# Unlimiformer: Long-Range Transformers with Unlimited Length Input
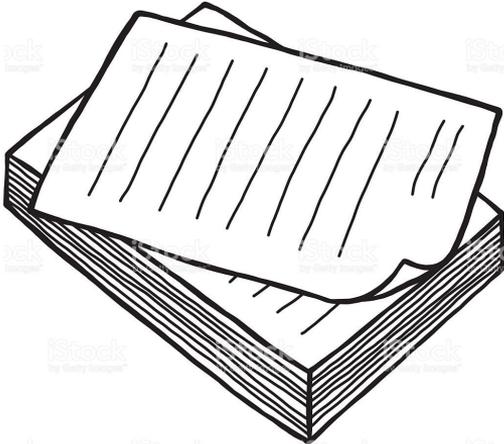
Amanda Bertsch     Uri Alon     Graham Neubig     Matt Gormley

# The inputs we'd like to work with keep getting bigger…

1,000 tokens

10,000 tokens

100,000 tokens

# …and our models don't scale that well

- Sparse attention
  - Pretraining is hugely expensive
  - Fixed maximum length
- Hierarchical summarization
  - Cascading errors
  - Can't see the big picture
- ???

100,000 tokens

The **length** of the context window is fixed… what about the **content**?

# Retrieval-augmented generation
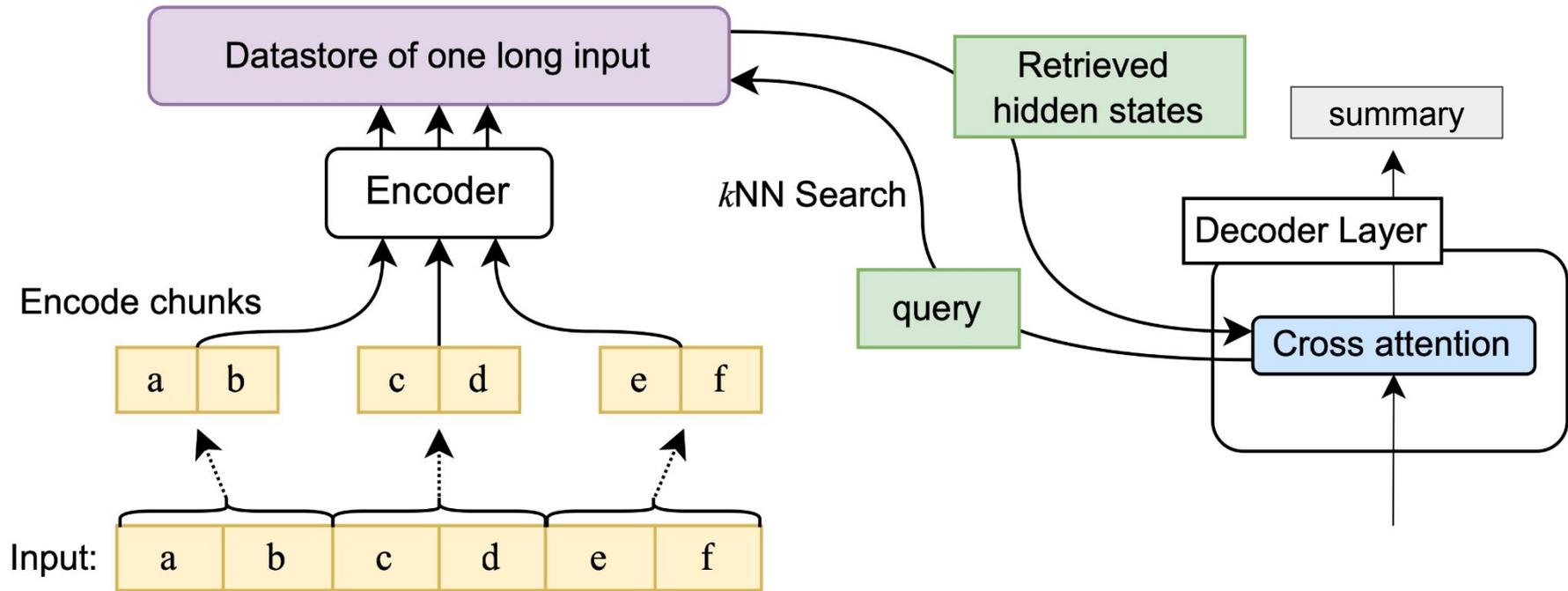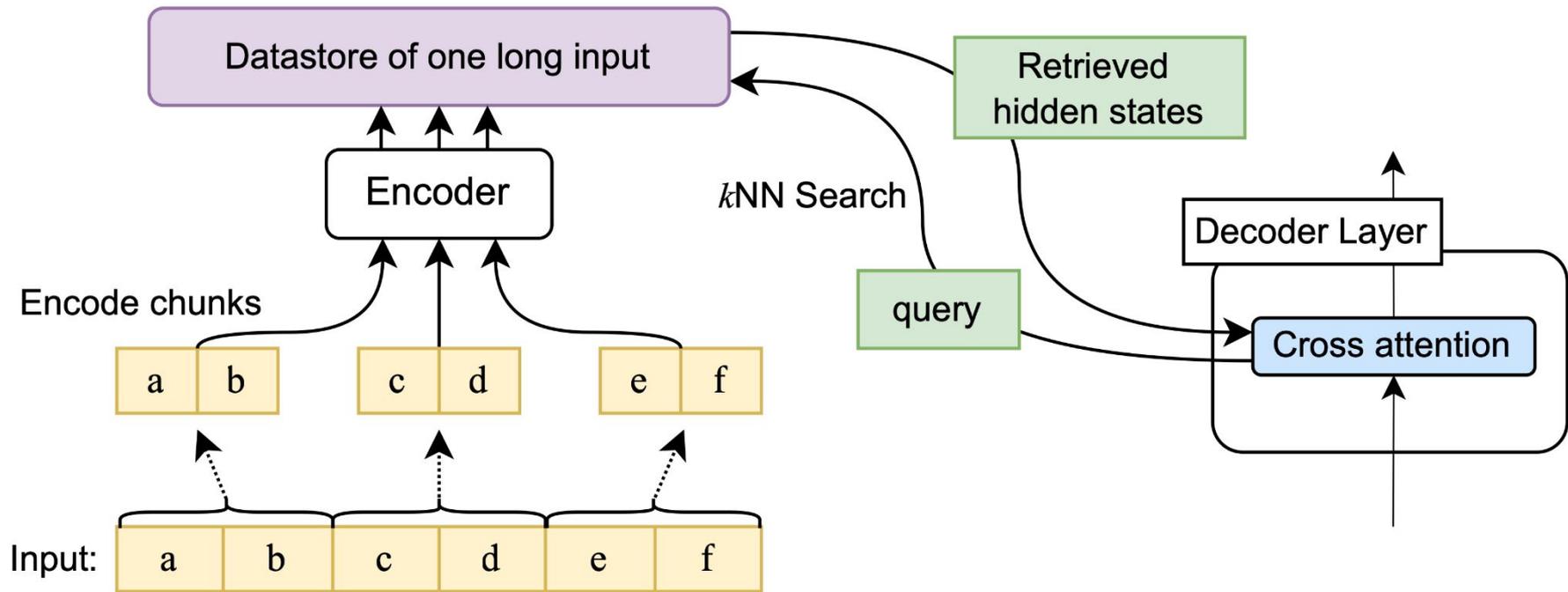
RETRO, Memorizing Transformers, etc:

- maintain a "base context" and augment with retrieved text
  - Unlimiformer has no "base context"
- add a layer (or a few layers) that cross attend to both external memory and the context
  - Unlimiformer cross attends only to external memory at every layer
- retrieve from set of relevant documents for QA or full pretraining corpus/recent examples for LM
  - Unlimiformer retrieves from the same long sequence
  - The datastore is static and unique for a single example

100,000 tokens
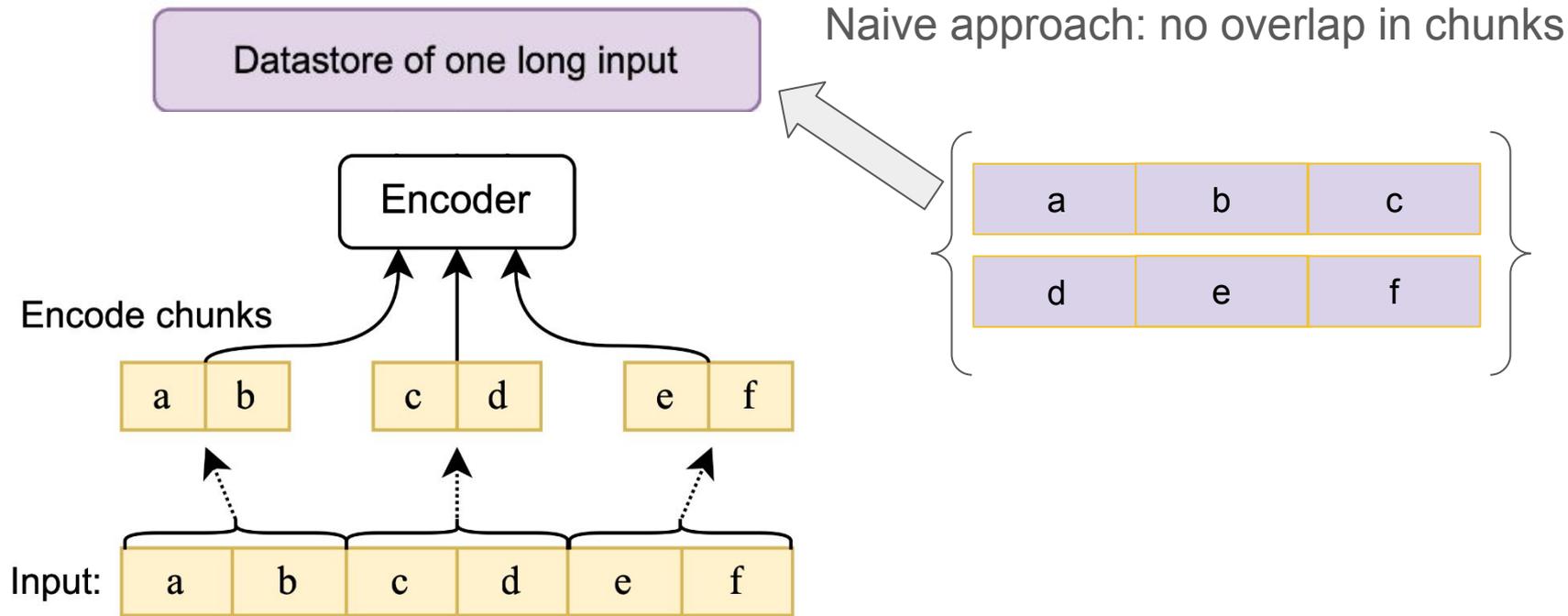
# Unlimiformer

# How do we do encoding?



Number of encoder passes: ⌈input len / encoder max len⌉

# How do we do encoding?

Datastore of one long input

Encoder

Encode chunks

| a | b | | c | d | | e | f |

Input:

| a | b | c | d | e | f |

Naive approach: no overlap in chunks

| a | b | c |

| d | e | f |

Number of encoder passes: ⌈input len / encoder max len⌉

7

# What about context?

embeddings with no left context:

| a | | d |
|---|---|---|

embeddings with left+right context:

| b | | e |
|---|---|---|

embeddings with no right context:

| c | | f |
|---|---|---|

# What about positional embeddings?

encoding:

| a | b | c |
|---|---|---|
| d | e | f |

positional embeddings:

| a | b | c | d | e | f |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 1 | 2 | 3 |

# How do we do encoding?

Datastore of one long input

Encoder

Encode chunks

a b    c d    e f

Input: a b c d e f

Overlapping chunks: all tokens in the middle of the input have left and right context!

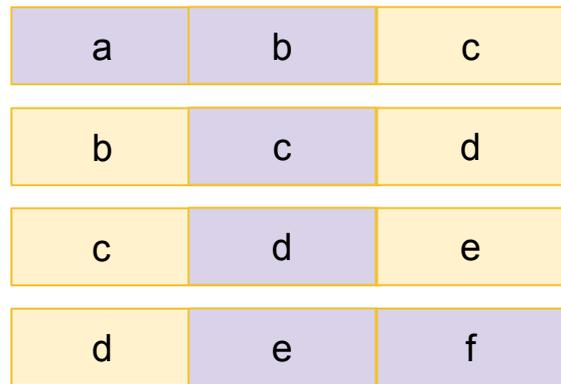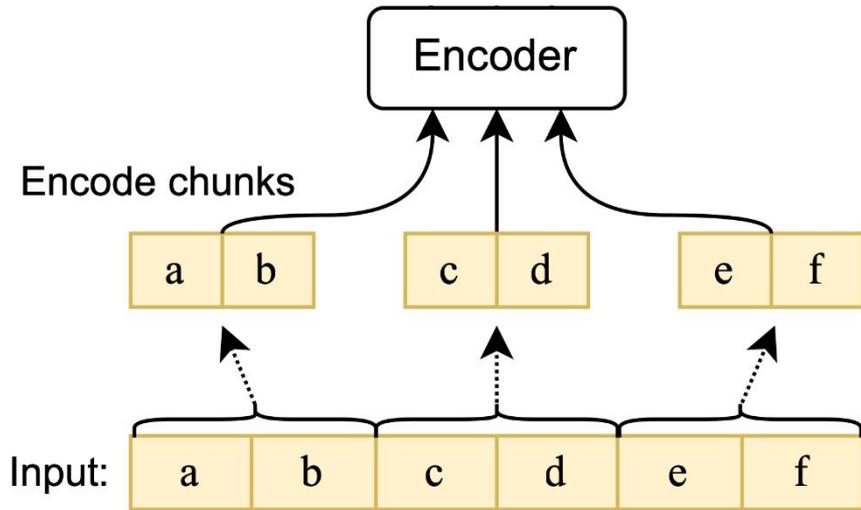| a | b | c |
|---|---|---|
| b | c | d |
| c | d | e |
| d | e | f |

in practice: use embeddings from middle half of window

Number of encoder passes: $\lceil$ input len / (0.5 * encoder max len) $\rceil$ - 1

10

# How do we do encoding?

Overlapping chunks: all tokens in the middle of the input have left and right context!

Encoder

Encode chunks

| a | b |   | c | d |   | e | f |

Input: | a | b | c | d | e | f |

| a | b | c |
| b | c | d |
| c | d | e |
| d | e | f |

in practice: use embeddings from middle half of window

Number of encoder passes: $\lceil$ input len / (0.5 * encoder max len) $\rceil$ - 1

# What about context?

embeddings with no left context:

a

embeddings with left+right context:

b    c

d    e

embeddings with no right context:

f

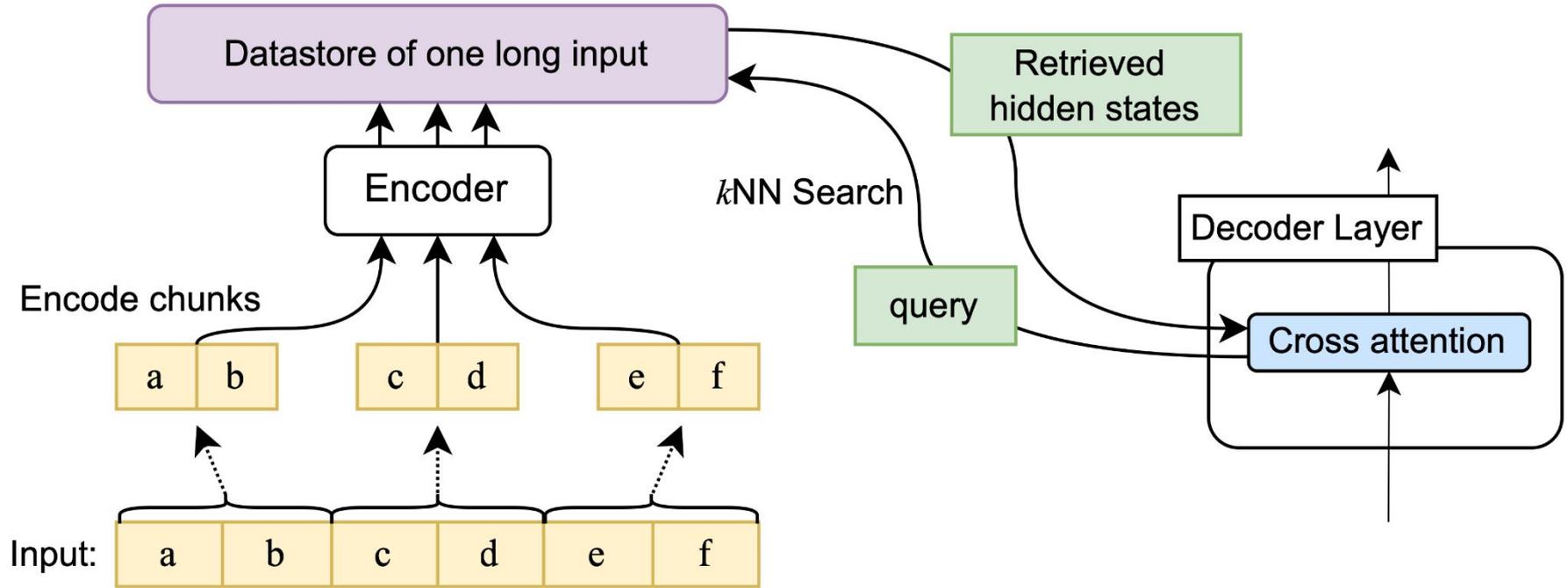# What about positional embeddings?

encoding:

| a | b | c |
|---|---|---|

| b | c | d |
|---|---|---|

| c | d | e |
|---|---|---|

| d | e | f |
|---|---|---|

positional embeddings:

| a | b | c | d | e | f |
|---|---|---|---|---|---|
| 1 | 2 | 2 | 2 | 2 | 3 |

also…

the decoder positional embeddings are unaffected!

# What is the datastore?

# How do we choose the context window?

# How do we choose the context window? cross-attention

decoder hidden state    encoder hidden state

layer specific
head specific

$$QK^T = (\boldsymbol{h}_d \boxed{W_q})\,(\boldsymbol{h}_e \boxed{W_k})^\top$$

Memorizing Transformers (Wu et al. ICLR'2022)
kept two datastores for each <layer,head> pair
Overall datastores: 2 X layers X heads

Project the query differently
for every layer/head

We can keep a **single** datastore of
the encoded hidden states

# How do we choose the context window?



Datastore of one long input $h_e$

Retrieved hidden states $h_r$

$k$NN Search

query

$\left( h_d W_q W_k^\top \right)$

Decoder Layer

Cross attention

$h_d$

Cross attention

$$Q = h_d W_q$$
$$K = h_r W_k$$
$$V = h_r W_v$$

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# How do we do efficient search?

Datastore of one long input



FAISS search:
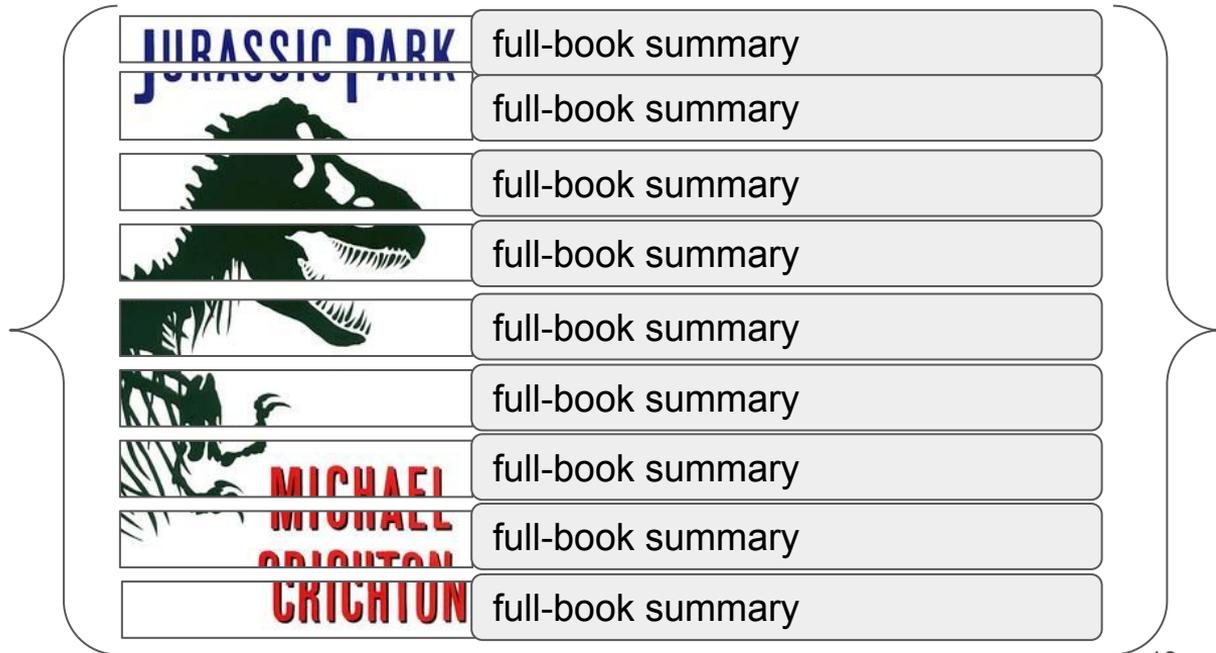- *Supports datastores on GPU, CPU, or disk*
- *Approximate*
- *Sublinear*

# Data augmentation (not Unlimiformer-specific!)

standard finetuning

JURASSIC PARK | full-book summary

chunked finetuning

JURASSIC PARK | full-book summary
| full-book summary
| full-book summary
| full-book summary
| full-book summary
| full-book summary
MICHAEL | full-book summary
| full-book summary
CRICHTON | full-book summary

# How do we train Unlimiformer?

Summarize:

Running example:
book summarization



117,645
words

# Normal training: truncating all inputs

During training:

During early stopping:

During test-time:

# Adding Unlimiformer after training

During training:
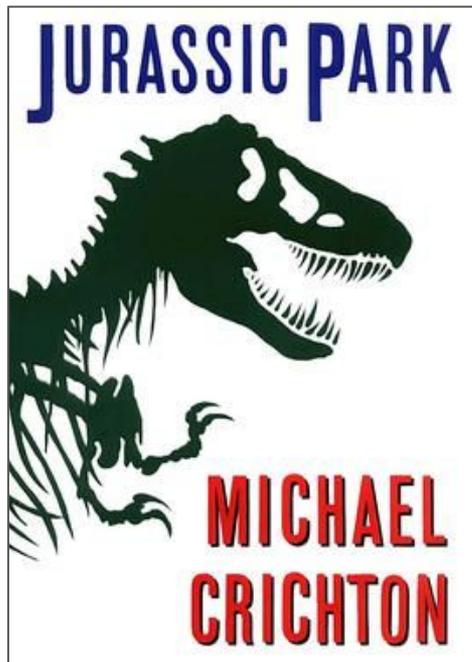
During early stopping:

During test-time:

22

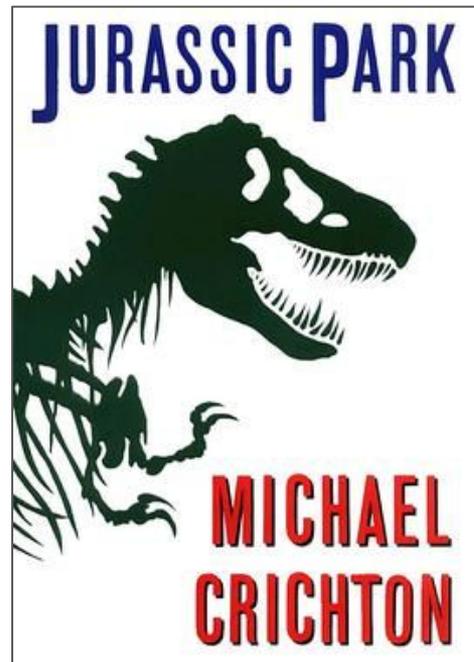# Low cost training: Unlimiformer-aware early stopping

During training:

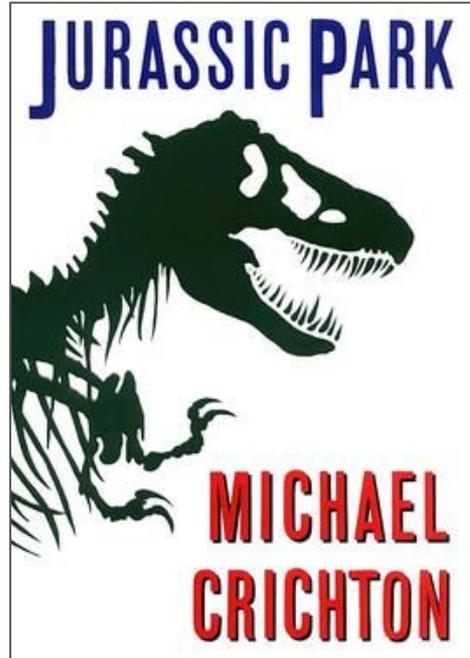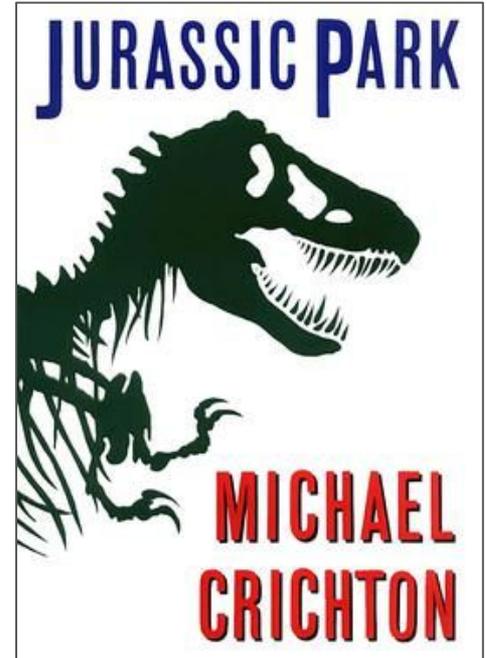During early stopping:

During test-time:

# Higher cost training methods
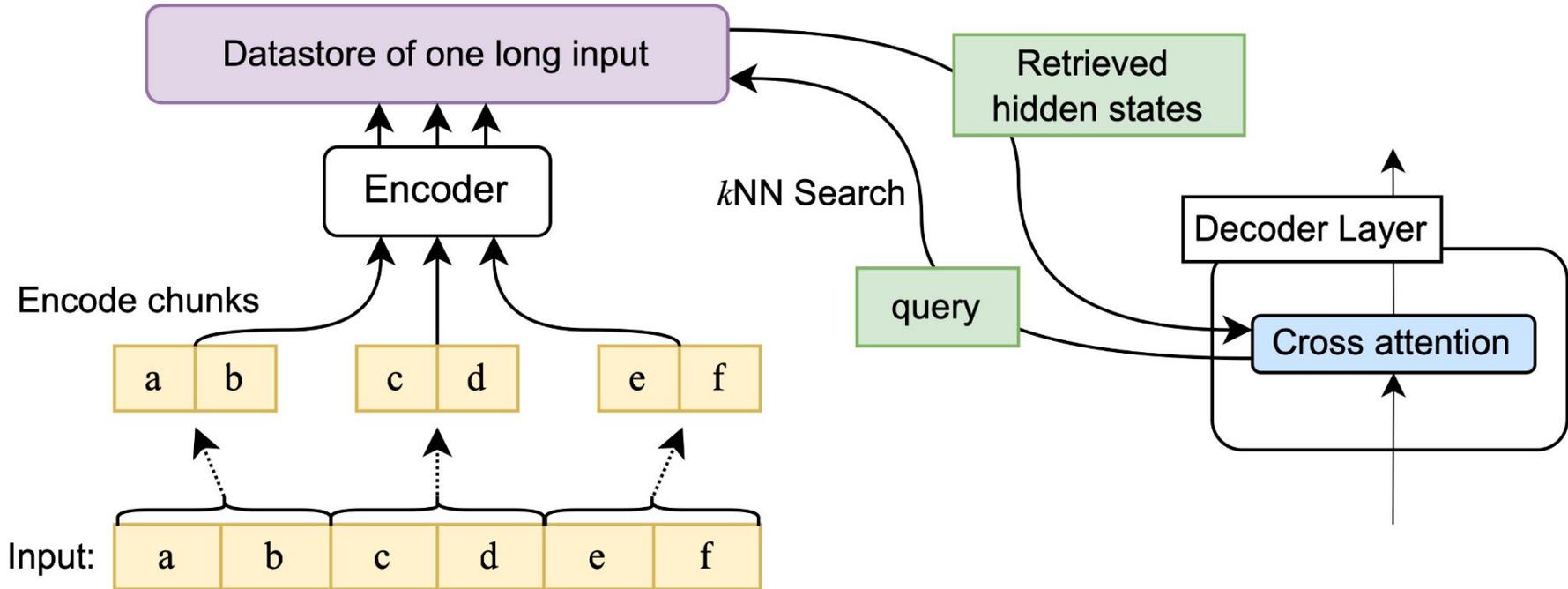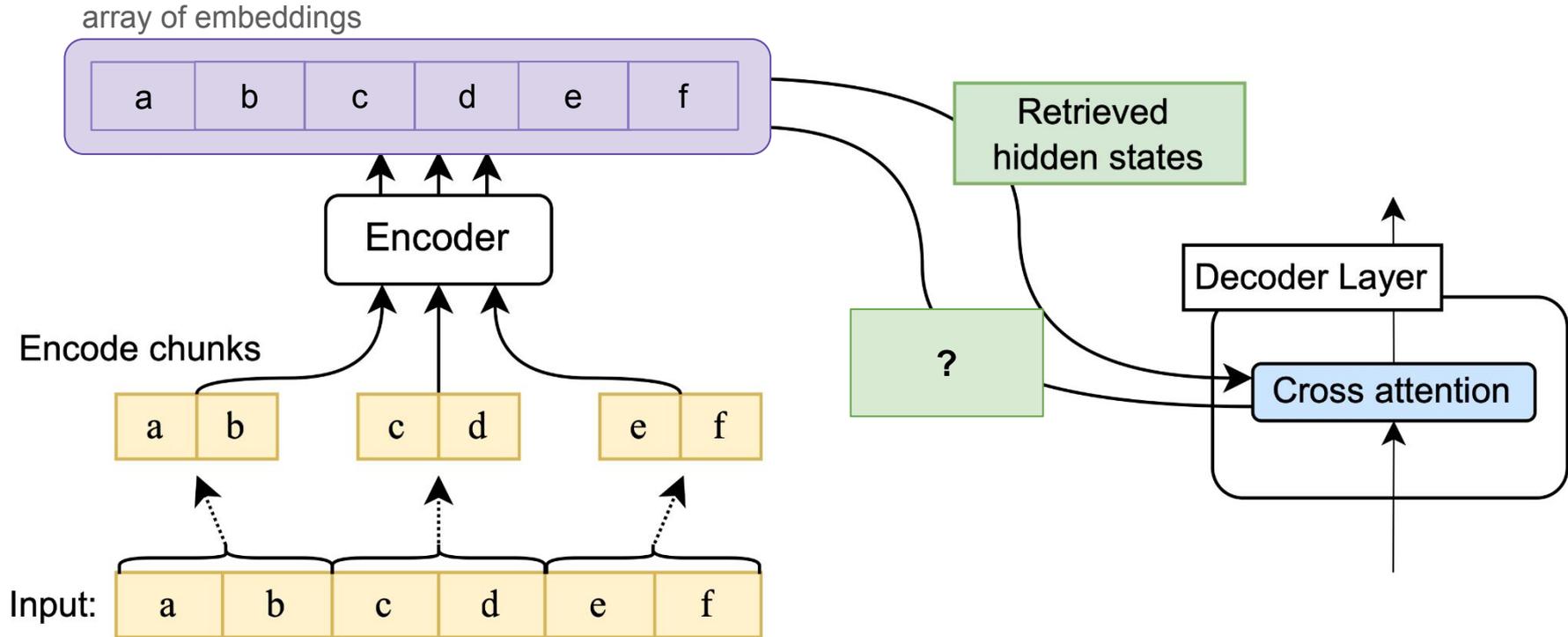
During training:

During early stopping:
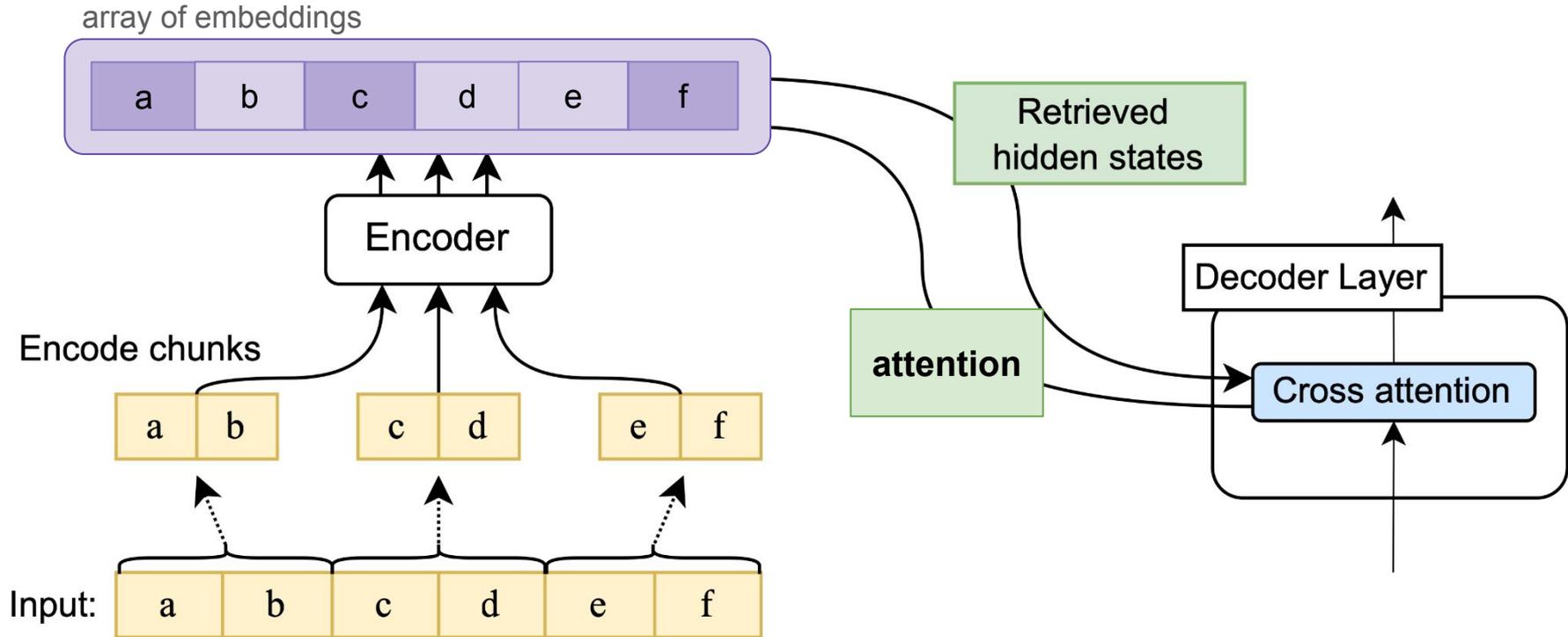
During test-time:



24

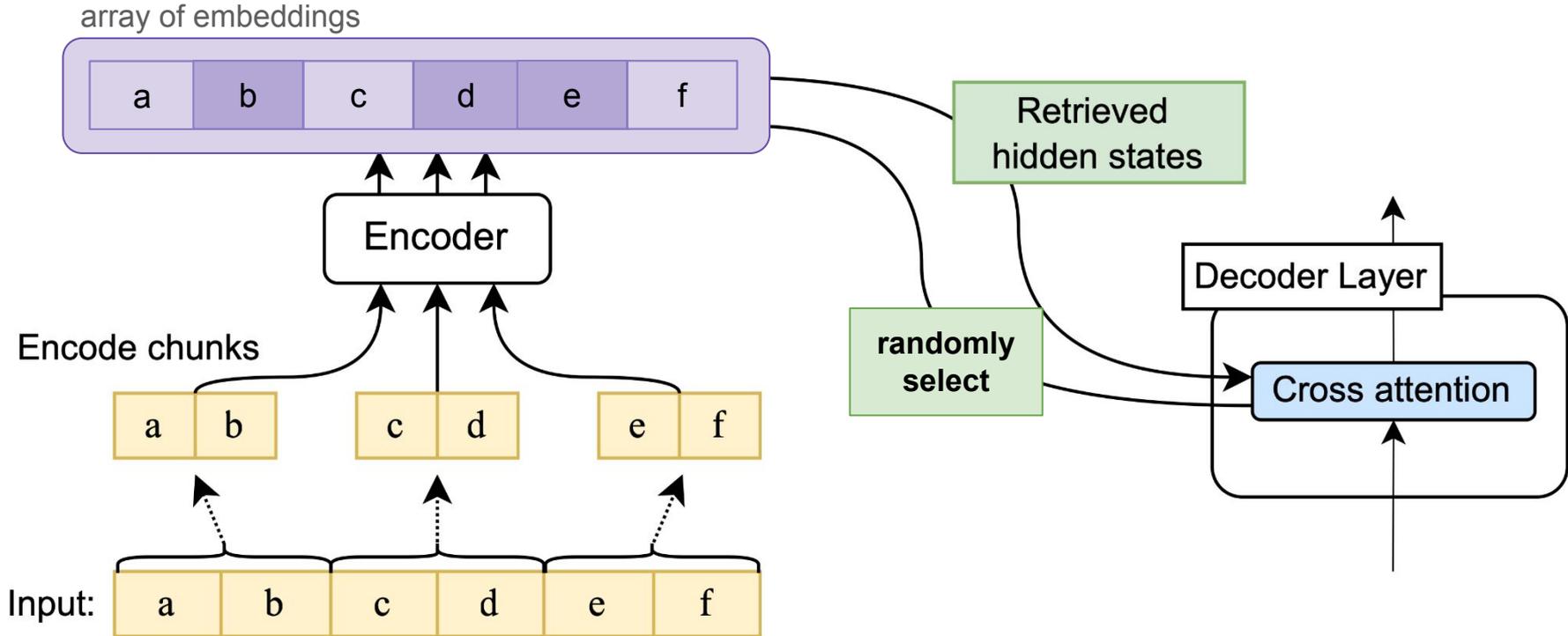# Higher cost training: which embeddings to backprop through?

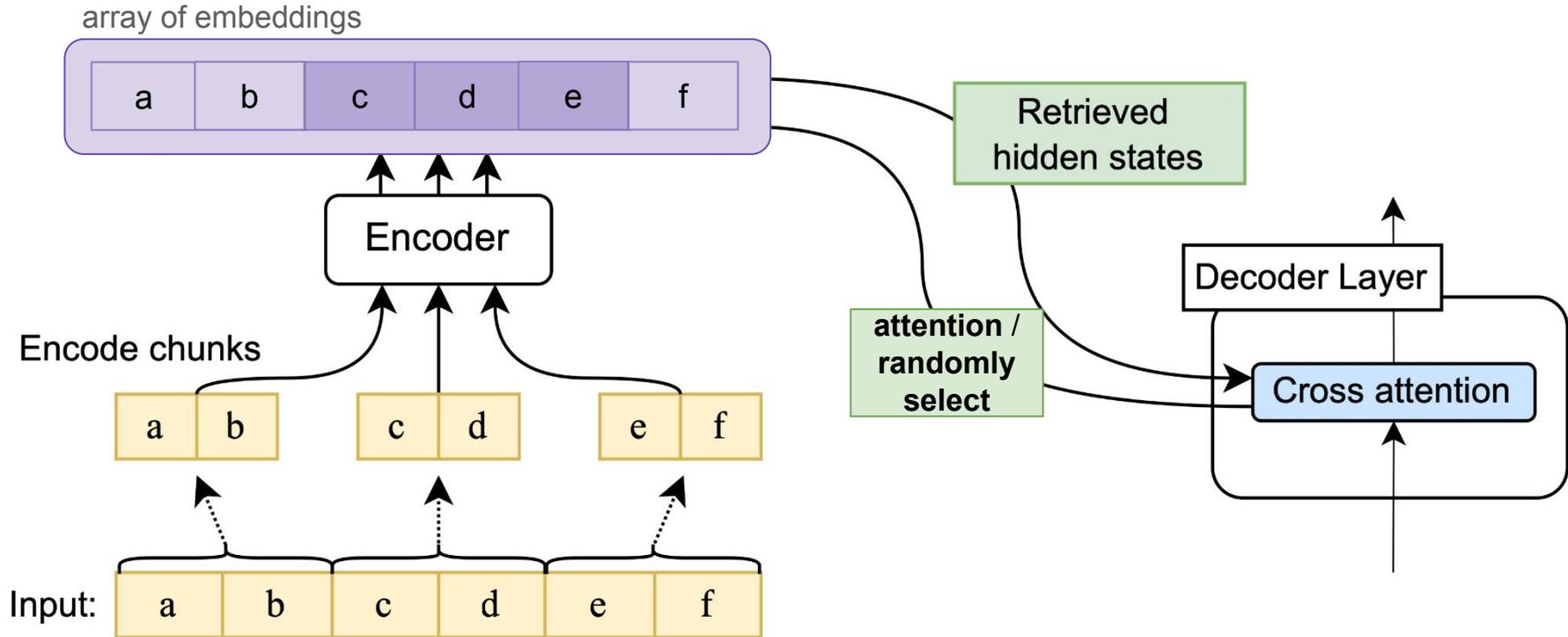# Higher cost training: which embeddings to backprop through?

# Higher cost training: retrieval training

# Higher cost training: random-encoded
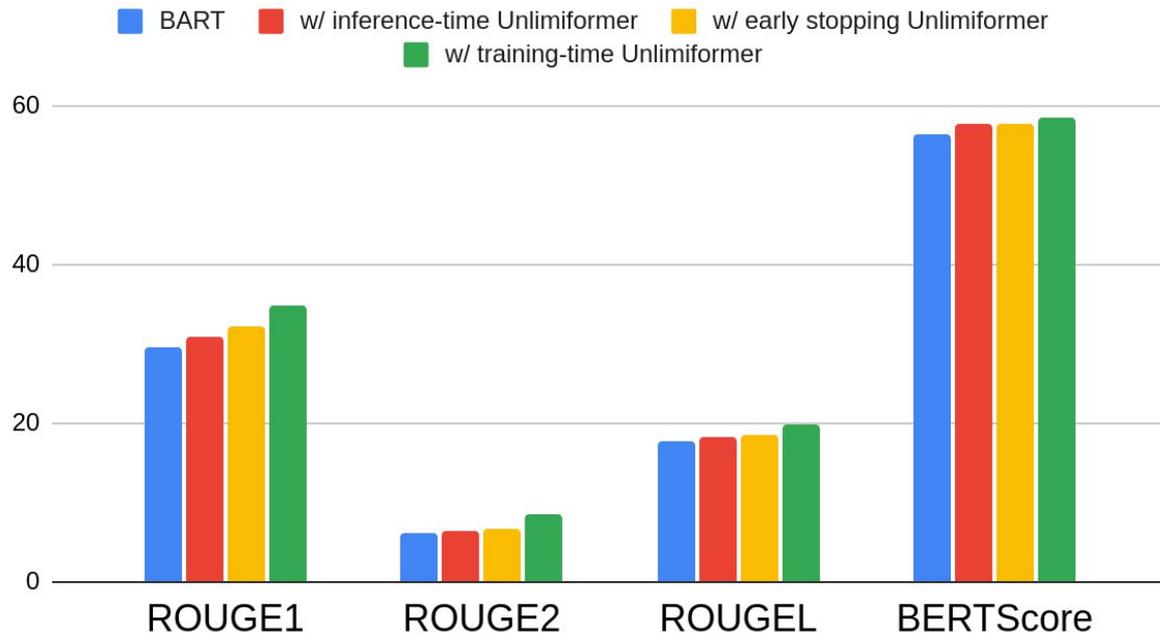
array of embeddings



28

# Higher cost training: alternating

# Results on SummScreen

**Domain:** TV screenplays

**Avg input length:** 8,987

**Avg output length:** 137

SummScreen

Legend: BART — w/ inference-time Unlimiformer — w/ early stopping Unlimiformer — w/ training-time Unlimiformer

Chart with categories: ROUGE1, ROUGE2, ROUGEL, BERTScore

# Results on GovReport

**Domain:** government reports

**Avg input length:** 9,616

**Avg output length:** 597



GovReport

Legend: BART, w/ inference-time Unlimiformer, w/ early stopping Unlimiformer, w/ training-time Unlimiformer

Categories: ROUGE1, ROUGE2, ROUGEL, BERTScore

model

# Comparison to other long-range methods [GovReport]

**Domain:** government reports

**Avg input length:** 9,616

**Avg output length:** 597

Long-range methods



Legend: BART, SLED, Memorizing Transformers, Unlimiformer

Categories: ROUGE1, ROUGE2, ROUGEL, BERTScore

model

# Results on BookSum

**Domain:**
public-domain
novels

**Avg input
length:** 143,301

**Avg output
length:** 1,294

BookSum



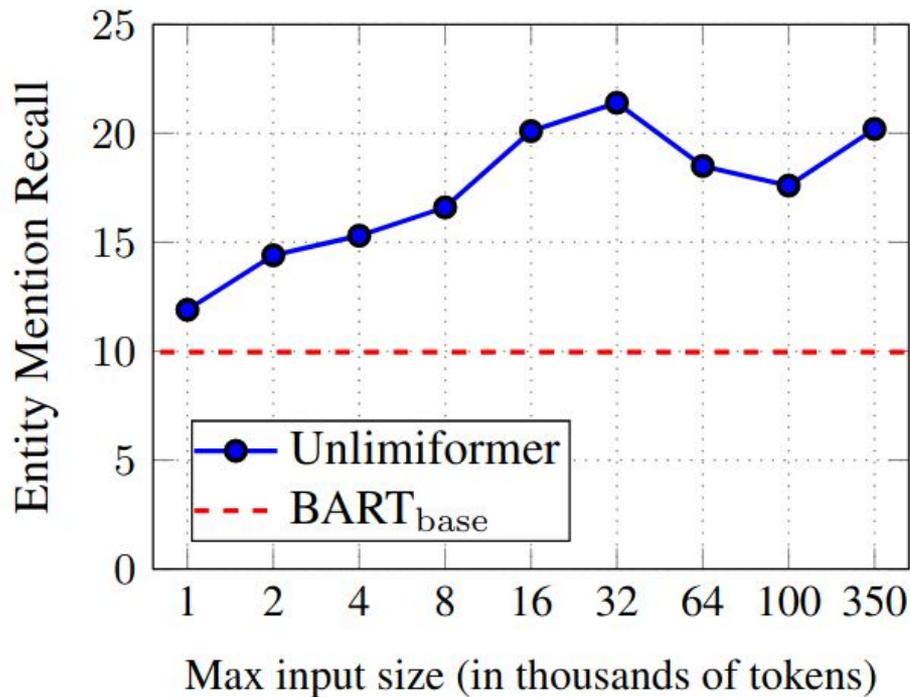Legend: PRIMERA, w/ inference-time Unlimiformer, w/ early stopping Unlimiformer, w/ training-time Unlimiformer
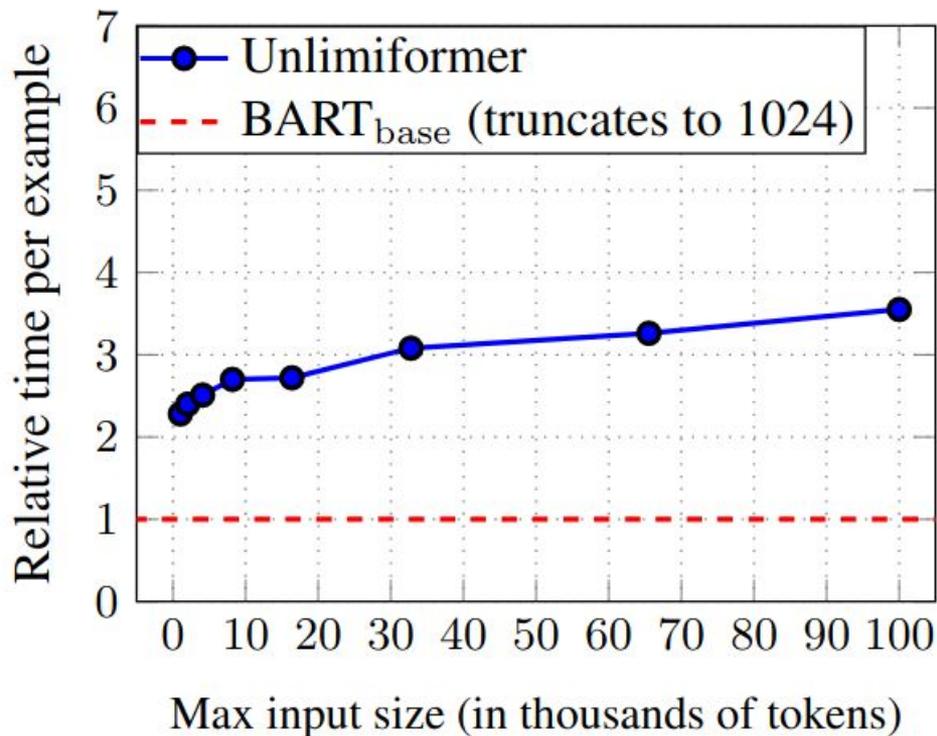
33

# EntMent

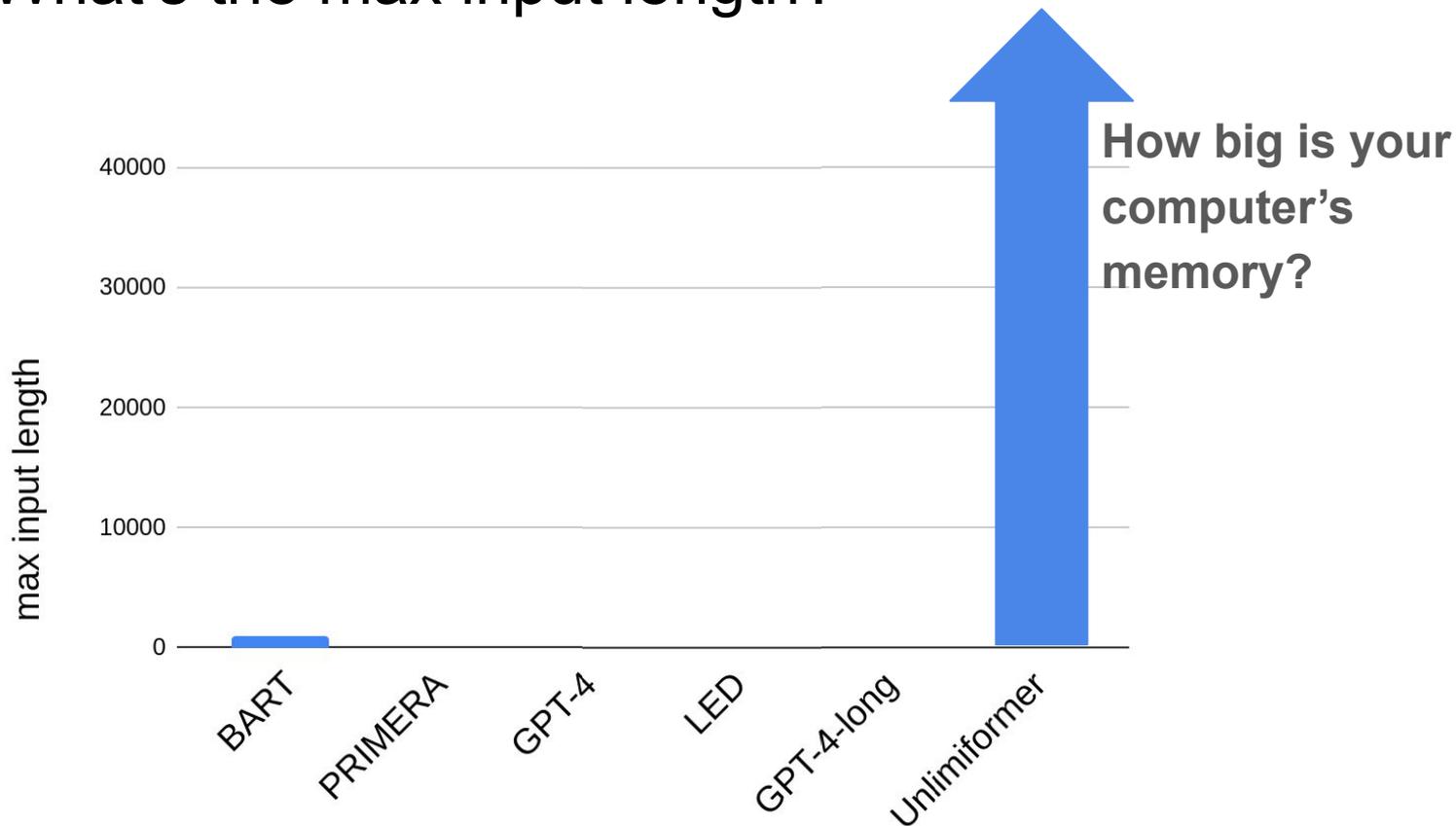Idea: **important to include entities from the gold summary**

# Computational cost

Additional cost from:

- Encoding additional input
- Datastore construction
- Datastore search

# What's the max input length?

**How big is your computer's memory?**



max input length

40000

30000

20000

10000

0

BART    PRIMERA    GPT-4    LED    GPT-4-long    Unlimiformer

# What (could be) next?

- Decoder-only models with Unlimiformer: LLaMA and Falcon
- Multi-doc summarization with Unlimiformer


- Better evaluation for long text
- Generation of long text
- Training to include *all* input

questions?

Amanda Bertsch    Uri Alon    Graham Neubig    Matt Gormley