

RETRO

Improving Language Models by Retrieving from Trillions of Tokens
(Sebastian Borgeaud et al. 2022)

Pseudo-Lab | 신서현

Background

Parametric

모델 파라미터에 모든 정보 저장
학습 데이터 → 모델 내부화된 기억

장점: 모델이 정보 내재화
단점: 메모리 사용량 증가

Non-parametric

외부 데이터베이스 사용 (e.g., 키-값 쌍)
필요 시 정보 검색

장점: 실시간 데이터 업데이트 가능
단점: 검색 속도 중요

Background

Self-attention

동일한 입력 내에서 정보를 상호 참조

Cross-attention

두 개 이상의 서로 다른 소스 간 정보를 연결

RETRO : Abstract Introduction

autoregressive language modeling

previous method more parameters , more memory

increasing the number of parameters in Transformer models

⇒ tremendous increase in training **energy cost**

⇒ resulted in a generation of **dense LLMs** with 100+ billion parameters

alternate path **RETRO**; Retrieval Enhanced TRansfOrmers

└ semi-parametric approach
: 모델이 **외부 지식 DB**에 접근해 필요한 지식을 활용

RETRO : 2.1. Training dataset

Introduction

for both training and retrieval data
use a multi-lingual version of **MassiveText** (Rae et al., 2021)
: multiple sources, multi-lingual

tokenize the dataset by SentencePiece
(128,000 tokens)

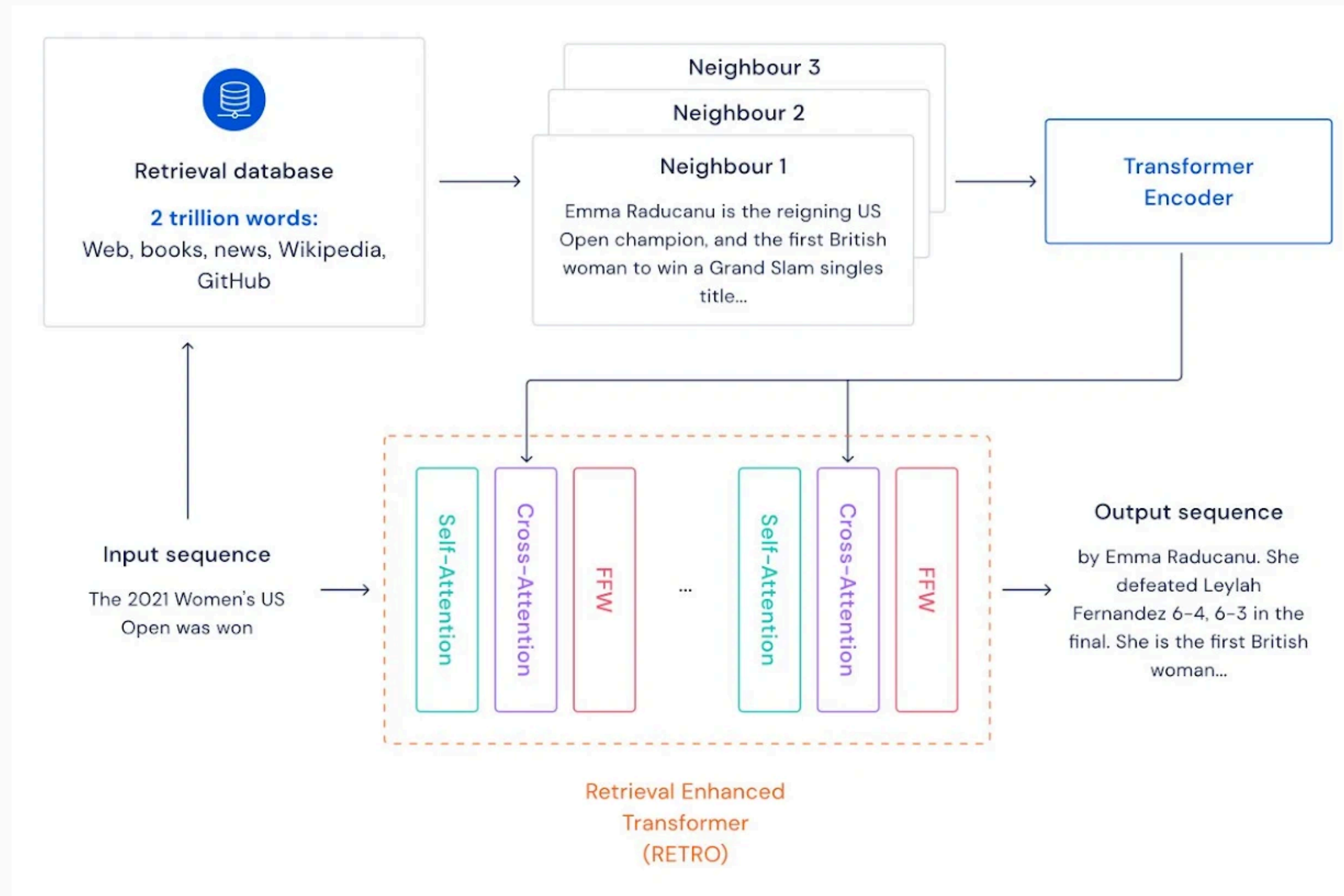
training retrieval : 600B tokens
evaluation retrieval : 1.75T tokens

To limit test set leakage - **MinHash**
compute the 13-gram Jaccard similarity
(train-test documents) using the MinHash scheme
remove all training documents with high similarity (0.8 or higher)
to a validation or test set document

Source	Language	Token count (M)	Documents	Sampling weight
Web	En	483,002	604,938,816	0.314
	Ru	103,954	93,004,882	0.033
	Es	95,762	126,893,286	0.033
	Zh	95,152	121,813,451	0.033
	Fr	59,450	76,612,205	0.033
	De	57,546	77,242,640	0.033
	Pt	44,561	62,524,362	0.033
	It	35,255	42,565,093	0.033
	Sw	2,246	1,971,234	0.0044
	Ur	631	455,429	0.0011
Books	En	3,423,740	20,472,632	0.25
News	En	236,918	397,852,713	0.1
Wikipedia	En	3,977	6,267,214	0.0285
	De	2,155	3,307,818	0.003
	Fr	1,783	2,310,040	0.003
	Ru	1,411	2,767,039	0.003
	Es	1,270	2,885,013	0.003
	It	1,071	2,014,291	0.003
	Zh	927	1,654,772	0.003
	Pt	614	1,423,335	0.003
	Ur	61	344,811	0.0001
	Sw	15	58,090	0.0004
Github	-	374,952	142,881,832	0.05
Total	-	5,026,463	1,792,260,998	1

deepmind blog : Figure 1

Introduction



<https://deepmind.google/discover/blog/improving-language-models-by-retrieving-from-trillions-of-tokens>

RETRO : 2.2. Retrieval-enhanced autoregressive token models

Method

retrieval-enhanced autoregressive language model

use a **chunked cross-attention module** to incorporate the retrieved text
time complexity linear in the amount of retrieved data

**n-token-long example X
into a sequence of l
chunks of size m = n/l**
n=2048, m=64

chunked Trillion 단위의 database에서 retrieval
저장 공간과 계산을 효율화하기 위해 chunked token을 기본단위로 사용

$$C_1 \triangleq (x_1, \dots, x_m), \dots, C_l \triangleq (x_{n-m+1}, \dots, x_n) \in \mathbb{V}^m$$

cross-attention 각 청크 C에 대해 Retrieval 데이터를 검색하고, 이를 Cross-Attention으로 통합

Method

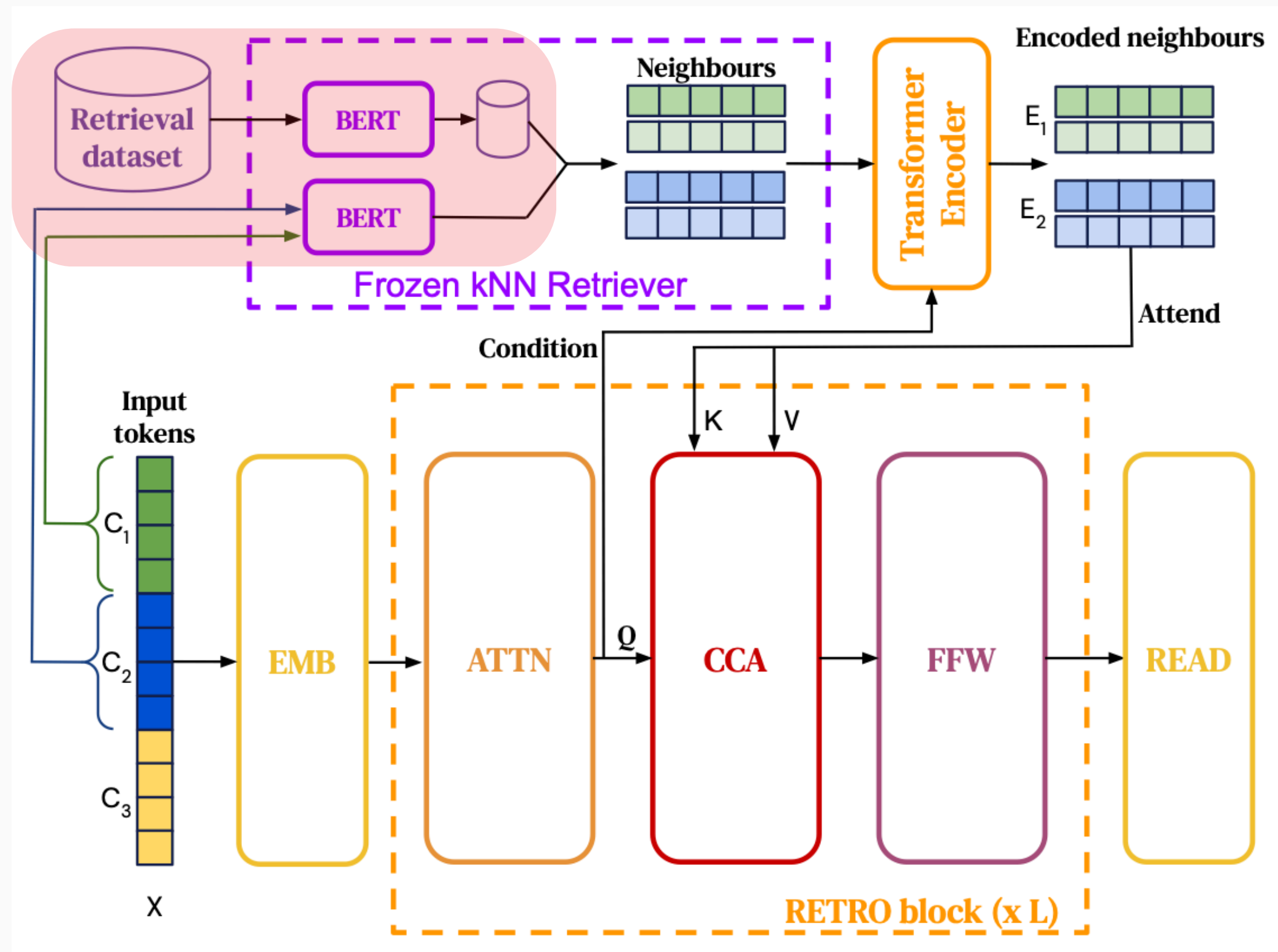
retrieval database

key : frozen BERT embeddings

value : raw chunks of text tokens

key	value
key 1	key에 해당하는 raw text
	next text chunk

고정된 모델(frozen BERT) 사용
→ 학습 중 전체 데이터베이스에 대해
임베딩을 다시 계산할 필요 없음.



RETRO : 2.3. Nearest neighbour retrieval Method

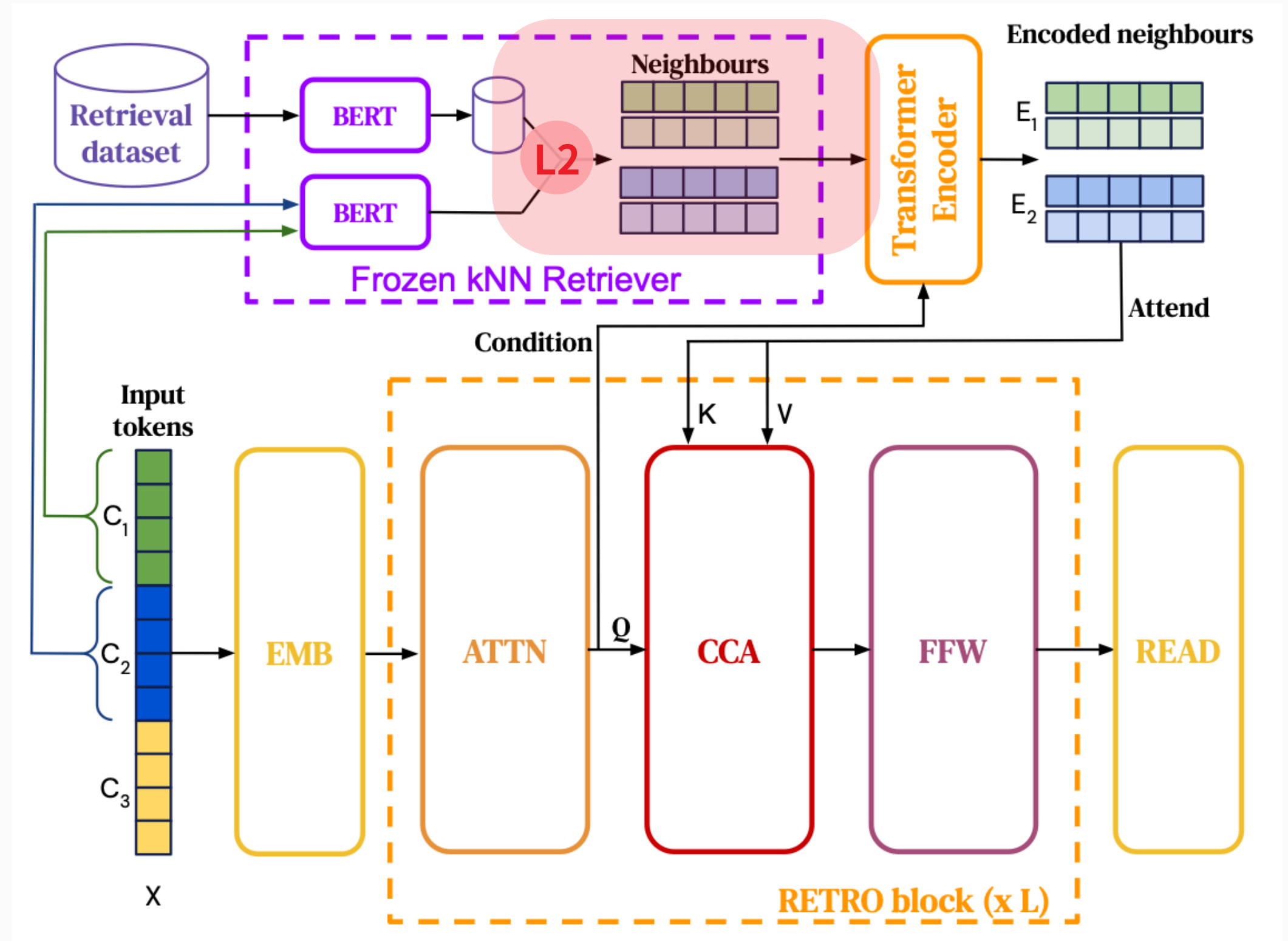
input chunk

L2 distance

$$d(C, N) = ||\text{BERT}(C) - \text{BERT}(N)||_2^2$$

nearest neighbours

Transformer encoder



Method

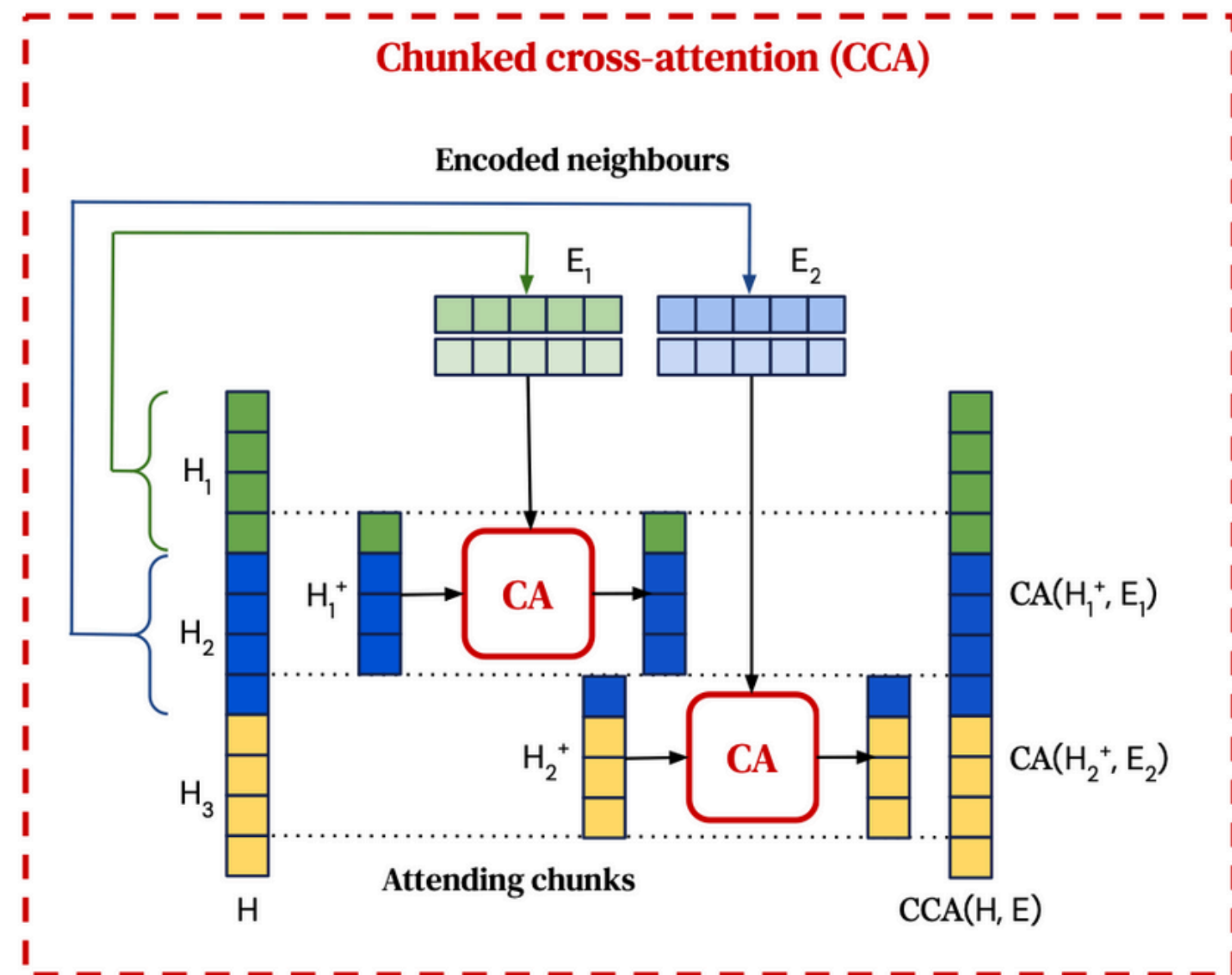
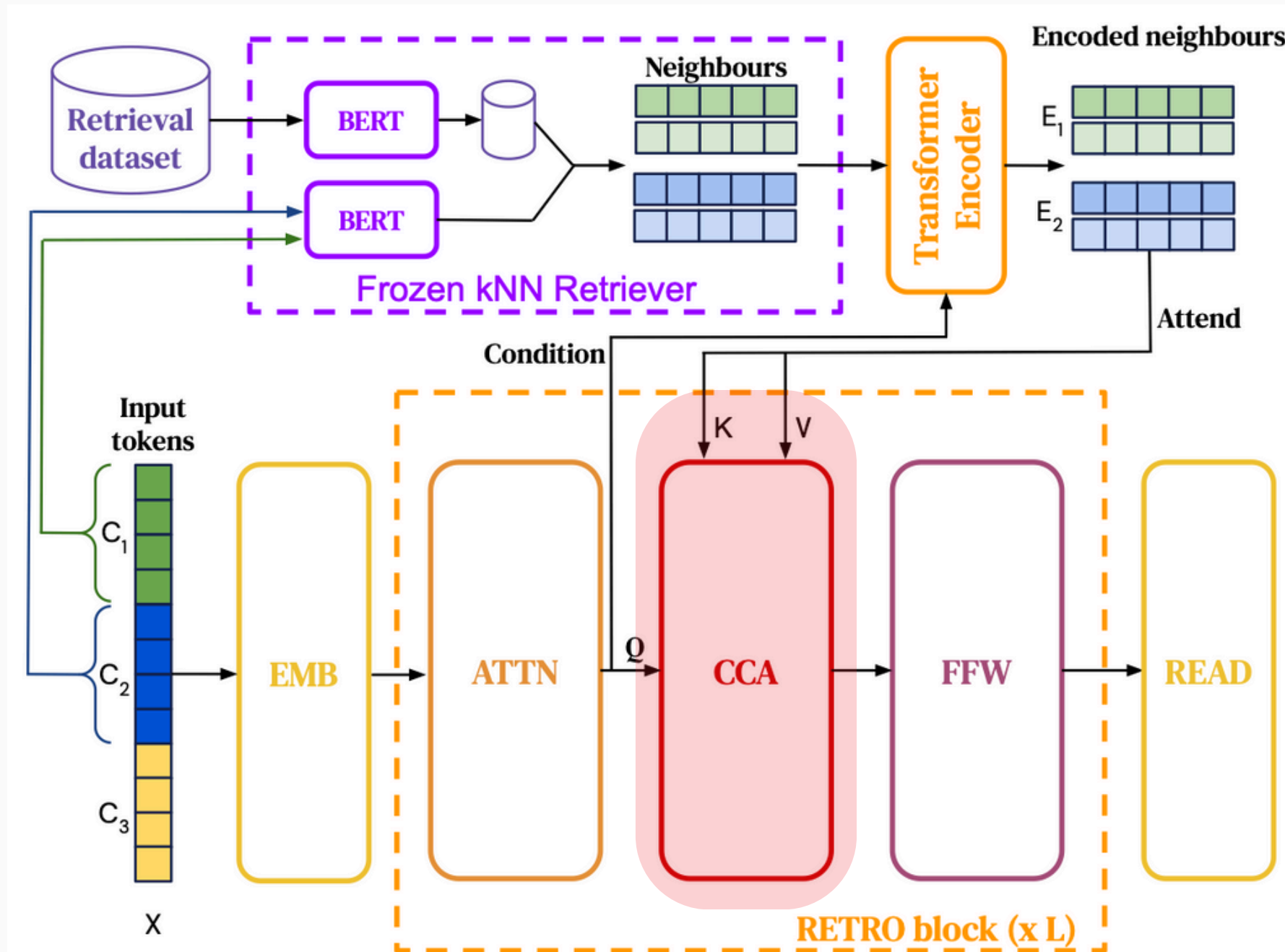
LM architecture

$$\text{LM}(H) \triangleq \text{FFW}(\text{ATTN}(H))$$

RETRO model architecture

$$\text{RETRO}(H, E) \triangleq \text{FFW}(\text{CCA}(\text{ATTN}(H), E))$$

Method



Method

Split decoder input into $l - 1$ attending chunks
maintain autoregressivity
(dependency propagation)

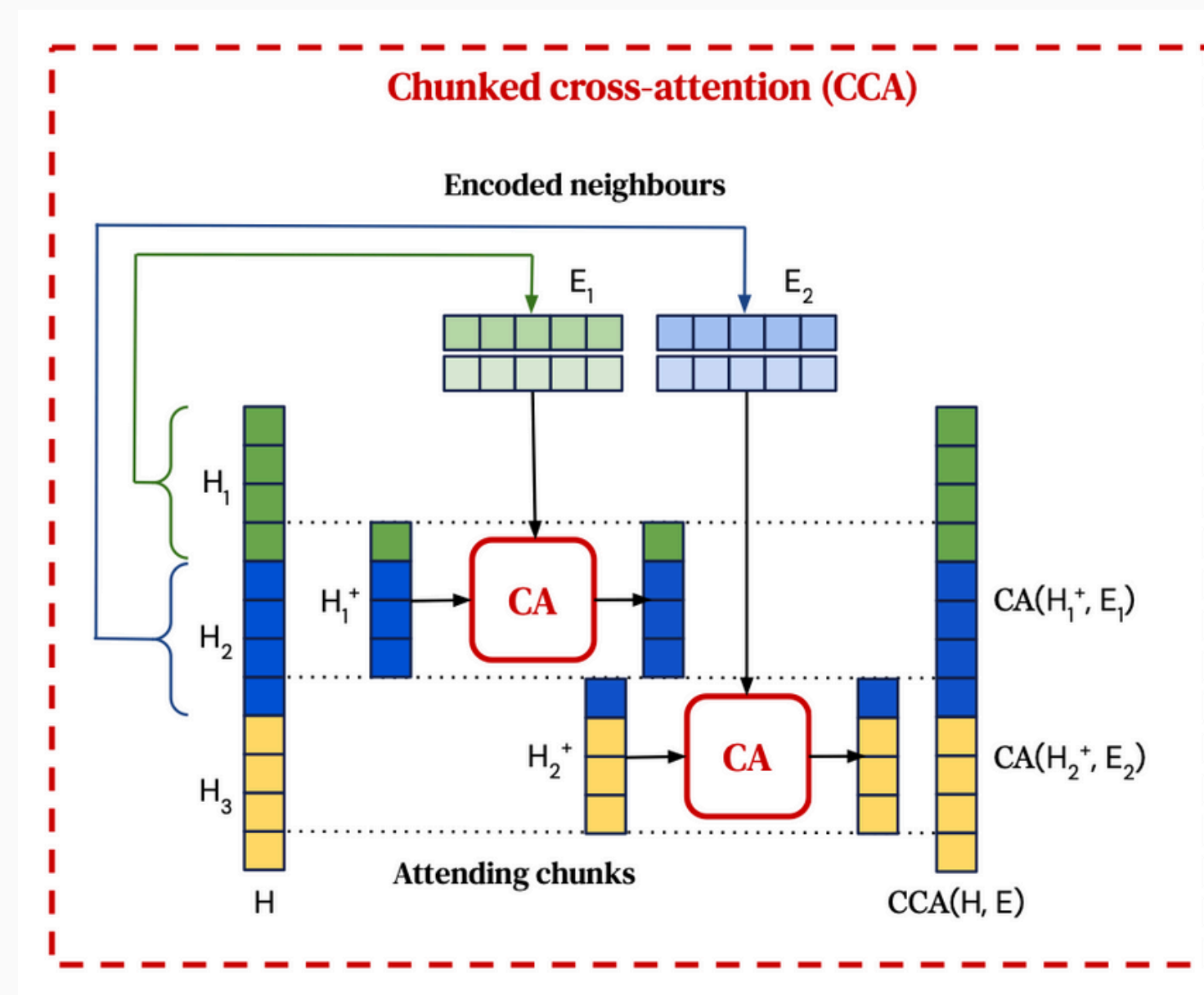
Chunked Cross-Attention

$$\left(H_u^+ \triangleq (h_{u \ m+i-1})_{i \in [1, m]} \in \mathbb{R}^{m \times d} \right)_{u \in [1, l-1]}$$

최종 출력 활성화 함수

$$CCA(H, E)_{u \ m+i-1} \triangleq CA(h_{u \ m+i-1}, E_u).$$

$$CA(h, Y) \triangleq \text{softmax}(Y K Q^T h) Y V,$$



H = input 활성화

E = retrieved neighbor

Q, K, V 는 각각 쿼리(query), 키(key), 값(value)

$Y \cdot K$ = 검색된 데이터 Y 의 키를 계산

Results

Datasets evaluated on

- ✓ C4
- ✓ Wikitext103
- ✓ Curation Corpus
- ✓ Lambada,
- ✓ The Pile

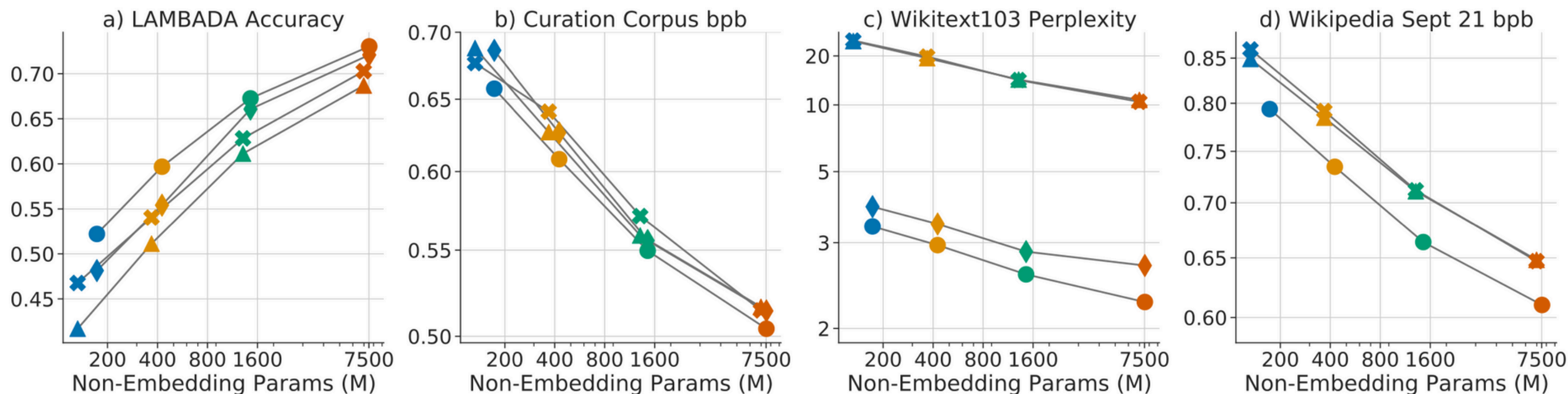
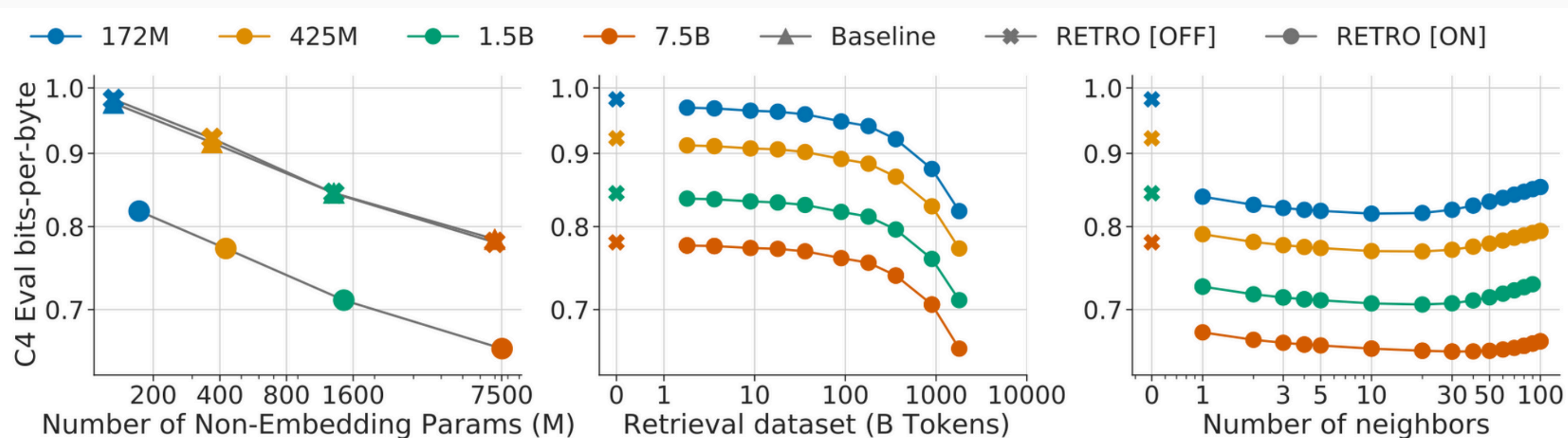
Metric : bits-per-byte(bpb)

$$\text{perplexity} = 2^{\frac{|\text{bytes}|}{|\text{words}|} \cdot \text{bpb}}$$

tokenizer와 무관하게 모델의 언어
모델링 성능을 측정

BPB가 낮을수록 모델 성능이 좋음

Results

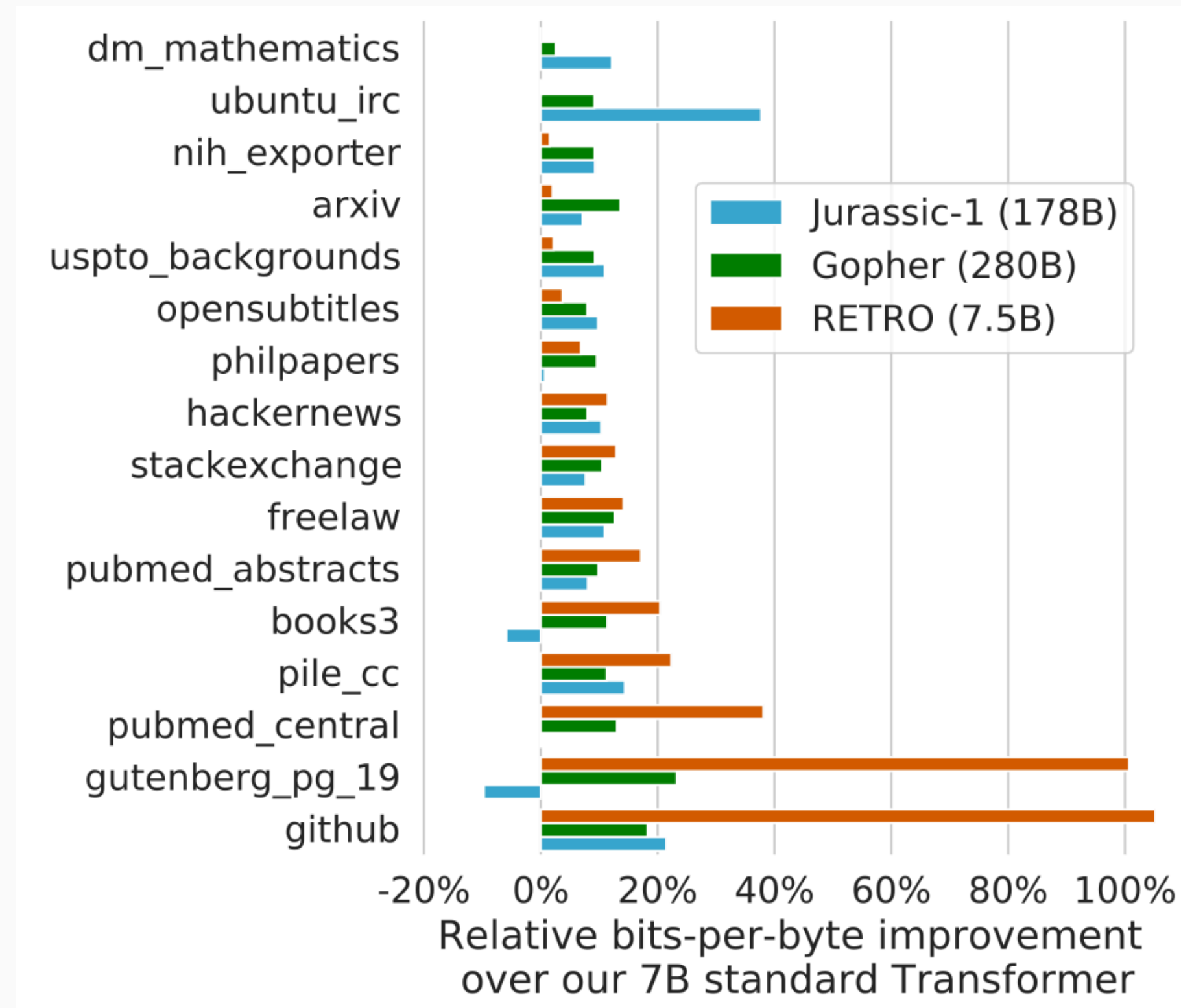


Results

The Pile dataset

Jurassic-1 : 178B

Gopher : 280B



Results

Table 16. **Daniel Radcliffe**, from Wikitext103Valid, retrieval data from c4. The chunks C_2 and C_3 are almost entirely retrieved from neighbours $[N_1, F_1]$ and $[N_2, F_2]$ respectively, up to formatting differences, which dramatically reduces the loss for these tokens. This example illustrates that when training data leaks into evaluation sets despite deduplication, our RETRO model can directly exploit this leakage.

<div>C_u colored by loss difference $L_{\text{RETRO}[\text{Off}]} - L_{\text{RETRO}} \leq -0.5, = 0, \geq 0.5$</div>	<div>C_u colored by LCP with $\text{RET}(C_u - 1)$ LCP = 0, 1, 2, 3, 4, ≥ 5</div>	<div>$[N_u^1, F_u^1]$ colored by LCP with C_{u+1} LCP = 0, 1, 2, 3, 4, ≥ 5</div>	<div>$[N_u^2, F_u^2]$ colored by LCP with C_{u+1} LCP = 0, 1, 2, 3, 4, ≥ 5</div>
<div>= Daniel Radcliffe =Daniel Jacob Radcliffe (born 23 July 1989) is an English actor who rose to prominence as the title character in the Harry Potter film series. He made his acting debut at 10 years of age in BBC One's 1999 television film David Copperfield, followed by his cinematic debut</div>	<div>= Daniel Radcliffe = Daniel Jacob Radcliffe (born 23 July 1989) is an English actor who rose to prominence as the title character in the Harry Potter film series. He made his acting debut at 10 years of age in BBC One's 1999 television film David Copperfield, followed by his cinematic debut</div>	<div>Daniel Jacob Radcliffe (born 23 July 1989) is an English actor who rose to prominence as the title character in the Harry Potter film series. He made his acting debut at 10 years of age in BBC One's 1999 television film David Copperfield, followed by his cinematic debut in 2001's The Tailor of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New York productions of Equus, and</div>	<div>Daniel Jacob Radcliffe (born 23 July 1989) is an English actor who rose to prominence as the title character in the Harry Potter film series. He made his acting debut at 10 years of age in BBC One's 1999 television movie David Copperfield, followed by his film debut in 2001's The Tailor of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New York productions of Equus, and in the</div>
<div>in 2001's The Tailor of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New</div>	<div>in 2001's The Tailor of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New</div>	<div>in 2001's The Tailor of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New York productions of Equus, and in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities</div>	<div>of Panama. At age 11, he was cast as Harry Potter in the first Harry Potter film, and starred in the series for 10 years until the release of the eighth and final film in 2011. Radcliffe began to branch out to stage acting in 2007, starring in the London and New York productions of Equus, and in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities, including Demelza House Children's</div>
<div>York productions of Equus, and in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities,</div>	<div>York productions of Equus, and in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities,</div>	<div>York productions of Equus, and in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities, including Demelza House Children's Hospice and The Trevor Project. He also made public service announcements for the latter. In 2011, he was awarded the Trevor Project's "Hero Award." Sources disagree about Radcliffe's personal wealth; he was reported to have earned £1 million for the first Harry Potter</div>	<div>in the 2011 Broadway revival of the musical How to Succeed in Business Without Really Trying. He starred in the 2012 horror film The Woman in Black, and played beat poet Allen Ginsberg in the 2013 independent film Kill Your Darlings. He has contributed to many charities, including Demelza House Children's Hospice and The Trevor Project. He also made public service announcements for the latter. In 2011, he was awarded the Trevor Project's "Hero Award."</div>
<div>including Demelza House Children's Hospice and The Trevor Project for suicide prevention among LGBTQ youth, which gave him its Hero Award in 2011. = Early life = Radcliffe was born in West London, England. He is the only child of Alan George Radcliffe, a literary agent, and</div>	<div>including Demelza House Children's Hospice and The Trevor Project for suicide prevention among LGBTQ youth, which gave him its Hero Award in 2011. = Early life = Radcliffe was born in West London, England. He is the only child of Alan George Radcliffe, a literary agent, and</div>		

2024.11.18

감사합니다

Pseudo-Lab | 신서현