# FOUNDATION MODELS BOOST LOW-LEVEL PERCEPTUAL SIMILARITY METRICS

*Abhijay Ghildyal*[1] , *Nabajeet Barman*[2] , *Saman Zadtootaghaj*[2]

[1]Department of Computer Science, Portland State University, USA, abhijay@pdx.edu
[2]Sony Interactive Entertainment (PlayStation), {Nabajeet.Barman,Saman.Zadtootaghaj}@sony.com

## ABSTRACT

For full-reference image quality assessment (FR-IQA) using deep-learning approaches, the perceptual similarity score between a distorted image and a reference image is typically computed as a distance measure between features extracted from a pretrained CNN or more recently, a Transformer network. Often, these intermediate features require further fine-tuning or processing with additional neural network layers to align the final similarity scores with human judgments. So far, most IQA models based on foundation models have primarily relied on the final layer or the embedding for the quality score estimation. In contrast, this work explores the potential of utilizing the intermediate features of these foundation models, which have largely been unexplored so far in the design of low-level perceptual similarity metrics. We demonstrate that the intermediate features are comparatively more effective. Moreover, without requiring any training, these metrics can outperform both traditional and state-of-the-art learned metrics by utilizing distance measures between the features. Code: `https://github.com/abhijay9/ZS-IQA`

***Index Terms***— FR-IQA, Zero-shot, Foundation Models

## 1. INTRODUCTION

It is well established that neural networks trained on ImageNet [1], using architectures like AlexNet [2], VGG [3], EfficientNet [4], and others, provide effective metrics for assessing low-level perceptual similarity [5, 6]. Recent advancements reveal that large vision models like CLIP [7] and DINO [8], originally designed for high-level tasks such as image recognition and multimodal understanding, are also effective in assessing mid to low-level perceptual similarity [9, 10, 11]. For example, Fu *et al.* [10] utilize the final embeddings of CLIP's [7] ViT vision encoder and DINO to compute cosine distance, providing a similarity score between two images, and find application in tasks like image retrieval and synthesis. Other previous works [11, 9] also focus on the final embedding rather than the *intermediate features*. The embedding is the final representation produced by the model, typically from the last layer, after processing the input image. These embeddings capture high-level semantic information and are used for tasks such as image-text

matching, retrieval, and classification. Embeddings serve as global representations that summarize the entire input. In contrast, *intermediate features* refer to activations within the model prior to the computation of the final embedding. These activations capture localized, lower-level information, including features like edges, textures, and various patterns. Consequently, for developing low-level perceptual similarity metrics, past approaches have consistently relied on the *intermediate features* of models [5, 12, 13]. Based on these observations, we pose the question: *Could utilizing the intermediate features of large foundation models result in the development of more effective low-level perceptual similarity metrics?* To address this question, this study explores several aspects, with the following key contributions:
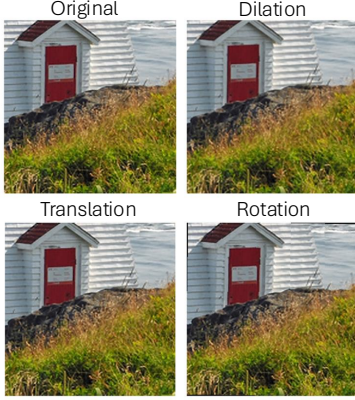
1. We compare whether embeddings or *intermediate features* are more effective for developing a more accurate and robust low-level perceptual similarity metric.
2. Through evaluations on various datasets and across different distribution and distance measures, we demonstrate that foundation models like DINO and CLIP variants yield more accurate and robust metrics.

## 2. RELATED WORK

Traditional metrics rely on pixel-level comparisons. PSNR measures direct pixel value differences, while SSIM [15] and its variants such as MS-SSIM [16], and FSIMc [17] leverage the understanding of HVS and measure differences in luminance, contrast, and structure across local and global image regions, enhancing robustness to misalignments.

In contrast, deep learning-based metric such as LPIPS [5] uses *intermediate features* extracted from pre-trained networks like AlexNet [2] and VGG [3] to assess perceptual similarity, aligning more closely with human judgments. The performance is further improved by fine-tuning on datasets specifically tailored for human perceptual similarity.

The LPIPS metric processes the *intermediate features* of the two images using 1x1 convolution layers to produce weighted feature maps, and subsequently calculates the $l_2$ distance between these maps. However, recent research shows that using the same CNN-based backbones, such as VGG and EfficientNet, and directly comparing the *intermediate features* with distribution comparison measures like Wasserstein

**Fig. 1**: Geometric transformations. Despite imperceptible changes, a metric's rank predictions often fluctuate.

| Model | Type | | TID2013 | | | | PIPAL | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | emb. | feats. | Ori. | Tra. | Dil. | Rot. | Ori. | Tra. | Dil. | Rot. |
| CLIP-RN50 [7] | ✓ | | 0.518 | 0.510 | 0.499 | 0.527 | 0.320 | 0.317 | 0.315 | 0.319 |
| | | ✓ | 0.604 | 0.616 | 0.674 | 0.638 | 0.581 | 0.538 | 0.567 | 0.570 |
| CLIP-ConvNext [7] | ✓ | | 0.632 | 0.625 | 0.697 | 0.654 | 0.479 | 0.477 | 0.474 | 0.482 |
| | | ✓ | 0.658 | 0.665 | 0.753 | 0.717 | 0.598 | 0.552 | 0.587 | 0.595 |
| CLIP-ViT-B [7] | ✓ | | 0.693 | 0.670 | 0.748 | 0.683 | 0.531 | 0.512 | 0.526 | 0.515 |
| | | ✓ | 0.758 | 0.702 | 0.807 | 0.784 | 0.623 | 0.544 | 0.597 | 0.593 |
| DINOv1-ViT-B [8] | ✓ | | 0.783 | 0.717 | 0.792 | 0.810 | 0.622 | 0.597 | 0.615 | 0.631 |
| | | ✓ | 0.786 | 0.730 | 0.804 | 0.814 | 0.637 | 0.607 | 0.629 | 0.643 |
| DINOv2-ViT-B [14] | ✓ | | 0.727 | 0.672 | 0.745 | 0.747 | 0.514 | 0.484 | 0.508 | 0.523 |
| | | ✓ | 0.722 | 0.679 | 0.763 | 0.766 | 0.573 | 0.523 | 0.558 | 0.580 |

**Table 1**: Comparison of SRCC scores for embeddings (emb.) and intermediate features (feats.) using Cosine distance across the TID2013 and PIPAL datasets with various CLIP and DINO backbones. In each subgroup, best model results are underlined, with overall best results highlighted in red and second-best in blue.

Distance (WSD), Jensen-Shannon Divergence (JSD), and Symmetric Kullback-Leibler Divergence (SKLD) can yield improved perceptual similarity scores [18]. These enhanced metrics, based on distribution measures, align with human perceptual similarity without any additional training.

Other metrics such as DISTS [12] adopt a design similar to the LPIPS(VGG) [5] but use measures akin to SSIM rather than $l_2$ distance between feature maps, thereby improving their metric's robustness. They also incorporate $l_2$-pooling layers, which blur *intermediate features* to enhance DISTS's robustness. Another recent work investigates aliasing in *intermediate features* resulting from imperceptible geometric misalignments [13]. By examining various neural network components, including max-pool, stride, and others, they aligned their LPIPS-based metric's sensitivity with human perception of an imperceptible misalignment between images.

Recent metrics [10, 9] that employ foundation models like CLIP [7] and DINO [14] as backbones have mainly focused on using embeddings, rather than *intermediate features*. Although [19] proposed utilizing *intermediate features* for CLIP variants, their evaluation and analysis were restricted to the TID2013 dataset. In contrast, our study offers a more comprehensive analysis by examining performance across three different datasets, including the larger PIPAL dataset, which encompasses both synthetic and real algorithmic distortions. We also evaluate various DINO and CLIP backbone architectures, utilizing a range of distribution and distance measures, and compare our results with existing low-level FR-IQA methods to demonstrate that using *intermediate features* in VLMs are also effective for improved robustness.

## 3. EXPERIMENTS AND RESULTS

In our study, we evaluate the accuracy of several low-level perceptual similarity metrics by comparing their performance using correlation measures like PLCC, SRCC, and KRCC. Our experiments primarily focus on two key aspects: demon-strating that *intermediate features* outperform embeddings for the low-level perceptual similarity task (Table 1) and showing that CLIP and DINO models serve as superior backbones, generating more robust *intermediate features* for enhanced low-level perceptual similarity metrics (Table 2 and Table 3).

For the CLIP and DINO metrics with ViT backbones, we evaluate results using a square sliding window of size 224 with a stride of 200. For comparisons against previous state-of-the-art training-free metrics using VGG and EfficientNet, we use the authors' code [18] and perform evaluations without downsampling the input image (due to lack of information on the exact downsampling factor used by the authors for each dataset). Therefore, for uniformity, all our experiments are conducted on full-sized images without downsampling.

To investigate the robustness of these models, we adopt the experimental design from [12], with the modification of applying geometric transformations to the distorted image while keeping the reference image unchanged. This approach better simulates real-world scenarios where distortions are more likely to be present in the image being analyzed. Specifically, we shift the distorted image horizontally to the right by 1% of the pixels for translation, increase its scale by 1% for dilation, and rotate it clockwise by 1 degree for rotation. A sample of this process is illustrated in Figure 1.

**Embeddings vs. Intermediate Features.** Table 1 demonstrates that for all models, using *intermediate features* (feats.) consistently results in better performance than using embeddings (emb.), suggesting that *intermediate features* capture more relevant information for this task. This improvement is especially significant for DINOv1. Among the CLIP models, CLIP-ViT-B with using *intermediate features* shows the best performance. DINOv1-ViT-B stands out as the best overall performer across both datasets and geometric distortion types, particularly when using *intermediate features*. This indicates that DINOv1-ViT-B is highly effective for the low-level perceptual similarity task. Moreover, the results from Table 1 align with those in Table 3, where using *intermediate features*,

| Model | Dist. | LIVE | | | | | | | | TID2013 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Original | | Translation | | Dilation | | Rotation | | Original | | Translation | | Dilation | | Rotation | |
| | | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC | PLCC | SRCC |
| *Distribution Measures* | | | | | | | | | | | | | | | | | |
| VGG [3] | | 0.914 | 0.936 | 0.456 | 0.429 | 0.714 | 0.697 | 0.646 | 0.619 | 0.751 | 0.780 | 0.231 | 0.285 | 0.667 | 0.656 | 0.523 | 0.531 |
| EfficientNet [4] | SKLD | 0.800 | 0.896 | 0.808 | 0.816 | 0.869 | 0.867 | 0.868 | 0.862 | 0.577 | 0.647 | 0.692 | 0.652 | 0.718 | 0.683 | 0.733 | 0.689 |
| CLIP-ViT-B [7] | [18] | 0.938 | 0.965 | 0.895 | 0.894 | 0.936 | 0.942 | 0.937 | 0.940 | 0.730 | 0.765 | 0.768 | 0.717 | 0.843 | 0.815 | 0.849 | 0.819 |
| DINOv1 [8] | | 0.949 | 0.964 | 0.893 | 0.883 | 0.920 | 0.914 | 0.943 | 0.941 | 0.770 | 0.768 | 0.769 | 0.709 | 0.828 | 0.788 | 0.848 | 0.813 |
| VGG [3] | | 0.910 | 0.928 | 0.436 | 0.405 | 0.690 | 0.667 | 0.620 | 0.587 | 0.752 | 0.774 | 0.209 | 0.268 | 0.649 | 0.639 | 0.500 | 0.513 |
| EfficientNet [4] | JSD | 0.863 | 0.905 | 0.834 | 0.839 | 0.875 | 0.870 | 0.875 | 0.867 | 0.642 | 0.649 | 0.719 | 0.669 | 0.733 | 0.689 | 0.744 | 0.694 |
| CLIP-ViT-B [7] | [18] | 0.852 | 0.959 | 0.869 | 0.879 | 0.925 | 0.936 | 0.932 | 0.934 | 0.670 | 0.764 | 0.742 | 0.710 | 0.809 | 0.806 | 0.838 | 0.809 |
| DINOv1 [8] | | 0.951 | 0.964 | 0.889 | 0.879 | 0.917 | 0.911 | 0.941 | 0.940 | 0.774 | 0.768 | 0.765 | 0.704 | 0.827 | 0.785 | 0.847 | 0.812 |
| VGG [3] | | 0.916 | 0.933 | 0.457 | 0.428 | 0.714 | 0.692 | 0.649 | 0.619 | 0.764 | 0.780 | 0.257 | 0.323 | 0.693 | 0.680 | 0.552 | 0.560 |
| EfficientNet [4] | WSD | 0.837 | 0.897 | 0.780 | 0.775 | 0.879 | 0.864 | 0.868 | 0.852 | 0.625 | 0.656 | 0.629 | 0.586 | 0.736 | 0.685 | 0.741 | 0.684 |
| CLIP-ViT-B [7] | [18] | 0.950 | 0.964 | 0.903 | 0.903 | 0.935 | 0.941 | 0.936 | 0.942 | 0.735 | 0.744 | 0.747 | 0.702 | 0.823 | 0.795 | 0.824 | 0.794 |
| DINOv1 [8] | | 0.954 | 0.965 | 0.899 | 0.888 | 0.926 | 0.920 | 0.943 | 0.943 | 0.774 | 0.765 | 0.778 | 0.716 | 0.832 | 0.791 | 0.848 | 0.811 |
| *Other Distance Measures* | | | | | | | | | | | | | | | | | |
| CLIP-ViT-B [7] | $l_2$ | 0.947 | 0.968 | 0.899 | 0.909 | 0.924 | 0.947 | 0.919 | 0.946 | 0.831 | 0.793 | 0.797 | 0.727 | 0.842 | 0.821 | 0.844 | 0.824 |
| DINOv1 [8] | | 0.937 | 0.964 | 0.896 | 0.891 | 0.906 | 0.914 | 0.915 | 0.940 | 0.824 | 0.792 | 0.794 | 0.720 | 0.821 | 0.790 | 0.834 | 0.813 |
| CLIP-ViT-B [7] | Cos. | 0.937 | 0.965 | 0.869 | 0.871 | 0.908 | 0.936 | 0.899 | 0.924 | 0.817 | 0.758 | 0.778 | 0.702 | 0.840 | 0.807 | 0.833 | 0.784 |
| DINOv1 [8] | | 0.933 | 0.965 | 0.896 | 0.898 | 0.911 | 0.941 | 0.911 | 0.943 | 0.820 | 0.786 | 0.800 | 0.730 | 0.824 | 0.804 | 0.833 | 0.814 |

**Table 2**: Performance comparison of different backbone architectures on LIVE and TID2013 datasets. Within each subgroup, better model results are underlined, with the overall best performing model highlighted in red and the second-best in blue.

compared to embeddings (emb.), results in better correlation scores for both CLIP-ViT-B and DINOv1[1].

**Comparison across different backbones.** As demonstrated in Table 2 and Table 3, the ViT-B backbone models DINOv1 and CLIP-ViT-B generally perform the best across all three datasets, LIVE, TID2013, and PIPAL, and all geometric distortions, indicating higher consistency with human perceptual judgments. In contrast, VGG shows the lowest performance, particularly in robustness under geometric transformations.

**Comparison across distribution and distance measures.** Based on the results in Table 2 and Table 3, the different distribution measures, i.e., SKLD, JSD, and WSD, do not notably affect the relative performance rankings of the models. Among the training-free models, the DINOv1 and CLIP-ViT-B backbones outperform the VGG and Efficient backbones, which had previously achieved state-of-the-art results using different distribution measures in a recent study [18]. The result is consistent across all three datasets, i.e., LIVE, TID2013, and PIPAL, demonstrating that the DINOv1 and CLIP-ViT-B features are both more accurate and robust.

We also explore two other distance measures: $l_2$, and cosine distance (Cos.). Since, SKLD, JSD, and WSD include a weighted Euclidean norm, we evaluate the use of $l_2$ alone. As observed in Table 2 and Table 3, the results using different distribution measures are almost the same as those obtained with $l_2$, suggesting that the distribution measures are dominated by the added Euclidean norm and that the adaptive weighting strategy requires further refinement. Additionally, previous studies show that cosine distance between embeddings aligns well with perceptual similarity scores, so we compute it between *intermediate features* as well.

In summary, although CLIP-ViT-B performs better than DINOv1 on the LIVE and TID datasets with $l_2$ distance, DINOv1 exhibits significantly better performance on the larger PIPAL dataset. For cosine distance, SKLD, JSD, and WSD measures, DINOv1 consistently outperforms CLIP-ViT-B across all datasets. *Thus, we recommend DINOv1 for low-level perceptual similarity tasks.*

**Comparison against existing metrics.** Since none of these metrics were trained on the PIPAL Training dataset [21], which includes a wide range of synthetic and computer vision algorithm-based distortions, it serves as a good test set for comparing all methods. Based on the results presented in Table 3, we can conclude that FSIMc is the best performing metric among the considered traditional methods. For the learned methods, we chose the most robust and high-performing metrics. Among the learned models we evaluated, ST-LPIPS(AlexNet) demonstrated superior performance. This might be due to the fact that LPIPS and ST-LPIPS are trained on the BAPPS dataset [5], which includes a range of distortions from computer vision algorithms like super-resolution, frame interpolation, deblurring, and colorization. In contrast, DISTS is trained on the KADID-10k dataset [22], which mainly includes synthetic distortions.

Despite not being fine-tuned for low-level perceptual similarity tasks, training-free metrics with CLIP-ViT-B and DINOv1 backbones perform well. Among these, DINOv1 stands out as the top performer among all models, particularly in its handling of various geometric distortions.

**Comparison against ImageNet-ViT.** In terms of $l_2$ distance, CLIP-ViT and DINOv1-ViT outperform ImageNet-ViT [20] significantly. However, for cosine distance CLIP-ViT and ImageNet-ViT perform similarly while DINOv1 has a significantly better performance. Both DINOv1 and ImageNet-ViT

---

[1]In the rest of the paper, *intermediate features* (feats.) are used, especially in Tables 2 and 3, unless the embedding (emb.) is specifically indicated.

| | | PIPAL | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Dist. | Original | | | Translation | | | Dilation | | | Rotation | | |
| | | PLCC | SRCC | KRCC | PLCC | SRCC | KRCC | PLCC | SRCC | KRCC | PLCC | SRCC | KRCC |
| *Traditional* | | | | | | | | | | | | | |
| PSNR | | 0.403 | 0.407 | 0.276 | -0.065 | -0.079 | -0.052 | 0.064 | 0.038 | 0.026 | -0.024 | -0.042 | -0.028 |
| MS-SSIM [16] | - | 0.562 | 0.562 | 0.397 | -0.003 | -0.012 | -0.008 | 0.318 | 0.227 | 0.153 | 0.186 | 0.108 | 0.072 |
| FSIMc [17] | | 0.609 | 0.589 | 0.416 | 0.209 | 0.190 | 0.127 | 0.368 | 0.328 | 0.220 | 0.292 | 0.254 | 0.170 |
| *Learned* | | | | | | | | | | | | | |
| DISTS [12] | | 0.580 | 0.579 | 0.407 | 0.560 | 0.558 | 0.390 | 0.582 | 0.575 | 0.403 | 0.546 | 0.539 | 0.375 |
| LPIPS(AlexNet) [5] | | 0.586 | 0.588 | 0.412 | 0.567 | 0.567 | 0.394 | 0.539 | 0.532 | 0.366 | 0.574 | 0.573 | 0.398 |
| LPIPS(VGG) [5] | - | 0.611 | 0.588 | 0.415 | 0.548 | 0.527 | 0.366 | 0.588 | 0.561 | 0.392 | 0.584 | 0.558 | 0.390 |
| ST-LPIPS(AlexNet) [13] | | 0.628 | 0.631 | 0.448 | 0.625 | 0.628 | 0.445 | 0.618 | 0.614 | 0.434 | 0.627 | 0.630 | 0.447 |
| ST-LPIPS(VGG) [13] | | 0.578 | 0.580 | 0.406 | 0.572 | 0.573 | 0.401 | 0.557 | 0.550 | 0.382 | 0.577 | 0.580 | 0.407 |
| *Training-free* | | | | | | | | | | | | | |
| VGG [18] | | 0.536 | 0.563 | 0.392 | 0.233 | 0.228 | 0.152 | 0.401 | 0.411 | 0.277 | 0.353 | 0.350 | 0.235 |
| EfficientNet [18] | SKLD | 0.504 | 0.503 | 0.354 | 0.453 | 0.418 | 0.289 | 0.481 | 0.452 | 0.314 | 0.474 | 0.439 | 0.305 |
| CLIP-ViT-B [18] | [18] | 0.600 | 0.616 | 0.438 | 0.564 | 0.557 | 0.388 | 0.598 | 0.597 | 0.421 | 0.607 | 0.596 | 0.420 |
| DINOv1 | | 0.627 | 0.639 | 0.458 | 0.632 | 0.615 | 0.437 | 0.643 | 0.633 | 0.452 | 0.681 | 0.660 | 0.476 |
| VGG [18] | | 0.530 | 0.552 | 0.383 | 0.214 | 0.208 | 0.139 | 0.387 | 0.391 | 0.263 | 0.335 | 0.328 | 0.220 |
| EfficientNet [18] | JSD | 0.520 | 0.506 | 0.356 | 0.464 | 0.428 | 0.296 | 0.488 | 0.454 | 0.316 | 0.482 | 0.444 | 0.308 |
| CLIP-ViT-B [18] | [18] | 0.494 | 0.610 | 0.444 | 0.528 | 0.545 | 0.390 | 0.549 | 0.588 | 0.426 | 0.590 | 0.584 | 0.422 |
| DINOv1 | | 0.630 | 0.638 | 0.457 | 0.633 | 0.613 | 0.436 | 0.645 | 0.632 | 0.451 | 0.682 | 0.660 | 0.475 |
| VGG [18] | | 0.558 | 0.582 | 0.407 | 0.311 | 0.305 | 0.205 | 0.458 | 0.465 | 0.316 | 0.421 | 0.418 | 0.282 |
| EfficientNet [18] | WSD | 0.525 | 0.503 | 0.354 | 0.405 | 0.365 | 0.251 | 0.488 | 0.443 | 0.308 | 0.473 | 0.423 | 0.293 |
| CLIP-ViT-B [18] | [18] | 0.614 | 0.627 | 0.446 | 0.578 | 0.571 | 0.400 | 0.609 | 0.608 | 0.430 | 0.619 | 0.607 | 0.429 |
| DINOv1 | | 0.633 | 0.641 | 0.460 | 0.637 | 0.618 | 0.439 | 0.647 | 0.636 | 0.454 | 0.683 | 0.661 | 0.476 |
| ImageNet-ViT [20] | | 0.469 | 0.444 | 0.306 | 0.239 | 0.212 | 0.143 | 0.371 | 0.335 | 0.228 | 0.345 | 0.308 | 0.209 |
| CLIP-RN50 | | 0.567 | 0.558 | 0.391 | 0.498 | 0.488 | 0.335 | 0.550 | 0.540 | 0.375 | 0.546 | 0.537 | 0.373 |
| CLIP-ConvNext | $l_2$ | 0.122 | 0.137 | 0.091 | 0.109 | 0.121 | 0.081 | 0.099 | 0.109 | 0.073 | 0.111 | 0.124 | 0.082 |
| CLIP-ViT-B | | 0.649 | 0.619 | 0.440 | 0.600 | 0.563 | 0.393 | 0.634 | 0.601 | 0.424 | 0.634 | 0.602 | 0.425 |
| DINOv1 | | 0.683 | 0.639 | 0.458 | 0.664 | 0.616 | 0.438 | 0.679 | 0.633 | 0.452 | 0.698 | 0.660 | 0.475 |
| DINOv2 | | 0.446 | 0.416 | 0.286 | 0.340 | 0.300 | 0.203 | 0.403 | 0.363 | 0.247 | 0.371 | 0.336 | 0.228 |
| ImageNet-ViT [20] | | 0.657 | 0.616 | 0.438 | 0.612 | 0.557 | 0.391 | 0.641 | 0.592 | 0.418 | 0.645 | 0.598 | 0.423 |
| CLIP-RN50 | | 0.610 | 0.598 | 0.426 | 0.564 | 0.552 | 0.387 | 0.609 | 0.587 | 0.415 | 0.618 | 0.595 | 0.422 |
| CLIP-ConvNext | | 0.603 | 0.581 | 0.413 | 0.565 | 0.538 | 0.377 | 0.598 | 0.567 | 0.401 | 0.600 | 0.570 | 0.404 |
| CLIP-ViT-B | | 0.663 | 0.623 | 0.444 | 0.590 | 0.544 | 0.379 | 0.641 | 0.597 | 0.421 | 0.638 | 0.593 | 0.418 |
| CLIP-ViT-B emb. | Cos. | 0.522 | 0.531 | 0.369 | 0.508 | 0.512 | 0.354 | 0.517 | 0.526 | 0.365 | 0.509 | 0.515 | 0.357 |
| R-CLIP$_T$ [11] ViT-B emb. | | 0.385 | 0.547 | 0.383 | 0.382 | 0.528 | 0.367 | 0.386 | 0.455 | 0.309 | 0.384 | 0.479 | 0.331 |
| DINOv1 | | 0.677 | 0.637 | 0.456 | 0.654 | 0.607 | 0.431 | 0.671 | 0.629 | 0.448 | 0.683 | 0.643 | 0.460 |
| DINOv1 emb. | | 0.663 | 0.622 | 0.444 | 0.644 | 0.597 | 0.423 | 0.658 | 0.615 | 0.437 | 0.672 | 0.631 | 0.451 |
| DINOv2 | | 0.619 | 0.573 | 0.406 | 0.576 | 0.523 | 0.366 | 0.608 | 0.558 | 0.393 | 0.625 | 0.580 | 0.410 |

**Table 3**: Performance comparison of different metrics on the PIPAL dataset. Within each subgroup, better model results are underlined, with the overall best performing model highlighted in red and the second-best in blue.

are trained on the ImageNet dataset, highlighting the importance of this pre-training in producing better features. Furthermore, DINOv1 employs a self-distillation training method, where the student model learns to align its image representations with those of the teacher network across various augmentations, without requiring labeled data. As a result, for downstream tasks, the pretrained features from DINOv1 have been shown to outperform those from ImageNet-ViT, which is trained with supervision on ImageNet. We observed a similar trend in our work for low-level perceptual similarity.

**Embeddings of Adversarially Robust-CLIP.** There is growing interest in using adversarially trained models for perceptual similarity [23, 24]. A recent study found that cosine distance between Robust-CLIP embeddings outperforms that between CLIP embeddings. [11]. Consequently, we compared the embeddings of R-CLIP$_T$-ViT and CLIP-ViT. The results presented in Table 3 show that R-CLIP exhibits superior performance in terms of SRCC and KRCC, although it

lags behind in terms of PLCC scores. The pretrained adversarially robust CLIP is less sensitive to translation distortions but not as effective against dilation and rotation. At present, it remains inconclusive whether adversarial training improves both accuracy and geometric robustness. In the future, we plan to further explore adversarially robust backbones for the low-level perceptual similarity task.

## 4. CONCLUSION

In this study, we performed a comprehensive evaluation of different distance measures between *intermediate features* of various foundation models as training-free metrics for low-level perceptual similarity estimation. Our findings highlight that metrics using *intermediate features* outperform those using embeddings, with DINOv1 being the overall top performer. Future work will focus on fine-tuning these features on perceptual similarity datasets for improved results.

# 5. REFERENCES

[1] Olga Russakovsky, Jia Deng, Hao Su, et al., "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, pp. 211–252, 2015.

[2] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.

[3] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[4] Mingxing Tan and Quoc Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019, vol. 97, pp. 6105–6114.

[5] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 586–595.

[6] Manoj Kumar, Neil Houlsby, Nal Kalchbrenner, and Ekin Dogus Cubuk, "Do better imagenet classifiers assess perceptual similarity better?," *Transactions on Machine Learning Research*, 2022.

[7] Alec Radford, Jong Wook Kim, Chris Hallacy, et al., "Learning transferable visual models from natural language supervision," in *International conference on machine learning*, 2021, pp. 8748–8763.

[8] Mathilde Caron, Hugo Touvron, Ishan Misra, et al., "Emerging properties in self-supervised vision transformers," in *International Conference on Computer Vision*, 2021, pp. 9650–9660.

[9] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy, "Exploring clip for assessing the look and feel of images," in *AAAI Conference on Artificial Intelligence*, 2023, vol. 37, pp. 2555–2563.

[10] Stephanie Fu, Netanel Tamir, Shobhita Sundaram, et al., "Dreamsim: Learning new dimensions of human visual similarity using synthetic data," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[11] Francesco Croce, Christian Schlarmann, Naman Deep Singh, and Matthias Hein, "Adversarially robust CLIP models induce better (robust) perceptual metrics," in *ICML Wkshp. on Foundation Models in the Wild*, 2024.

[12] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P. Simoncelli, "Image quality assessment: Unifying structure and texture similarity," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[13] Abhijay Ghildyal and Feng Liu, "Shift-tolerant perceptual similarity metric," in *European Conference on Computer Vision*, 2022.

[14] M Oquab, T Darcet, T Moutakanni, et al., "DINOv2: Learning robust visual features without supervision," *Transactions on Machine Learning Research*, 2024.

[15] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[16] Zhou Wang, Eero P Simoncelli, and Alan C Bovik, "Multiscale structural similarity for image quality assessment," in *The Asilomar Conference on Signals, Systems & Computers*. 2003, vol. 2, pp. 1398–1402, IEEE.

[17] Lin Zhang, Lei Zhang, Xuanqin Mou, and David Zhang, "FSIM: a feature similarity index for image quality assessment," *IEEE Transactions on Image Processing*, vol. 20, no. 8, pp. 2378–2386, 2011.

[18] X Liao, X Wei, M Zhou, Z Li, and S Kwong, "Image quality assessment: Measuring perceptual degradation via distribution measures in deep feature spaces," *IEEE Transactions on Image Processing*, 2024.

[19] Pablo Hernández-Cámara, Jorge Vila-Tomás, Jesus Malo, and Valero Laparra, "Measuring human-CLIP alignment at different abstraction levels," in *ICLR Wkshp. on Representational Alignment*, 2024.

[20] Alexey Dosovitskiy, Lucas Beyer, A Kolesnikov, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *International Conference on Learning Representations*, 2021.

[21] J Gu, H Cai, H Chen, X Ye, J Ren, and C Dong, "Pipal: a large-scale image quality assessment dataset for perceptual image restoration," in *European Conference on Computer Vision*, 2020, pp. 633–651.

[22] Hanhe Lin, Vlad Hosu, and Dietmar Saupe, "Kadid-10k: A large-scale artificially distorted iqa database," in *International Conference on Quality of Multimedia Experience*, 2019, pp. 1–3.

[23] Abhijay Ghildyal and Feng Liu, "Attacking perceptual similarity metrics," *Transactions on Machine Learning Research*, 2023.

[24] Sara Ghazanfari, Siddharth Garg, Prashanth Krishnamurthy, Farshad Khorrami, and Alexandre Araujo, "R-LPIPS: An adversarially robust perceptual similarity metric," in *ICML Wkshp. on New Frontiers in Adversarial Machine Learning*, 2023.