

Lecture 07: Efficiency Strategies for Large Language Models

Notes

- Lab 2 is due next Friday.
- Think about the project, discuss with me during office hours or after class.
- No quiz today.
- Midterm
 - Oct 29, in class.
 - Will cover materials up to lecture 8 (Oct 22)
 - Will send out a coverage by this weekend
- Final presentation
 - Virtual
 - Dec 16 and Dec 17
 - Final report due at Dec 19



Recap

- Distillation
- Neural architecture search (NAS)
- Low-rank factorization
- Dynamic / Conditional Computing



Topics

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM



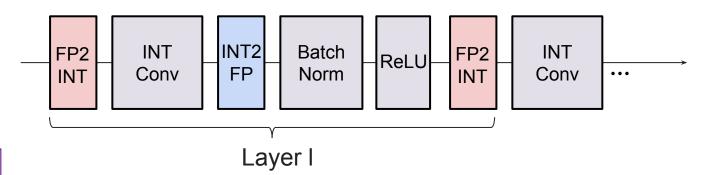
Topics

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM



Quantization

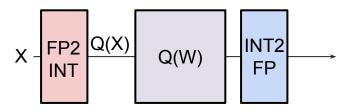
- Quantization on activation needs to be performed dynamically. This will introduce additional compute overhead.
- Also the activation will pass the nonlinear functions, which are usually very sensitive to quantization error, so dequantization is required to convert back to FP 16/32.
- For large models with substantial computational demand, even when input activations are dynamically quantized, this approach can still significantly reduce overall processing latency.





Quantization

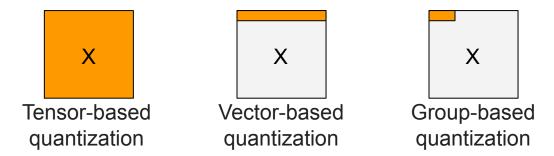
- Quantization on activation needs to be performed dynamically. This will introduce additional compute overhead.
- Also the activation will pass the nonlinear functions, which are usually very sensitive to quantization error, so dequantization is required to convert back to FP 16/32.
- For large models with substantial computational demand, even when input activations are dynamically quantized, this approach can still significantly reduce overall processing latency.





Granularity of Quantization

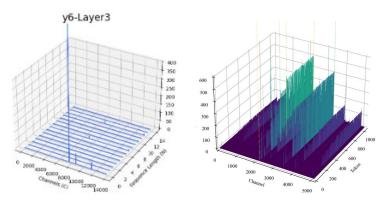
- The activation and weight can be quantized with different granularity:
 - Tensor-based quantization
 - Vector-based quantization
 - Group-based quantization
- A higher quantization granularity will lead to a lower quantization error and a higher hardware implementation cost.





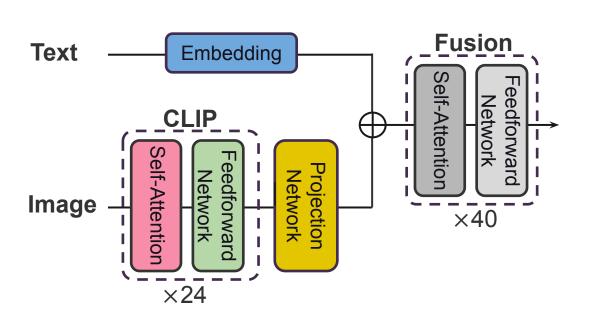
Types of Outlier

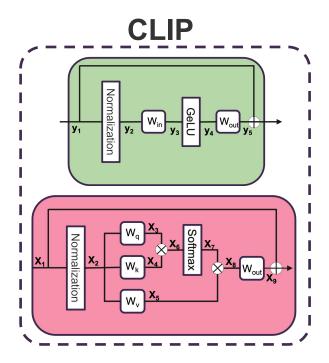
- Massive Activation:
 - For an activation matrix A, an massive activation is an element Aij within it that satisfies:
 - Aij > $\eta \times \text{mean}(|A|)$
 - Aij > γ
 - \circ For example, η=300, γ=50
- Channelwise Outlier:
 - \circ mean(Ai) > $\eta \times std(A) + mean(|A|)$
 - \circ std(Ai) < β
 - \circ For example, η=3, β=0.6





CLIP Model



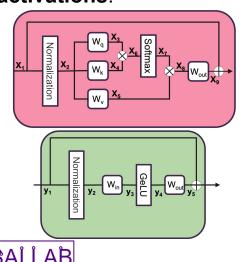


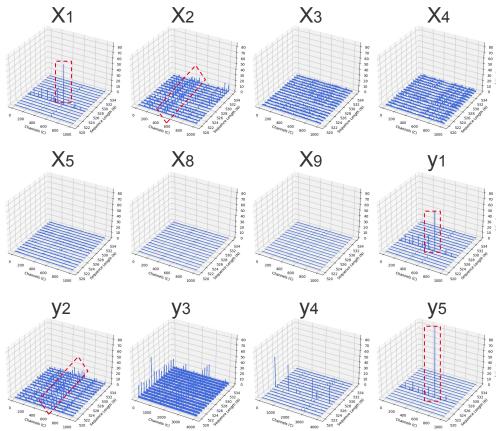


Radford, Alec, et al. "Learning transferable visual models from natural language supervision." *International conference on machine learning*. PmLR, 2021.

Outlier Study: CLIP Activations

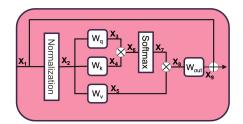
 Outliers with large magnitudes appear at positions x1, y1, and y5, referred to as massive activations.



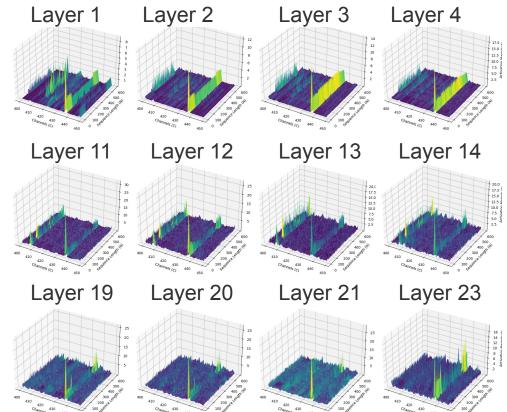


Outlier Study: CLIP Activations

 3D plots of x2 across layers.

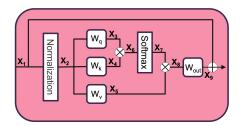


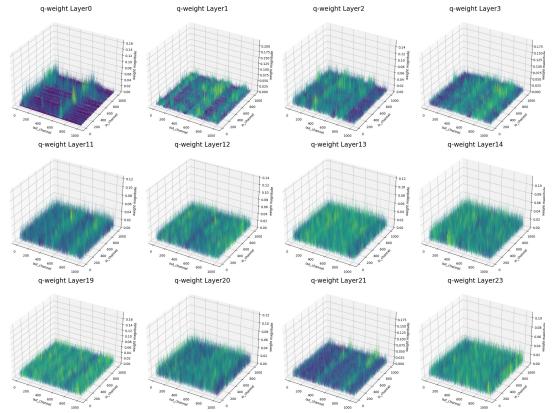
 x2 exhibits channel wise outlier





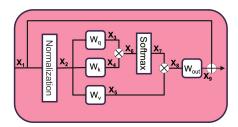
Wq across CLIP layers.

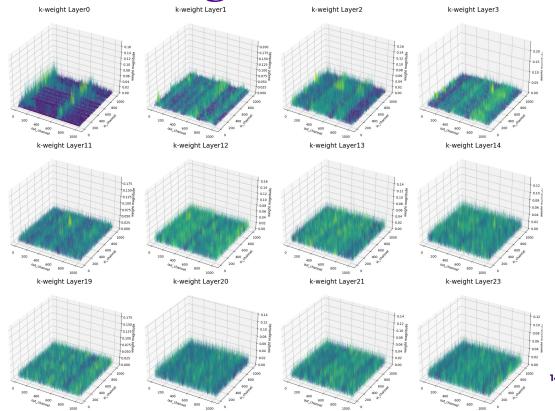






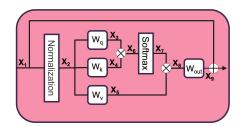
Wk across CLIP layers.

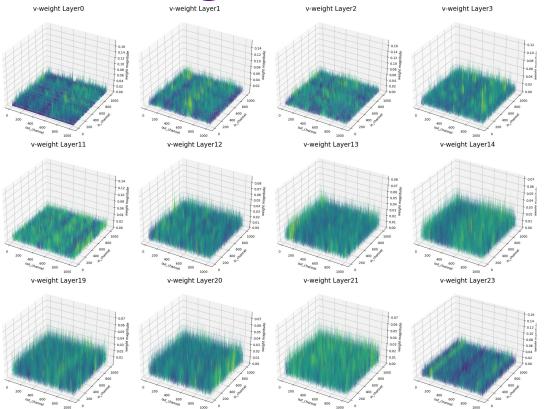






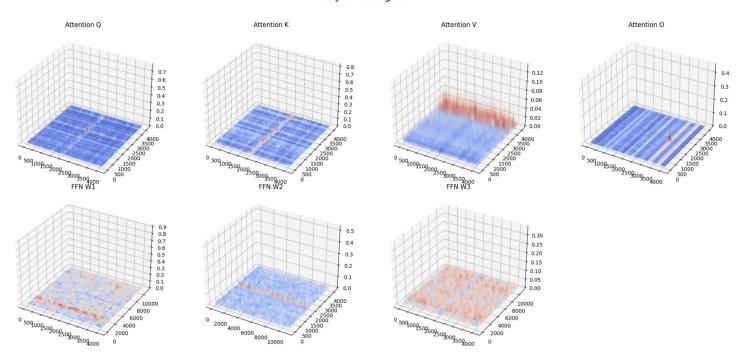
W_V across CLIP layers.



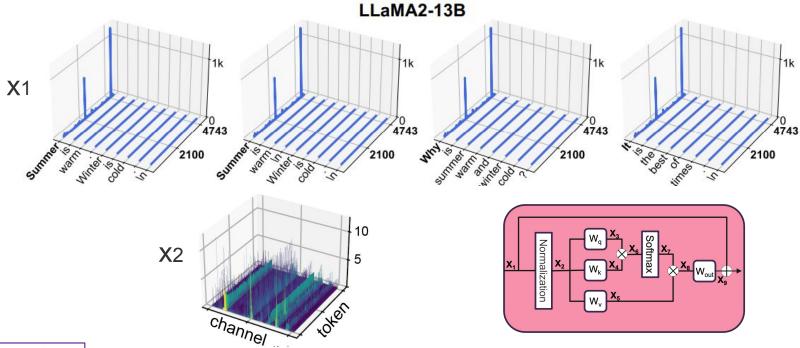




Layer 0 Weights



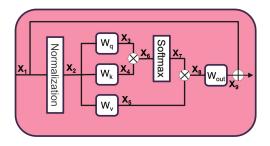
Outlier Study: LLaMA Activations





Why Massive Activations Exist?

- These massive activations are not random spikes or errors, but rather learned, input-agnostic constants that function as implicit bias terms within the model.
- The paper argues that LLMs use massive activations to inject bias into the self-attention computation. These large, fixed activations essentially allow certain tokens to always attract disproportionately high attention.
- No clear conclusion has been reached.





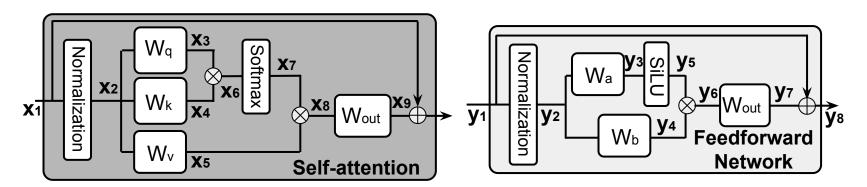
Impact of Massive Activation

Intervention	LLaMA3.2-3B		LLaMA3.1-8B		LLaMA2-13B		GPT-2		Qwen2.5-7B	
	WikiText	C4	WikiText	C4	WikiText	C4	WikiText	C4	WikiText	C4
Original	5.567	10.790	6.941	9.046	4.355	6.405	14.795	19.460	6.520	11.773
TMAs to mean at y7	1124111.75	21046.82	21281.49	1301562.25	1301562.25	6469.42	14.841	19.560	71216.17	66588.86
TMAs to zeroes at y7	1138151.23	21951.41	21601.10	1302018.53	1309211.61	7128.32	14.911	19.928	71835.61	67518.35

- The truncation of massive activation will cause the significant accuracy degradation of the LLM.
- Massive activations also occur in other types of foundation models that utilize attention-based architectures.

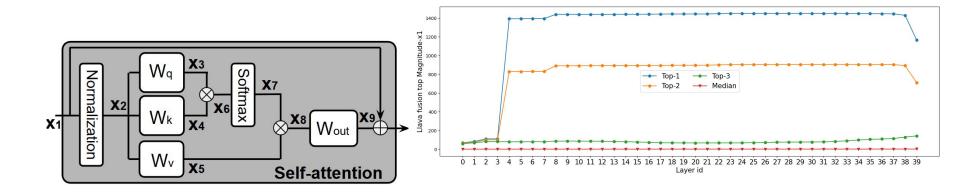


- y7 of layer 3 first produce massive activation (MA), the outlier exists within the first text token.
- x1 and y1 in the following layers contain MA, which may reach into the thousands.
- The MA then propagates through residual link across layers.
- Channel wise outliers exists within x2, y2.



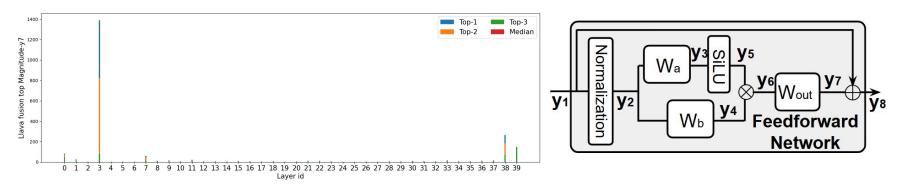


Top magnitude of x1 across layers.





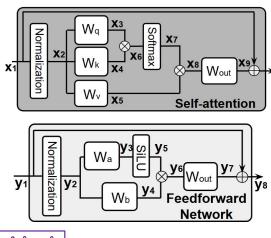
Top magnitude of y7 across layers.

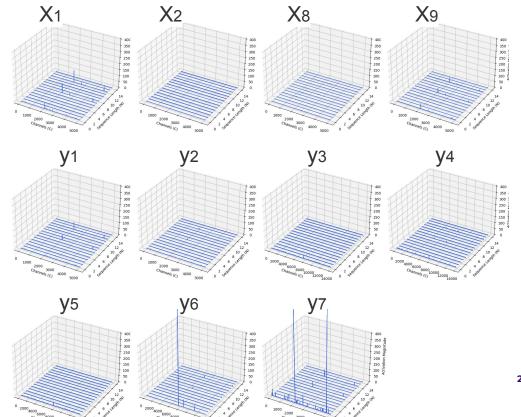


 Truncating the massive activation that caused by residual connection will not lead to large accuracy drop.



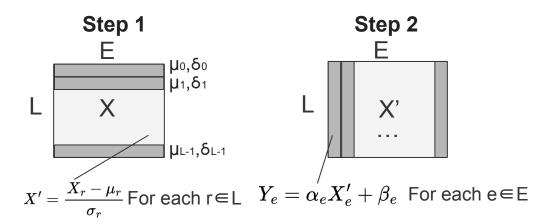
- Distribution of activation within layer 3.
- MA first appears in y6 of token=0.

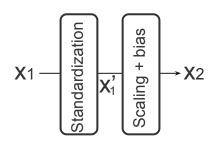




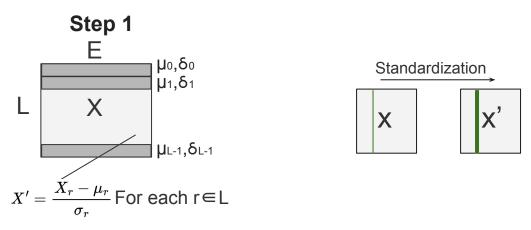


- X1 already has some channelwise outlier, but not significant enough.
- Partial channelwise outlier forms after standardization.
- The scaling process greatly contributes to part of the channelwise outliers.



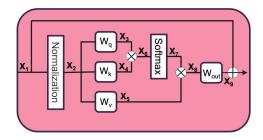




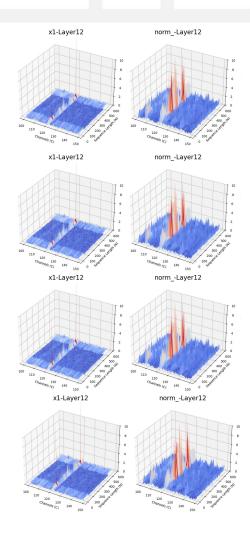


When the standard deviation is low, channel-wise outliers become more pronounced.

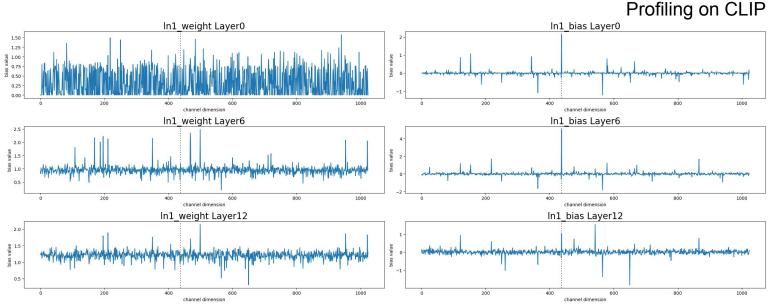




- X1 typically already contains some channelwise outliers.
- Following standardization, the outlier's magnitude becomes more pronounced.

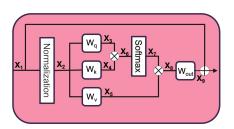




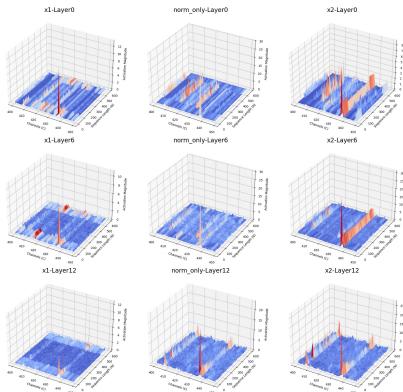


 The scaling and bias factors within the normalization layer also contribute to the channelwise outlier.

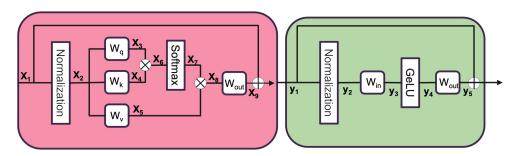




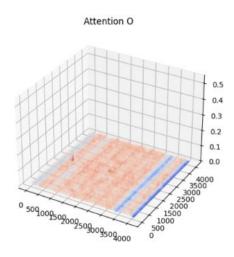
 The scaling and bias factors within the normalization layer also contribute to the channelwise outlier.





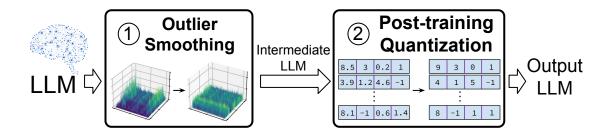


- The question then arises: why does x1 already exhibit some channelwise outliers?
- This is attributed to Wout, which results in y5 having some minor channelwise outliers.
- Additionally, the residual link carries intermediate activations with channelwise outliers, which are also generated by the Wout from preceding layers.





Quantization Strategy for LMs



 When performing post-training quantization on a LLM, it's common to include a step of outlier smoothing prior to the quantization process.



Topics

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM

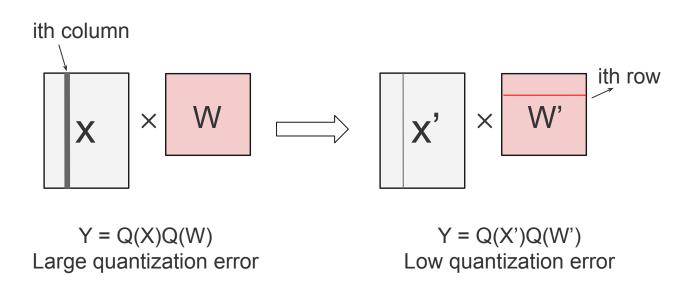


Topics

- Large Model Data Distribution
- Large Model Quantization
 - Smoothing Techniques
 - Quantization Techniques
- Large Model Pruning
- Low-rank Decomposition for LLM

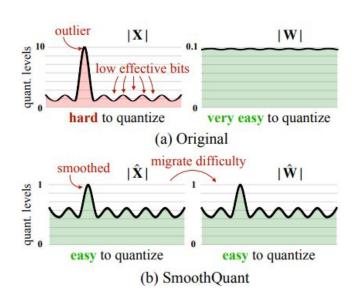


SmoothQuant





SmoothQuant

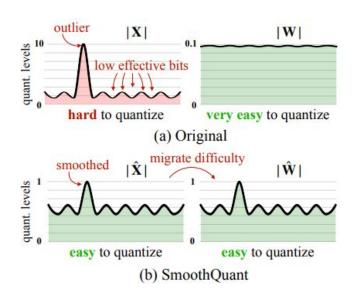


- The intermediate results within LLM usually have a lot of outliers.
- SmoothQuant smooths the activation outliers by offline migrating the quantization difficulty from activations to weights with a mathematically equivalent transformation.

$$\mathbf{Y} = (\mathbf{X} \operatorname{diag}(\mathbf{s})^{-1}) \cdot (\operatorname{diag}(\mathbf{s})\mathbf{W}) = \hat{\mathbf{X}}\hat{\mathbf{W}}$$



SmoothQuant



$$\mathbf{s}_j = \max(|\mathbf{X}_j|)^{\alpha} / \max(|\mathbf{W}_j|)^{1-\alpha}$$

• α is a hyperparameter.

$$\mathbf{Y} = (\mathbf{X} diag(\mathbf{s})^{-1}) \cdot (diag(\mathbf{s})\mathbf{W}) = \hat{\mathbf{X}} \hat{\mathbf{W}}$$

 α is determined from the calibration dataset.



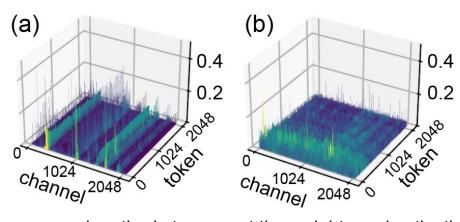
Activation-Aware Weight Quantization (AWQ)

$$\mathbf{s}^* = \operatorname*{arg\,min}_{\mathbf{s}} \mathcal{L}(\mathbf{s})$$

$$\mathcal{L}(\mathbf{s}) = \|Q(\mathbf{W} \cdot \mathrm{diag}(\mathbf{s}))(\mathrm{diag}(\mathbf{s})^{-1} \cdot \mathbf{X}) - \mathbf{W}\mathbf{X}\|$$

 AWQ improves the performance of smoothquant by making "s" learnable.

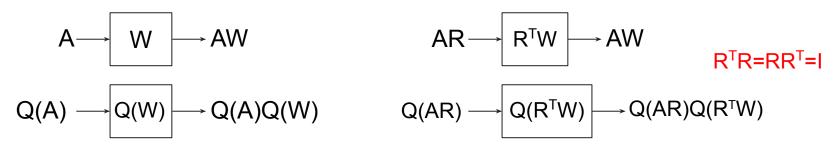
PPL↓			Llama-2		LLaMA				
		7B	13B	70B	7B	13B	30B	65B	
FP16	=	5.47	4.88	3.32	5.68	5.09	4.10	3.53	
INT3 g128	RTN GPTQ GPTQ-R AWQ	6.66 6.43 6.42 6.24	5.52 5.48 5.41 5.32	3.98 3.88 3.86 3.74	7.01 8.81 6.53 6.35	5.88 5.66 5.64 5.52	4.88 4.88 4.74 4.61	4.24 4.17 4.21 3.95	
INT4 g128	RTN GPTQ GPTQ-R AWQ	5.73 5.69 5.63 5.60	4.98 4.98 4.99 4.97	3.46 3.42 3.43 3.41	5.96 6.22 5.83 5.78	5.25 5.23 5.20 5.19	4.23 4.24 4.22 4.21	3.67 3.66 3.66 3.62	



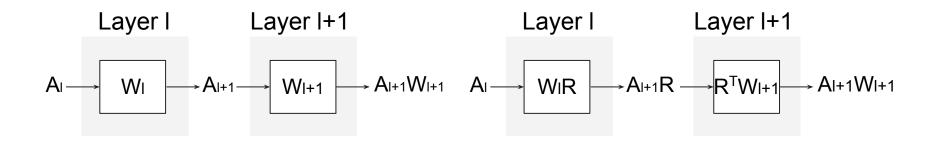
- QuaRot introduces a novel methods to convert the weights and activation of LLM.
- After conversion, most of the outliers within the activation and weights are removed.
- This conversion introduces almost no additional cost during the inference.



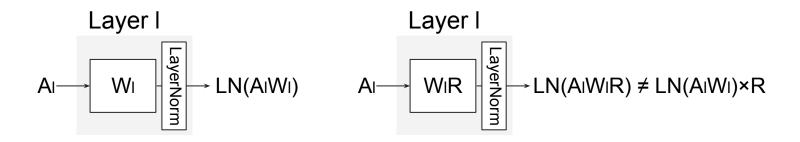
- Assume Y = AW, where A may have outliers, quantizing A and W as Q(A) and Q(W) could result in increased quantization error. Consequently, Q(A)Q(W) may differ significantly from AW.
- With QuaRot, a orthogonal matrix is applied to eliminate the outliers within A.



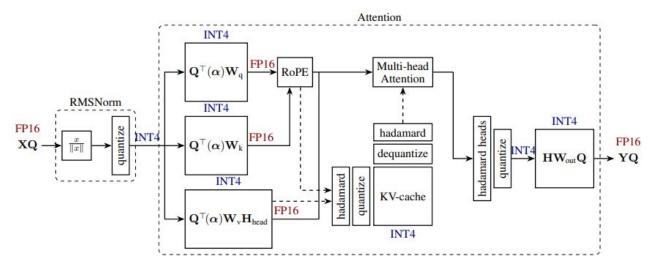
 R¹W can be computed offline, AR can be generated by modifying the weight matrices of the last layer.



 R^TW can be computed offline, AR can be generated by modifying the weight matrices of the last layer.



 However, if a nonlinear function follows the linear layers, the Hadamard transform must be computed dynamically.



- For some of the layers, the conversion needs to be performed online.
- We can use Hadamard matrix, which consists of only 1 and -1 to facilitate the matrix multiplications.



DuQuant



- Permutation matrix is another types of matrix which cause no change on the output.
- Duquant combines three types of computational invariance operation, including the scalar-based, rotation-based and permutation-based operations for mitigating the outliers.

DuQuant

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad P = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \quad A' = AP = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 3 & 2 & 1 \\ 6 & 5 & 4 \\ 9 & 8 & 7 \end{bmatrix}$$

$$P^T = egin{bmatrix} 0 & 0 & 1 \ 0 & 1 & 0 \ 1 & 0 & 0 \end{bmatrix}$$

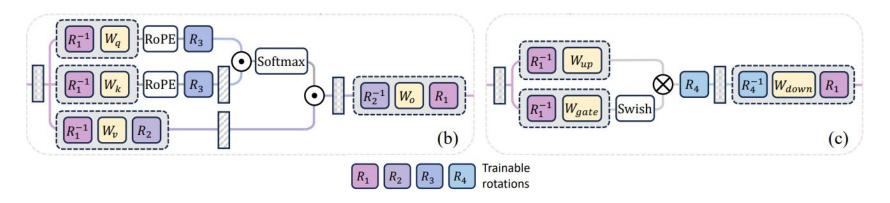
$$P^T = egin{bmatrix} 0 & 0 & 1 \ 0 & 1 & 0 \ 1 & 0 & 0 \end{bmatrix} \hspace{1cm} A_{ ext{recovered}} = A'P^T = egin{bmatrix} 3 & 2 & 1 \ 6 & 5 & 4 \ 9 & 8 & 7 \end{bmatrix} egin{bmatrix} 0 & 0 & 1 \ 0 & 1 & 0 \ 1 & 0 & 0 \end{bmatrix} = egin{bmatrix} 1 & 2 & 3 \ 4 & 5 & 6 \ 7 & 8 & 9 \end{bmatrix}$$

Quantization group





SpinQuant



 $\underset{R \in \mathcal{M}}{\operatorname{arg\,min}} \, \mathcal{L}_Q(R_1, R_2 \mid W, X)$

- SpinQuant optimizes (or learns) the rotation matrices to obtain the minimal changes on the training loss.
- We have to ensure the rotational matrix still satisfies the orthogonal property → Cayley Optimization.

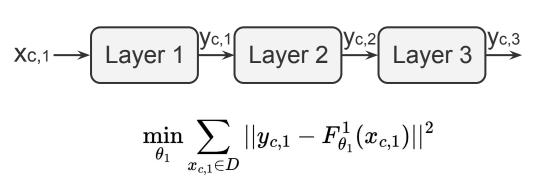


Topics

- Large Model Data Distribution
- Large Model Quantization
 - Smoothing Techniques
 - Quantization Techniques
- Large Model Pruning
- Low-rank Decomposition for LLM



Sequential Update



- We use a calibration dataset and profile some data xc, yc on the first layer.
- After that, we use these data to train the optimal quantization hyperparameters, smoothing hyperparameters.



Calibration Dataset

Xc,1
$$\longrightarrow$$
 Layer 1 $\xrightarrow{y_{c,1}}$ Layer 2 $\xrightarrow{y_{c,2}}$ Layer 3 $\xrightarrow{y_{c,3}}$ $x_{c,2} = F_{\theta_1^*}^1(x_{c,1})$

 After that, we regenerate the calibration dataset using the new weight values in layer 1, results in (x_{c,2}, y_{c,2}).



AdaQuant

$$\begin{split} \left(\hat{\Delta}_{w}, \hat{\Delta}_{x}, \hat{V}\right) &= \\ \underset{\Delta_{w}, \Delta_{x}, V}{\operatorname{arg \, min}} ||WX - Q_{\Delta_{w}}(W')Q_{\Delta_{x}}(X)||^{2} \end{split}$$

$$\begin{pmatrix} \hat{\Delta}_{w_l}, \hat{\Delta}_{x_l}, \hat{V}_l \end{pmatrix}
= \underset{\Delta_{w_l}, \Delta_{x_l}, V_l}{\operatorname{arg min}} ||W_l X_l - Q_{\Delta_{w_l}}(W'_l) \cdot Q_{\Delta_{x_l}}(X_l^q)||^2
X_q = \sigma(Q_{\Delta_{w_{l-1}}}(W_{l-1} + V_{l-1}) \cdot Q_{\Delta_{x_l}}(X_{l-1}^q)),$$

AdaQuant

Sequential AdaQuant

A small calibration dataset is used to find the optimal scale and V.

Pruning Criteria: Training Loss Change

$$egin{aligned} \mathcal{L}(w') &= \mathcal{L}(w) +
abla \mathcal{L}(w)^T (w-w') \ &+ rac{1}{2} (w-w')^T
abla^2 \mathcal{L}(w) (w-w') \end{aligned}$$

 When reflecting on each individual value, the pruning criteria becomes:

$$\mathcal{L}(0) - \mathcal{L}(w) = rac{d\mathcal{L}(w)}{dw}w + rac{1}{2}rac{d^2\mathcal{L}(w)}{dw^2}w^2$$

The gradient is usually estimated to zero.



Optimal Brain Surgeon (OBS)

Propose ways to prune the weight, and finetune the remaining weights.

$$\mathcal{L}(w') - \mathcal{L}(w) pprox rac{1}{2} (w - w')^T H(w - w') \quad where \,\, H =
abla^2(w)$$

If we choose to prune qth weight, the problem we will try to solve is:

$$\min_{s} rac{1}{2} s^T H s \quad where \,\, s = w - w'$$

- s.t. $e_q^T s = w_q^{}$ eq is the one-hot vector with qth element equaling 1 and 0 otherwise
- Solve this problem, we have:

$$s^* = rac{-w_q}{[H^{-1}]_{qq}} H^{-1} e_q \qquad L_q = rac{w_q^2}{2[H^{-1}]_{qq}}$$



How to update rest weights How to select the weight

Optimal Brain Quantization (OBQ)/GPTQ

$$\operatorname{argmin}_{\widehat{\mathbf{W}}_{\ell}} \quad ||\mathbf{W}_{\ell} \mathbf{X}_{\ell} - \widehat{\mathbf{W}}_{\ell} \mathbf{X}_{\ell}||_{2}^{2} \quad \text{s.t.} \quad \mathcal{C}(\widehat{\mathbf{W}}_{\ell}) > C.$$

- Each element of W is pruned one-by-one, until the target sparsity ratio is achieved.
- The same idea can be applied to perform quantization, where each value within W is rounded to its neared quantized value, and rest of the weight will adjusted accordingly.

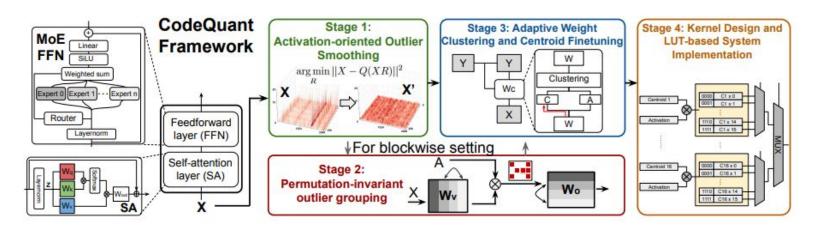
$$w_p = \operatorname{argmin}_{w_p} \frac{w_p^2}{[\mathbf{H}^{-1}]_{pp}}, \quad \boldsymbol{\delta_p} = -\frac{w_p}{[\mathbf{H}^{-1}]_{pp}} \cdot \mathbf{H}_{:,p}^{-1}, \quad w_p = \operatorname{argmin}_{w_p} \frac{(\operatorname{quant}(w_p) - w_p)^2}{[\mathbf{H}^{-1}]_{pp}}, \quad \boldsymbol{\delta_p} = -\frac{w_p - \operatorname{quant}(w_p)}{[\mathbf{H}^{-1}]_{pp}} \cdot \mathbf{H}_{:,p}^{-1}$$



Frantar, Elias, and Dan Alistarh. "Optimal brain compression: A framework for accurate post-training quantization and pruning." *Advances in Neural Information Processing Systems* 35 (2022): 4475-4488.

Frantar, Elias, et al. "Gptq: Accurate post-training quantization for generative pre-trained transformers." *arXiv preprint arXiv:2210.17323* (2022).

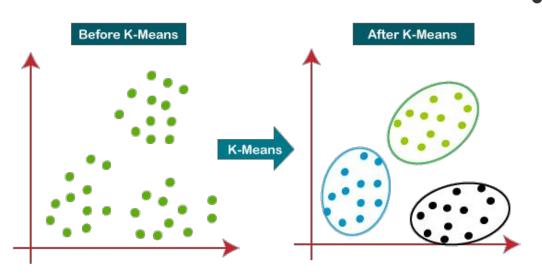
CodeQuant



 Compared to quantization, the clustering operation exhibits greater tolerance to outliers.



LUT-based LLM Inference



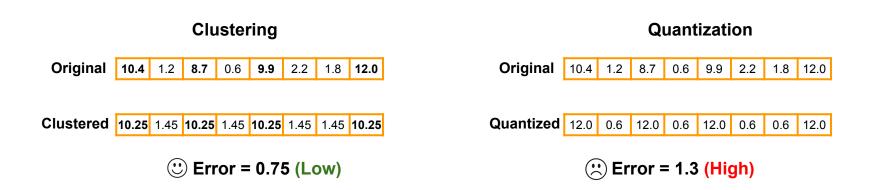
Input of 4 bits, weight with K centroids.

Input centroid	Weight centroid	Results			
a1	w1	0011			
a2	w1	0010			
•••					
ak	wk	1010			



 The multiplication and addition operations can be performed using lookup table (LUT).

Advantage of Clustering



Clustering can be used to hand excessive outlier.

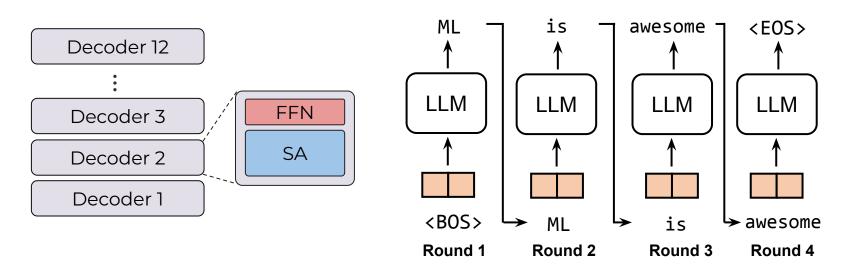


Topics

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM



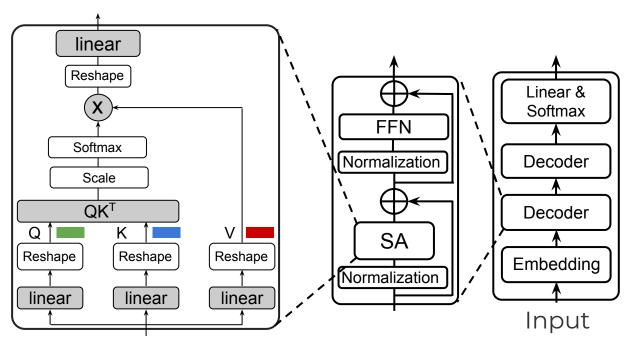
Transformers as a Generative AI Tool



Each token is generated in an autoregressive manner.



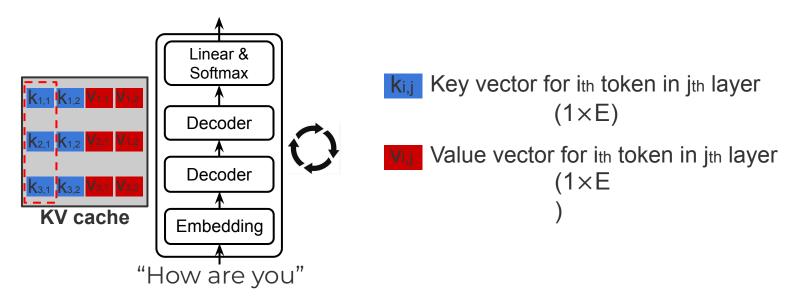
Transformers as a Generative AI Tool





We need to buffer the v and k for later usage.

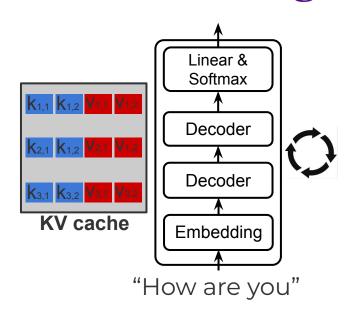
LLM: Prefilling

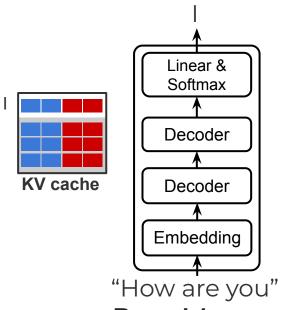


 During the prefilling stage, LLM processes the entire prompt, or context tokens jointly, saving the KV vectors into the memory.



LLM: Decoding



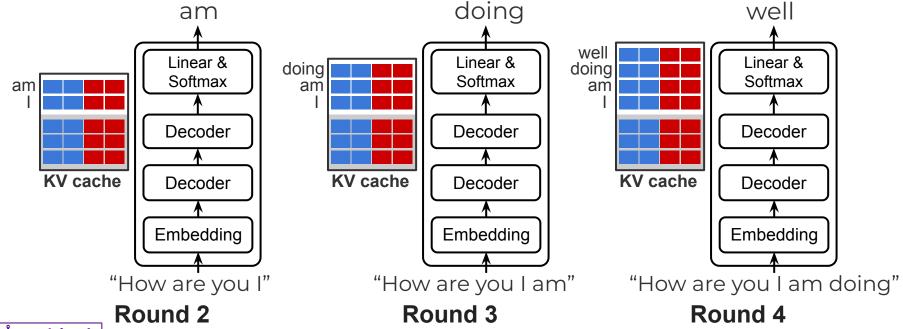


Round 1

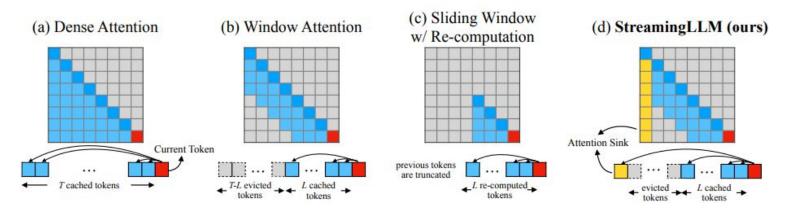
During the decoding stage, LLM generates the responses in an autoregressive way.



LLM: Decoding



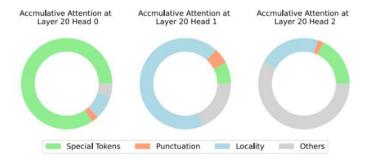
Streaming-LLM



- We observe an interesting phenomenon, namely attention sink, that keeping the KV of initial tokens will largely recover the performance of window attention.
- There are strong attention scores towards initial tokens as a "sink" even if they are not semantically important.



Pruning on Large Models: KV Cache Pruning

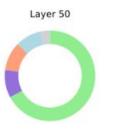


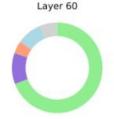






- We show the attention score of each token.
- Different attention heads usually have different importance scores on KV vectors.
- The importance of KV vectors also varies across layers.



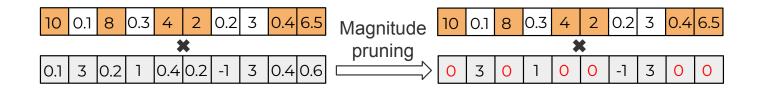






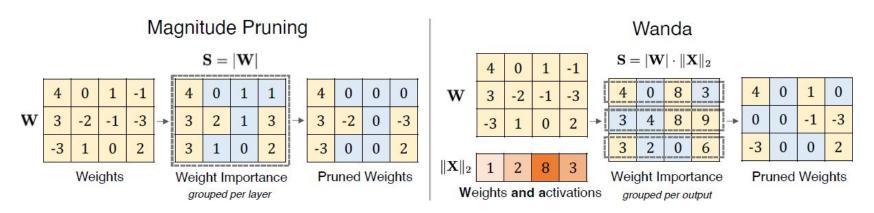
Drawbacks of Magnitude Pruning

• The major drawback of magnitude based pruning is that it does not consider the impact of the input when making the pruning decision.





LLM Pruning: Wanda



- Prune the weights by considering the input statistics.
- For each weight, if the corresponding input's magnitude is large, the output will also be large.
- Need some samples for calibration.



LLM Pruning: Wanda

			LLaMA			LLaMA-2			
Method	Weight Update	Sparsity	7B	13B	30B	65B	7B	13B	70B
Dense	-	0%	5.68	5.09	4.77	3.56	5.12	4.57	3.12
Magnitude	×	50%	17.29	20.21	7.54	5.90	14.89	6.37	4.98
SparseGPT	✓	50%	7.22	6.21	5.31	4.57	6.51	5.63	3.98
Wanda	×	50%	7.26	6.15	5.24	4.57	6.42	5.56	3.98
Magnitude	×	4:8	16.84	13.84	7.62	6.36	16.48	6.76	5.54
SparseGPT	1	4:8	8.61	7.40	6.17	5.38	8.12	6.60	4.59
Wanda	×	4:8	8.57	7.40	5.97	5.30	7.97	6.55	4.47
Magnitude	×	2:4	42.13	18.37	9.10	7.11	54.59	8.33	6.33
SparseGPT	✓	2:4	11.00	9.11	7.16	6.28	10.17	8.32	5.40
Wanda	X	2:4	11.53	9.58	6.90	6.25	11.02	8.27	5.16

Table 3: WikiText perplexity of pruned LLaMA and LLaMA-2 models. Wanda performs competitively against prior best method SparseGPT, without introducing any weight update.



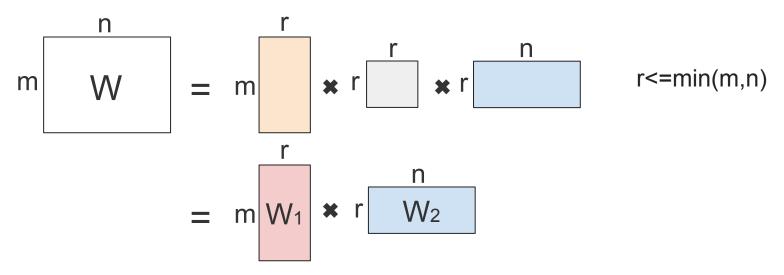
Topics

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM



Low Rank Optimization for DNN Efficiency

Weight tensors can be decomposed into:



- We can train the W₁ and W₂ in the DNN instead of W.
- Less storage is required.



Singular Value Decomposition

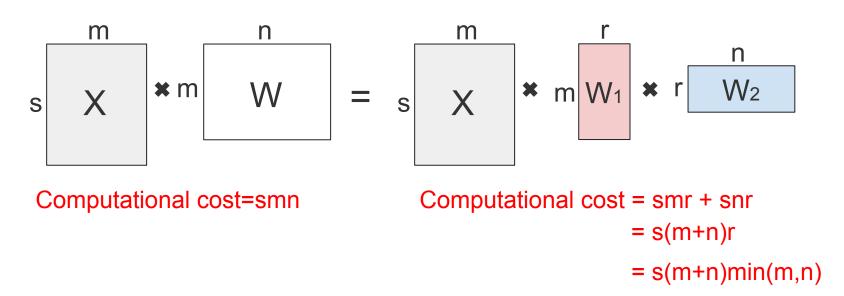
- W: Input matrix
 - m**×**n matrix
- V (n*r matrix) contains the orthonormal eigenvectors of W^TW
- U (m*r matrix) contains the orthonormal eigenvectors of WW^T
- R: Singular value matrix
 - o r∗r diagonal matrix, r is the rank of W

Before: mn After: mr + rn = r(m+n) = min(m,n)(m+n) assume W is full rank

 Without rank truncation, the number of parameters increases.



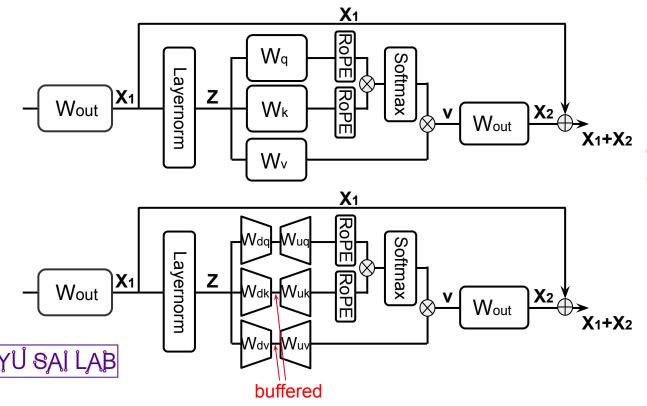
Singular Value Decomposition



Without rank truncation, the computational cost will increase.

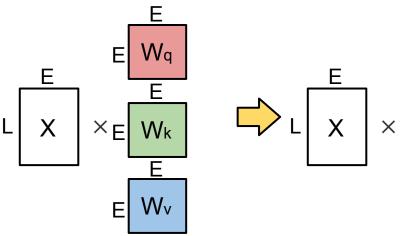


Singular Value Decomposition



SVD decomposition can potentially save the MAC operations, memory storage and KV cache size.

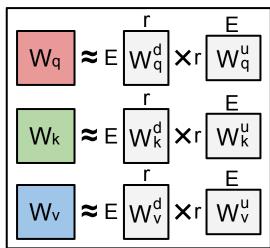
QSVD



Total MACs: 3LE²

Total weight parameters: 3E²

Total cache size: 2LE



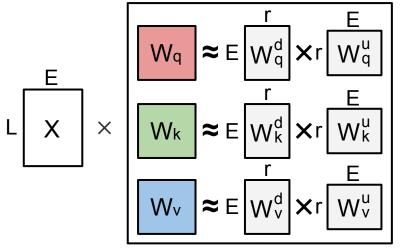
Total MACs: 6LrE

Total weight parameters: 6rE

Total cache size: 2rL



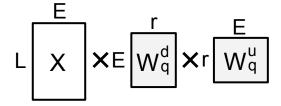
QSVD



Total MACs: 6LrE

Total weight parameters: 6rE

Total cache size: 2rL



Total MACs: 6LrE

Total weight parameters: 6rE

Total cache size: 2rL



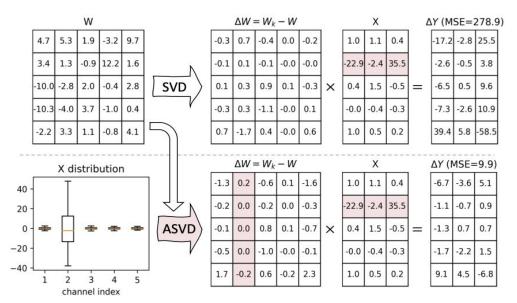
ASVD

 ASVD changes the objective considering impact of activation x

$$L = ||WX - SVD(W)X||$$

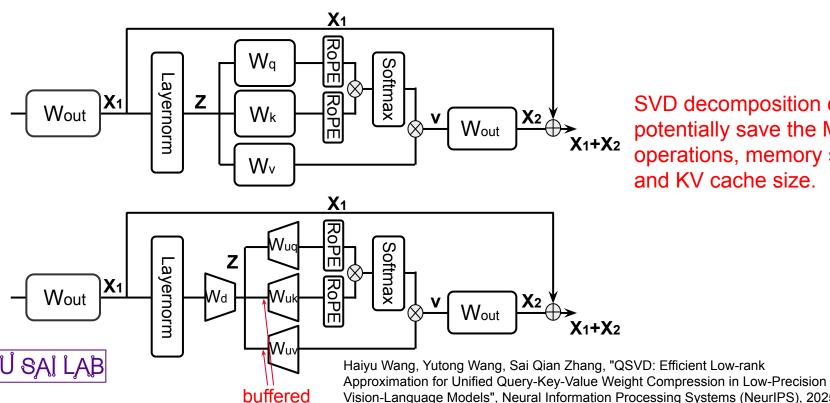
$$L = ||WX - SVD(WS)S^{-1}X||$$
S is a diagonal matrix

 To mitigate impact of channel wise outlier within X, we scale the input activation to mitigate the outliers.



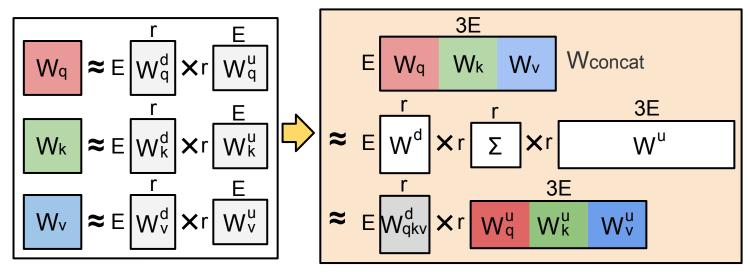


Yuan, Zhihang, et al. "ASVD: Activation-Aware Singular Value Decomposition for Compressing Large Language Models." 2025.



SVD decomposition can potentially save the MAC operations, memory storage and KV cache size.

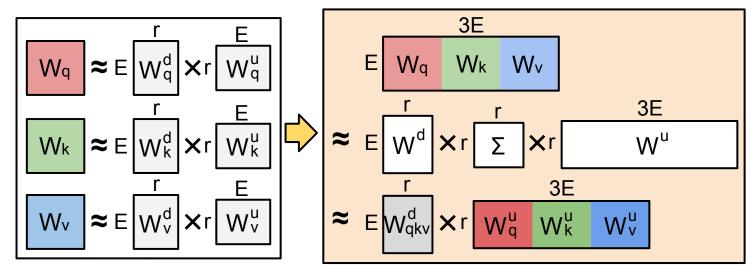
Vision-Language Models", Neural Information Processing Systems (NeurIPS), 2025.



Concat and SVD



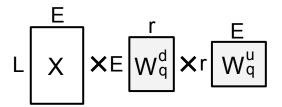
$$[W_q, W_k, W_v] = W_{concat} \approx W^d \times \Sigma \times W^u$$



Down/up projection

$$\begin{aligned} W_{qkv}^d &= W_r^d \Sigma_r^{1/2}, \ [W_q^u, \ W_k^u, \ W_v^u] = \Sigma_r^{1/2} W_r^u \\ [W_q, \ W_k, \ W_v] &= W_{qkv}^d \times [W_q^u, \ W_k^u, \ W_v^u] \end{aligned}$$

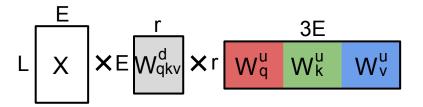




Total MACs: 6LrE

Total weight parameters: 6rE

Total cache size: 2rL

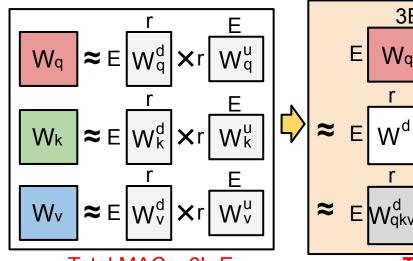


Total MACs: 4LrE

Total weight parameters: 4rE

Total cache size: rL

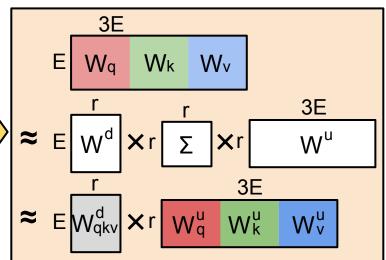




Total MACs: 6LrE

Total weight parameters: 6rE

Total KV cache size: 2rL

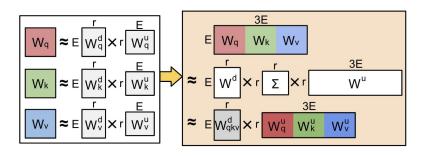


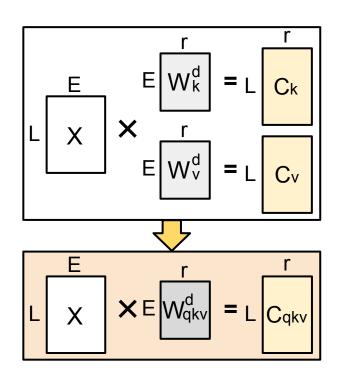
Total MACs: 4LrE

Total weight parameters: 4rE

Total KV cache size: rL

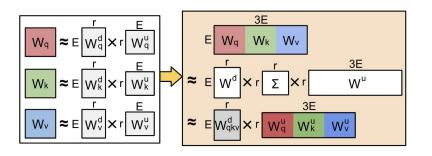


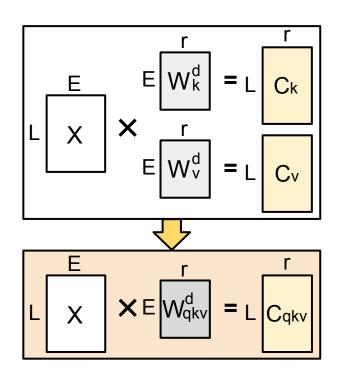




Shared latent

$$[Q,\ K,\ V] = XW_{qkv}^d \times [W_q^u,\ W_k^u,\ W_v^u] = C_{qkv} \times [W_q^u,\ W_k^u,\ W_v^u]$$





Reconstruction



$$K = C_{qkv}W_k^u, \ V = C_{qkv}W_v^u$$

How to assign the rank r to each layer?

$$rac{L(\sigma + \Delta\sigma) - L(\sigma)}{\Delta\sigma} pprox rac{dL}{d\sigma} \;
ightharpoons \; \mathbb{E}_D(|rac{dL}{d\sigma}|)$$

The importance of each layer can be expressed as:

$$\sum_i \mathbb{E}_D(|rac{dL}{d\sigma_i}|)$$

 Given a total rank budget R, we can allocate the rank for each layer in proportion to the importance score.



QSVD Performance

	Method	ScienceQA-IMG ↑							
		Acc.	Hw cost	Acc.	Hw cost	Acc.	Hw cost	Acc.	Hw cost
SmolVLM 2B	ASVD SVDLLM	53.84% 65.89%	$R_1:100\%$ $R_2:50.0\%$	7.88% 34.61%	$R_1:90.0\%$ $R_2:42.5\%$	0.69% 9.07%	$R_1:80.0\%$ $R_2:35.0\%$	0.10% 3.02%	$R_1:70.0\%$ $R_2:27.5\%$
	QSVD-noQ	83.78%	R_1 :100% R_2 :37.5%	81.70%	R_1 :90.0% R_2 :33.75%	79.57%	R_1 :80.0% R_2 :30.0%	77.64%	R_1 :70.0% R_2 :26.25%
	FP16	Accuracy: 84,53%							
LLaVA-Next 7B	ASVD SVDLLM	50.72% 65.94%	$R_1:63.3\%$ $R_2:22.5\%$	47.15% 66.14%	$R_1:60.0\%$ $R_2:20.0\%$	40.26% 64.90%	$R_1:56.7\%$ $R_2:17.5\%$	25.73% 62.87%	$R_1:53.3\%$ $R_2:15.0\%$
	QSVD-noQ	69.91%	R_1 :60.0% R_2 :22.5%	68.22%	R_1 :53.3% R_2 :20.0%	67.03%	R_1 :46.7% R_2 :17.5%	65.15%	R_1 :40.0% R_2 :15.0%
	FP16	Accuracy: 69.51%							
LLaVA-v1.5 13B	ASVD SVDLLM	64.70% 71.44%	$R_1:63.3\%$ $R_2:22.5\%$	56.92% 71.44%	$R_1:60.0\%$ $R_2:20.0\%$	46.50% 71.29%	$R_1:56.7\%$ $R_2:17.5\%$	42.79% 70.50%	$R_1:53.3\%$ $R_2:15.0\%$
	QSVD-noQ	71.79%	R_1 :60.0% R_2 :22.5%	71.74%	R_1 :53.3% R_2 :20.0%	71.74%	R ₁ :46.7% R ₂ :17.5%	70.80%	R ₁ :40.0% R ₂ :15.0%
	FP16	Accuracy: 71.78%							

 QSVD achieves a comparable and even an better performance than the original model.



Problem of SVD

The problem we have to solve can be formulated as this:

$$egin{argmin} argmin ||W-W'||^2 \ W' = DU \ D \in \mathcal{R}^{E imes r} \ U \in \mathcal{R}^{r imes E} \ \end{array}$$

 However, we know that each element of W has different impact on the final accuracy.



Weighted SVD

The problem we have to solve can be formulated as this:

$$egin{argmin} argmin \sum_{ij} ||a_{ij}W_{ij} - a_{ij}W_{ij}'|| \ W' = DU \ D \in \mathcal{R}^{E imes r} \ U \in \mathcal{R}^{r imes E} \ \end{array}$$

- This problem has no closed-form solution, therefore we can use the deep-learning method to solve this, train D and U such that the objective function is minimized.
- For importance score, we can use the following formula to evaluate.

$$L(w)-L(w')pprox rac{dL}{dw}(w-w')$$

