

Life-IQA: Boosting Blind Image Quality Assessment through GCN-enhanced Layer Interaction and MoE-based Feature Decoupling

Supplementary Material

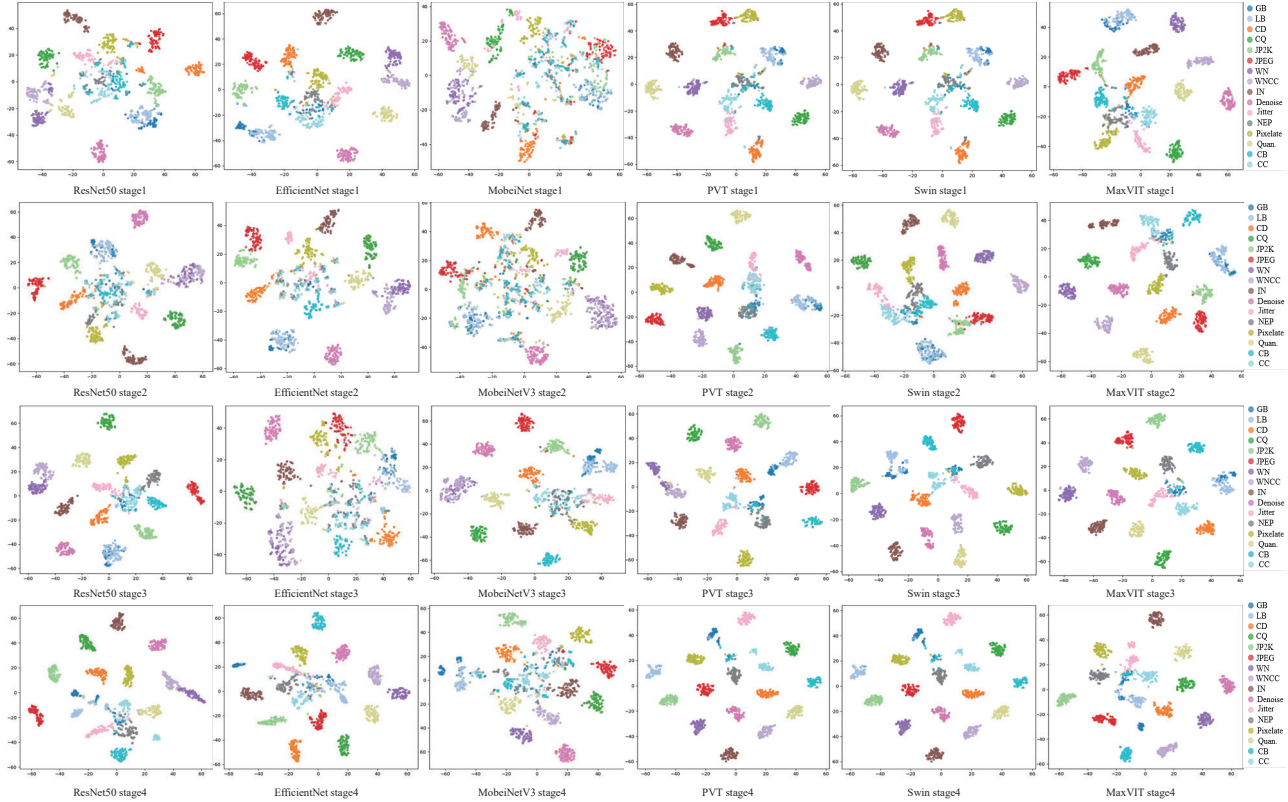


Figure 10. t-SNE visualization of multi-scale features from different pretrained models on the KADID-10k dataset.

Method	CSIQ		KADID-10K	
	SROCC	PLCC	SROCC	PLCC
Stage2 \iff Stage1	0.951	0.959	0.915	0.918
Stage3 \iff Stage1	0.960	0.965	0.937	0.940
Stage3 \iff Stage2	0.962	0.968	0.938	0.938
Stage4 \iff Stage1	0.958	0.965	0.930	0.930
Stage4 \iff Stage2	0.964	0.968	0.937	0.939
Life-IQA	0.966	0.971	0.940	0.943

Table 6. Cross-layer pairing study: deep and shallow interactions for BIQA

503 6. Cross-layer pairing study

504 Table 6 varies only the interaction pair in the GCN-
505 enhanced layer interaction, while the decoder depth, heads,

embedding size, query count, training recipe, and parameter budget are kept fixed. A consistent pattern emerges:

- 506 • **Deep with deep performs best.** Using Stage4 with Stage3 (Life-IQA) achieves 0.966/0.971 on CSIQ and 0.940/0.943 on KADID-10K. 507
- 508 • **Deep with shallow is weaker.** Relative to Stage4 with Stage1, the deep-deep pairing gains +0.008/+0.006 SROCC/PLCC on CSIQ and +0.010/+0.013 on KADID-10K. 509 510
- 511 • **Deep with mid sits in between.** Compared with Stage4 with Stage2, the deep-deep pairing improves by +0.002/+0.003 on CSIQ and +0.003/+0.004 on KADID-10K; compared with Stage3 with Stage2, the gains are +0.004/+0.003 on CSIQ and +0.002/+0.005 on KADID-10K. 512 513 514
- 515 • **Shallow with shallow is the weakest.** The Stage2 \iff Stage1 pairing yields 0.951/0.959 on CSIQ and 0.915/0.918 on KADID-10K, underperforming 516 517 518 519 520 521 522 523

524 the deep–deep setting by $-0.015/-0.012$ (CSIQ) and
525 $-0.025/-0.025$ (KADID-10K).

526 The trend holds for both rank correlation (SROCC) and lin-
527 ear correlation (PLCC), indicating that interactions between
528 semantically compatible deep representations yield more
529 reliable quality cues, whereas shallow–shallow coupling
530 provides limited semantic separation and thus the weakest
531 performance.

532 7. t-SNE analysis of stage-wise features

533 Figure 10 visualizes features from four stages of several
534 pretrained backbones on KADID-10K (25 distortion types).
535 A clear monotonic trend is observed across all architec-
536 tures: from stage1 to stage4, clusters become progressively
537 more compact with larger inter-class margins, while cross-
538 class overlap is reduced. Early layers exhibit scattered
539 manifolds with substantial mixing among distortion cate-
540 gories, indicating dominance of low-level textures and color
541 edges. Mid-level layers begin to separate families of ar-
542 tifacts (e.g., blur, noise, compression), yet boundaries re-
543 main porous. Deep layers form well-delimited groups with
544 sparse interstitial points, suggesting that higher-level repre-
545 sentations encode distortion-aware semantics that are more
546 linearly separable. The pattern holds for both convolutional
547 (ResNet50, EfficientNet, MobileNetV3) and Transformer-
548 style models (PVT, Swin, MaxViT), with the latter gener-
549 ally showing clearer margins at deeper stages. These obser-
550 vations support using deep features for query formation and
551 shallower deep features as keys/values in the proposed in-
552 teraction module, as deep–deep pairing exploits the stronger
553 class separability of later stages while retaining comple-
554 mentary detail from preceding stages.