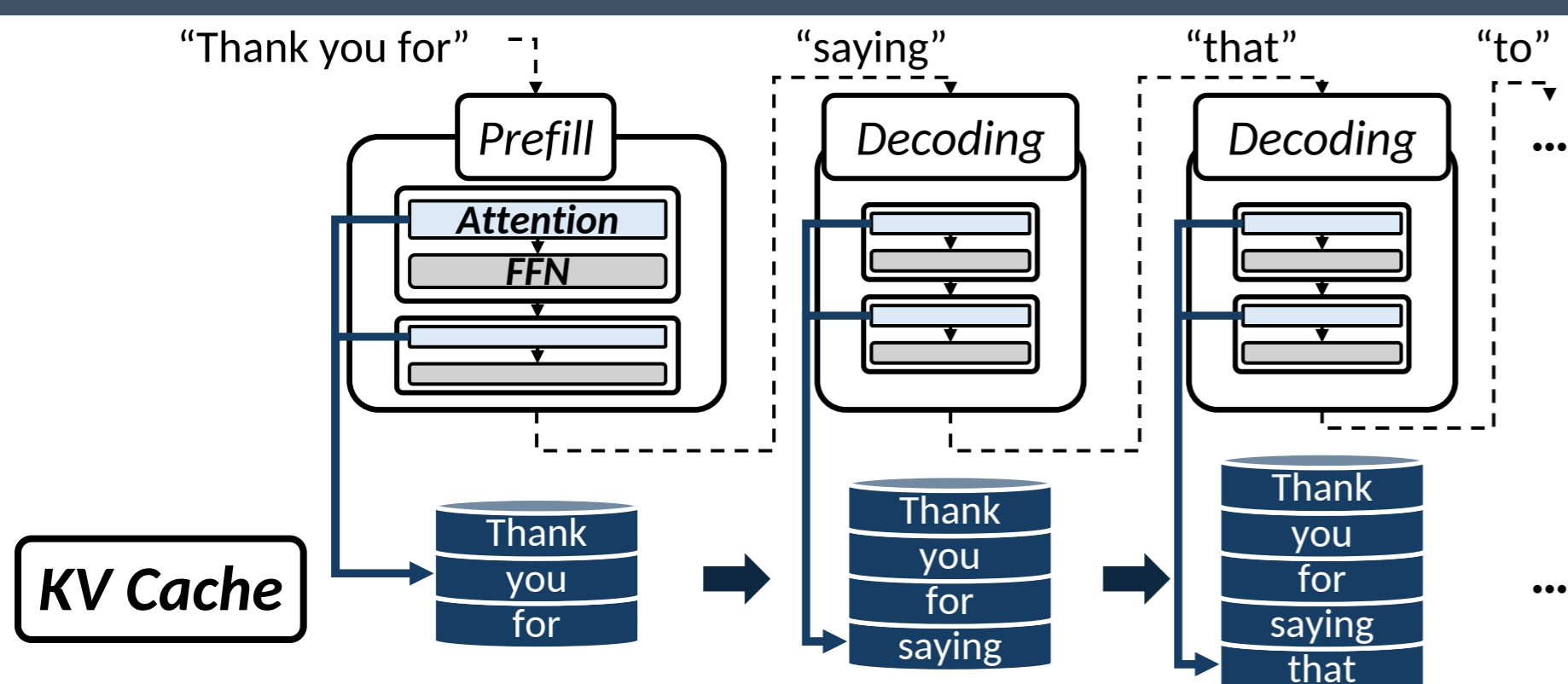# InfiniGen: Efficient Generative Inference of Large Language Models with Dynamic KV Cache Management

Wonbeom Lee*, Jungi Lee*, Junghwan Seo, Jaewoong Sim
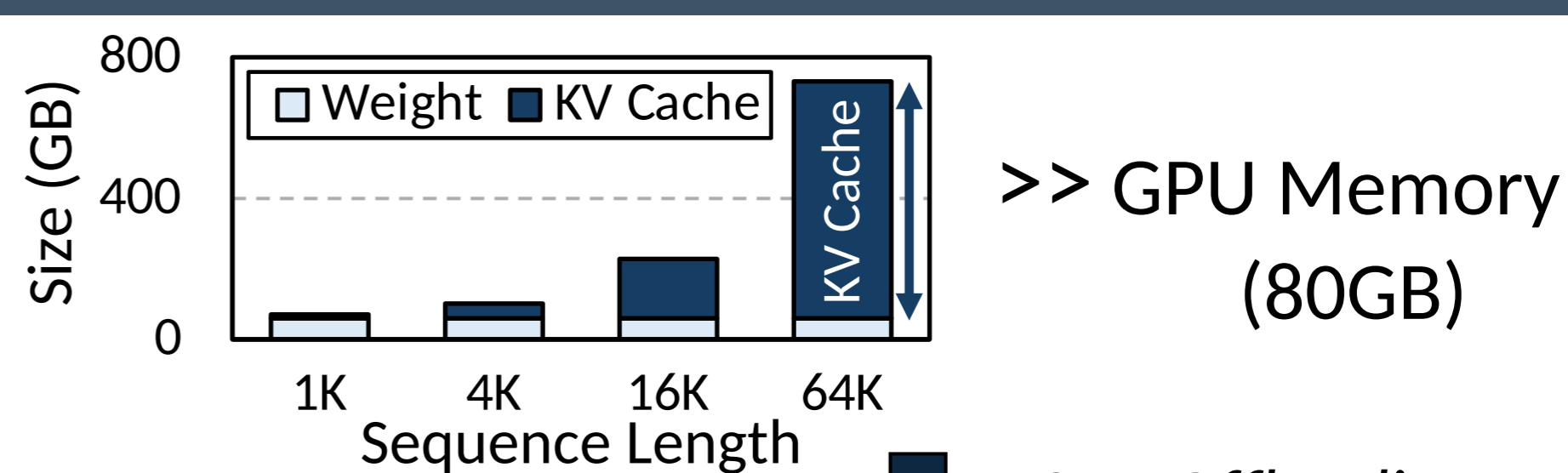
Seoul National University

## Motivation

### LLM Inference & KV Cache

"Thank you for" — *Prefill* — *Attention* / *FFN*
"saying" — *Decoding*
"that" — *Decoding*
"to" ...

**KV Cache**

Thank you for → Thank you for saying → Thank you for saying that ...

LLMs **memoize** the keys and values of the preceding tokens in memory to generate a new token that aligns well with the context.

### KV Cache Size & Transfer Overhead



Weight / KV Cache — Size (GB) vs Sequence Length (1K, 4K, 16K, 64K)

>> GPU Memory (80GB)

*CPU Offloading*

| Attention | FFN | **KV Cache Transfer** (CPU → GPU) | Attention | FFN | *time* |

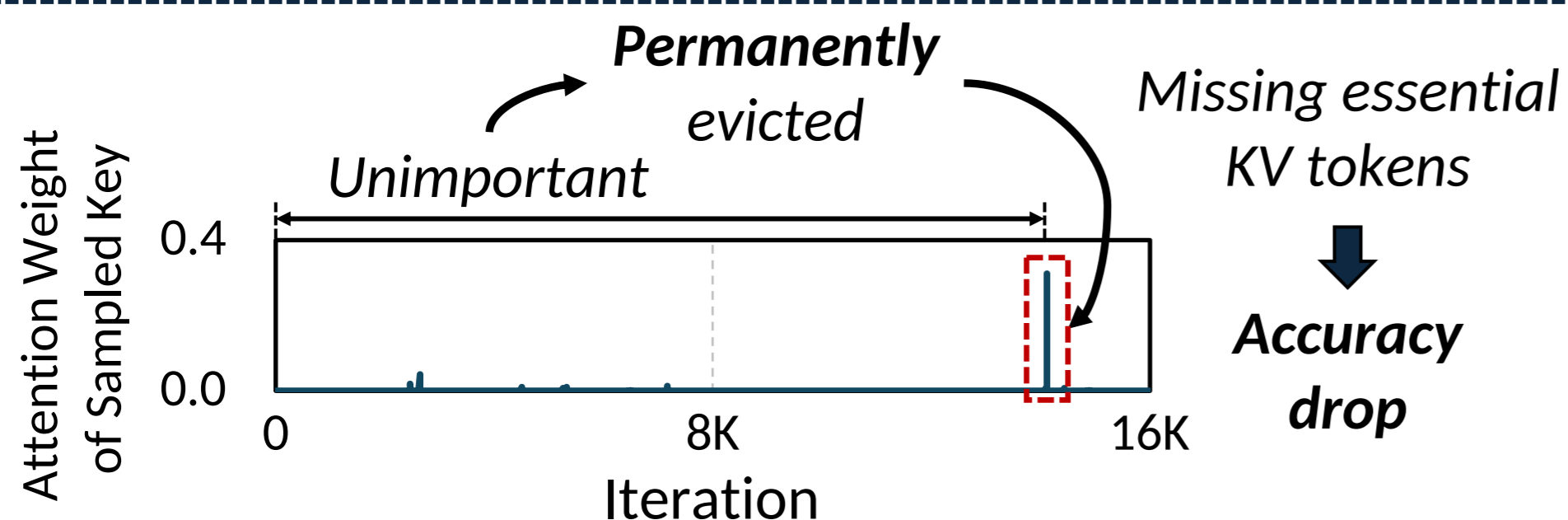Offloading can enable longer sequences exploiting large CPU memory, but lead to a **significant slowdown** due to the **limited PCIe bandwidth**.
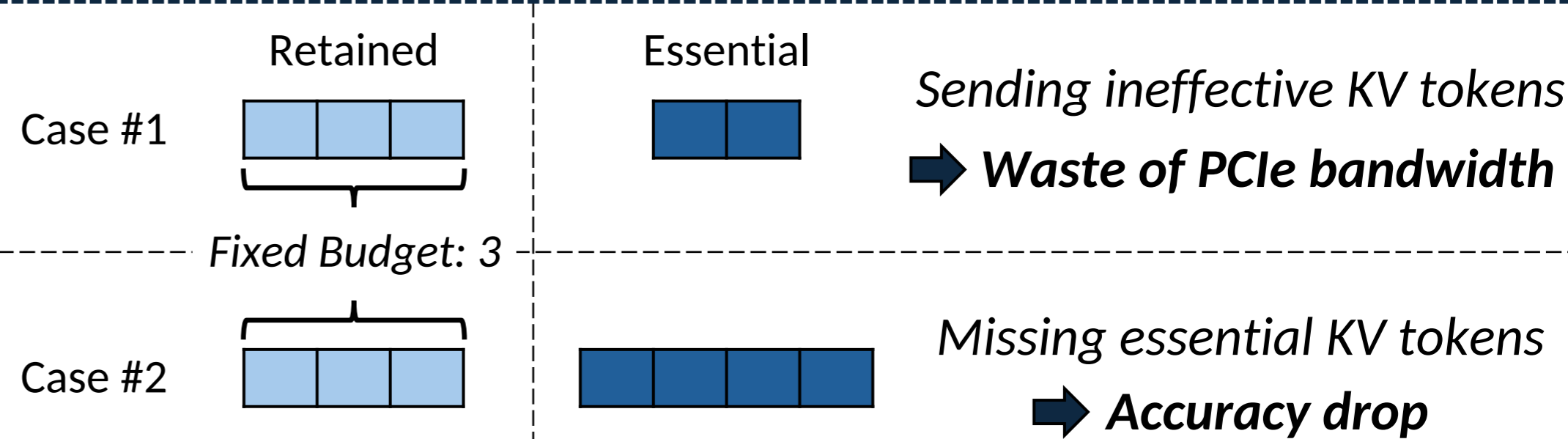
## Challenges

### Dynamic Nature of Attention Patterns across Iterations

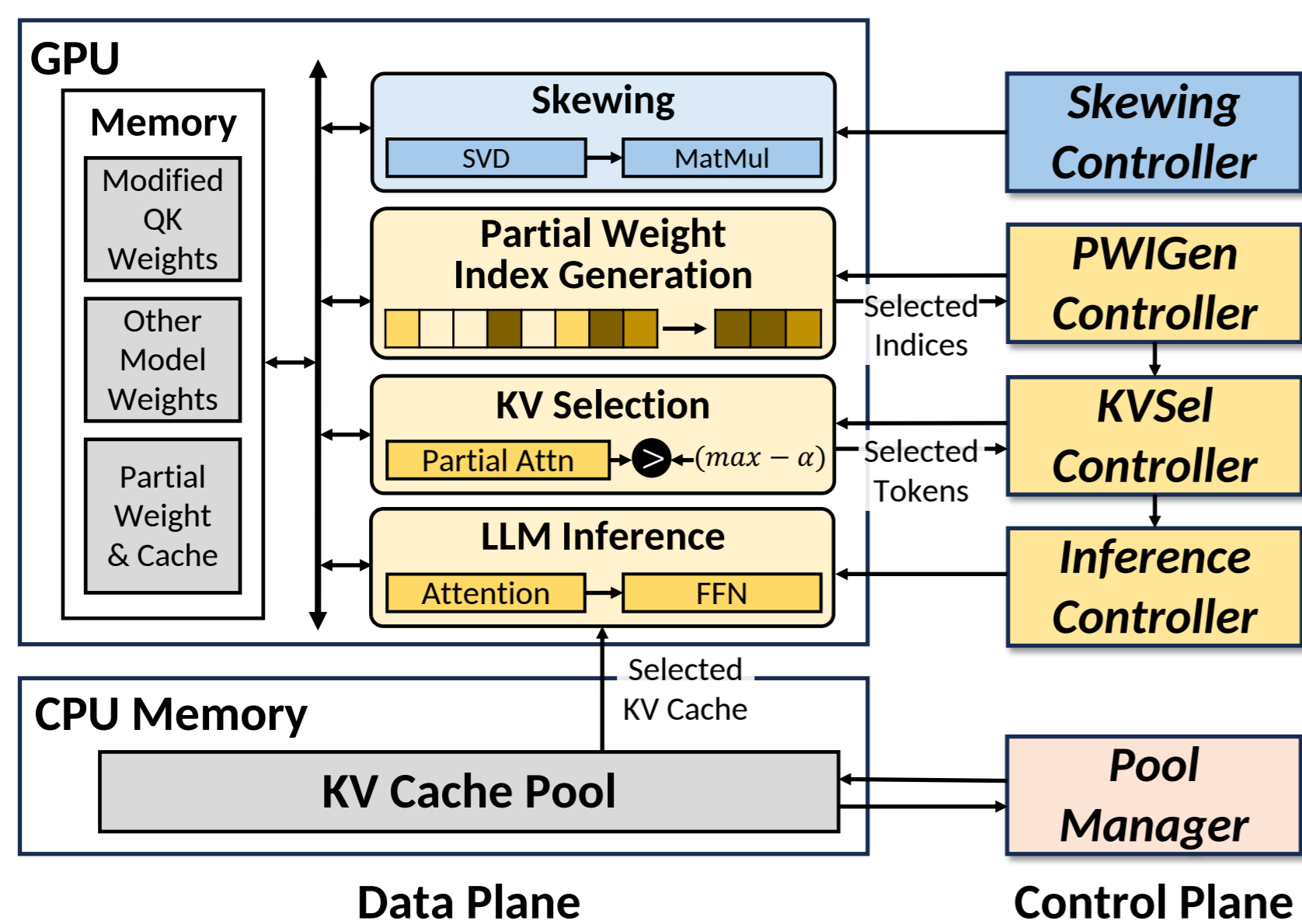Prior approaches **permanently evict** unimportant tokens.



**Permanently evicted**

*Unimportant*

*Missing essential KV tokens*

→ **Accuracy drop**

Attention Weight of Sampled Key (0.0, 0.4) vs Iteration (0, 8K, 16K)

### Adjusting the Number of KV across Layers and Queries

Prior approaches use a **fixed KV cache size budget**.

Retained — Essential

Case #1 — *Sending ineffective KV tokens* → **Waste of PCIe bandwidth**

*Fixed Budget: 3*

Case #2 — *Missing essential KV tokens* → **Accuracy drop**

## InfiniGen

### InfiniGen Design



**GPU** — Memory: Modified QK Weights, Other Model Weights, Partial Weight & Cache

Skewing: SVD / MatMul — *Skewing Controller*
Partial Weight Index Generation — *PWIGen Controller* — Selected Indices
KV Selection: Partial Attn ⊗ $(max - \alpha)$ — *KVSel Controller* — Selected Tokens
LLM Inference: Attention → FFN — *Inference Controller*

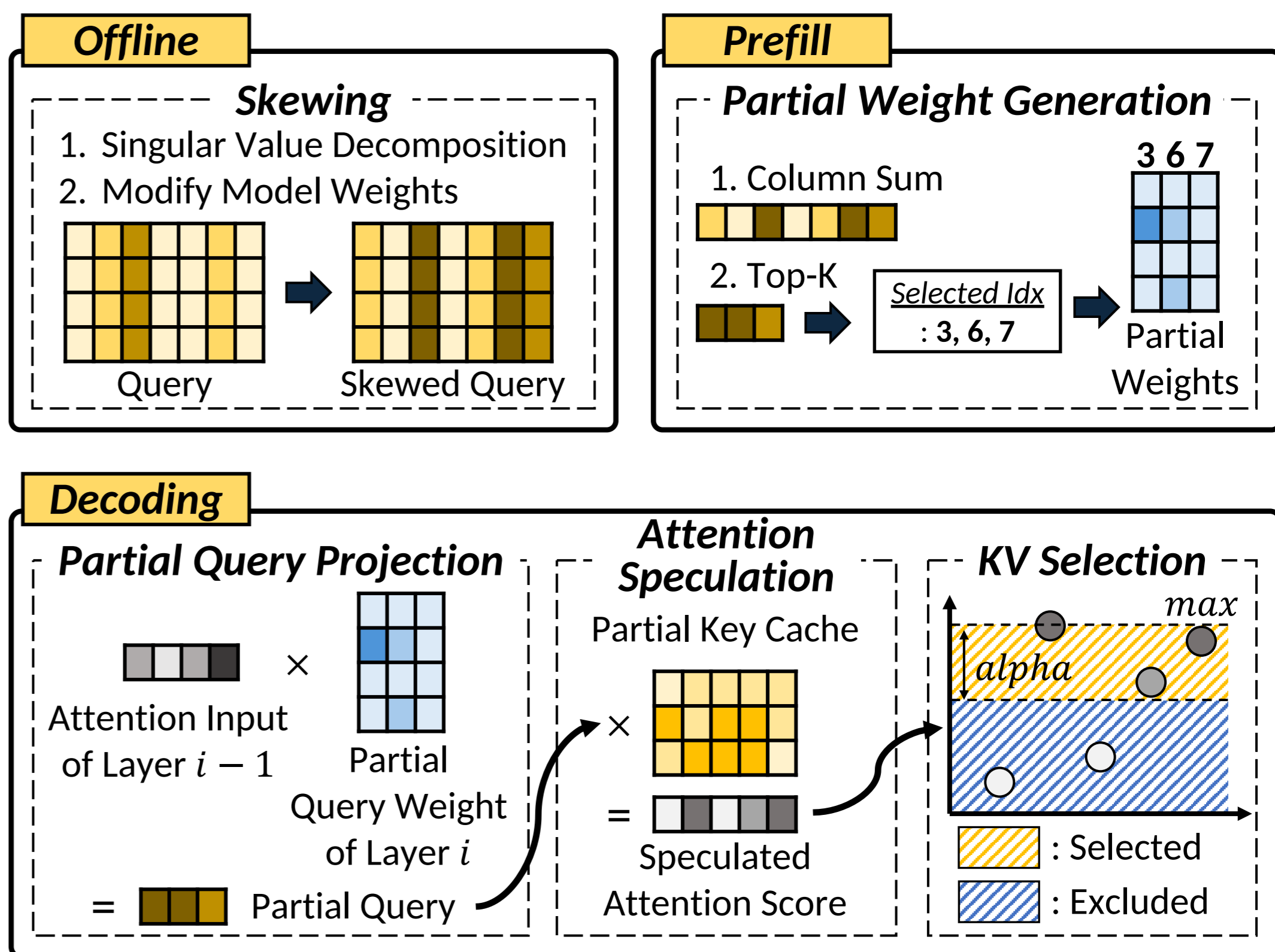**CPU Memory** — KV Cache Pool — Selected KV Cache — *Pool Manager*
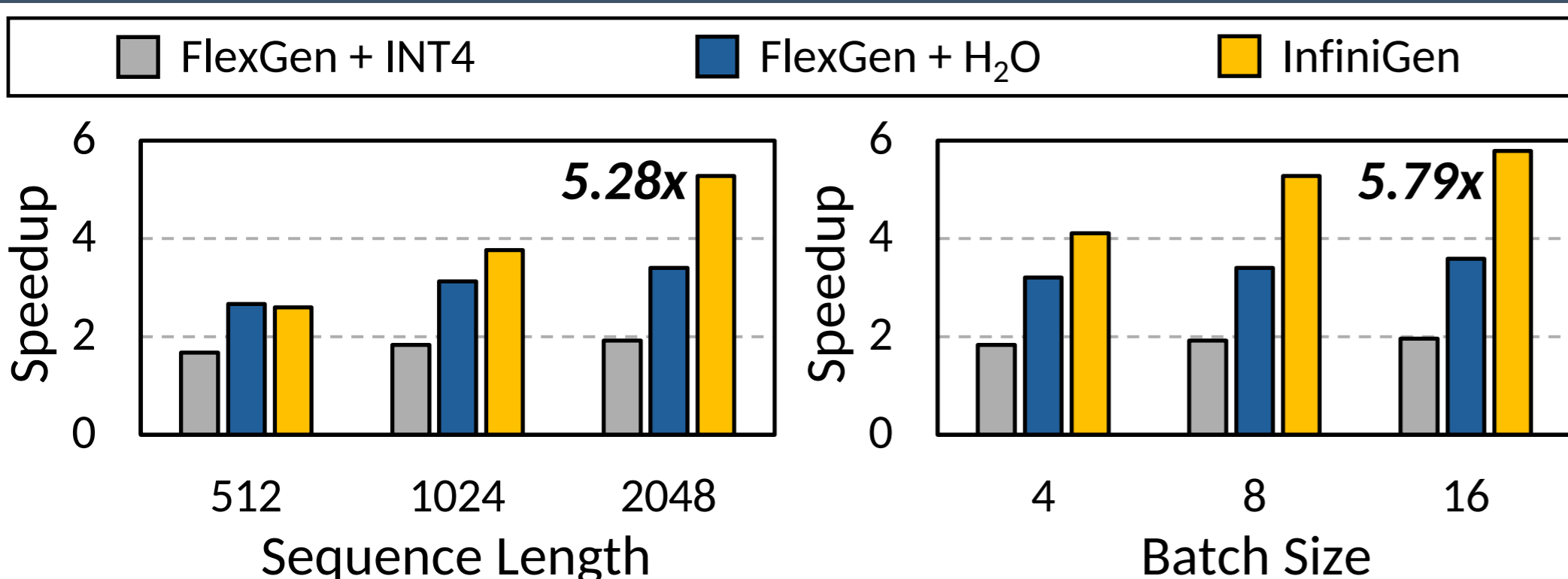
**Data Plane** — **Control Plane**

1. (Offline) Modify model weights in *Skewing Controller*
2. Speculate and prefetch KV tokens in *PWIGen/KVSel Controller*
3. Manage KV cache pool on CPU memory in *Pool Manager*
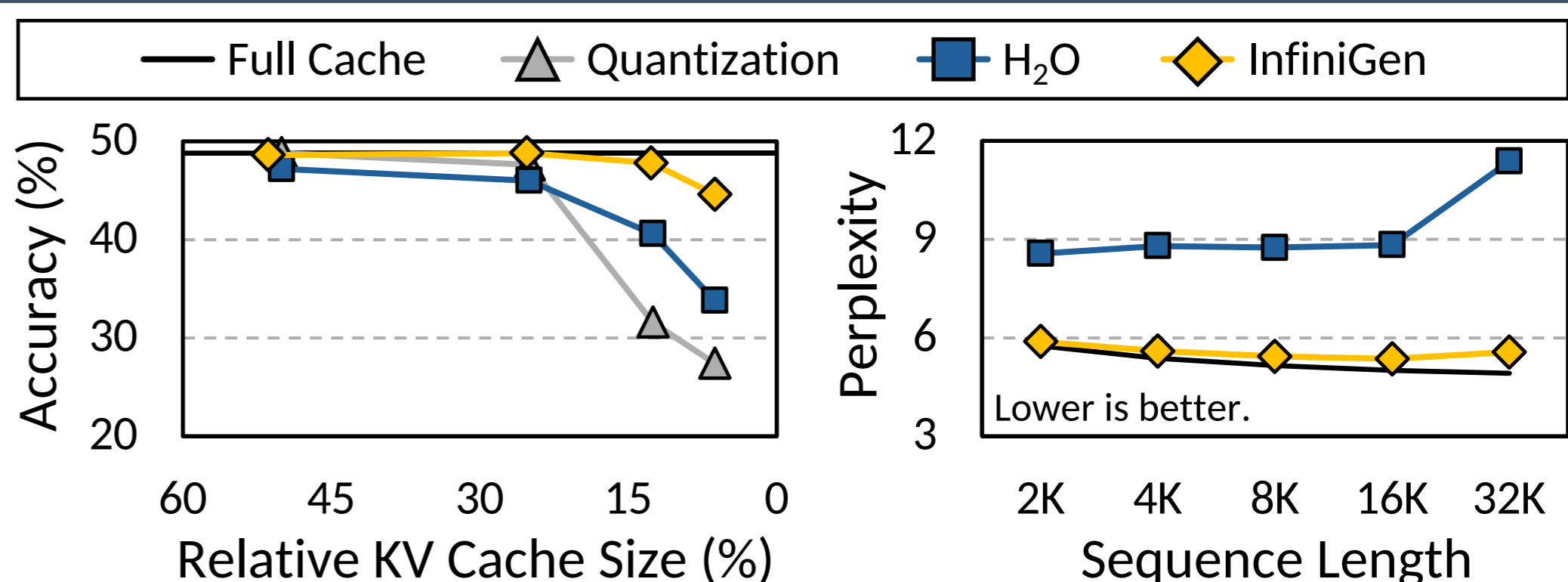
### Efficient KV Cache Prefetching

**Offline** — *Skewing*
1. Singular Value Decomposition
2. Modify Model Weights
Query → Skewed Query

**Prefill** — *Partial Weight Generation*
1. Column Sum — **3 6 7**
2. Top-K — *Selected Idx* : 3, 6, 7 — Partial Weights

**Decoding**

*Partial Query Projection*
Attention Input of Layer $i - 1$ × Partial Query Weight of Layer $i$ = Partial Query

*Attention Speculation*
Partial Key Cache × = Speculated Attention Score

*KV Selection*
$max$ / $alpha$ — : Selected / : Excluded

## Results

### Speedup over Modern Offloading-based Inference System



FlexGen + INT4 / FlexGen + $H_2O$ / InfiniGen

Speedup vs Sequence Length (512, 1024, 2048) — **5.28x**

Speedup vs Batch Size (4, 8, 16) — **5.79x**

→ *InfiniGen shows a shorter inference latency!*

### Accuracy



Full Cache / Quantization / $H_2O$ / InfiniGen

Accuracy (%) vs Relative KV Cache Size (%) (60, 45, 30, 15, 0)

Perplexity vs Sequence Length (2K, 4K, 8K, 16K, 32K) — Lower is better.

→ *InfiniGen consistently offers better accuracy!*