# Keyformer: KV Cache reduction through key tokens selection for Efficient Generative Inference

**Muhammad Adnan**, Akhil Arunkumar, Gaurav Jain, Prashant J. Nair, Ilya Soloveychik, Purushotham Kamath
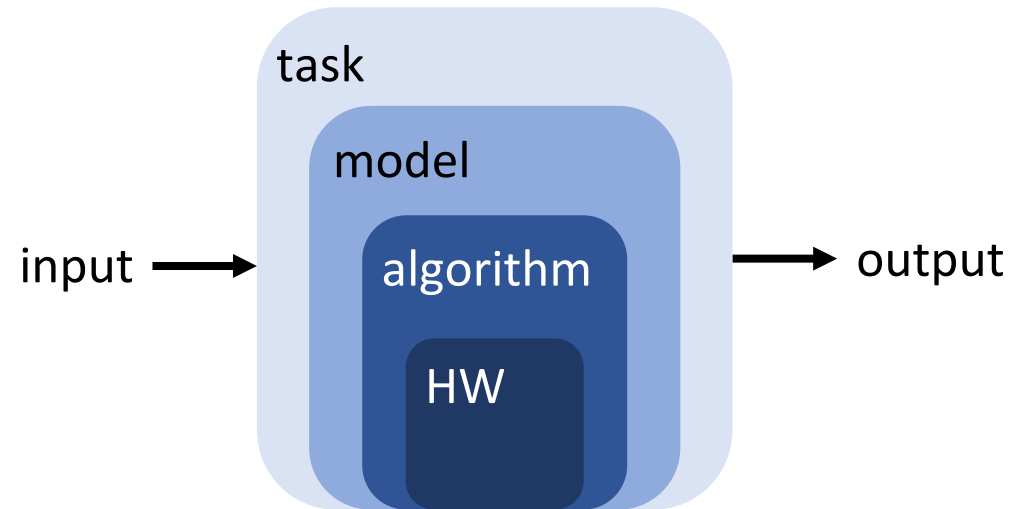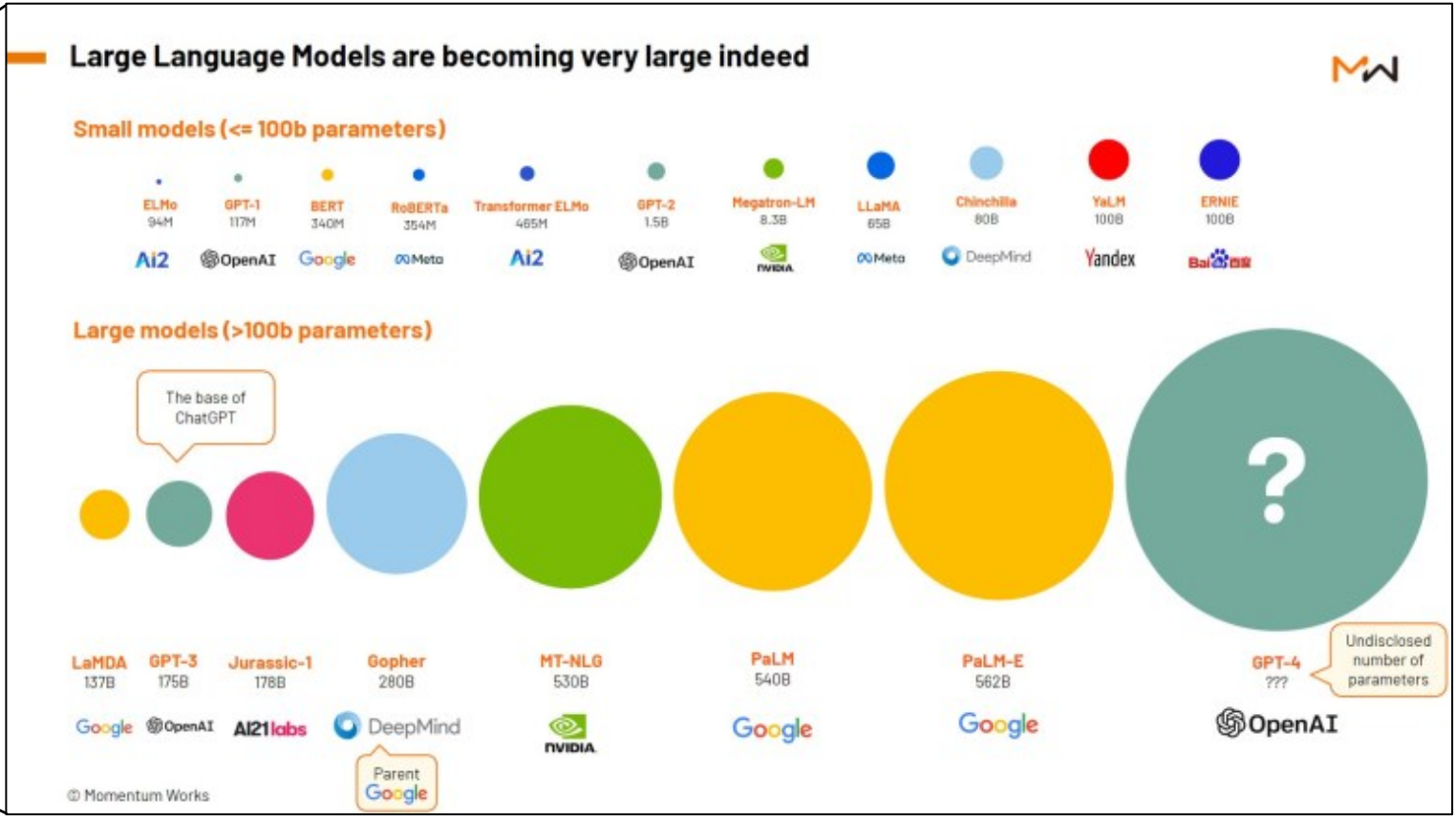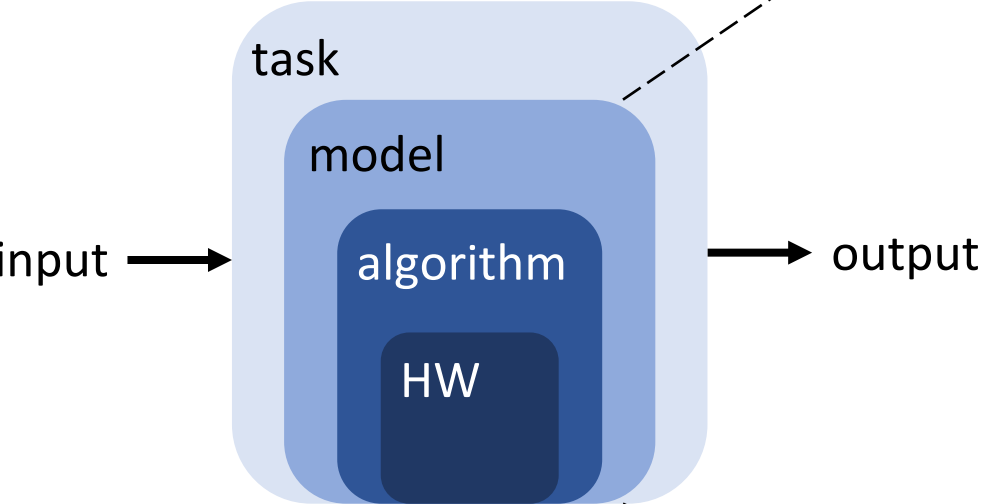
MLSys, 2024

THE UNIVERSITY OF BRITISH COLUMBIA

d-Matrix

# Background

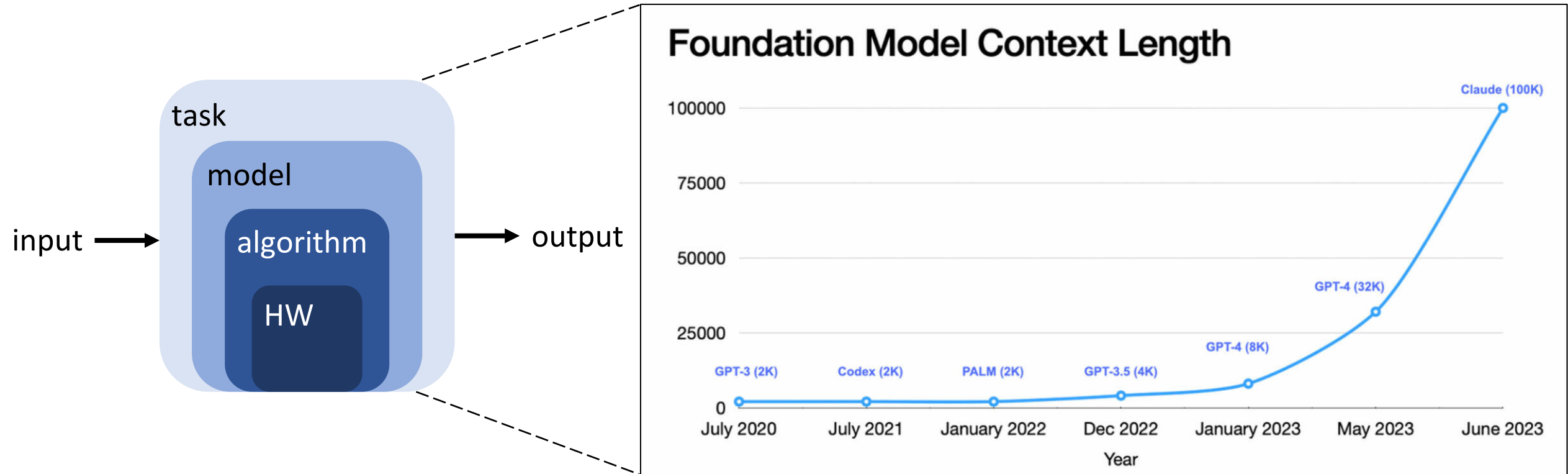Large Language Models (LLMs) abstraction

input →

task
model
algorithm
HW

→ output

# Background

## Models are getting bigger

task

model

input →

algorithm

HW

→ output

**Large Language Models are becoming very large indeed**

**Small models (<= 100b parameters)**

| ELMo | GPT-1 | BERT | RoBERTa | Transformer ELMo | GPT-2 | Megatron-LM | LLaMA | Chinchilla | YaLM | ERNIE |
|---|---|---|---|---|---|---|---|---|---|---|
| 94M | 117M | 340M | 354M | 465M | 1.5B | 8.3B | 65B | 80B | 100B | 100B |
| Ai2 | OpenAI | Google | Meta | Ai2 | OpenAI | nVIDIA | Meta | DeepMind | Yandex | Baidu |

**Large models (>100b parameters)**

The base of ChatGPT

Undisclosed number of parameters

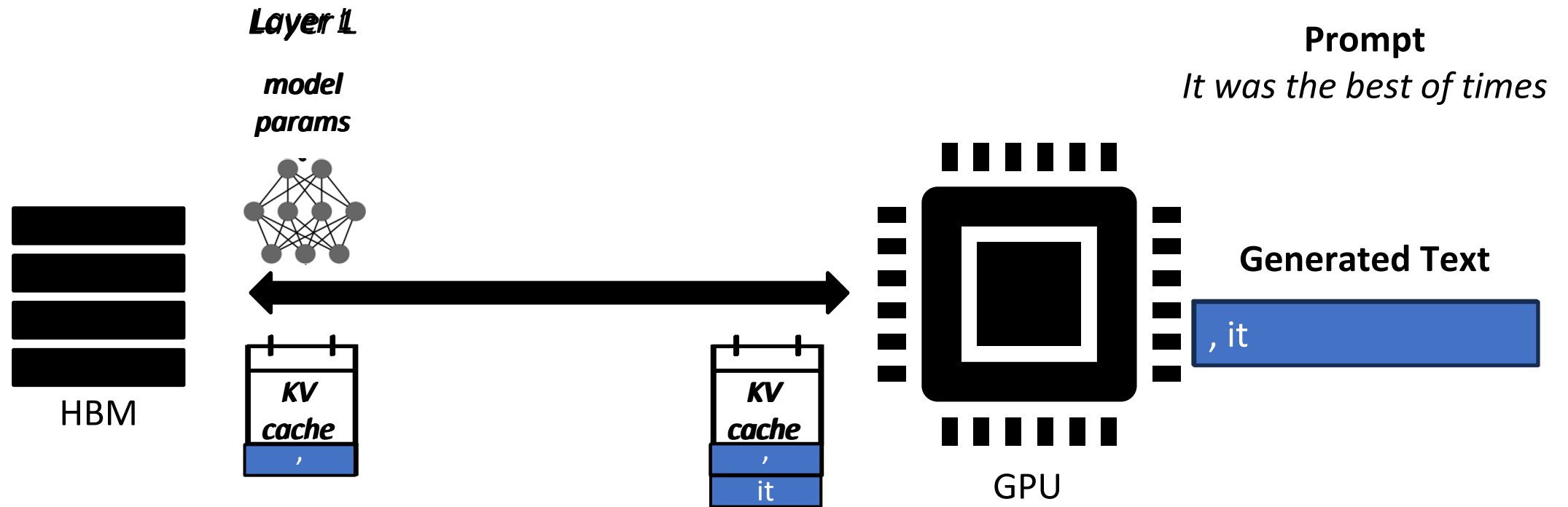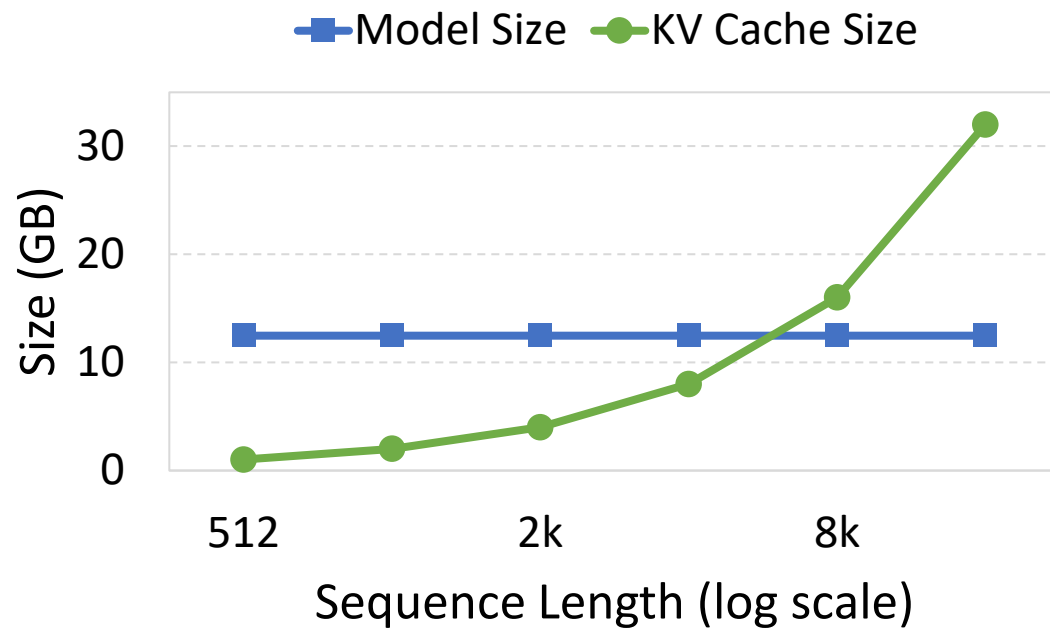| LaMDA | GPT-3 | Jurassic-1 | Gopher | MT-NLG | PaLM | PaLM-E | GPT-4 |
|---|---|---|---|---|---|---|---|
| 137B | 175B | 178B | 280B | 530B | 540B | 562B | ??? |
| Google | OpenAI | AI21labs | DeepMind | nVIDIA | Google | Google | OpenAI |

Parent Google

© Momentum Works

# Background

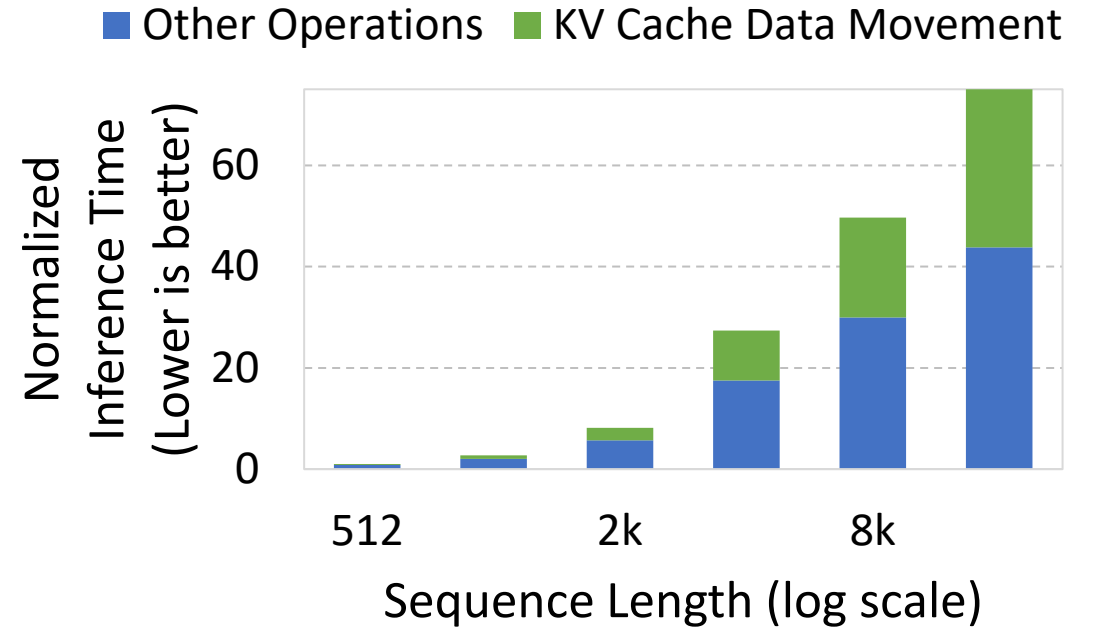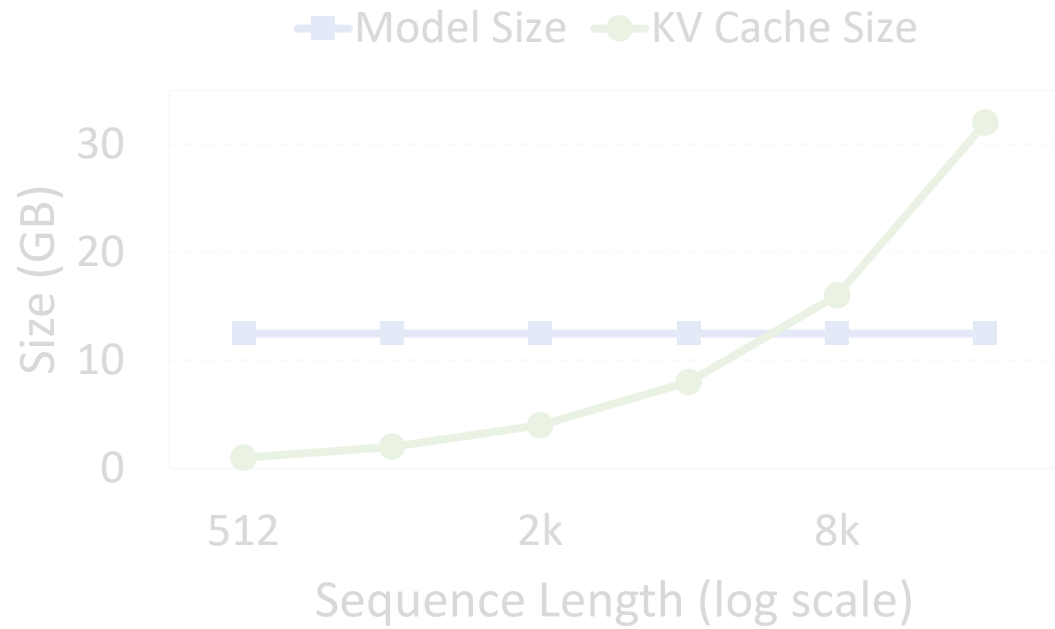Tasks require longer context length
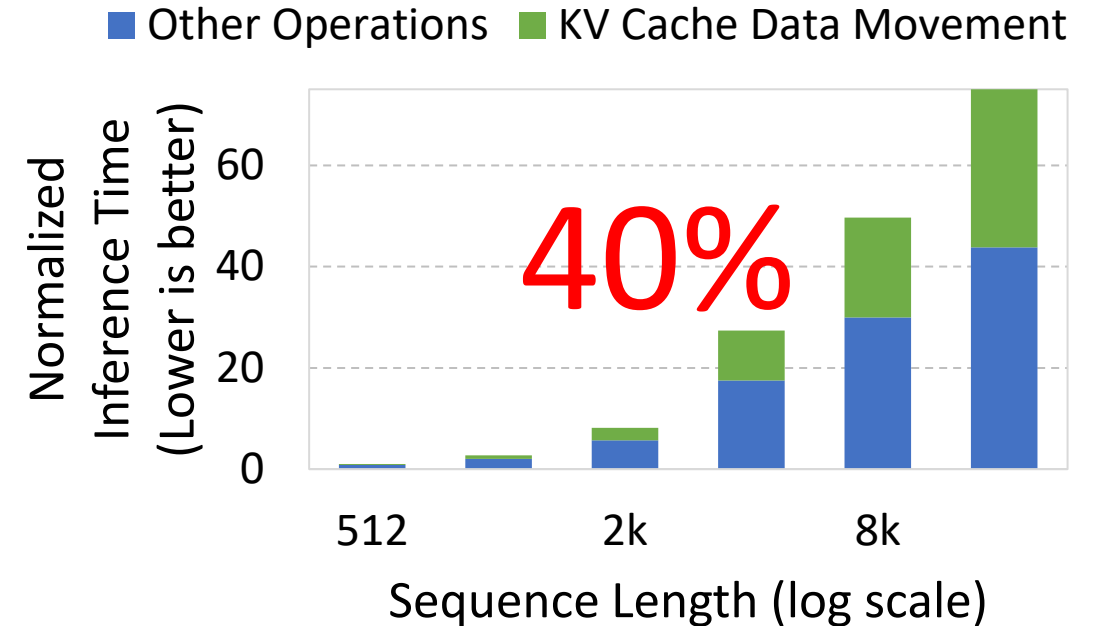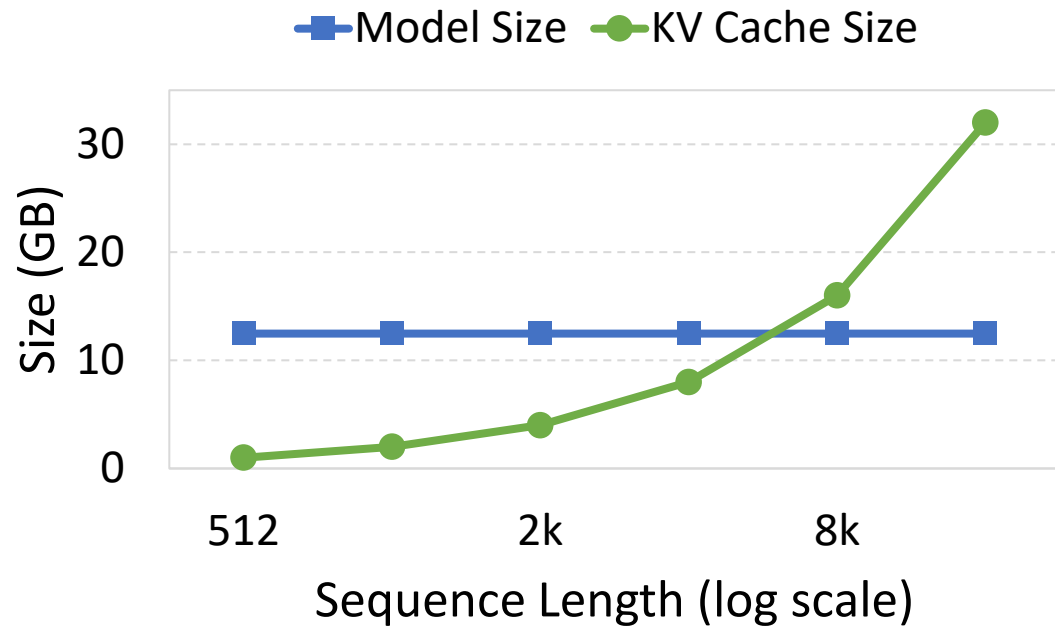
# Off-chip Data Movement

# Problem: Capacity, Bandwidth

# Problem: Inference Time

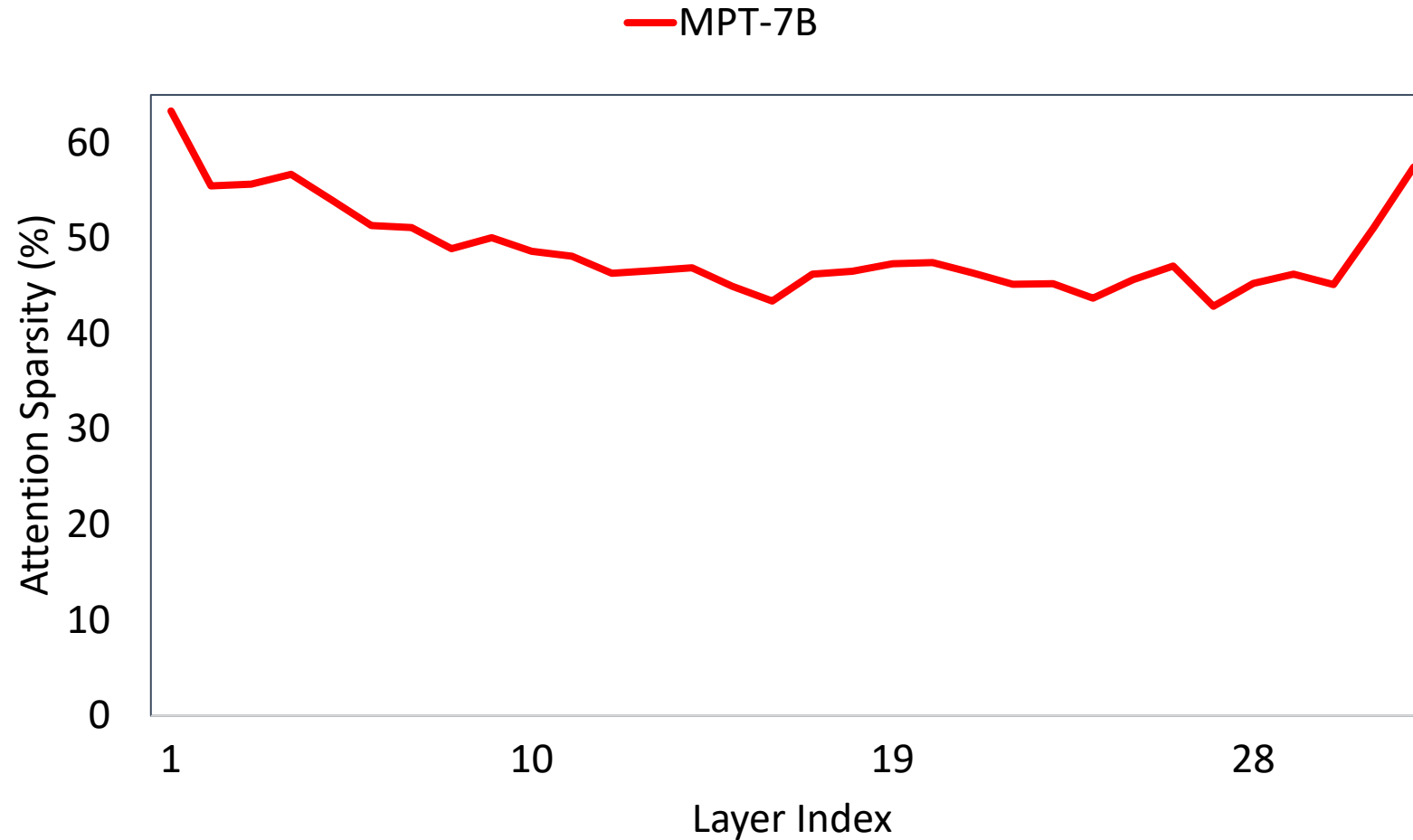# Problem: Capacity, Bandwidth & Time



Profiling with MPT-7B model using varying Sequence Length (50% context + 50% text generation) using NVIDIA A100 (80GB) GPU and synthetic data.
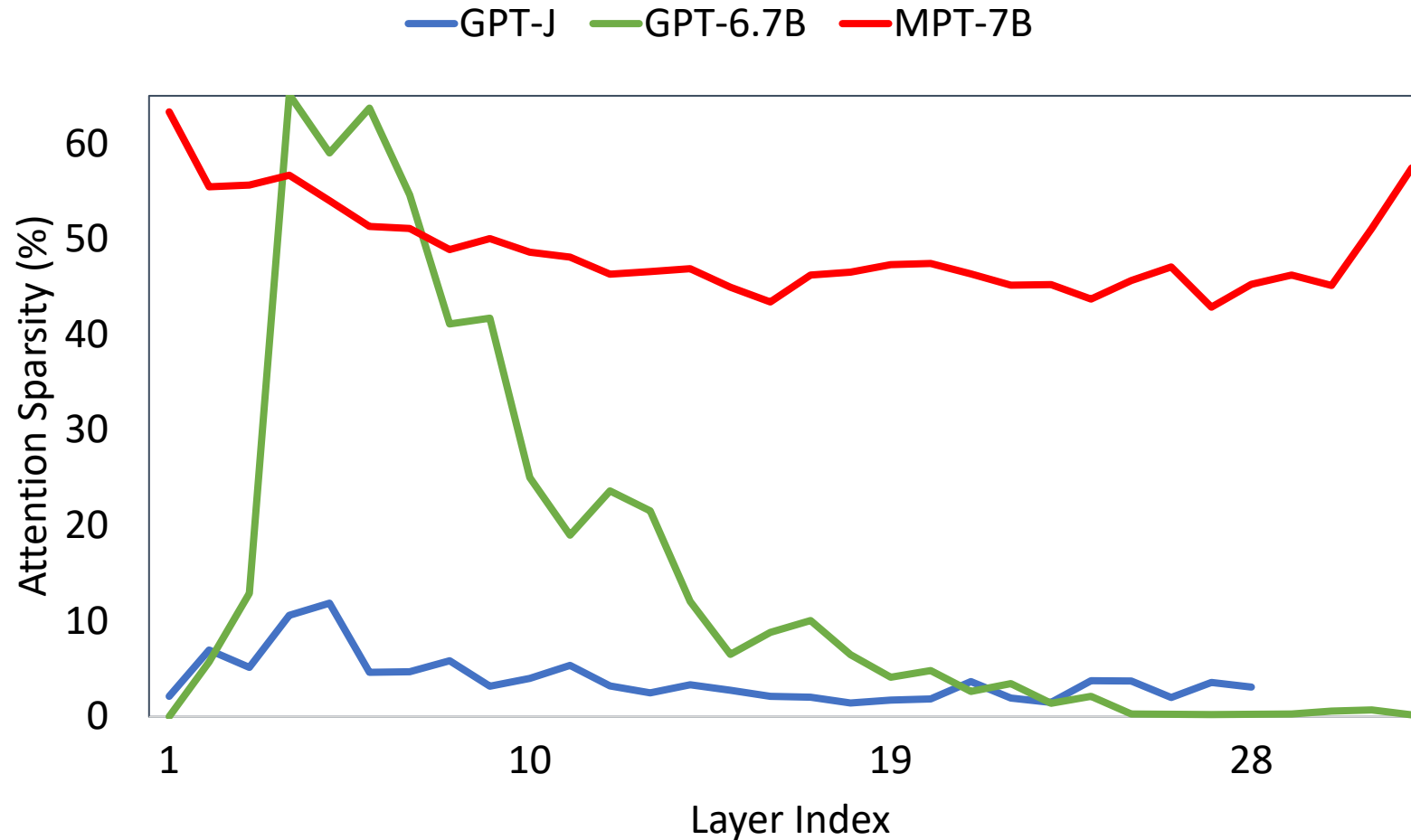
# Goals

KV cache reduction without accuracy drop

KV cache reduction on the fly without any retraining or finetuning
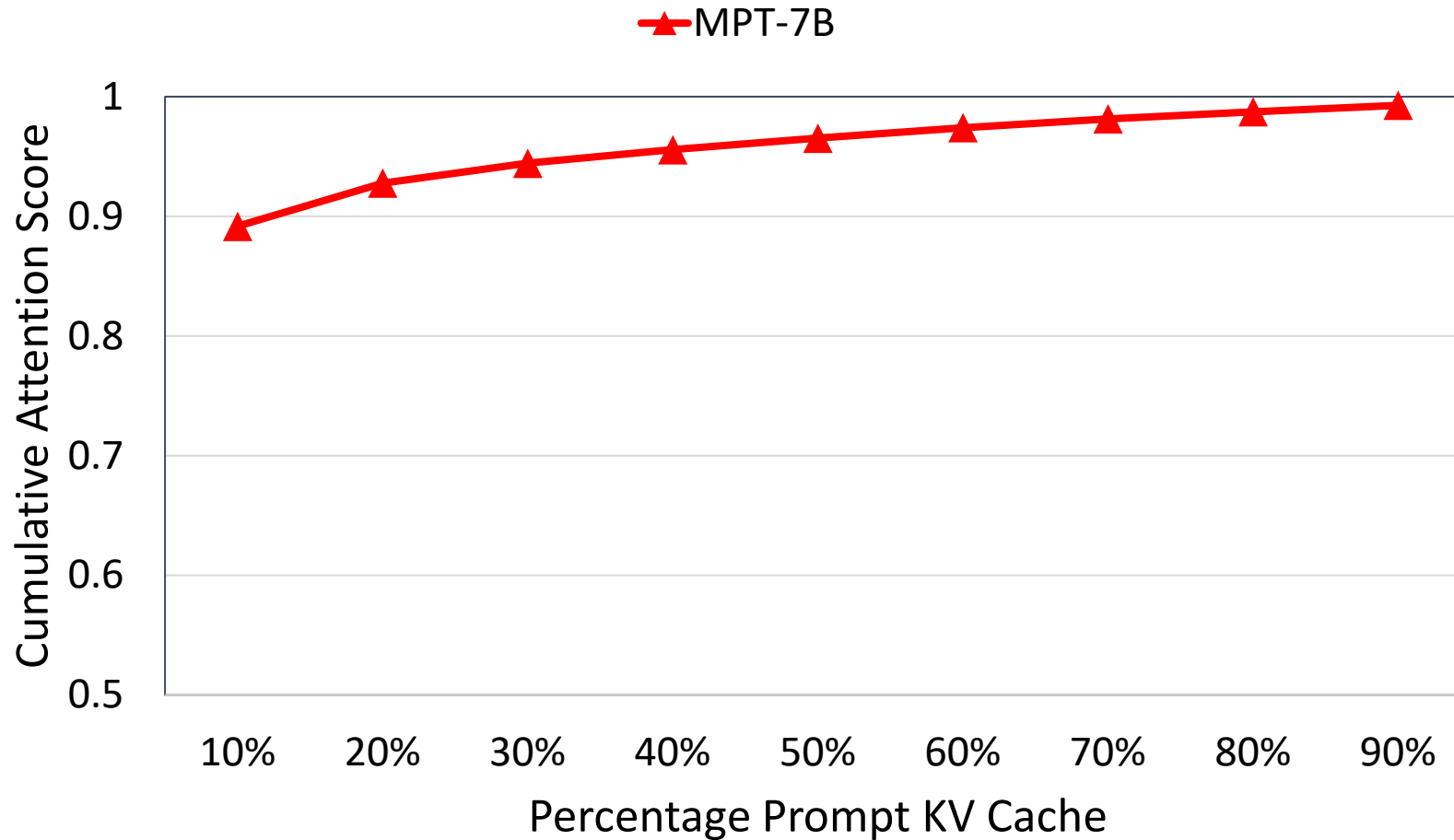
# Key Insight – I: Sparsity
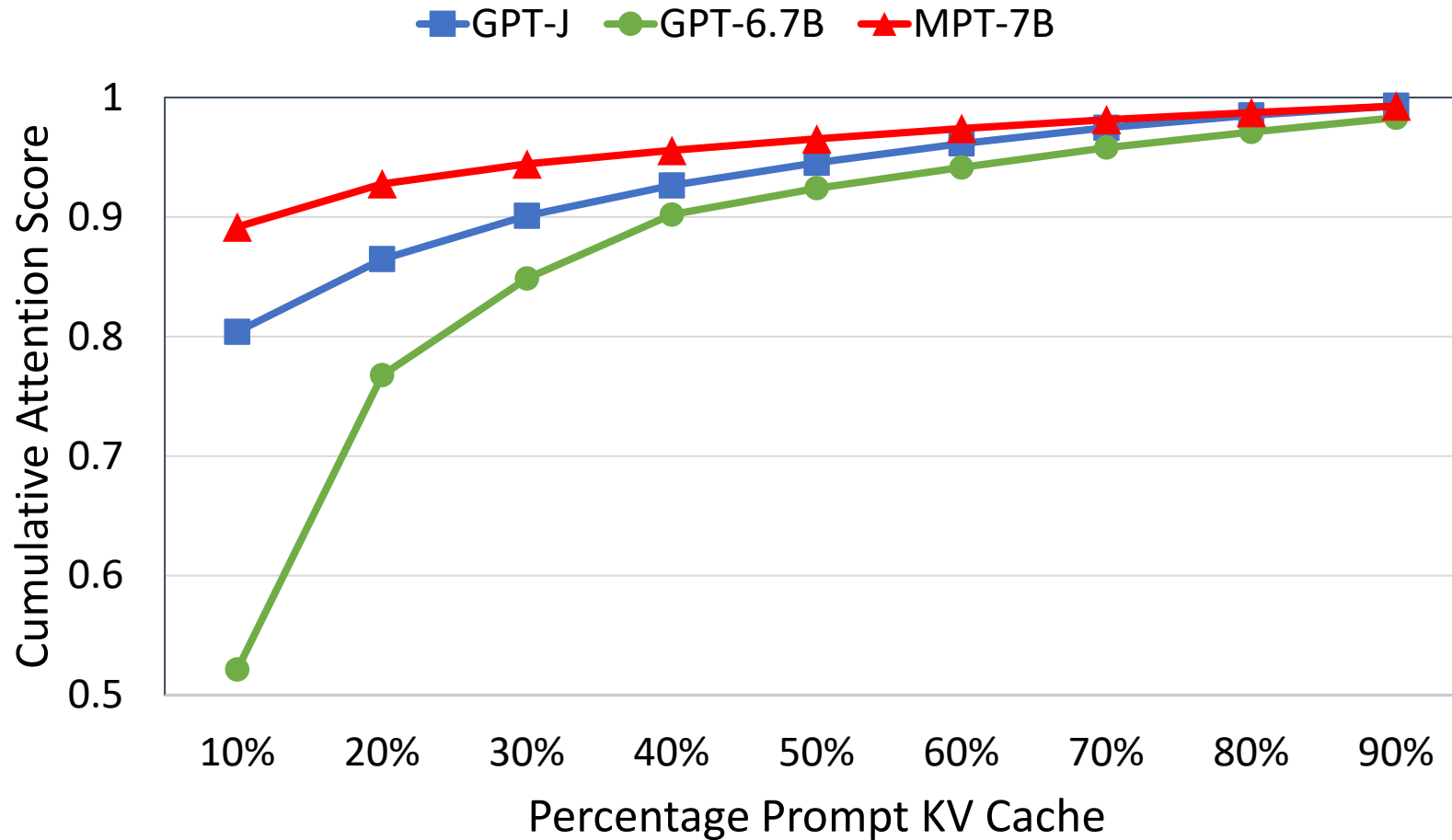
# Key Insight – I: Sparsity



Average Attention Sparsity using CNN/DailyMail dataset for Summarization Task.
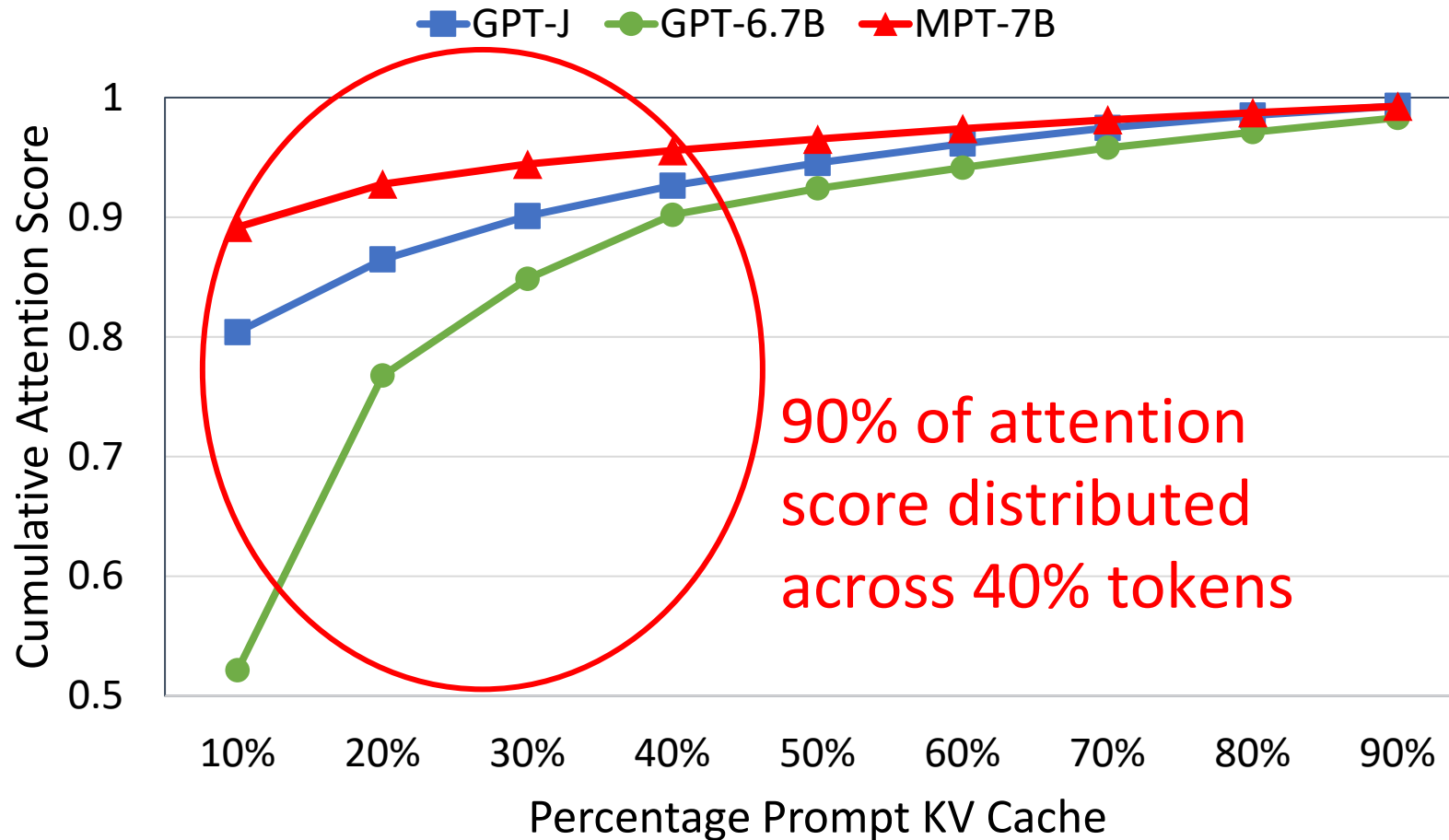
# Key Insight – II: Attention Distribution

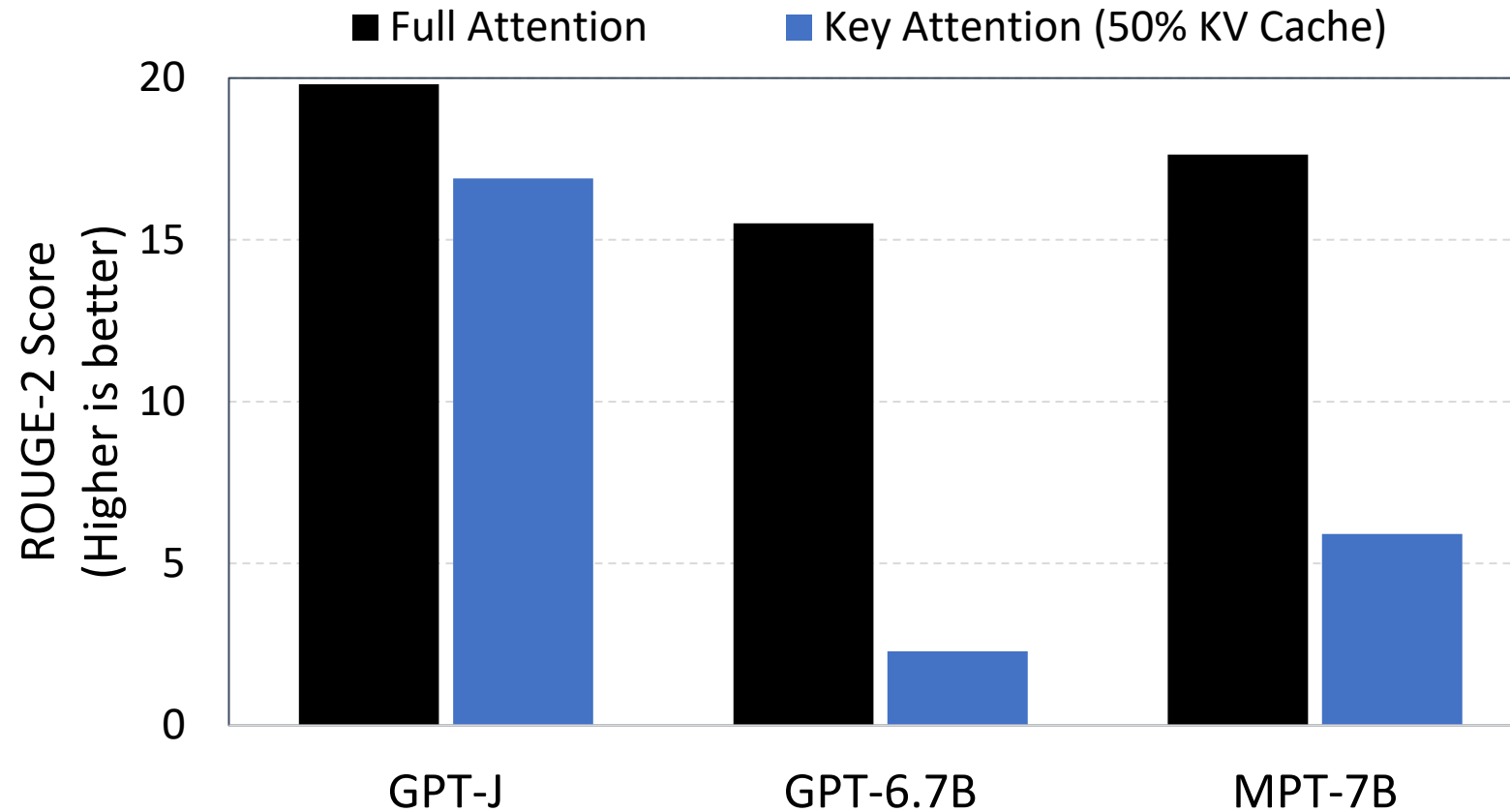# Key Insight – II: Attention Distribution

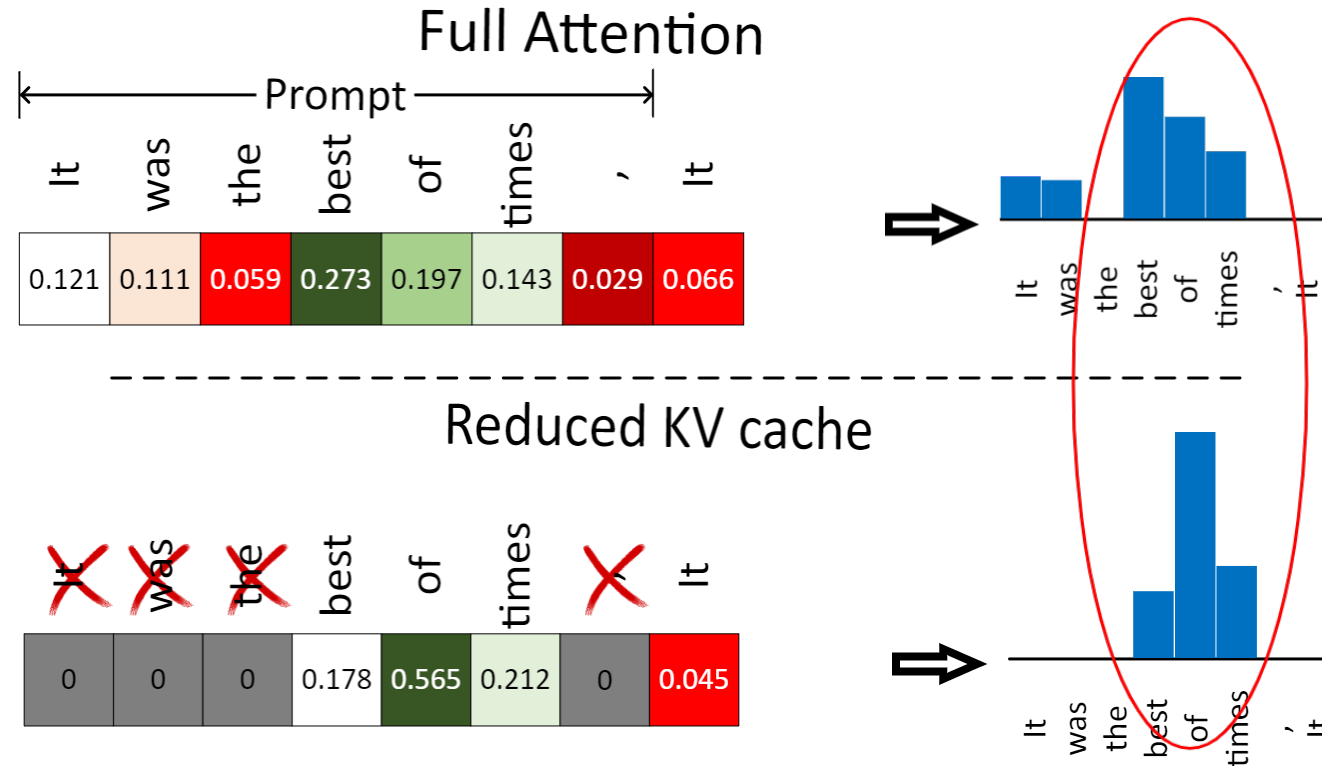# Key Insight – II: Attention Distribution



Attention Score Distribution using CNN/DailyMail dataset for Summarization Task.

# Limitation – Key Tokens



Accuracy Drop with 50% KV Cache for Summarization Task with CNN/DailyMail dataset.
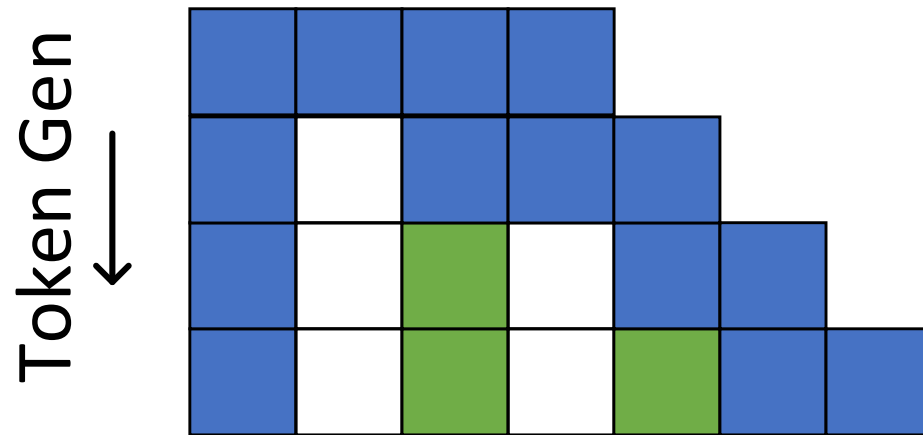
# Key Insight – III: Change in Score Distribution



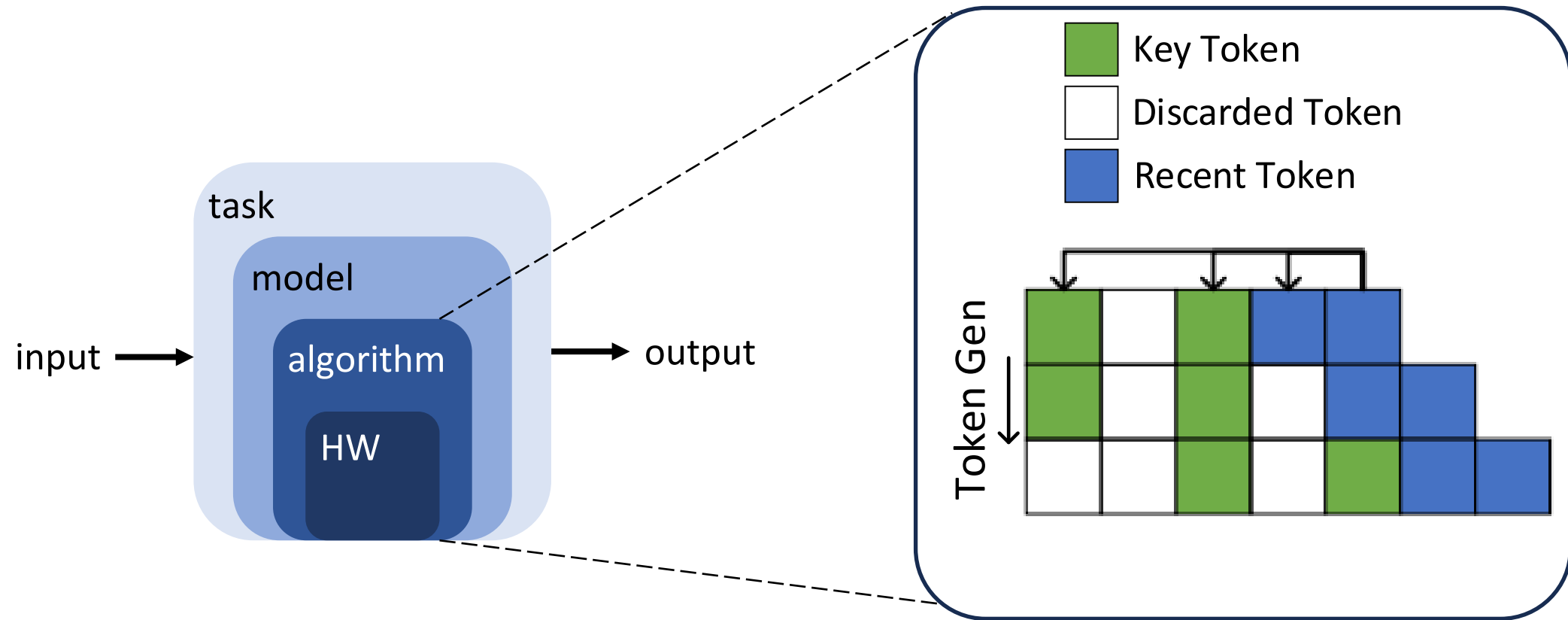MPT-7B model layer 1 and Head 1 stats.

# Challenges



Key Tokens Identification

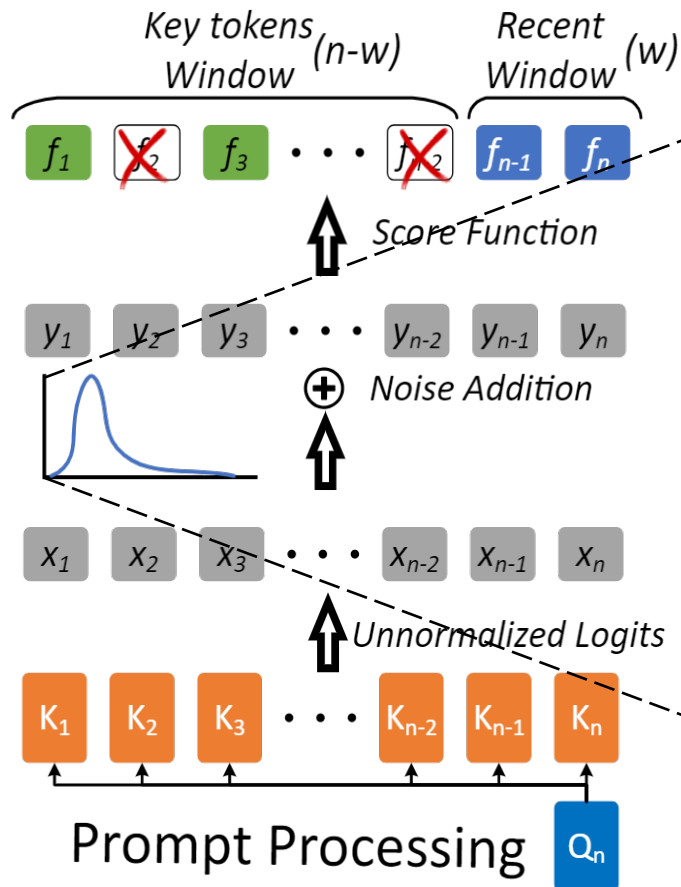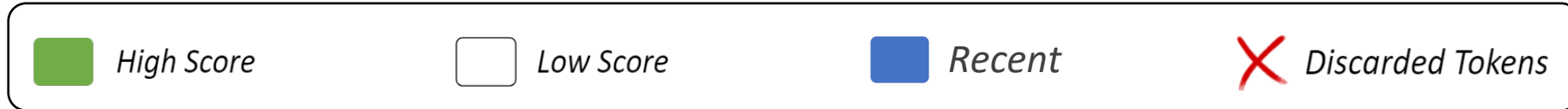Key Token    Discarded Token

Token Gen
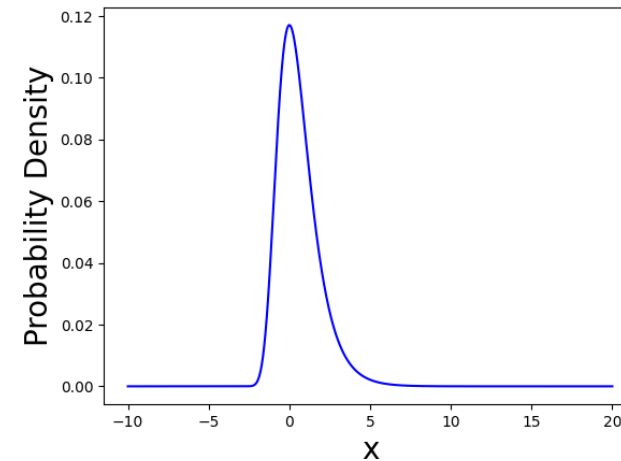
Discarded Tokens Utilization

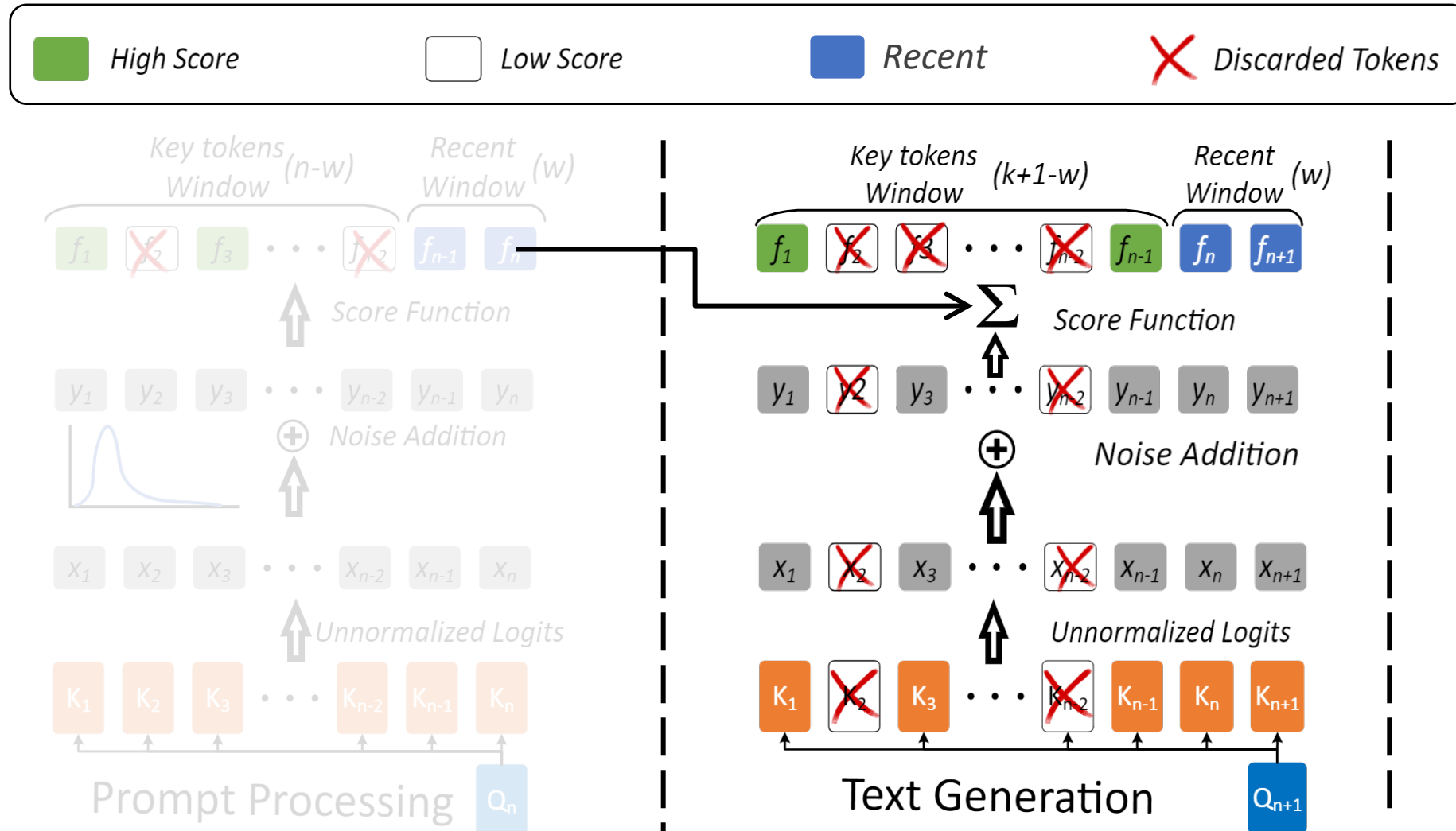# Keyformer – Regularization based Key Token
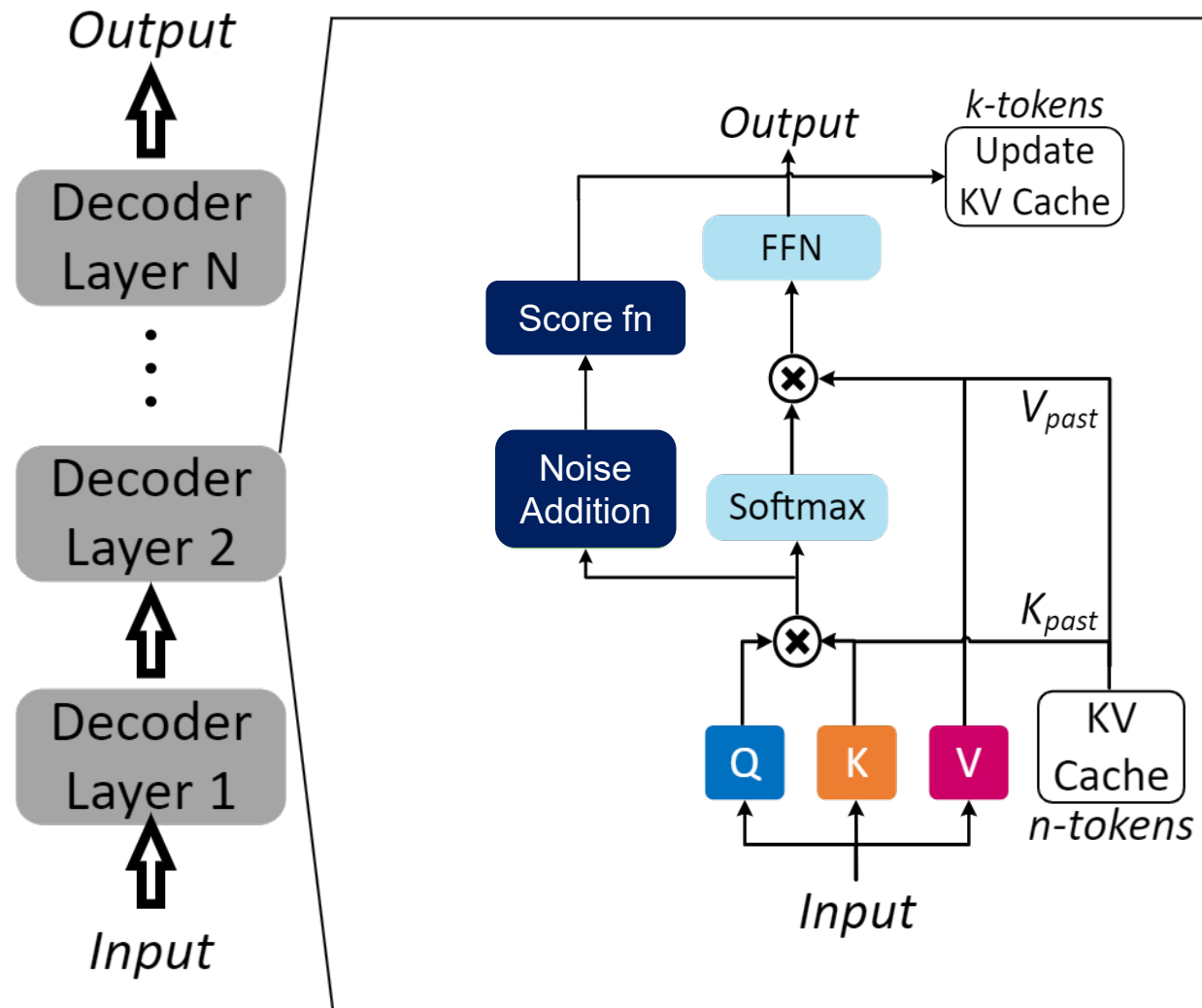
# Keyformer - KV Cache Reduction



- Bias towards Initial tokens
- Long tail Distribution

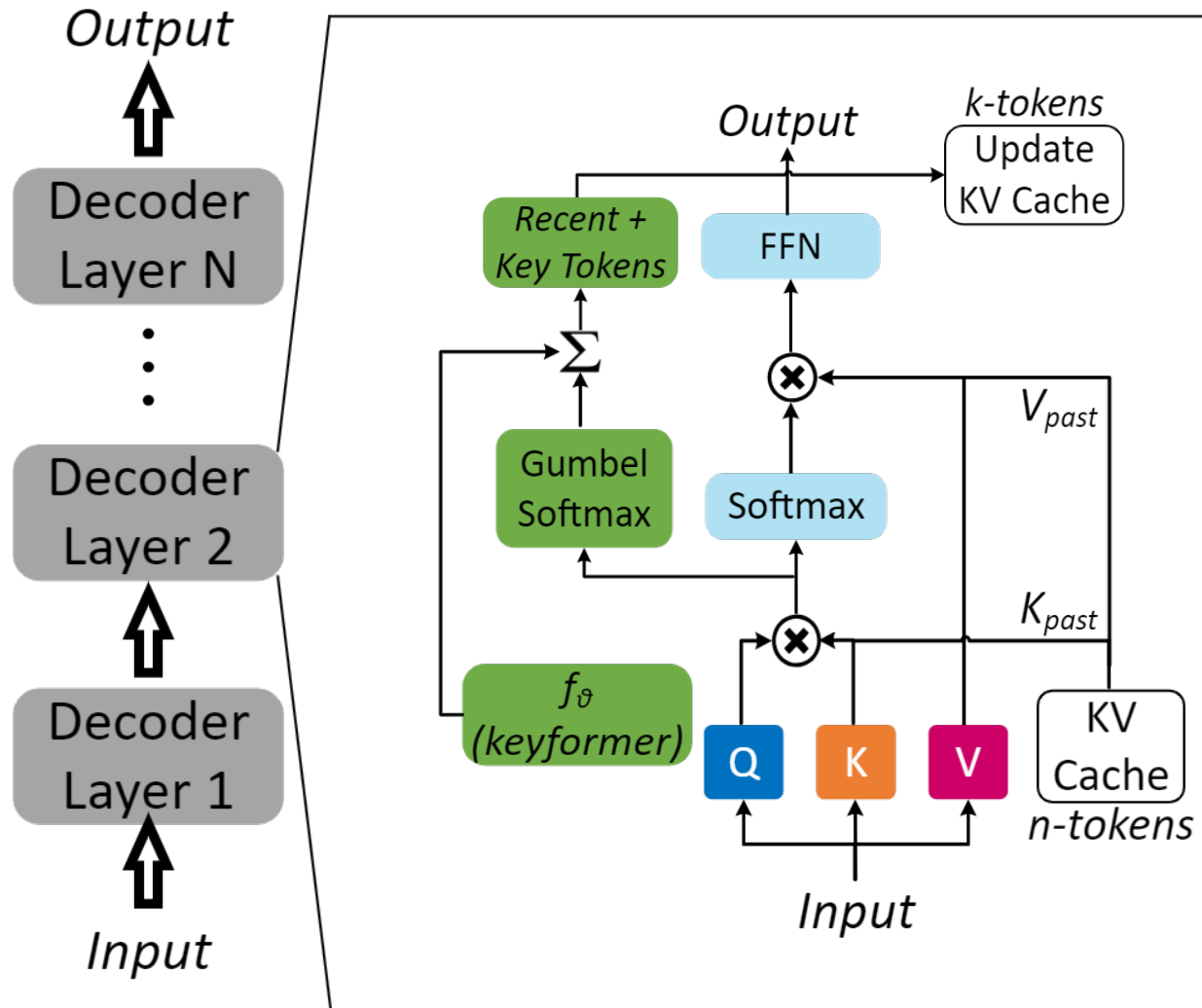# Keyformer - KV Cache Reduction

# Keyformer - Decoder
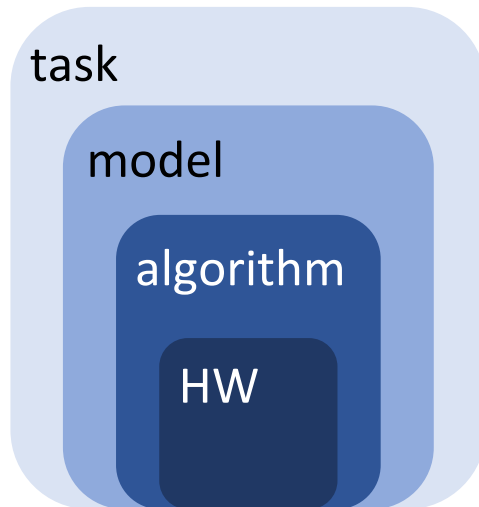
# Keyformer - Decoder

# Evaluation

task
model
algorithm
HW

Summarization
Conversation,
COPA (Commonsense Reasoning)
OpenBookQA (Commonsense Reasoning)
Winogrande (Language Understanding)
PIQA (Commonsense Reasoning)

GPT-J (finetuned – Summarization) - *RoPE*
MPT- chat (finetuned – Dialogue) - *ALiBi*
Cerebras GPT (pretrained) - *Absolute*
MPT (pretrained) - *ALiBi*

Full Attention
Window Attention
$H_2O$

NVIDIA A100 (80GB)

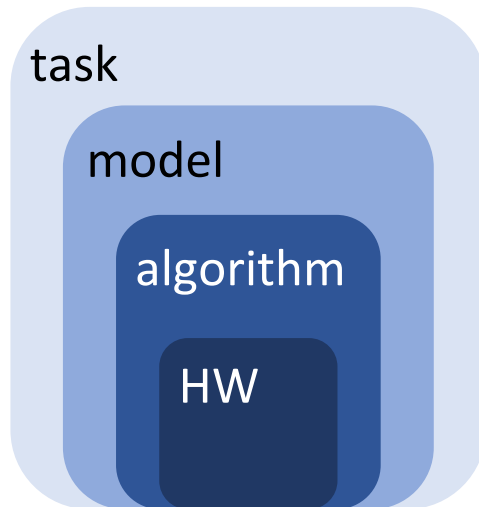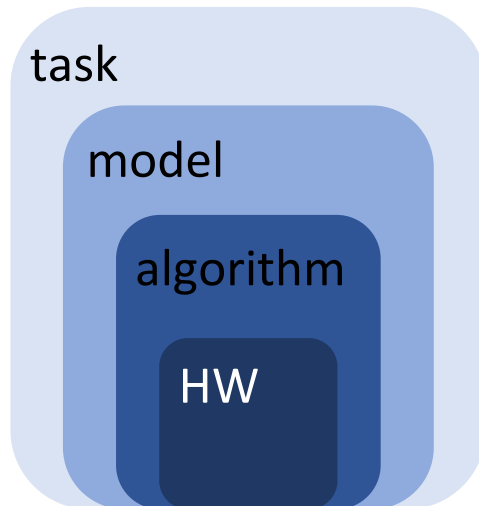# Evaluation

task

model

algorithm

HW

Summarization
Conversation,
COPA (Commonsense Reasoning)
OpenBookQA (Commonsense Reasoning)
Winogrande (Language Understanding)
PIQA (Commonsense Reasoning)

GPT-J (finetuned – Summarization) - *RoPE*
MPT- chat (finetuned – Dialogue) - *ALiBi*
Cerebras GPT (pretrained) - *Absolute*
MPT (pretrained) - *ALiBi*

Full Attention
Window Attention
$H_2O$

NVIDIA A100 (80GB)
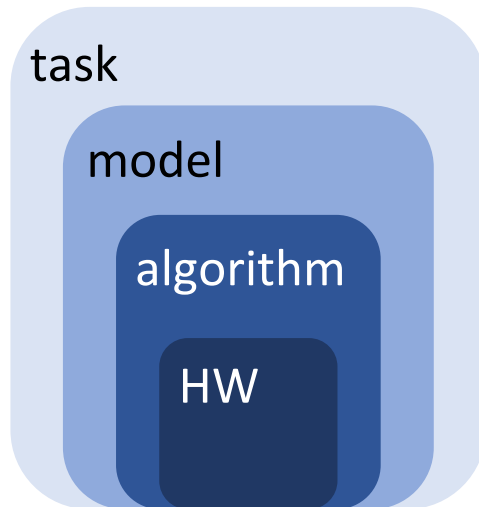
# Evaluation

task
model
algorithm
HW

Summarization
Conversation,
COPA (Commonsense Reasoning)
OpenBookQA (Commonsense Reasoning)
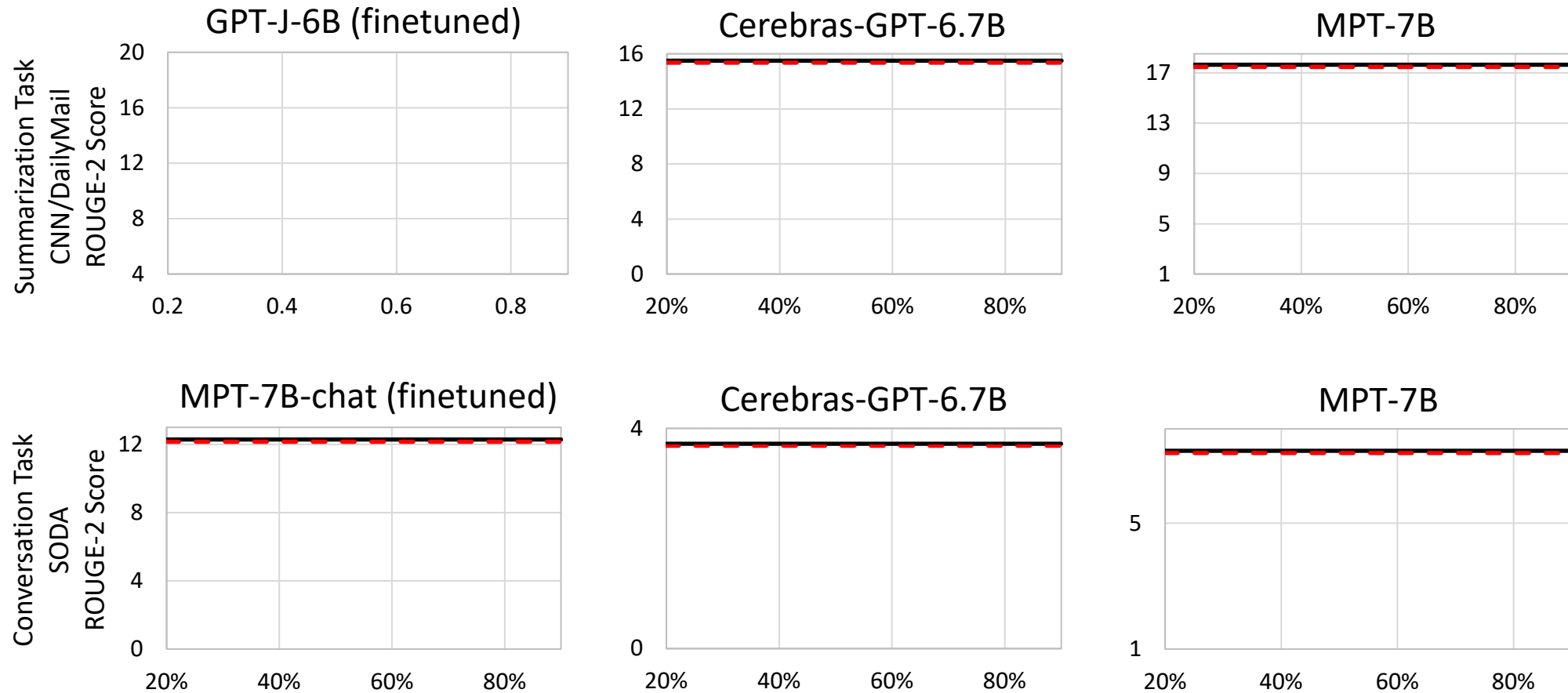Winogrande (Language Understanding)
PIQA (Commonsense Reasoning)

GPT-J (finetuned – Summarization) - *RoPE*
MPT- chat (finetuned – Dialogue) - *ALiBi*
Cerebras GPT (pretrained) - *Absolute*
MPT (pretrained) - *ALiBi*

Full Attention
Window Attention
$H_2O$

NVIDIA A100 (80GB)

# Evaluation

task
model
algorithm
HW

Summarization
Conversation,
COPA (Commonsense Reasoning)
OpenBookQA (Commonsense Reasoning)
Winogrande (Language Understanding)
PIQA (Commonsense Reasoning)

GPT-J (finetuned – Summarization) - *RoPE*
MPT- chat (finetuned – Dialogue) - *ALiBi*
Cerebras GPT (pretrained) - *Absolute*
MPT (pretrained) - *ALiBi*

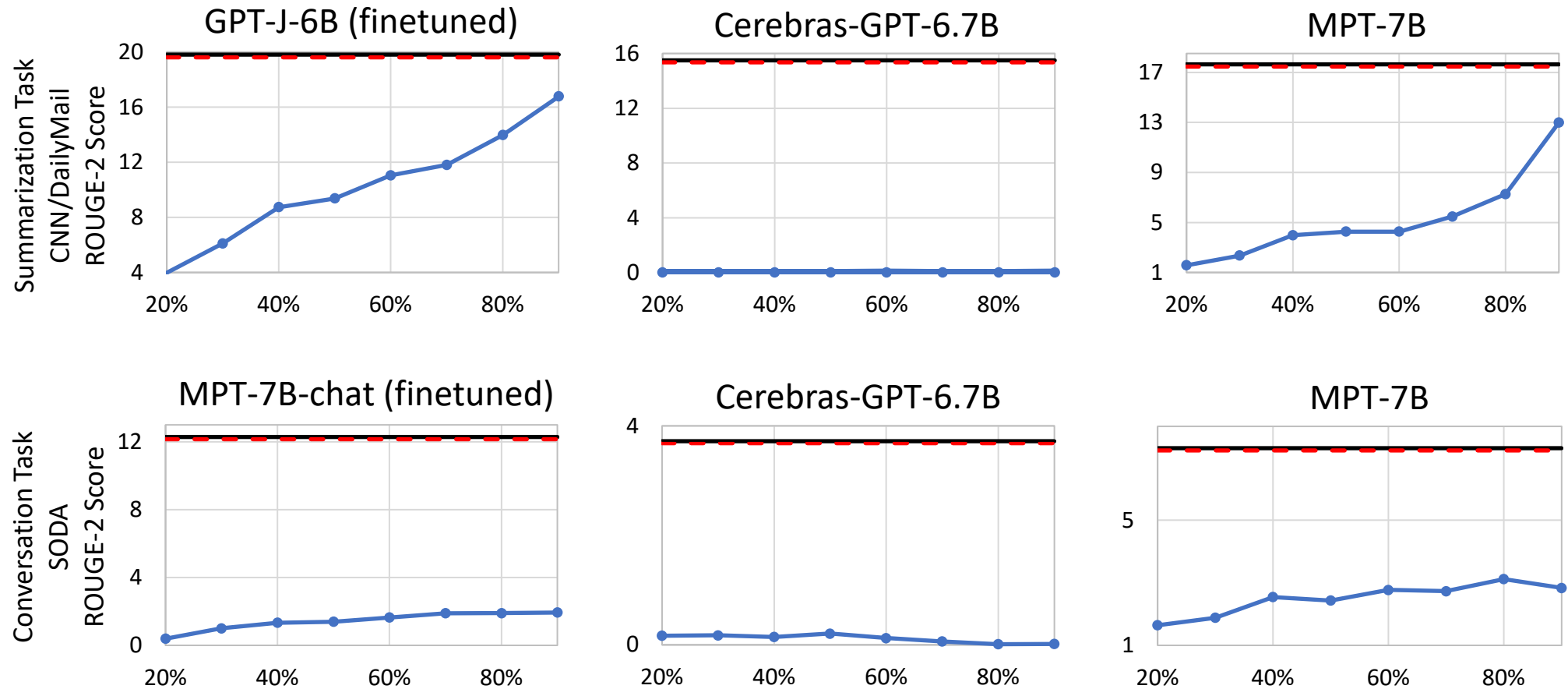Full Attention
Window Attention
$H_2O$

NVIDIA A100 (80GB)

# Accuracy Tradeoff



Legend: Full Attention — 99% Accuracy — Window Attention — H2O — Keyformer
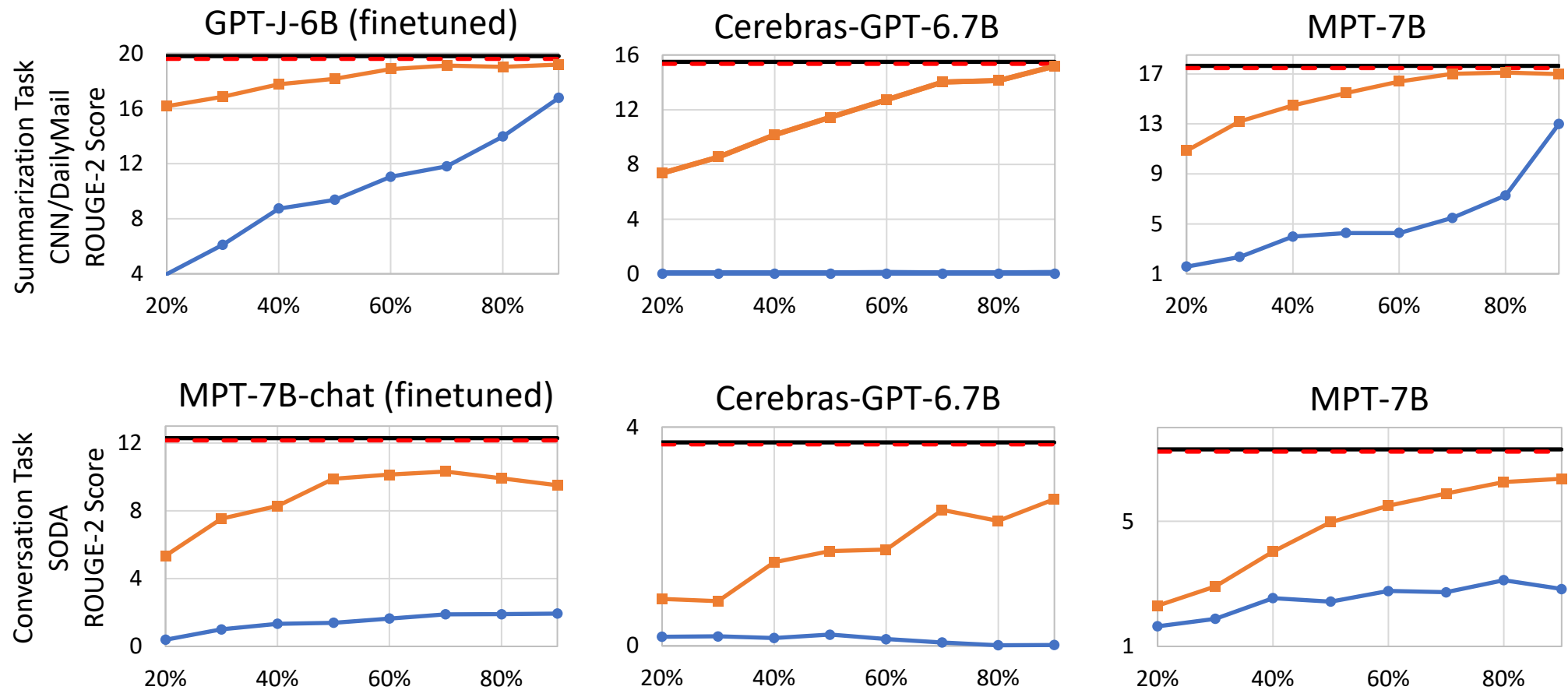
# Accuracy Tradeoff

# Accuracy Tradeoff



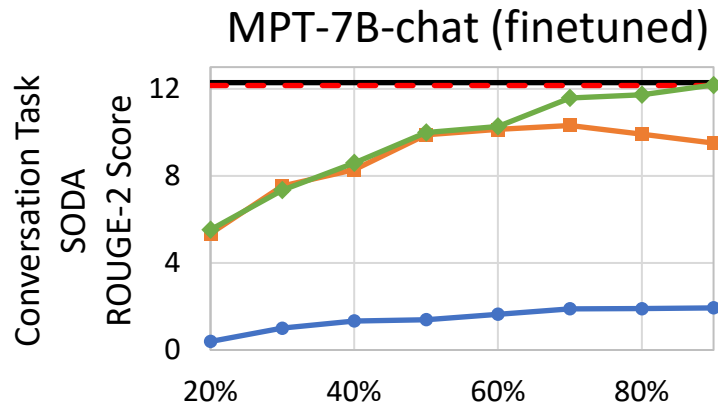Full Attention — 99% Accuracy — Window Attention — H2O — Keyformer

# Accuracy Tradeoff



30

# Accuracy Tradeoff

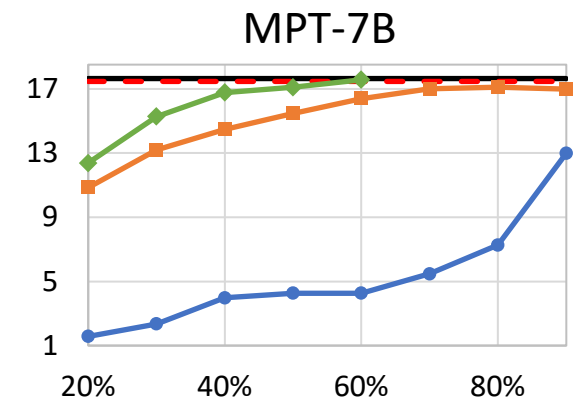—Full Attention    – – 99% Accuracy    —Window Attention    —H2O    —Keyformer

**Summarization Task → 70% of prompt KV cache**
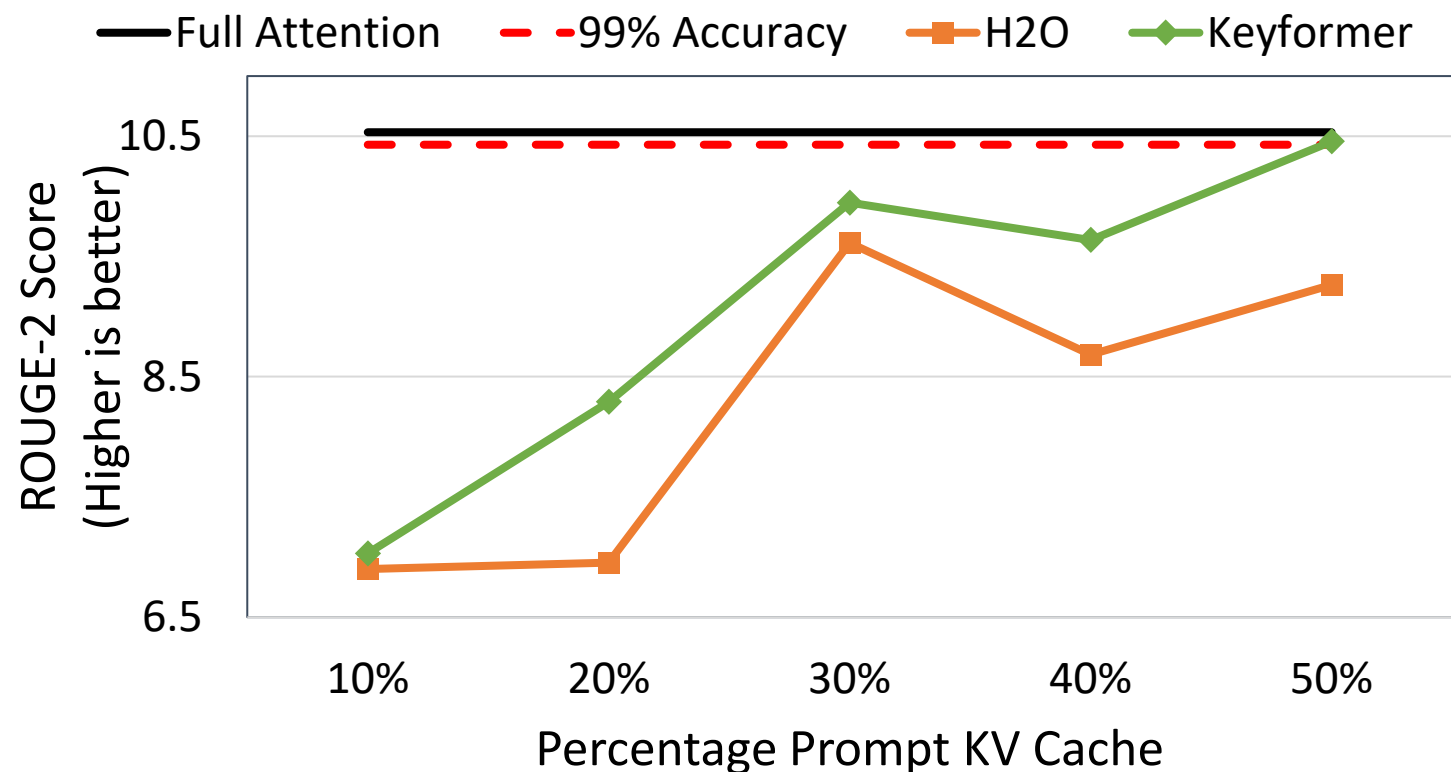
**Conversation Task → 90% of prompt KV cache**

MPT-7B-chat (finetuned)    Cerebras-GPT-6.7B    MPT-7B

# Long Context
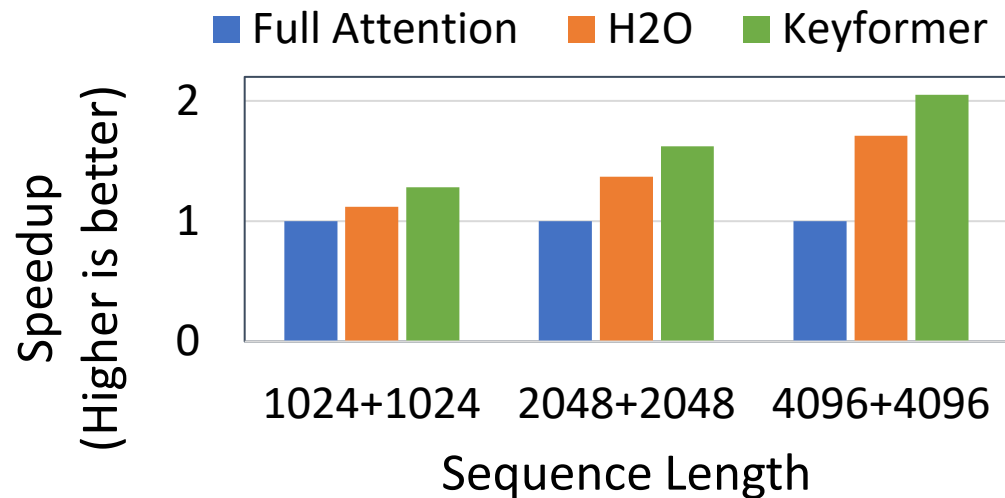
**Model** : MPT-7B-storywriter (65k seq len)

**Dataset** : Gov Reports (Mean Context: 9k)

# Long Context

**Model**   : MPT-7B-storywriter (65k seq len)
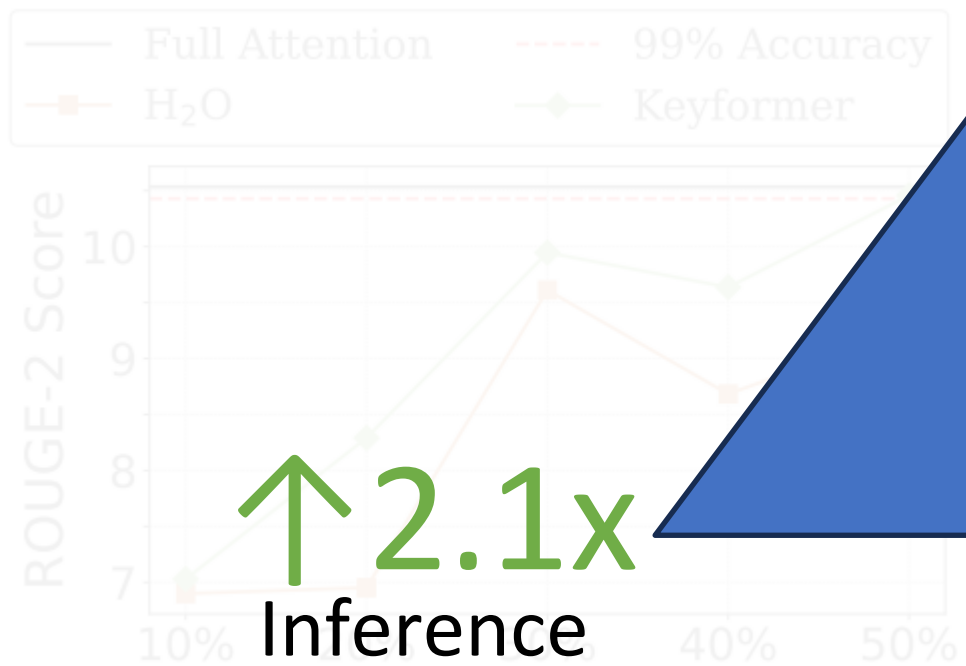**Dataset** : Gov Reports (Mean Context: 9k)



Speedup (Higher is better) vs Sequence Length — Full Attention, H2O, Keyformer

### Generation Throughput (tokens/sec)

| Sequence Length | Full Attention Original cache | H₂O 90% KV cache | Keyformer 50% KV cache |
|---|---|---|---|
| 1024 + 1024 | 24.9 | 27.8 | 32.0 |
| 2048 + 2048 | 15.0 | 20.5 | 24.3 |
| 4096 + 4096 (BS=1) | 8.3 | 14.1 | 17.0 |
| 4096 + 4096 (BS=2) | OOM | OOM | 19.85 |

# Long Context



> 99%
Accuracy

↑2.1x
Inference
Speedup

↑2.4x
Token
Generation
Throughput

# Conclusion

- LLMs → Inherent sparsity within attention
- Sparsity → KV cache reduction
- Token Discarding → Effect on attention score distribution
- Keyformer → Regularization based Key token Identification

**> 99%**
Accuracy

**↑2.1x**
Inference
Speedup

**↑2.4x**
Token
Generation
Throughput

# Questions