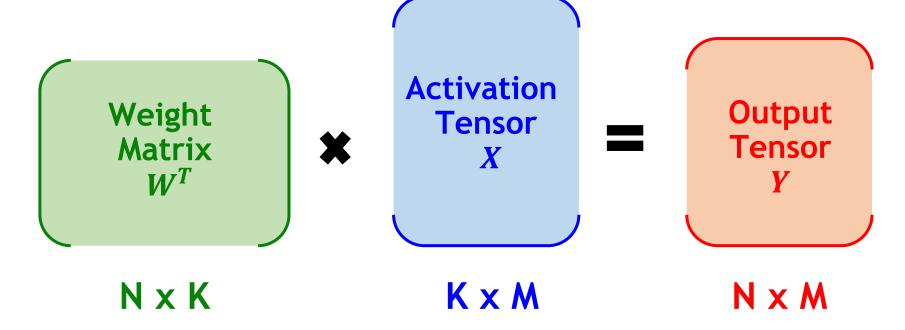
General Matrix Multiplication (GEMM)

$$\mathbf{Y} = \mathbf{W}^T \mathbf{X}$$

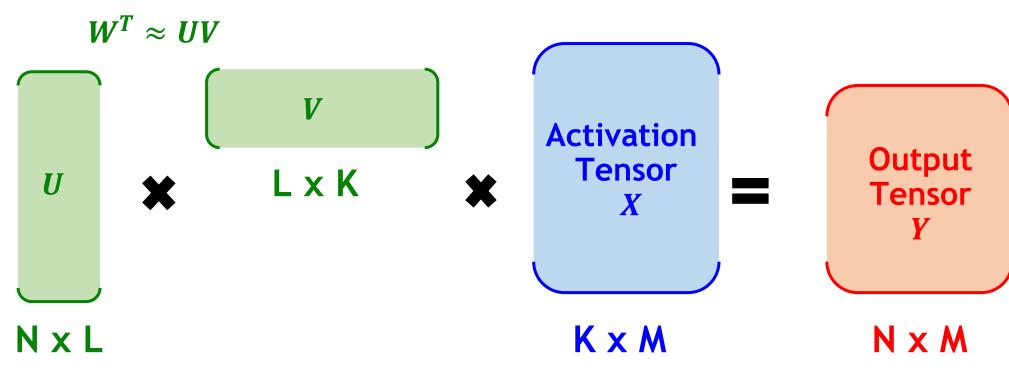
- trainable parameters: trainable parameters:
- model size: inference cost X



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

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

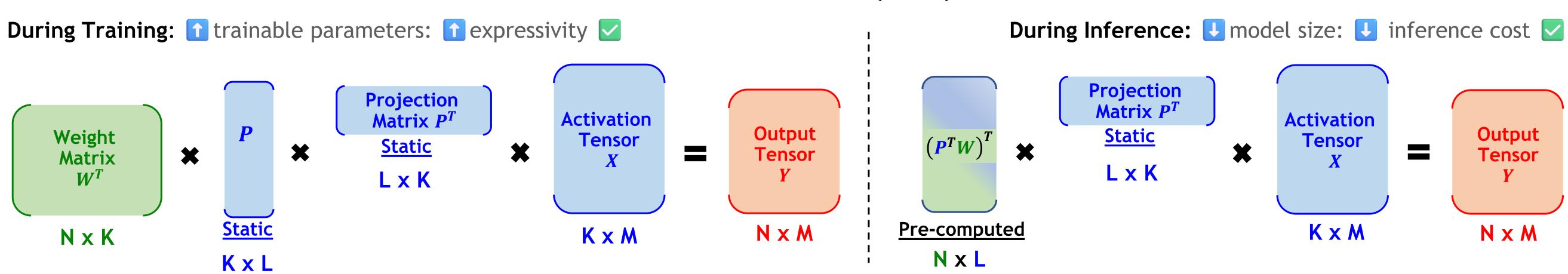
 $Y \approx UVX$



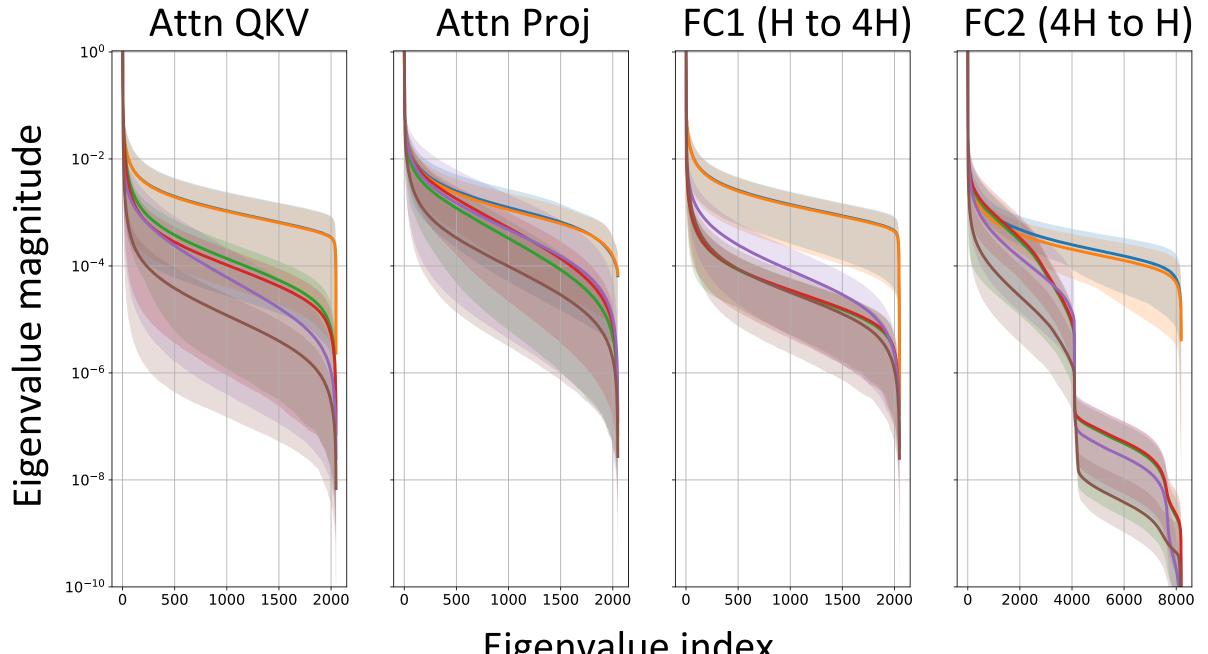
K = input dimension, N= output dimension, M = batch * sequence L = intermediate dimension

ESPACE: Activation Dimensionality Reduction Using Static Projections

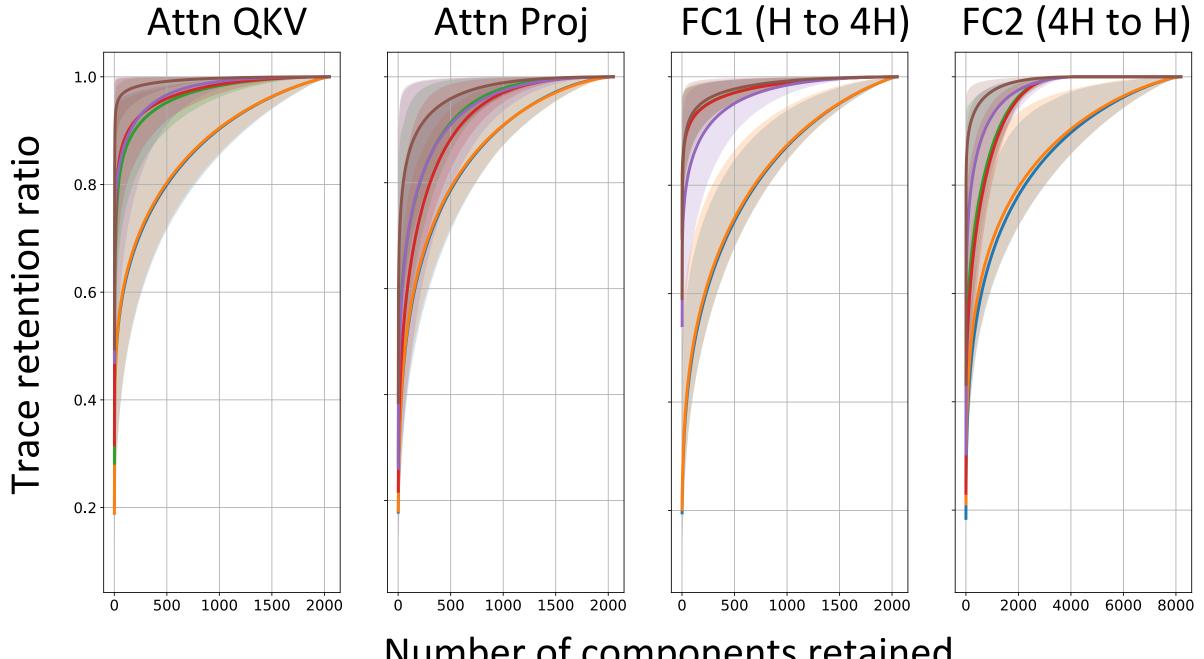
$$Y \approx W^T(PP^TX) = (P^TW)^T(P^TX)$$



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



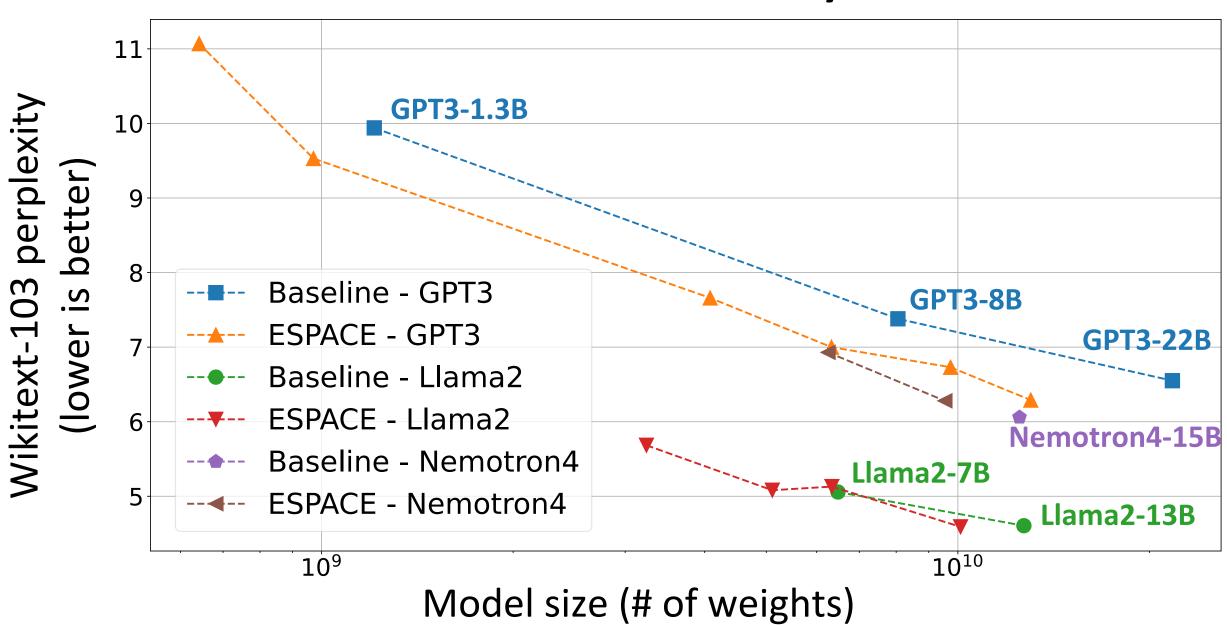
Eigenvalue index



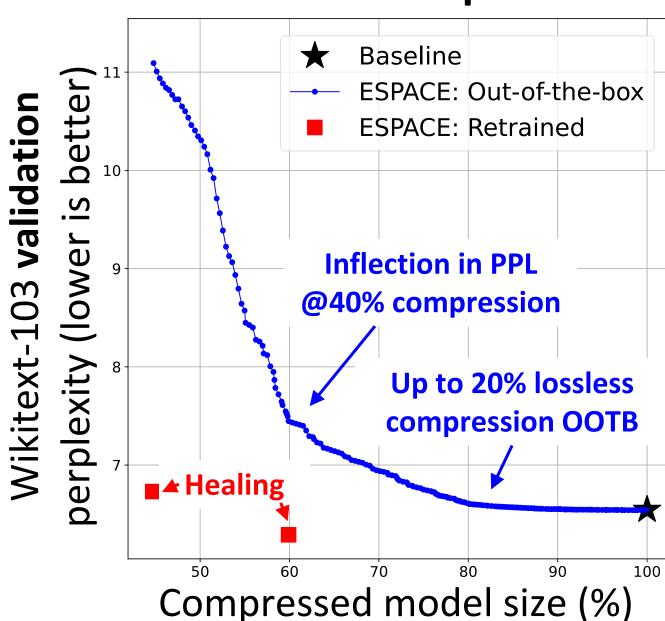
Number of components retained

— GO-MSE — GO-MSE w/ L_2 -normalization — NL-MSE — NL-MSE with L_2 -normalization - MSE - NMSE

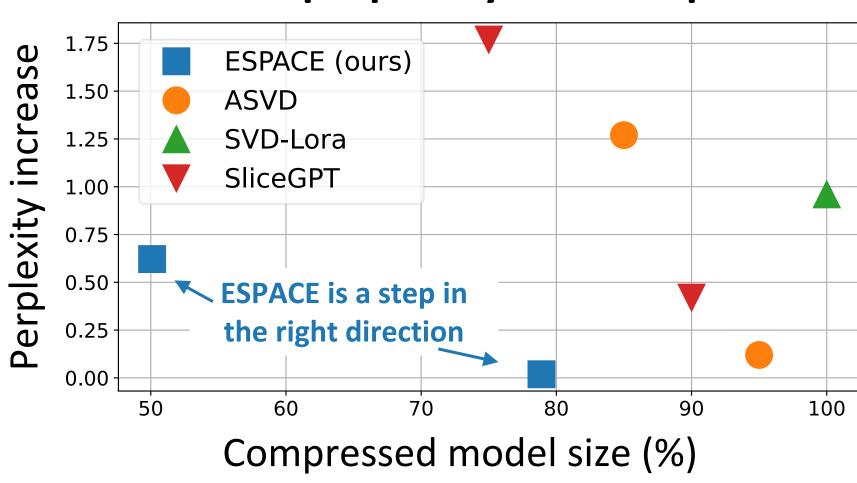
LLM size vs accuracy



GPT3-22B Compression



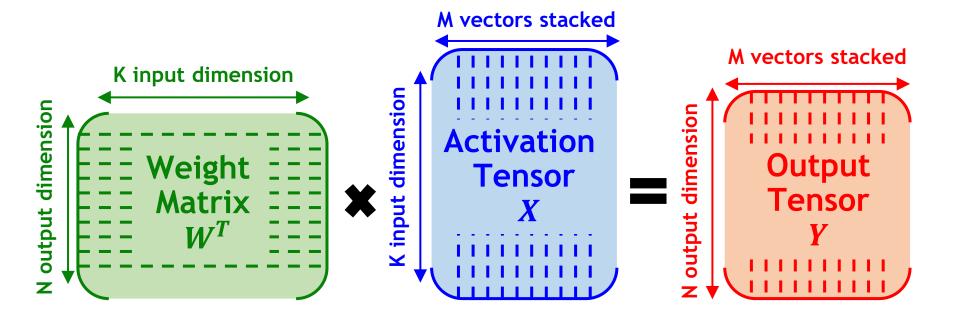
Llama2-7B perplexity and compression

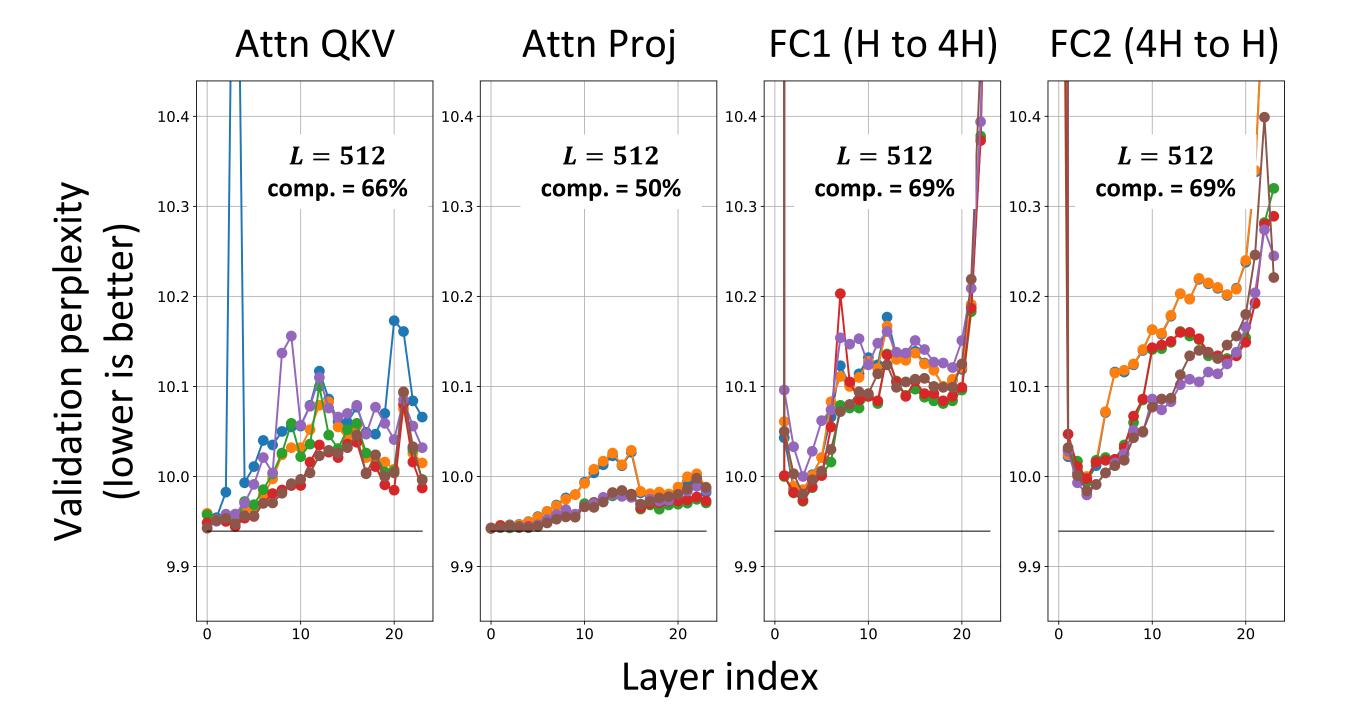


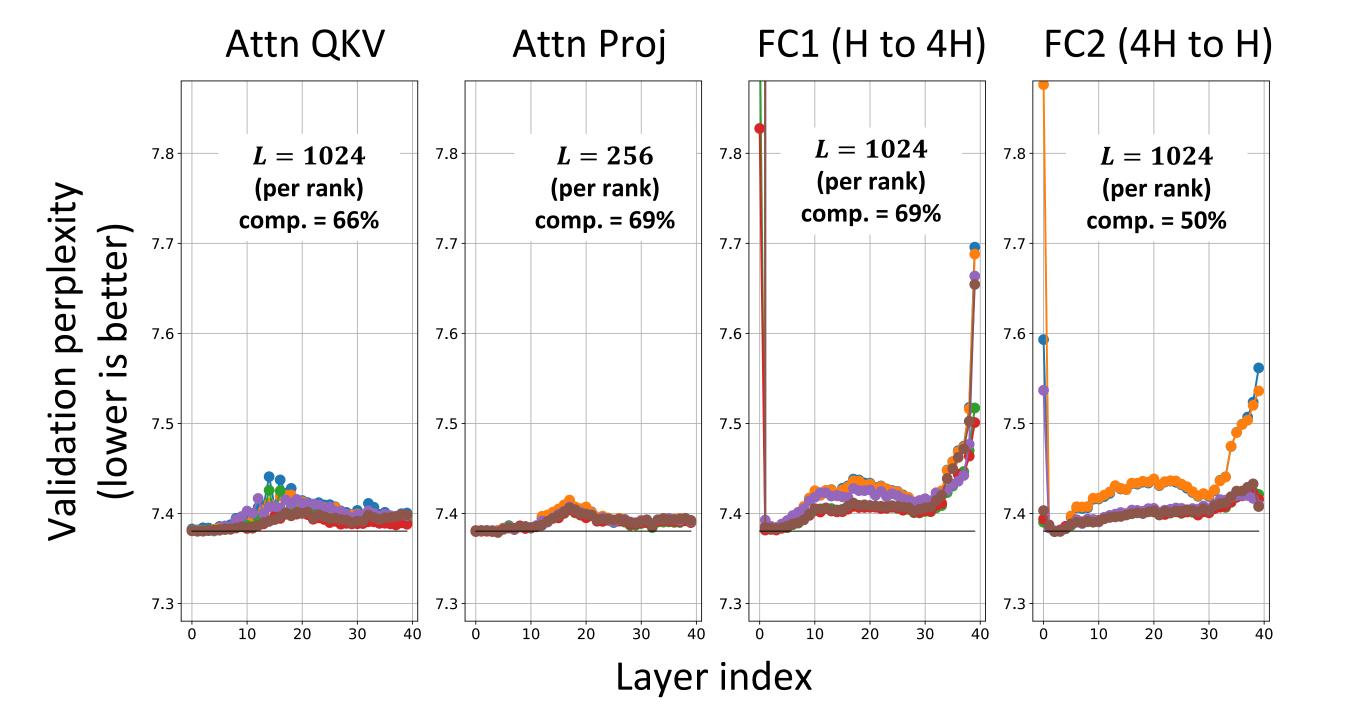
General Matrix Multiplication (GEMM)

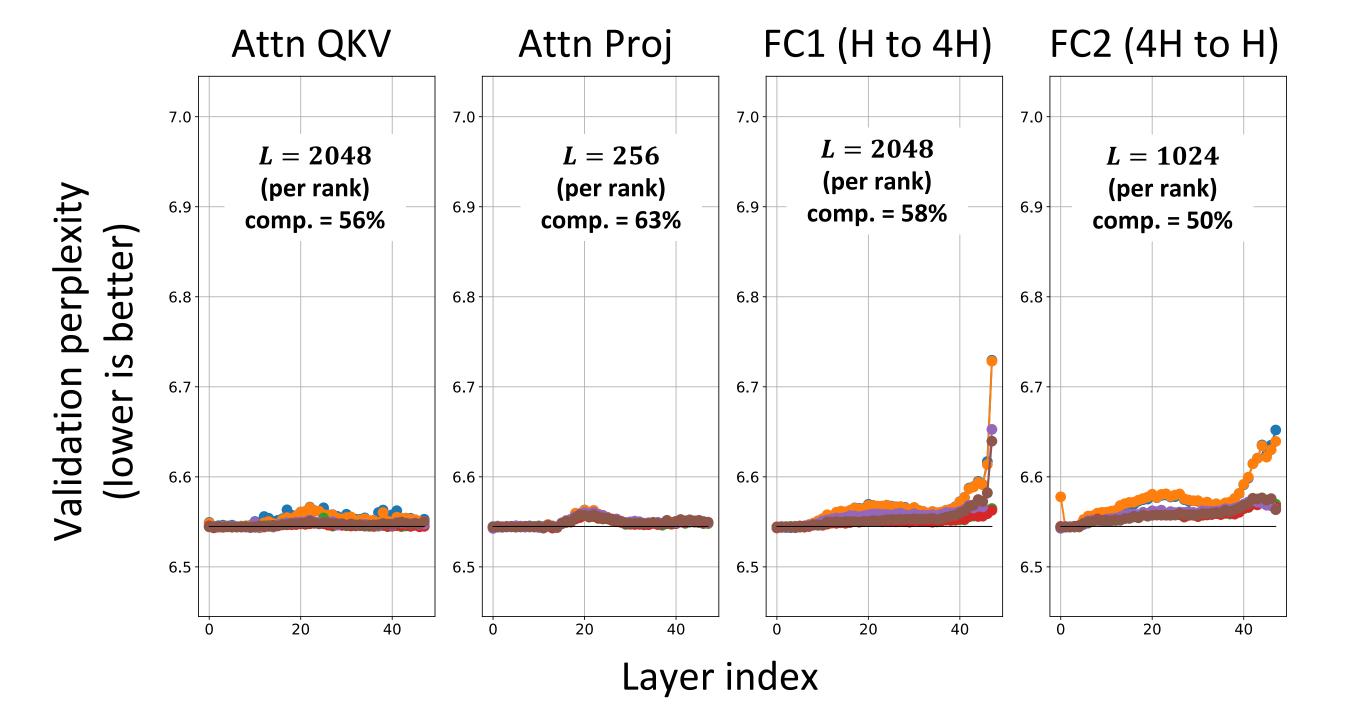
$$\mathbf{Y} = \mathbf{W}^T \mathbf{X}$$

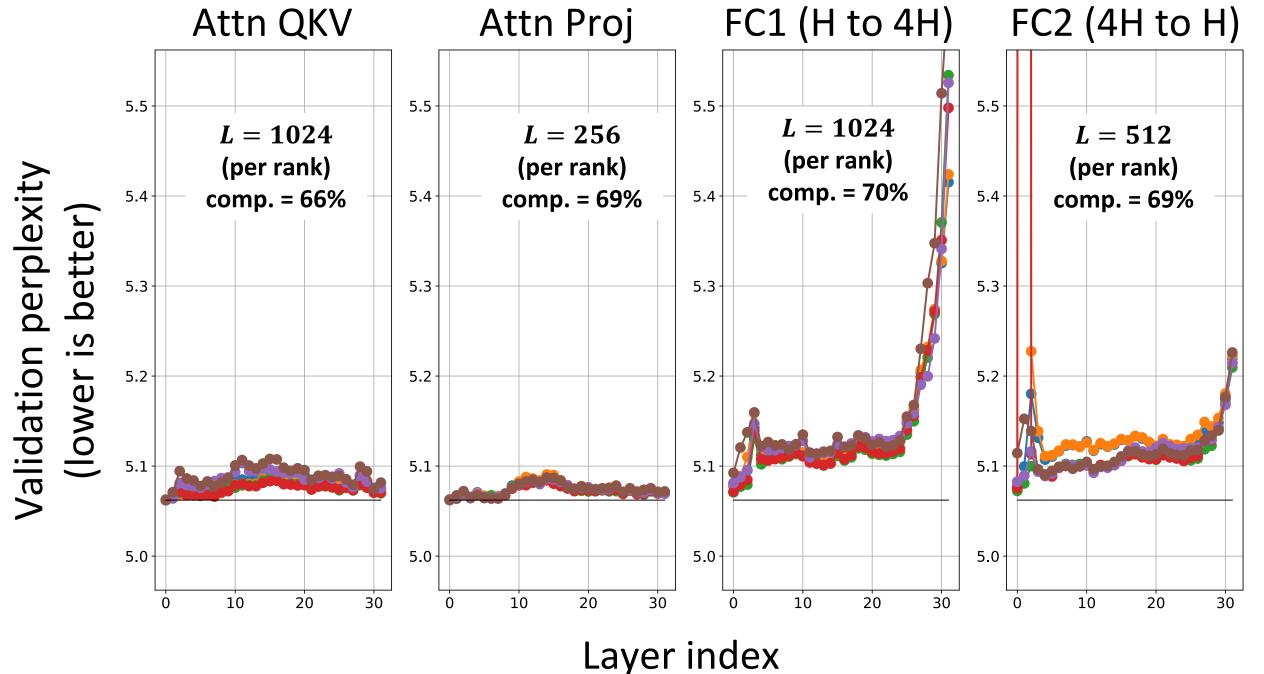
- trainable parameters: expressivity model size: inference cost

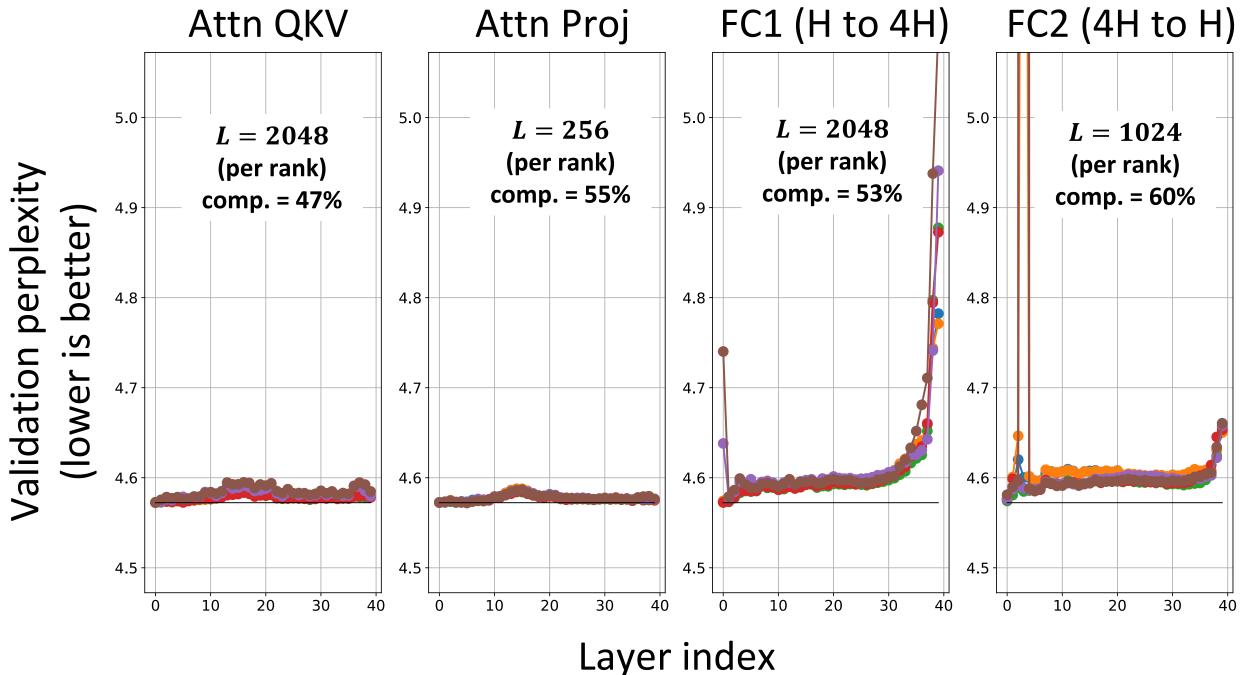


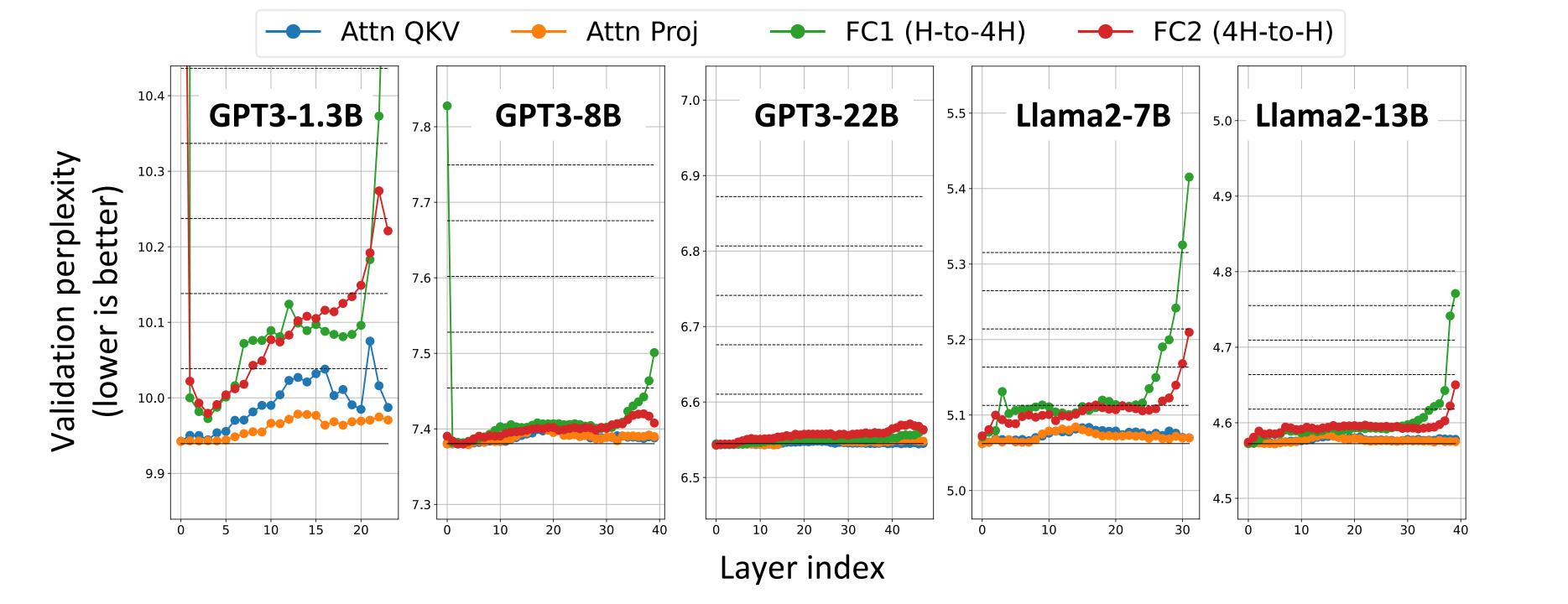


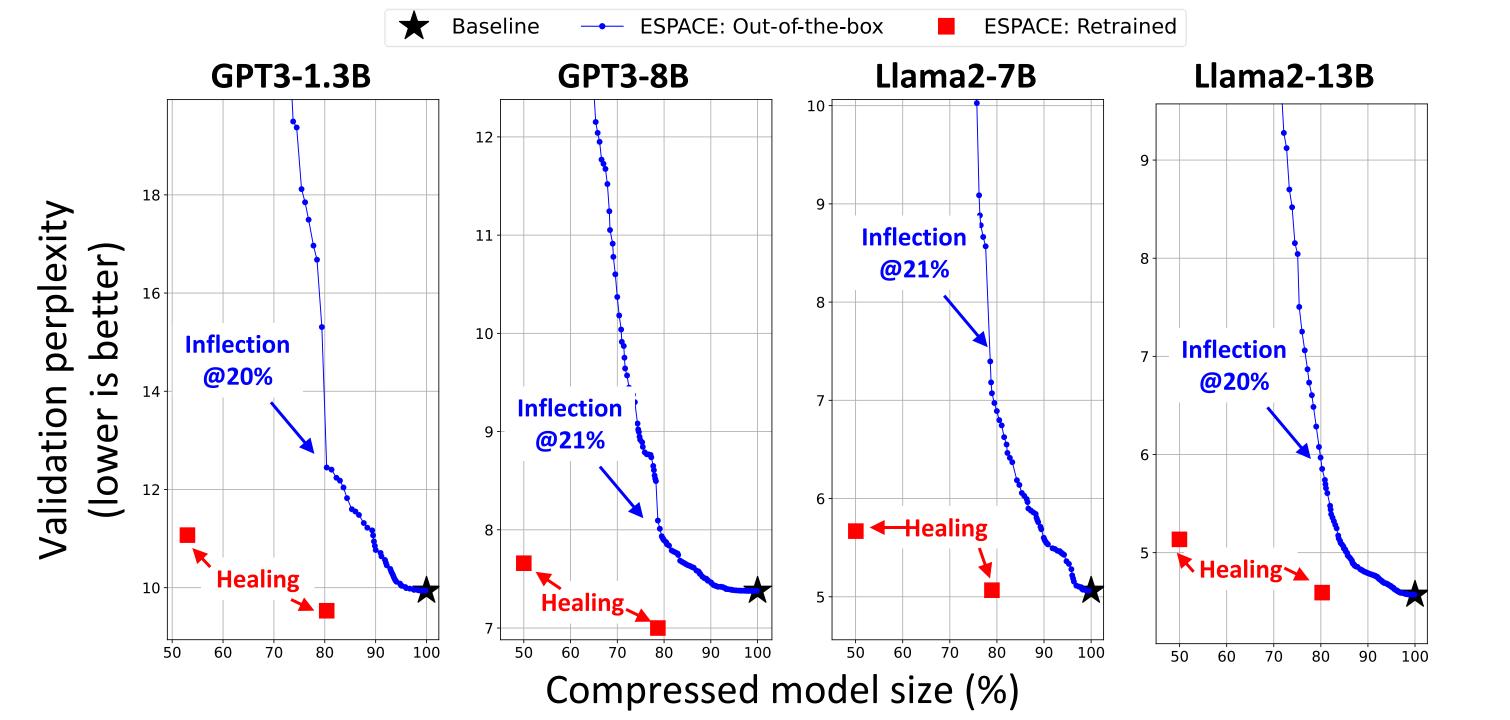




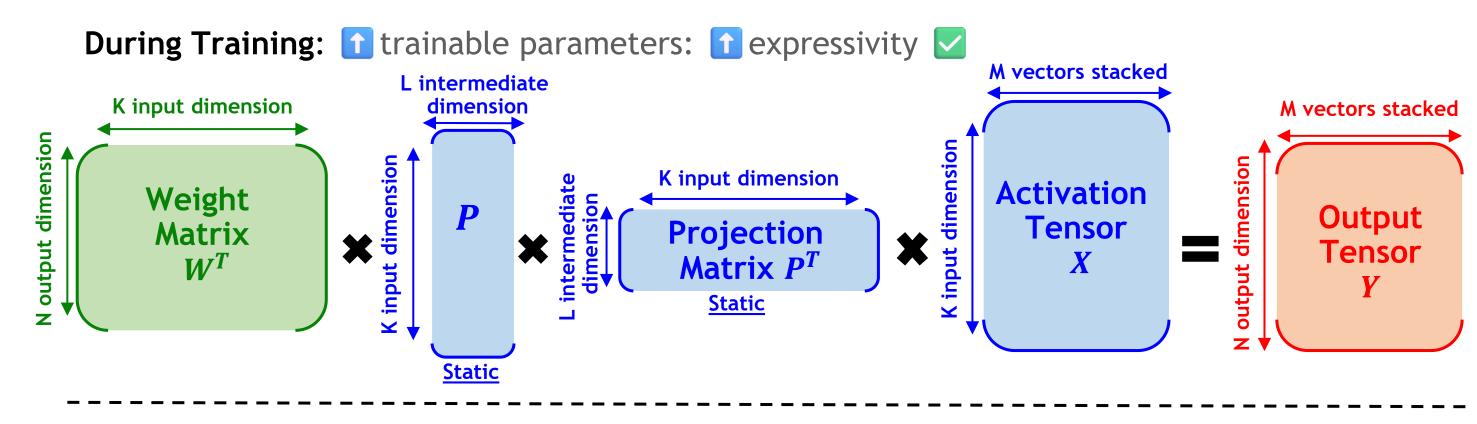








Activation projection does not interfere with weight learnability in training, yet leads to model compression at inference

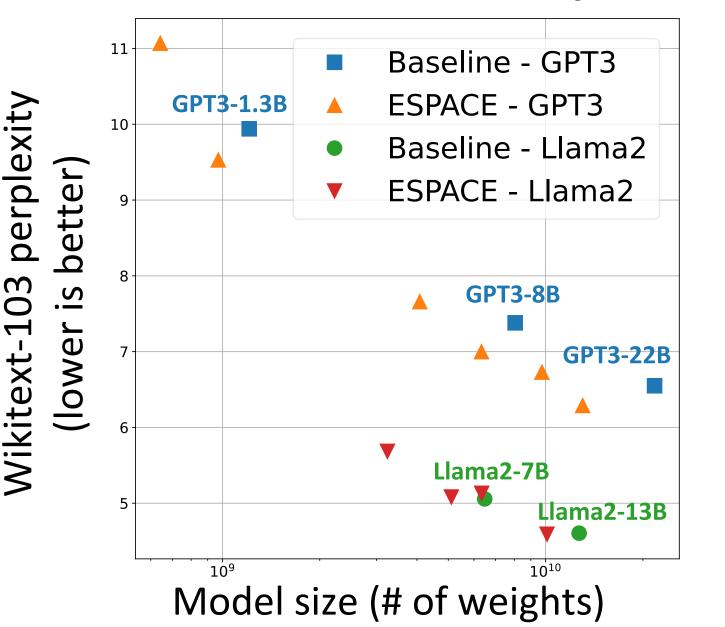


During Inference:

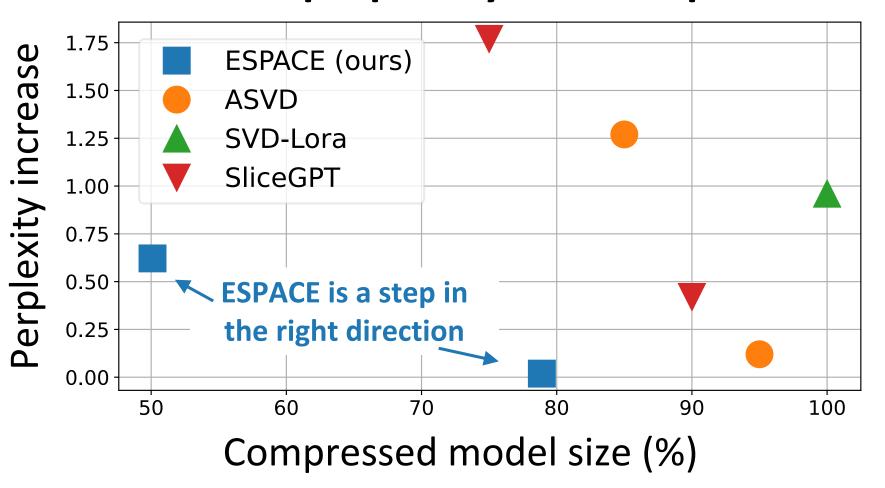
☐ model size: ☐ inference cost ☐

☐ M vectors stacked dimension

LLM size vs accuracy



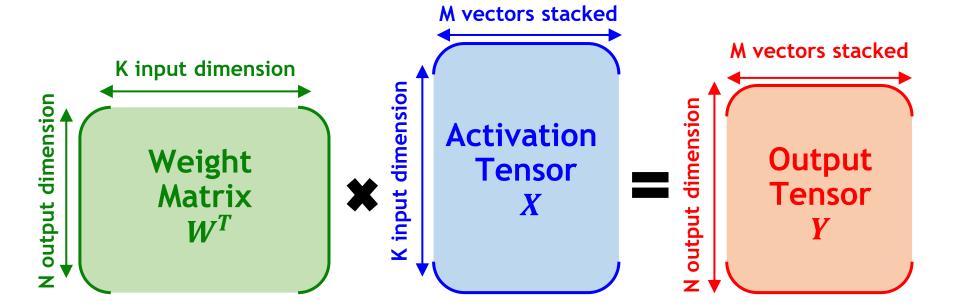
Llama2-7B perplexity and compression



General Matrix Multiplication (GEMM)

$$\mathbf{Y} = \mathbf{W}^T \mathbf{X}$$

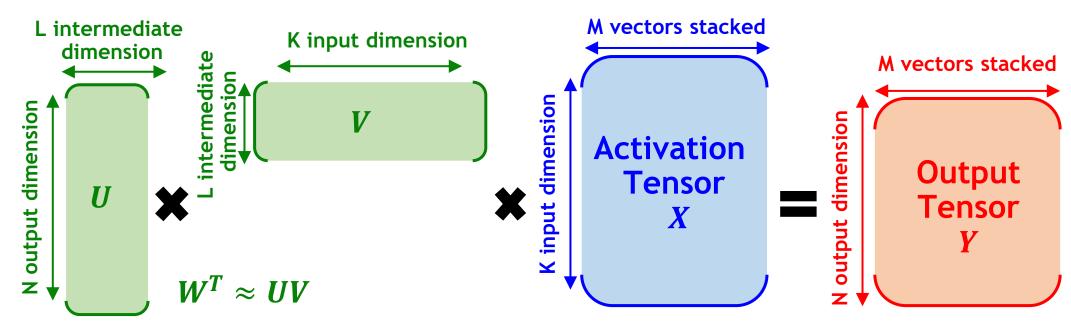
- trainable parameters: texpressivity model size: tinference cost



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

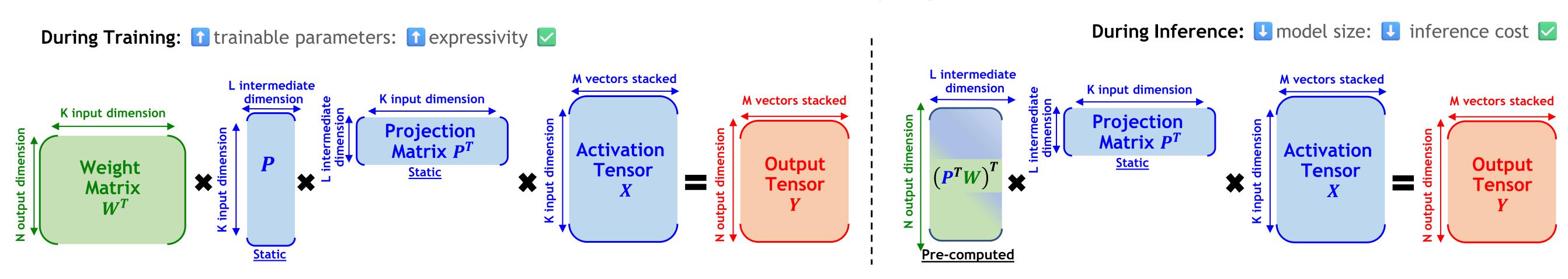
 $Y \approx UVX$

- U trainable parameters: U expressivity XU model size: U inference cost ✓



ESPACE: Activation Dimensionality Reduction Using Static Projections

$$\mathbf{Y} \approx \mathbf{W}^T (\mathbf{P} \mathbf{P}^T \mathbf{X}) = (\mathbf{P}^T \mathbf{W})^T (\mathbf{P}^T \mathbf{X})$$



LLM size vs accuracy

