# Page Attention

Philipp Krähenbühl, UT Austin

# Training and Generation

## Vanilla Generation

| | Training | Training - Checkpointing | Generation |
|---|---|---|---|
| Peak Memory | $O(NL)$ | $O(NL^{\frac{1}{2}})$ | $O(N)$ |
| Runtime | $O(N^2L)$ | $O(2\,N^2L)$ | $O(N^3L)$ |
| # forward calls | 1 | 1 | N |

N

Output Text

Detokenizer

Layer

Layer

Layer

Layer

Layer

Layer

L

Tokenizer

Input Text

# Generation

- Step 1: Tokenize

- Step 2: $N \times$ Forward

- Step 3: Detokenize

Output Text

Detokenizer

Layer

Layer

Layer

Layer

Layer

Layer

Tokenizer

Input Text

# Vanilla Attention



Attention

$$q_{i-2} \qquad q_{i-1} \qquad q_i$$

$$\ldots \qquad \ldots \qquad k_{i-2}, v_{i-2} \qquad k_{i-1}, v_{i-1} \qquad k_i, v_i$$
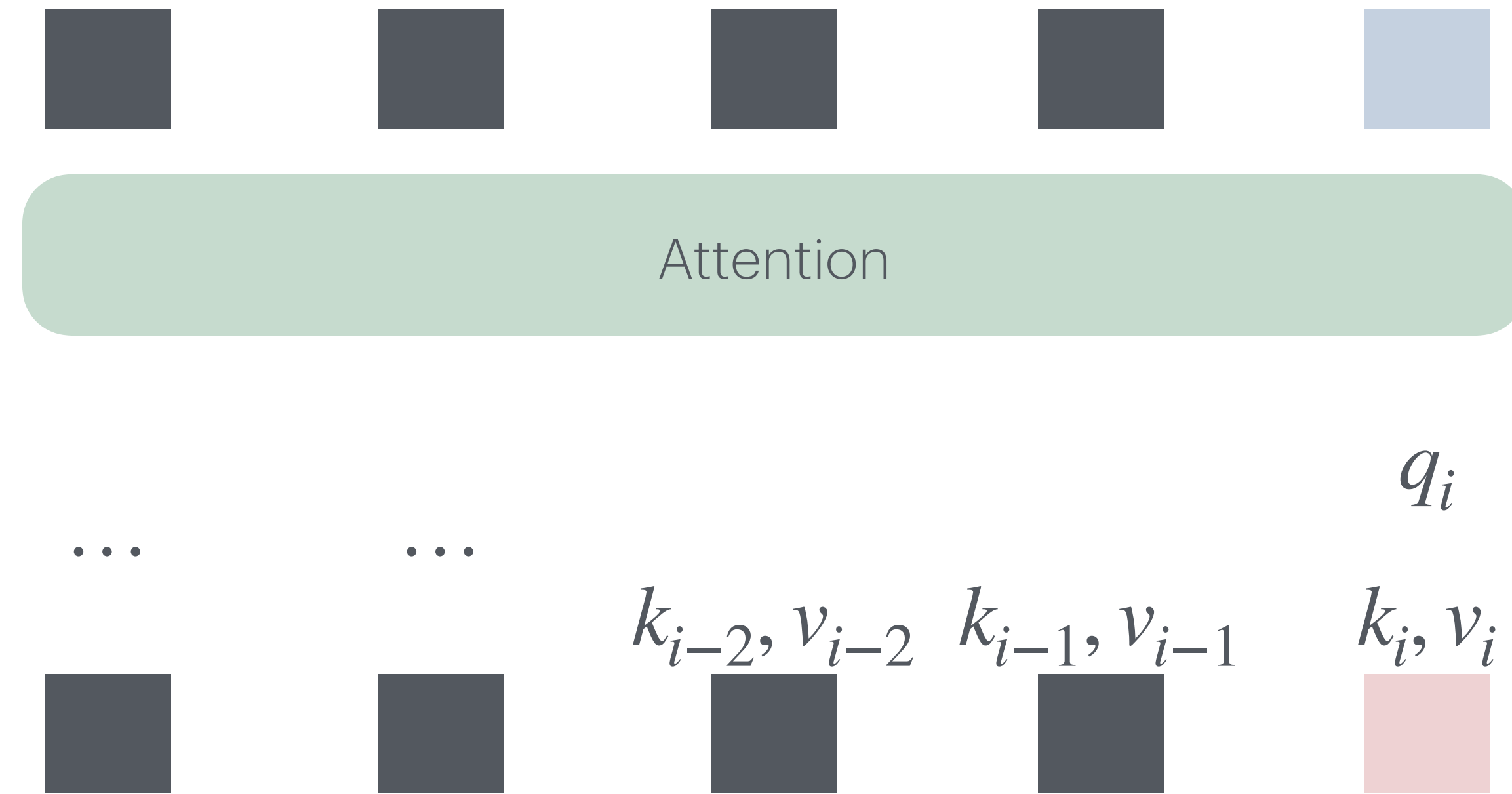
# Paged Attention

# Paged Attention - Analysis

- Cache keys and value

  - $O(NL)$ memory

- Less computation

  - $O(NL)$ per output token
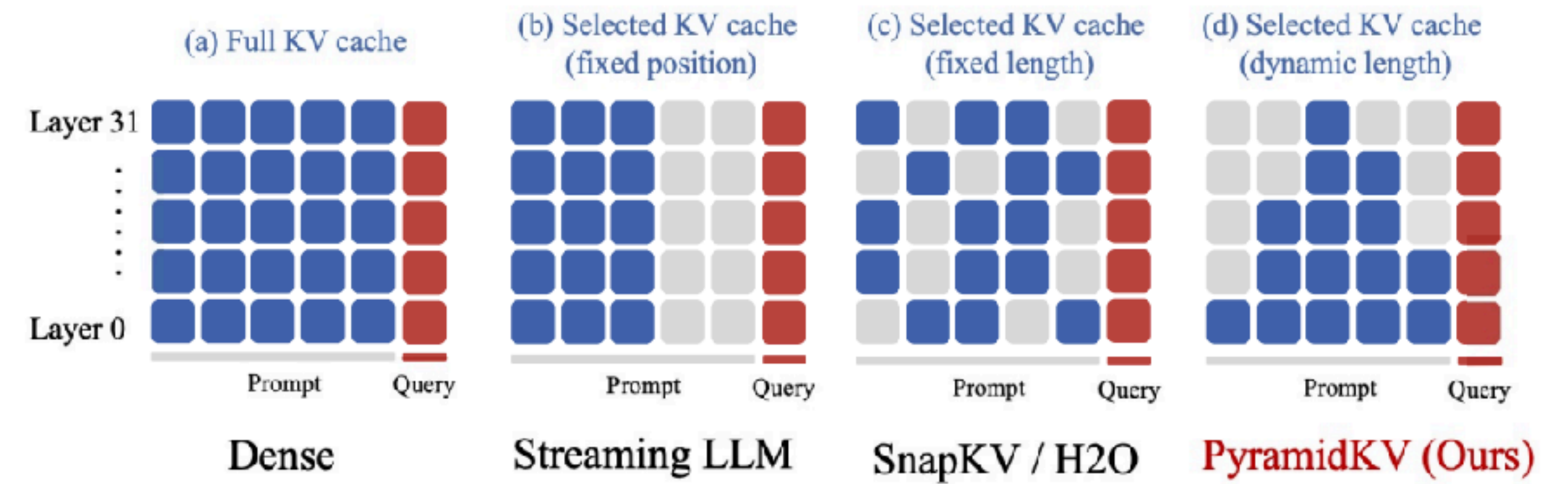
Attention

$$\dots \quad \dots \quad k_{i-2}, v_{i-2} \quad k_{i-1}, v_{i-1} \quad \begin{matrix} q_i \\ k_i, v_i \end{matrix}$$

Efficient Memory Management for Large Language Model Serving with PagedAttention, Kwon etal 2023

# Training and Generation
## Paged Attention

| | Training | Training - Checkpointing | Generation | Paged Attention |
|---|---|---|---|---|
| Peak Memory | $O(NL)$ | $O(NL^{1/2})$ | $O(N)$ | $O(NL)$ |
| Runtime | $O(N^2L)$ | $O(2\,N^2L)$ | $O(N^3L)$ | $O(N^2L)$ |
| # forward calls | 1 | 1 | N | N |



N

Output Text

Detokenizer

Layer

Layer

Layer

Layer

Layer

Layer

L

Tokenizer

Input Text

# Open Problem



(a) Full KV cache (b) Selected KV cache (fixed position) (c) Selected KV cache (fixed length) (d) Selected KV cache (dynamic length)

Dense — Streaming LLM — SnapKV / H2O — PyramidKV (Ours)

- A more efficient KV-Cache

  - Group Query Attention

  - Pruning

  - Low-rank representations?

- Connection to state-space models



Multi-head — Grouped-query — Multi-query

Values / Keys / Queries

PyramidKV: Dynamic KV Cache Compression based on Pyramidal Information Funneling, Cai etal 2024
GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints, Ainslie etal 2023

# References

- [1] Efficient Memory Management for Large Language Model Serving with PagedAttention, Kwon et al 2023. ([link](link))

- [2] PyramidKV: Dynamic KV Cache Compression based on Pyramidal Information Funneling, Cai et al 2024. ([link](link))

- [3] GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints, Ainslie et al 2023. ([link](link))