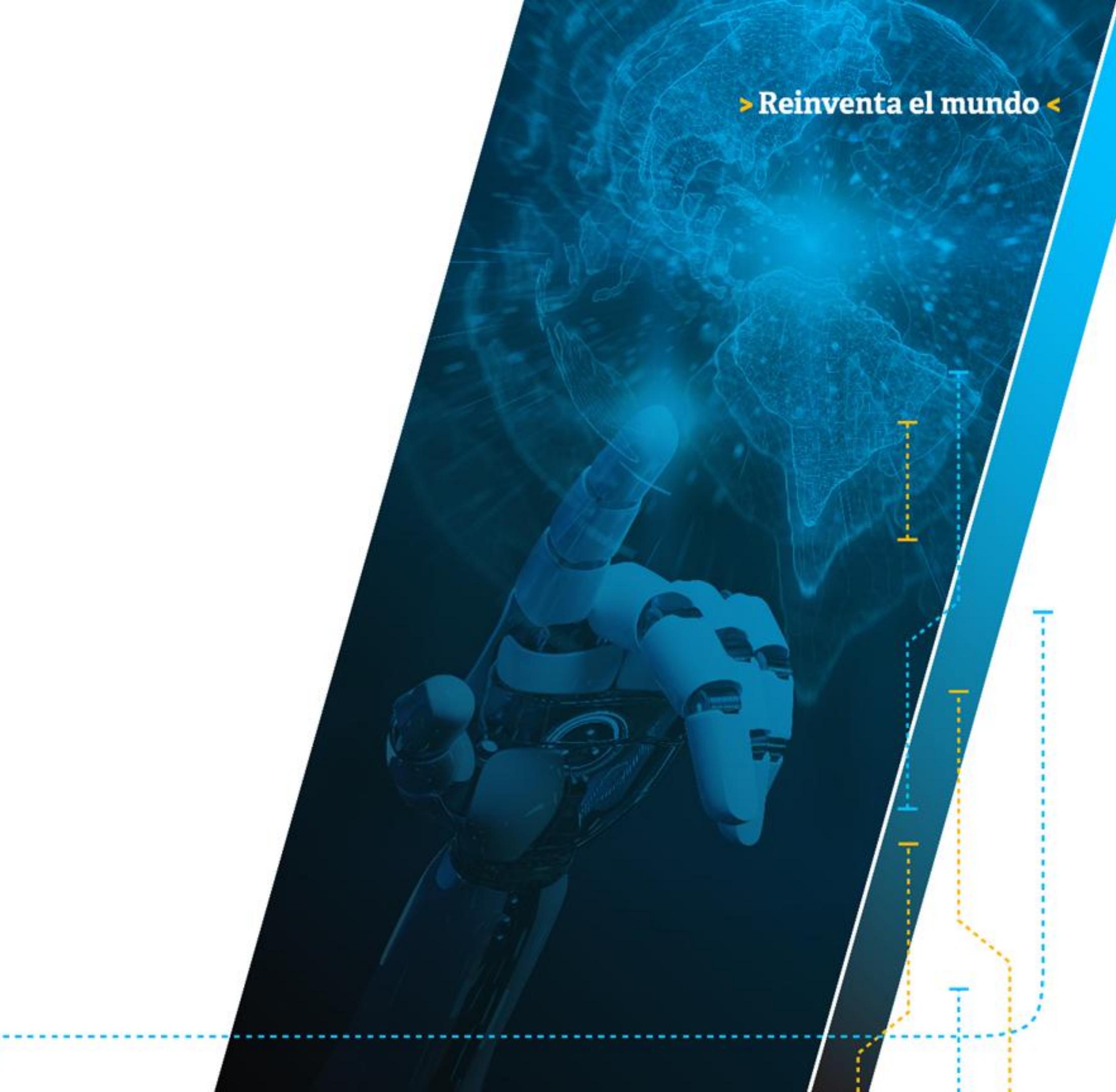
Multi-Token Prediction (MTP)



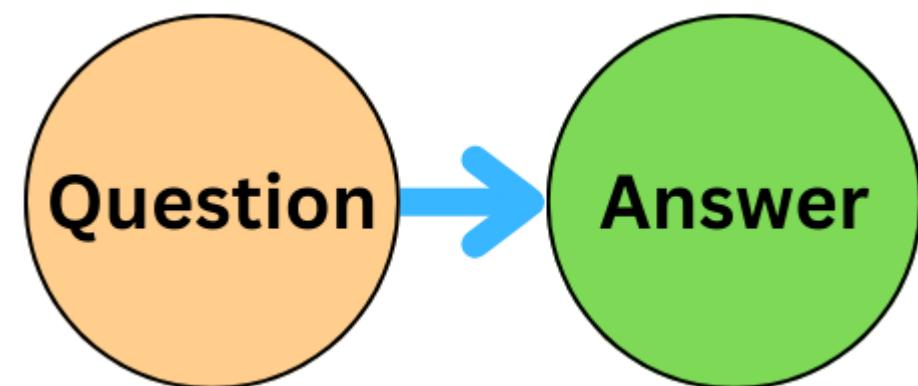
3.



Prom*pting*



Zero-Shot *Prompting*



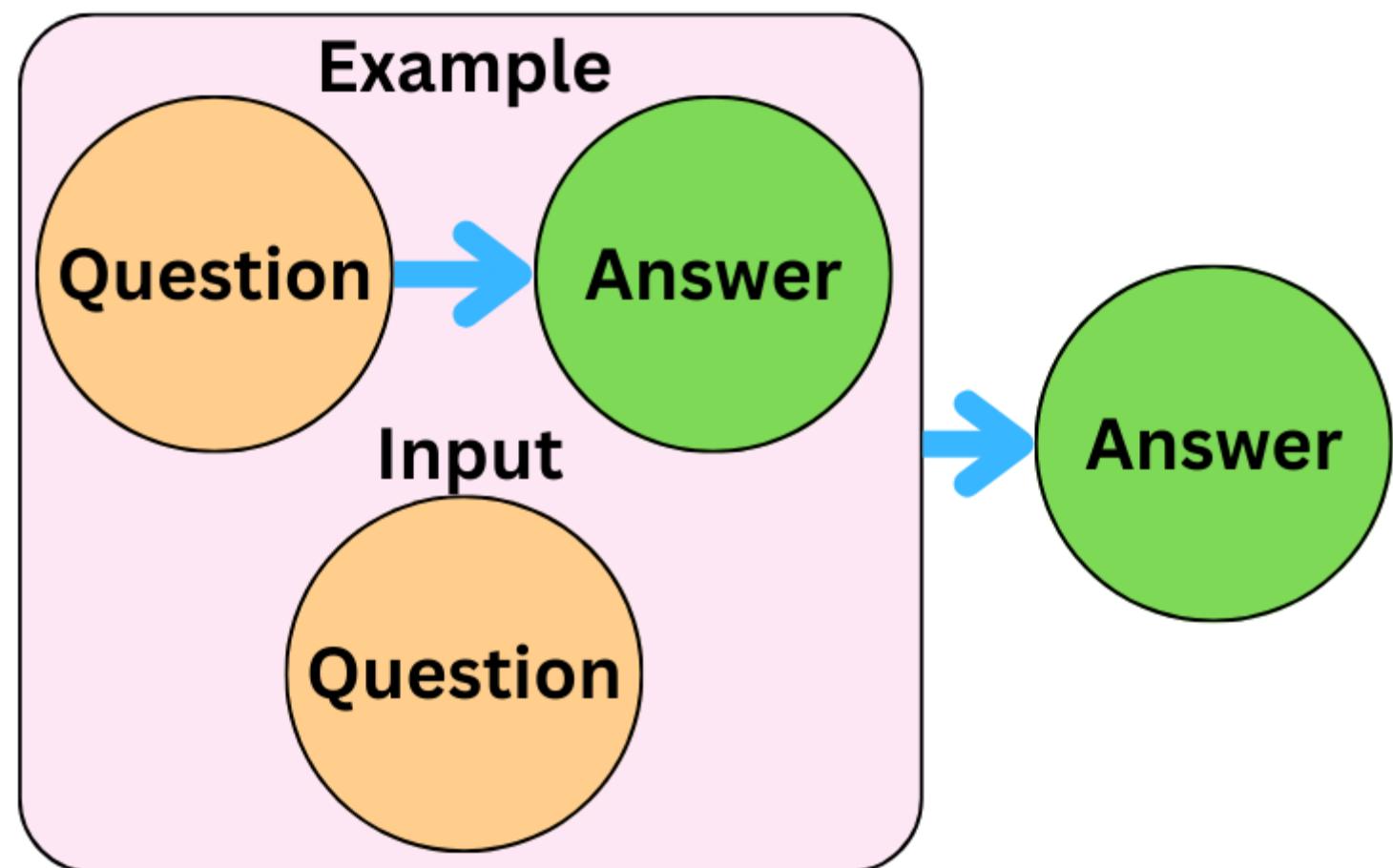
Write a Python function that takes a list of numbers and returns the average but exclude the highest and lowest numbers from the calculation.

Translate this sentence to French: '{input}'

Summarize the following paragraph in one sentence:
'{input}'



Few-Shot *Prompting*



Here is how to solve analogies: 'Man is to Woman as King is to Queen.' -> Queen. 'Sun is to Day as Moon is to Night.' -> Night. Now solve this analogy: '{input}'

Here is how to categorize words semantically: 'Apple, Banana, Cherry' -> Fruits. 'Dog, Cat, Horse' -> Animals. Now categorize these words: '{input}'

Here is how to write a project plan:
'Project: Website Redesign
Objective: Revamp the company's website to improve user experience and SEO rankings.
Timeline: 3 months

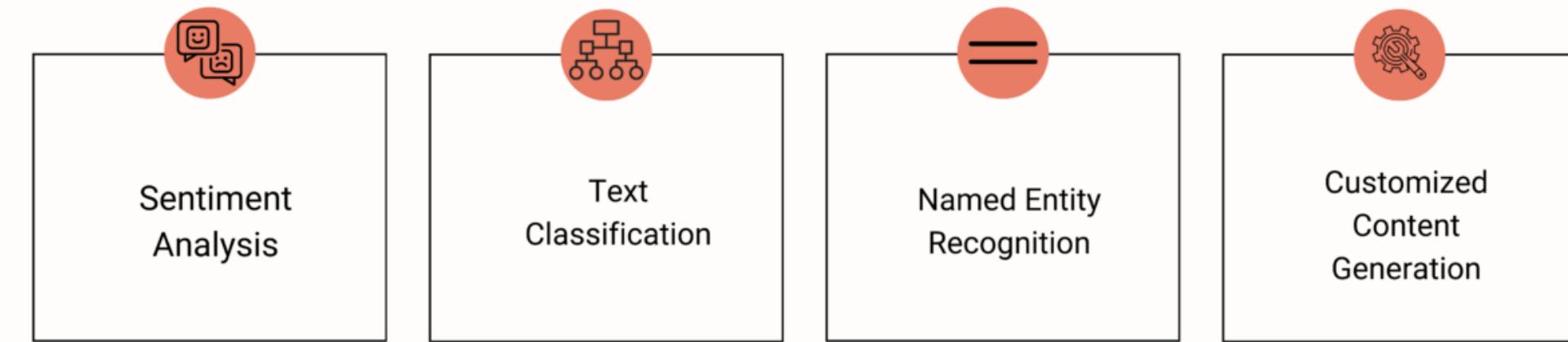
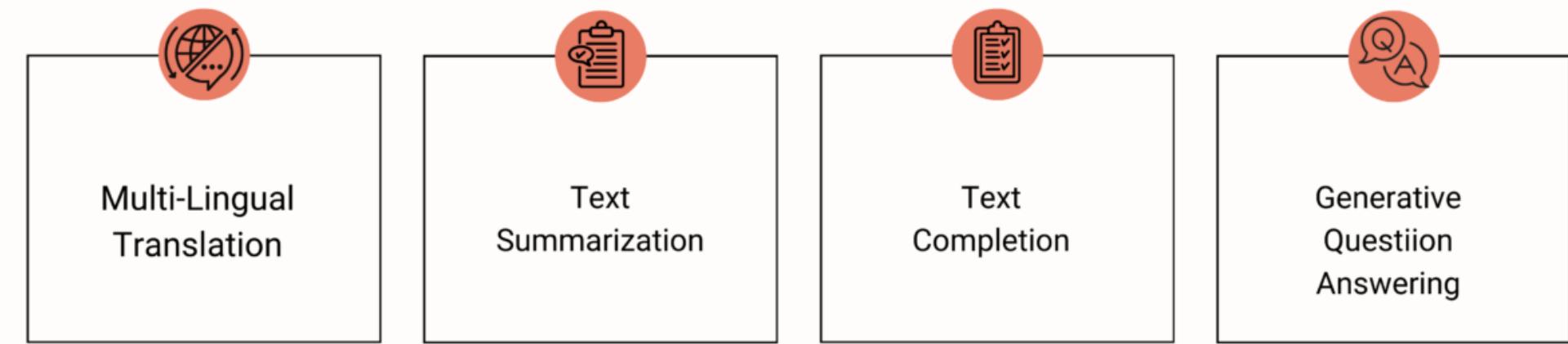
Key Milestones:

1. Wireframe Design (Month 1)
2. Development Phase (Month 2)
3. Testing and Launch (Month 3)

Now, create a project plan based on this goal: '{input}'

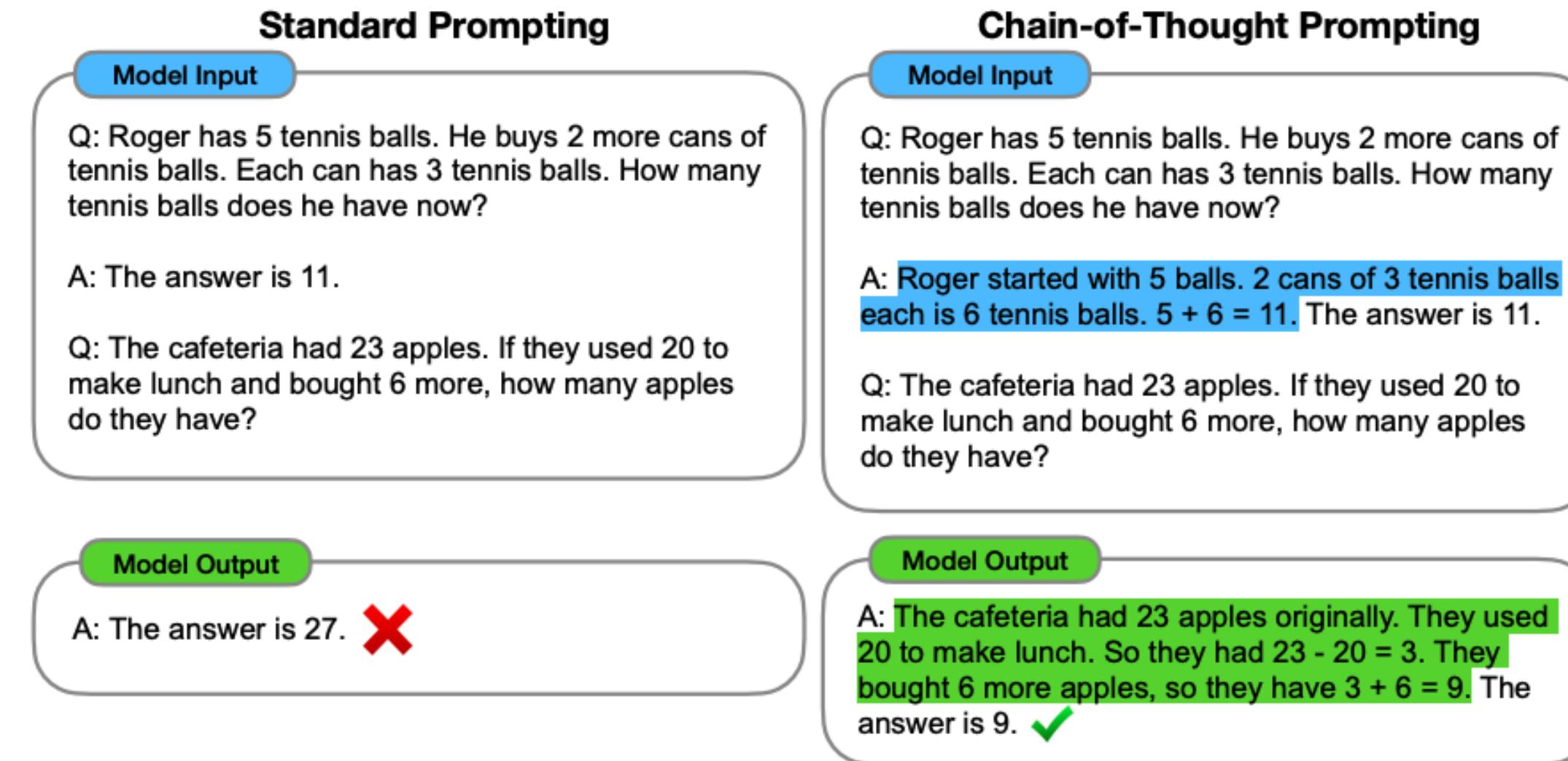


Prompting



Chain-of-Thought (CoT) Prompting

Few-shot-CoT



Chain-of-Thought (CoT) Prompting

Zero-shot-CoT

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?

A: Let's think step by step.

Rationale Generation

LLM

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?

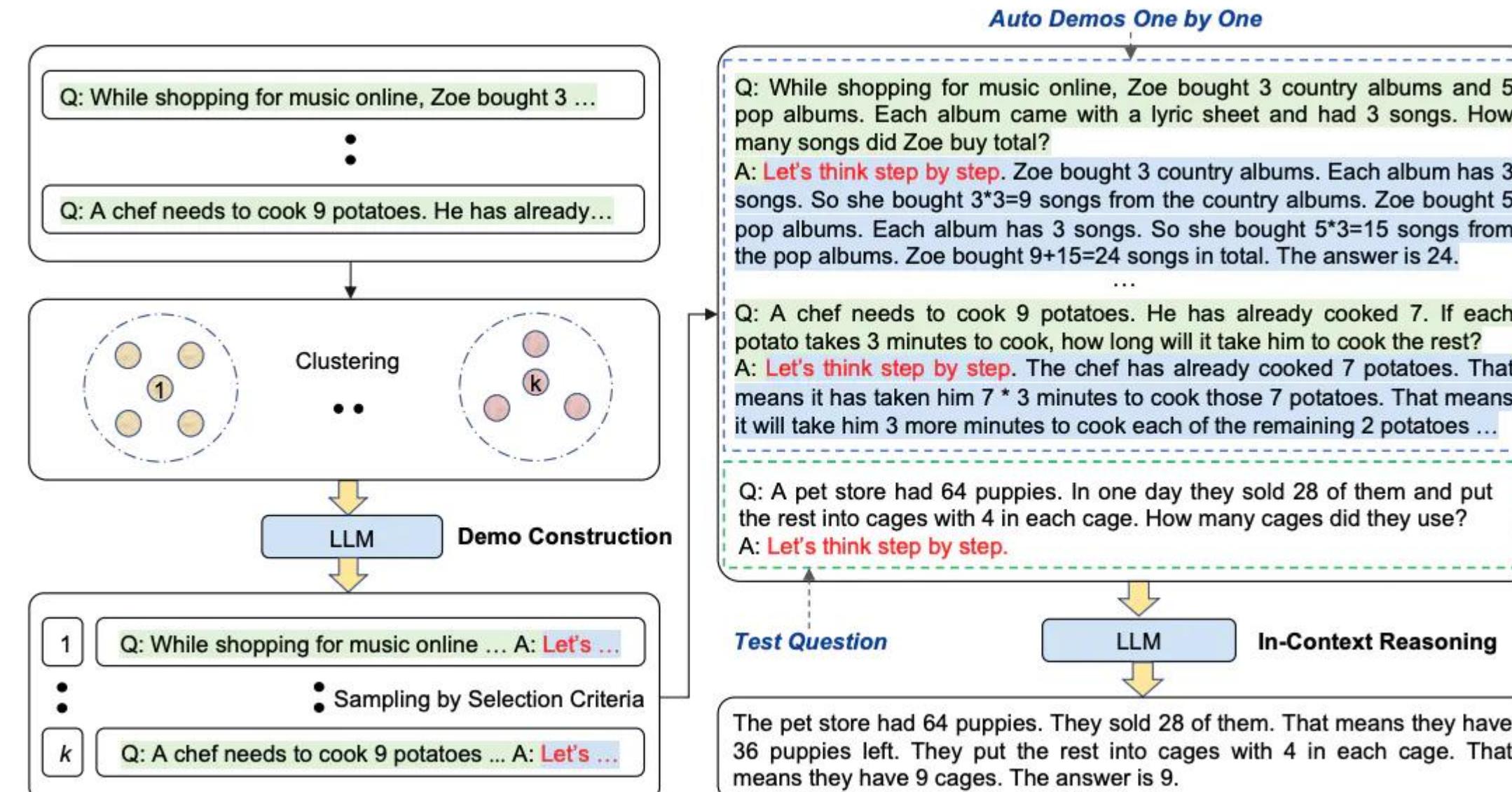
A: Let's think step by step. There are 64 puppies. 28 of them were sold. This leaves 36 puppies. Each cage has 4 puppies, so we need 9 cages. Therefore, the answer (arabic numerals) is

Answer Extraction

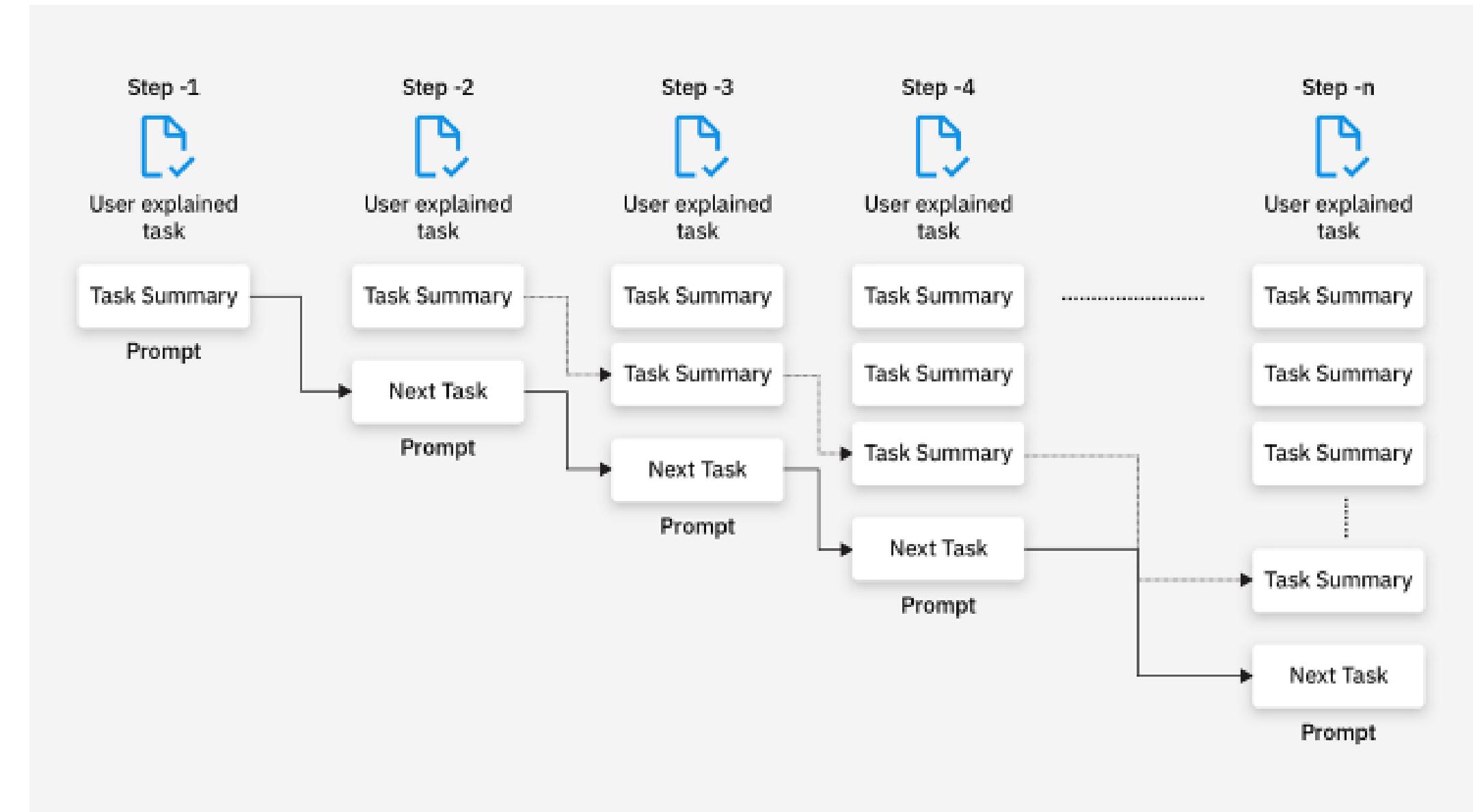
LLM

9.

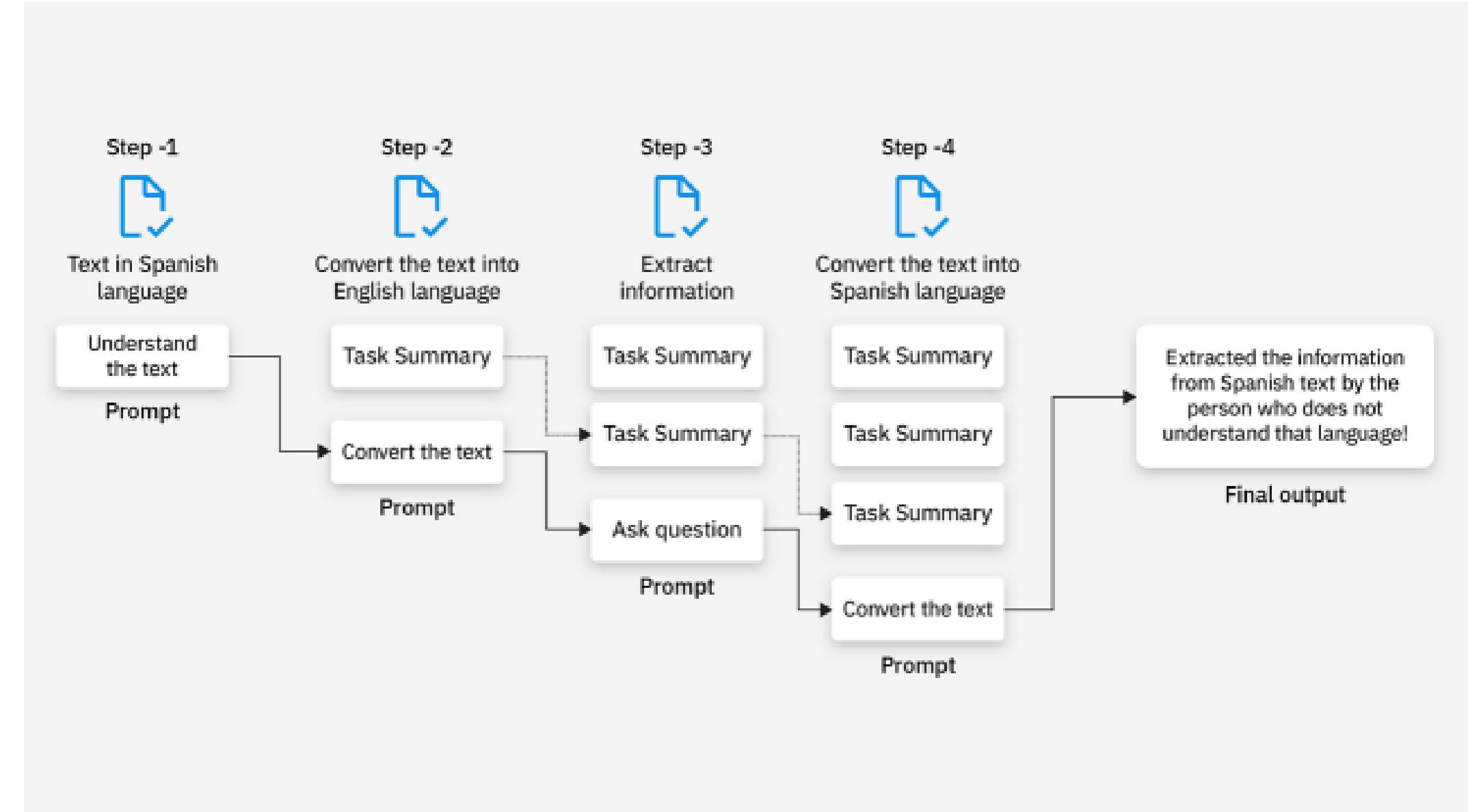
Automatic chain of thought (Auto-CoT)



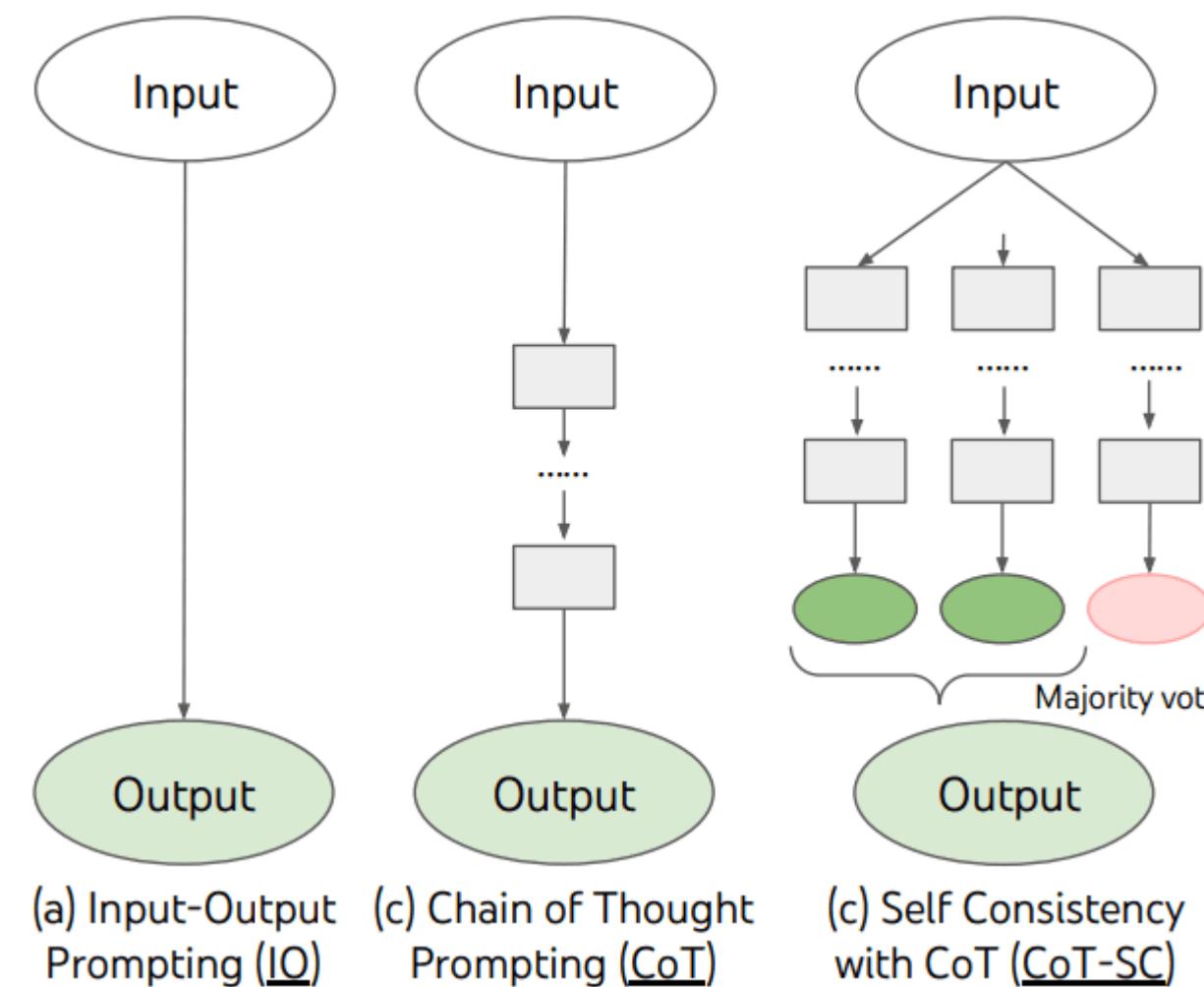
Prompt Chaining



Prompt Chaining



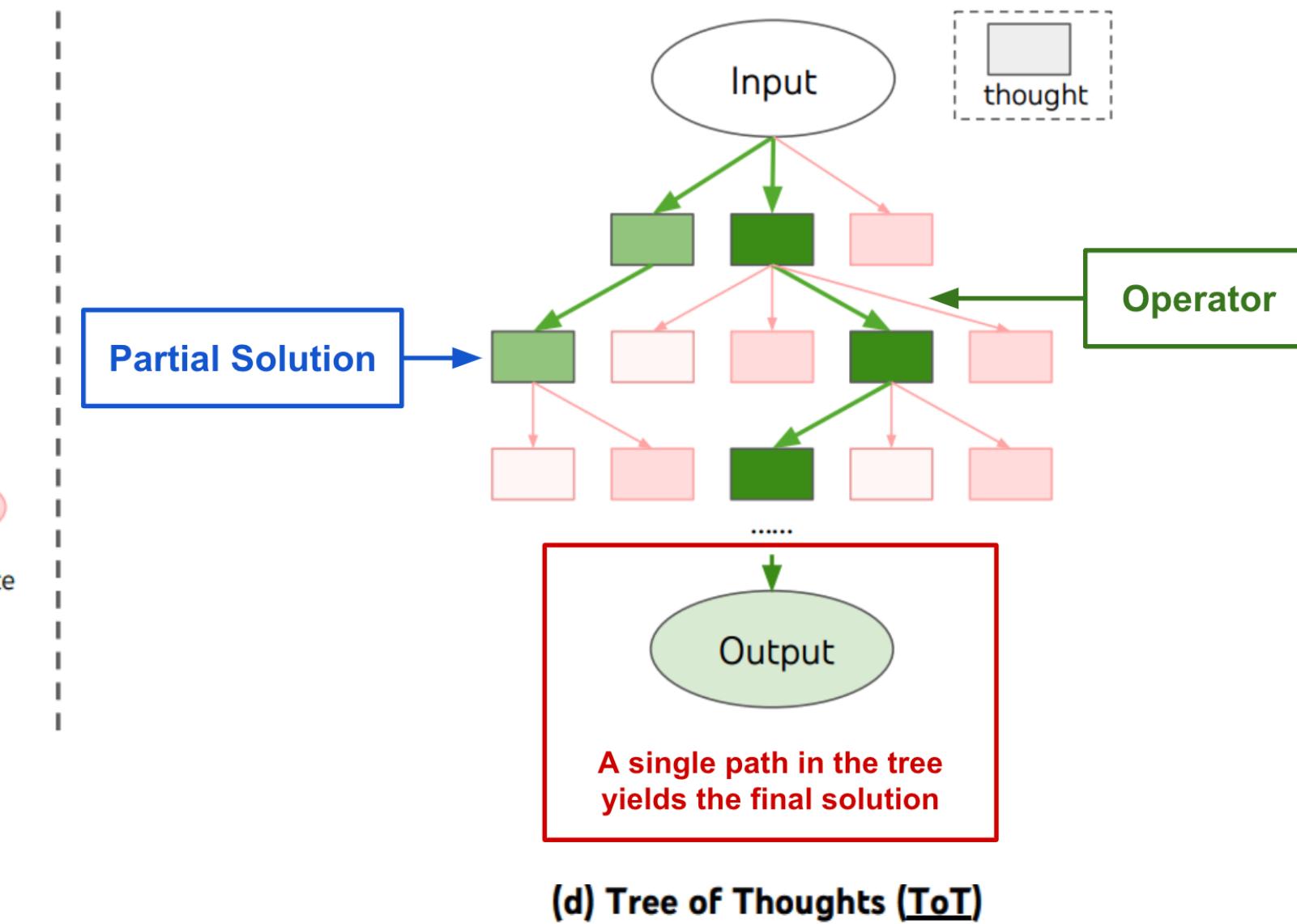
Tree of Thoughts (ToT)



(a) Input-Output
Prompting (IO)

(b) Chain of Thought
Prompting (CoT)

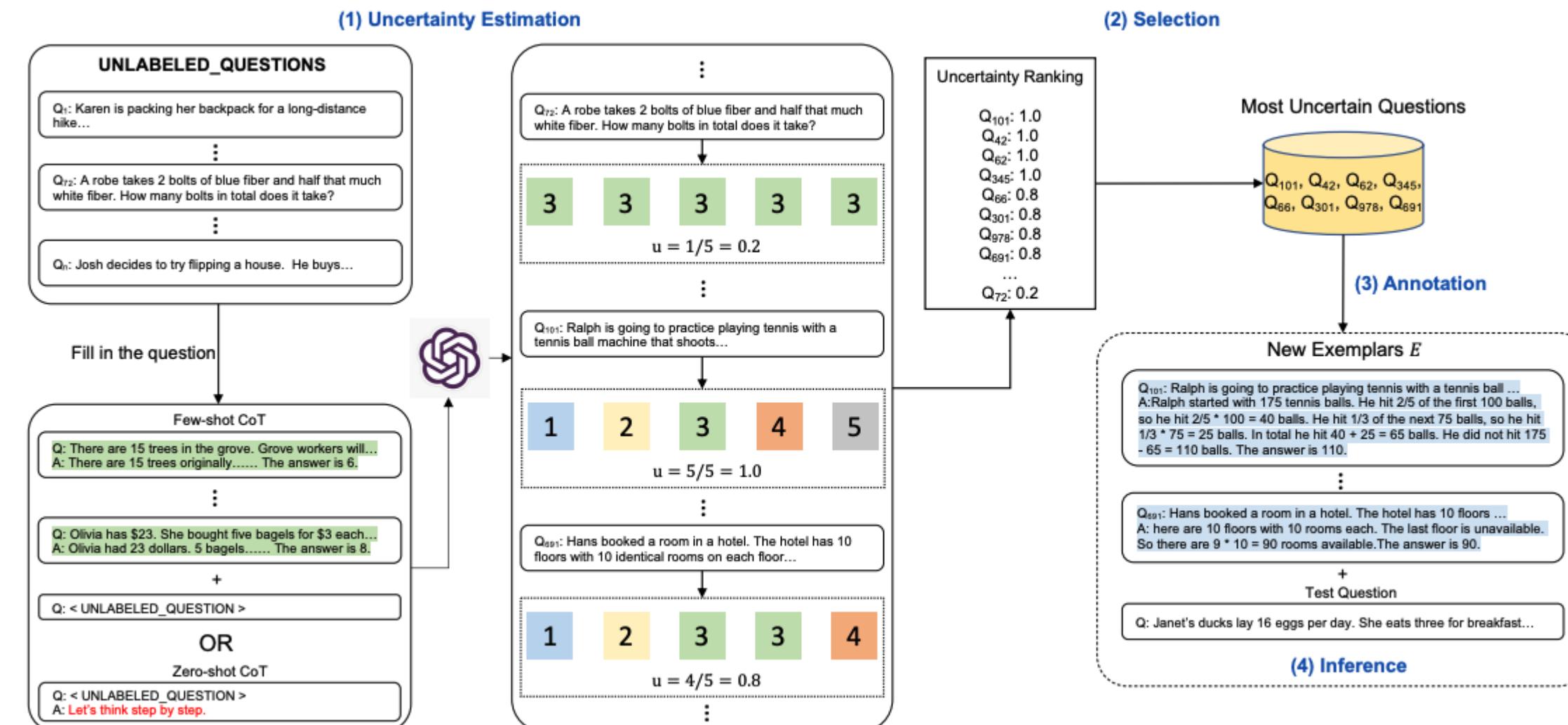
(c) Self Consistency
with CoT (CoT-SC)



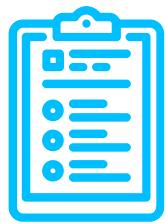
(d) Tree of Thoughts (ToT)



Active-Prompt



4.



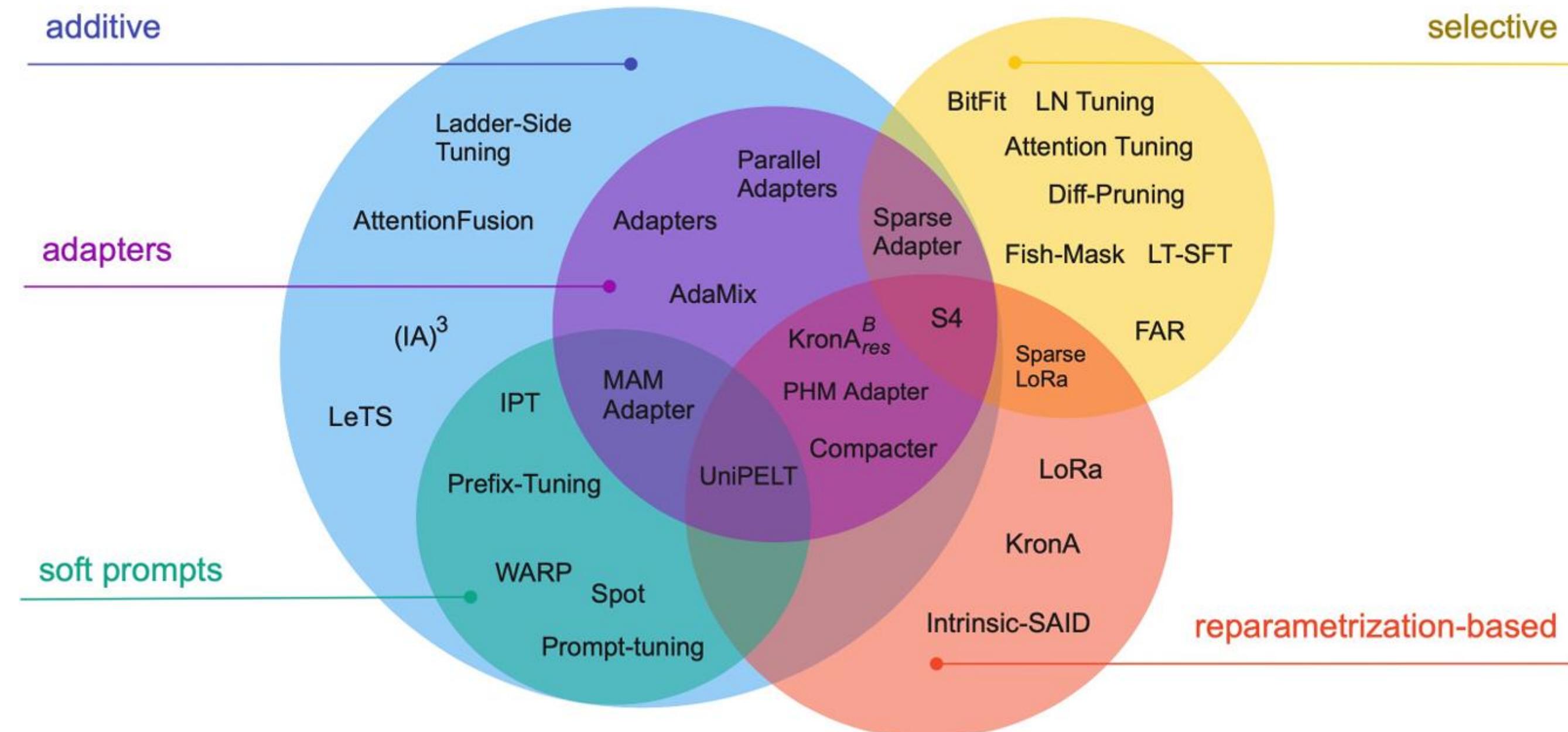
LLM fine-*tuning*

TRANSFORMATEC

> Reinventa el mundo <



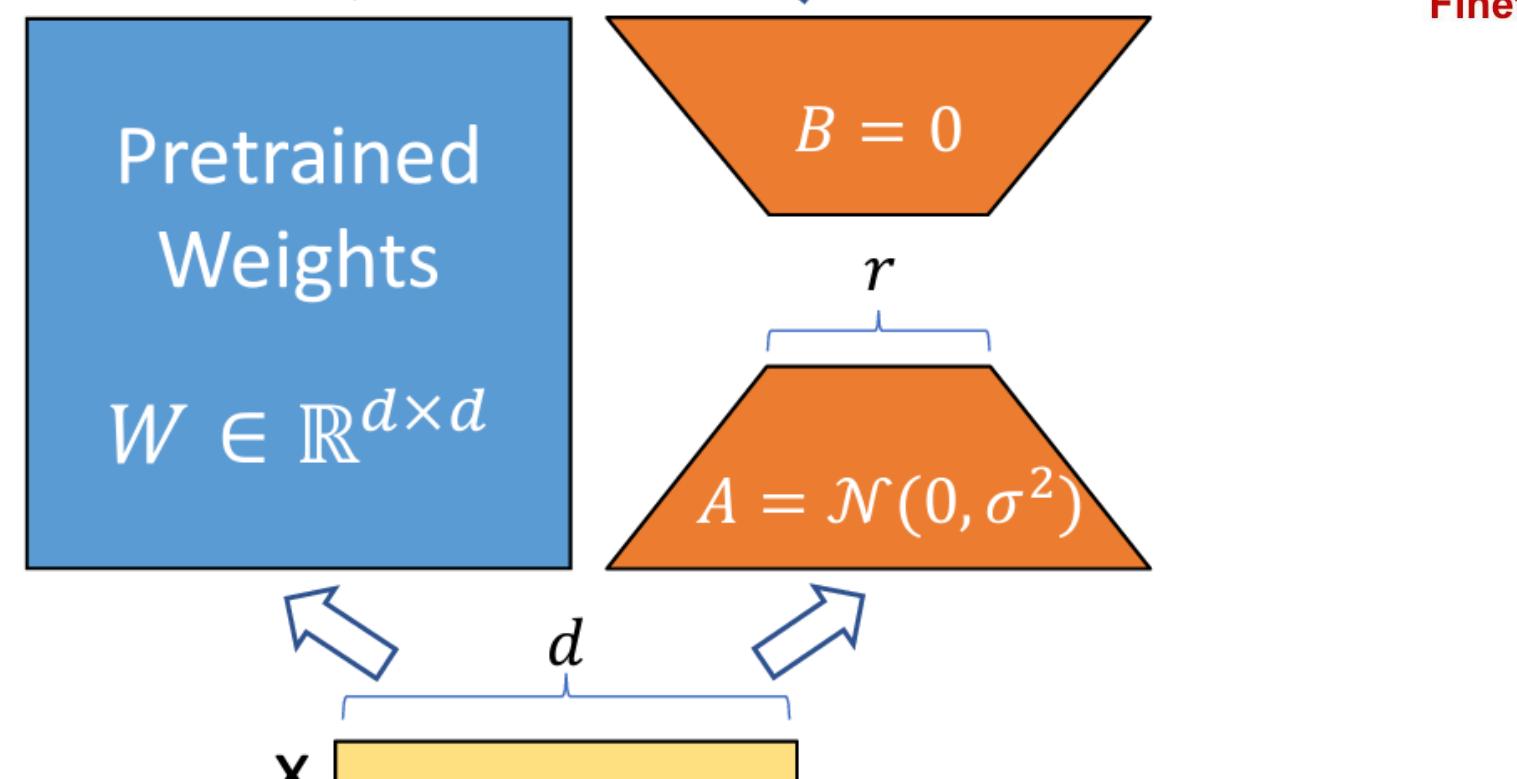
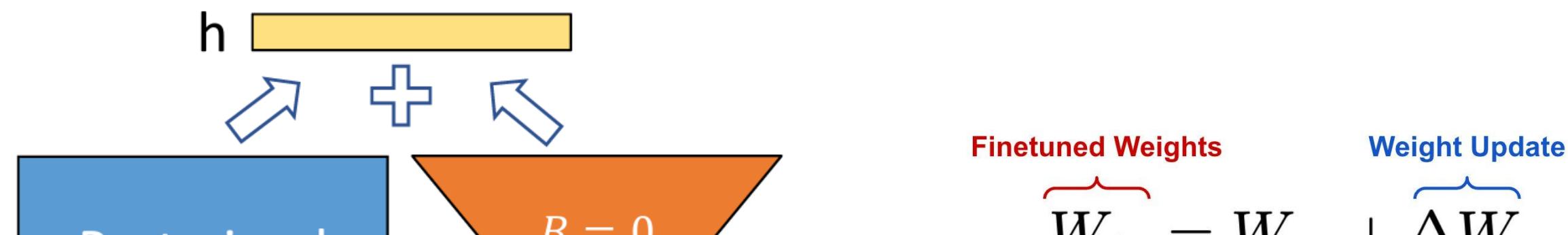
Natural language processing



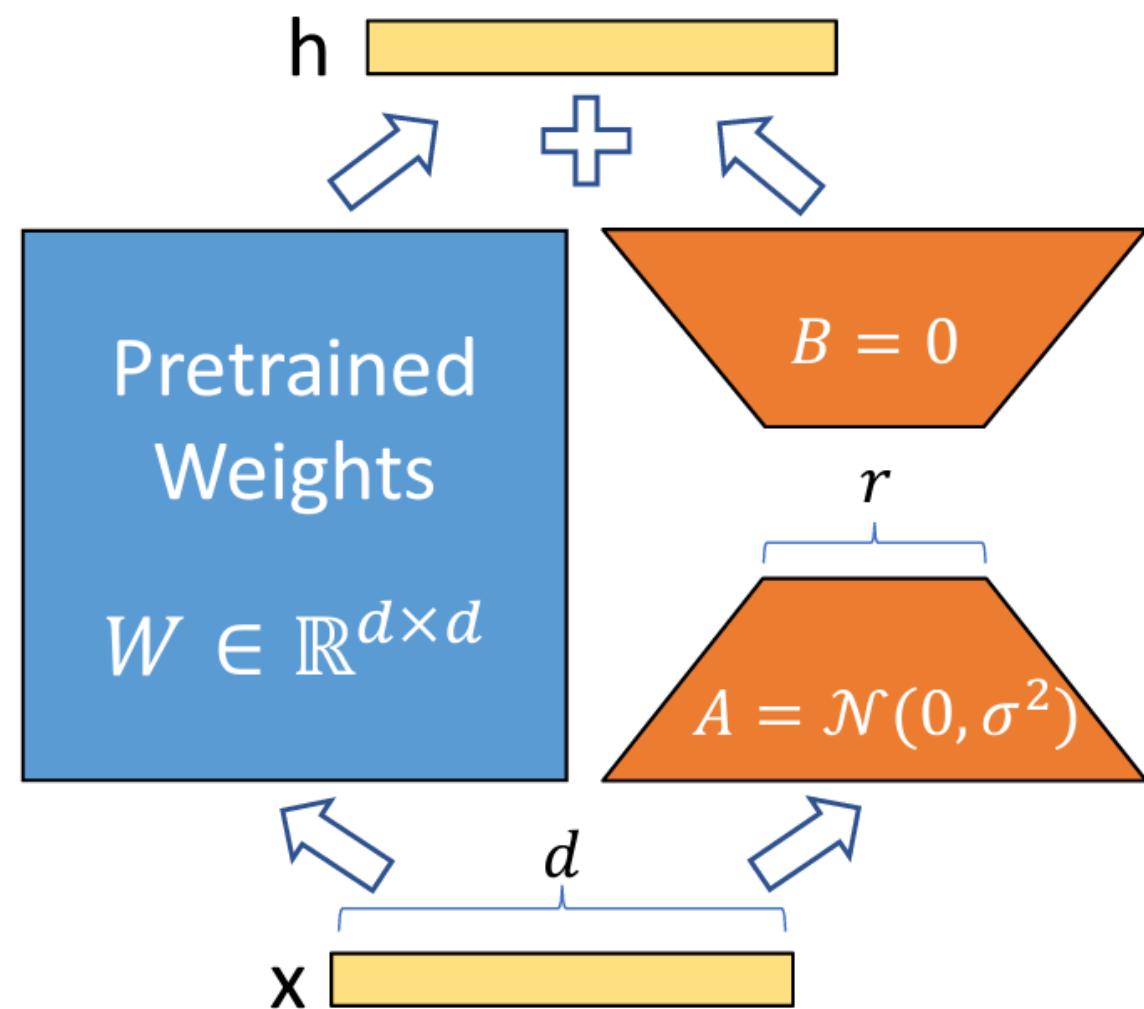
<https://arxiv.org/pdf/2303.15647>

TRANSFORMATEC

LoRA



LoRA

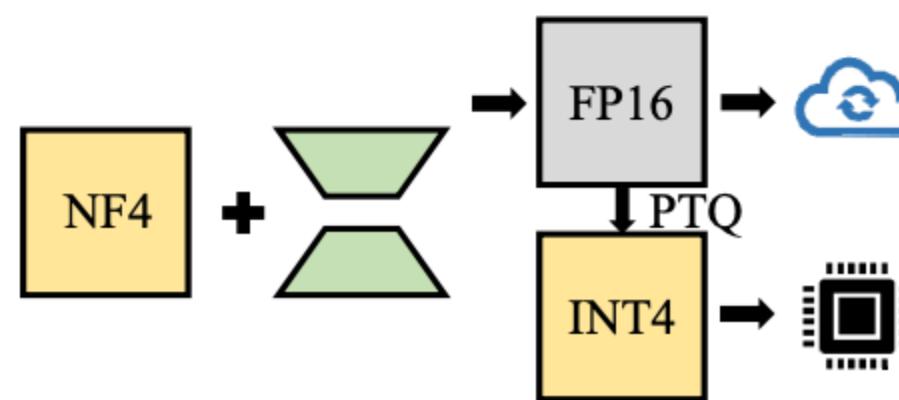
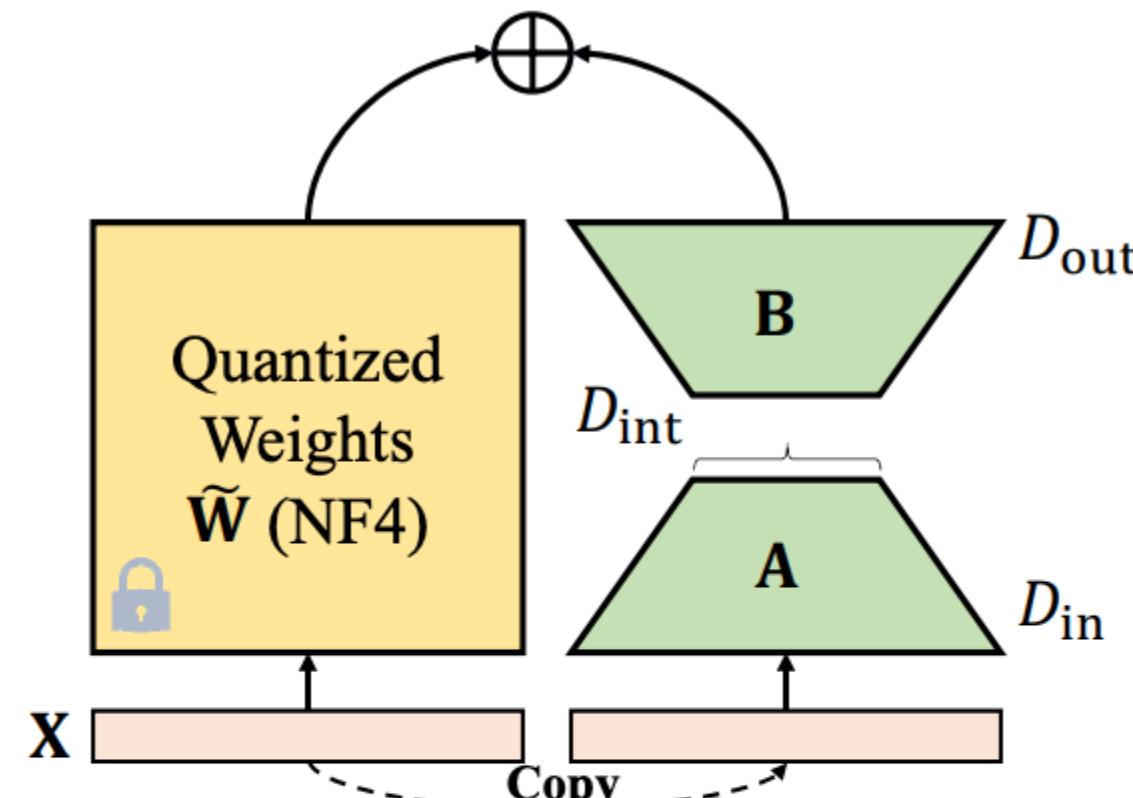


$$W_{\text{ft}} = W_{\text{pt}} + \Delta W = W_{\text{pt}} + \underbrace{AB}_{\text{Low Rank}}$$

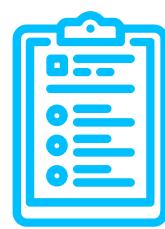
where $W_{\text{ft}}, W_{\text{pt}}, \Delta W, AB \in \mathbb{R}^{d \times d}$
 and $\underbrace{A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}}_{\text{Low Rank}}$



QLoRA



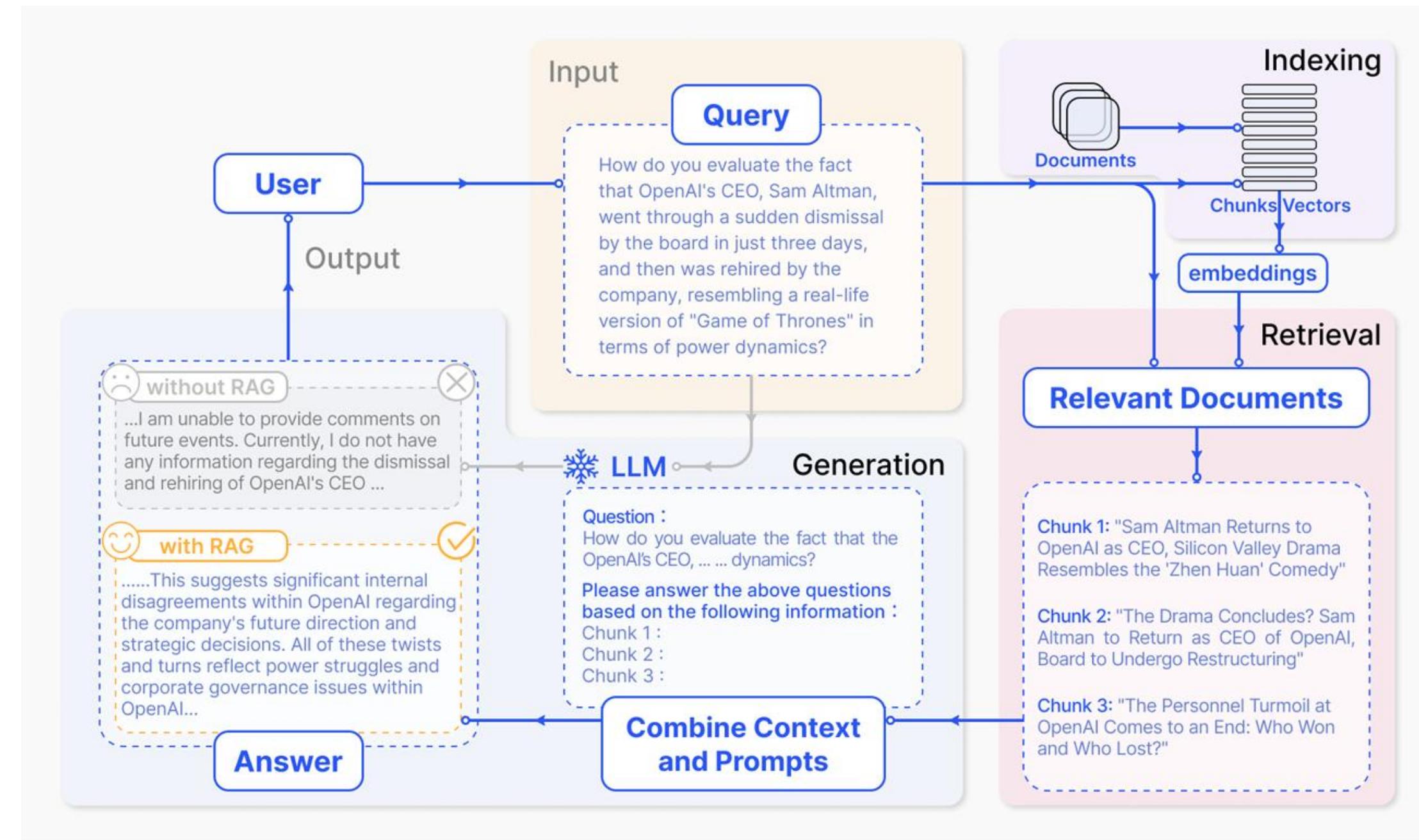
5.



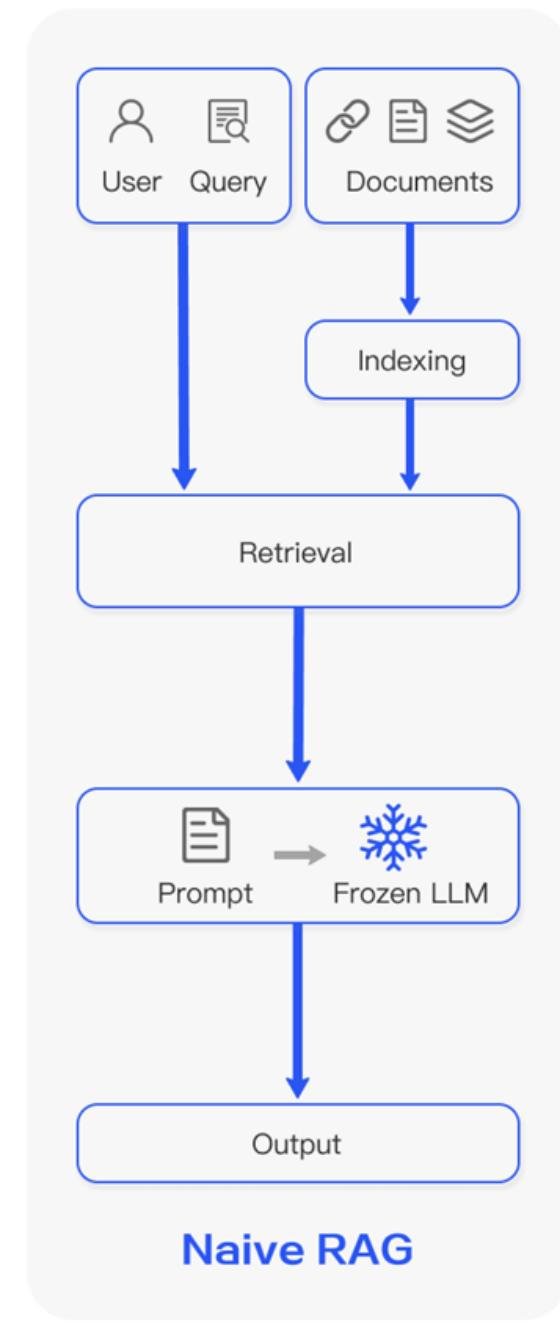
Retrieval Augmented Generation (RAG)



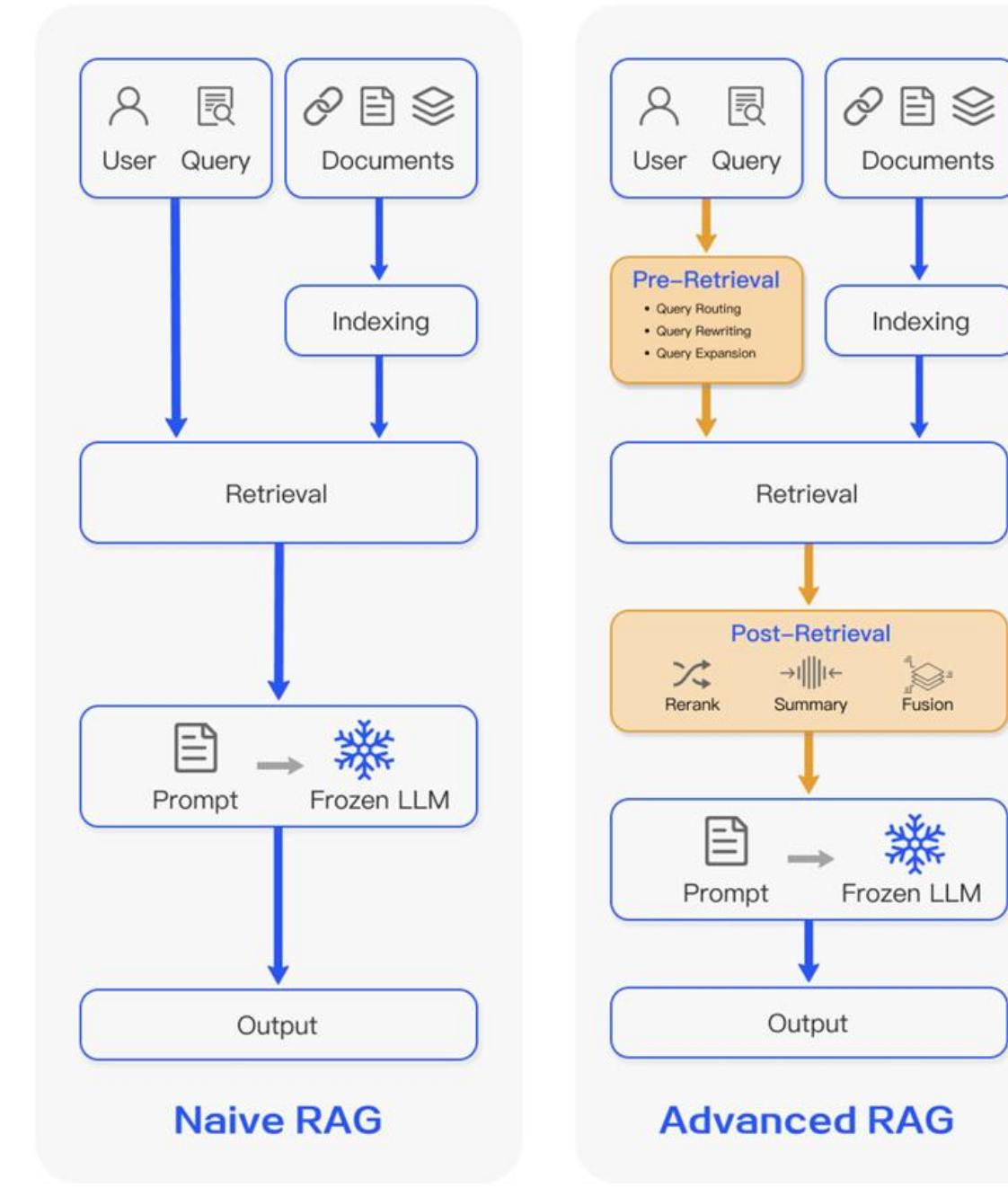
Retrieval Augmented Generation (RAG)



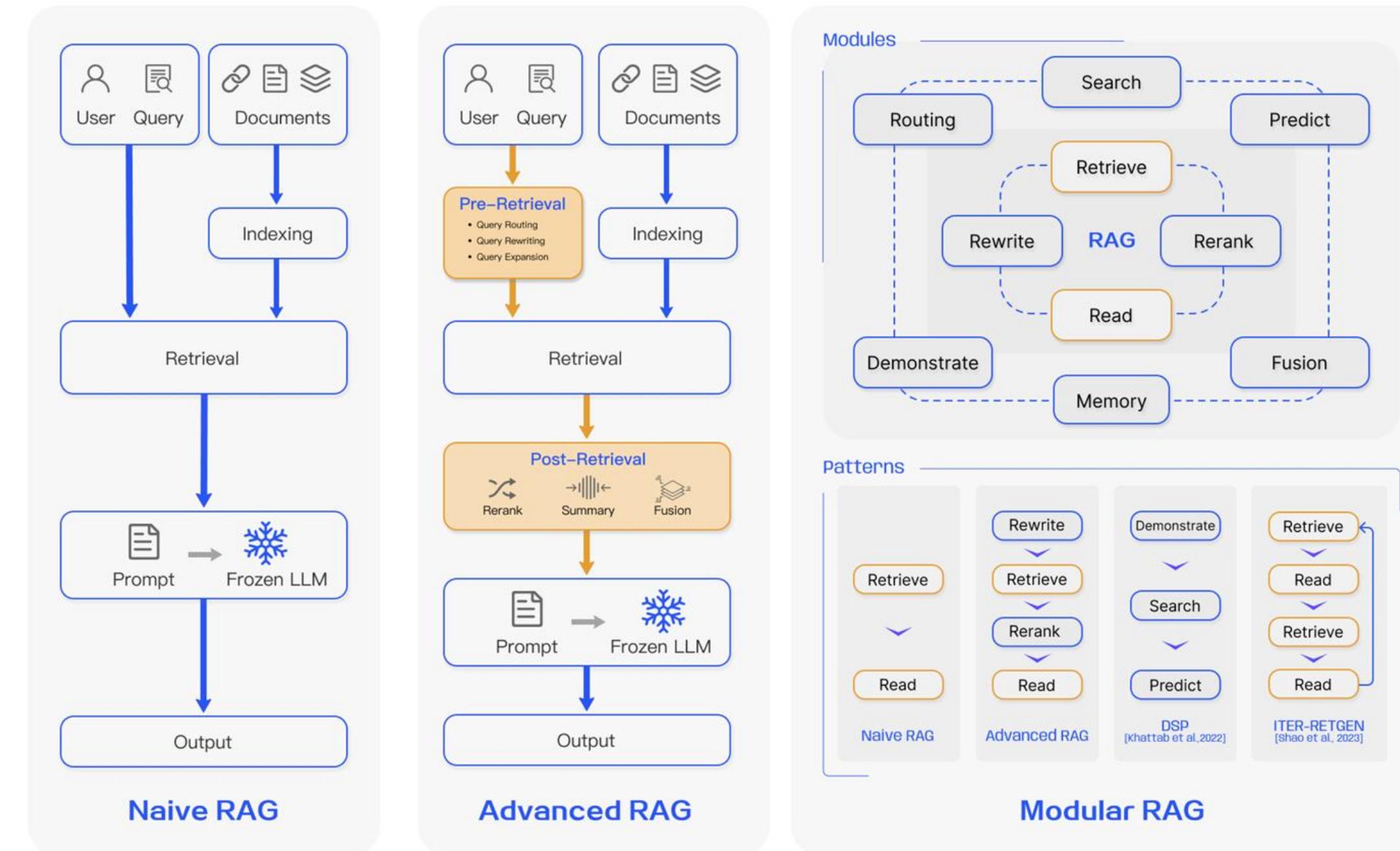
Retrieval Augmented Generation (RAG)



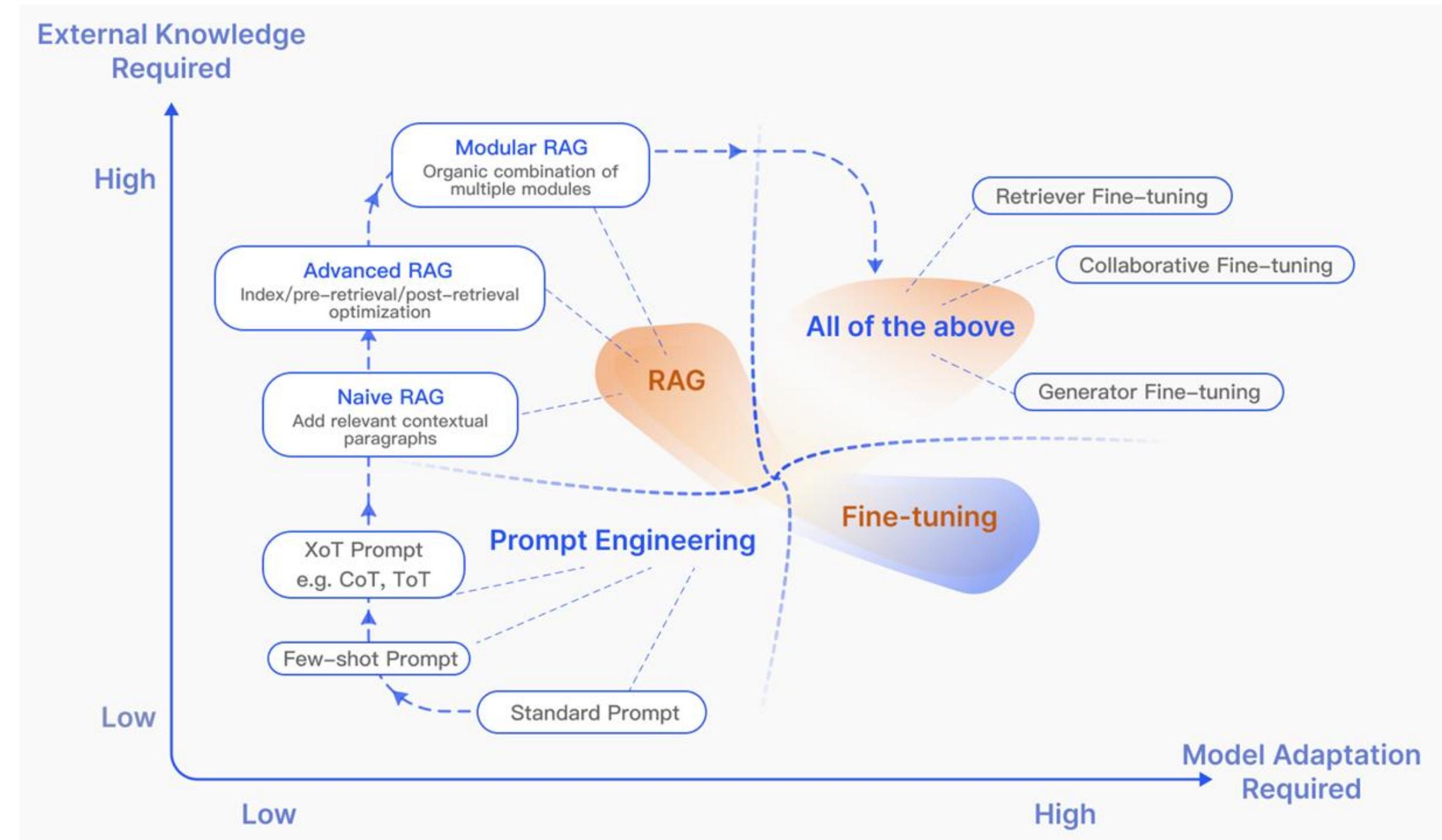
Retrieval Augmented Generation (RAG)



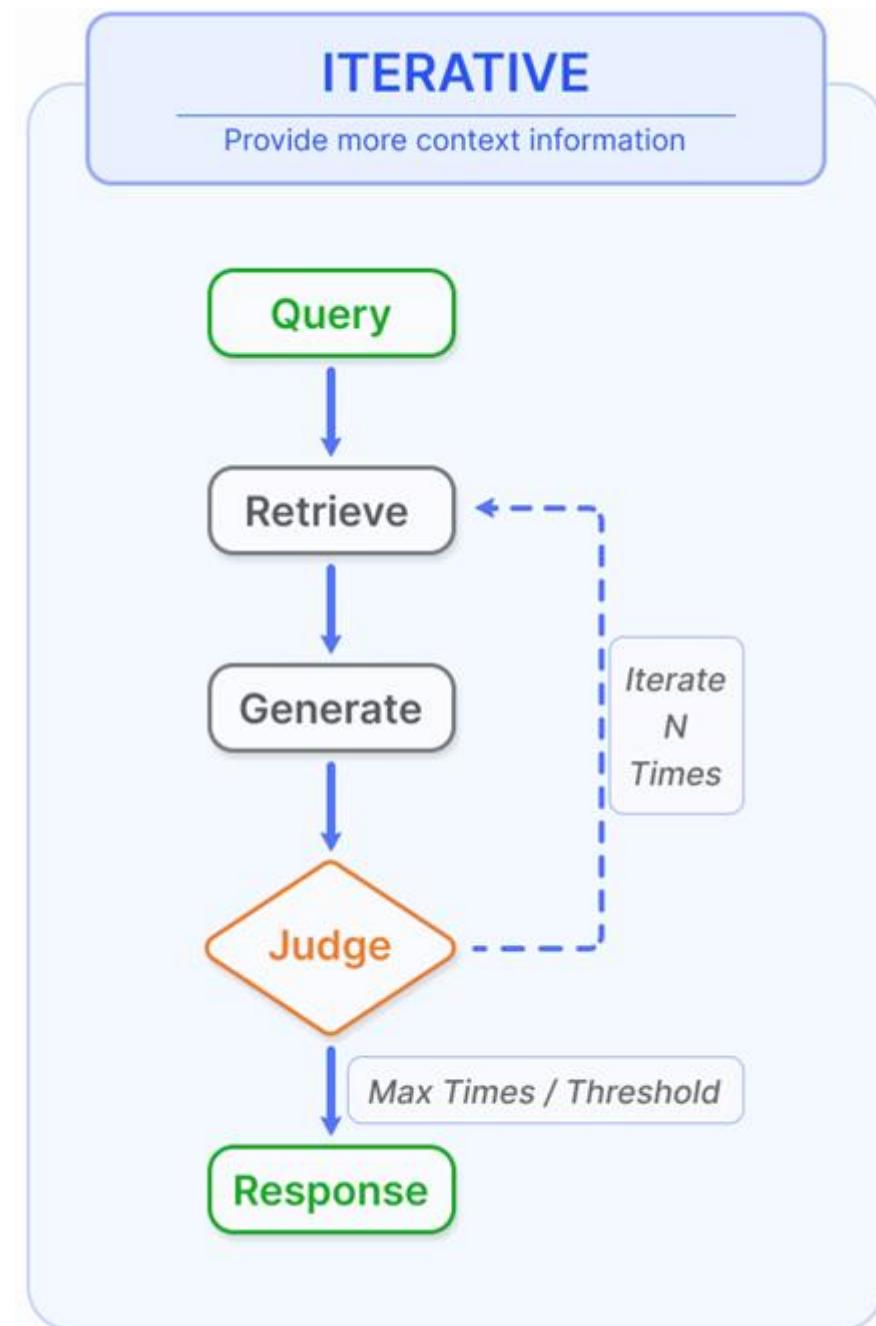
Retrieval Augmented Generation (RAG)



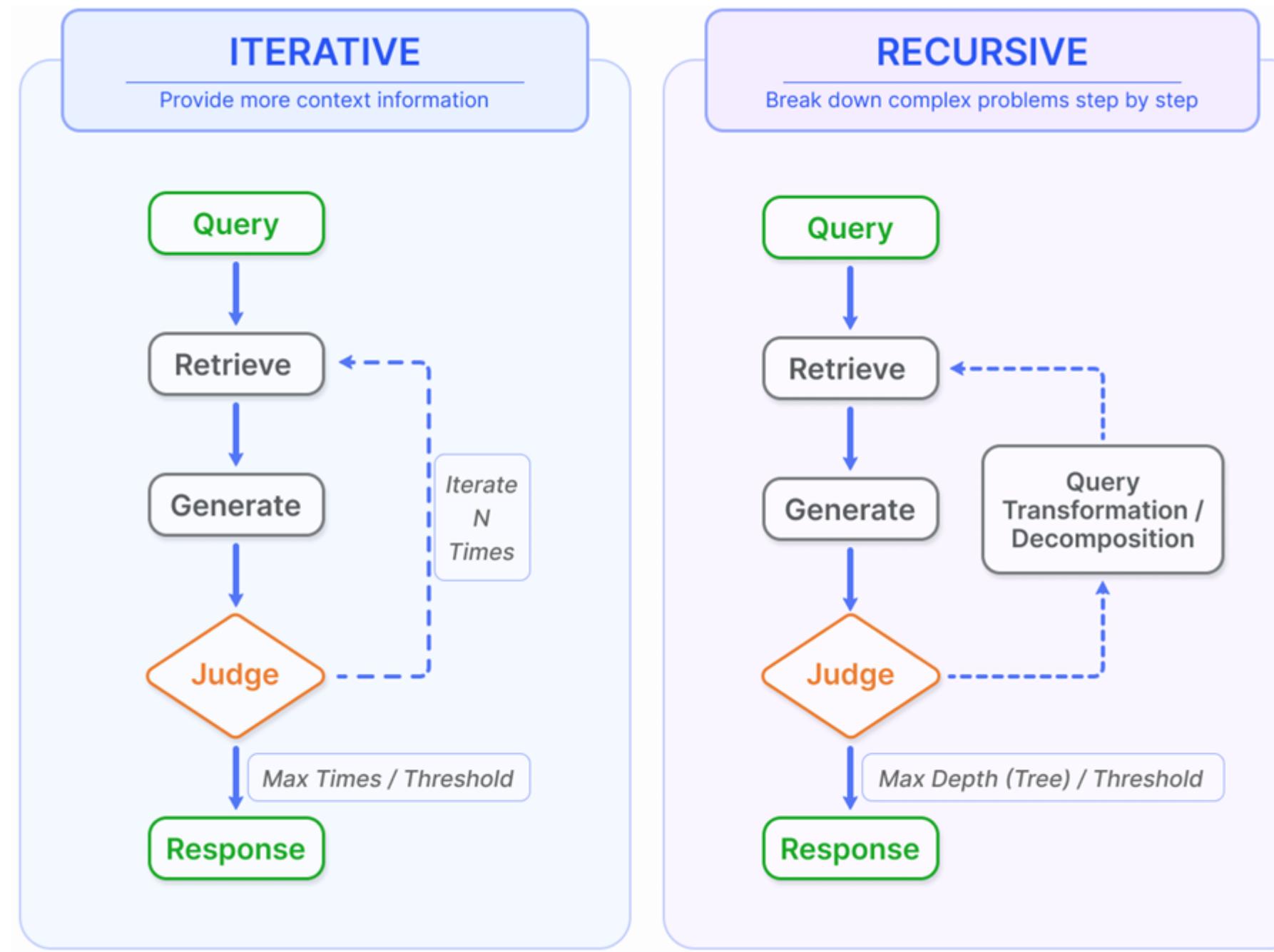
Retrieval Augmented Generation (RAG)



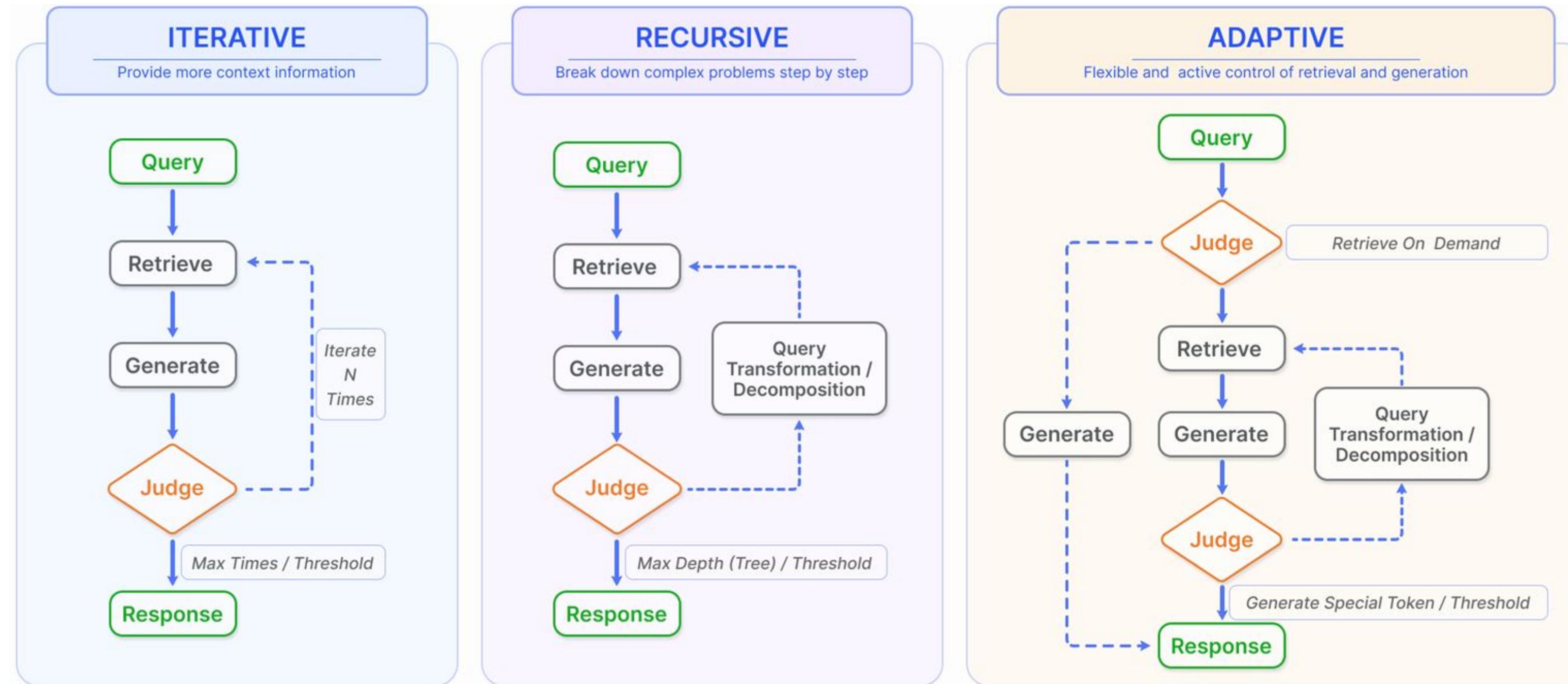
Retrieval Augmented Processes



Retrieval Augmented Processes



Retrieval Augmented Processes



GRACIAS

Victor Flores Benites