

PACEAttention: Principled and Adaptive Feature Compression-Expansion Grounded in the Geometry of MCR²



Xiaojie Yu¹ Haibo Zhang^{2,*} Jeremiah D. Deng^{1,*} Lizhi Peng^{3,4}
¹University of Otago, ²University of New South Wales, ³Quancheng Laboratory, ⁴University of Jinan
¹ Contact: haibo.zhang@unsw.edu.au, jeremiah.deng@otago.ac.nz



1. Background

The maximal coding rate reduction (MCR²) objective provides a principled framework for learning low-dimensional subspace representations Z by promoting intra-class compactness (i.e., feature compression) and inter-class separation (i.e., feature expansion). White-box deep models can be derived by unrolling optimization steps of this objective implied by the MCR² gradient, as the approximation of the gradient may deviate from the intended compression objective.

Our goal aims to optimize the following MCR²-like objective, where $U_{[K]}$ are K trainable subspace matrices.

$$\Delta R(Z, \Pi) = R(Z) - R_c(Z | U_{[K]})$$

$$= \underbrace{\frac{1}{2} \log \det \left(I + \frac{d}{n\epsilon^2} ZZ^T \right)}_{\text{Expansion term}} - \underbrace{\frac{1}{2} \sum_{k=1}^K \log \det \left(I + \frac{p}{n\epsilon^2} (U_k^T Z)^T (U_k^T Z) \right)}_{\text{Compression term}}$$

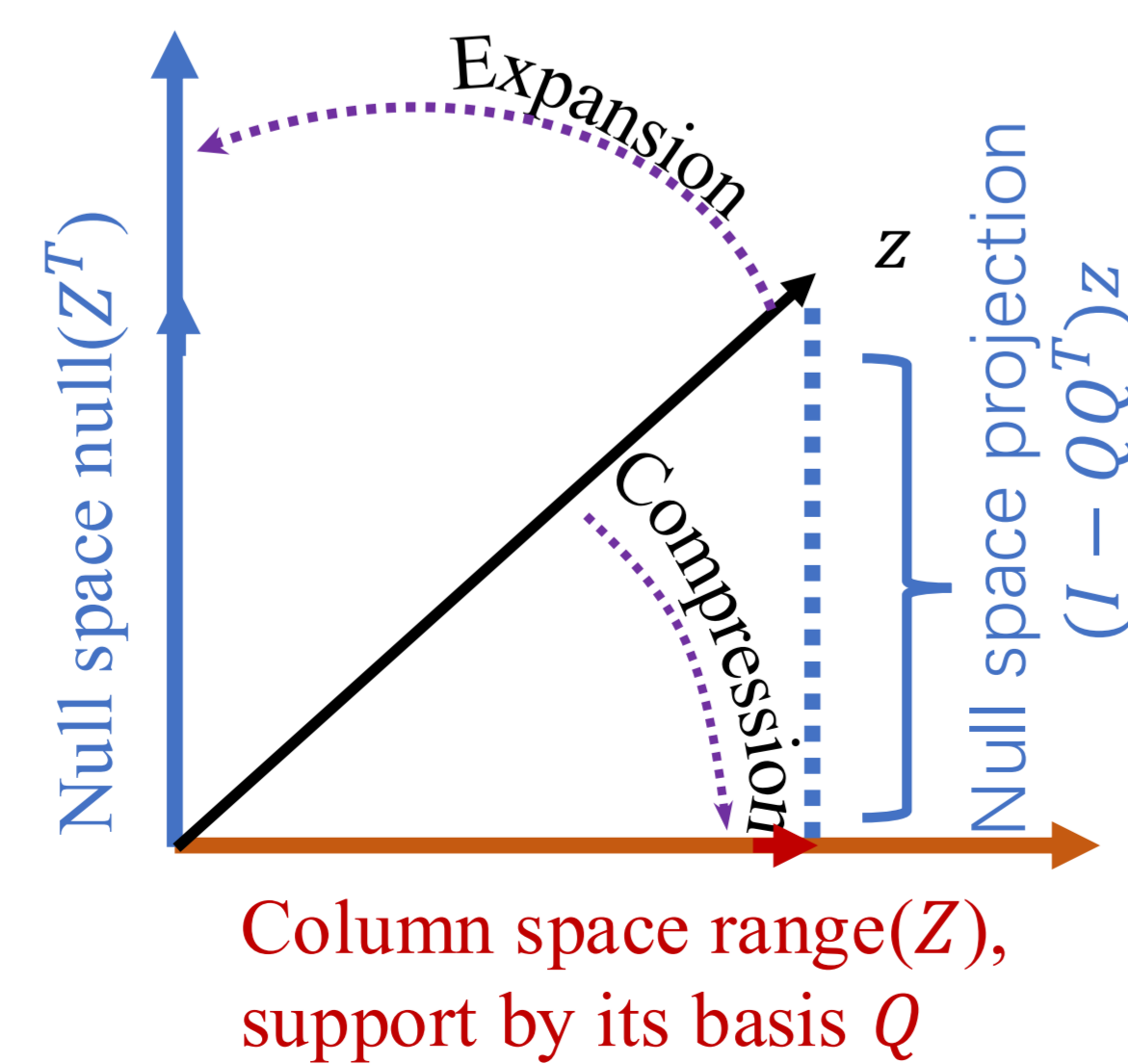
2. Key Contributions

- ▶ Principled feature updates derived from the geometry of MCR² gradient.
- ▶ Using randomization to efficiently approximate subspace structures.
- ▶ Linear complexity.
- ▶ Learnable α, β balance expansion and compression per layer.
- ▶ Heads learn fine-grained object-region memberships.

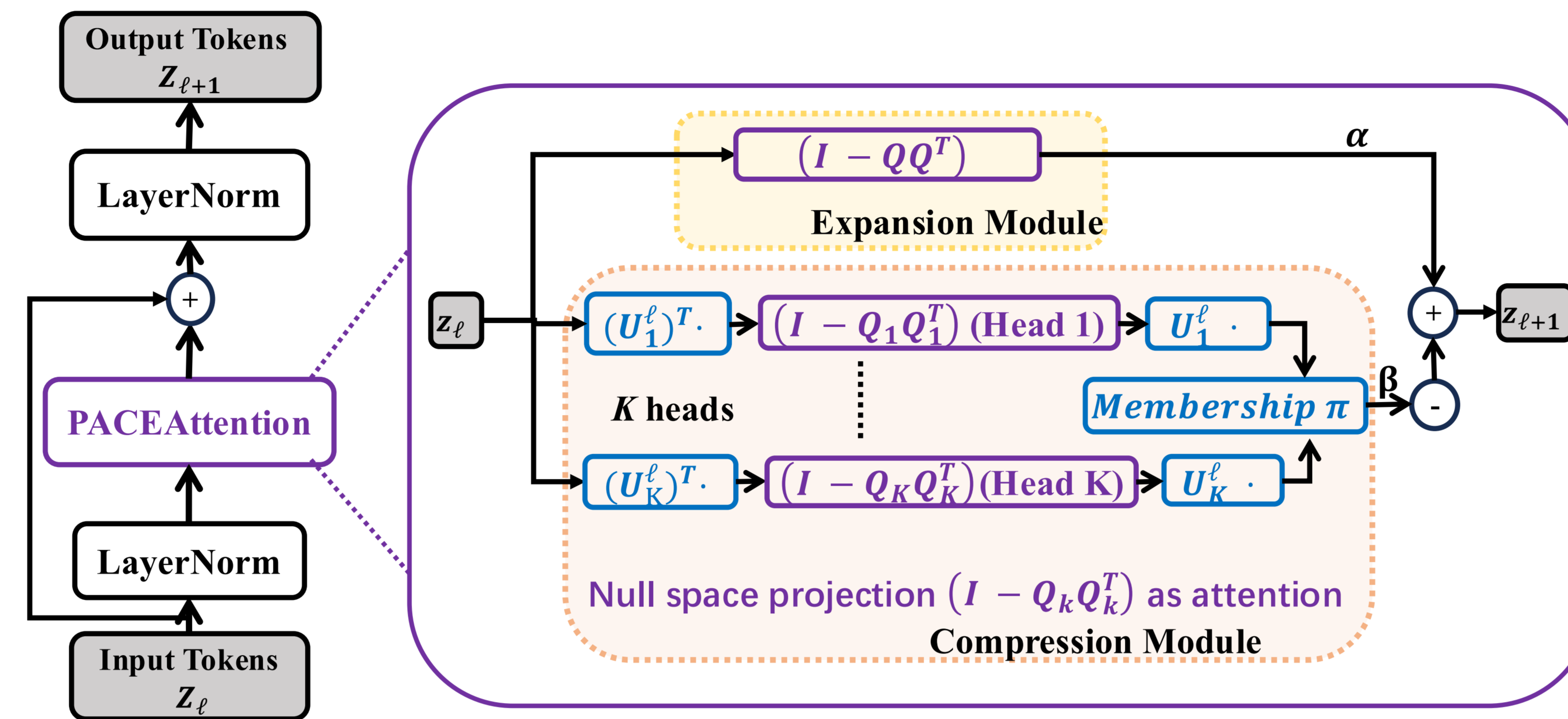
3. Geometric Intuition

We exploit the geometric interpretation of the gradient of $\log \det(I + ZZ^T)$, a term included in MCR²-like objective. To increase or decrease the value of $\log \det(I + ZZ^T)$, we can enlarge or reduce the dimensionality of column space $\text{range}(Z)$. Therefore, let Q be the orthonormal basis of $\text{range}(Z)$. We arrive at two core ideas:

- ▶ For **feature expansion**, adding the null space projection $(I - QQ^T)z$ moves a token z toward the null space, hence increasing the column space dimensionality and promoting separability.
- ▶ For **feature compression**, subtracting null space projection moves the token z toward the column space, hence decreasing the column space dimensionality.



4. PACEAttention Layer



$$z_{l+1} = z_l + \underbrace{\alpha(I - QQ^T)z_l}_{\text{Expansion}} - \underbrace{\beta \sum_{k=1}^K \pi_k U_k (I - Q_k Q_k^T) U_k^T z_l}_{\text{Compression}}$$

Expansion: Increase global coding rate $R(Z)$ by expanding toward null-space directions.

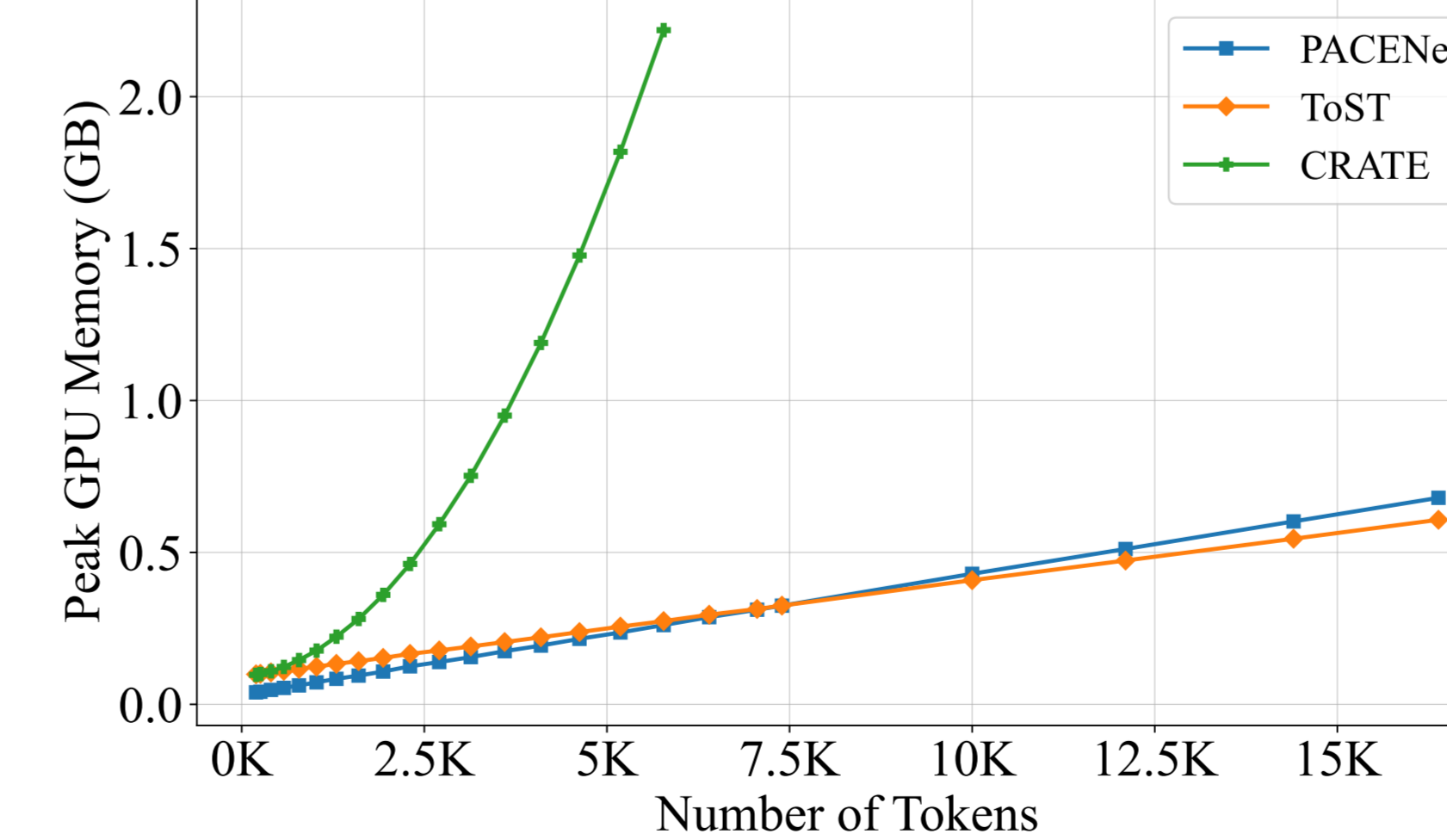
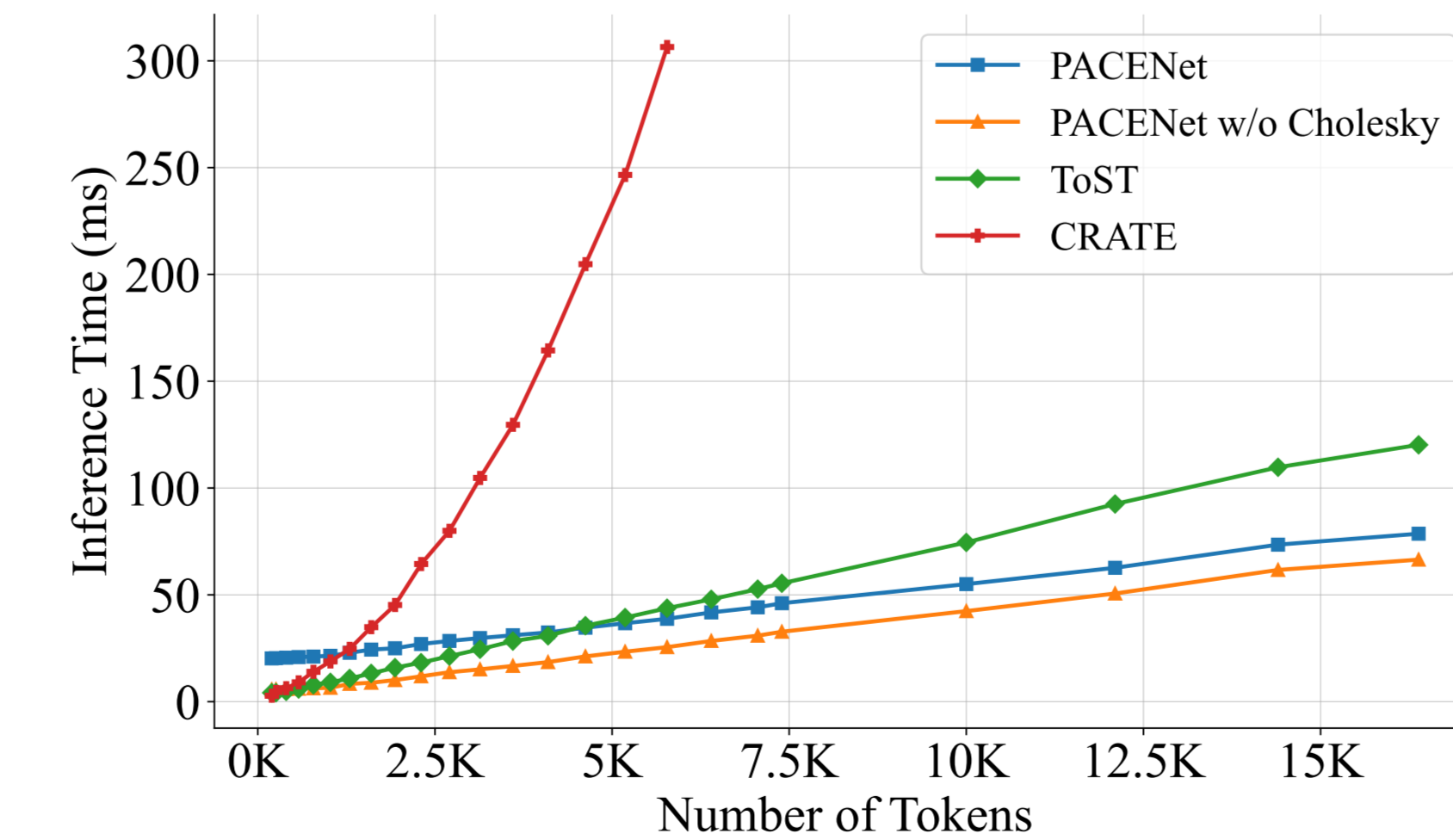
Compression: Reduce $R_c(Z | U_{[K]})$ by subtracting group-wise null-space components.

Self-adaptive dynamics: α controls expansion; β controls compression.

5. Efficient Subspace Approximation based on Randomization

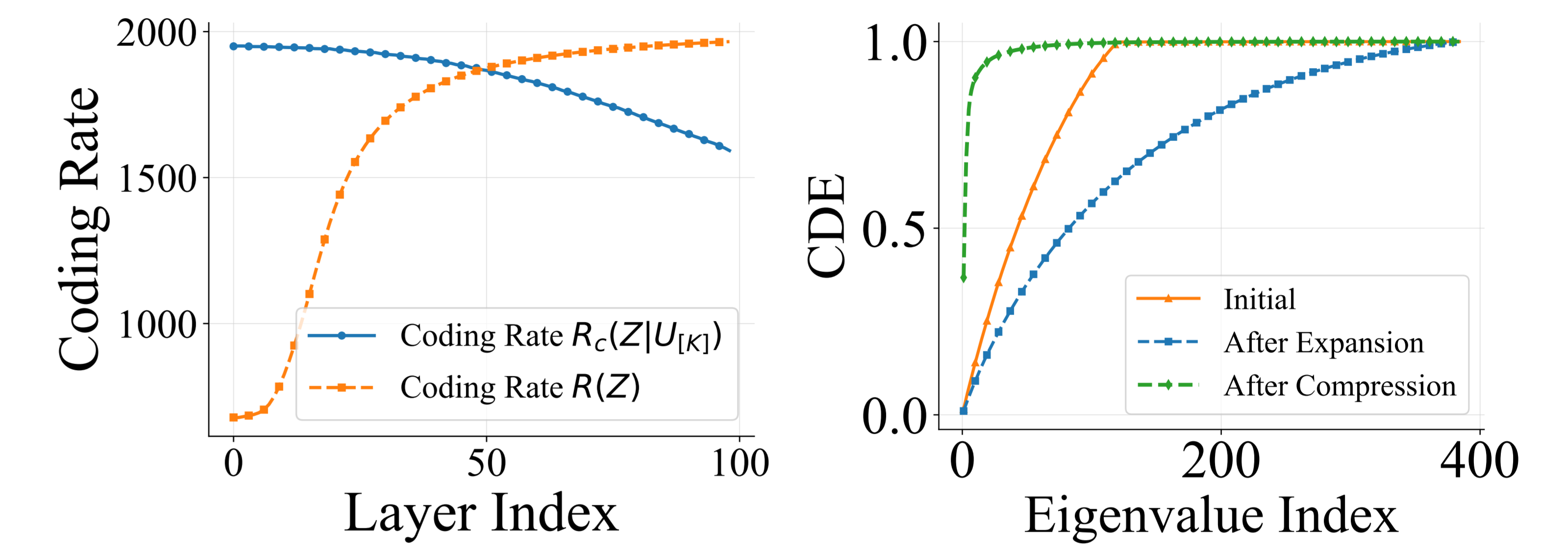
$$Y = Z\Omega, \quad Q = \text{orth}(Y)$$

Random matrix Ω captures the column space.
Cholesky decomposition extract orthonormal basis efficiently.



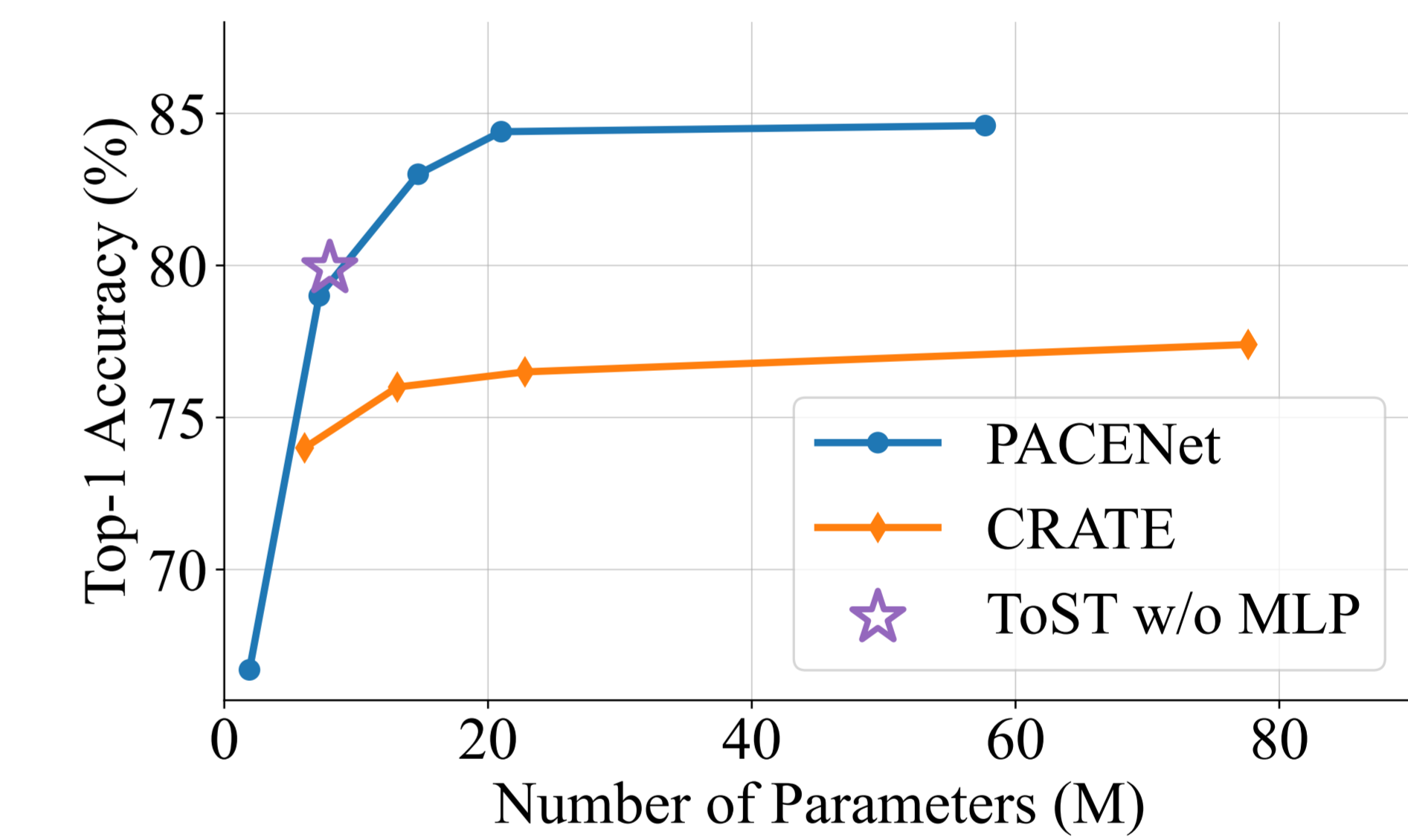
Scaling. Linear complexity in the number of tokens.

6. Verification on Toy Data

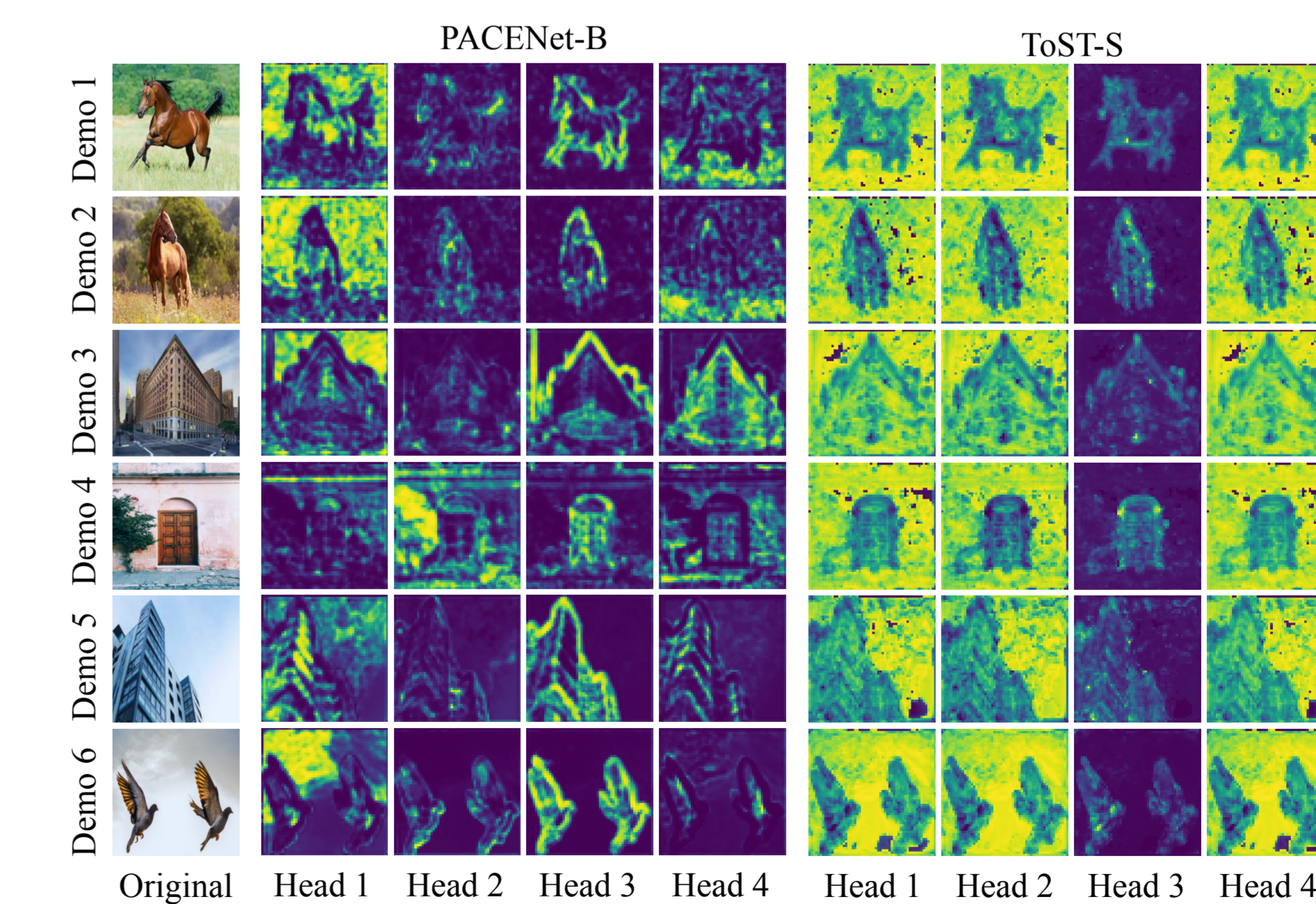


- ▶ Expansion increases $R(Z)$ layer by layer.
- ▶ Compression decreases $R_c(Z | U_{[K]})$.
- ▶ Eigen-spectrum confirms effective rank expansion/compression.

7. Top-1 on ImageNet-Real



8. Interpretability



Different heads attend to distinct object parts, enabling fine-grained structural modeling.