

# Improving Shape Bias in Learnable Geometric Moment Representations

Anonymous WACV **Algorithms Track** submission

Paper ID \*\*\*\*\*

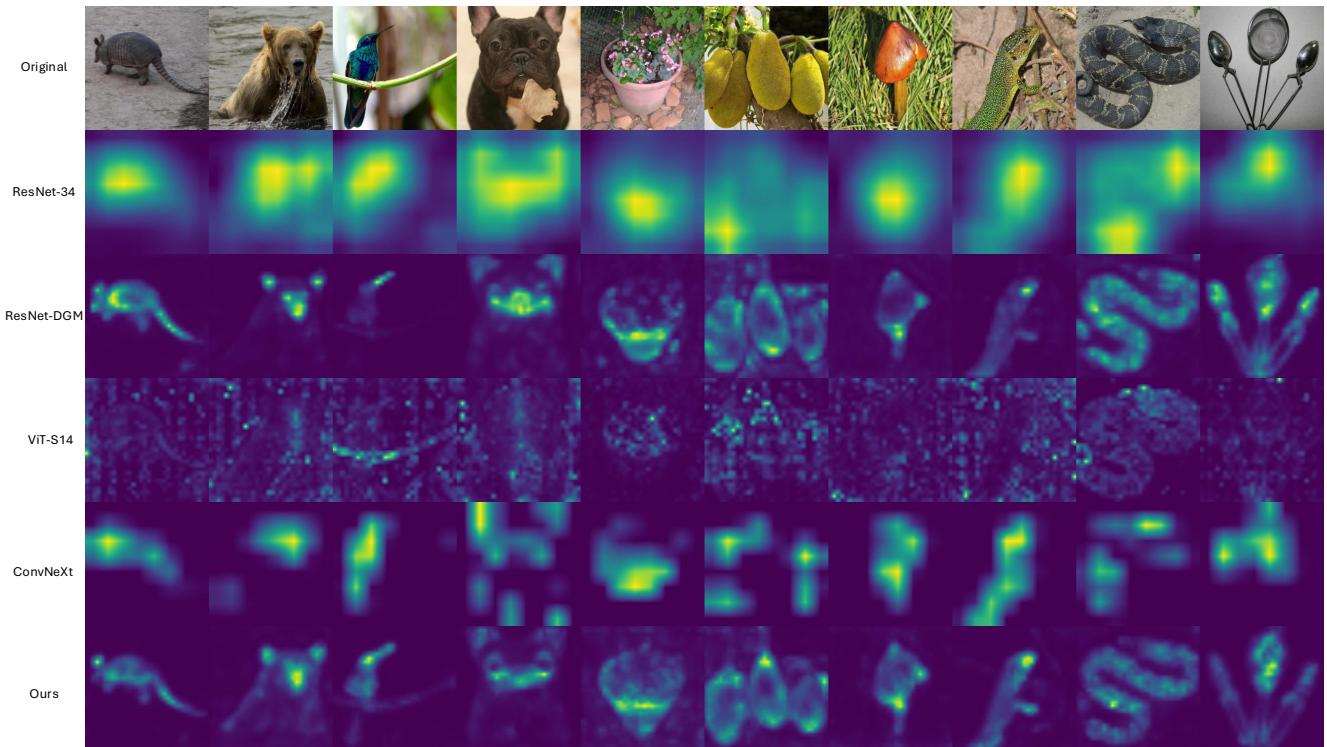


Figure 1. Feature visualization comparison across ResNet-34, ResNet-DGM, ViT-s14 (DINOv2), ConvNeXt-Base, and our method (ConvNeXt-DGM). Standard baselines (ResNet, ViT, ConvNeXt) tend to produce diffuse activations that fail to capture accurate object boundaries. In contrast, the DGM variations demonstrate superior shape alignment and localization.

## Abstract

001 *Deep Geometric Moments (DGMs) have been shown to en-*  
 002 *code shape-aware representations in image features. In*  
 003 *this work, we revisit the DGM framework through the lens*  
 004 *of ConvNeXt, a modern convolutional network (ConvNet)*  
 005 *architecture. By leveraging features extracted from Con-*  
 006 *vNeXt, we improve classification accuracy while further*  
 007 *strengthening the geometric shape-awareness of DGMs.*  
 008 *Our results demonstrate that features from modern Con-*  
 009 *vNet backbones serve as compatible stems for training Deep*  
 010 *Geometric Moments, and that the learned representations*  
 011 *remain tightly aligned with object geometry while exhibit-*

012 *ing robustness to visual perturbations. We quantitatively*  
 013 *characterize this shape awareness using geometric metrics*  
 014 *such as the Hausdorff distance and the Average Symmet-*  
 015 *ric Surface Distance (ASSD) complemented by Intersec-*  
 016 *tion over Union (IoU) to assess regional overlap. Furthermore,*  
 017 *we conduct an extensive analysis of the invariance of the*  
 018 *learned feature representations under diverse image per-*  
 019 *turbations, including changes in rotation, brightness, color,*  
 020 *and scale. We posit that these shape-aligned features offer*  
 021 *significant value not only for traditional computer vision*  
 022 *tasks, such as object detection and image segmentation, but*  
 023 *also for modern efficient training-free image and video edit-*

024 *ing methods.*

## 025 1. Introduction

026 Learning visual representations that are sensitive to object  
 027 shape while remaining robust to appearance variations is a  
 028 long-standing goal in computer vision. Shape-aware features  
 029 are particularly valuable for tasks that require precise  
 030 geometric reasoning, such as image segmentation, depth  
 031 estimation, and 3D understanding, where semantic consistency  
 032 must be maintained despite changes in color, illumination,  
 033 viewpoint, or scale. [10, 18]. Beyond traditional  
 034 recognition pipelines, recent works have further demonstrated  
 035 that intermediate feature representations such as attention  
 036 maps or activation responses can be directly manipulated  
 037 for image editing, video guidance, and controllable generation,  
 038 placing additional importance on the interpretability and geometric  
 039 alignment of learned features.

040 Deep Geometric Moments (DGM) were introduced as a  
 041 mechanism for encouraging shape-awareness in deep neural  
 042 networks by explicitly modeling geometric properties  
 043 of feature responses [17]. The original study demonstrated  
 044 that DGM-based representations tend to align more closely  
 045 with object structure, enabling improved spatial correspondences  
 046 and robustness to certain transformations. However,  
 047 this foundational work relied exclusively on ResNet [8, 13],  
 048 a conventional backbone that predates more modern architectures.  
 049 Consequently, the interaction between these geometric  
 050 properties and advanced network designs remains  
 051 underexplored, and it is unclear whether the observed  
 052 benefits persist when utilizing stronger feature extractors.

053 In parallel, the landscape of visual representation learning  
 054 has evolved significantly with the introduction of Vision  
 055 Transformers (ViTs) [2, 4, 7, 11, 14, 19], which have  
 056 demonstrated strong scalability and performance across a  
 057 wide range of vision tasks. While ViTs excel at modeling  
 058 long-range dependencies, ConvNets [6] retain important  
 059 advantages due to strong inductive biases such as locality and  
 060 translation equivariance. [6] These properties often lead to  
 061 better generalization in limited-data regimes and reduced  
 062 computational overhead. ConvNeXt revisits these design  
 063 principles through Transformer-inspired training and architectural  
 064 choices, demonstrating that ConvNets can achieve  
 065 performance competitive with ViTs while maintaining their  
 066 inherent inductive strengths.

067 Motivated by these developments, we revisit Deep Geometric  
 068 Moments using ConvNeXt [12, 20] as the feature  
 069 backbone. Our objective is twofold: first, to assess whether  
 070 DGM remains effective when paired with a high-capacity  
 071 ConvNet architecture; and second, to conduct a systematic  
 072 analysis of the shape-aware features learned in this setting.  
 073 By integrating DGM with ConvNeXt, we aim to investigate  
 074 the specific contribution of the geometric moment for-

mulation relative to the underlying feature extractor and to gain further insight into the properties of the representations DGM induces.

To this end, we conduct a comprehensive empirical study of ConvNeXt-based DGM models. We quantitatively evaluate shape awareness by measuring the alignment between learned feature responses and object masks using metrics such as mean Intersection-over-Union (mIoU), Hausdorff distance, and Average Symmetric Surface Distance (ASSD). To further probe the robustness and invariance of the representations, we perform controlled perturbation experiments spanning variations in color, brightness, scale, and orientation. These analyses provide insight into the stability of DGM features against common appearance changes and their ability to preserve geometric consistency under transformations.

Our results indicate that ConvNeXt serves as a compatible feature backbone for Deep Geometric Moments, yielding representations that align more closely with object geometry than the baseline model. Furthermore, the observed robustness to visual perturbations demonstrates that shape-aware mechanisms complement advanced ConvNet architectures by effectively enhancing geometric consistency.

In summary, this work provides a re-examination of Deep Geometric Moments by integrating them with advanced convolutional backbones and conducting a detailed representational analysis. Our findings confirm that DGM remains a powerful tool for inducing shape bias, bridging the gap between strong feature extractors and geometric interpretability.

## 2. Preliminary

### 2.1. Deep Geometric Moments

Geometric moments have long been used in computer vision as a compact and informative representation of shape, capturing spatial distributions of mass while facilitating invariance to geometric transformations. Classical moment formulations encode object geometry through weighted integrals of pixel coordinates, enabling descriptors that are sensitive to shape while being robust to appearance variations such as color or texture. Formally, consider a feature map  $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$  extracted from a backbone network. In a standard formulation, the general geometric moment of order  $(p, q)$  for a given channel is defined as:

$$M_{pq} = \sum_{x=1}^W \sum_{y=1}^H x^p y^q \phi(F_{x,y}) \quad (1)$$

where  $(x, y)$  represents spatial coordinates and  $\phi(\cdot)$  denotes a mapping function applied to the feature activations. By selecting specific moment orders and normalization schemes, these moments can explicitly encode geometric properties

123 such as orientation, scale, and higher-order shape characteristics.  
 124 However, determining the optimal set of orders ( $p, q$ ) prior is non-trivial.  
 125 Deep Geometric Moments (DGM) address this by embedding geometric structure directly into  
 126 the representation learning process. Rather than manually  
 127 prescribing fixed integer orders, our method fixes the dimensionality of the representation but learns the underlying  
 128 basis functions end-to-end. Consequently, we simplify the  
 129 notation by dropping the order subscripts  $p, q$  and indexing  
 130 moments solely by the feature channel  $c$ . In this formulation,  
 131 the DGM module projects extracted image features onto a learnable coordinate basis. The moment  $m^c$  for the  
 132  $c$ -th channel is defined as:  
 133

$$136 \quad m^c = \frac{1}{N \times N} \sum_{x=1}^N \sum_{y=1}^N g^c(x, y) f^c(x, y) \quad (2)$$

137 where  $N \times N$  denotes the spatial resolution of the feature  
 138 map,  $f^c(x, y)$  represents the feature activation at location  
 139  $(x, y)$  extracted by the backbone CNN, and  $g^c(x, y)$  is a  
 140 learnable 2D polynomial function serving as the coordinate  
 141 basis. To further account for variations in object location,  
 142 size, and pose, we explicitly learn affine parameters that dy-  
 143 namically deform the coordinate grid  $g^c(x, y)$  during the  
 144 moment computation, thereby ensuring geometric invari-  
 145 ance. For our specific implementation, we set the spatial  
 146 dimension  $N = 32$  (yielding a  $32 \times 32$  grid) and maintain  
 147 a fixed channel dimensionality of  $C = 256$ .

## 148 2.2. ConvNeXt

149 ConvNeXt revisits the standard convolutional design by  
 150 systematically integrating architectural patterns that con-  
 151 tribute to the success of Vision Transformers [15]. As il-  
 152 lustrated in Figure 2, the core building block shifts from  
 153 the traditional ResNet bottleneck to an inverted bottleneck  
 154 structure. This design decouples spatial mixing from chan-  
 155 nel mixing, employing large-kernel depthwise convolutions  
 156 (increasing from  $3 \times 3$  to  $7 \times 7$ ) to significantly expand the  
 157 effective receptive field. Additionally, the network adopts  
 158 Layer Normalization [1] and GELU activations [9], re-  
 159 placing the conventional Batch Normalization and ReLU  
 160 activation layers found in ResNet. These adjustments  
 161 align the network’s processing dynamics with Transformer-  
 162 style blocks while retaining the computational efficiency of  
 163 sliding-window convolutions. In the context of this study,  
 164 ConvNeXt offers a distinct benefit by yielding discriminative  
 165 feature representations that preserve the inherent spatial  
 166 structure of standard convolutional networks. Unlike ViTs,  
 167 which process token sequences that may lose local adjac-  
 168 ency information, ConvNeXt outputs dense, grid-aligned  
 169 activation maps. This preservation of spatial layout is cru-  
 170 cial, as DGM computes geometric moments directly from  
 171 the feature map coordinates to enforce shape constraints.

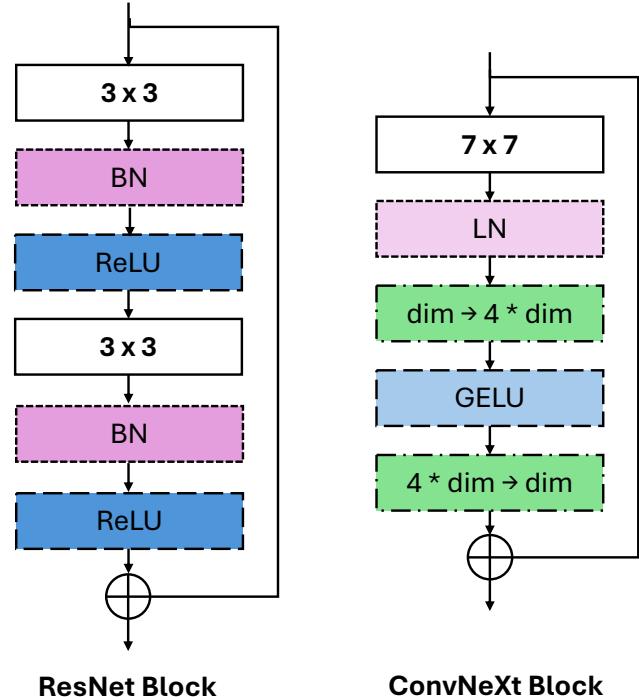


Figure 2. Composition of ResNet and ConvNeXt blocks. Red depicts normalization layers, Green is linear layer with input and output dimensions, Blue is activation functions.

## 172 2.3. Combining DGM and ConvNeXt

We combine these frameworks to evaluate the behavior of Deep Geometric Moments within a more powerful architectural setting. ConvNeXt provides a strong, spatially grounded feature substrate that is technically compatible with the moment integration process of DGM. By employing this advanced backbone, we aim to investigate the interaction between explicit geometric constraints and stronger feature representations. This approach allows us to verify how DGM adapts to improved feature extractors and whether it continues to enhance representation alignment beyond what is achieved by the backbone alone.

## 184 3. Experimental Setup

In this section, we discuss how our experiments are designed. We explain the evaluation metrics used for our analysis, following the model architecture and its variations.

### 188 3.1. Metrics

Our primary objective is to evaluate the capacity of learned features to preserve object geometry and shape structure. While the original DGM study relied primarily on qualitative visualizations to demonstrate shape awareness, we extend this analysis by establishing a rigorous quantitative

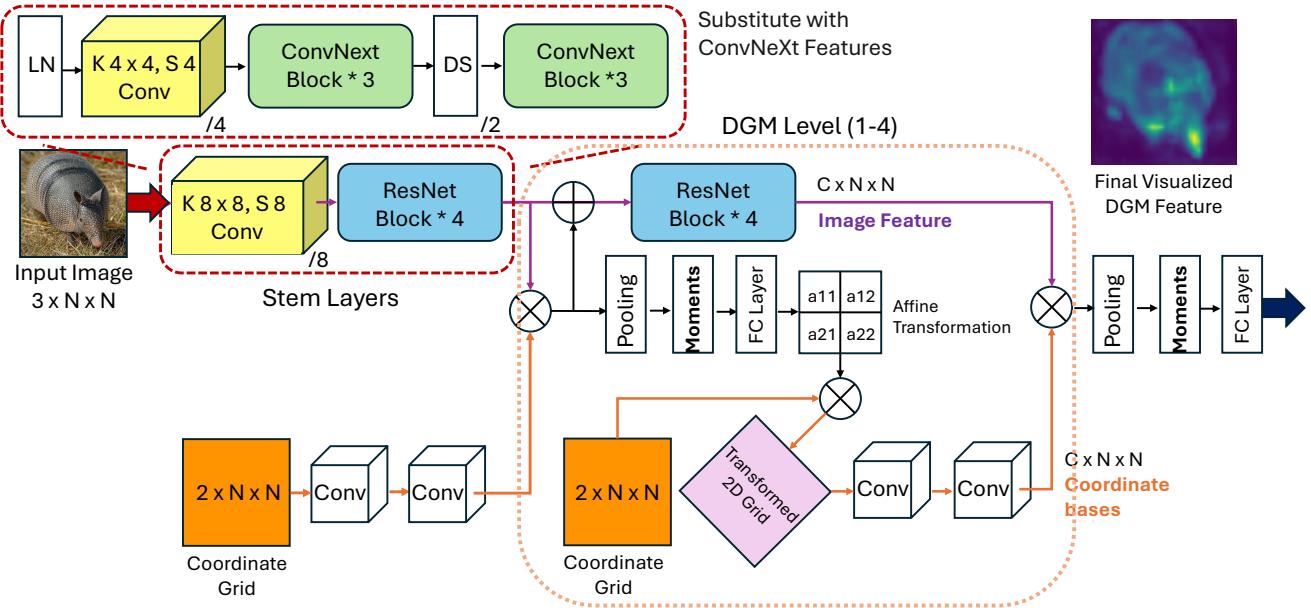


Figure 3. DGM model architecture. DGM is learned by dual path processing; upper image features, lower grid bases. We substitute the stem layers for the image features with ConvNeXt.

194 framework. We adopt the following three metrics to  
 195 systematically measure the geometric alignment of each  
 196 model's image features.  
 197

### Intersection over Union (IoU)

The Intersection over Union (IoU), also known as the Jaccard Index, is a standard metric for quantifying the spatial overlap between two distinct regions. It is calculated by dividing the area of the intersection by the area of the union of the two sets. Formally, given a predicted binary mask  $A$  (derived from the feature response) and the ground truth object mask  $B$ , the IoU is defined as:

$$\text{IoU}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|} \quad (3)$$

where  $|\cdot|$  denotes the cardinality (area) of the set. The resulting score ranges from 0 to 1, where 1 indicates a perfect spatial match. Unlike detection tasks that often apply IoU to bounding boxes, we compute the mean IoU (mIoU) directly on pixel-wise segmentation masks to accurately capture the alignment of the feature activation with the interior shape of the object.

### Hausdorff Distance

The Hausdorff Distance (HD) measures the maximum distance from a point in one set to the nearest point in another set, effectively quantifying the worst-case mismatch between two boundaries. Formally, given two sets of boundary points  $A$  and  $B$  (representing the predicted and

ground-truth contours), the Hausdorff distance is defined as:

$$H(A, B) = \max \left\{ \sup_{a \in A} \inf_{b \in B} d(a, b), \sup_{b \in B} \inf_{a \in A} d(a, b) \right\} \quad (4)$$

where  $d(a, b)$  denotes the Euclidean distance between points  $a$  and  $b$ . While IoU focuses on the volumetric overlap of the object's interior, the Hausdorff distance supplements this by explicitly evaluating the geometric precision of the object boundaries.

However, because the standard HD is determined by the single largest deviation, it is highly sensitive to outliers. To mitigate this, we adopt the 95th percentile Hausdorff Distance (HD95), which excludes the most extreme anomalies to provide a more robust metric of shape consistency.

### Average Symmetric Surface Distance (ASSD)

The Average Symmetric Surface Distance (ASSD) is a boundary-based metric that evaluates the global alignment between two surfaces. ASSD computes the average distance between all points on the predicted and ground-truth boundaries. It is defined as:

$$\text{ASSD}(A, B) = \frac{\sum_{a \in S(A)} d(a, S(B)) + \sum_{b \in S(B)} d(b, S(A))}{|S(A)| + |S(B)|} \quad (5)$$

where  $S(A)$  and  $S(B)$  represent the set of surface points (contour pixels) for the ground truth  $A$  and prediction  $B$ , respectively. The term  $d(x, S(Y)) = \min_{y \in S(Y)} \|x - y\|$  denotes the shortest Euclidean distance from a point  $x$  to the surface set  $S(Y)$ . Crucially, this metric incorporates

246 two directional components to ensure symmetry. By aver-  
 247 eraging these bidirectional distances, ASSD penalizes both  
 248 error types equally, providing a holistic measure of bound-  
 249 ary adherence that complements the worst-case sensitivity  
 250 of the Hausdorff distance.

### 251 3.2. Model

252 We detail our model architecture in Figure 3. The origi-  
 253 nal DGM implementation utilized a convolutional stem to  
 254 downsample the input by a factor of 8, followed by a se-  
 255 quence of four ResNet blocks. In this work, we replace  
 256 this backbone with a pretrained ConvNeXt encoder. Specif-  
 257 ically, the image is processed through two stages, each com-  
 258 prising three ConvNeXt blocks, configured to match the  
 259 spatial downsampling ratios and output dimensions of the  
 260 original architecture. These extracted features subsequently  
 261 serve as the input for the DGM module, where the geo-  
 262 metric learning formulation remains identical to the origi-  
 263 nal method. All models are trained for 100 epochs using a  
 264 cosine decay scheduler with an initial learning rate of 0.1.

## 265 4. Experimental Results

266 In this section, we outline our experimental analysis and re-  
 267 sults, with a primary focus on the shape alignment capabili-  
 268 ties of DGM features and their robustness to perturbations.

Table 1. ImageNet-1K classification accuracy using different ConvNeXt models as stem feature extractors. We newly train ResNet34-DGM following [17]. The first two columns denote the number of channels after each stage.

Model	Stage 1	Stage 2	Accuracy (%)	Params (M)
Tiny	96	192	76.26 (+1.34)	20.86
Base	128	256	76.76 (+1.84)	21.82
Large	192	384	77.25 (+2.33)	24.51
X-Large	256	512	77.36 (+2.44)	28.25
ResNet34-DGM	-	-	74.93	24.35

269

### 270 4.1. Model Variations

271 We experiment with ConvNeXt model variants ranging  
 272 from Tiny to X-Large. The architecture is structured in four  
 273 stages, beginning with a stem that reduces the input resolu-  
 274 tion by a factor of four. Between subsequent stages, down-  
 275 sampling layers further reduce the spatial dimension by  
 276 half. To accommodate our DGM training pipeline, which  
 277 operates on a fixed  $32 \times 32$  grid, we fix the input resolution  
 278 to  $256 \times 256$  and extract features from the output of Stage 2.  
 279 With this setup, the cumulative downsampling transforms  
 280 the  $256 \times 256$  input precisely into the required  $32 \times 32$  fea-  
 281 ture map. The model variants differ only in channel width  
 282 at this extraction point; consequently, we omit ConvNeXt-  
 283 Small, as it is structurally identical to the Tiny variation up

284 to Stage 2. As summarized in Table 1, we trained four DGM  
 285 variants utilizing different ConvNeXt backbones. For com-  
 286 parison, we also retrained the ResNet34-DGM baseline [17]  
 287 under the same settings. All models were trained on the  
 288 ImageNet-1K classification task [5]. We observed that  
 289 every ConvNeXt-based variant outperformed the ResNet-  
 290 DGM baseline, with performance consistently improving as  
 291 the model size increased. However, since our primary fo-  
 292 cus is on geometric analysis rather than maximizing down-  
 293 stream classification accuracy, we utilize the ConvNeXt-  
 294 Base variant for the detailed experiments presented in the  
 remainder of this paper.

Table 2. Evaluation of geometry-shape alignment metrics on various feature extractors. GD denotes GradCAM. (\*) Denotes an outlier: visual inspection confirms that ResNet-GD produces over-expanded masks, inflating mIoU despite poor shape fidelity. Consequently, we adopt **ConvNext-Base-DGM** as our experimented model for the remainder of this study.

Model	mIoU ( $\uparrow$ )	mHD95 ( $\downarrow$ )	mASSD ( $\downarrow$ )
ResNet-GD	<b>0.4635*</b>	128.73	58.97
ResNet-DGM	0.3282	133.91	42.41
ConvNext-GD	0.3158	131.02	46.38
ConvNext-Tiny-DGM	<b>0.4118</b>	128.61	41.66
<b>ConvNext-Base-DGM</b>	0.4003	<b>126.70</b>	<b>40.09</b>
ConvNext-Large-DGM	0.3972	127.33	40.57
DINOv2 (ViT-S/14)	0.1975	146.02	48.64
DINOv2 (ViT-L/14)	0.2223	141.02	46.74

295

### 296 4.2. Shape Alignment Evaluation

297 We evaluate the geometric shape alignment of the learned  
 298 feature representations. As detailed in Section 3.1, we quan-  
 299 tify this alignment by comparing generated feature masks  
 300 against ground-truth object masks using three metrics: Inter-  
 301 section over Union (IoU), 95th percentile Hausdorff Dis-  
 302 tance (HD95), and Average Symmetric Surface Distance  
 303 (ASSD). Feature masks are min-max normalized and bina-  
 304 rized using a fixed empirical threshold of 0.15. For baseline  
 305 comparisons, we derive masks for the standard ConvNeXt-  
 306 Base using GradCAM [16], and for DINO models [3] by  
 307 averaging attention maps across all heads. Through exten-  
 308 sive analysis, we observe that region-based metrics like IoU  
 309 can be misleading when applied to coarse feature extrac-  
 310 tors (e.g., standard ResNet), which often produce diffuse,  
 311 over-expanded masks. To address this limitation, we incor-  
 312 porate boundary-aware metrics to ensure contour precision.  
 313 Among these, we find ASSD to be particularly robust. As  
 314 shown in Table 2, our ConvNeXt-Base DGM significantly  
 315 outperforms the baselines across all three metrics.

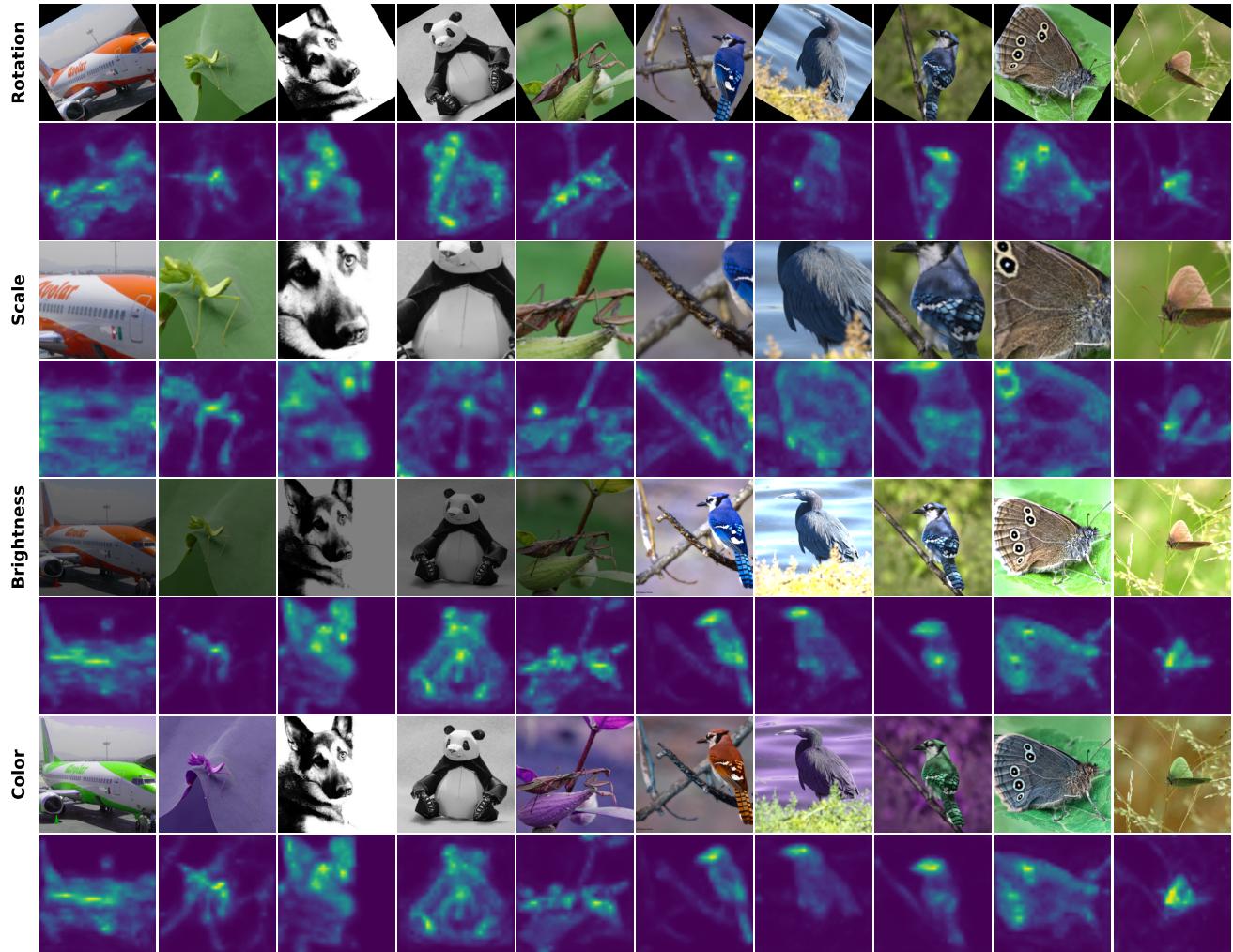


Figure 4. Qualitative comparison of model robustness against geometric (Angle, Scale) and photometric (Brightness, Color) perturbations. Our model maintains consistent shape alignment and structural integrity.

316

### 4.3. Invariance of Shape Alignment

317 To highlight the invariance of the learned geometric  
 318 moments, we conducted experiments using diverse perturba-  
 319 tions, including rotation, photometric shifts, and scale varia-  
 320 tions. For rotation, we evaluated six distinct angles rang-  
 321 ing from  $45^\circ$  to  $240^\circ$ . For color, we perturbed specific RGB  
 322 channel intensities (increase/decrease) and performed chan-  
 323 nel swapping. For scale, we tested performance on enlarged  
 324 inputs. Table 3 demonstrates that while standard ConvNet  
 325 features (ConvNeXt baseline) become inconsistent under  
 326 these transformations, our model maintains high structural  
 327 consistency. Figure 4 provides qualitative examples of this  
 328 stability. Specifically, we observe that the baseline often  
 329 suffers from catastrophic degradation when texture statis-  
 330 tics shift, suggesting an over-reliance on local appearance.  
 331 In contrast, our method’s performance remains stable, con-

firmed that the learned geometric moments effectively cap-  
 332 ture intrinsic shape properties that are independent of pho-  
 333 tometric or viewpoint variations.

### 4.4. Performance of Different Shape Categories

We conduct an ablation study to analyze how geometric  
 336 alignment varies across different semantic categories. We  
 337 selected 16 ImageNet classes that exhibit a wide range of  
 338 performance gaps between the baseline and our method. As  
 339 shown in Figure 5, a compelling trend emerges: our model  
 340 achieves significantly higher gains on animate objects (e.g.,  
 341 animals) compared to inanimate objects, even for species  
 342 with complex, diverse morphologies such as butterflies and  
 343 crocodiles. Taken together with Figure 6, these results es-  
 344 tablish that our model offers superior shape preservation  
 345 and alignment generality.

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

Table 3. Robustness evaluation under geometric (Angle, Scale) and photometric (Color Perturbation, Color Swap) transformations. We compare the stability of the default ConvNext-Base model versus ConvNext-DGM across mIoU, mHD95, and mASSD metrics.

Angle (°)	mIoU (↑)		mHD95 (↓)		mASSD (↓)	
	Base	Ours	Base	Ours	Base	Ours
45	0.3558	0.4364	128.32	123.08	44.72	37.68
60	0.3478	0.4468	131.60	125.65	46.25	38.70
90	0.3434	0.4751	137.32	133.53	46.58	40.05
120	0.3448	0.4413	132.31	127.98	46.35	39.59
180	0.3465	0.4734	134.27	131.31	45.95	39.77
240	0.3633	0.4419	125.31	127.66	43.94	39.57

Color Perturb	mIoU (↑)		mHD (↓)		mASSD (↓)	
	Base	Ours	Base	Ours	Base	Ours
R boost	0.3502	0.4730	133.46	127.10	45.77	38.09
R reduce	0.3532	0.4809	132.56	123.28	45.48	36.82
G boost	0.3506	0.4666	133.33	128.12	45.92	38.80
G reduce	0.3471	0.4667	134.17	127.22	46.03	38.19
B boost	0.3530	0.4811	132.49	124.82	45.27	37.31
B reduce	0.3532	0.4727	131.79	123.96	45.15	37.03

Color Swap	mIoU (↑)		mHD (↓)		mASSD (↓)	
	Base	Ours	Base	Ours	Base	Ours
R ↔ G	0.4669	0.3428	129.04	137.01	38.94	47.22
G ↔ B	0.4671	0.3455	129.11	135.47	38.82	46.69
B ↔ R	0.4730	0.3441	126.38	135.83	38.05	46.71

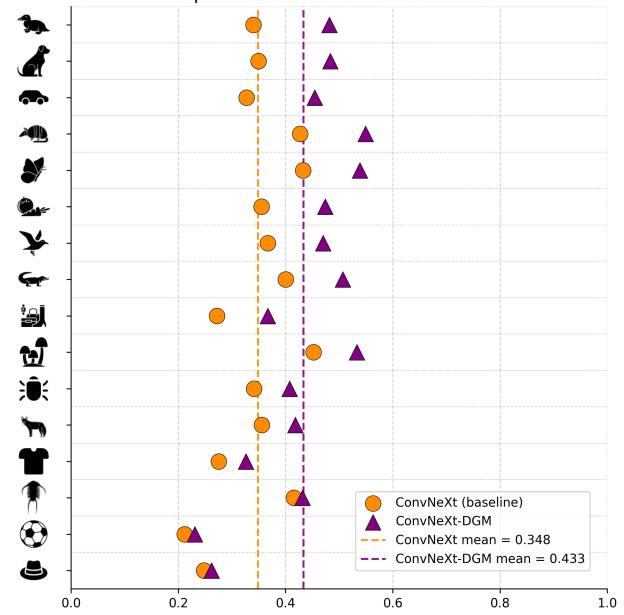
  

Scale	mIoU (↑)		mHD (↓)		mASSD (↓)	
	Base	Ours	Base	Ours	Base	Ours
1.25	0.3351	0.4801	147.34	136.62	48.09	40.39
1.5	0.3252	0.4753	156.57	144.34	50.48	42.48
2.0	0.3066	0.4614	171.47	155.91	54.59	45.62

## 347 5. Conclusion

348 By revisiting Deep Geometric Moments through the lens of  
 349 modern convnet architectures like ConvNeXt, we have syn-  
 350 thesized the inductive bias of geometric theory with rep-  
 351 resentational power of deep learning. Our work quantita-  
 352 tively validates that DGMs learn representations that are  
 353 both tightly aligned with object geometry through exten-  
 354 sive evaluation using geometric metrics such as Hausdorff  
 355 distance and ASSD and shows robustness to diverse visual  
 356 perturbations including rotation, scaling, and photometric  
 357 shifts. This establishes a distinct “shape bias” that contrasts  
 358 sharply with the texture-reliance often observed in standard  
 359 convnets, suggesting that DGMs can offer a robust foun-  
 360 dation for the next generation of geometry-sensitive down-  
 361 stream applications, providing critical capabilities for tasks  
 362 ranging from controllable image manipulation and video  
 363 editing workflows.

Mean IoU per class — ConvNeXt vs ConvNeXt-DGM



Mean ASSD per class — ConvNeXt vs ConvNeXt-DGM

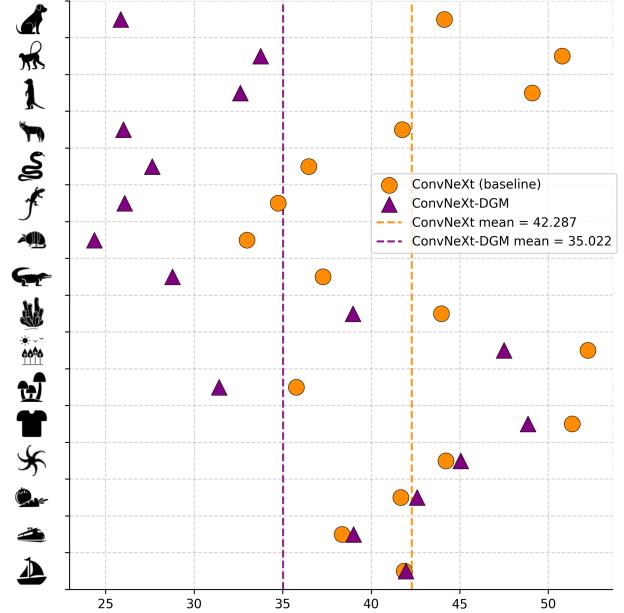


Figure 5. Per-class performance comparison on 16 curated categories from the ImageNet dataset. The chart illustrates the IoU and ASSD scores for the Base model versus DGM (Ours). While improvement margins vary depending on class complexity, DGM demonstrates consistent gains across diverse semantic categories.

## 364 6. Limitation and Future Works

In this study, we constrain our models to approximately 25M parameters to maintain computational efficiency and prioritize algorithmic interpretability. We hypothesize that scaling the input grid size offers significant potential, as

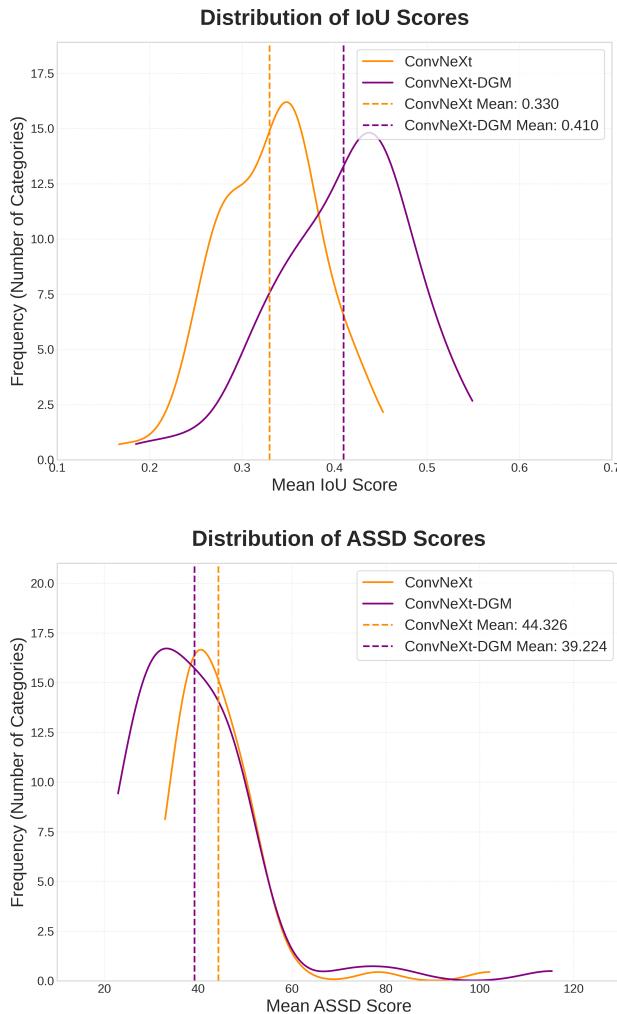


Figure 6. Distribution of IoU and ASSD scores on the ImageNet dataset. Compared to the baseline, DGM (Ours) shows a distinct distributional shift towards higher IoU and lower ASSD values, indicating superior generalization and more consistent geometric alignment across the dataset.

369 finer resolutions theoretically allow for the calculation of  
 370 more precise geometric moments—a direction we reserve  
 371 for future exploration. Additionally, we do not conduct a  
 372 direct comparison with massive self-supervised ViTs (e.g.,  
 373 DINOv3). While we acknowledge the exceptional repre-  
 374 sentational power of such foundation models, our results  
 375 demonstrate that robust shape awareness is attainable in sig-  
 376 nificantly more compact architectures without relying on  
 377 large-scale pre-training.

## 378 References

- 379 [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hin-  
 380 ton. Layer normalization. *arXiv preprint arXiv:1607.06450*,  
 381 2016. 3

- [2] Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are  
 382 transformers more robust than cnns? *Advances in neural*  
 383 *information processing systems*, 34:26831–26843, 2021. 2
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou,  
 384 Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerg-  
 385 ing properties in self-supervised vision transformers. In *Pro-  
 386 ceedings of the IEEE/CVF international conference on com-  
 387 puter vision*, pages 9650–9660, 2021. 5
- [4] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr  
 388 Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter  
 389 Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdul-  
 390 mohsin, et al. Scaling vision transformers to 22 billion pa-  
 391 rameters. In *International conference on machine learning*,  
 392 pages 7480–7512. PMLR, 2023. 2
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,  
 393 and Li Fei-Fei. Imagenet: A large-scale hierarchical image  
 394 database. In *2009 IEEE conference on computer vision and*  
 395 *pattern recognition*, pages 248–255. Ieee, 2009. 5
- [6] Josip Djolonga, Jessica Yung, Michael Tschannen, Rob  
 396 Romijnders, Lucas Beyer, Alexander Kolesnikov, Joan  
 397 Puigcerver, Matthias Minderer, Alexander D’Amour, Dan  
 398 Moldovan, et al. On robustness and transferability of convo-  
 399 lutional neural networks. In *Proceedings of the IEEE/CVF*  
 400 *Conference on Computer Vision and Pattern Recognition*,  
 401 pages 16458–16468, 2021. 2
- [7] Alexey Dosovitskiy. An image is worth 16x16 words:  
 402 Transformers for image recognition at scale. *arXiv preprint*  
*arXiv:2010.11929*, 2020. 2
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.  
 403 Deep residual learning for image recognition. In *Pro-  
 404 ceedings of the IEEE conference on computer vision and pattern*  
*405 recognition*, pages 770–778, 2016. 2
- [9] D Hendrycks. Gaussian error linear units (gelus). *arXiv*  
*406 preprint arXiv:1606.08415*, 2016. 3
- [10] Sangmin Jung, Utkarsh Nath, Yezhou Yang, Giulia Pedrielli,  
 407 Joydeep Biswas, Amy Zhang, Hassan Ghasemzadeh, and  
 408 Pavan Turaga. Guiding diffusion with deep geometric mo-  
 409 ments: Balancing fidelity and variation. *arXiv preprint*  
*arXiv:2505.12486*, 2025. 2
- [11] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng  
 410 Zhang, Stephen Lin, and Baining Guo. Swin transformer:  
 411 Hierarchical vision transformer using shifted windows. In  
 412 *Proceedings of the IEEE/CVF international conference on*  
*413 computer vision*, pages 10012–10022, 2021. 2
- [12] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feicht-  
 414 enhofer, Trevor Darrell, and Saining Xie. A convnet for the  
 415 2020s. In *Proceedings of the IEEE/CVF conference on com-  
 416 puter vision and pattern recognition*, pages 11976–11986,  
 417 2022. 2
- [13] Thao Nguyen, Maithra Raghu, and Simon Kornblith. Do  
 418 wide and deep networks learn the same things? uncover-  
 419 ing how neural network representations vary with width and  
 420 depth. *arXiv preprint arXiv:2010.15327*, 2020. 2
- [14] Maithra Raghu, Thomas Unterthiner, Simon Kornblith,  
 421 Chiyuan Zhang, and Alexey Dosovitskiy. Do vision trans-  
 422 formers see like convolutional neural networks? *Advances*  
 423 *in neural information processing systems*, 34:12116–12128,  
 424 2021. 2

- 440 [15] Nataniel Ruiz, Sarah Bargal, Cihang Xie, Kate Saenko, and  
 441 Stan Sclaroff. Finding differences between transformers and  
 442 convnets using counterfactual simulation testing. *Advances*  
 443 *in Neural Information Processing Systems*, 35:14403–14418,  
 444 2022. 3
- 445 [16] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das,  
 446 Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra.  
 447 Grad-cam: Visual explanations from deep networks via  
 448 gradient-based localization. In *Proceedings of the IEEE in-*  
 449 *ternational conference on computer vision*, pages 618–626,  
 450 2017. 5
- 451 [17] Rajhans Singh, Ankita Shukla, and Pavan Turaga. Improving  
 452 shape awareness and interpretability in deep networks using  
 453 geometric moments. In *Proceedings of the IEEE/CVF Con-*  
 454 *ference on Computer Vision and Pattern Recognition*, pages  
 455 4159–4168, 2023. 2, 5
- 456 [18] Rajhans Singh, Ankita Shukla, and Pavan Turaga. Poly-  
 457 nomial implicit neural representations for large diverse  
 458 datasets. In *Proceedings of the IEEE/CVF Conference*  
 459 *on Computer Vision and Pattern Recognition*, pages 2041–  
 460 2051, 2023. 2
- 461 [19] Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L  
 462 Griffiths. Are convolutional neural networks or transformers  
 463 more like human vision? *arXiv preprint arXiv:2105.07197*,  
 464 2021. 2
- 465 [20] Kirill Vishniakov, Zhiqiang Shen, and Zhuang Liu. Convnet  
 466 vs transformer, supervised vs clip: Beyond imagenet accu-  
 467 racy. *arXiv preprint arXiv:2311.09215*, 2023. 2