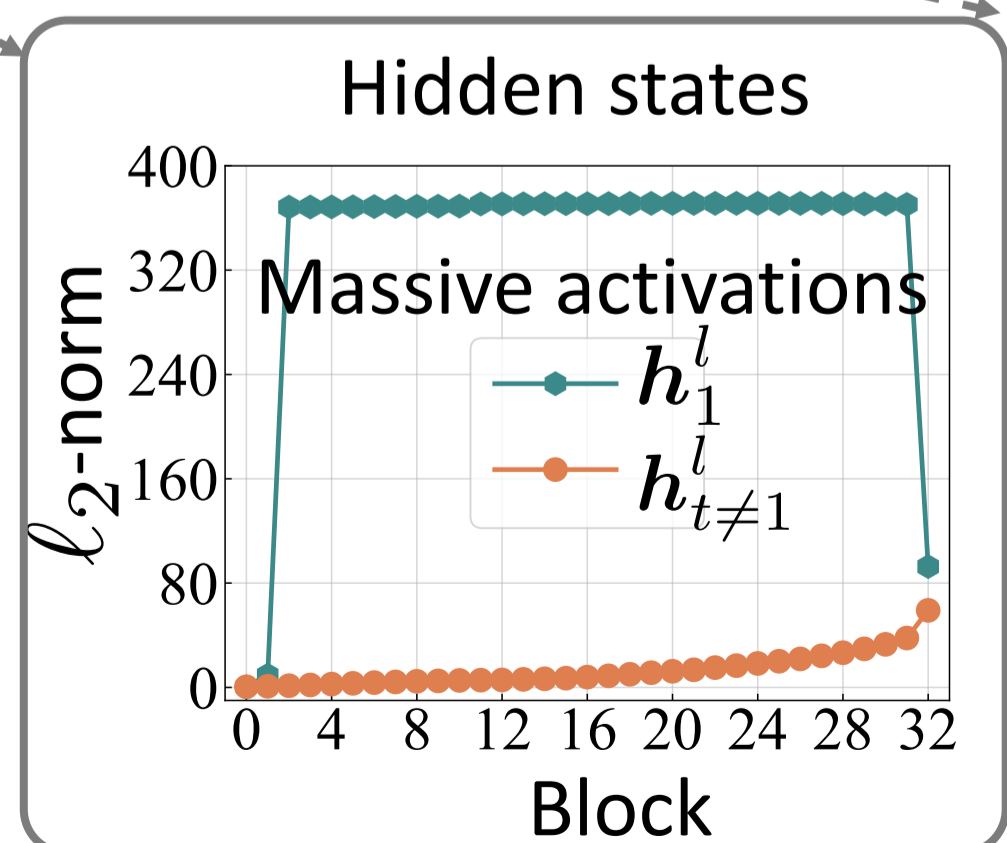
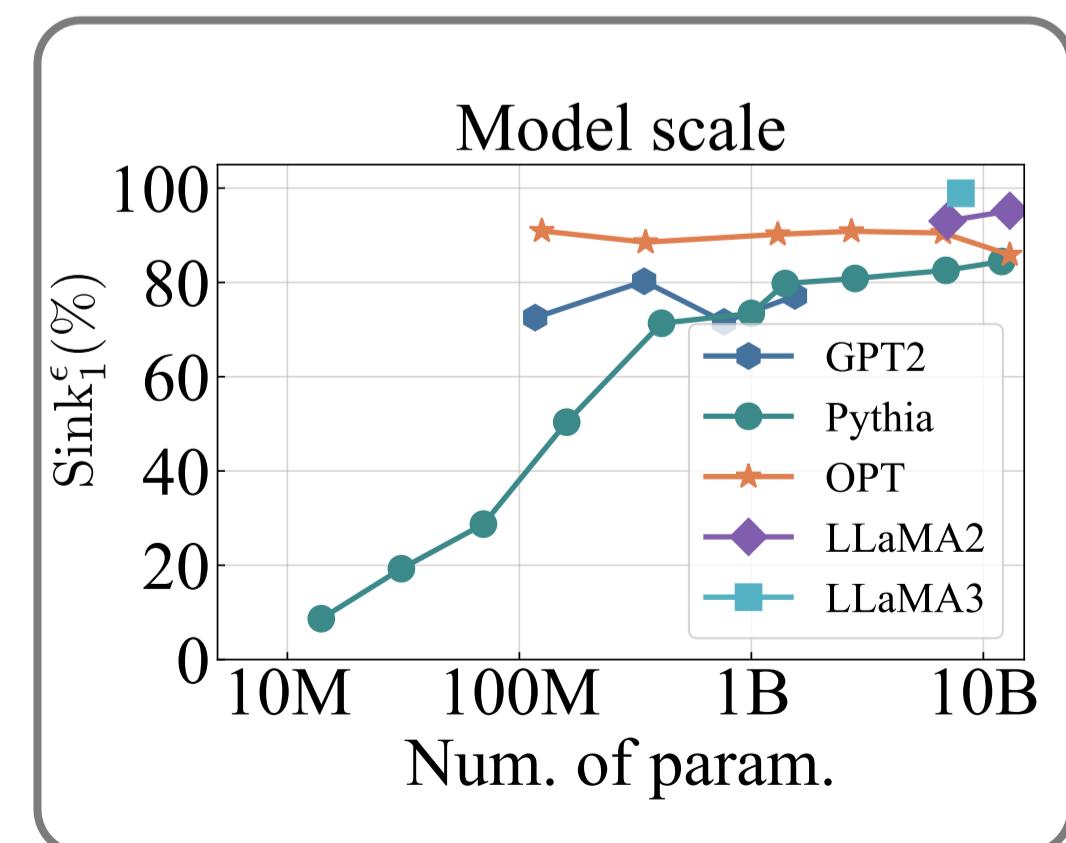
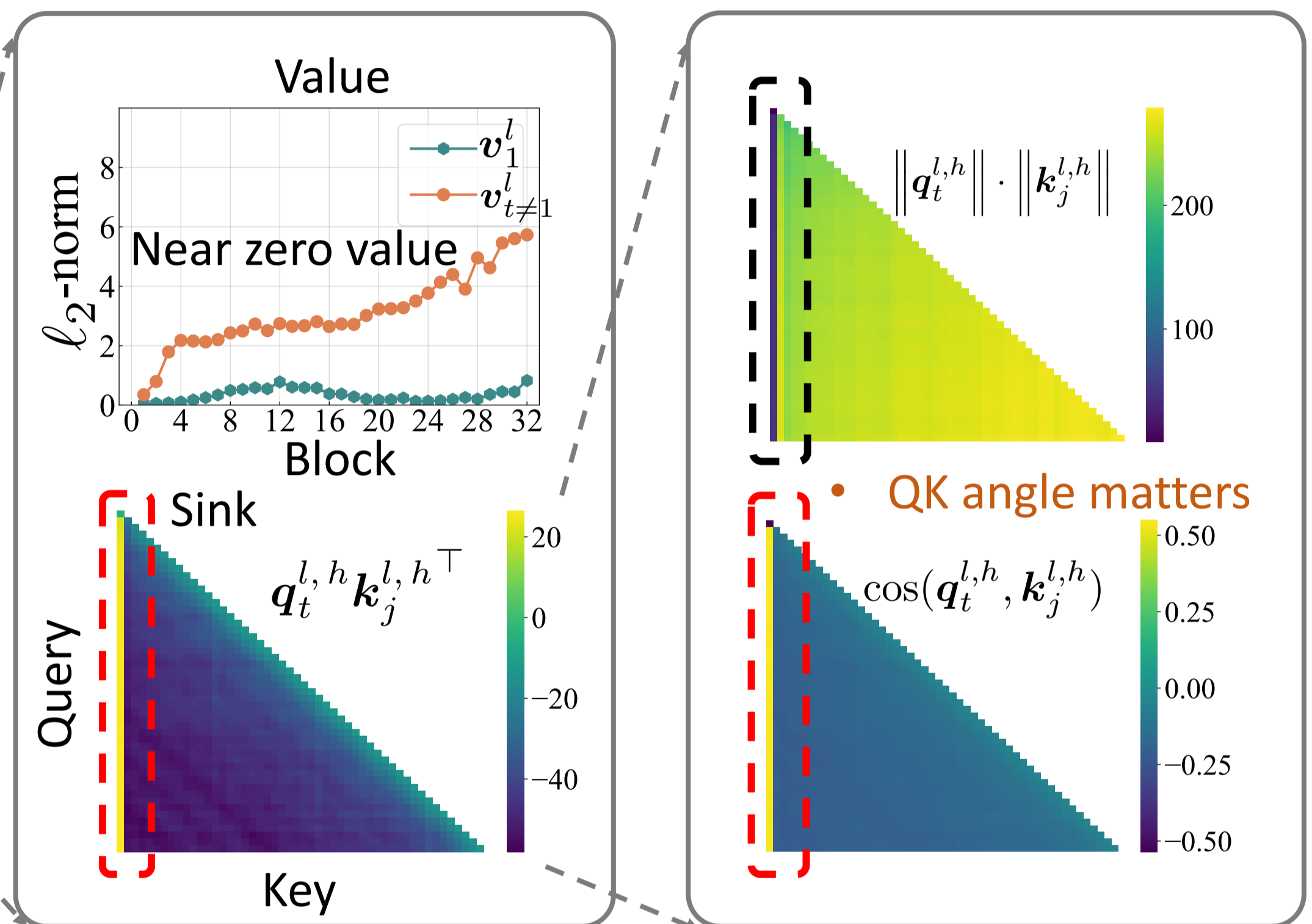
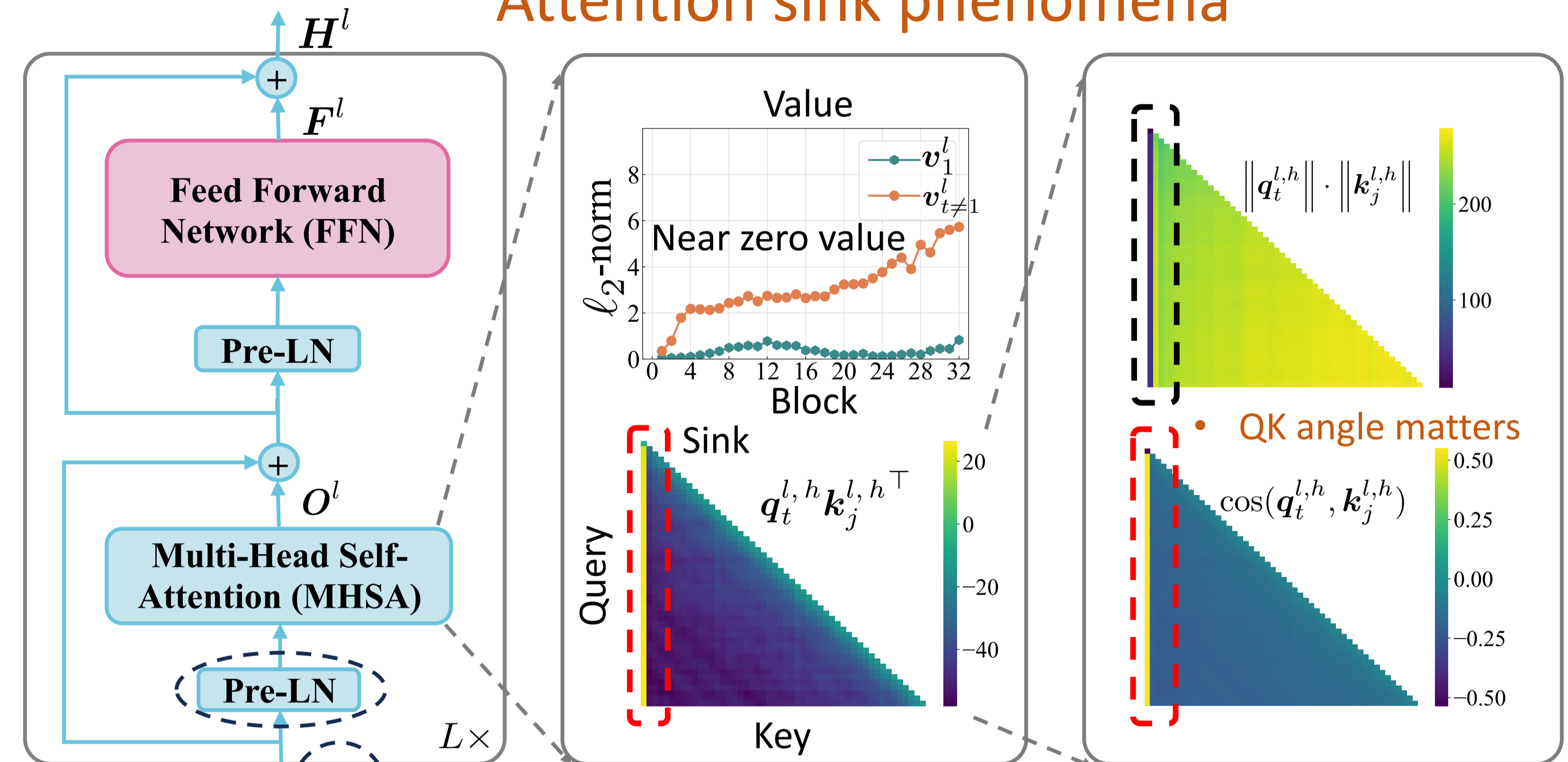




Xiangming Gu^{1,2}, Tianyu Pang¹, Chao Du¹, Qian Liu¹, Fengzhuo Zhang^{1,2}, Cunxiao Du¹, Ye Wang², Min Lin¹

¹Sea AI Lab ²National University of Singapore

Attention sink phenomena



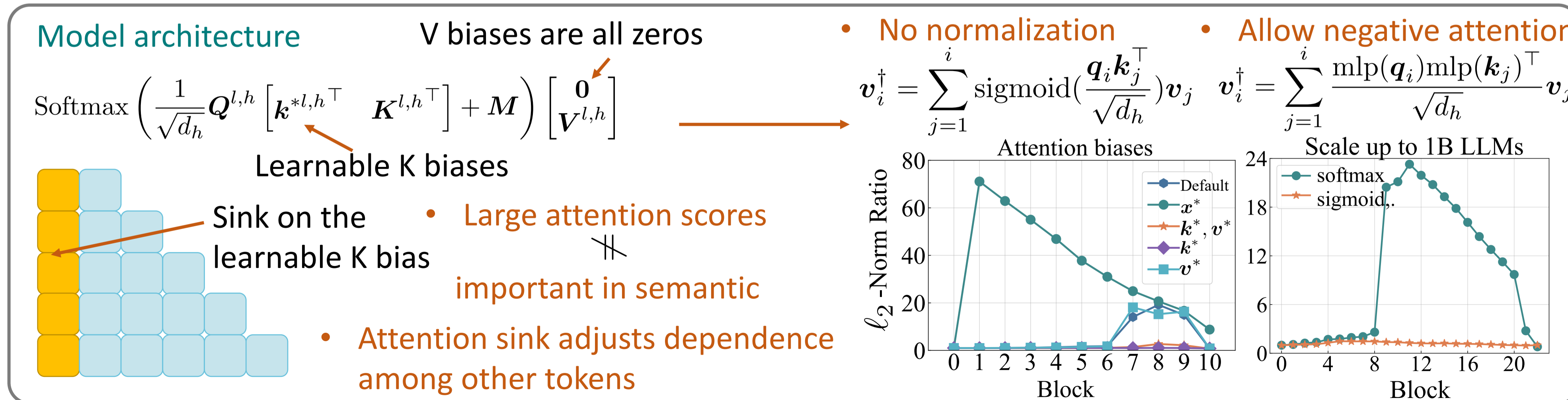
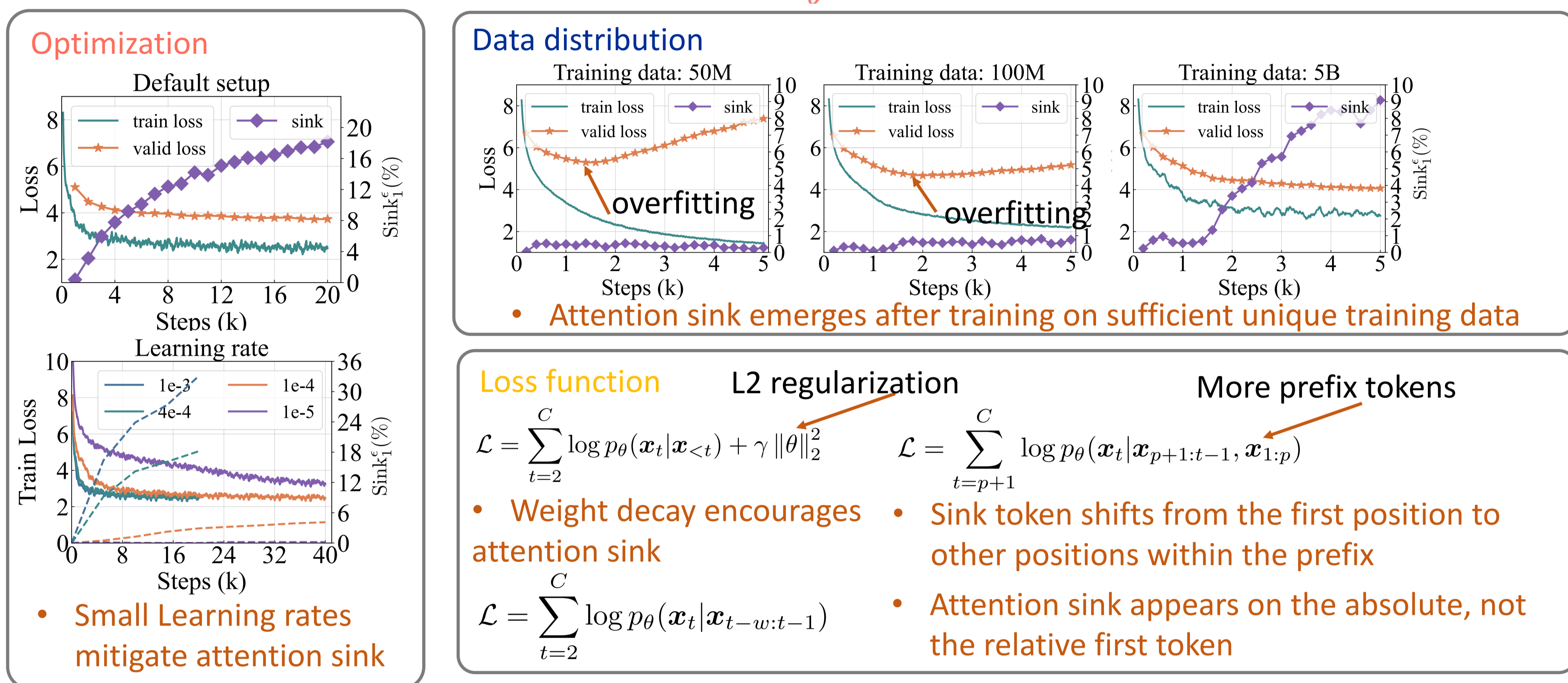
Key of the first token is distributed in a different manifold

$$\mathbf{k}_t^{l,h} = \text{LN}(\mathbf{h}_t^{l-1}) \mathbf{W}_K^{l,h} \mathbf{R}_{\Theta, -t}$$

$$\text{LN}(\mathbf{h}) = \frac{\mathbf{h}}{\sqrt{\frac{1}{d} \sum_{i=1}^d h_i^2}} \odot \mathbf{g}$$

• Attention sink even exists in Pythia-14M

LM pre-training $\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\mathcal{L}(p_{\theta}(\mathbf{X}))]$



Find more interesting conclusions in our paper!