



**ICLR**

International Conference On  
Learning Representations

# **KV Cache Transform Coding for Compact Storage in LLM Inference**

Konrad Staniszewski & Adrian Łańcucki

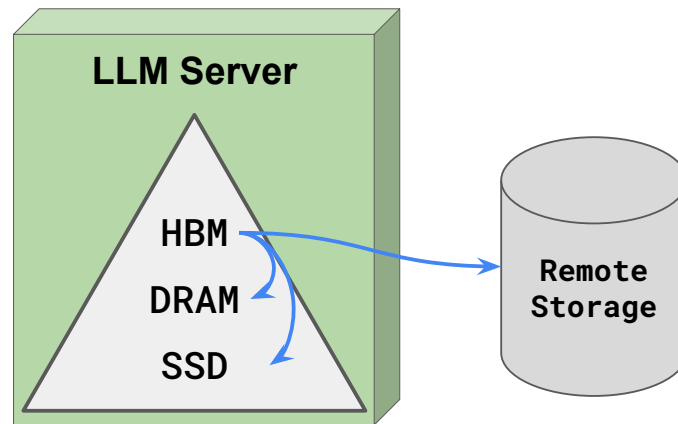
Rio de Janeiro, Brazil, 2026

# Problem Setting

## Compression of KV Cache for Storage

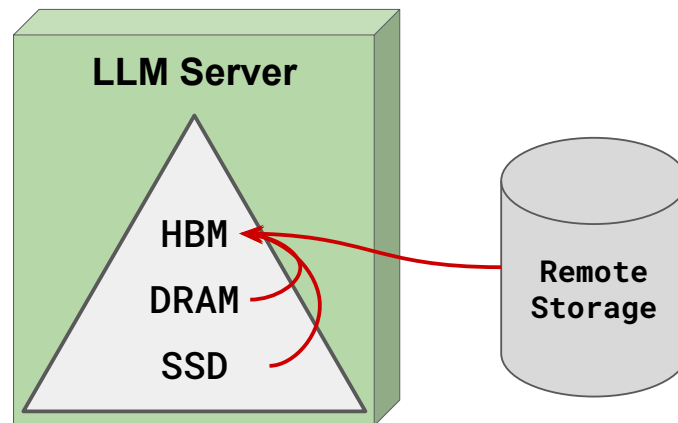
```
Assistant
>How to X?
  One can use Y.
>
```

*Awaiting for  
user/tool  
response*



```
Assistant
>How to X?
  One can use Y.
>Implement Y.
  Thinking...
```

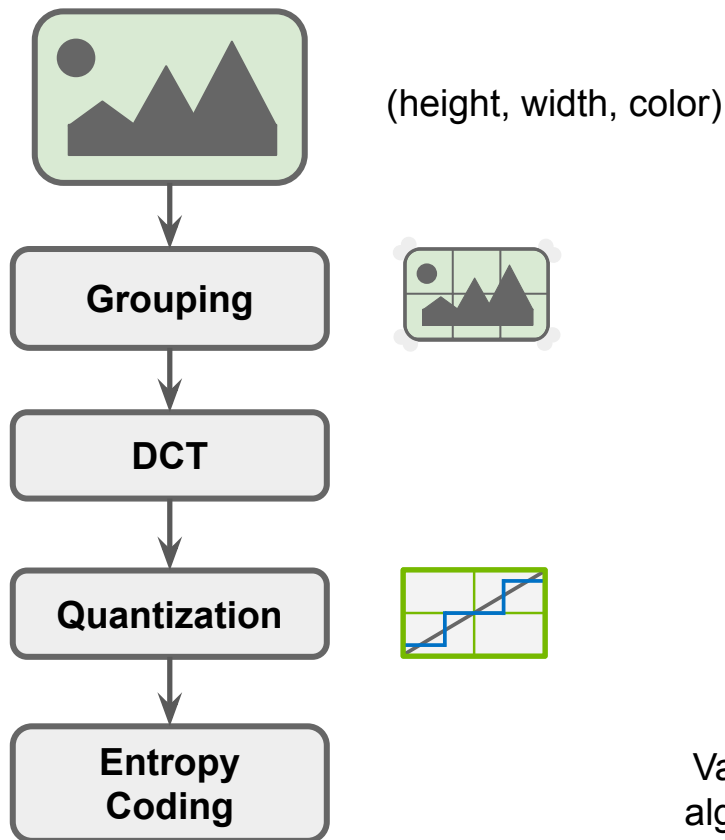
*Processing next  
request*



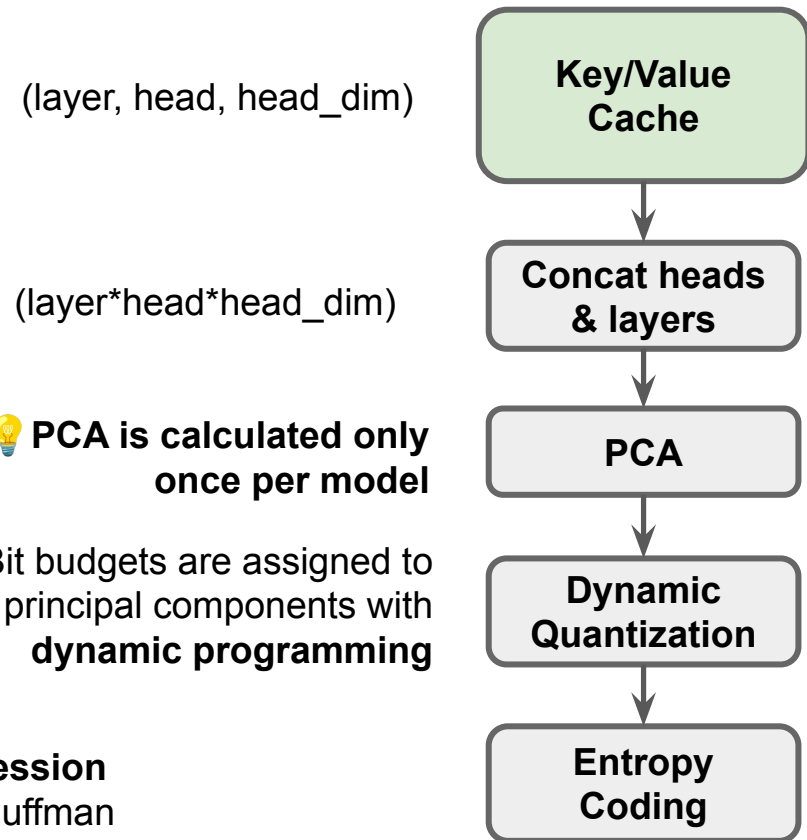
# Our Approach - KVTC

## Transform Coding

### JPEG:



### KV Cache Transform Coding:



Variety of **lossless compression** algorithms: **ANS**, Deflate, Huffman Coding

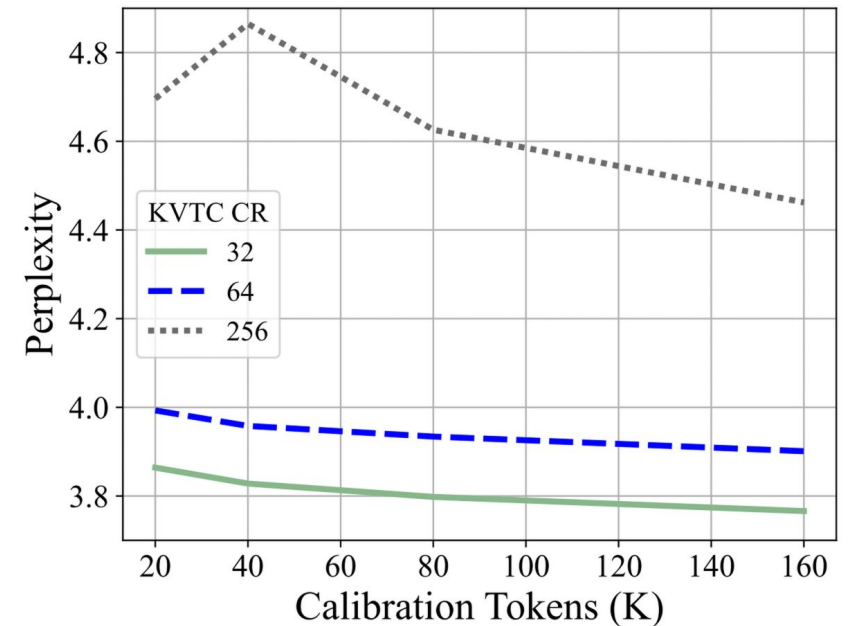
# Generalization Study

More Calibration Data -> Better Performance

KVTC generalizes well for distant calibration data

Method	Calibration Data	CR	GSM	MMLU	QASPER	LITM	RULER-VT
Mistral NeMo 12B							
Vanilla		1	61.9 <sub>1.3</sub>	64.5 <sub>0.4</sub>	38.4	99.5 <sub>0.1</sub>	99.8 <sub>0.2</sub>
GEAR 2bit	-	5	59.8 <sub>1.4</sub>	64.0 <sub>0.4</sub>	38.6	96.9 <sub>0.2</sub>	99.4 <sub>0.3</sub>
KIVI 2bit		5	59.7 <sub>1.4</sub>	64.3 <sub>0.4</sub>	38.2	91.9 <sub>0.3</sub>	98.3 <sub>0.4</sub>
kv <sub>TC</sub> 16x	FineWeb + OpenR1Math	17-20	63.5 <sub>±0.9</sub>	64.7 <sub>±0.3</sub>	37.2 <sub>±0.4</sub>	99.6 <sub>±0.3</sub>	99.7 <sub>±0.0</sub>
	Python	17-20	59.4 <sub>1.4</sub>	65.2 <sub>0.4</sub>	37.4	99.9 <sub>0.0</sub>	99.7 <sub>0.2</sub>
	C	17-20	55.0 <sub>1.4</sub>	65.4 <sub>0.4</sub>	38.0	99.2 <sub>0.1</sub>	99.3 <sub>0.3</sub>
	Assembly	17-20	58.9 <sub>1.4</sub>	65.2 <sub>0.4</sub>	37.0	99.9 <sub>0.0</sub>	99.7 <sub>0.2</sub>

Careful calibration unlocks higher compression ratios



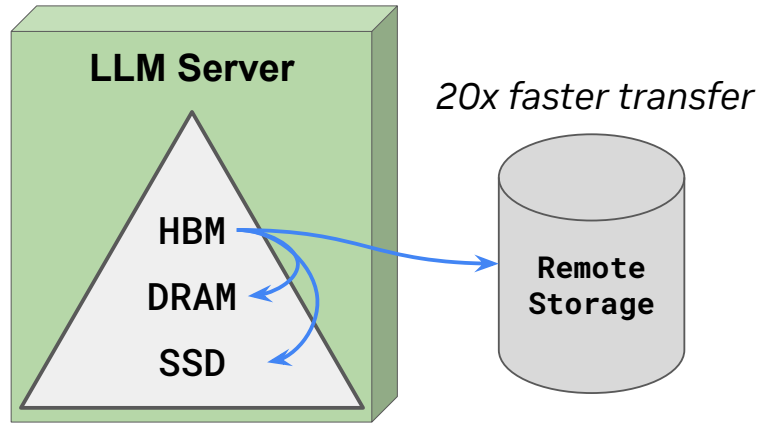
# Results

## Comparison with Other Methods

	Vanilla	GEAR <sub>2-bit</sub>	KIVI <sub>2-bit</sub>	H <sub>2</sub> O	TOVA	xKV	FP8	kvtC <sub>8x</sub>	kvtC <sub>16x</sub>	kvtC <sub>32x</sub>	kvtC <sub>64x</sub>
	Llama 3.1 8B										
CR	1	5	5	8	8	1-5	2	9-10	18-22	34-44	60-88
GSM8K	<b>56.8</b>	52.8	52.8	54.3	54.5	<b>56.6</b>	55.2	<b>57.0</b>	<b>56.9</b>	<b>57.8</b>	<b>57.2</b>
MMLU	<b>60.5</b>	<b>59.6</b>	<b>59.6</b>	44.3	44.8	<b>59.5</b>	<b>60.1</b>	<b>59.8</b>	<b>60.1</b>	<b>60.6</b>	<b>60.7</b>
QASPER	<b>40.4</b>	<b>40.4</b>	39.1	34.3	38.6	35.6	<b>40.8</b>	<b>40.1</b>	<b>40.7</b>	<b>39.4</b>	37.8
LITM	<b>99.4</b>	96.9	88.8	20.2	1.2	<b>99.9</b>	<b>99.4</b>	<b>99.3</b>	<b>99.3</b>	<b>99.1</b>	90.2
RULER-VT	<b>99.8</b>	<b>99.8</b>	<b>98.9</b>	50.4	<b>99.7</b>	<b>99.8</b>	<b>99.9</b>	<b>99.1</b>	<b>99.1</b>	<b>98.9</b>	95.9

Task	Vanilla	kvtC <sub>8x</sub>	kvtC <sub>16x</sub>
Llama 3.3 70B Instruct			
MATH-500	75.6 <sub>±1.92</sub>	73.2 <sub>±1.98</sub>	73.2 <sub>±1.98</sub>
NIAH	100.0	100.0	100.0
LITM	100.0	100.0	100.0
Qwen 2.5 R1 7B			
AIME24	50.9 <sub>±4.9</sub>	52.5 <sub>±3.6</sub>	50.9 <sub>±6.8</sub>
AIME25	40.8 <sub>±4.3</sub>	40.8 <sub>±5.2</sub>	38.3 <sub>±5.5</sub>
LCB	36.7	36.5	31.6
Qwen3.5 4B			
$\tau^2$ -Bench Telecom	95.6	96.5	93.9
NIAH	100.0	100.0	100.0

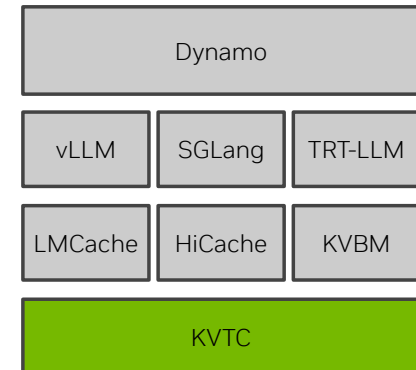
# Performance & Deployment



*20x compression increases  
KV Cache retention time 20x*

- Compression can be performed async on CPU / DPU
- GPU can be utilized for decompression

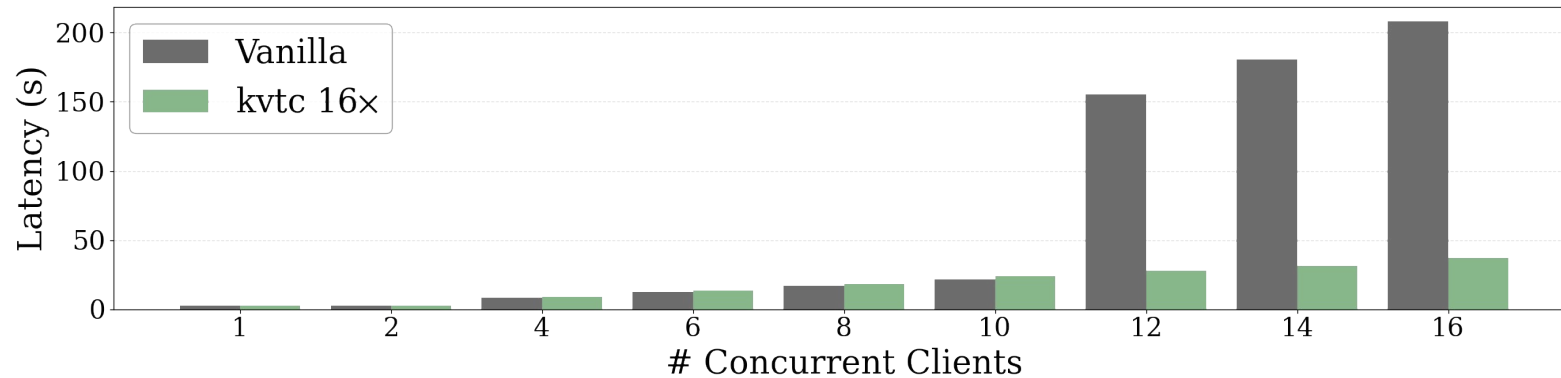
- KVTC can act as a **separate compression layer** in the inference stack
- Easy to integrate - **no alterations** to vLLM/SGLang/...
- Input: KV cache Output: Compressed data



# Performance

## vLLM with LMCache Example

**Llama 3.3 70B FP8 TP2**  
**~64K prefill ~100 tokens per turn**  
256GB DRAM KV Cache storage space





**ICLR**

International Conference On  
Learning Representations

# **KV Cache Transform Coding for Compact Storage in LLM Inference**

Konrad Staniszewski & Adrian Łańcucki

Rio de Janeiro, Brazil, 2026