

**ACL 2023 Tutorial:**

# **Retrieval-based Language Models and Applications**

Akari Asai, Sewon Min, Zexuan Zhong, Danqi Chen

University of Washington, Princeton University

발표자: 송선영

2024/01/23

# Retrieval-based language model

---

	What do retrieve?	How to use retrieval?	When to retrieve?
REALM (Guu et al 2020)	Text chunks	Input layer	Once
Retrieve-in-context LM (Shi et al 2023, Ram et al 2023)	Text chunks	Input layer	Every n tokens
RETRO (Borgeaud et al. 2021)	Text chunks	Intermediate layers	Every n tokens
kNN-LM (Khandelwal et al. 2020)	Tokens	Output layer	Every token
FLARE (Jiang et al. 2023)	Text chunks	Input layer	Every n tokens (adaptive)
Adaptive kNN-LM (He et al 2021, Alon et al 2022, etc)	Tokens	Output layer	Every n tokens (adaptive)
Entities as Experts (Fevry et al. 2020), Mention Memory (de Jong et al. 2022)	Entities or entity mentions	Intermediate layers	Every entity mentions
Wu et al. 2022, Bertsch et al. 2023, Rubin & Berant. 2023	Text chunks <b>from the input</b>	Intermediate layers	Once or every n tokens

# Retrieval-based language model

---

	What do retrieve?	How to use retrieval?	When to retrieve?
REALM (Guu et al 2020)	Text chunks	Input layer	Once
Retrieve-in-context LM (Shi et al 2023, Ram et al 2023)	Text chunks	Input layer	Every n tokens
RETRO (Borgeaud et al. 2021)	Text chunks	Intermediate layers	Every n tokens
kNN-LM (Khandelwal et al. 2020)	Tokens	Output layer	Every token
FLARE (Jiang et al. 2023)	Text chunks	Input layer	Every n tokens ( <i>adaptive</i> )
Adaptive kNN-LM (He et al 2021, Alon et al 2022, etc)	Tokens	Output layer	Every n tokens ( <i>adaptive</i> )
Entities as Experts (Fevry et al. 2020), Mention Memory (de Jong et al. 2022)	Entities or entity mentions	Intermediate layers	Every entity mentions
Wu et al. 2022, Bertsch et al. 2023, Rubin & Berant. 2023	Text chunks <i>from the input</i>	Intermediate layers	Once or every n tokens

# REALM (Guu et al. 2020)

- A masked language model

$x$  = World Cup 2022 was the last before the increase to [MASK] in the 2026 tournament.



# Retrieval-in-context LM (Ram et al. 2023, Shi et al. 2023)

- An auto-regressive language model

$x$  = World Cup 2022 was the last with 32 teams, before the increase to

World Cup 2022 was the last with 32 teams, before the increase to

↓  
**Retrieval**  
↓

\* Can use multiple text blocks too (see the papers!)

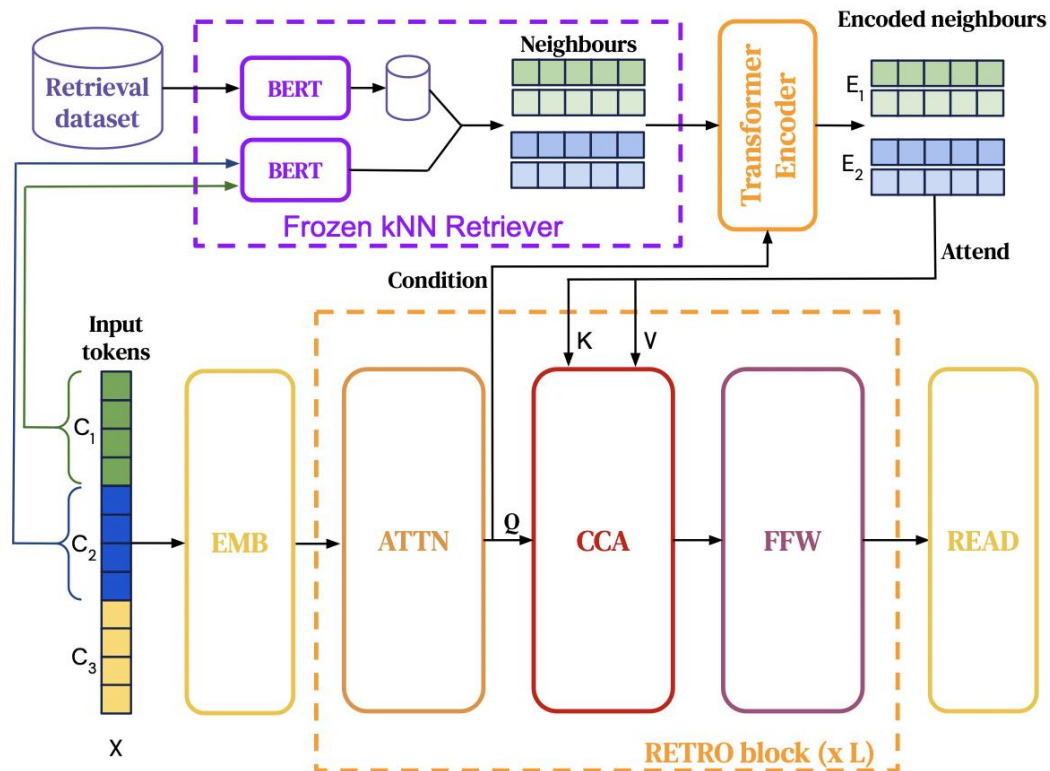
FIFA World Cup 2026 will expand to 48 teams. World Cup 2022 was the last with 32 teams, before the increase to

↓  
**LM**  
↓

48 in the 2026 tournament.

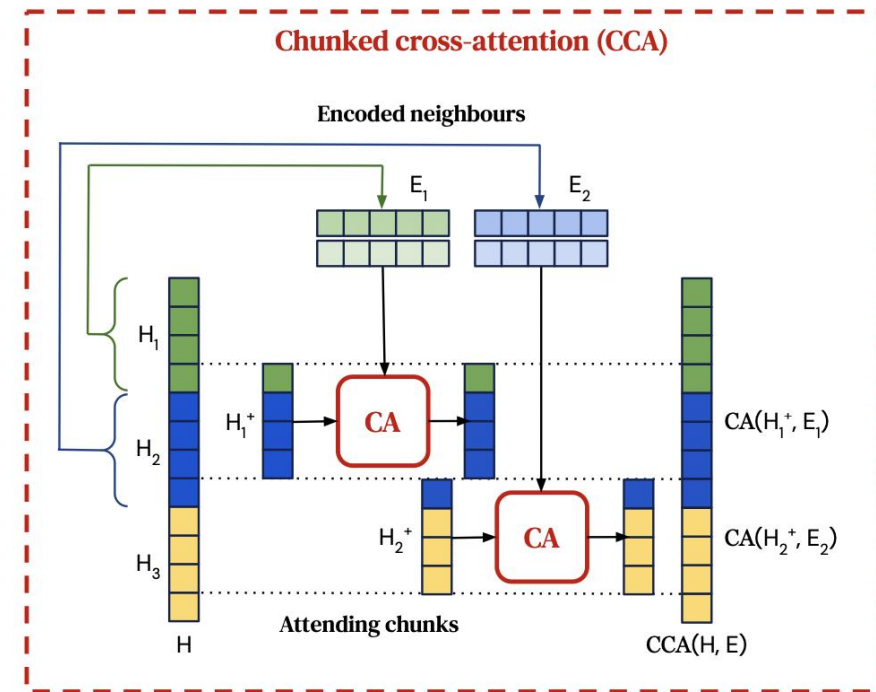
# RETRO (Borgeaud et al. 2022)

- Using Chunked cross-attention (CCA)







$x$  = World Cup 2022 was ~~the last~~ with 32 teams, ~~before~~ the increase to


$x_1$   $x_2$   $x_3$



# kNN-LM (Khandelwal et al. 2020)

- The size of the datastore = # of tokens in the corpus (not vocab size)

Training Contexts $c_i$	Targets $v_i$	Representations $k_i = f(c_i)$
Obama was senator for	Illinois	
Barack is married to	Michelle	
Obama was born in	Hawaii	
...	...	...
Obama is a native of	Hawaii	

Test Context $x$	Target	Representation $q = f(x)$
Obama's birthplace is	?	

Which tokens in a datastore are close to the next token?

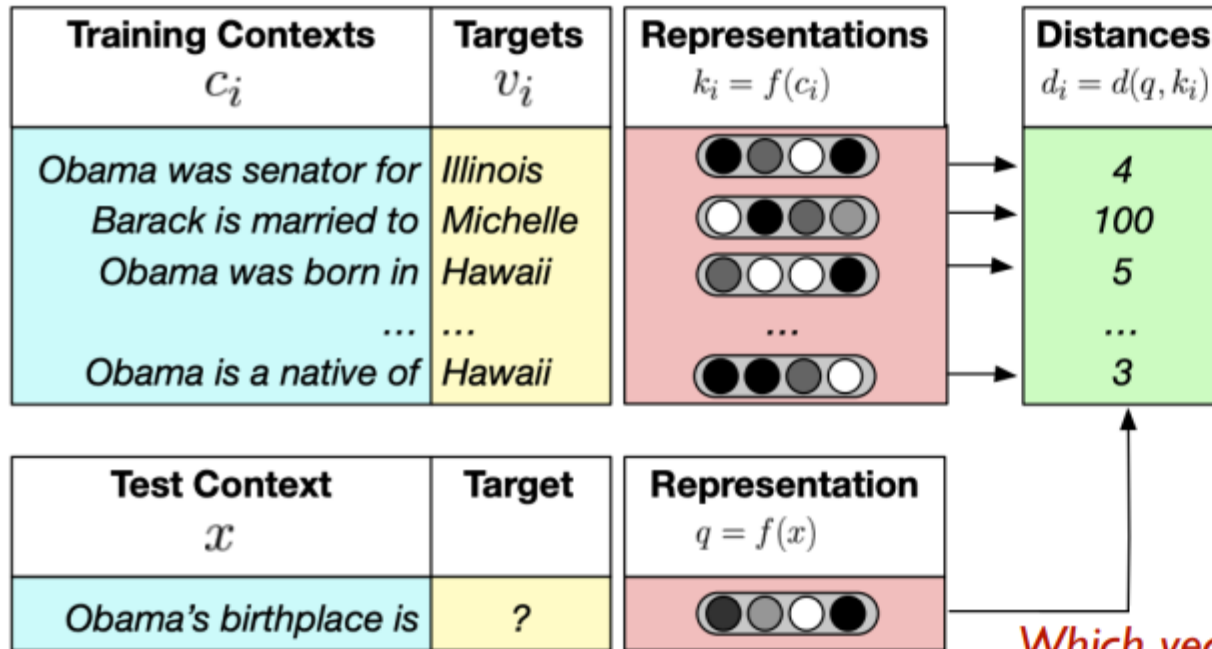
=

Which prefixes in a datastore are close to the prefix we have?

=

Which vectors in a datastore are close to the vector we have?

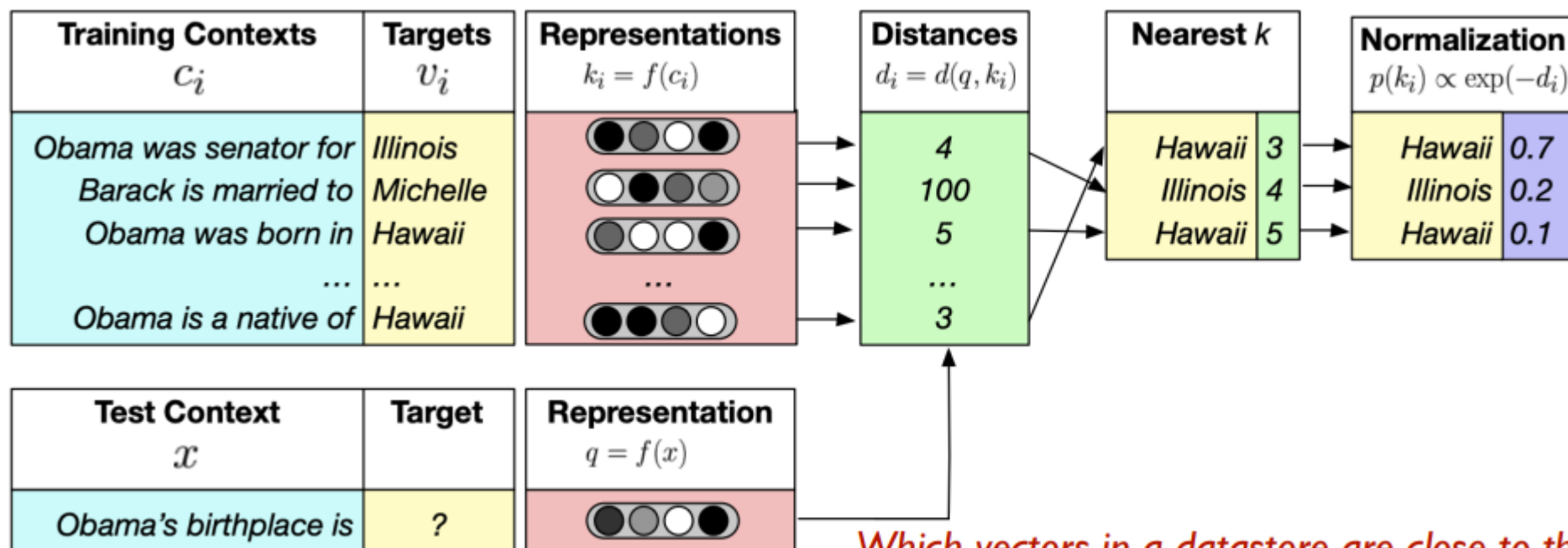
# kNN-LM (Khandelwal et al. 2020)



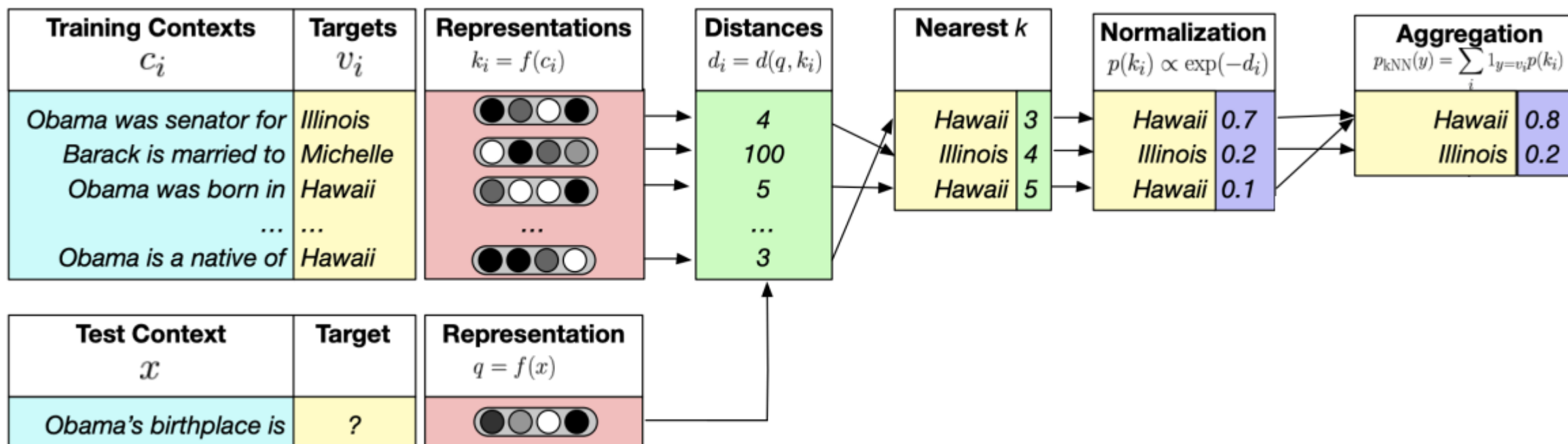
Which vectors in a datastore are close to the vector we have?



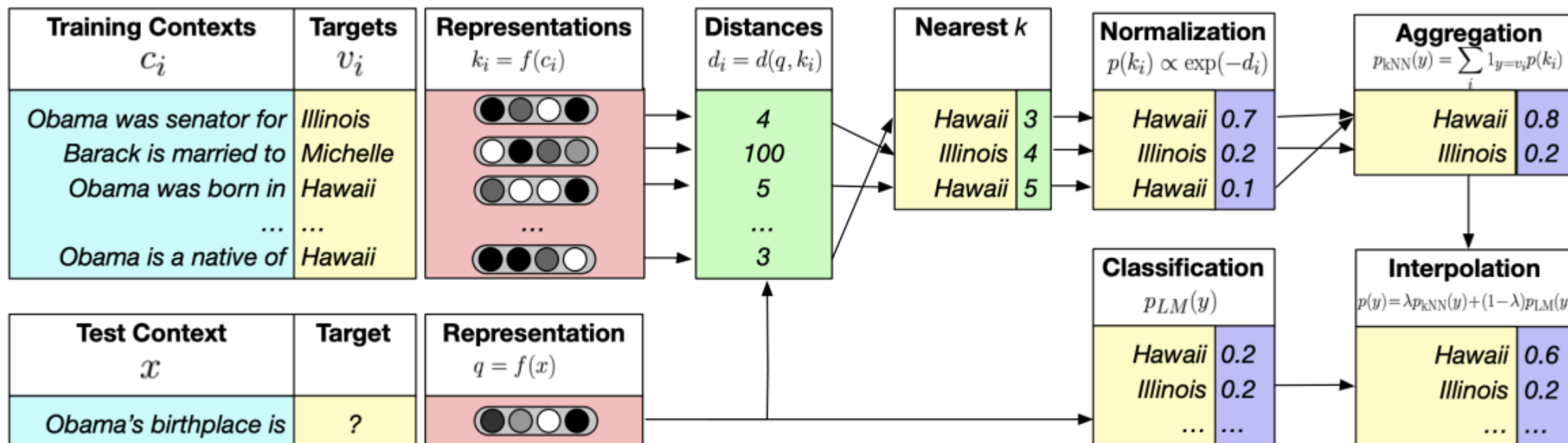
# kNN-LM (Khandelwal et al. 2020)



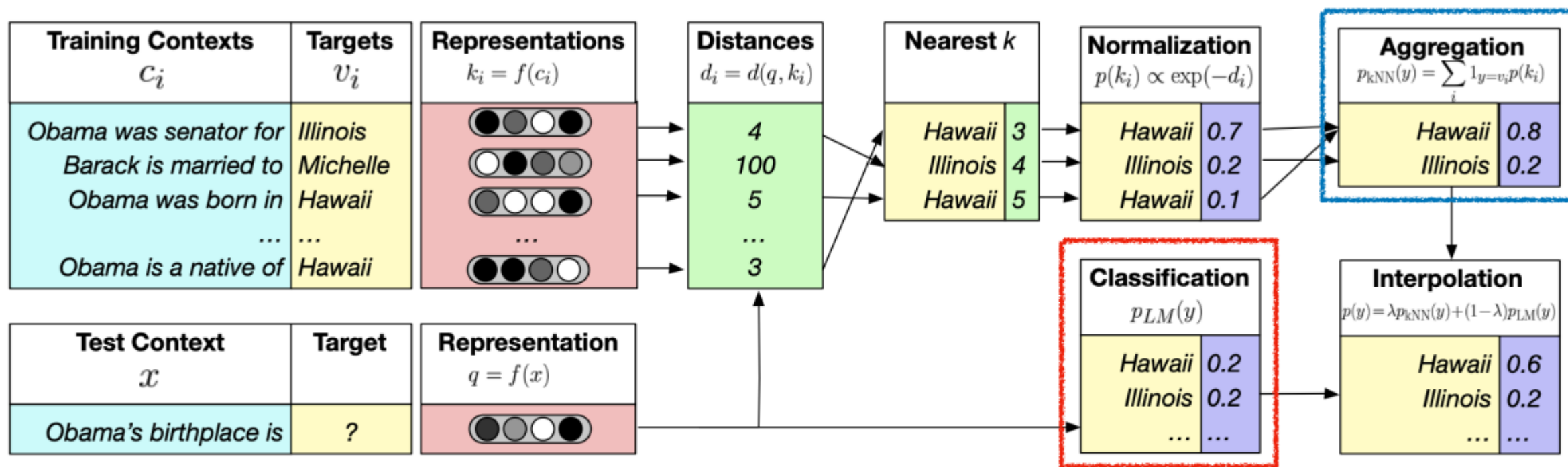
# kNN-LM (Khandelwal et al. 2020)



# kNN-LM (Khandelwal et al. 2020)



# kNN-LM (Khandelwal et al. 2020)



$\lambda$ : hyperparameter

$$P_{kNN-LM}(y|x) = (1 - \lambda)P_{LM}(y|x) + \lambda P_{kNN}(y|x)$$

# Adaptive kNN-LM (He et al. 2021, Alon et al. 2022)

- Adaptive retrieval of **tokens**

kNN-LM (Khandelwal et al. 2020)	Tokens	Output layer	Every token
Adaptive kNN-LM (He et al 2021, Alon et al 2022, etc)	Tokens	Output layer	Every n tokens ( <i>adaptive</i> )

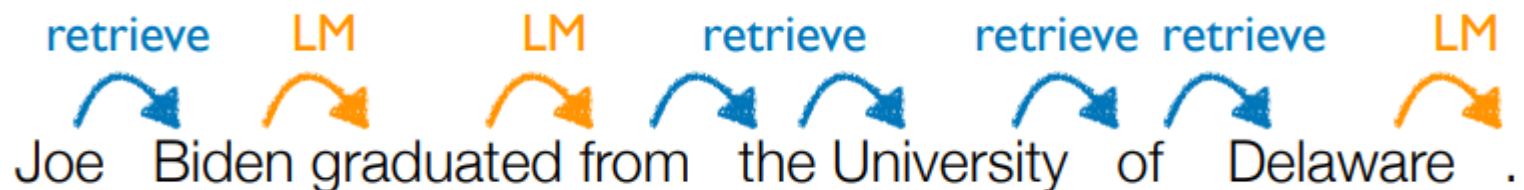
# Adaptive kNN-LM (He et al. 2021, Along et al. 2022)

---

kNN-LM

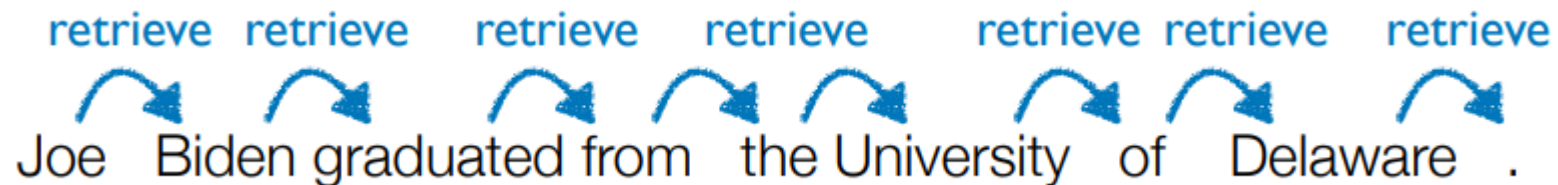


Adaptive  
kNN-LM

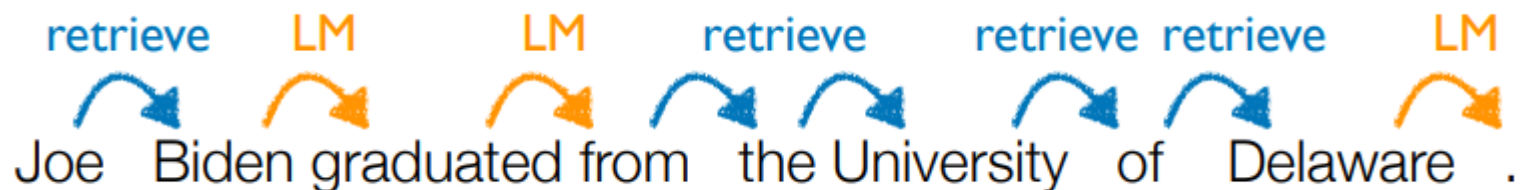


# Adaptive kNN-LM (He et al. 2021, Along et al. 2022)

kNN-LM



Adaptive  
kNN-LM



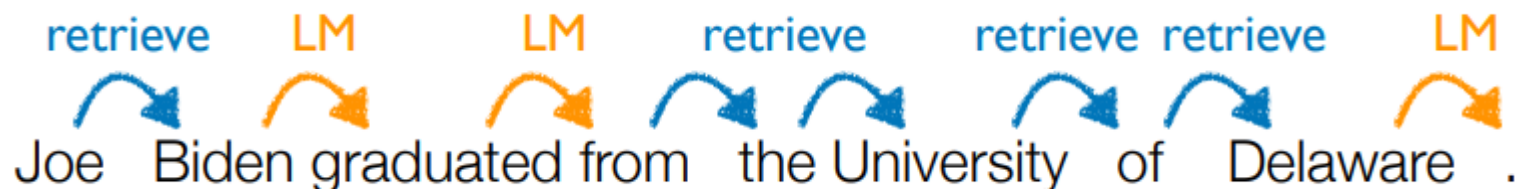
$$P_{\text{kNN-LM}}(y|x) = (1 - \lambda(x))P_{\text{LM}}(y|x) + \lambda(x)P_{\text{kNN}}(y|x)$$

# Adaptive kNN-LM (He et al. 2021, Along et al. 2022)

kNN-LM



Adaptive kNN-LM



$$P_{kNN-LM}(y|x) = (1 - \lambda(x))P_{LM}(y|x) + \lambda(x)P_{kNN}(y|x)$$

$$\mathcal{L} = \frac{1}{T} \sum_t [\log p(w_t|c_t; \lambda_\theta(c_t)) - a \cdot \lambda_\theta(c_t)]$$



# FLARE (Jiang et al. 2023)

- Adaptive retrieval of **chunks**

Retrieve-in-context LM (Shi et al 2023, Ram et al 2023)

Text chunks

Input layer

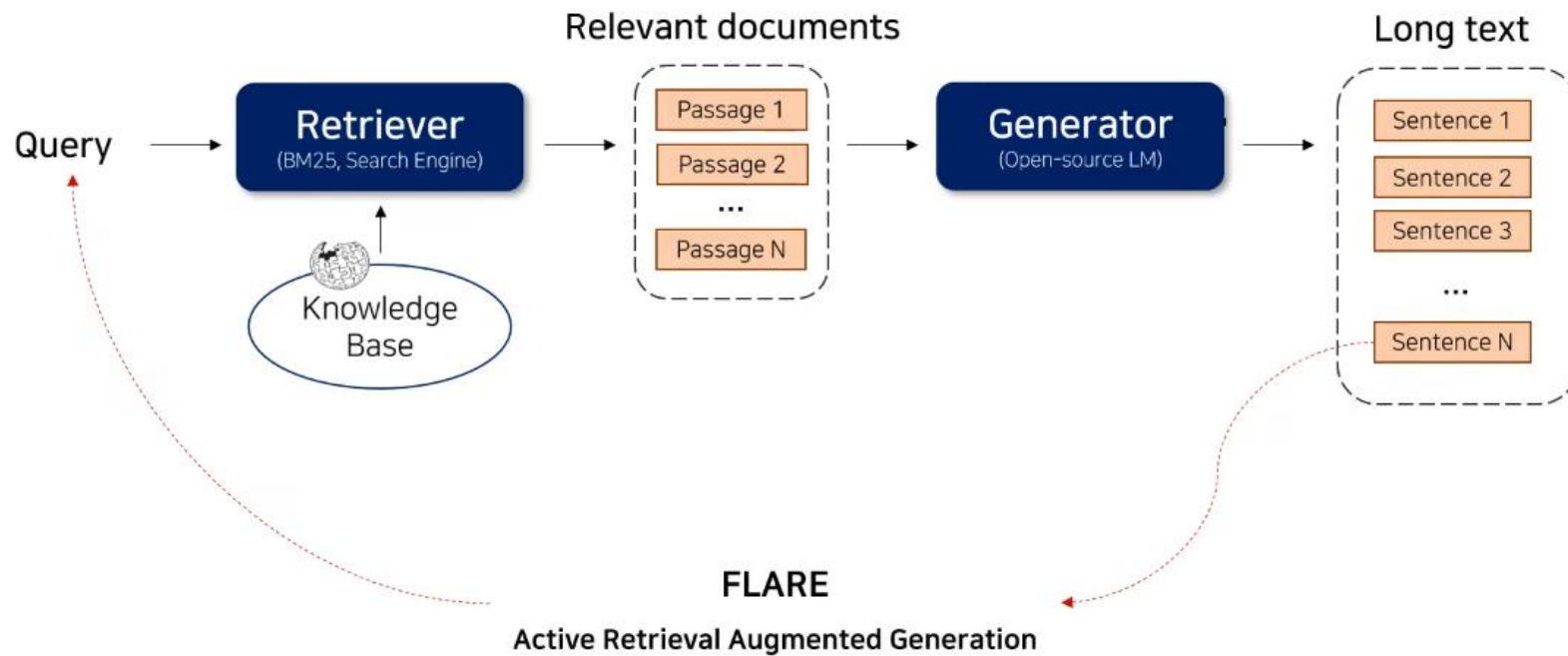
Every n tokens

FLARE (Jiang et al. 2023)

Text chunks

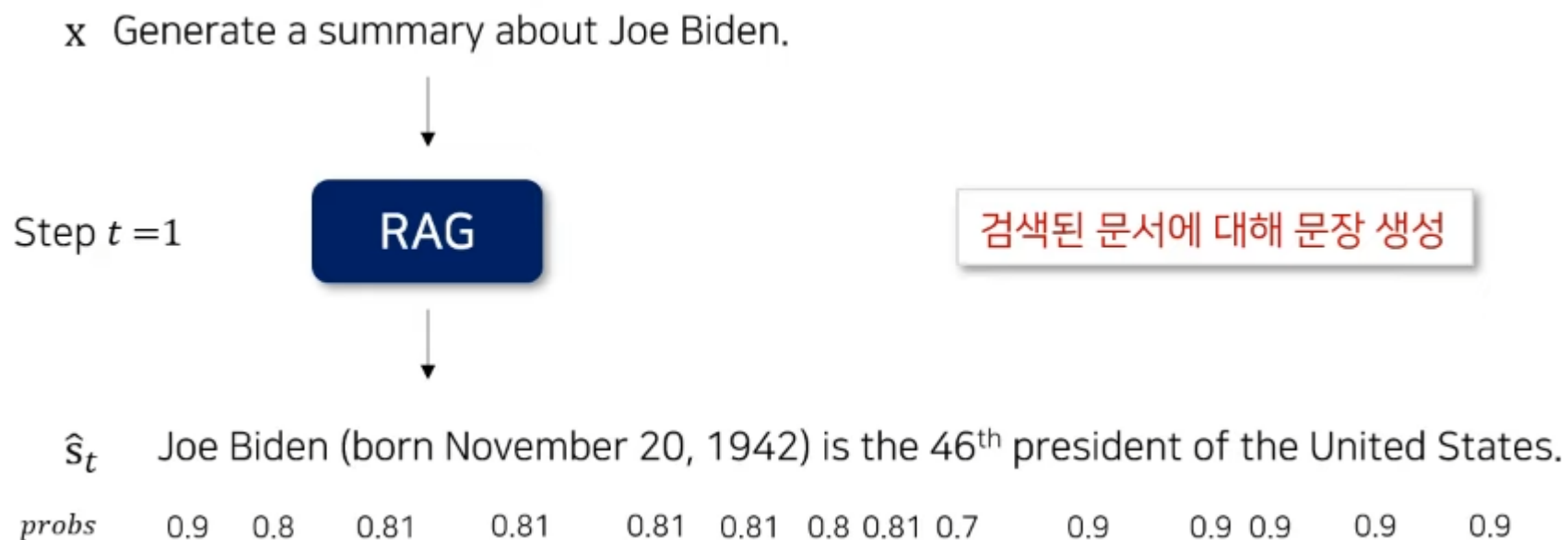
Input layer

Every n tokens  
(adaptive)



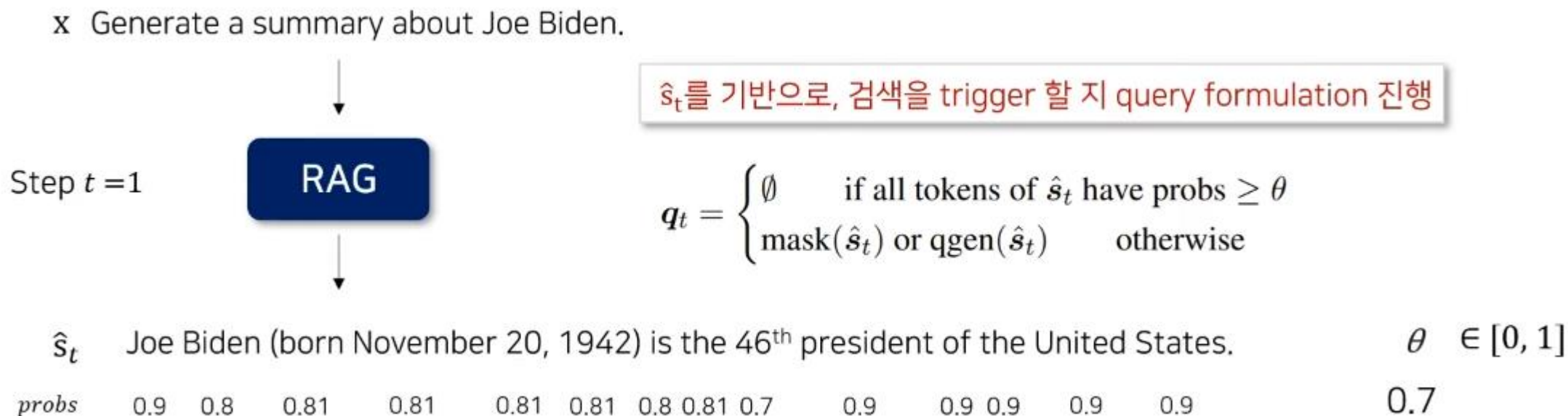
# FLARE (Jiang et al. 2023)

---



# FLARE (Jiang et al. 2023)

- Query formulation



# FLARE (Jiang et al. 2023)

- Query formulation

1. Mask (implicit)
2. Qgen (explicit)

x Generate a summary about Joe Biden.

$Y_1$  Joe Biden (born November 20, 1942) is the 46<sup>th</sup> president of the United States.

Step  $t = 2$

RAG

$\hat{s}_t$ 를 기반으로, 검색을 trigger 할 지 query formulation 진행

$$q_t = \begin{cases} \emptyset & \text{if all tokens of } \hat{s}_t \text{ have probs } \geq \theta \\ \text{mask}(\hat{s}_t) \text{ or qgen}(\hat{s}_t) & \text{otherwise} \end{cases}$$

$\hat{s}_t$  Joe Biden attended the University of Pennsylvania, where he earned a law degree.  $\theta$

probs    0.9   0.8   0.81   0.74   0.62   0.7   0.58   0.9   0.9   0.9   0.7   0.5   0.73   0.7

# FLARE (Jiang et al. 2023)

## 1. Implicit query by masking

$\hat{s}_t$  Joe Biden attended the University of Pennsylvania, where he earned a law degree.  $\beta$

*probs* 0.9 0.8 0.81 0.74 0.62 0.7 0.58 0.9 0.9 0.9 0.7 0.5 0.73 0.8

Joe Biden attended , where he earned .

RAG

# FLARE (Jiang et al. 2023)

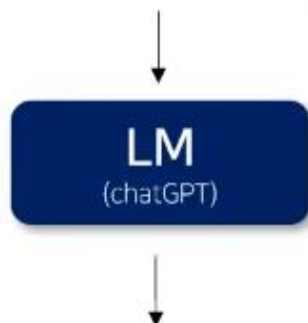
## 2. Explicit query by question generation (qgen)

$\hat{s}_t$  Joe Biden attended the University of Pennsylvania, where he earned a law degree.  $\beta$

*probs* 0.9 0.8 0.81 0.74 0.62 0.7 0.58 0.9 0.9 0.9 0.7 0.5 0.73 0.8

Ask a question to which the answer is "the University of Pennsylvania"

Ask a question to which the answer is "a law degree"



What university did Joe Biden attend?

What degree did Joe Biden earn?

$\hat{s}_t$  에서 probs가 낮은 span  $z$  를 대상으로 explicit query 생성

### Prompt 3.2: zero-shot question generation

User input  $x$ .

Generated output so far  $y_{\leq t}$ .

Given the above passage, ask a question to which the answer is the term/entity/phrase " $z$ ".

# FLARE (Jiang et al. 2023)

---

## 2. Explicit query by question generation

What university did Joe Biden attend?

What degree did Joe Biden earn?

RAG

```
graph TD; Q1[What university did Joe Biden attend?]; Q2[What degree did Joe Biden earn?]; RAG[RAG]; Q1 --> RAG; Q2 --> RAG; RAG --> A[He graduated from the University of Delaware in 1965 with a Bachelor of Arts in history and political science];
```

$s_2$  He graduated from the University of Delaware in 1965 with a Bachelor of Arts in history and political science

# Entities as Experts (Fevry et al. 2020)

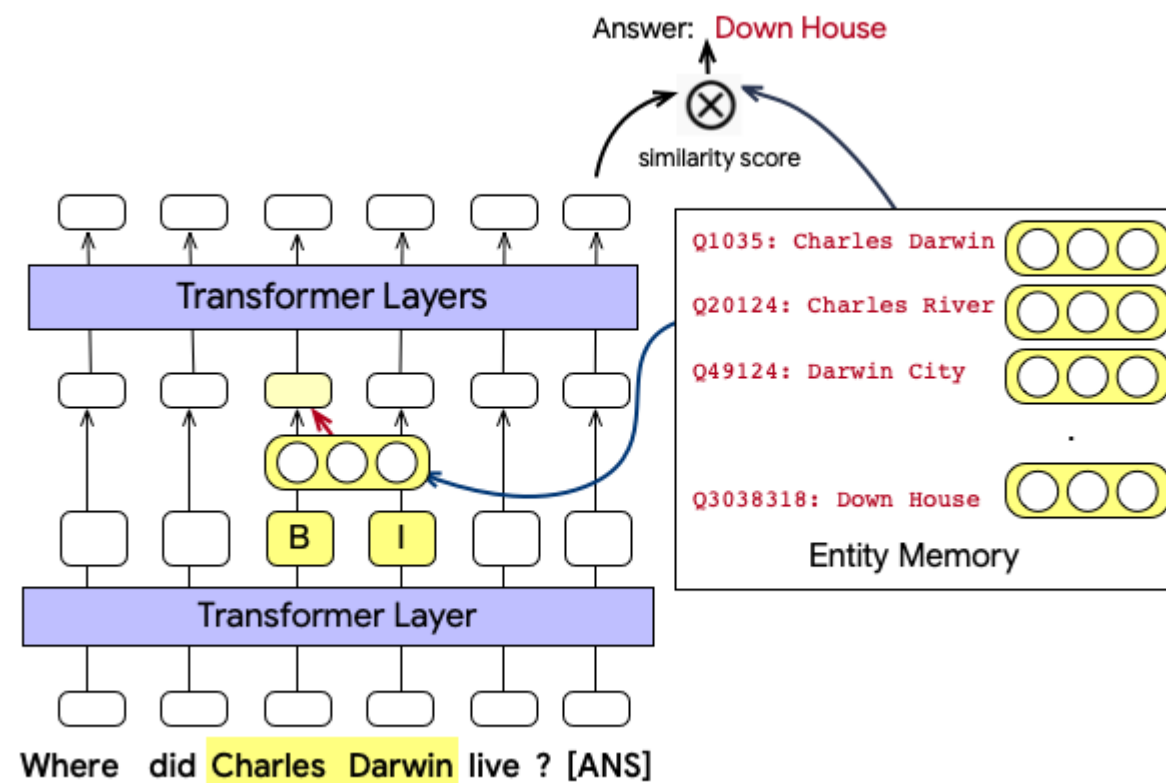
- Entities as Experts (EAE)

Entities as Experts (Fevry et al. 2020), Mention Memory (de Jong et al. 2022)

Entities or entity mentions

Intermediate layers

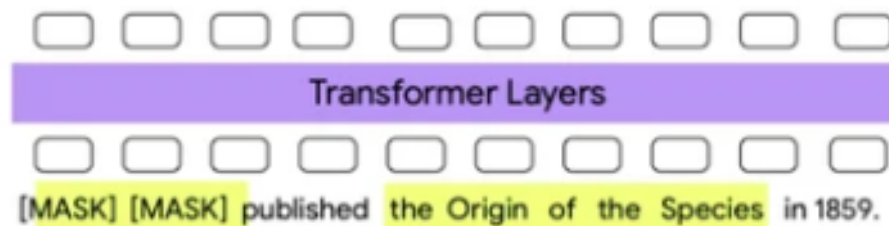
Every entity mentions





# Entities as Experts (Fevry et al. 2020)

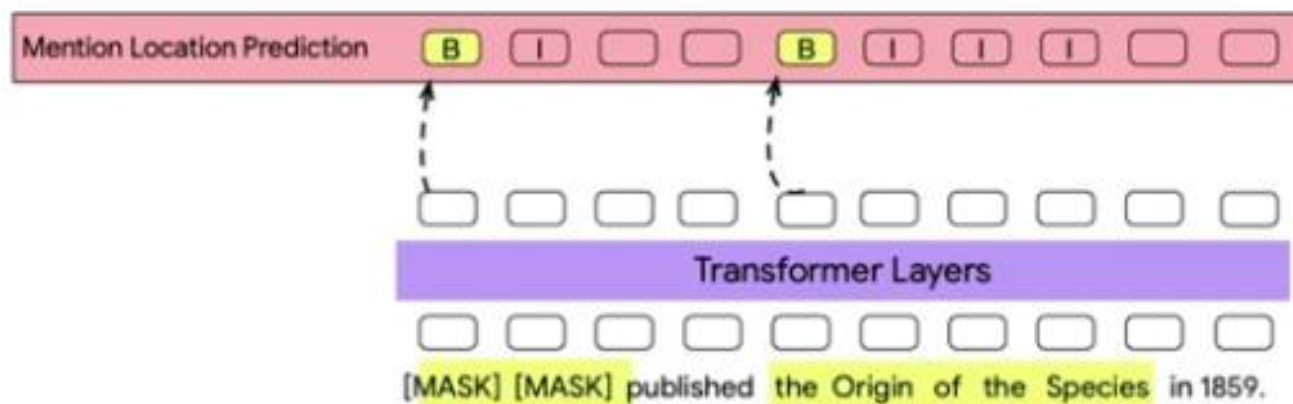
---



"Charles Darwin published the Origin of the Species in 1859."

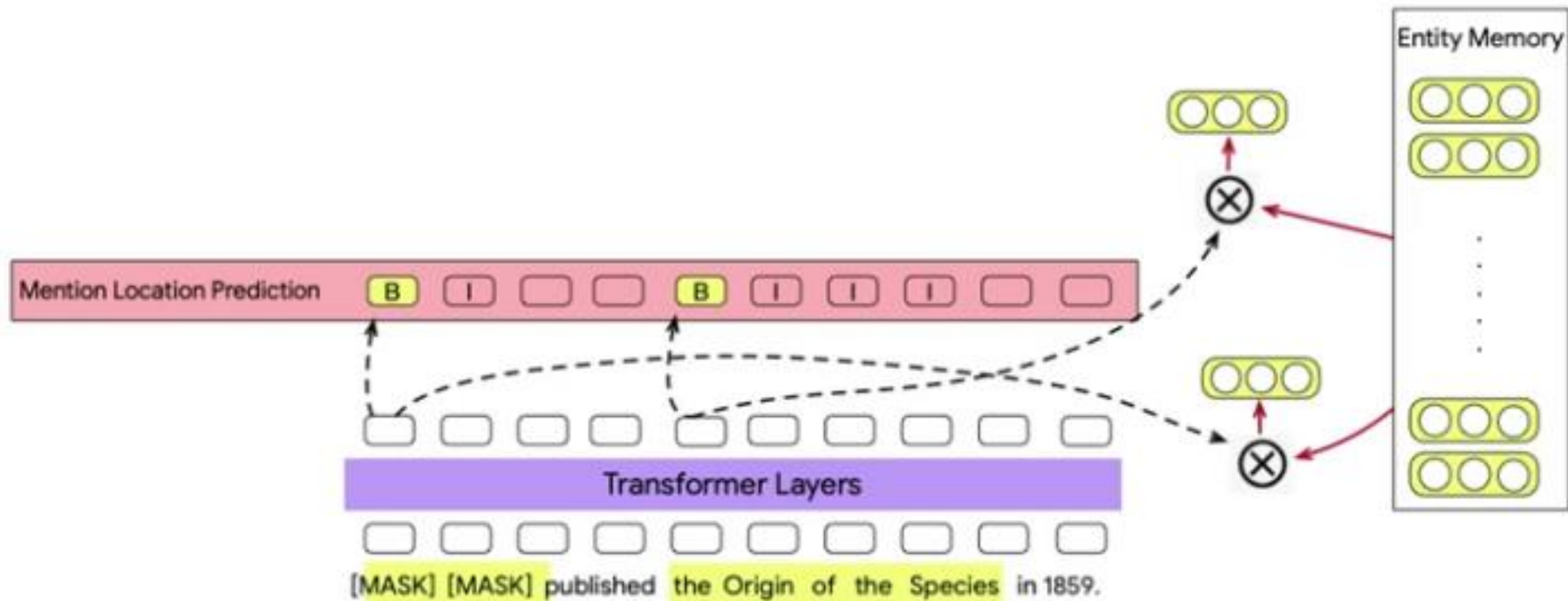
# Entities as Experts (Fevry et al. 2020)

---



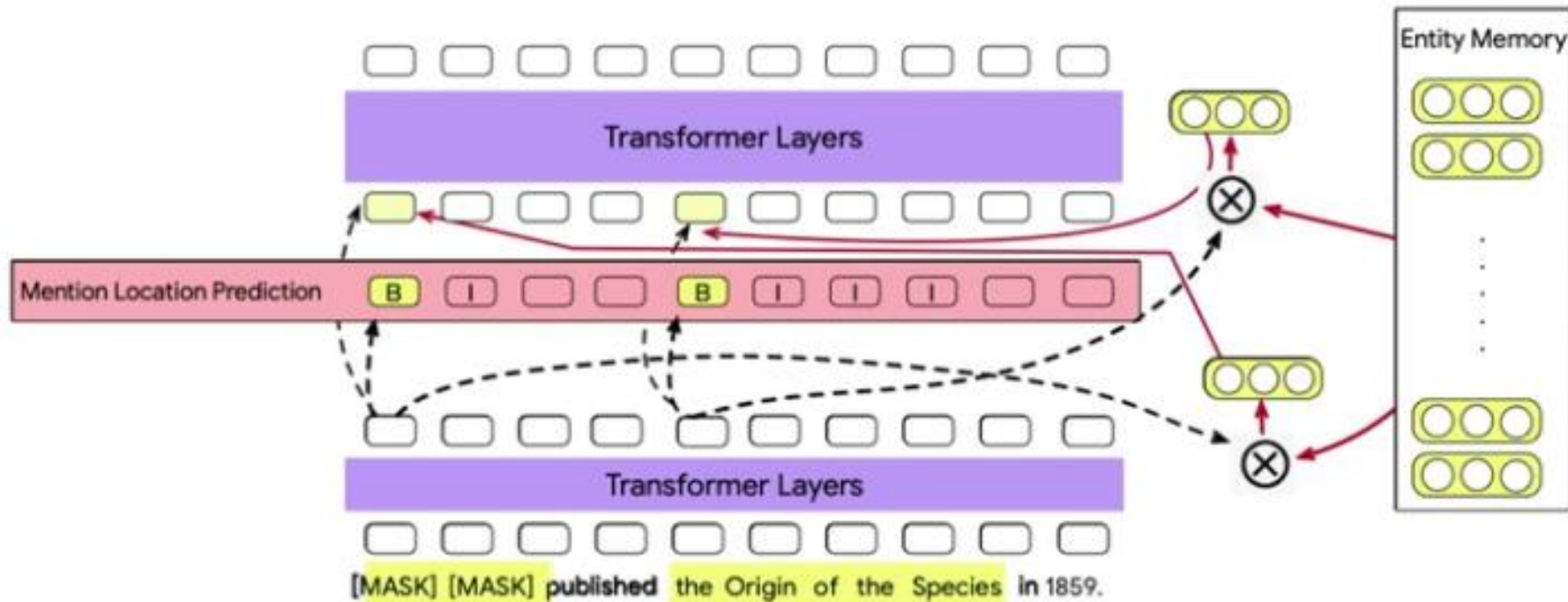
"Charles Darwin published the Origin of the Species in 1859."

# Entities as Experts (Fevry et al. 2020)



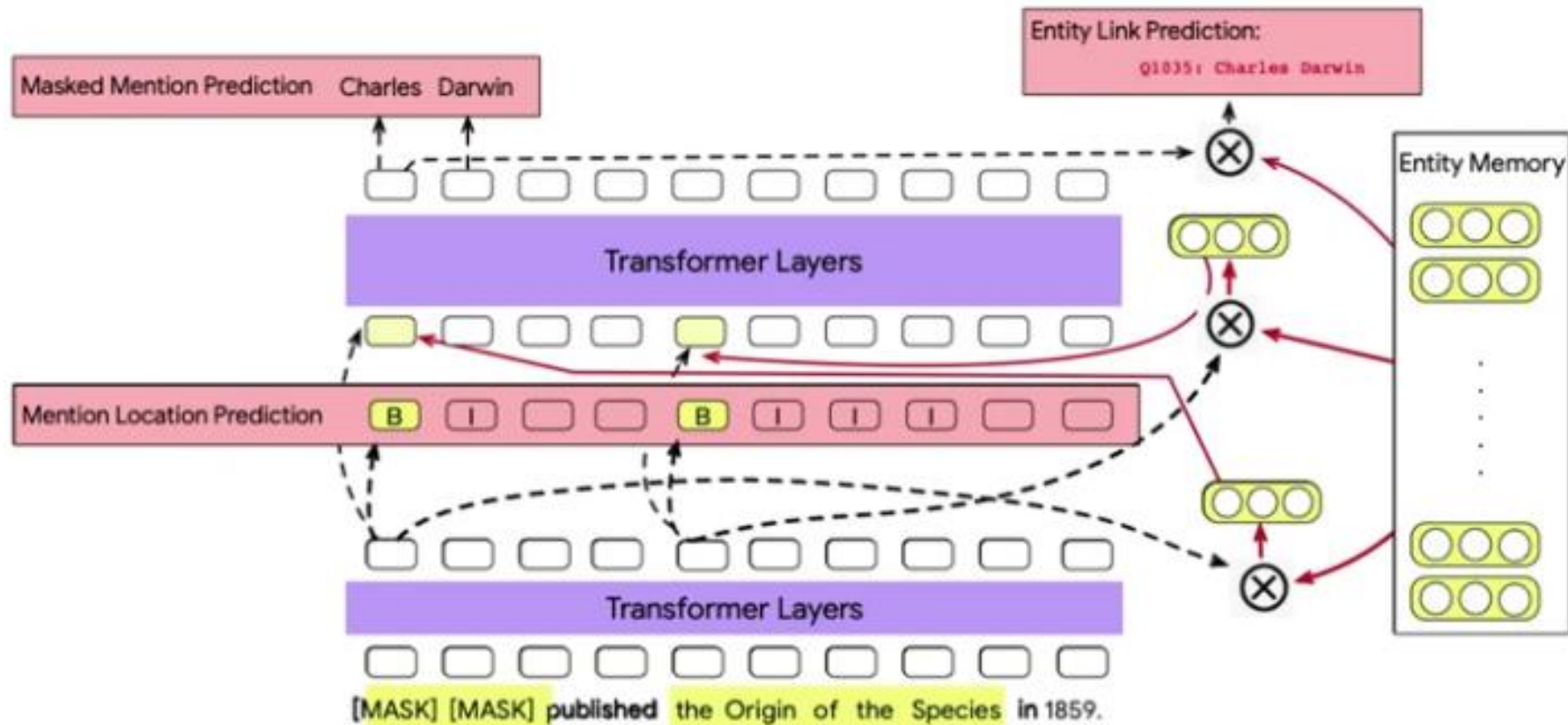
"Charles Darwin published the Origin of the Species in 1859."

# Entities as Experts (Fevry et al. 2020)



"Charles Darwin published the Origin of the Species in 1859."

# Entities as Experts (Fevry et al. 2020)



"Charles Darwin published the Origin of the Species in 1859."

# Retrieval for long-range LM (Wu et al. 2022)

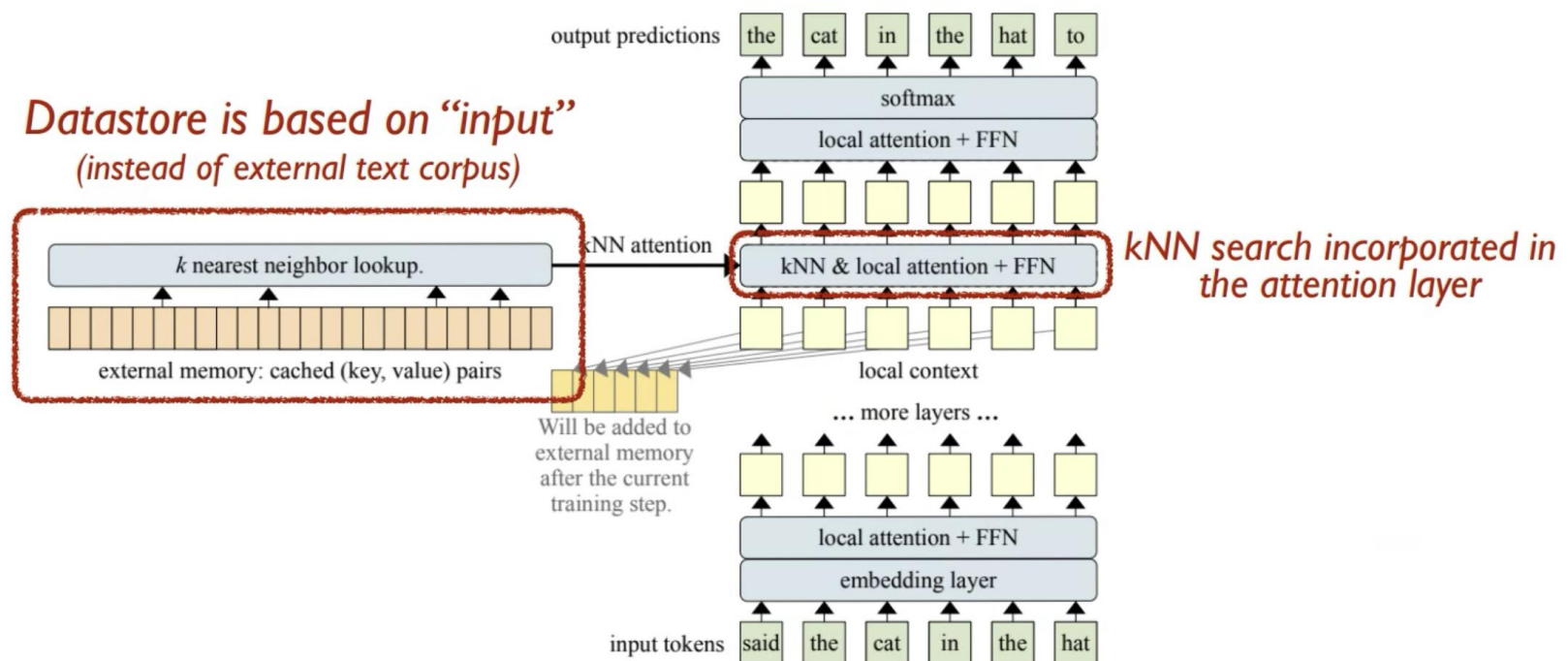
- It is hard to process long input (entire book) because transformer have a sequence length limit

Wu et al. 2022, Bertsch et al. 2023,  
Rubin & Berant. 2023

Text chunks *from*  
*the input*

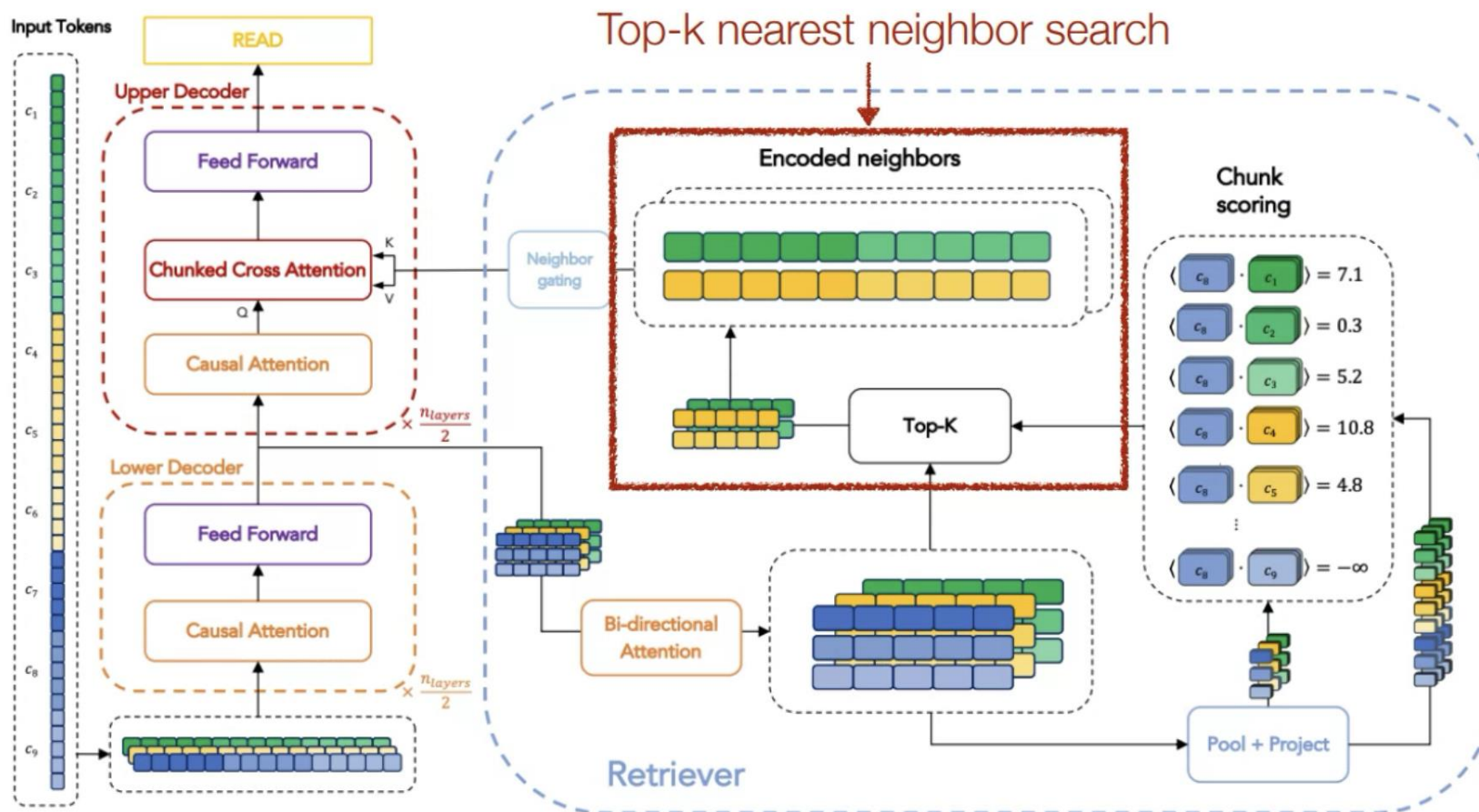
Intermediate layers

Once or every n  
tokens



# Retrieval for long-range LM (Rubin & Brent et al. 2023)

- This is from the RETRO



# Thank You

---

감사합니다.