

Skip-Attention: Improving Vision Transformers by Paying Less Attention



Shashanka Venkataramanan*, Amir Ghodrati*, Yuki M. Asano, Fatih Porikli, Amirhossein Habibian

Improving efficiency in ViTs

- Transformers are state-of-the-art in most image and video tasks.
- ViTs are **computationally expensive**.
 - Compute and memory grows **quadratically** ($N \times N$). 14x14 for image classification to 128x128 = 16K for image denoising.
- How to improve the efficiency of vision transformers?

We leverage the **high correlation** across MSA blocks to improve **efficiency** of ViTs

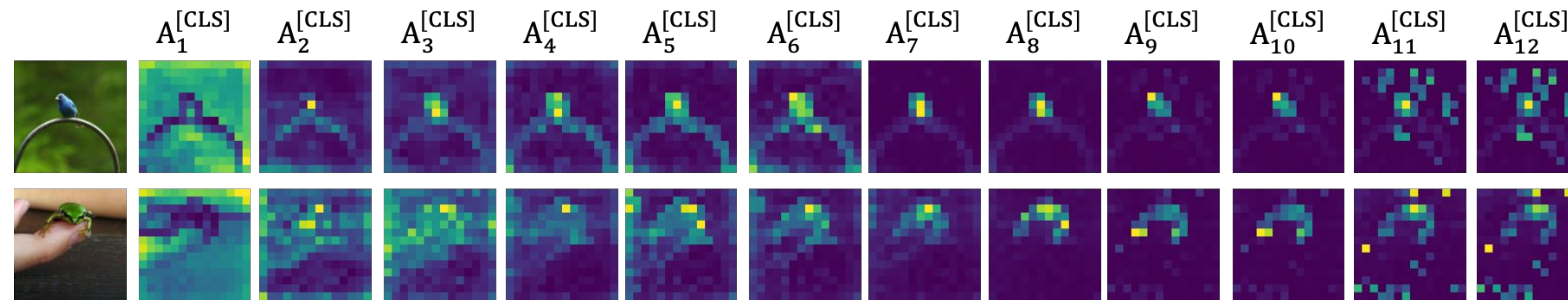
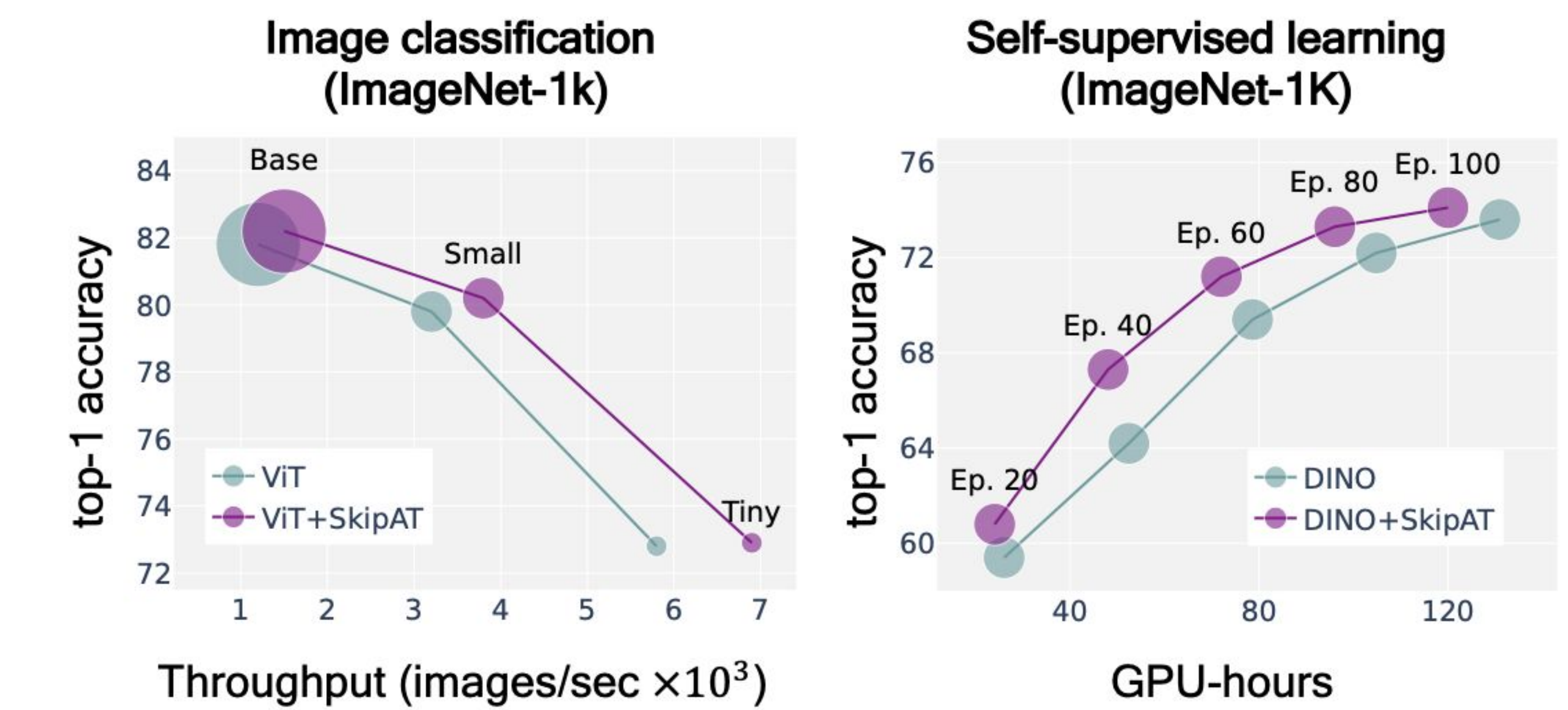
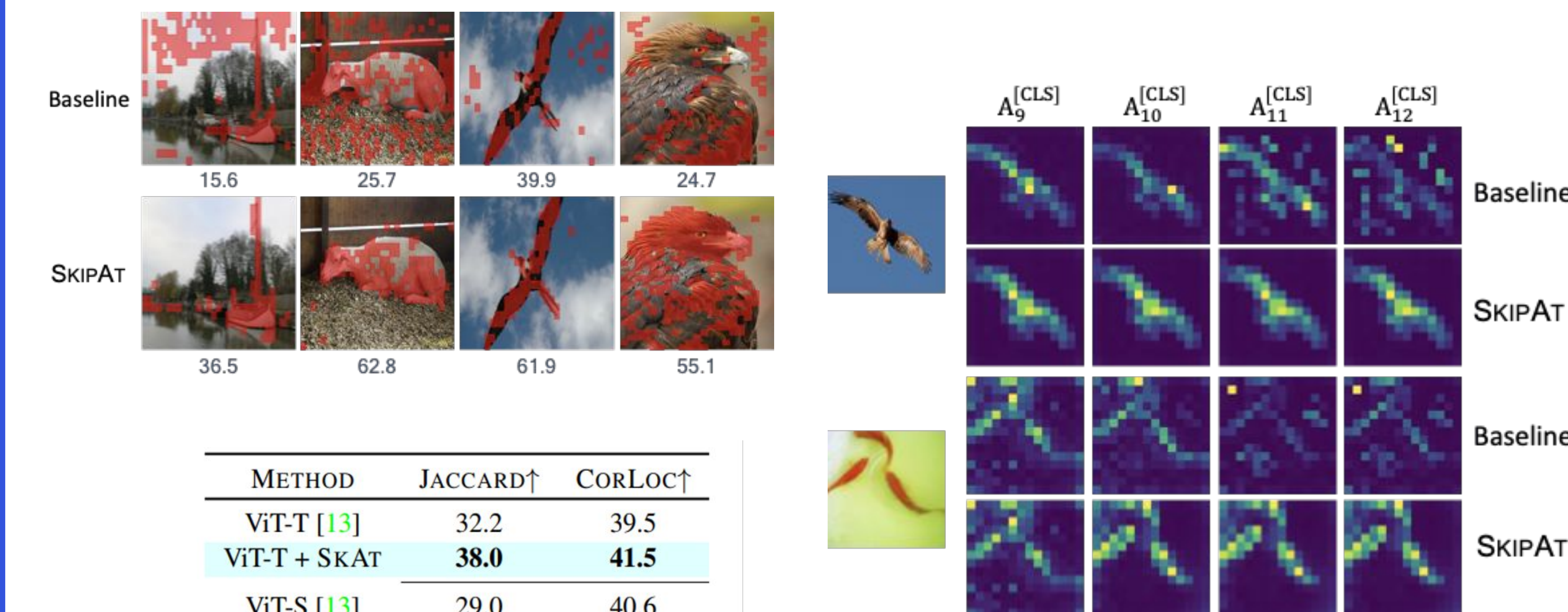


Image classification

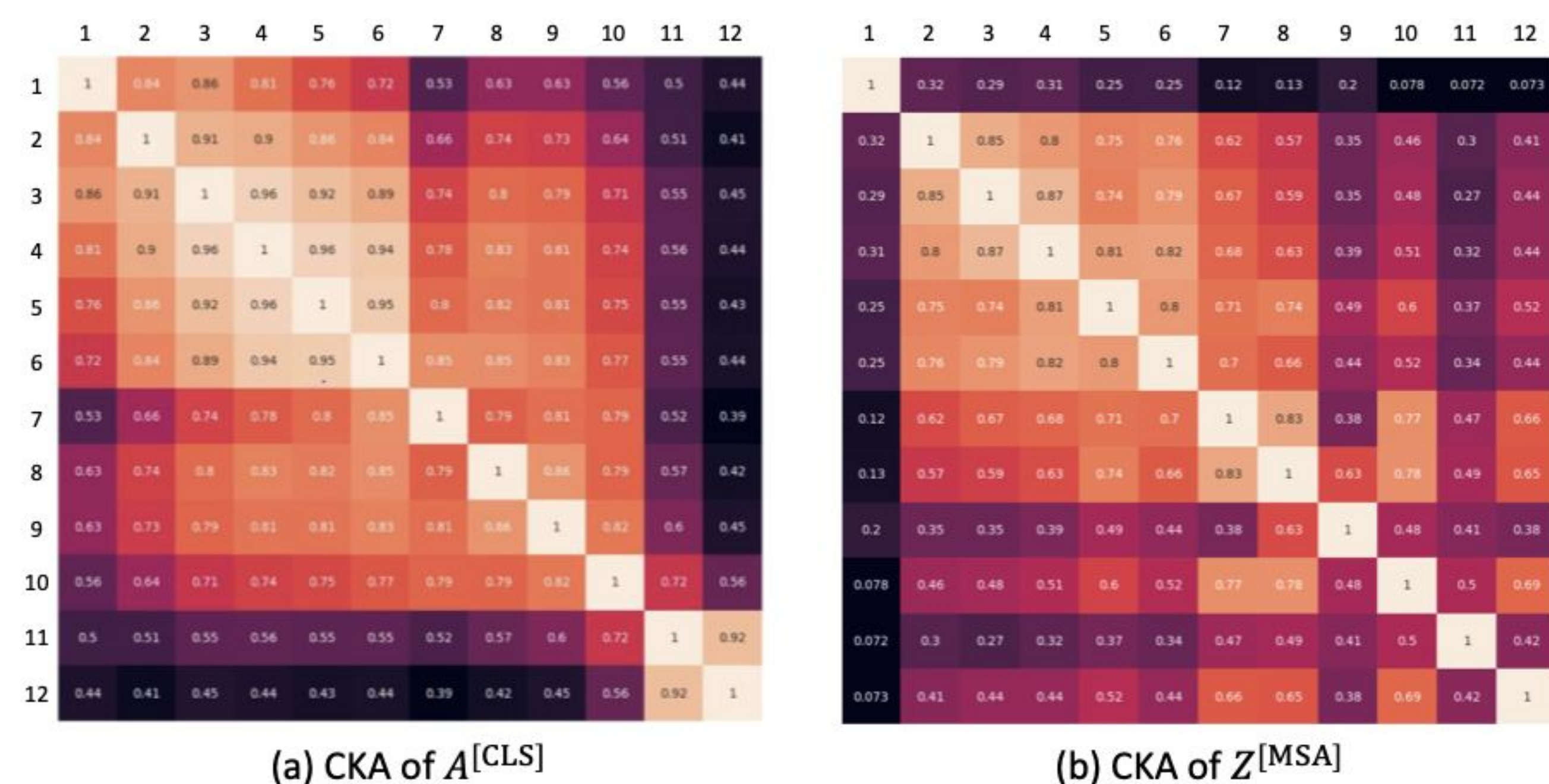


Unsupervised Object Discovery

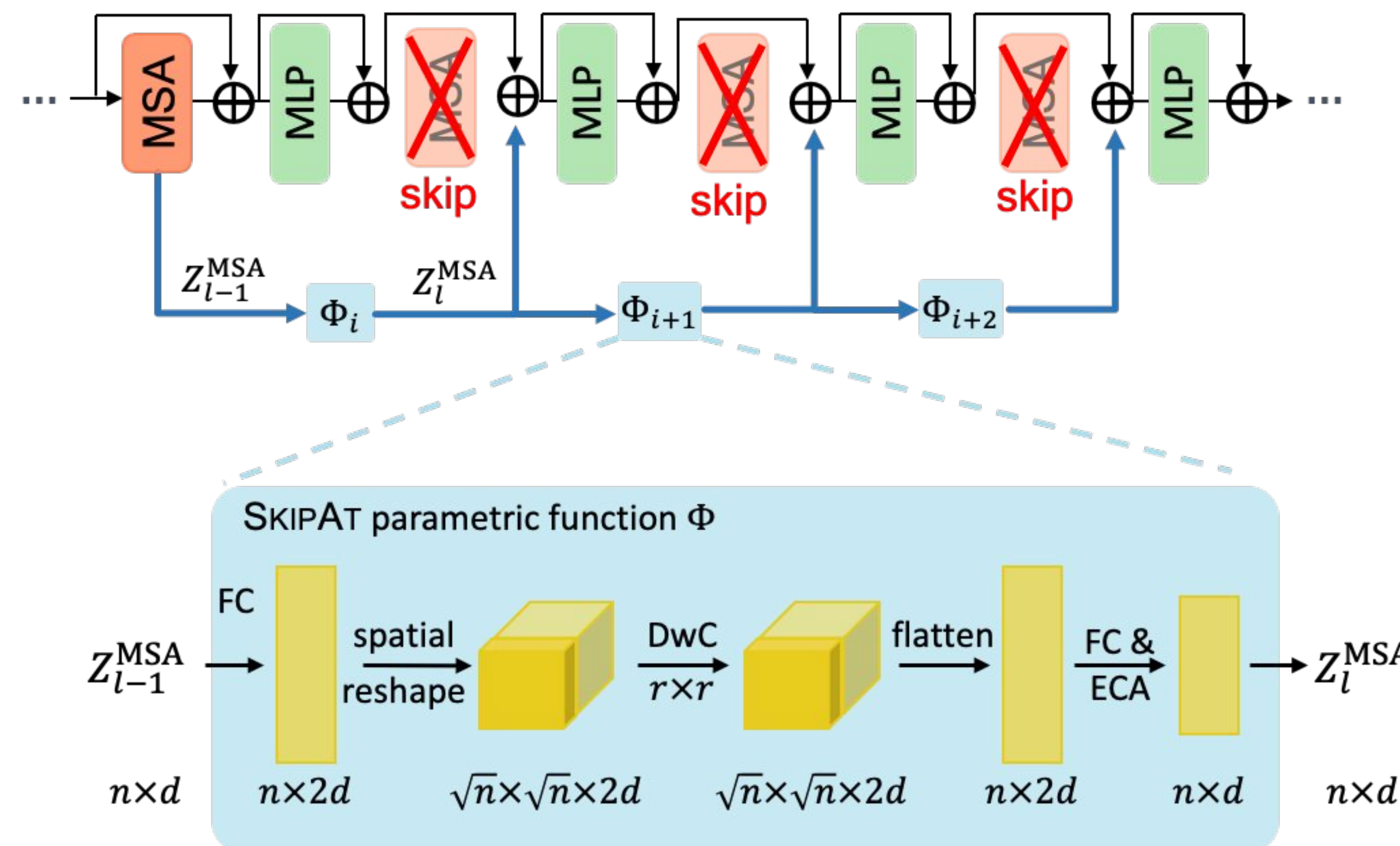


CKA analysis

- High correlation** between attention matrices across layers.
- Same for output of MSA blocks (features Z).



Approximating MSA with SkipAT parametric function



$$Z_l \leftarrow Z_l^{MSA} + Z_{l-1}$$

$$Z_l \leftarrow \text{MLP}(Z_l) + Z_l$$

$$Z_l \leftarrow \Phi(Z_{l-1}^{MSA}) + Z_{l-1}$$

$$Z_l \leftarrow \text{MLP}(Z_l) + Z_l$$

Additional Downstream Tasks

