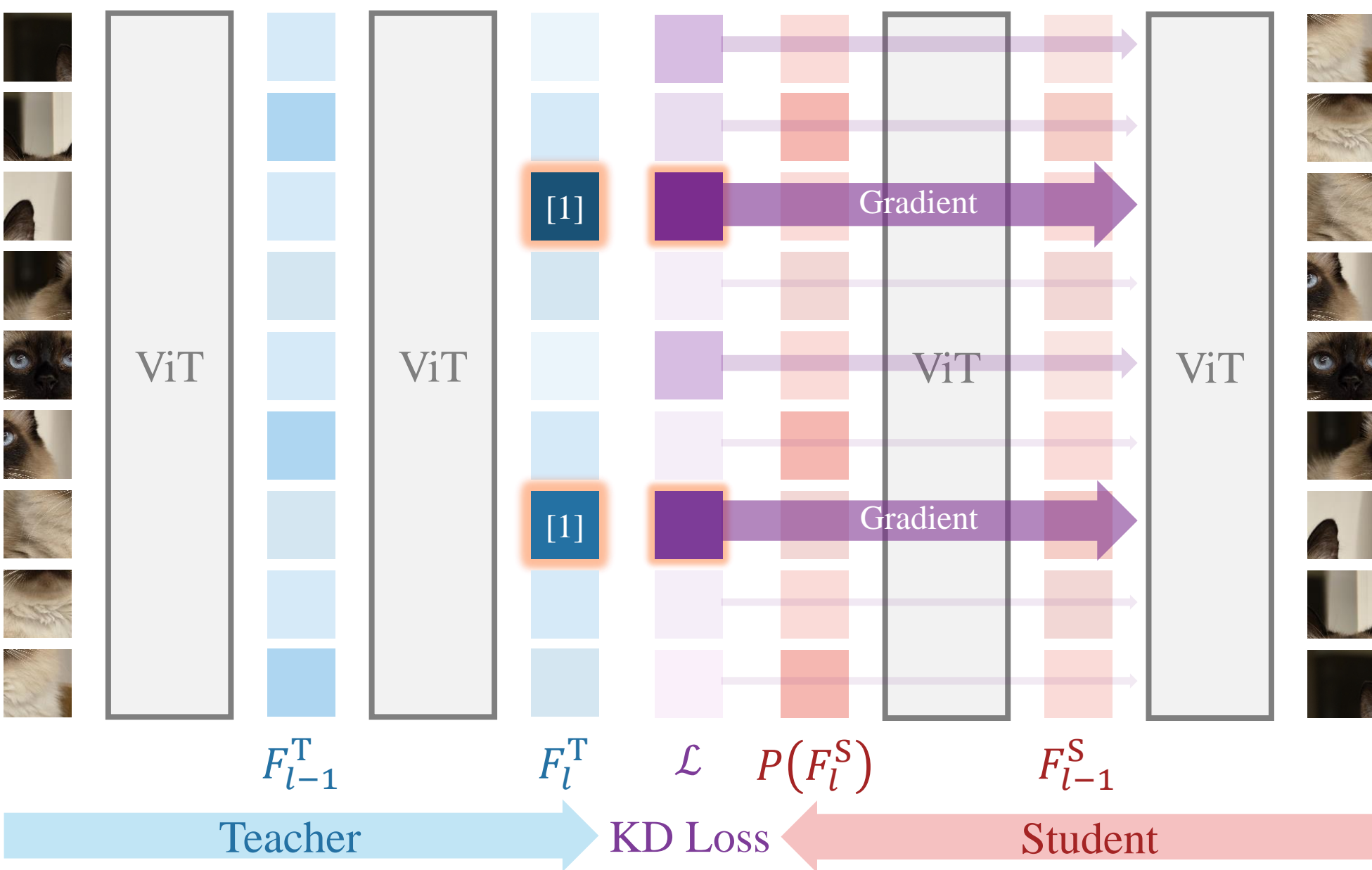




I. Gradient Dominance



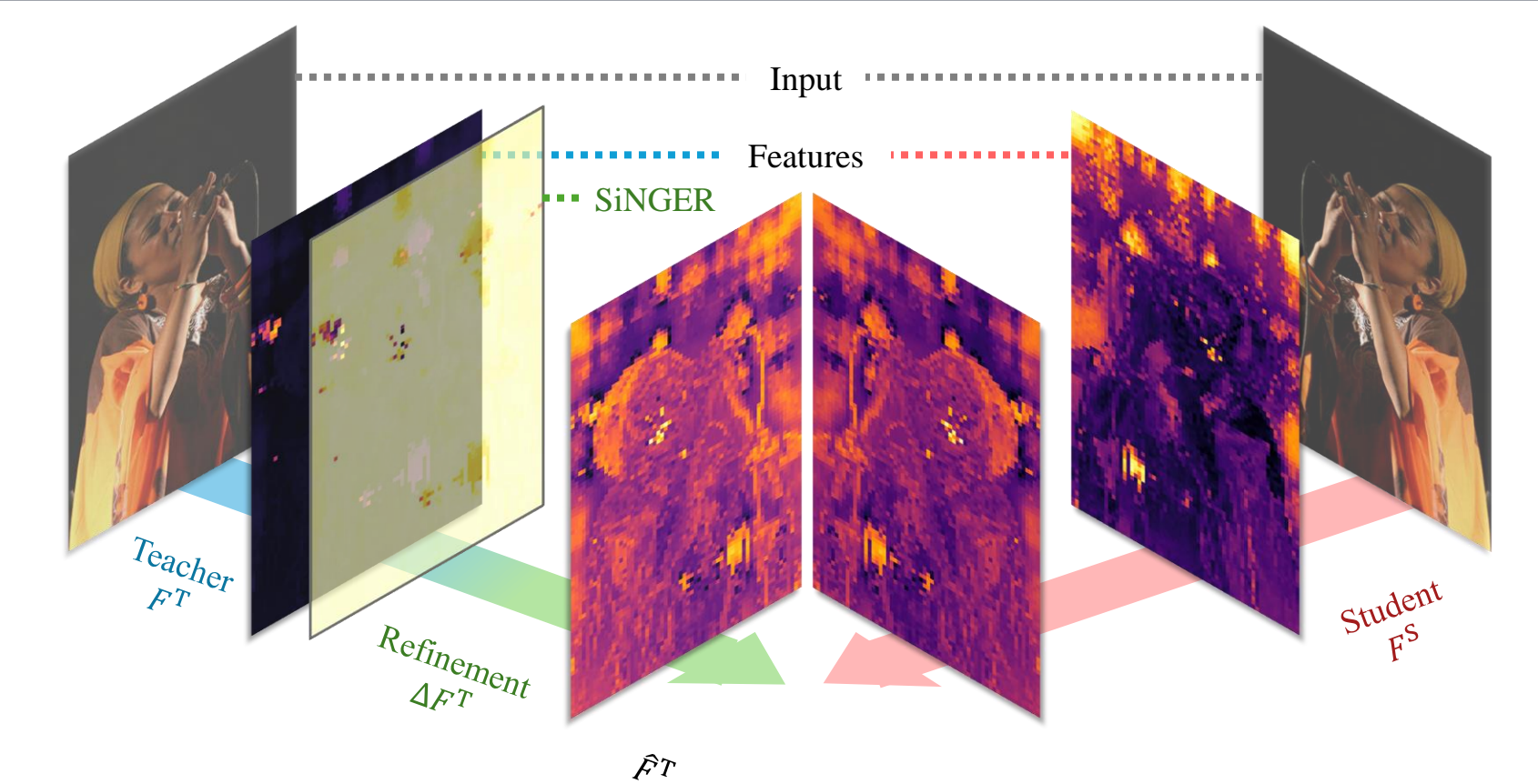
$$\mathcal{L} = \frac{1}{N} \sum_i \|P_l(F_{l,i}^S) - F_{l,i}^T\|_2^2 \quad \dots \text{Distance-based loss is proportional to the norm}$$

$$\nabla_{P_l(F_{l,i}^S)} \mathcal{L} = \frac{2}{N} \sum_i P_l(F_{l,i}^S) - F_{l,i}^T \quad \dots \text{thus, the outliers dominate the gradient portion, lead to representation degradation.}$$

✧ Contributions ✧

- ✓ Propose a distillation framework SiNGER that refines teacher signals via a lightweight adapter with nullspace initialization to guide effective perturbations.
- ✓ Analyze a fundamental limitation of naïve ViT distillation, showing degraded transfer on downstream benchmarks along with qualitative evidence.
- ✓ Provide extensive ablation studies to analyze the contribution of each component in SiNGER and validate the robustness of our framework.
- ✓ Demonstrate through extensive experiments that our method exceeds baseline performance across tasks and produces more interpretable feature maps.

II. Training Objectives

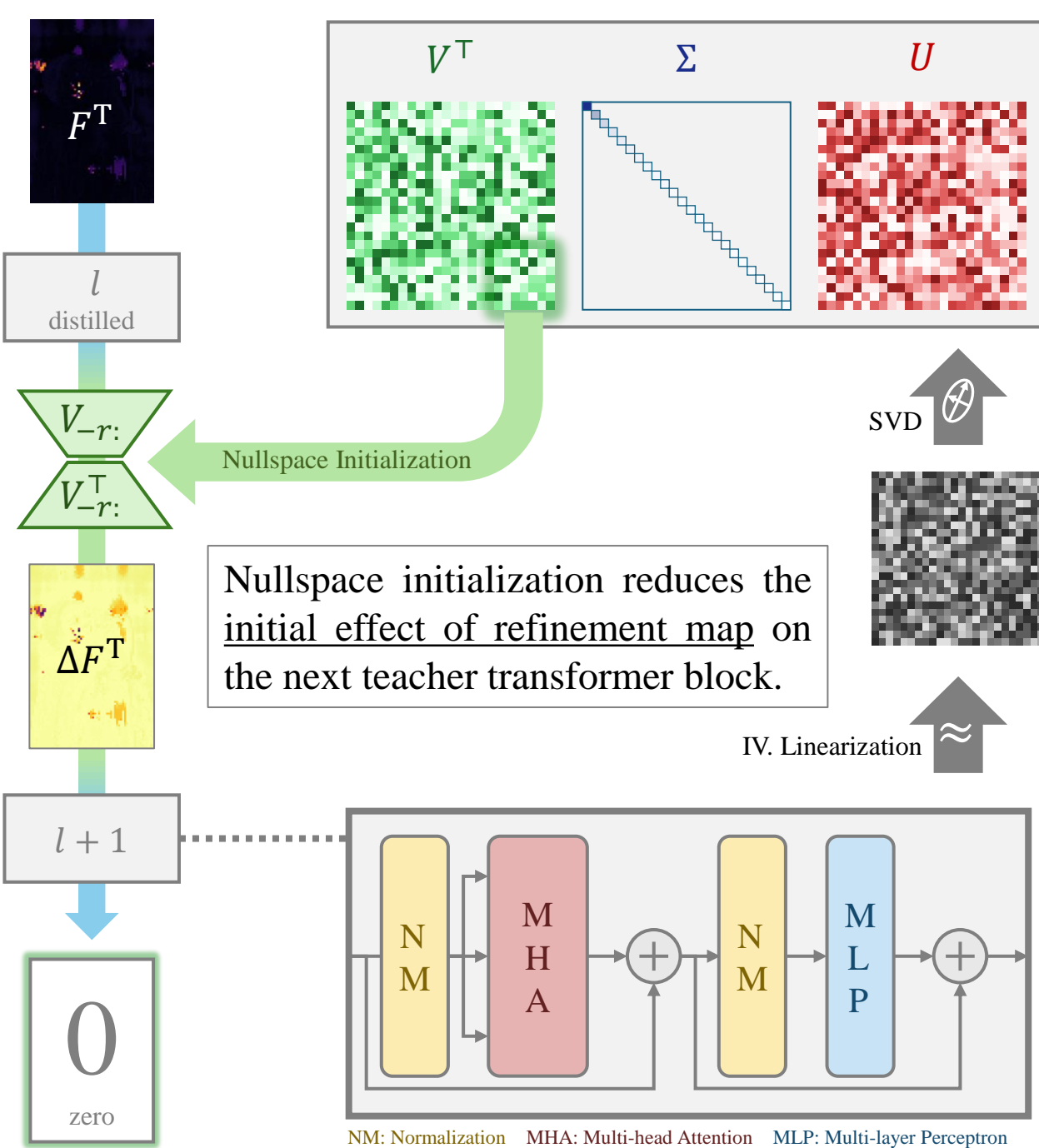


$$\mathcal{L}_{\text{outlier}} = \sum_{l \in \mathcal{D}} \frac{1}{|\mathcal{O}_l|} \sum_{i \in \mathcal{O}_l} (\|\hat{F}_{l,i}^T\|_2 - q_{\alpha,l})$$

$$\mathcal{L}_{\text{info},l} = \begin{cases} \text{MSE}(G(\hat{F}_{l+1}^T), G(F_{l+1}^T)), & l \in l_{\text{inter}}, \\ \text{MSE}(G(\hat{F}_l^T), G(F_l^T)), & l = l_{\text{final}}. \end{cases}$$

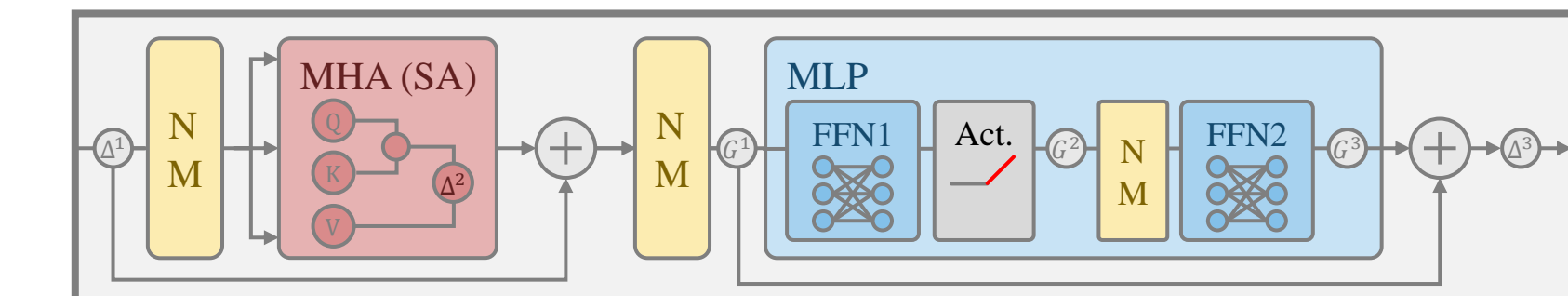
G : Gram matrix

III. Adapter Initialization



IV. Linearization

- To enable SVD on inherently non-linear Transformer blocks, we utilize a linear surrogate ($\tilde{W} = W_{\text{FFN1}} W_{\text{FFN2}}$) focusing on the FFN sub-layer, which our empirical analysis identifies as the primary driver of norm inflation.
- Verification against the full Jacobian baseline confirms that this FFN-centric proxy is highly robust, preserving the teacher's functional behavior with negligible output deviation.



Metric	SiNGER	Jacobian	Metric	mean	median	p95
L_2 (\downarrow)	0.169	0.191	$G_{\text{FFN1}}(G^1, G^2)$	0.0719	0.0714	0.0844
Cosine sim (\uparrow)	0.9787	0.9564	$G_{\text{FFN2}}(G^2, G^3)$	0.7988	0.7945	0.8571
CKA (\uparrow)	0.9975	0.9947	$\Delta_{\text{FFN}}(\Delta^1, \Delta^2)$	0.1871	0.1852	0.2210
			$\Delta_{\text{SA}}(\Delta^2, \Delta^3)$	0.1417	0.1415	0.1611

Table 8. The output deviation of the non-linear block when inputs are perturbed along null directions computed by ours versus the full Jacobian. Table 7 in Appendix A. Identification of norm inflation. The p95 metric refers to the norm at the top 5% quantile. G : rhs/lhs, Δ : (rhs - lhs)/lhs.

V. Results

Teacher	Student	Distillation	IN-val top-1 (\uparrow)	ADE-20K mIoU (\uparrow)	NYUd-v2 RMSE (\downarrow)	iNat2019 top-1 (\uparrow)	Domain Shift top-1 (\uparrow)	Fine-Grained top-1 (\uparrow)
ViT-L [2]	ViT-T [2]	FitNet [4]	<u>62.43</u>	<u>18.73</u>	<u>1.0093</u>	<u>40.02</u>	<u>32.32</u>	<u>62.48</u>
		ViTKD [5]	5.07	11.92	1.1903	23.69	2.08	33.52
		SiNGER	70.59	21.76	0.9406	41.11	38.87	64.61
		Δ	+8.16	+3.03	+0.0687	+1.09	+6.55	+2.13
DeiT-III-L [3]	DeiT-III-B [3]	FitNet [4]	60.00	<u>26.79</u>	<u>1.1625</u>	50.04	31.53	<u>74.78</u>
		ViTKD [5]	<u>66.54</u>	19.58	1.2525	30.78	<u>35.77</u>	58.36
		SiNGER	79.37	29.47	1.1514	55.50	49.14	57.41
		Δ	+12.83	+2.68	+0.0111	+3.46	+13.37	+0.63

Table 1. Multi-task evaluation results. Δ rows indicate the performance gains of SiNGER, computed against the best-performing baseline among the distilled students (underlined).

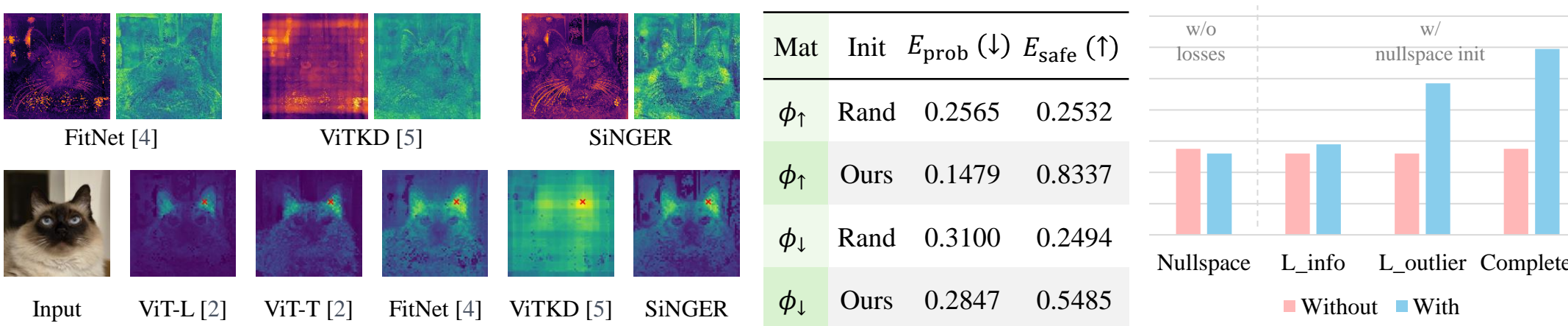


Figure 2. Qualitative analysis. Row 1: KD method comparison. Row 2: Feature map comparison. Table 5a. Ablation study on nullspace initialization ($l = 17$). Figure from Table 4. Ablation study on the nullspace initialization and loss.