# ESPACE: DIMENSIONALITY REDUCTION OF ACTIVATIONS FOR MODEL COMPRESSION
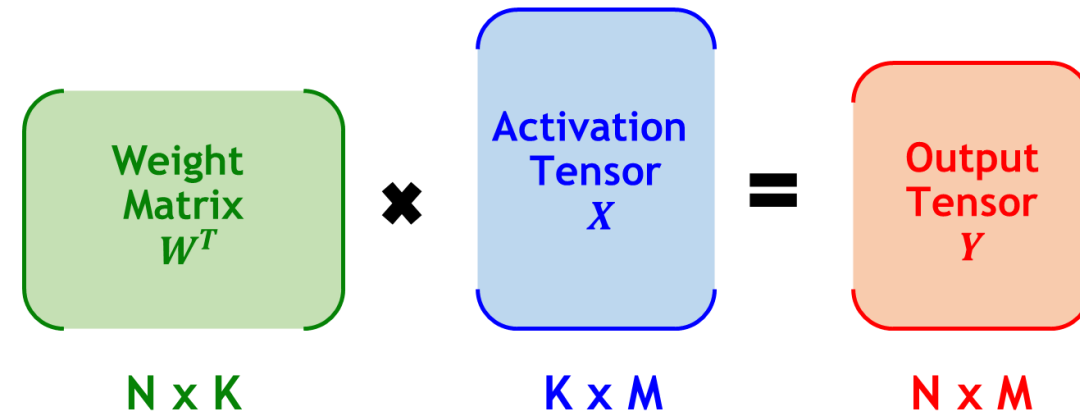
*Charbel Sakr* and Brucek Khailany

# COMPRESSING MODELS USING MATRIX FACTORIZATION

**General Matrix Multiplication (GEMM)**

$$Y = W^T X$$

⬆ trainable parameters: ⬆ expressivity ✅
⬆ model size: ⬆ inference cost ❌

| Weight Matrix $W^T$ | ✖ | Activation Tensor $X$ | = | Output Tensor $Y$ |
|---|---|---|---|---|
| N x K | | K x M | | N x M |

**K = input dimension, N= output dimension, M = batch $*$ sequence**
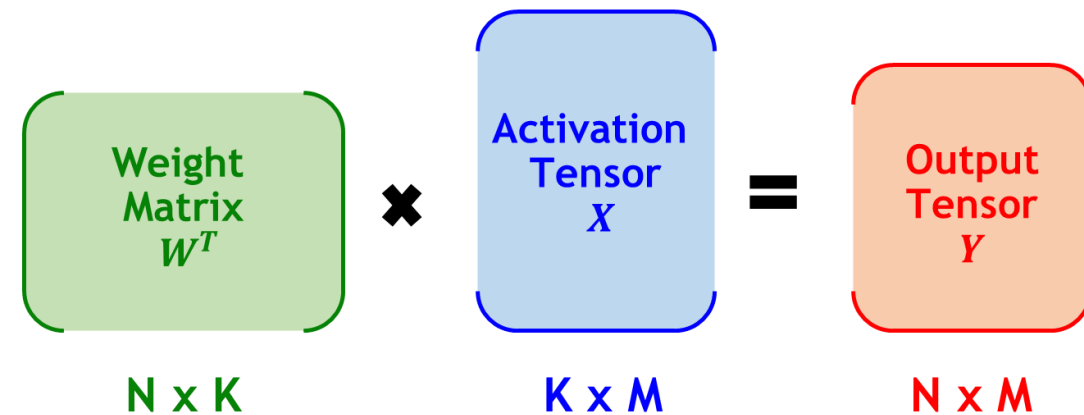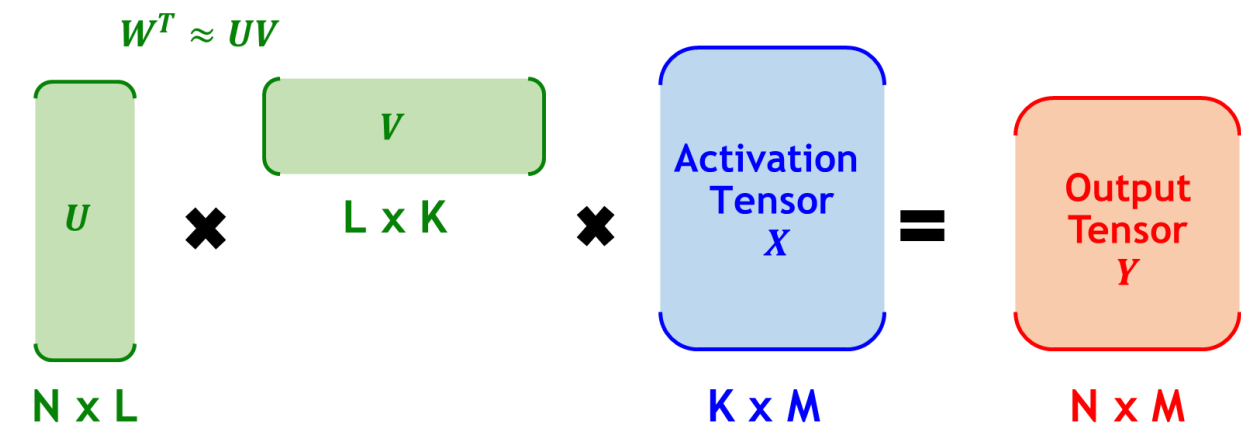
▸ Matrix multiplications between weights and activations

▸ Pervasive operation in LLM serving and training

▸ Weight matrix contains the model parameters

    ▸ it stores its representative power

▸ Activation tensor is generated on the fly for every input

# COMPRESSING MODELS USING MATRIX FACTORIZATION

**General Matrix Multiplication (GEMM)**

$$Y = W^T X$$

⬆️ trainable parameters: ⬆️ expressivity ✅
⬆️ model size: ⬆️ inference cost ❌



N x K      K x M      N x M

**K = input dimension, N= output dimension, M = batch * sequence**

**GEMM with Weight Decomposition**
**(e.g., truncated SVD)**

$$Y \approx UVX$$

⬇️ trainable parameters: ⬇️ expressivity ❌
⬇️ model size: ⬇️ inference cost ✅



N x L      K x M      N x M

**K = input dimension, N= output dimension, M = batch * sequence**
**L = intermediate dimension**

---

▶ Matrix multiplications between weights and activations

▶ Pervasive operation in LLM serving and training

▶ Weight matrix contains the model parameters

    ▶ it stores its representative power
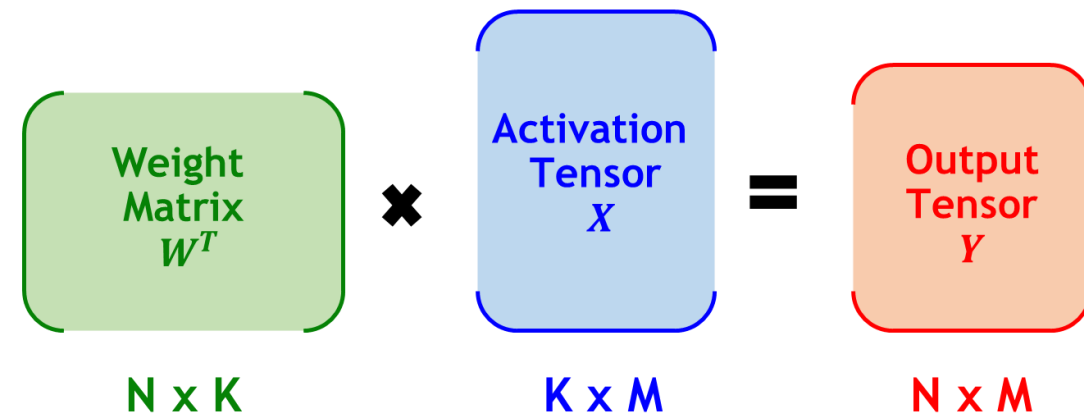
▶ Activation tensor is generated on the fly for every input

▶ Weight matrix can be approximated using factorization

    ▶ e.g. a truncated SVD, where we hope to factorize $W^T \approx UV$

▶ To achieve compression, we need $L \ll K$

    ▶ Since we now have two matrices instead of one

▶ Such decomposition has two limitations

    ▶ fewer weights → lesser representative power

    ▶ Breakage of weight structure → optimizer cannot be recovered

# COMPRESSING MODELS USING MATRIX FACTORIZATION

**General Matrix Multiplication (GEMM)**

$$Y = W^T X$$

⬆️ trainable parameters: ⬆️ expressivity ✅
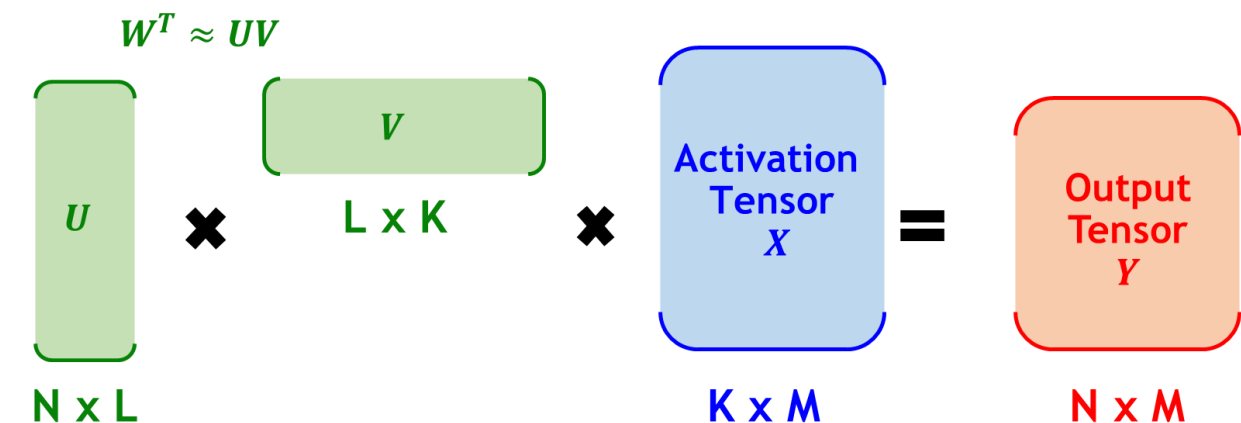⬆️ model size: ⬆️ inference cost ❌



**Weight Matrix** $W^T$ ✖️ **Activation Tensor** $X$ = **Output Tensor** $Y$

N x K          K x M          N x M

K = input dimension, N= output dimension, M = batch ∗ sequence

**GEMM with Weight Decomposition**
**(e.g., truncated SVD)**

$$Y \approx UVX$$

⬇️ trainable parameters: ⬇️ expressivity ❌
⬇️ model size: ⬇️ inference cost ✅

$W^T \approx UV$



$U$ ✖️ $V$ / L x K ✖️ **Activation Tensor** $X$ = **Output Tensor** $Y$

N x L          K x M          N x M

K = input dimension, N= output dimension, M = batch ∗ sequence
L = intermediate dimension

▶ Matrix multiplications between weights and activations

▶ Pervasive operation in LLM serving and training

▶ Weight matrix contains the model parameters

　　▶ it stores its representative power

▶ Activation tensor is generated on the fly for every input

▶ **Can we decompose activations instead?**

▶ Weight matrix can be approximated using factorization

　　▶ e.g. a truncated SVD, where we hope to factorize $W^T \approx UV$

▶ To achieve compression, we need $L \ll K$

　　▶ Since we now have two matrices instead of one

▶ Such decomposition has two limitations

　　▶ fewer weights → lesser representative power

　　▶ Breakage of weight structure → optimizer cannot be recovered

# ESPACE: OUR PROPOSAL TO DECOMPOSE ACTIVATIONS

$$
\begin{array}{ccccccccc}
\boxed{\begin{array}{c}\textbf{Weight}\\\textbf{Matrix}\\W^T\end{array}} & \times & \boxed{P} & \times & \boxed{\begin{array}{c}\textbf{Projection}\\\textbf{Matrix } P^T\end{array}} & \times & \boxed{\begin{array}{c}\textbf{Activation}\\\textbf{Tensor}\\X\end{array}} & = & \boxed{\begin{array}{c}\textbf{Output}\\\textbf{Tensor}\\Y\end{array}}\\
\text{N x K} & & \begin{array}{c}\underline{\text{Static}}\\\text{K x L}\end{array} & & \begin{array}{c}\underline{\text{Static}}\\\text{L x K}\end{array} & & \text{K x M} & & \text{N x M}
\end{array}
$$

- ▸ Rather than computing: $Y = W^T X$

- ▸ let's insert a small matrix $P$ and its transpose $P^T$

  - ▸ this *static and orthonormal* matrix will be related to the activation tensor $X$

  - ▸ It is likely that activation tensors have more redundancies as well (long sequences, repeated tokens, etc..)

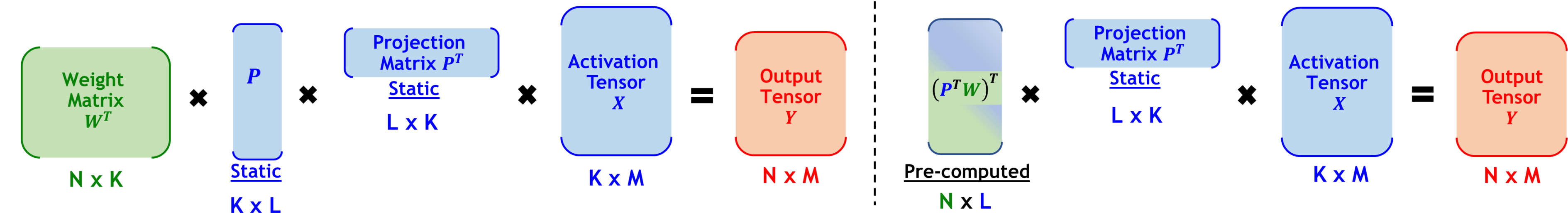- ▸ we will compute: $\widehat{Y} = W^T P P^T X$

# ESPACE: EIGEN STATIC PRINCIPAL ACTIVATION COMPONENT ESTIMATION

## LLM compression via activation dimensionality reduction

$$Y \approx W^T(PP^TX) = (P^TW)^T(P^TX)$$

**During Training:** ⬆️ trainable parameters: ⬆️ expressivity ✅          **During Inference:** ⬇️ model size: ⬇️ inference cost ✅

| Weight Matrix $W^T$ | × | $P$ | × | Projection Matrix $P^T$ Static L x K | × | Activation Tensor $X$ | = | Output Tensor $Y$ |
|---|---|---|---|---|---|---|---|---|
| N x K | | Static K x L | | | | K x M | | N x M |

| $(P^TW)^T$ | × | Projection Matrix $P^T$ Static L x K | × | Activation Tensor $X$ | = | Output Tensor $Y$ |
|---|---|---|---|---|---|---|
| Pre-computed N x L | | | | K x M | | N x M |

**K = input dimension, N= output dimension, M = batch ∗ sequence, L = intermediate dimension**

▸ We made an approximation $\widehat{X} = PP^TX \approx X \rightarrow$ will use continuous training to help model adapt to this approximation

▸ Activation projection does not interfere with weight learnability $\rightarrow$ the whole weight matrix $W^T$ is preserved

   ▸ we didn't reduce the number of learnable parameters and we can load the optimizer state

▸ At inference, this leads to compression when $L \ll K$

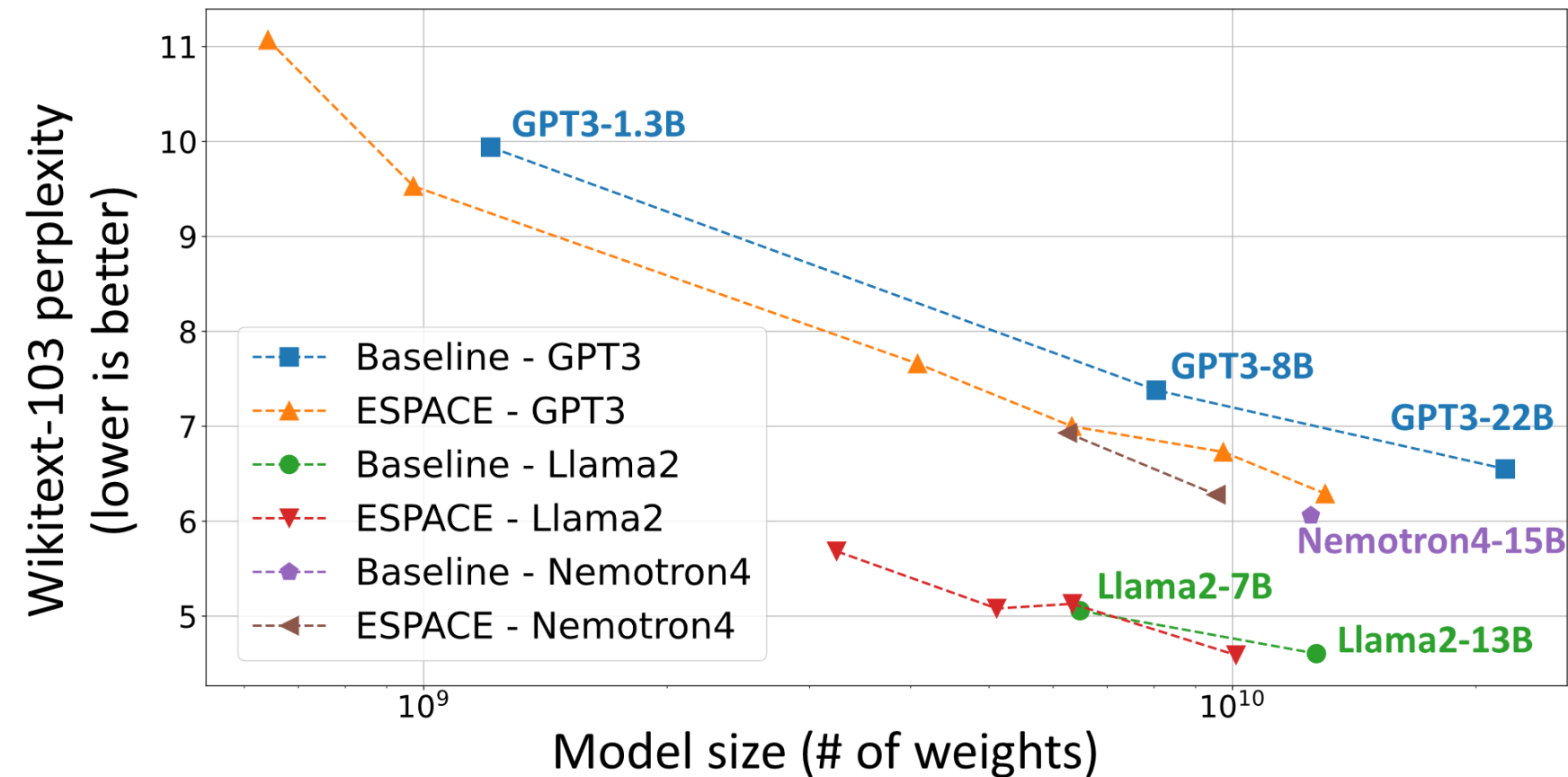   ▸ even during training, we end up computing less and observed a small end-to-end speed-up

# ESPACE: EIGEN STATIC PRINCIPAL ACTIVATION COMPONENT ESTIMATION

## Results and research roadmap

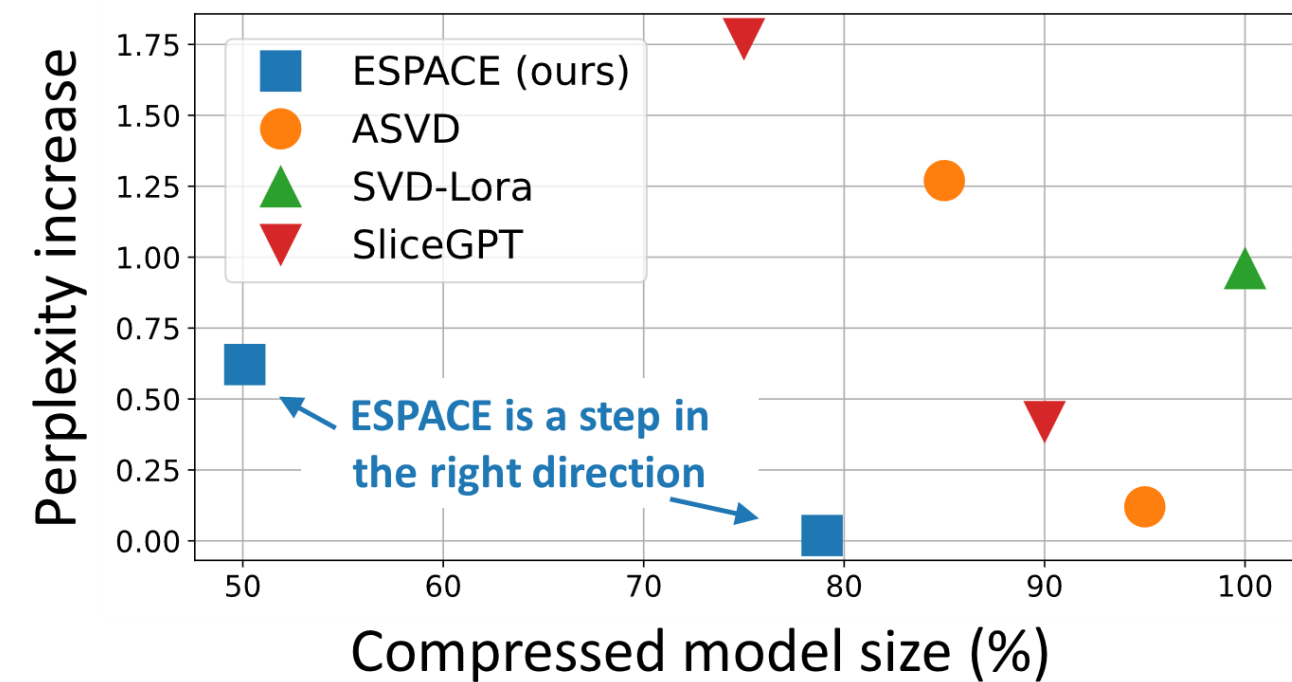**ESPACE strictly improves the pareto frontier of model size vs accuracy.**
**It also significantly improves on the SOTA of tensor decomposition method**

### LLM size vs accuracy



### Comparison to Prior Work
Tensor Decompositions of Weights

**Llama2-7B perplexity and compression**



▶ Benefits during inference:

  ▶ 20%-to-50% compression in # parameters

  ▶ less weight*activation math: 35%-to-45% reduction in GEMM latency

▶ Future work and research roadmap

  ▶ Combining ESPACE with other compression techniques, accelerate pre-training, optimize HW for back-to-back GEMMs

**NVIDIA.**