# Cache Me If You Must: Adaptive Key-Value Quantization for Large Language Models

Myeongjun Lee, Jeonghwan Ahn Department of Electrical Engineering POSTECH, Korea

### Contents

- 1. Preliminary study
- 2. Background and related work
- 3. Method
- 4. Experiments
- 5. Discussion



## 1. Preliminary study

- What is the KV Cache?
- Long Sequence
- Challenge in KV Cache compression

- ✓ Transformer-based models generate the next token (word) in an autoregressive manner based on previous inputs
- ✓ They predict the next word using the previously input text information.



[1] https://medium.com/@joaolages/kv-caching-explained-276520203249



- Repeated computations are performed during the auto-regressive generation process
  - recalculating the same previous token attention at each generation step
  - Computational cost !!!



Figure 1: The Transformer - model architecture.

[1] Ashish Vaswani et. al., "Attention Is All You Need", arxiv:1706.03762 [2] https://medium.com/@ioaolages/kv-caching-explained-276520203249 POSTECH

Zoom-in! (simplified without Scale and Softmax)

1. Preliminary study

### What is the KV-Cache?

- By caching the previous Keys an Values, it is possible to focus on only calculating the attention for the new token
- ✓ The matrices obtained with KV caching are small → fast matrix multiplications
- ✓ faster inference & to avoid wasteful re-computation



[1] https://medium.com/@joaolages/kv-caching-explained-276520203249

1. Preliminary study

### Long Sequence

- ✓ KV caching comes with its own pitfalls
- For longer context lengths, the KV cache becomes the dominant memory bottleneck



[Model size and activation memory size estimates for different LLaMA models]

BS	Model	Model Size (GB) $16 \rightarrow 2\text{-bit}$	K 32K	V Cache S 128K	Size w/ Diff 1M	f. Seq Len (GB) 10M (16 $\rightarrow$ 2-bit)
1	7B 13B 30B 65B	$12.6 \to 1.6 \\ 24.1 \to 3.0 \\ 60.3 \to 7.5 \\ 121.1 \to 15.1$	16 25 49 80	64 100 195 320	512 800 1560 2560	$\begin{array}{c} 4883 \to 610 \\ 7629 \to 954 \\ 14877 \to 1860 \\ 24414 \to 3052 \end{array}$
4	7B 13B 30B 65B	$\begin{array}{c} 12.6 \rightarrow 1.6 \\ 24.1 \rightarrow 3.0 \\ 60.3 \rightarrow 7.5 \\ 121.1 \rightarrow 15.1 \end{array}$	64 100 195 320	256 400 780 1280	2048 3200 6240 10240	$\begin{array}{c} 19531 \rightarrow 2441 \\ 30518 \rightarrow 3815 \\ 59509 \rightarrow 7439 \\ 97656 \rightarrow 12207 \end{array}$

[1] Coleman Hopper et. al., "KVQuant: Towards 10 Million Context Length LLM Inference with KV Cache Quantization", NuerIPS 2024

## 2. Background and related work

- Challenge in KV Cache compression
- Approach to KV Cache quantization



### **Challenges in KV-Cache Compression**

#### ✓ Both compression and decompression speeds

Dynamically adding new entries to the caches as well as decoding them at inference

### ✓ Inherent structure in caches

- The existence of attention sinks
- Large outlier values

- ✓ Quantization Granularity
- ✓ Error Handling
- ✓ Cross Layer merging



2. Background

#### Quantization Granularity

- An approach to determining at what scope or level the scale factor should be set for quantization
- Examples include Global, Group-wise, Channel-wise, Token-wise, hybrid granularity



Figure 2: Magnitude of key and value cache for Llama-2-13B and Falcon-7B. We observe (1) for key cache, there are a few channels whose magnitudes are very large. (2) for value cache, there is no obvious outlier pattern.

[1] Zirui Liu et. al., "KIVI A Tuning-Free Asymmetric 2bit Quantization for KV Cache", arxiv:2402.02750

#### POSTECH

### Error Handling

- Approaches to minimizing and correcting quantization error
- Most methods maintain a window of recent historical KV cache in full precision to preserve accuracy



Figure 1: The distributions of activations at the input to the FFN block in LLAMA2-7B model, in the tenth layer. Left: using the default configuration as downloaded from Hugging Face. Right: after processing using QuaRot. The processed distribution has no outliers, leading to superior quantization.

[1] Saleh Ashkboos et. al., "QuaRot: Outlier-Free 4Bit Inference in Rotated LLMs", NuerIPS 2024

#### POSTECH

2. Background

#### Cross Layer merging

 Combining the KV caches from multiple layers into a single shared cache to reduce memory usage



(b) KV Cache Sharing Strategy Searching



(c) Comparison between MiniCache and previous methods

[1] Yifei Yang et. al., "KVSHARER: EFFICIENT INFERENCE VIA LAYER-WISE DISSIMILAR KV CACHE SHARING", arxiv:2410.18517 [2] Akide Liu et. al., "MiniCache: KV Cache Compression in Depth Dimension for Large Language Models", arxiv:2405.14366

POSTECH

- Core research idea
- Inter-Layer dependencies
- AQUA-KV Calibration
- AQUA-KV Inference

### **Core research idea**

- Remove predictable information from the KV cache and cache only the residuals after quantization
- Predict the current layer's KV cache based on the Key and Value of the previous layer





*Figure 3.* An intuitive scheme of the AQUA-KV inference. Only the quantized residuals are saved for each block.

3. Method

#### Why Analyze Inter-Layer Dependencies?

- To build a predictor, the Key/Value of the current layer must be predictable
- AQUA-KV assumes a structure where the Key/Value is predicted linearly across layers
- To validate this hypothesis, the inter-layer dependency is experimentally verified across diverse inputs



#### ✓ Linear probe model

- Train linear regression model
- Prediction and interdependency analysis of i-th layer keys or values across various inputs

### ✓ Various inputs

Previous layer keys and values, adjacent tokens, and different vector types

#### ✓ Analysis of Layer Dependencies

- Inter Layer : cached vectors at different layer or tokens
- Intra Layer : (key & value) vectors within the same layer
- Extracting with the simple predictor models

#### Explained variance ratio

- A metric that quantifies how well the model explains the variance of the actual values.
- A higher value indicates better information reconstruction from the source input.
- Ensures fair comparison across layers with different scales (e.g., size and distribution of Keys and Values).

$$ext{EVR} = 1 - rac{ ext{Var}(y - \hat{y})}{ ext{Var}(y)}$$

- If a predictor captures 90% of the variance, it means that the subsequent quantization only needs to capture the remaining 10% of variance
- This would mean that the resulting quantization will also have roughly 10 times smaller error
- 1bit and 2bit quantizer usually explain 0.75 and 0.89 variance respectively

#### POSTECH

#### ✓ What Do the Results Tell Us?

- Keys and Values can be predicted based on inter-layer dependency
- For attention keys, using just one previous already achieves errors similar to 2bit quantization
- For values, the dependency on previous layer is also strong
- Strong dependencies between key and values within the same layer



*Figure 2.* Mean Explained Variance Ratios by linear probes from previous blocks (L), tokens (T) and role on Llama-3.2-3B.



### **AQUA-KV**

Exploits inter- and intra-layer dependencies to improve quantization accuracy

### Predictors

- Per layer
  - Key predictor, Value Predictor •
- Sequential Prediction
  - •
- Reflect the way predictors are used
  - For training: reconstructed cache entries as input •
- **Key Predictor** 
  - previous layer keys
- Value Predictor
  - previous layer keys + current layer keys
  - Cannot predict in both direction within same layer
    - $V_{L,rec} \rightarrow K_L, K_{L,rec} \rightarrow V_L$



Algorithm 1 AQUA-KV Calibration
<b>Require:</b> model, data, quantization method $Q$
1: $\mathbf{X} \leftarrow \texttt{model.input_embeddings}$ (data)
2: $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \emptyset$
3: predictors $\leftarrow \{\}$
4: for $i=1,\ldots,$ model.num_layers do
5: block $\leftarrow$ model.transformer_layers[i]
6: $(\mathbf{K}, \mathbf{V}) \leftarrow \texttt{block.get_attention_kv}(\mathbf{X})$
7: $\mathbf{X} \leftarrow \texttt{block}(\mathbf{X})$
8: <b>if</b> $\mathbf{K}_{old} = \emptyset$ and $\mathbf{V}_{old} = \emptyset$ then
9: $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}, \mathbf{V}$
10: continue
11: end if
12: # Predict keys from past keys
13: $f_{\text{kev}} \leftarrow \arg\min_f \ f(\mathbf{K}_{\text{old}}) - \mathbf{K}\ _2^2$
14: $\mathbf{K}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{K} - f_{\text{kev}}(\mathbf{K}_{\text{old}})))$
15: $\mathbf{K}_{\text{rec}} \leftarrow f_{\text{key}}(\mathbf{K}_{\text{old}}) + \mathbf{K}_{\text{residue}}$
16: # Predict values from past values and current keys
17: $f_{\text{value}} \leftarrow \arg\min_{f} \ f([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) - \mathbf{V}\ _2^2$
18: $\mathbf{V}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{V} - f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}])))$
19: $\mathbf{V}_{\text{rec}} \leftarrow f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) + \mathbf{V}_{\text{residue}}$
20: predictors[i] = $(f_{key}, f_{value})$
21: $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}_{\text{rec}}, \mathbf{V}_{\text{rec}}$
22: end for
23: return predictors



### ✓ AQUA-KV

• Exploits inter- and intra-layer dependencies to improve quantization accuracy

### ✓ Predictors

- Per layer
  - Key predictor, Value Predictor
- Sequential Prediction
  - Operate layer-by-layer
- Reflect the way predictors are used
  - For training: reconstructed cache entries as input
- Key Predictor
  - previous layer keys
- Value Predictor
  - previous layer keys + current layer keys
  - Cannot predict in both direction within same layer

• 
$$V_{L,rec} \rightarrow K_L$$
,  $K_{L,rec} \rightarrow V_L$ 

Values are harder to predict

Alge	orithm 1 AQUA-KV Calibration
Req	<b>uire:</b> model, data, quantization method $Q$
1:	$\mathbf{X} \leftarrow \texttt{model.input}_\texttt{embeddings}(\texttt{data})$
2:	$\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \emptyset$
3:	$predictors \leftarrow \{\}$
4:	for $i=1,\ldots,$ model.num_layers do
5:	$\texttt{block} \gets \texttt{model.transformer\_layers}[i]$
6:	$(\mathbf{K}, \mathbf{V}) \leftarrow \texttt{block.get_attention_kv}(\mathbf{X})$
7:	$\mathbf{X} \leftarrow \texttt{block}(\mathbf{X})$
8:	if $\mathbf{K}_{\mathrm{old}} = \emptyset$ and $\mathbf{V}_{\mathrm{old}} = \emptyset$ then
9:	$\mathbf{K}_{\mathrm{old}}, \mathbf{V}_{\mathrm{old}} \gets \mathbf{K}, \mathbf{V}$
10:	continue
11:	end if
12:	# Predict keys from past keys
13:	$f_{\text{kev}} \leftarrow \arg\min_{f} \ f(\mathbf{K}_{\text{old}}) - \mathbf{K}\ _2^2$
14:	$\mathbf{K}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{K} - f_{\text{key}}(\mathbf{K}_{\text{old}})))$
15:	$\mathbf{K}_{\text{rec}} \leftarrow f_{\text{key}}(\mathbf{K}_{\text{old}}) + \mathbf{K}_{\text{residue}}$
16:	# Predict values from past values and current keys
17:	$f_{\text{value}} \leftarrow \arg\min_{f} \ f([\mathbf{V}_{\text{old}}   \mathbf{K}_{\text{rec}}]) - \mathbf{V}\ _{2}^{2}$
18:	$\mathbf{V}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{V} - f_{\text{value}}([\mathbf{V}_{\text{old}}   \mathbf{K}_{\text{rec}}])))$
19:	$\mathbf{V}_{\text{rec}} \leftarrow f_{\text{value}}([\mathbf{V}_{\text{old}};\mathbf{K}_{\text{rec}}]) + \mathbf{V}_{\text{residue}}$
20:	predictors[i] = $(f_{\text{kev}}, f_{\text{value}})$
21:	$\mathbf{K}_{\mathrm{old}}, \mathbf{V}_{\mathrm{old}} \leftarrow \mathbf{K}_{\mathrm{rec}}, \mathbf{V}_{\mathrm{rec}}$
22:	end for
23:	return predictors



## **AQUA-KV** Calibration

3. Method

#### ✓ AQUA-KV



- Key Predictor
  - previous layer keys
- Value Predictor

POSTECH

- previous layer keys + current layer keys
- Cannot predict in both direction within same layer

• 
$$V_{L,rec} \rightarrow K_L$$
,  $K_{L,rec} \rightarrow V_L$ 

• Values are harder to predict

Algorithm 1 AQUA-KV Calibration **Require:** model, data, quantization method Q 1: X ← model.input\_embeddings(data) 2:  $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \emptyset$ 3: predictors  $\leftarrow$  {} 4: for  $i = 1, \ldots, model.num_layers do$  $block \leftarrow model.transformer_layers[i]$ 5:  $(\mathbf{K}, \mathbf{V}) \leftarrow block.get_attention_kv(\mathbf{X})$ 6:  $\mathbf{X} \leftarrow \texttt{block}(\mathbf{X})$ 7: if  $\mathbf{K}_{old} = \emptyset$  and  $\mathbf{V}_{old} = \emptyset$  then 8: 9:  $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}, \mathbf{V}$ continue 10: end if 11: # Predict keys from past keys 12:  $f_{\text{key}} \leftarrow \arg\min_{f} \|f(\mathbf{K}_{\text{old}}) - \mathbf{K}\|_2^2$ 13:  $\mathbf{K}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{K} - f_{\text{kev}}(\mathbf{K}_{\text{old}})))$ 14:  $\mathbf{K}_{\text{rec}} \leftarrow f_{\text{kev}}(\mathbf{K}_{\text{old}}) + \mathbf{K}_{\text{residue}}$ 15: # Predict values from past values and current keys 16:  $f_{\text{value}} \leftarrow \arg\min_{f} \|f([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) - \mathbf{V}\|_{2}^{2}$ 17: $\mathbf{V}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{V} - f_{\text{value}}(|\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}|)))$ 18:  $\mathbf{V}_{\text{rec}} \leftarrow f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) + \mathbf{V}_{\text{residue}}$ 19: predictors[i] =  $(f_{\text{key}}, f_{\text{value}})$ 20:  $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}_{\text{rec}}, \mathbf{V}_{\text{rec}}$ 21: 22: end for 23: return predictors

### ✓ AQUA-KV

POSTECH

Exploits inter- and intra-layer dependencies to improve quantization accuracy

### Lightweight Linear Regressors

- Fast training with closed-form solutions
  - $W^* = (X^T X)^{-1} X^T y$
  - ~ 4 hours training for Llama-3.1-70B
- Low compute & memory overhead
  - Modern LLMs: Grouped Query Attention (GQA)



Improves accuracy on recent tokens

[1] Ainslie, Joshua, et al. "Gqa: Training generalized multi-query transformer models from multi-head checkpoints." arXiv preprint arXiv:2305.13245 (2023)

- 1: X ← model.input\_embeddings(data) 2:  $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \emptyset$ 3: predictors  $\leftarrow$  {} 4: for  $i = 1, \ldots, model.num_layers do$  $block \leftarrow model.transformer_layers[i]$  $(\mathbf{K}, \mathbf{V}) \leftarrow block.get_attention_kv(\mathbf{X})$ 6:  $\mathbf{X} \leftarrow \texttt{block}(\mathbf{X})$ 7: if  $\mathbf{K}_{old} = \emptyset$  and  $\mathbf{V}_{old} = \emptyset$  then 8:  $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}, \mathbf{V}$ 9: continue 10: 11: end if # Predict keys from past keys 12:  $f_{\text{key}} \leftarrow \arg\min_{f} \|f(\mathbf{K}_{\text{old}}) - \mathbf{K}\|_2^2$ 13:  $\mathbf{K}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{K} - f_{\text{kev}}(\mathbf{K}_{\text{old}})))$ 14: 15:  $\mathbf{K}_{\text{rec}} \leftarrow f_{\text{kev}}(\mathbf{K}_{\text{old}}) + \mathbf{K}_{\text{residue}}$ # Predict values from past values and current keys 16:  $f_{\text{value}} \leftarrow \arg\min_{f} \|f([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) - \mathbf{V}\|_{2}^{2}$ 17: $\mathbf{V}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{V} - f_{\text{value}}(|\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}|)))$ 18:  $\mathbf{V}_{\text{rec}} \leftarrow f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) + \mathbf{V}_{\text{residue}}$ 19: predictors [i] =  $(f_{\text{kev}}, f_{\text{value}})$ 20:
- 21:  $\mathbf{K}_{old}, \mathbf{V}_{old} \leftarrow \mathbf{K}_{rec}, \mathbf{V}_{rec}$

Algorithm 1 AQUA-KV Calibration

**Require:** model, data, quantization method Q

- 22: **end for**
- 23: **return** predictors

### ✓ AQUA-KV

- Exploits inter- and intra-layer dependencies to improve quantization accuracy
- ✓ Lightweight Linear Regressors
  - Fast training with closed-form solutions
    - $W^* = (X^T X)^{-1} X^T y$
    - ~ 4 hours training for Llama-3.1-70B
  - Low compute & memory overhead
    - Modern LLMs: Grouped Query Attention (GQA)

### Quantization

- Only the residual parts that cannot be predicted
- Attention Sink
  - Keep the first 4 tokens uncompressed
- Recent token buffer
  - Up to 128 tokens uncompressed
  - Improves accuracy on recent tokens

Algorithm 1 AQUA-KV Calibration
<b>Require:</b> model, data, quantization method $Q$
1: $\mathbf{X} \leftarrow \texttt{model.input}_\texttt{embeddings}$ (data)
2: $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \emptyset$
3: predictors $\leftarrow$ {}
4: for $i=1,\ldots,$ model.num_layers do
5: block $\leftarrow$ model.transformer_layers[i]
6: $(\mathbf{K},\mathbf{V}) \leftarrow  extsf{block.get_attention_kv}(\mathbf{X})$
7: $\mathbf{X} \leftarrow  ext{block}(\mathbf{X})$
8: <b>if</b> $\mathbf{K}_{old} = \emptyset$ <b>and</b> $\mathbf{V}_{old} = \emptyset$ <b>then</b>
9: $\mathbf{K}_{\text{old}}, \mathbf{V}_{\text{old}} \leftarrow \mathbf{K}, \mathbf{V}$
10: continue
11: end if
12: # Predict keys from past keys
13: $f_{\text{key}} \leftarrow \arg\min_{f} \ f(\mathbf{K}_{\text{old}}) - \mathbf{K}\ _{2}^{2}$
14: $\mathbf{K}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{K} - f_{\text{key}}(\mathbf{K}_{\text{old}})))$
15: $\mathbf{K}_{\text{rec}} \leftarrow f_{\text{key}}(\mathbf{K}_{\text{old}}) + \mathbf{K}_{\text{residue}}$
16: # Predict values from past values and current keys
17: $f_{\text{value}} \leftarrow \arg\min_{f} \ f([\mathbf{V}_{\text{sld}}; \mathbf{K}_{\text{rec}}]) - \mathbf{V}\ _2^2$
18: $\mathbf{V}_{\text{residue}} \leftarrow Q^{-1}(Q(\mathbf{V} - f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}])))$
19: $\mathbf{V}_{\text{rec}} \leftarrow f_{\text{value}}([\mathbf{V}_{\text{old}}; \mathbf{K}_{\text{rec}}]) + \mathbf{V}_{\text{residue}}$
20: predictors[i] = $(f_{\text{key}}, f_{\text{value}})$
21: $\mathbf{K}_{old}, \mathbf{V}_{old} \leftarrow \mathbf{K}_{rec}, \mathbf{V}_{rec}$
22: end for
23: return predictors

### ✓ Inference

- Decode → Computation → Encode
  - Reconstruct KV-cache from a compressed representation

Algorithm 4 Inference with AQUA-KV Predictors Require: model, input, key\_cache, value\_cache, predictors 1:  $\mathbf{K}_{\text{past}}, \mathbf{V}_{\text{past}}, \mathbf{K}_{\text{inp-prev}}, \mathbf{V}_{\text{inp-prev}} \leftarrow \emptyset$ 2: X ← model.input\_embeddings(input) 3: for  $i = 0, \ldots, model.num_layers - 1$  do # Recover previously saved key-values 4:  $(\hat{\mathbf{K}}_{\text{past}}, \hat{\mathbf{V}}_{\text{past}}) \leftarrow \text{decode}(i, \text{key_cache}[i], \text{value_cache}[i], \hat{\mathbf{K}}_{\text{past}}, \hat{\mathbf{V}}_{\text{past}}, \text{predictors})$ 5: 6: # Run forward pass  $block \leftarrow model.transformer_layers[i]$ 7:  $(\mathbf{K}_{inp}, \mathbf{V}_{inp}) \leftarrow block.get_attention_kv(\mathbf{X})$ 8:  $\mathbf{X} \leftarrow block(\mathbf{X}, K = [\hat{\mathbf{K}}_{past}; \mathbf{K}_{inp}], V = [\hat{\mathbf{V}}_{past}, \mathbf{V}_{inp}])$  {Concatenate past and input keys/values} 9: # Compress new key-value entries 10:  $(\mathbf{K}_{\text{inp}}^{q}, \mathbf{V}_{\text{inp}}^{q}) \leftarrow \texttt{encode}(i, \hat{\mathbf{K}}_{\text{inp}}, \mathbf{V}_{\text{inp}}, \hat{\mathbf{K}}_{\text{inp-prev}}, \mathbf{V}_{\text{inp-prev}}, \texttt{predictors})$ 11:  $key_cache[i] \leftarrow [key_cache[i]; \mathbf{K}_{inp}^q] \{Concatenate key cache with new keys\}$ 12:  $value_cache[i] \leftarrow [value_cache[i]; \mathbf{V}_{inp}^q] \{Concatenate value cache with new values\}$ 13:  $(\hat{\mathbf{K}}_{\text{inp-prev}}, \hat{\mathbf{V}}_{\text{inp-prev}}) \leftarrow \texttt{decode}(i, \mathbf{K}_{\text{inp}}^q, \mathbf{V}_{\text{inp}}^q, \hat{\mathbf{K}}_{\text{inp-prev}}, \hat{\mathbf{V}}_{\text{inp-prev}}, \texttt{predictors})$ 14: 15: end for 16: return model.compute\_logits(X)

### ✓ Inference

- Decode → Computation → Encode
  - Inference combines the current token's KV with the reconstructed cache from decode step



### ✓ Inference

- Decode → Computation → Encode
  - · Only residual part that the predictor cannot predict is compressed and stored

Algorithm 4 Inference with AQUA-KV Predictors Require: model, input, key\_cache, value\_cache, predictors 1:  $\hat{\mathbf{K}}_{\text{past}}, \hat{\mathbf{V}}_{\text{past}}, \hat{\mathbf{K}}_{\text{inp-prev}}, \hat{\mathbf{V}}_{\text{inp-prev}} \leftarrow \emptyset$ 2: X ← model.input\_embeddings(input) 3: for  $i = 0, \ldots, model.num_layers - 1$  do # Recover previously saved key-values 4:  $(\hat{\mathbf{K}}_{\text{past}}, \hat{\mathbf{V}}_{\text{past}}) \leftarrow \text{decode}(i, \text{key_cache}[i], \text{value_cache}[i], \hat{\mathbf{K}}_{\text{past}}, \hat{\mathbf{V}}_{\text{past}}, \text{predictors})$ 5: # Run forward pass 6:  $block \leftarrow model.transformer_layers[i]$ 7:  $(\mathbf{K}_{inp}, \mathbf{V}_{inp}) \leftarrow block.get_attention_kv(\mathbf{X})$ 8:  $\mathbf{X} \leftarrow block(\mathbf{X}, K = [\hat{\mathbf{K}}_{past}; \mathbf{K}_{inp}], V = [\hat{\mathbf{V}}_{past}, \mathbf{V}_{inp}])$  {Concatenate past and input keys/values} 9: # Compress new key-value entries 10:  $(\mathbf{K}_{inp}^{q}, \mathbf{V}_{inp}^{q}) \leftarrow \texttt{encode}(i, \hat{\mathbf{K}}_{inp}, \mathbf{V}_{inp}, \hat{\mathbf{K}}_{inp\_prev}, \mathbf{V}_{inp\_prev}, \texttt{predictors})$ 11:  $key\_cache[i] \leftarrow [key\_cache[i]; \mathbf{K}_{inp}^{q}] \{Concatenate key cache with new keys\}$ 12:  $value_cache[i] \leftarrow [value_cache[i]; \mathbf{V}_{inp}^q] \{Concatenate value cache with new values\}$ 13:  $(\hat{\mathbf{K}}_{\text{inp-prev}}, \hat{\mathbf{V}}_{\text{inp-prev}}) \leftarrow \texttt{decode}(i, \mathbf{K}_{\text{inp}}^q, \mathbf{V}_{\text{inp}}^q, \hat{\mathbf{K}}_{\text{inp-prev}}, \hat{\mathbf{V}}_{\text{inp-prev}}, \texttt{predictors})$ 14: 15: end for

16: return model.compute\_logits $(\mathbf{X})$ 

## 4. Evaluation

- Evaluation Settings
- Detailed Evaluation
- Large-Scale Evaluation



### Predictor Calibration

- Dataset: RedPajama
  - Random 256 sequences with 8192 tokens
  - Use 32 of sequences as holdout for hyperparameter selection
  - Use remaining 224 for training the predictors

### Perplexity evaluation

Account for the effect of recent token buffers and attention sinks

### ✓ LongBench v1

- 14 English-language tasks
  - SamSum, 2WikiMQ, TREC, HotpotAQ, MultiNews, TriviaQA, QMSum, PsgCount, MFQA\_en, Musique, Qasper, PsgRetr, NarrativeQA, GovReport
- Without restricting the input 8192 tokens
  - To better explore the effectiveness of AQUA-KV on longer sequences

### Alternative Quantizer

- Addition of AQUA-KV to Quanto
  - Group size: 64, per token quantization
  - A significant improvement in both PPL and LongBench avg scores
- Addition of AQUA-KV to HIGGS
  - $d = 2, n \in \{16, 64, 256\}, 2-, 3-, 4$ -bit for each
  - Outperforms in both PPL and LongBench avg scores

### ✓ Layer Sharing

Sharing multiple layer pairs causes major accuracy drops

Config	Quant. Bits	Wiki2 PPL↓ (base model)	LongBench Avg.↑ (instruct model)
Uncompressed	16	6.98	44.47
Quanto-2b-gs64	2.50	21.56	33.59
KIVI-2b-gs128-r128	2.25	9.33	39.63
KVQuant-2b-s1%	2.33	9.43	20.56
HIGGS-2b-gs1024	2.02	7.47	43.25
KVSharer (1 pair)	15.43	7.45	36.81
KVSharer (4 pairs)	13.71	9.60	29.82
AQUA-KV (Quanto)	2.64	10.33	43.64
AQUA-KV (HIGGS)	2.16	7.03	44.26

### **Detailed Evaluation**

#### Predictor Architecture

- Linear regression (main proposal)
- Reduce-Rank Regression(with rank=256)
  - Slight increase in PPL
- MLP with two layers
  - Double hidden dimension + Layer normalization
  - Slightly better performance Image: Slightly better performance
  - A significant cost in size and speed
- Linear regressor + GPTQ
  - Effective for scenarios where minimal model size is crucial

Config	Quant. Bits	Wiki2 PPL↓ (base model)	LongBench Avg.↑ (instruct model)			
AQUA-K	V Predic	tor Architectu	re			
Linear (162 MiB) MLP (540 MiB) RRR (68 MiB) GPTQ (41 MiB)	2.16 2.16 2.16 2.16 2.16	7.03 7.03 7.22 7.03	44.26 44.61 44.30 44.19			



### **Detailed Evaluation**

#### ✓ First Layer Keys & Values

- Use 4-bit HIGGS quantization
  - Since the first layer is not compressed by AQUA-KV predictors

Config	Quant. Bits	Wiki2 PPL↓ (base model)	LongBench Avg.↑ (instruct model)							
AQUA-KV 1st Layer Keys & Values										
Keep in BF16	2.59	7.03	44.26							
HIGGS 4 bit	2.16	7.03	44.26							
HIGGS 3 bit	2.13	7.03	44.28							
HIGGS 2 bit	2.09	7.05	44.41							

### ✓ Rotary Positional Embedding (RoPE)

- Natural dilemma about whether to apply predictors and quantizers
  - Better in Pre-RoPE
  - Simple linear models are not rotation-equivariant

Config	Quant. Bits	Wiki2 PPL↓ (base model)	LongBench Avg.↑ (instruct model)
Before RoPE HIGGS	2.16	7.03	44.26
After RoPE HIGGS	2.16	7.05	44.13
Before RoPE Quanto gs64	2.64	10.90	44.54
After RoPE Quanto gs64	2.64	10.33	43.64



#### ✓ Ablation Analysis

Config	Quant.	Wiki2 PPL↓	LongBench Avg.↑							
	Bits	(base model)	(instruct model)							
Ablation Analysis										
AQUA-KV (default)	2.16	7.03	44.26							
w/o 16-bit Attn. Sink	2.16	7.15	44.13							
w/o V predictor	2.09	7.06	43.91							
w/o K predictor	2.09	7.50	42.92							
w/o pre-RoPE	2.16	7.05	44.13							

#### ✓ Evaluation of AQUA-KV with HIGGS backbone

Particularly noticeable for extreme 2-bit compression

	Quant	Cache size	WikiText-2 PPL↓ (base model)				LongBench Average↑ (instruct)					
Method	Juant.	CiR 70R		Llama 3.x	ĸ	Qwe	n 2.5		Llama 3.x		Qwe	n 2.5
	Dits	GIB, 70B	3B	8B	70B	3B	7B	3B	8B	70B	3B	7B
Uncompressed	16	40	6.98	5.61	2.54	7.14	6.13	44.61	48.13	52.92	38.80	46.82
AQUA-KV	2.09	5.7	7.03	5.72	2.62	7.20	6.17	44.30	47.77	52.79	38.31	46.43
HIGGS	2.02	5.1	7.47	5.89	2.77	7.93	8.08	42.80	47.37	52.18	30.92	25.97
KIVI	2.25	5.6	9.34	7.37	3.06	9.05	7.02	39.64	46.28	52.45	28.66	32.78
KVQuant	2.33	5.6	9.43	6.64	3.28	_		20.56	37.17	46.14	_	—
AQUA-KV	3.06	8.1	6.98	5.64	2.55	7.15	6.14	44.37	48.10	52.81	38.77	46.81
HIGGS	3.02	7.6	7.05	5.66	2.57	7.26	7.20	44.41	47.86	52.56	31.85	14.61
KIVI	3.05	7.7	7.87	6.04	2.87	7.63	6.37	41.40	46.98	52.87	30.37	32.63
KVQuant	3.33	8.3	7.26	5.84	2.75	_	—	41.40	46.42	50.74	_	—
AQUA-KV	4.02	10.9	6.98	5.61	2.54	7.15	6.14	44.48	48.10	52.95	38.92	46.77
HIGGS	4.02	10.6	7.01	5.62	2.55	7.16	6.88	44.41	48.07	52.97	32.11	11.54
KIVI	4.25	10.6	7.03	5.64	2.61	7.17	6.14	43.11	47.57	52.88	31.50	33.40
KVQuant	4.33	10.8	7.04	5.65	2.58	_	—	43.62	47.77	52.89	—	—

#### ✓ Evaluation of AQUA-KV with HIGGS backbone

Particularly noticeable for extreme 2-bit compression

	Quant	Cache size	1	WikiText-2	2 PPL↓ (b	ase mode	el)		LongBenc	h Average↑	(instruct)	
Method	bite	CiR 70R		Llama 3.3	ĸ	Qwe	en 2.5		Llama 3.x		Qwe	en 2.5
	Dits	GIB, 70B	3B	8B	70B	3B	7B	3B	8B	70B	3B	7B
Uncompressed	16	40	6.98	5.61	2.54	7.14	6.13	44.61	48.13	52.92	38.80	46.82
AQUA-KV	2.09	5.7	7.03	5.72	2.62	7.20	6.17	44.30	47.77	52.79	38.31	46.43
HIGGS	2.02	5.1	7.47	5.89	2.77	7.93	8.08	42.80	47.37	52.18	30.92	25.97
KIVI		KV for 2 hit	compr	ossion	is roug	hly og	uivalon	t to the	3_hit HI	CCS ha	solino a	uantizor
KVQuant			compre	531011	is roug	my eq	uivaicii		5-01111	665 <i>ba</i> .	senne q	uantizei
AQUA-KV	3.06	8.1	6.98	5.64	2.55	7.15	6.14	44.37	48.10	52.81	38.77	46.81
HIGGS	3.02	7.6	7.05	5.66	2.57	7.26	7.20	44.41	47.86	52.56	31.85	14.61
KIVI	3.05	7.7	7.87	6.04	2.87	7.63	6.37	41.40	46.98	52.87	30.37	32.63
KVQuant	3.33	8.3	7.26	5.84	2.75	_	—	41.40	46.42	50.74	—	—
AQUA-KV	4.02	10.9	6.98	5.61	2.54	7.15	6.14	44.48	48.10	52.95	38.92	46.77
HIGGS	4.02	10.6	7.01	5.62	2.55	7.16	6.88	44.41	48.07	52.97	32.11	11.54
KIVI	4.25	10.6	7.03	5.64	2.61	7.17	6.14	43.11	47.57	52.88	31.50	33.40
KVQuant	4.33	10.8	7.04	5.65	2.58	—	_	43.62	47.77	52.89	—	_

# 5. Discussion

- Summary
- Limitations & further research
- Our Observed Limitations



## Conclusion

### ✓ Summary

- Analysis of KV cache structure
  - Identified compression opportunities via mutual information
- AQUA-KV compression framework
  - Exploit inter- and intra-layer dependencies
  - Compatible with any quantization

### Limitations & further research

- Inference speed
  - Depends heavily on the backbone quantizer
  - Cannot be faster than the baseline
  - For Llama-70B

POSTPEH

- HIGSS: 5.91 tokens/s
- AQUA-KV: 5.76 tokens/s

**Require:** layer\_index, K, V, reconstructed keys and values of the previous layer  $\hat{\mathbf{K}}_{\text{prev}}, \hat{\mathbf{V}}_{\text{prev}}$ , predictors 1: if layer\_index = 0 then 2: return K. V 3: end if 4: f<sub>key</sub>, f<sub>value</sub> ← predictors[layer\_index] 5:  $\mathbf{K}_{\text{pred}} \leftarrow f_{\text{kev}}(\hat{\mathbf{K}}_{\text{prev}})$  $\leftarrow O(\mathbf{K} - \mathbf{K})$  $\hat{\mathbf{K}} \leftarrow O^{-1}(\mathbf{K}_a) + \mathbf{K}_{\text{pred}}$ red / Jvalue ([V prev, R]) {Concatenate and  $\hat{\mathbf{K}}$ 9:  $\mathbf{V}_a \leftarrow Q(\mathbf{V} - \mathbf{V}_{\text{pred}})$ 10: return  $\mathbf{K}_{q}, \mathbf{V}_{q}$ **Redundant computations** Algorithm 3 decode Require: layer\_index,  $\mathbf{K}_q$ ,  $\mathbf{V}_q$ , reconstructed keys and values of the previous layer  $\hat{\mathbf{K}}_{\text{prev}}$ ,  $\hat{\mathbf{V}}_{\text{prev}}$ , predictors 1: if layer\_index = 0 then 2: return  $\mathbf{K}_{q}, \mathbf{V}_{q}$ 3: end if  $f_{key}(\hat{\mathbf{K}}_q) + f_{key}(\hat{\mathbf{K}}_{prev})$ and predictor expects both previous reconstructed  $\hat{\mathbf{V}}_{prev}$  and current reconstructed  $\hat{\mathbf{K}}$ . 7:  $\hat{\mathbf{V}} \leftarrow Q^{-1}(\mathbf{V}_q) + f_{\text{value}}([\hat{\mathbf{V}}_{\text{prev}}; \hat{\mathbf{K}}])$ 8: return K. Ŷ

Algorithm 2 encode

Unoptimized pytorch can be improved with specialized libraries.

5. Discussion

### Conclusion

### Our Observed Limitations

Predictor overhead	Method	Quant. bits	Cache size GiB, 70B
<ul> <li>Additional memory and compute costs</li> </ul>	Uncompressed	16	40
<ul> <li>No evaluation is provided for non-GQA architectures</li> </ul>	AQUA-KV	2.09	5.7
	HIGGS	2.02	5.1
<ul> <li>Although claimed to be minimal due to GQA</li> </ul>	KIVI	2.25	5.6
	KVQuant	2.33	5.6

- Suggesting the justification may be incidental rather than fundamental.
- Limited Parallelism
  - Current key usage for value predictor restricts parallel execution
  - Unclear whether the performance gain justifies the added cost
- The Lack of Evalutaion
  - There is no EVR results for  $[V_{L-1}; K_L]$
  - only  $[KV_{L-1}; K_L]$  and  $[KV_{L-1}]$  are evaluated
  - What if using  $[KV_{L-1}]$  yields comparable performance with only minor inference-time overhead?

# Thank you



### ✓ Encode

Only residual part that the predictor cannot predict is compressed and stored

```
      Algorithm 2 encode

      Require: layer_index, K, V, reconstructed keys and values of the previous layer \hat{K}_{prev}, \hat{V}_{prev}, predictors

      1: if layer_index = 0 then

      2: return K, V

      3: end if

      4: f_{key}, f_{value} \leftarrow predictors[layer_index]

      5: K_{pred} \leftarrow f_{key}(\hat{K}_{prev})

      6: K_q \leftarrow Q(K - K_{pred})

      7: \hat{K} \leftarrow Q^{-1}(K_q) + K_{pred}

      8: V_{pred} \leftarrow f_{value}([\hat{V}_{prev}; \hat{K}]) {Concatenate \hat{V}_{prev} and \hat{K}}

      9: V_q \leftarrow Q(V - V_{pred})

      10: return K_q, V_q
```

#### ✓ Decode

Reconstruct KV-cache from a compressed representation

```
Algorithm 3 decode

Require: layer_index, \mathbf{K}_q, \mathbf{V}_q, reconstructed keys and values of the previous layer \hat{\mathbf{K}}_{\text{prev}}, \hat{\mathbf{V}}_{\text{prev}}, predictors

1: if layer_index = 0 then

2: return \mathbf{K}_q, \mathbf{V}_q

3: end if

4: f_{\text{key}}, f_{\text{value}} \leftarrow \text{predictors[layer_index]}

5: \hat{\mathbf{K}} \leftarrow Q^{-1}(\mathbf{K}_q) + f_{\text{key}}(\hat{\mathbf{K}}_{\text{prev}})

6: # Value predictor expects both previous reconstructed \hat{\mathbf{V}}_{\text{prev}} and current reconstructed \hat{\mathbf{K}}.

7: \hat{\mathbf{V}} \leftarrow Q^{-1}(\mathbf{V}_q) + f_{\text{value}}([\hat{\mathbf{V}}_{\text{prev}}; \hat{\mathbf{K}}])

8: return \hat{\mathbf{K}}, \hat{\mathbf{V}}
```

Appendix

### **Auto-regressive PPL evaluation**

```
Appendix.
```

```
if sequence_index < num_sequences_without_padding:
    shift_logits = lm_logits[:, :-1, :].contiguous()
    shift_labels = input_ids[:, 1:]
    loss_fct = nn.CrossEntropyLoss()
    loss = loss_fct(shift_logits.view(-1, shift_logits.size(-1)), shift_labels.view(-1))
    total_nll += loss.float() * shift_labels.numel()
    total_tokens += shift_labels.numel()
else:
    raise RuntimeError
```

```
ppl = torch.exp(total_nll / total_tokens)
return ppl.item()
```