# InfiniGen

# Efficient Generative Inference of Large Language Models with Dynamic KV Cache Management

**Wonbeom Lee**[†]   Jungi Lee[†]   Junghwan Seo   Jaewoong Sim

**Seoul National University**

[†]Co-first Authors

ARTIFACT EVALUATED usenix ASSOCIATION AVAILABLE

ARTIFACT EVALUATED usenix ASSOCIATION FUNCTIONAL

ARTIFACT EVALUATED usenix ASSOCIATION REPRODUCED
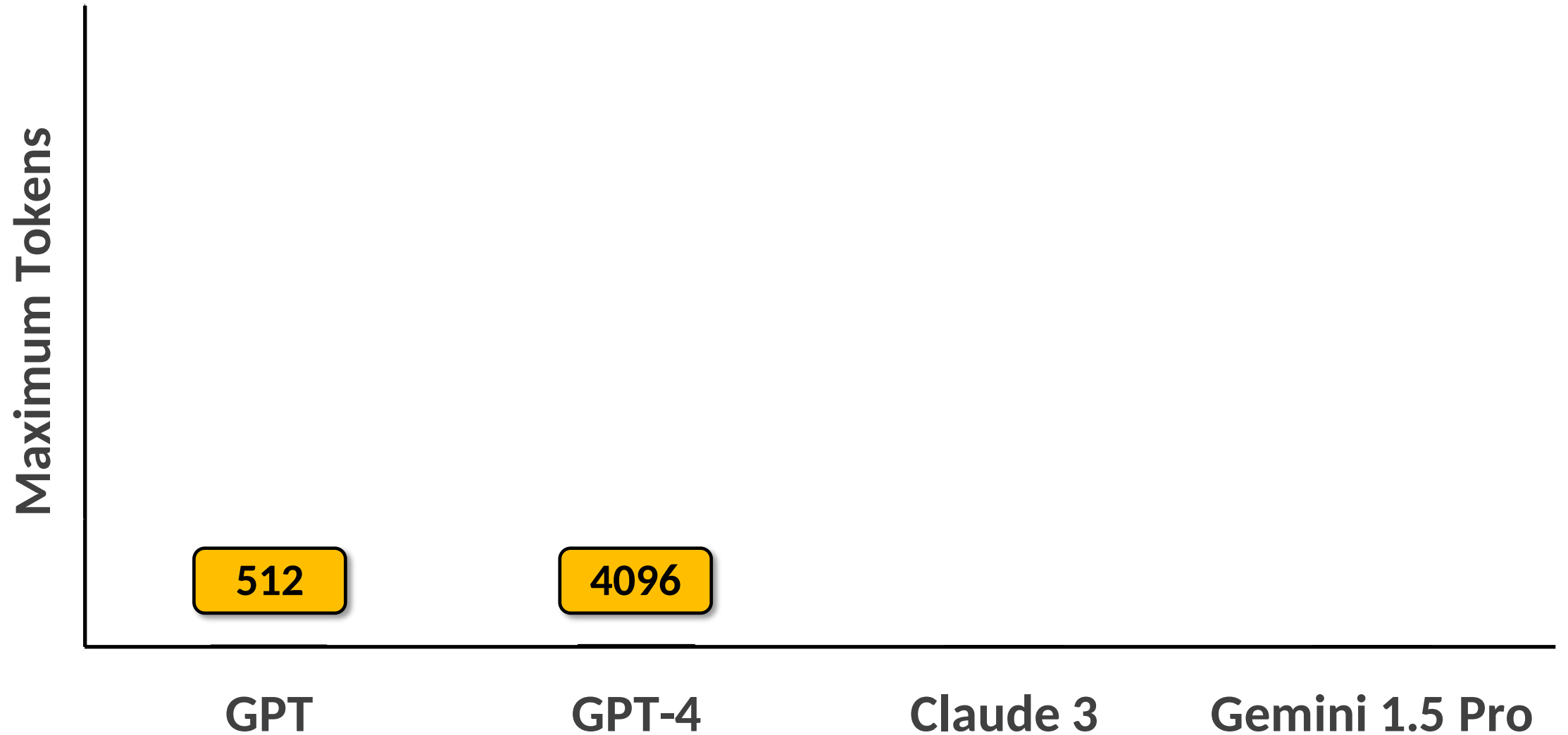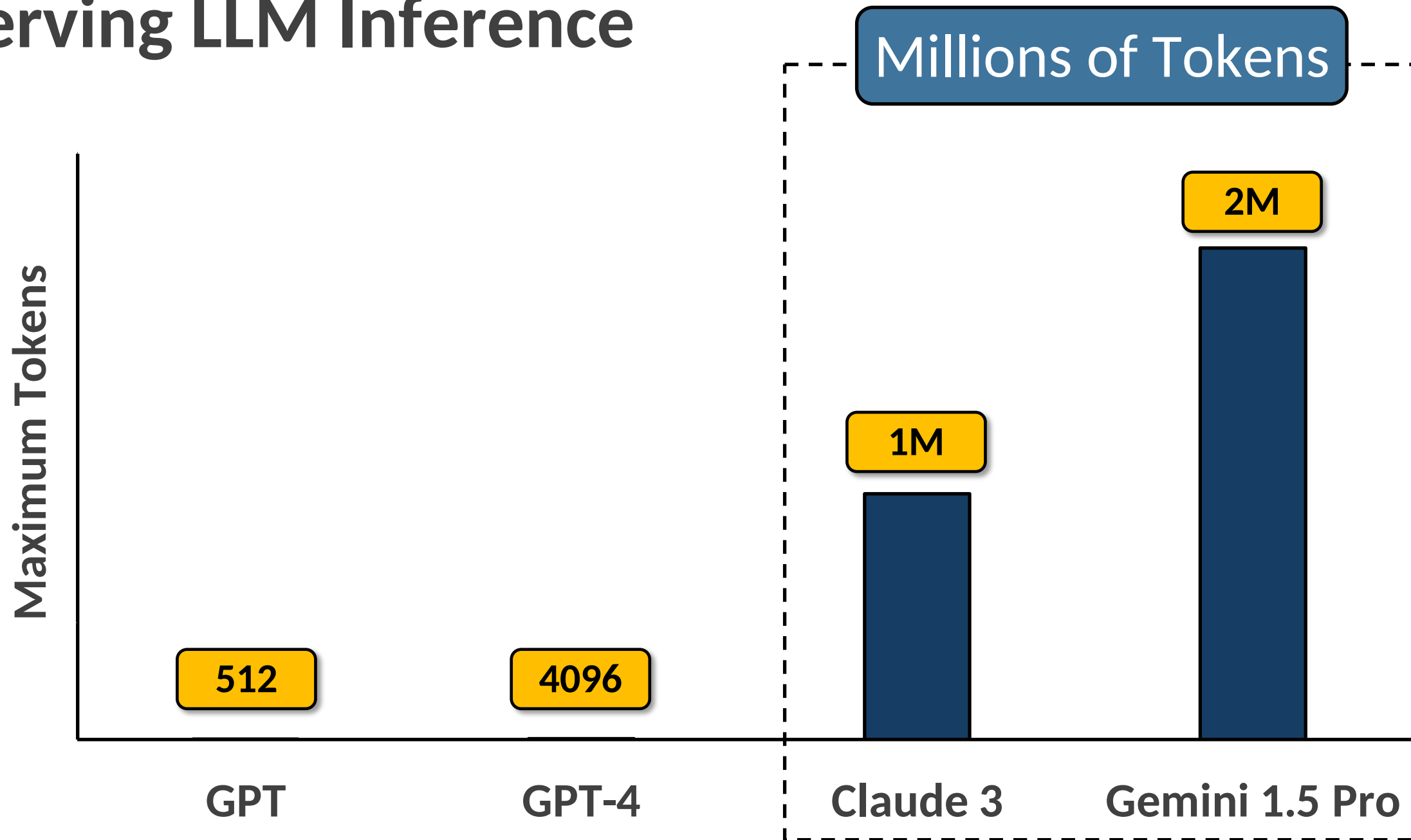
# Outline

- **LLM Inference & KV Cache**

- **Prior Approaches & Limitations**

- **InfiniGen: Dynamic KV Cache Management**
  - Speculative KV Prefetching
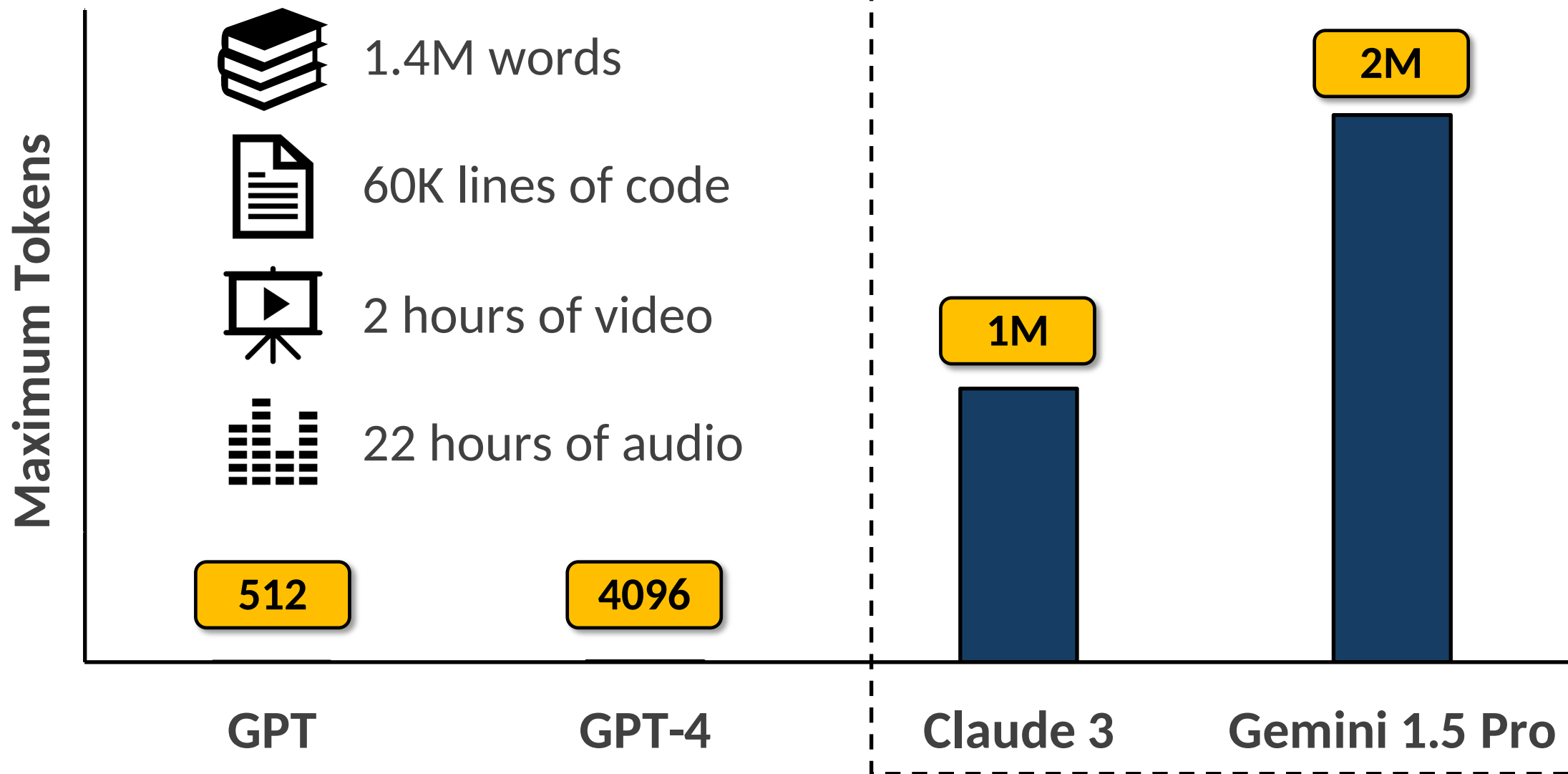  - Key/Query Skewing

- **Evaluation**

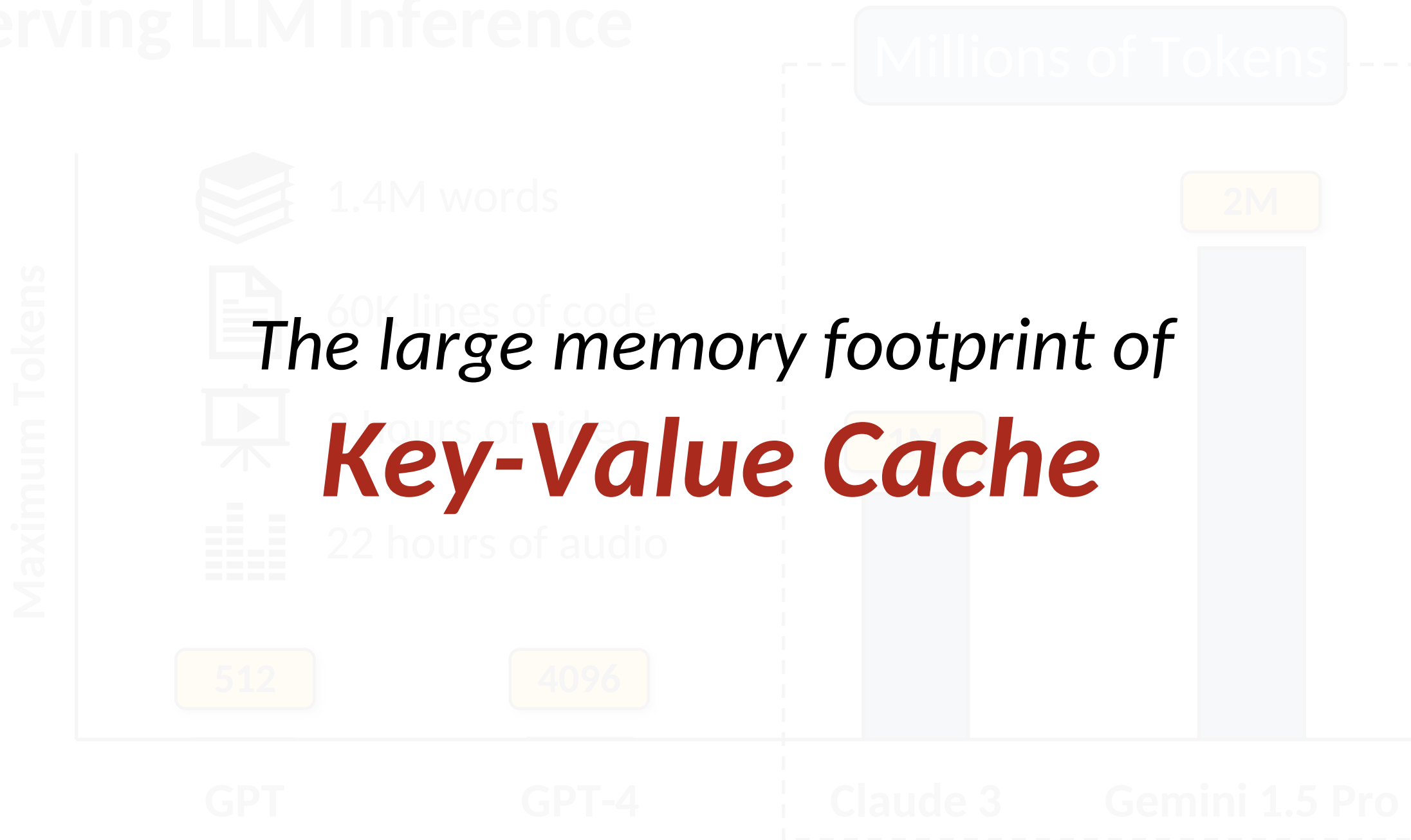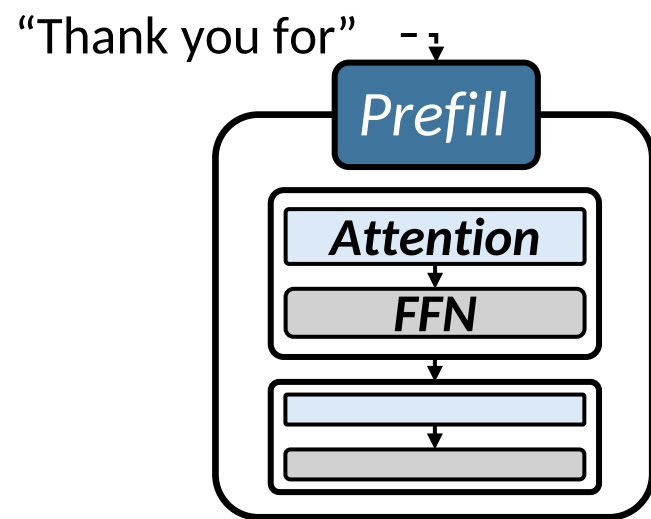- **Conclusion**

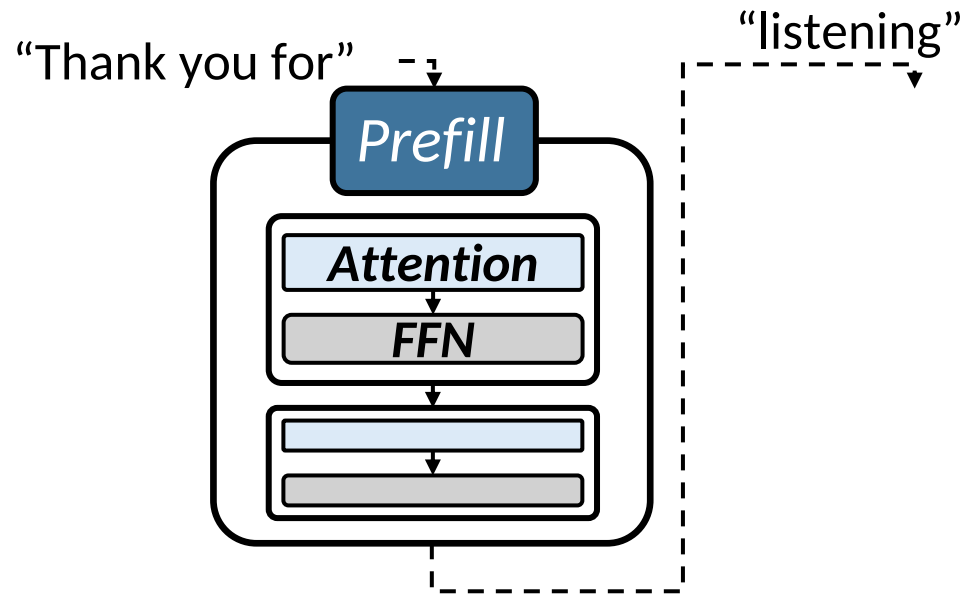# Serving LLM Inference

# Serving LLM Inference

# Serving LLM Inference



Millions of Tokens

- 1.4M words
- 60K lines of code
- 2 hours of video
- 22 hours of audio

Maximum Tokens

| GPT | GPT-4 | Claude 3 | Gemini 1.5 Pro |
| 512 | 4096 | 1M | 2M |

*The large memory footprint of*
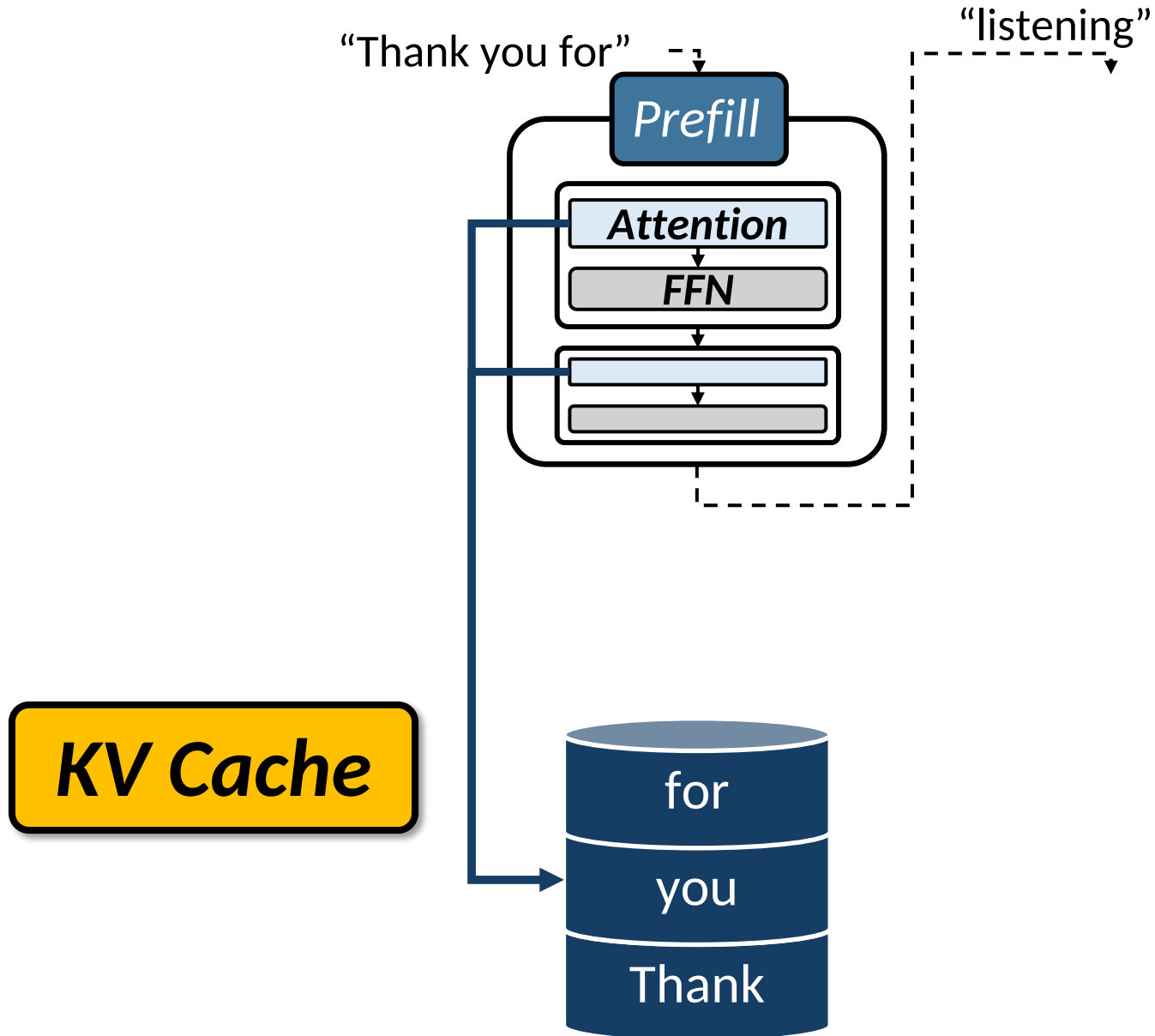***Key-Value Cache***
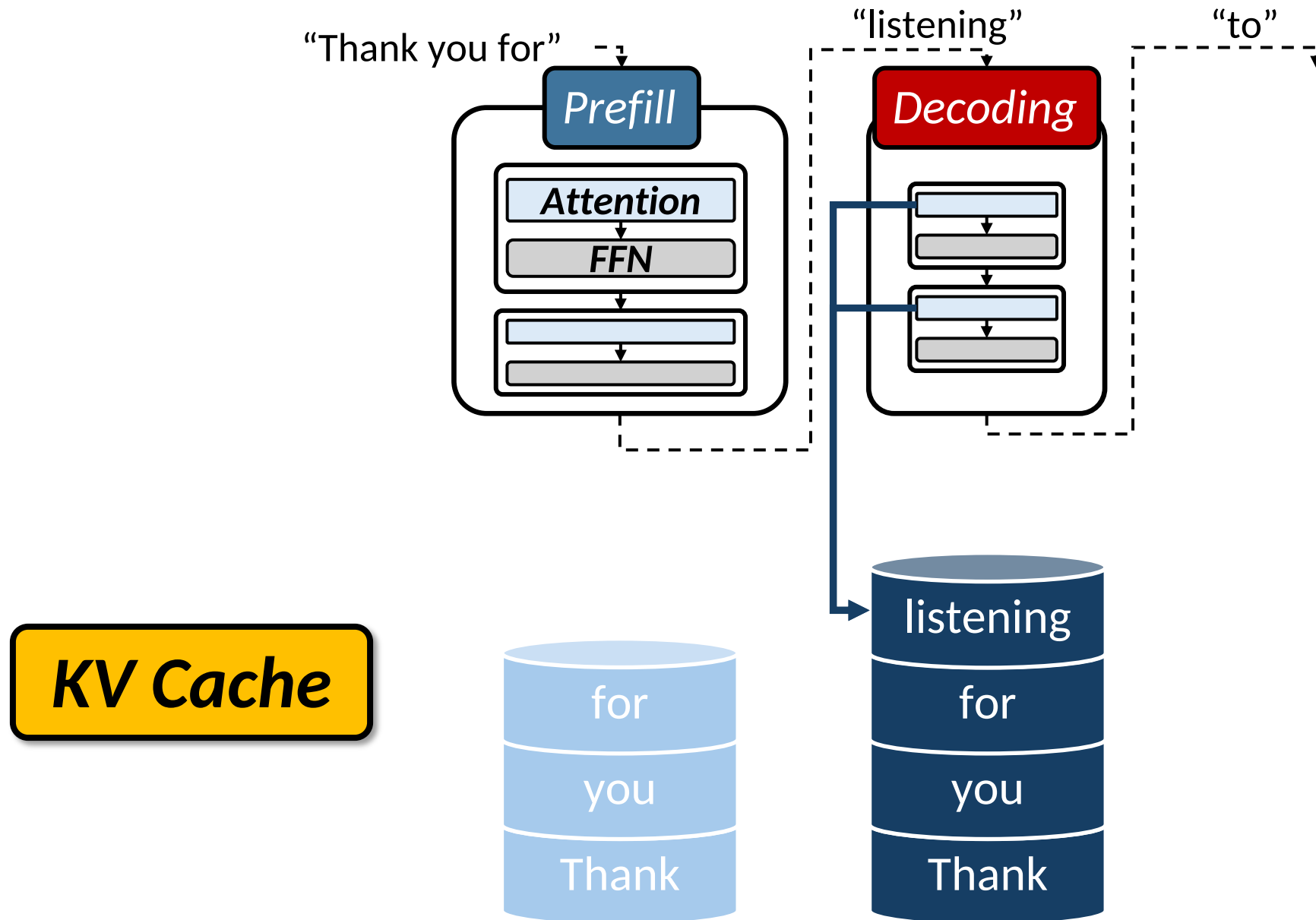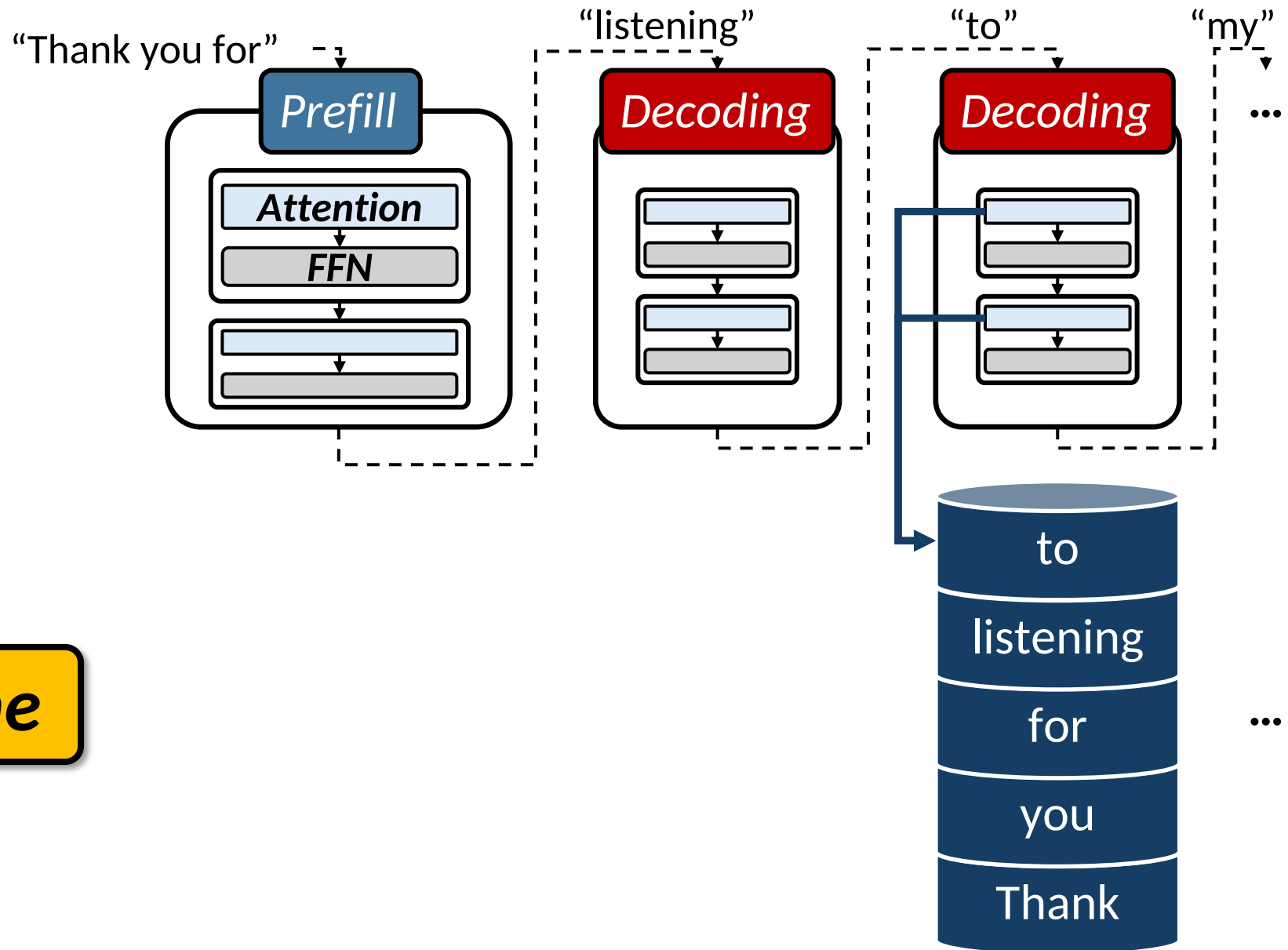
# KV Cache in LLM Inference

"Thank you for"

# KV Cache in LLM Inference
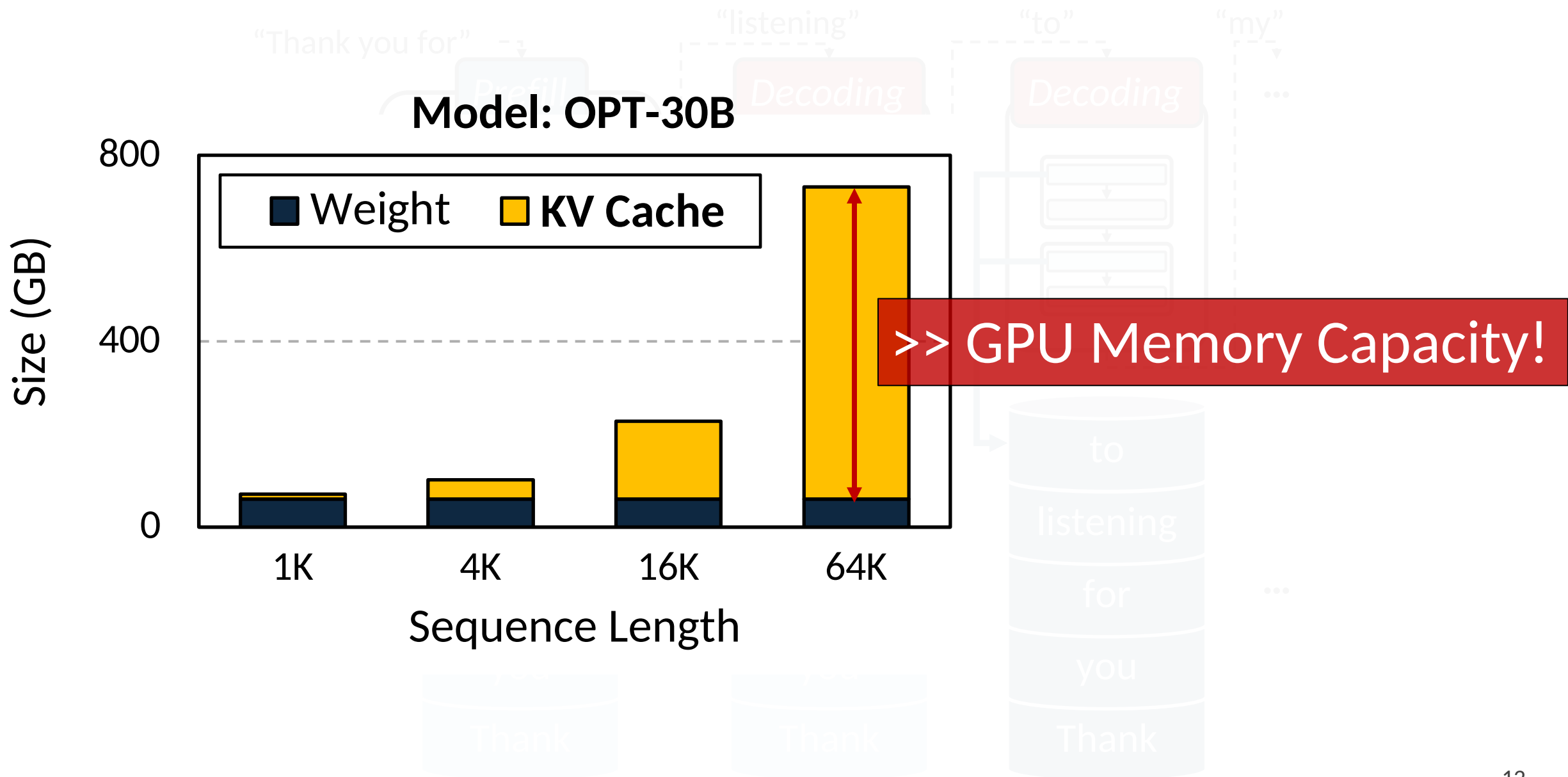
# KV Cache in LLM Inference

# KV Cache in LLM Inference

# KV Cache in LLM Inference

# KV Cache in LLM Inference



**Model: OPT-30B**

>> GPU Memory Capacity!

# Prior Approaches for KV Cache Problem

**1. Quantization**



Compress the **KV cache**
by converting it into **a low-bit data format**
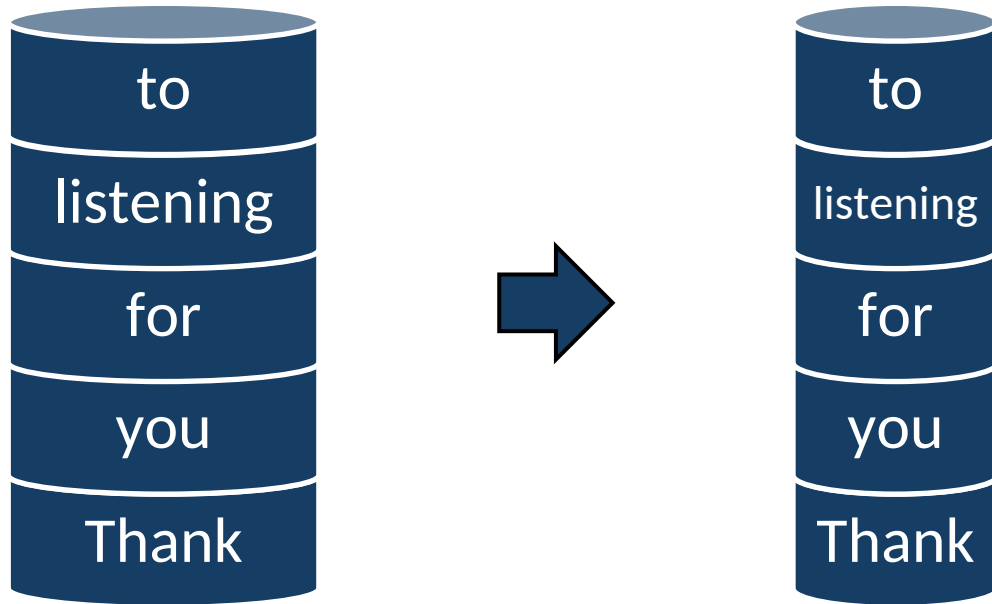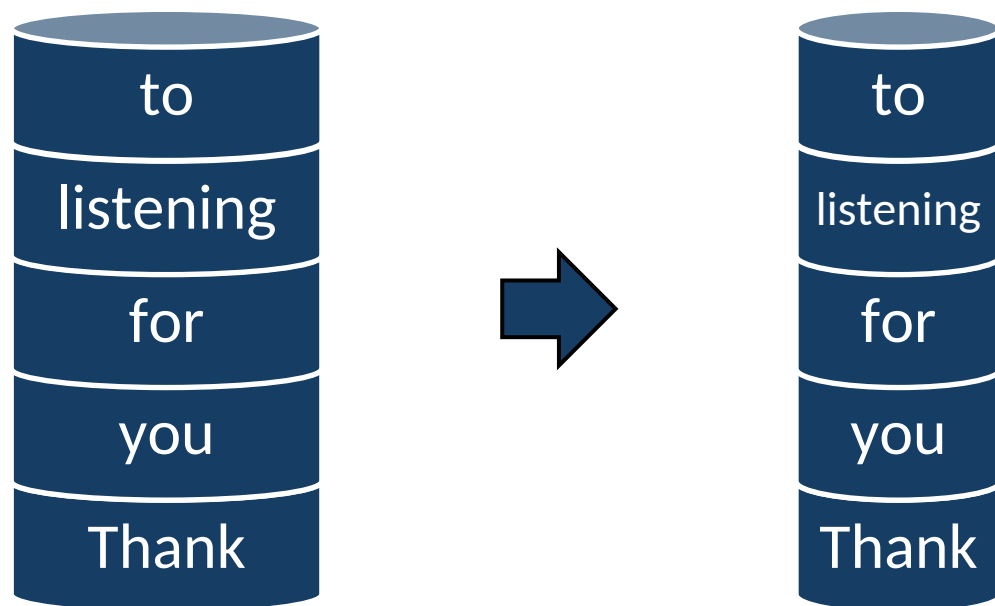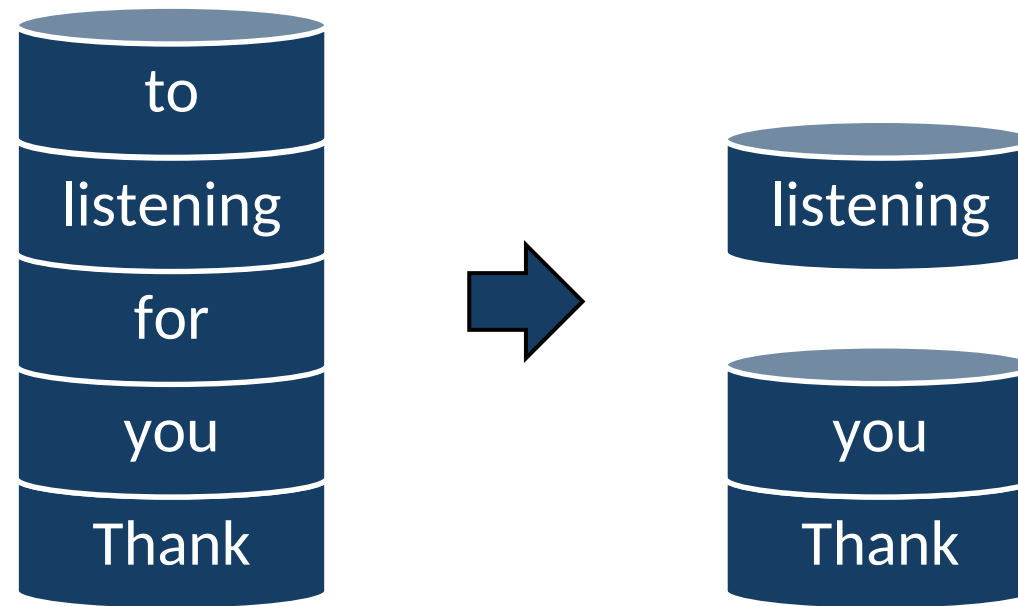
# Prior Approaches for KV Cache Problem

## 1. Quantization



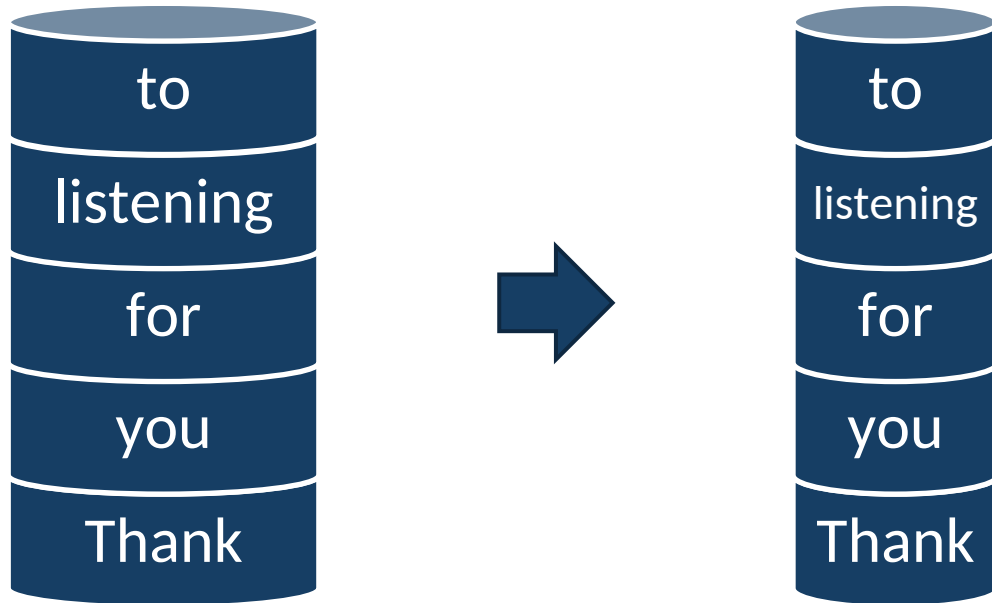Compress the **KV cache** by converting it into **a low-bit data format**

## 2. Eviction



**Permanently eliminate** **unimportant tokens** to keep the KV cache size in check

# Limitations of Prior Approaches

## 1. Quantization

to
listening
for
you
Thank

➡

to
listening
for
you
Thank

Compress the **KV cache**
by converting it into **a low-bit data format**
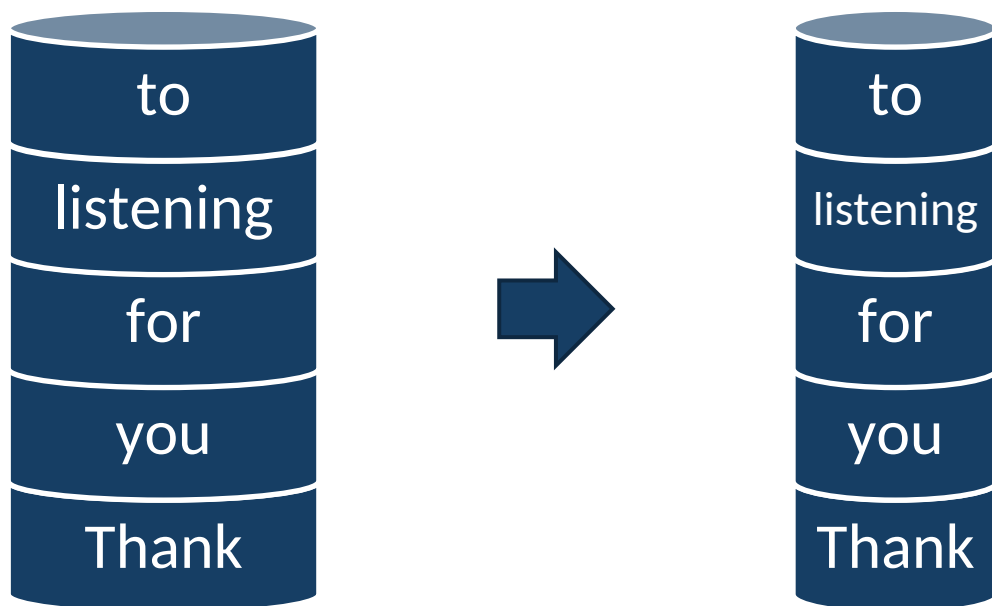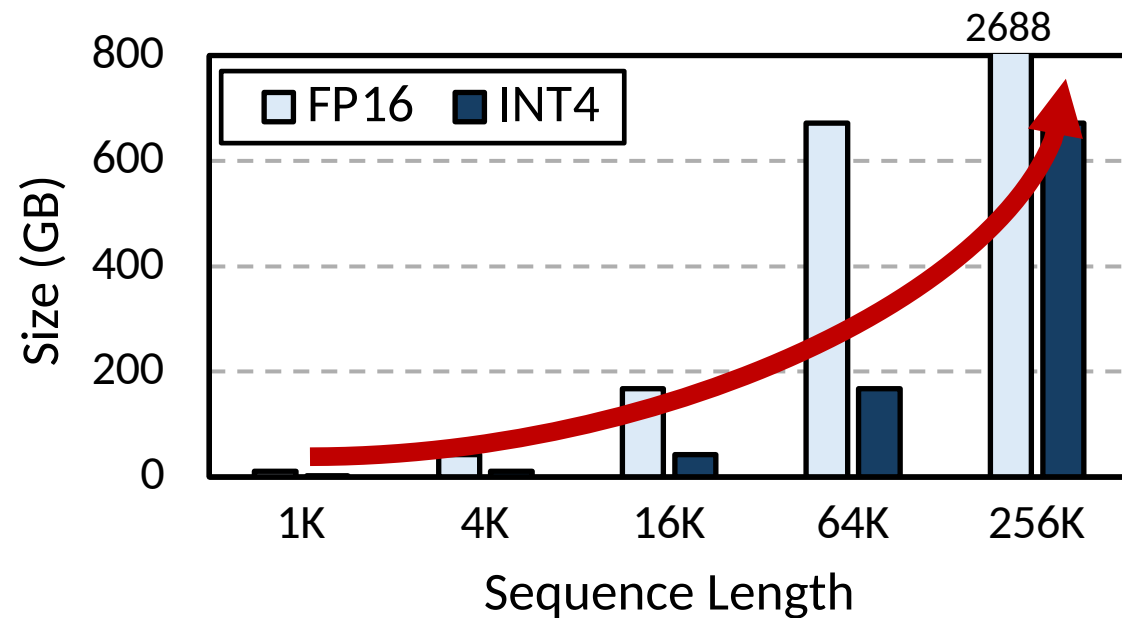
## 2. Eviction

that
saying
for
you
Thank

➡

saying

you

**Problem 1**

KV cache access still
**linearly increases** with the sequence length

# Limitations of Prior Approaches

## 1. Quantization

to
listening
for
you
Thank

→

to
listening
for
you
Thank

Compress the **KV cache**
by converting it into **a low-bit data format**

## 2. Eviction

2688

Size (GB)

□ FP16 ■ INT4

800
600
400
200
0

1K    4K    16K    64K    256K

Sequence Length

**Problem 1**

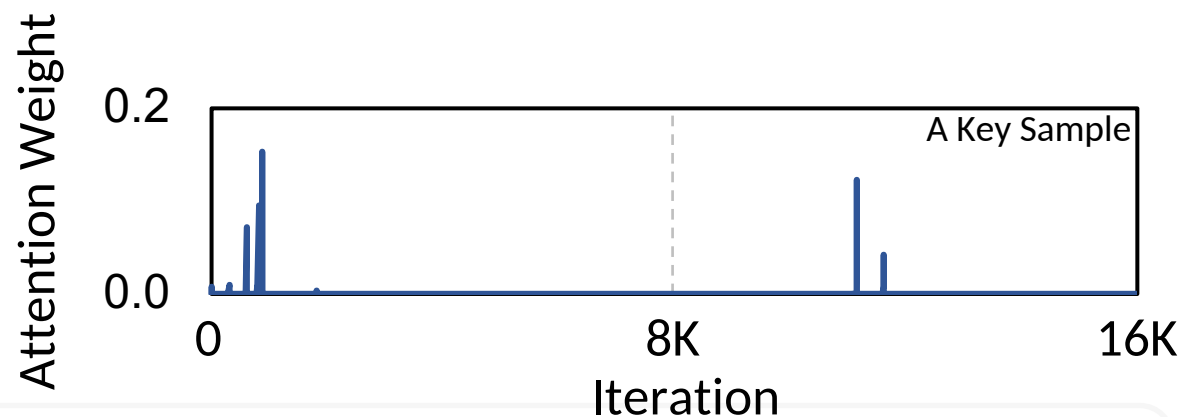KV cache access still
**linearly increases** with the sequence length
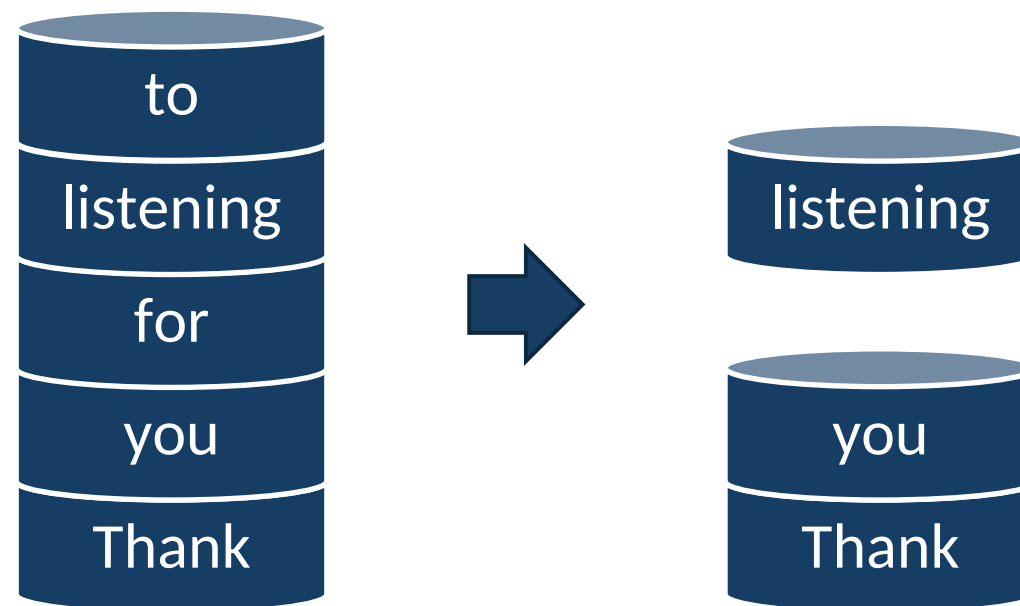
# Limitations of Prior Approaches

## 1. Quantization

**Problem 2**

Permanent elimination of tokens can lead to an **accuracy drop**



Attention Weight

0.2

0.0

0    8K    16K

Iteration

A Key Sample

Compress the KV Cache
by quantizing into Low-bit Data Format

## 2. Eviction



to
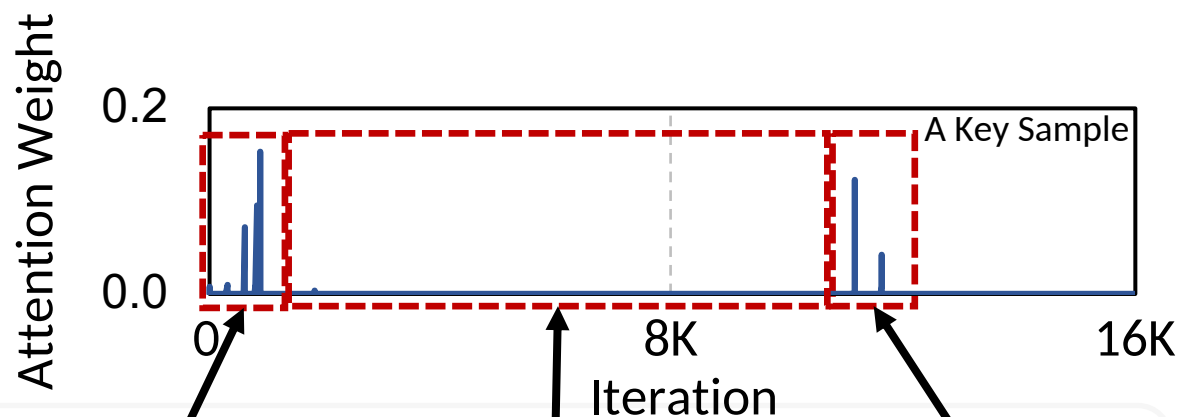listening
for
you
Thank

→

listening

you
Thank

**Permanently eliminate** unimportant tokens
to keep the KV cache size in check

# Limitations of Prior Approaches

## 1. Quantization

**Problem 2**

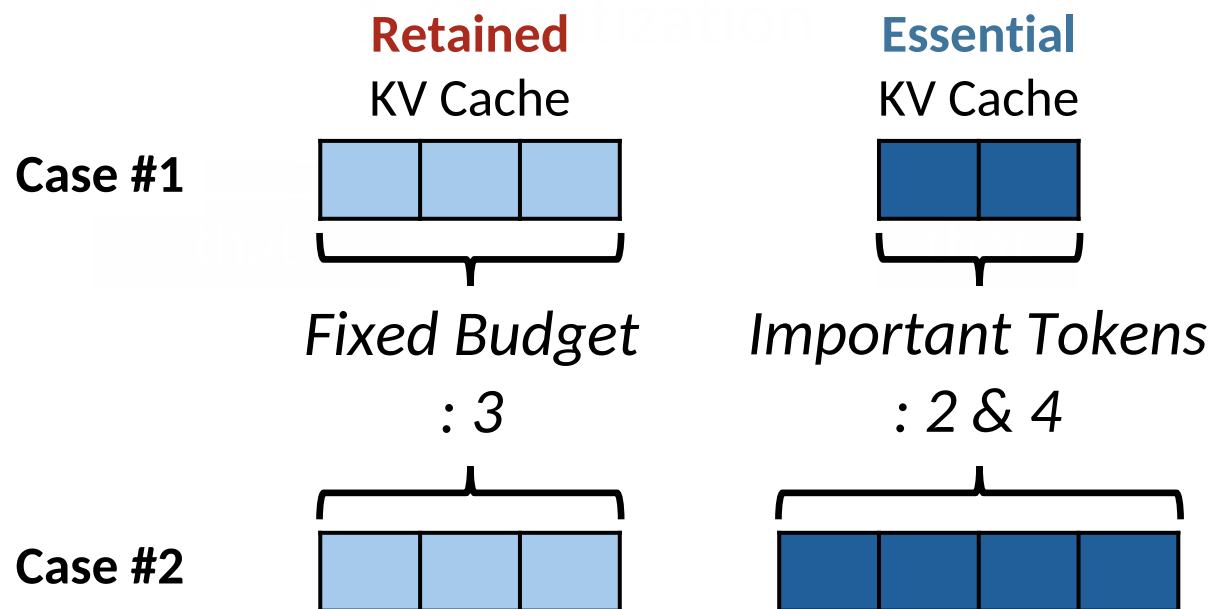Permanent elimination of tokens can lead to an **accuracy drop**



Attention Weight

0.2

0.0

0  8K  16K
Iteration

A Key Sample

*Important*  *Unimportant*  *Important*

Compress the KV Cache

## 2. Eviction



to
listening
for
you
Thank

→

listening

you
Thank

**Permanently eliminate** unimportant tokens to keep the KV cache size in check

# Limitations of Prior Approaches

**Retained**
KV Cache

**Essential**
KV Cache

**2. Eviction**

Case #1

*Fixed Budget
: 3*

*Important Tokens
: 2 & 4*

Case #2

to
listening
for
you
Thank

→

listening
you
Thank

**Problem 3**

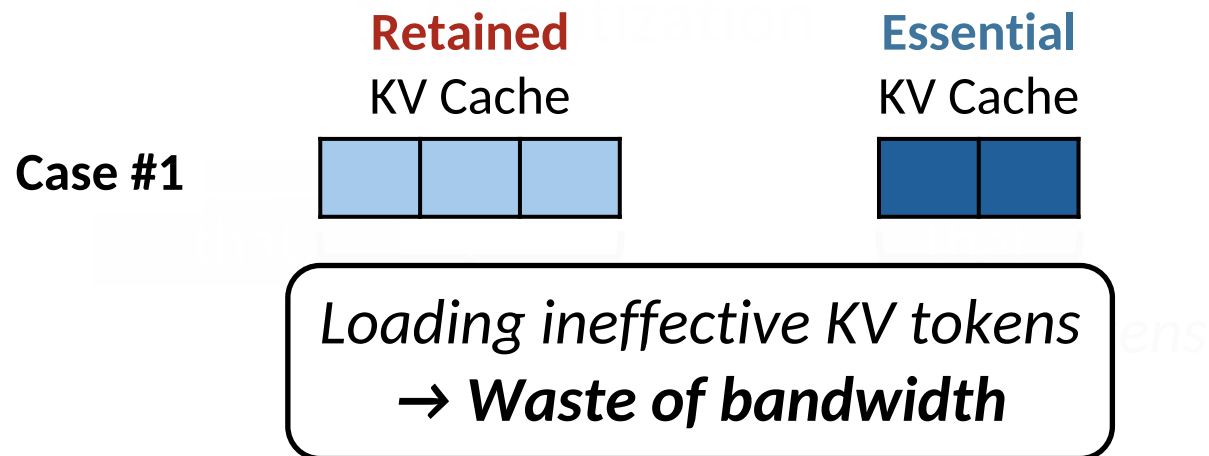Fixed KV cache budget
can lead to **subpar performance**

**Permanently eliminate** unimportant tokens
to keep the KV cache size in check

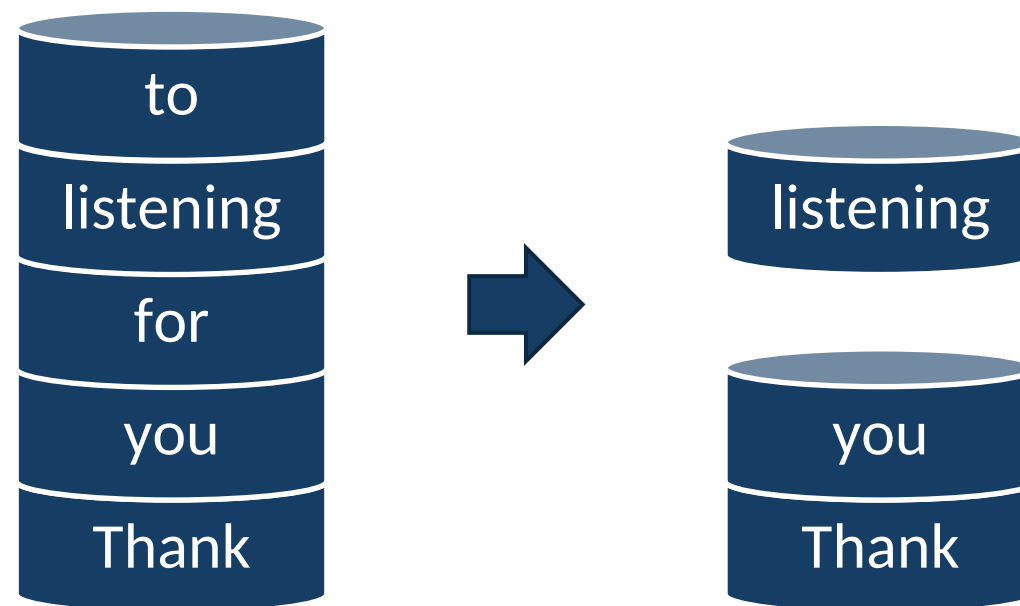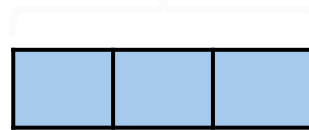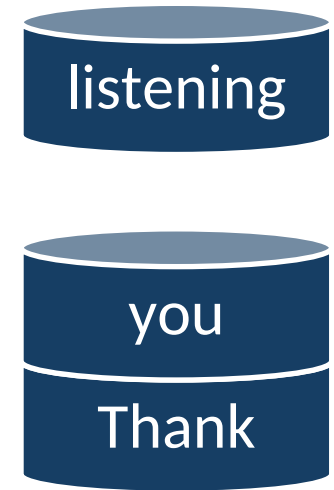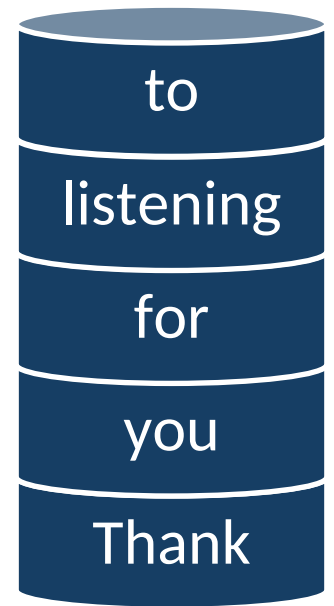# Limitations of Prior Approaches

**Case #1**

**Retained** KV Cache

**Essential** KV Cache

*Loading ineffective KV tokens* → **Waste of bandwidth**

**Problem 3**

Fixed KV cache budget can lead to **subpar performance**

## 2. Eviction

to
listening
for
you
Thank

→

listening

you
Thank

**Permanently eliminate** unimportant tokens to keep the KV cache size in check

# Limitations of Prior Approaches

**Retained**
KV Cache

**Essential**
KV Cache

Case #1

*Missing essential KV tokens*
**→ Accuracy drop**

**Case #2**

**Problem 3**

Fixed KV cache budget
can lead to **subpar performance**

## 2. Eviction



to
listening
for
you
Thank

→

listening

you
Thank

**Permanently eliminate** unimportant tokens
to keep the KV cache size in check

# Prior Approaches

**Problem 1**

KV cache access still linearly increases
with the sequence length

**Problem 2**

Permanent elimination of tokens
can lead to an accuracy drop

**Problem 3**

Fixed KV cache budget
can lead to subpar performance

# Prior Approaches

*Problem 1*

KV cache access still linearly increases
with the sequence length

**Not a scalable nor effective solution
in an era of millions of tokens!**

*Problem 3*

Fixed KV cache budget
can lead to subpar performance

# Outline

- LLM Inference & KV Cache

- Prior Approaches & Limitations

- **InfiniGen: Dynamic KV Cache Management**
  - Speculative KV Prefetching
  - Key/Query Skewing

- Evaluation

- Conclusion

# InfiniGen: Key Direction

Problem 1

Memory access overhead still
linearly scales with the sequence length

## Exploit the abundant CPU memory capacity to manage the KV cache!

Problem 3

Fixed KV cache budget
can lead to subpar performance

# KV Cache Offloading

GPU

GPU Memory

KV Cache

Expensive & Small

# KV Cache Offloading



PCIe

**CPU**

**GPU**

CPU Memory

**KV Cache**

GPU Memory

Cheap & Large

Expensive & Small

# KV Cache Offloading



CPU

PCIe

GPU

CPU Memory

**KV Cache**

GPU Memory

Cheap & Large

Expensive & Small

# KV Cache Offloading



PCIe

**CPU**

*KV Cache*

CPU Memory

GPU Memory

Cheap & Large

Expensive & Small

# KV Cache Offloading



Latency

***Significant slowdown
due to the limited PCIe bandwidth***

# KV Cache Offloading

**Transfer Less**

**Transfer Early**

# KV Cache Offloading

**Load and compute only with a few important tokens**

**Prefetch essential KV entries in the preceding layer**

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching



GPU

Layer $i - 1$

**Predict** | Attention | FFN →

Layer $i$

**Predict** | Light Attention | FFN → ...

time

CPU

Prefetching

Corresponding KV Cache

Overlap

Small Data Transfer

**Challenge**

## Q: How to predict the important tokens?

# Speculative KV Prefetching



**Original Attention:** *Layer $i$*

## Query Projection

Attention Input $\times$ Query Weight

$=$ Query

## Attention

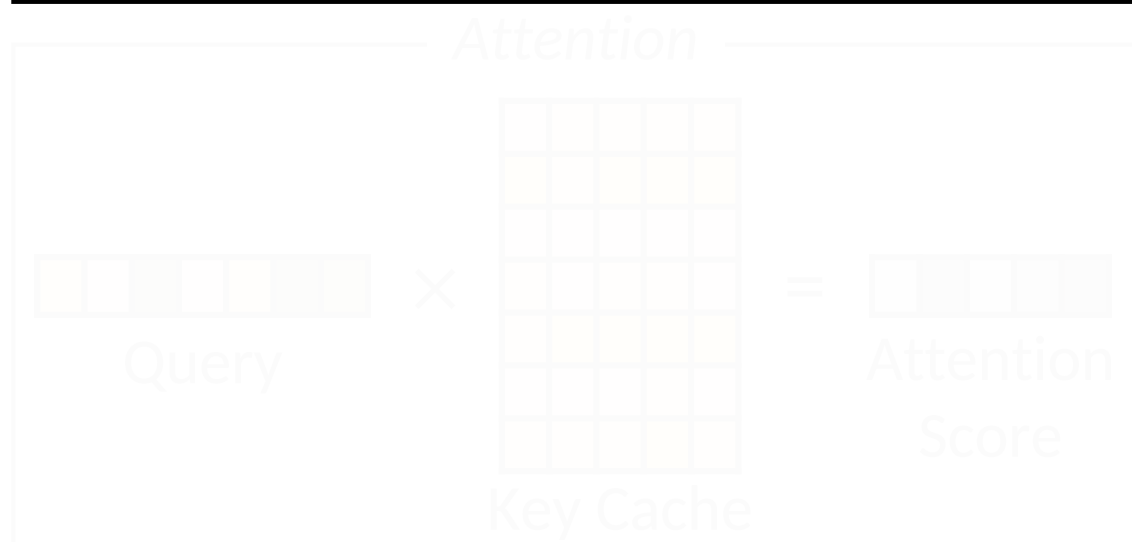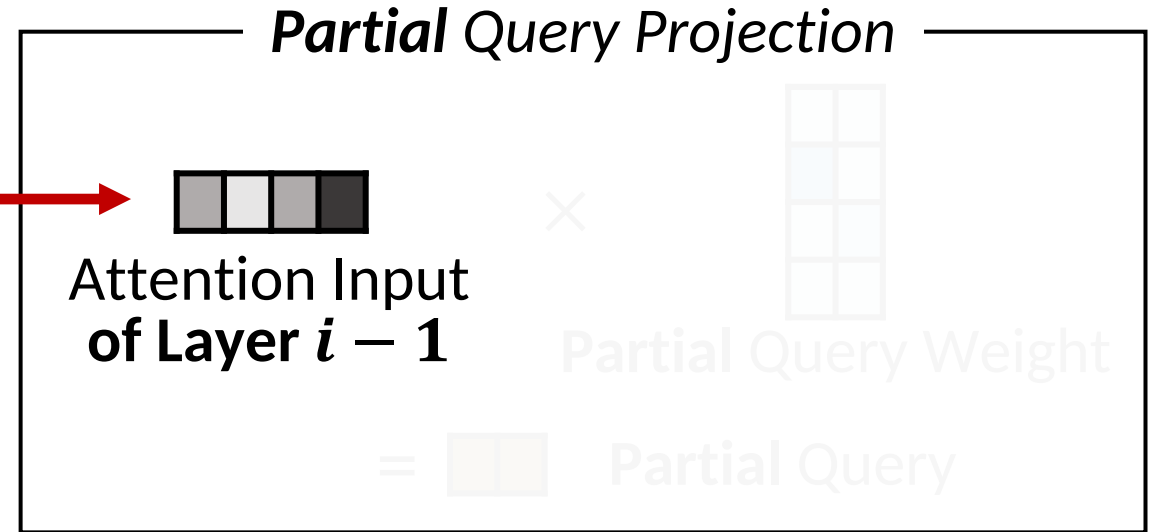Query $\times$ Key Cache $=$ Attention Score

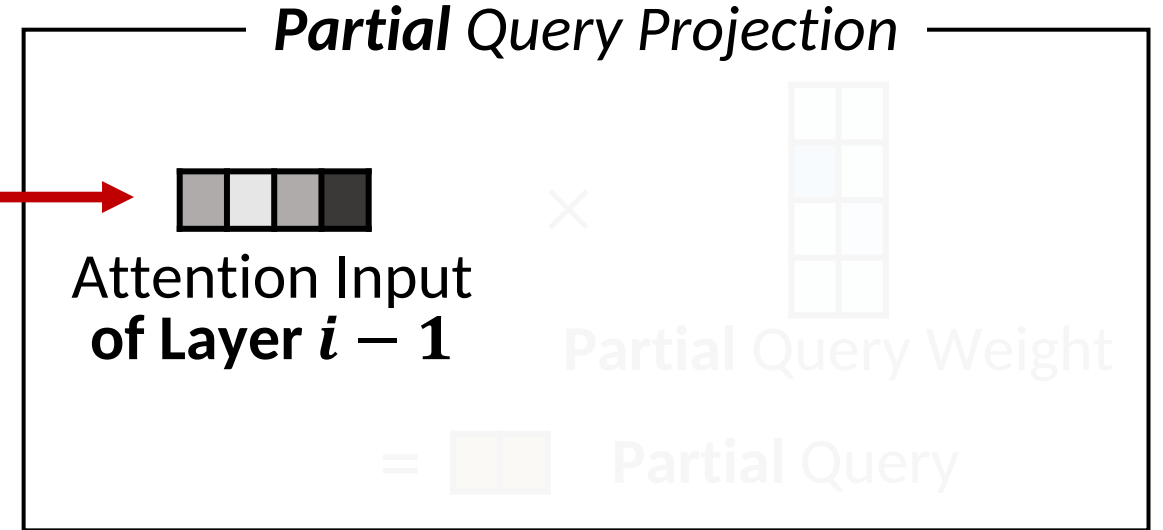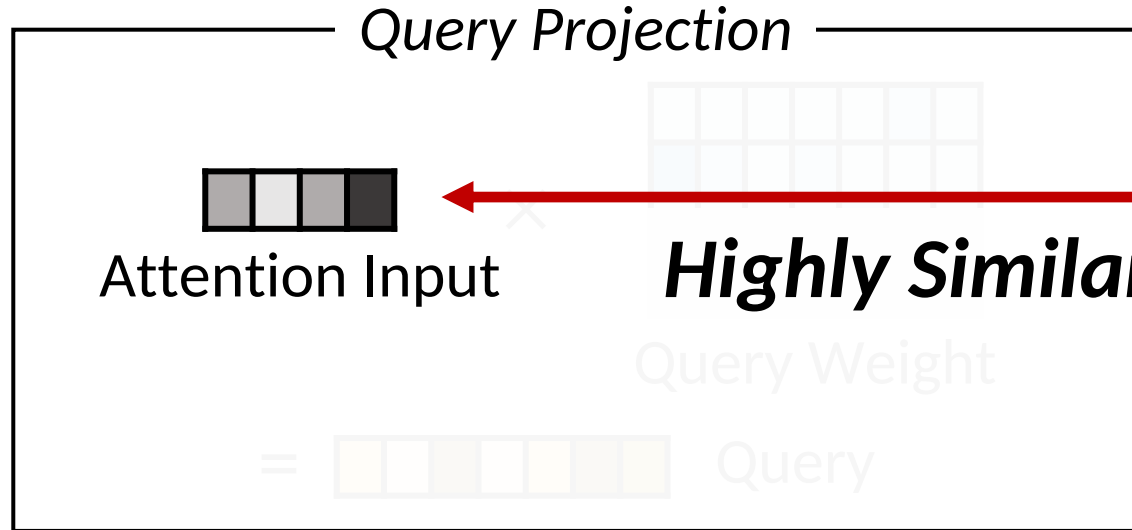**Minimal Rehearsal:** *Layer $i-1$*

## *Partial* Query Projection

Attention Input **of Layer $i-1$** $\times$ **Partial** Query Weight

$=$ **Partial** Query

## Attention *Speculation*

**Partial** Query $\times$ **Partial** Key Cache

$=$ **Speculated** Attention Score

# Speculative KV Prefetching

Minimal Rehearsal: Layer $i-1$

## Query Projection

Attention Input

Query Weight

Query

## Partial Query Projection

Attention Input
of Layer $i-1$

Partial Query Weight

Partial Query

## Attention

Query

Key Cache

Attention
Score

## Attention Speculation

Partial Query

Partial Key Cache

Speculated
Attention Score
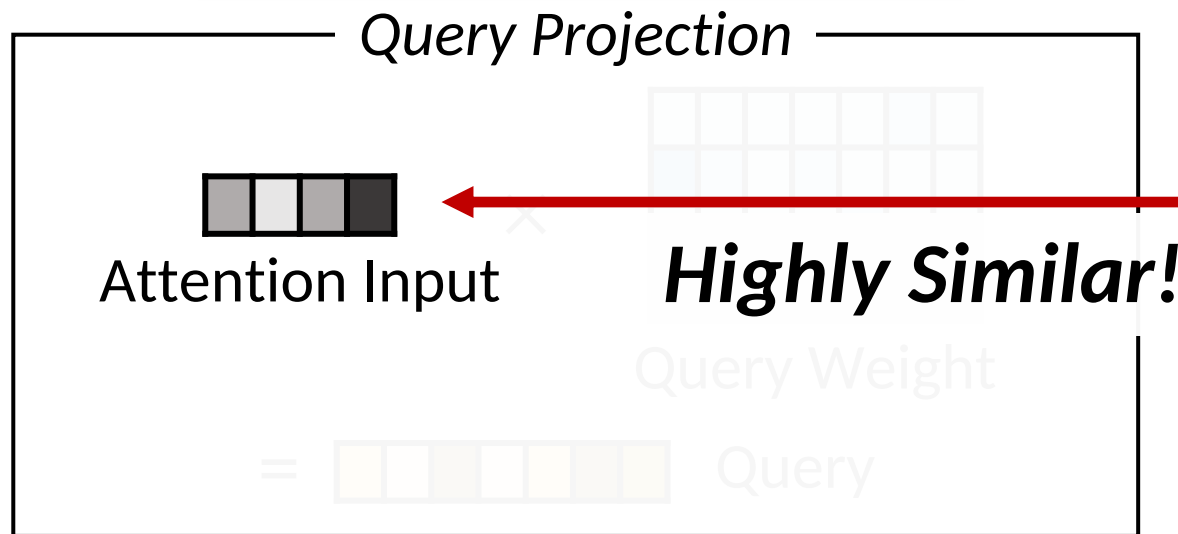
40

# Speculative KV Prefetching



**Original Attention:** *Layer $i$*

Query Projection

Attention Input

Query Weight

Attention

Key Cache

**Minimal Rehearsal:** *Layer $i-1$*

***Partial* Query Projection**

Attention Input
of Layer $i-1$

**Partial** Query Weight

Attention **Speculation**

**Partial** Key Cache

41

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

# Speculative KV Prefetching

**Partial** Query Projection

Attention Input
of Layer $i - 1$

**Partial** Query Weight

*Skewing*

**Partial** Query

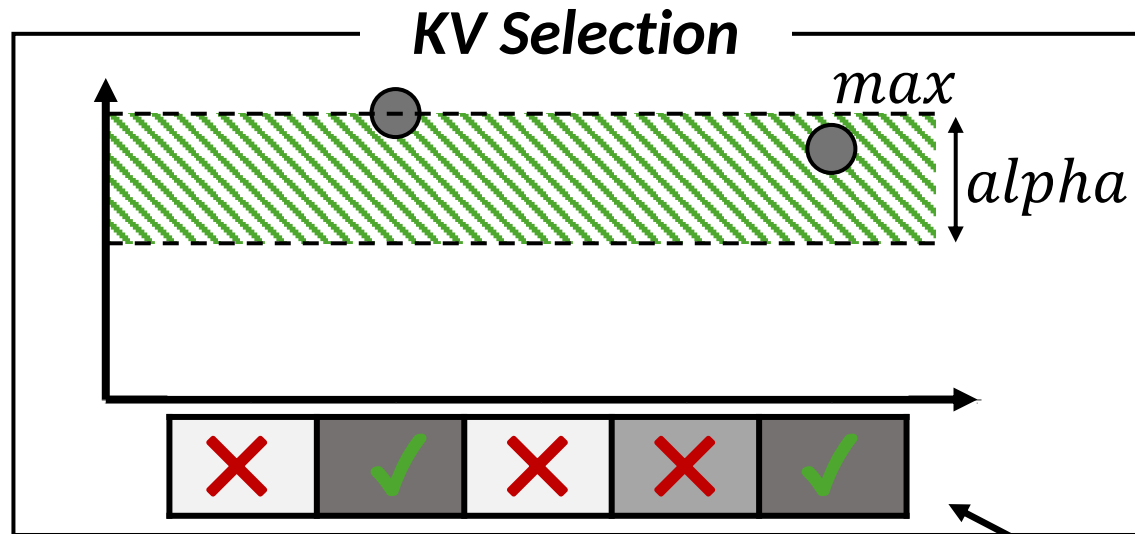KV Selection

$max$

Attention **Speculation**

**Partial** Query

**Partial** Key Cache

**Speculated**
Attention Score

49

# Key/Query Skewing

Query

| 1 | 4 | 1 | 4 | 3 |
|---|---|---|---|---|

Key

| 1 | 3 | 2 | 3 | 2 |
|---|---|---|---|---|

50

# Key/Query Skewing

Query

| 1 | 4 | 1 | 4 | 3 |

Partial Query

| 4 | 4 |

Key

| 1 | 3 | 2 | 3 | 2 |

Partial Key

| 3 | 3 |

51

# Key/Query Skewing

Query

| 1 | 4 | 1 | 4 | 3 |

×

Key

| 1 | 3 | 2 | 3 | 2 |

=    **33**

| 4 | 4 |

×

| 3 | 3 |

=    **24**

Partial Query

Partial Key

Attention Score

**9** ☹

# Key/Query Skewing

# Key/Query Skewing

# Key/Query Skewing

Before

**Offline modification** of the query and key weights using **singular value decomposition**

After

The **identical** computation result
$$(Q \times \boldsymbol{A}) \times (\boldsymbol{A^T} \times K^T) = Q \times K^T$$

# Key/Query Skewing

**Offline modification** of the query and key weights using **singular value decomposition**



(a) $Q = \mathbf{U}\Sigma\mathbf{V}^{\mathbf{T}}$

(b) $\widetilde{Q} = \mathbf{U}\Sigma(\mathbf{V}^{\mathbf{T}}A)$

# Outline

- LLM Inference & KV Cache

- Prior Approaches & Limitations

- InfiniGen: Dynamic KV Cache Management
  - Speculative KV Prefetching
  - Key/Query Skewing

- **Evaluation**

- Conclusion

# Experimental Setup

## Model
- **Open Pre-trained Transformer (OPT)**
  : 6.7B, 13B, 30B
- **Llama-2**
  : 7B, 13B

## Baseline
- KV Cache Offloading
  - CUDA Unified Virtual Memory (UVM)
  - **FlexGen**
- KV Cache Management Methods
  - **$H_2O$: Eviction-based**
  - Quantization

## Workload
- lm-evaluation-harness
- PG-19

### System Configuration

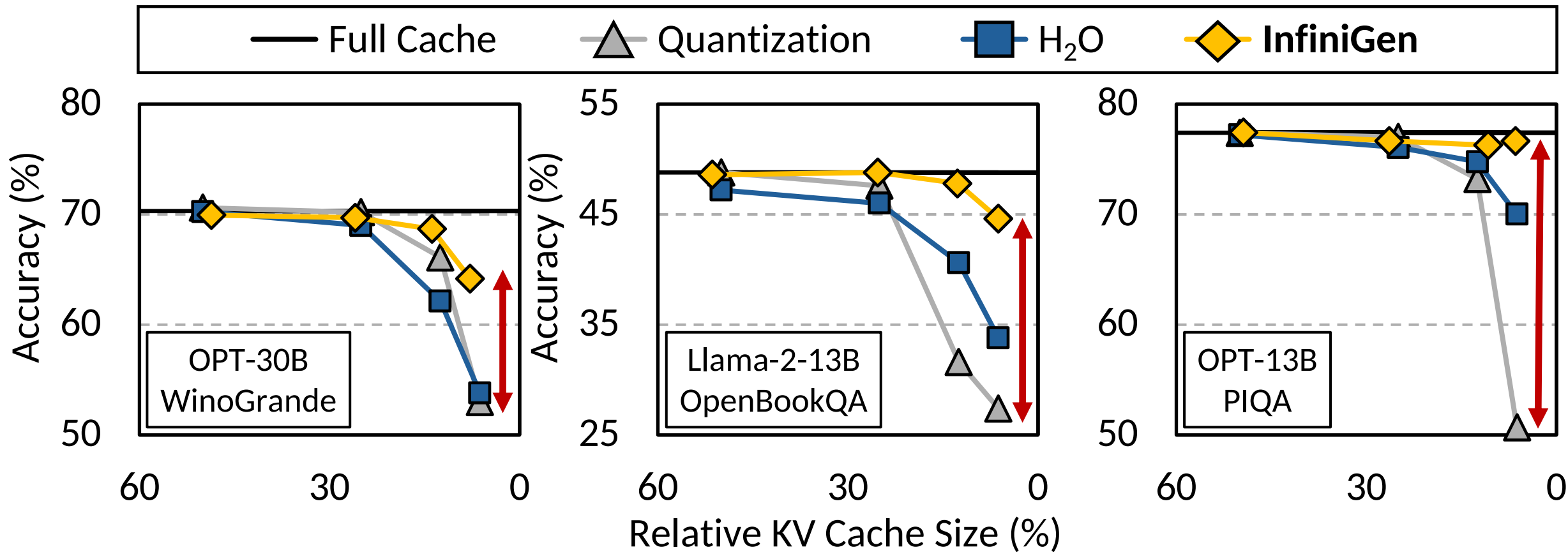| GPU | |
|---|---|
| GPU | NVIDIA RTX A6000 |
| GPU Memory Size | 48 GB |
| **CPU** | |
| CPU | Intel Xeon Gold 6136 |
| CPU Memory Size | 96GB |
| **Interconnect** | |
| PCIe Generation | 3.0 |
| Lane | 16 |

# Performance



**InfiniGen greatly improves the overall performance**
of a modern offloading-based inference system

# Performance



**InfiniGen improves performance with longer sequences**
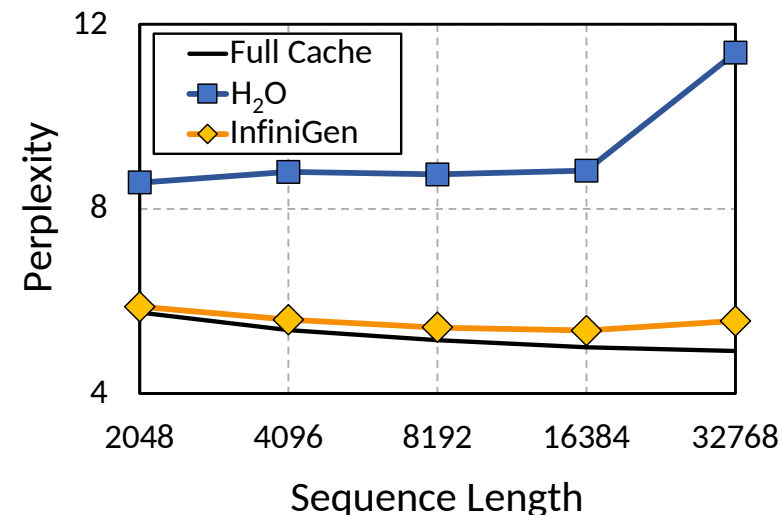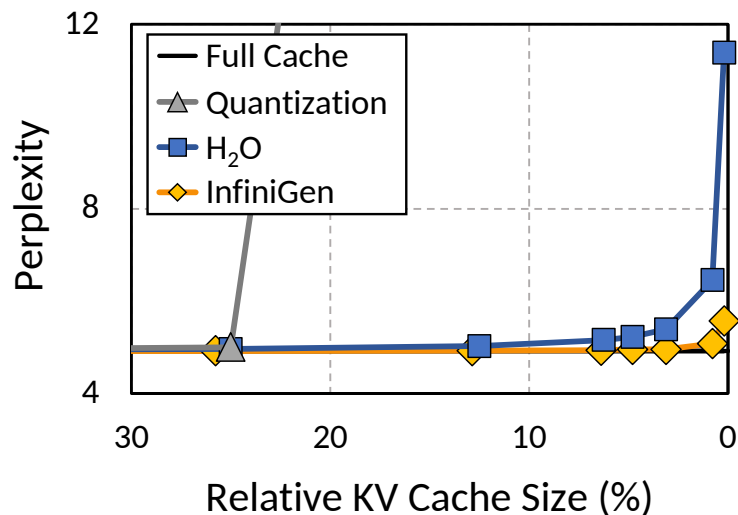while others lead to saturating speedups

# Accuracy



**InfiniGen offers substantially better model accuracy**
than other KV cache management methods

# More Details in Our Paper

- Long Sequences



- Key/Query Skewing

- Sensitivity Study and Overhead Analysis

- Latency across batch sizes and model sizes

- Accuracy with Other Combinations

- Others …

# Conclusion

**Problem**

- The large memory footprint of **the KV cache size** in LLM inference
- Existing methods show subpar performance

**Solution**: **InfiniGen**, a dynamic KV cache management framework

- **Speculative prefetching** of the essential KV cache
- **Skewing** query and key for efficient speculation

**Result**

- **InfiniGen** shows **3x faster** performance while **preserving model accuracy**
- It also shows **better scalability** than prior solutions! ☺

# Thank You!

## InfiniGen

Efficient Generative Inference
of Large Language Models
with Dynamic KV Cache Management

**Wonbeom Lee (wonbeom@snu.ac.kr)**



**Data Plane**　　**Control Plane**