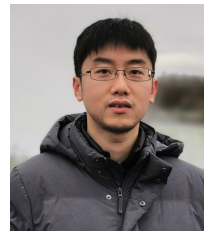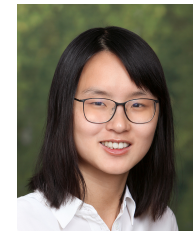# Training Language Models with Memory Augmentation

**Zexuan Zhong**   Tao Lei   Danqi Chen
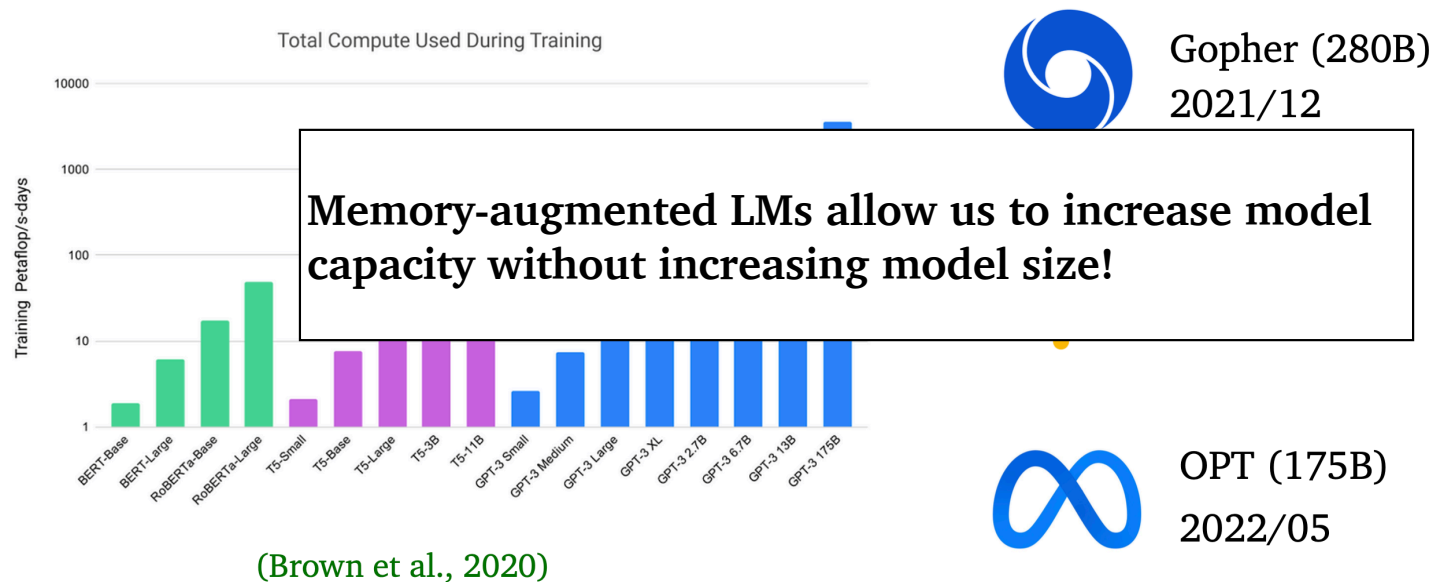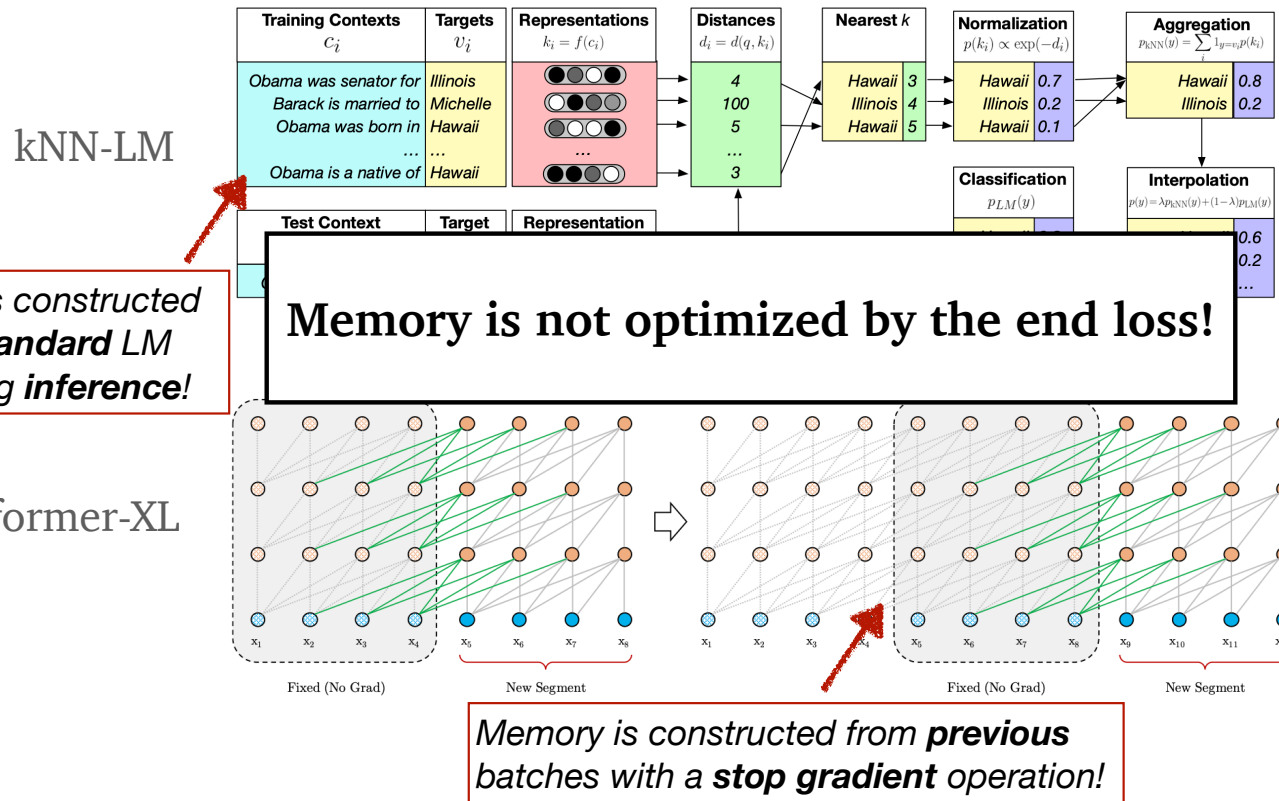
Princeton University

# Motivation

Large language models are very expensive to train/use!

Total Compute Used During Training

Gopher (280B)
2021/12

**Memory-augmented LMs allow us to increase model capacity without increasing model size!**

(Brown et al., 2020)

OPT (175B)
2022/05

# Prior works on memory-augmented LMs



kNN-LM

Transformer-XL

*Memory is constructed using a **standard** LM only during **inference**!*

**Memory is not optimized by the end loss!**

*Memory is constructed from **previous** batches with a **stop gradient** operation!*

Khandelwal et al., 2021. Generalization through memorization: Nearest neighbor language models.
Dai et al., 2019. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context

# Our approach: TRIME

Key idea: building a **training memory** from the **same training batch** on the fly
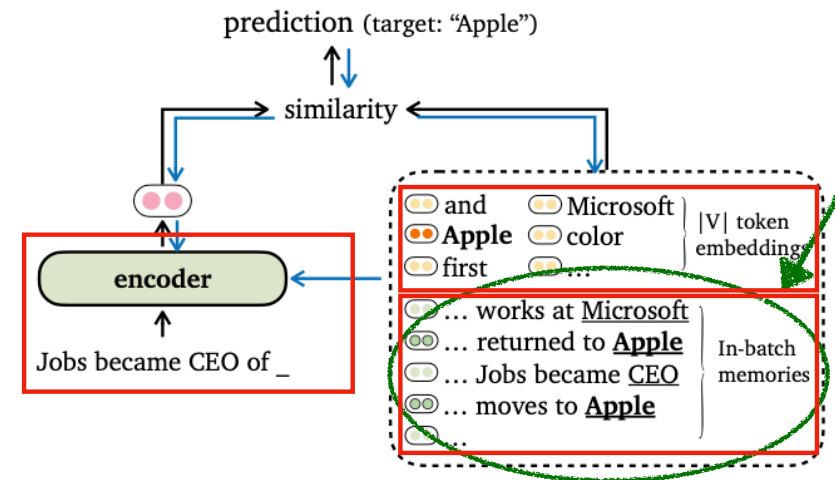


$c_t$: Jobs became CEO of ___

1) 🟠🟠 **Apple** (output embedding)

2) Other $c$ in the **training memory** that share the same next word as $x_t$

    🟢🟢 … returned to **Apple**

    🟢🟢 … moves to **Apple**

- ●● Target token's embedding
- ⊙⊙ Other token embeddings
- ◎◎ Positive in-batch memory
- ○○ Negative in-batch memory

↑ Forward pass   ↓ Back-propagation

prediction (target: "Apple")

similarity

encoder

Jobs became CEO of _

- ⊙● and      ⊙● Microsoft
- ●● **Apple**  ⊙● color
- ○● first     ⊙○ …

|V| token embeddings

- ○○ … works at <u>Microsoft</u>
- ◎◎ … returned to **Apple**
- ○○ … Jobs became <u>CEO</u>
- ◎◎ … moves to **Apple**
- ○○ …

In-batch memories

# Three TRIME language models



Current token    In memory    Not in memory

Segment len L

(a) Default batching
(b) Batching consecutive segments
(c) Batching lexically similar segments

Batch size B

Randomly sampled

Doc A
Doc B
Doc C

BM25 selected

**TrimeLM**
(local memory)

**TrimeLM$_{long}$**
(long-term memory)

**TrimeLM$_{ext}$**
(External memory)

*Compared to*

Vanilla LM

Transformer-XL

kNN-LM

5

# Experiments: WikiText-103
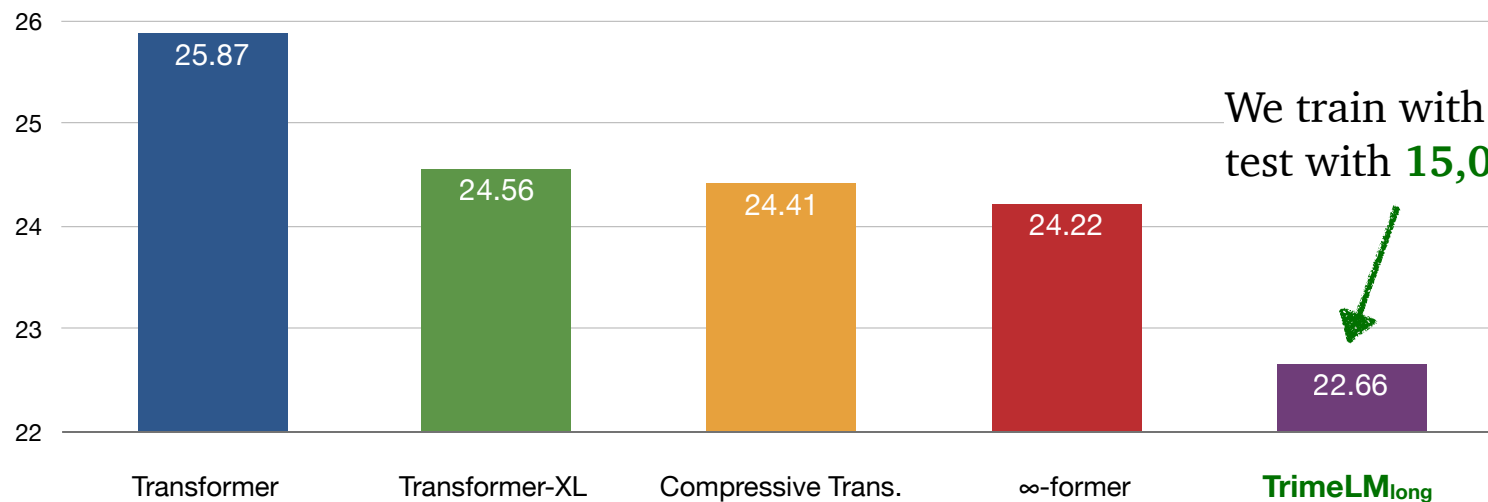## (247M model, L=3072)



Metric: perplexity

Khandelwal et al., 2021. Generalization through memorization: Nearest neighbor language models.
Dai et al., 2019. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context

# Experiments: WikiText-103

## (150M model, **L=150**)



We train with segment length **150**; test with **15,000** tokens!

Metric: perplexity

Dai et al., 2019. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context
Rae et al., 2019. Compressive Transformers for long-range sequence modeling.
Martins et al., 2022. ∞-former: Infinite memory Transformer.

# Summary: TRIME

**TRIME:** a **training objective** which leverages in-batch memories + three ways of memory construction

- adds very little computational overhead
- does not modify model architectures
- is compatible with other model architectures and techniques
- outperforms existing methods!

Check out more results on **domain adaptation**, **machine translation**, and **character-level LM** in the paper!

Paper: https://arxiv.org/pdf/2205.12674.pdf
Code & models: https://github.com/princeton-nlp/TRIME
Email: zzhong@cs.princeton.edu