
Evaluating Foundation Models’ 3D Understanding Through Multi-View Correspondence Analysis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper extends the Hummingbird framework with the Multi-View ImageNet
2 (MVIImgNet) dataset to evaluate how foundation model image encoders handle
3 in-context object segmentation under unseen camera angles. We group MVIImgNet
4 object views and construct memory banks from selected viewpoints. We assess
5 generalization by evaluating performance on held-out angles. Across six pretrained
6 Vision Transformer (ViT) models—CLIP, DINO, DINOV2, SigLIP2, C-RADIOv2,
7 and TIPS—we find that DINO-based models perform best: DINO leads when more
8 context viewpoints are available, while DINOV2 is strongest with fewer reference
9 views. These results highlight the benefits of contrastive pretraining for robust
10 performance across large viewpoint shifts.

11 1 Introduction

12 In recent years, we have seen the rise of foundation models [2, 20, 11, 23]. These large-scale pre-
13 trained models support many tasks, including visual understanding, which enables direct comparison
14 of their encoders.

15 Despite their impressive capabilities, state-of-the-art ViTs can suffer catastrophic failures when
16 objects are viewed from unusual angles, as demonstrated by recent angle-sensitivity studies [15].
17 However, most vision benchmarks emphasize single-view tasks or dense synthesis [10], leaving
18 segmentation robustness under camera rotations comparatively underexplored.

19 To address viewpoint variability, recent work has explored in-context learning (ICL) as a means to
20 adapt models to unseen views without retraining. ICL is the ability of a model to perform new tasks
21 by conditioning on a prompt [22]. An example is Hummingbird [1], a memory-augmented ViT for
22 ICL. Its encoder extracts dense image features and projects them into key–value pairs stored in a
23 dynamic memory. At query time, cross-attention over the keys assigns soft nearest-neighbor (NN)
24 weights to aggregate values and predict novel views in-context. This modular design permits the use
25 of any encoder.

26 In this work, we evaluate foundation model encoders for in-context object segmentation under novel
27 camera angles. We extend the Hummingbird framework with the MVIImgNet dataset [24], grouping
28 object views into angular bins. By constructing dynamic memory banks from selected viewpoints
29 and evaluating segmentation performance on held-out angles, we assess how well different encoders
30 generalize across viewpoint shifts. For clarity, some terms are followed by their variable name to
31 maintain cohesion with our codebase.

32 **2 Related work**

33 Recent developments in view generation have significantly improved the understanding of 3D scene
34 structure in computer vision, setting new benchmarks along the way. A notable example is NeRF [10],
35 which achieves state-of-the-art novel view synthesis by optimizing a continuous volumetric function
36 from sparse input images. While NeRF addresses generation, other studies investigate recognition
37 in 3D understanding. Ruan et al. [15] demonstrate that modern recognition models (ResNets [5],
38 ViTs [4], Swin [8], and Masked Autoencoders (MAEs) [6]) remain sensitive to viewpoint shifts in 3D
39 scenes. With adversarial training, they improve invariance to 3D viewpoint changes beyond standard
40 rotation-based augmentations.

41 In parallel, Shifman and Weiss [19] show that state-of-the-art encoders such as CLIP and DINOv2
42 change predictions under a single-pixel shift in up to 40% of cases, despite extensive augmentation
43 during training. Although their analysis centers on 2D translation invariance, the findings highlight
44 that robustness to geometric changes, whether in 2D or 3D, remains limited in current vision encoders.

45 Beyond model architectures, multi-view datasets such as MVIImgNet and PASCAL3D+ offer rich
46 multi-view annotations that support the evaluation of cross-view generalization in object recognition.
47 The rise of self-supervised Transformers has also opened new possibilities for downstream tasks such
48 as segmentation. In particular, DINOv2 [3] produces robust embeddings that transfer effectively to
49 prediction tasks, despite training without labels.

50 Nevertheless, a gap remains across these benchmarks and models. Most prior work emphasizes
51 detection, classification, or view-angle estimation, but not the measurement of 3D contextual under-
52 standing through multi-view segmentation. To our knowledge, no benchmark currently evaluates
53 a model’s ability to generalize segmentation across viewpoint shifts using dynamic memory. Our
54 proposed evaluation method, therefore, represents a novel contribution to this area.

55 **3 Methodology**

56 We evaluate the view generalization ability of frozen ViT models in semantic segmentation using a
57 non-parametric, retrieval-based framework. Our approach builds on Hummingbird [1], which applies
58 in-context learning (ICL) to vision tasks. Hummingbird supports evaluation of spatial perception and
59 semantic understanding, but does not analyze performance variation under viewpoint changes. We
60 address this by introducing a viewpoint binning protocol and a curated dataset subset for cross-view
61 robustness analysis.

62 **3.1 Encoders**

63 We evaluate six pretrained ViT models: CLIP [13], SigLIP2 [21], DINO [3], DINOv2 [12], C-
64 RADIOv2 [14, 7], and TIPS [9]. All models use a ViT-B/16 backbone, except DINOv2, which is
65 released only as ViT-B/14 and reported to outperform its ViT-B/16 counterpart. All encoders are
66 publicly released under open-source licenses (see Appendix C).

67 **3.2 Inference Pipeline**

68 We follow the Hummingbird [1] inference pipeline, where patch-level features from a frozen ViT are
69 stored in a memory bank with one-hot semantic labels. During inference, query features are matched
70 to the memory using FAISS with cosine similarity and $k = 30$ nearest neighbours. Hummingbird’s
71 cross-attention decoder aggregates the retrieved labels to produce segmentation predictions. The
72 framework is released under the MIT license (see Appendix C).

73 **3.3 Dataset**

74 Our study is based on MVIImgNet, a large-scale dataset with over 6.5 million frames across 238 object
75 categories [24]. Each frame includes a segmentation mask, camera extrinsics, and a reconstructed 3D
76 point cloud.

77 This dataset is chosen for the following three key reasons: (1) consistent object annotations across
78 wide viewpoint ranges, enabling view generalization analysis; (2) camera extrinsics that allow precise

79 angular binning; and (3) scale and diversity that capture a broad range of object appearances and
80 spatial configurations. It is released with CC BY-NC 4.0 for its code, and the dataset itself is
81 distributed under a custom Terms of Use provided by the authors (see Appendix C).

82 **Viewpoint binning.** To study viewpoint robustness, we discretize relative camera angles into seven
83 bins spanning $0^\circ - 90^\circ$ in 15° steps. Using COLMAP [16–18] extrinsics, we compute the relative
84 rotation R_{rel} between each frame and the canonical (first) frame of the instance¹. The angular
85 deviation θ is computed as:

$$\theta = \arccos\left(\frac{\text{trace}(R_{rel}) - 1}{2}\right). \quad (1)$$

86 For each object instance, we select one representative frame per bin by choosing the frame with the
87 smallest angular error relative to the bin center. Images and masks are stored per bin for downstream
88 use.

89 **Dataset construction.** To evaluate cross-view generalization, we curated a subset of MVIImgNet.
90 Our goal was to select categories with sufficient angular coverage and manageable size for repeated
91 memory-based evaluation.

92 The two criteria for selecting a category are: (1) the total zipped folder size must be 1–6 GB to ensure
93 feasible data loading and storage; and (2) at least one object instance must span all seven angular bins
94 with a maximum variance of 6 degrees. Categories that did not meet these criteria were excluded.
95 For example, category 23 (*laptop*) was excluded, as it did not have an instance that reached an angle
96 of 90° . This selection yielded 15 segmentation classes, excluding the background.

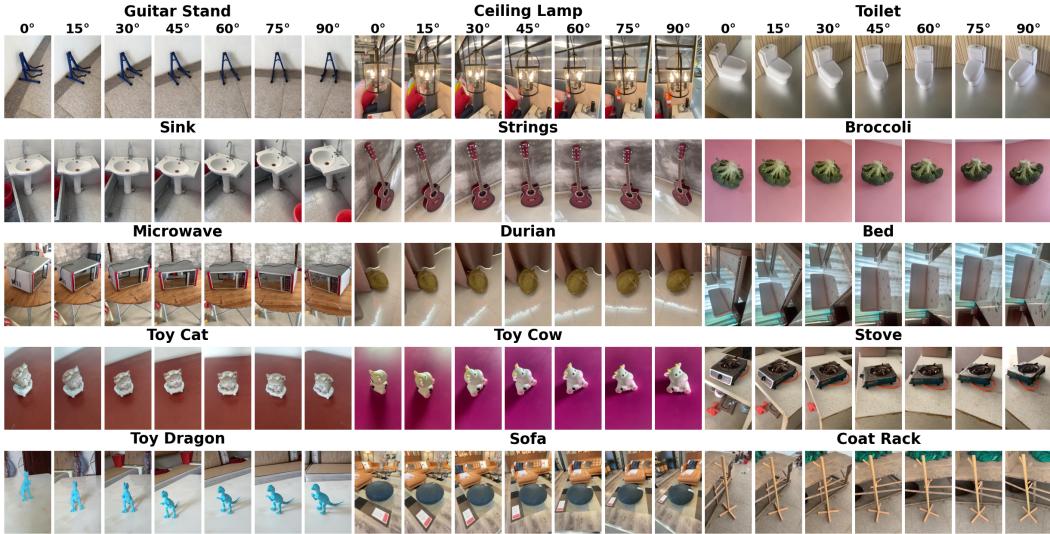


Figure 1: **Multi-view categories.** Each MVIImgNet category is shown across all viewpoint bins (0° – 90°) for the 15 selected classes.

97 For each category instance, we parsed the COLMAP camera pose data (`images.bin`) to extract
98 camera extrinsics and compute relative angles. For each valid instance, we selected the closest
99 frame to each angular bin center and organized the resulting RGB images and masks in a structured
100 directory grouped by category and angle. Additionally, we computed angular selection errors per bin
101 (see Table 5). Examples of multi-view objects from the resulting dataset are shown in Figure 1.

¹COLMAP is a Structure-from-Motion (SfM) and Multi-View Stereo (MVS) pipeline that estimates camera poses and 3D scene geometry from a set of images. It is commonly used to recover camera extrinsics and sparse 3D reconstructions from unstructured image sets.

102 **3.4 Model Input**

103 The models we use have different native input sizes: 224×224 for CLIP ViT-B/16 [13], DINO ViT-
104 B/16 [3] and C-RADIOv2 ViT-B/16-CPE [7, 14]; 512×512 for DINOv2 ViT-B/14 [12]; 504×504
105 for SigLIP2 B/16-512 [21]; and 448×448 for TIPS ViT-B/14-HR [9].

106 To ensure a fair comparison, the largest native input size across the models is used: 504×504
107 for models that use a patch size of 14 and 512×512 for the ones using a patch size of 16. The
108 evaluation of models with a bigger input size than their original resolution requires an additional step
109 of embedding interpolation. In our setup, we upsample the absolute positional embeddings using
110 bicubic interpolation to align with the new input resolution, as it provides smooth spatial transitions
111 and is widely used for resizing in vision models. All evaluations are performed with a batch size of 4.
112 The full configuration details are provided in Table 6.

113 **3.5 Metrics and Multi-view Evaluation**

114 The segmentation quality is measured by mean Intersection-over-Union (mIoU) across all semantic
115 classes, including the background. Each predicted segmentation mask consists of the background
116 class (jac_0) and the 15 object classes (jac_1 – jac_{15})², where jac_i refers to the Intersection
117 over Union (IoU) of class i .

118 We report three metrics per evaluation setting as follows: (1) class-wise IoU, (2) mIoU over all
119 16 classes (jac_{mean}), and (3) standard deviation (jac_{std}). Results are presented per model,
120 reference bins, and validation bin to analyze viewpoint robustness. Although all models perform well
121 on the background class (see Figure 4), we include it in the mIoU computation for completeness.

122 **3.6 Hardware and software**

123 All experiments were conducted on a high-performance computing cluster with $4 \times$ A100 GPUs,
124 72 CPU cores, and a maximum wall time of 30 hours per job. Due to memory constraints, we used
125 FAISS sharding to distribute the search index across GPUs and data parallelism to distribute the
126 computation of a batch size of 4. For reproducibility, we fixed the random seed of Python, NumPy,
127 PyTorch, and CUDA to 42.

128 **4 Reproducing the Hummingbird benchmark**

129 To validate our setup, we begin by reproducing the Hummingbird evaluation on Pascal VOC using a
130 publicly available implementation. Our reproduction results are shown in Table 8, obtained with the
131 same methodology and configuration described in the reference implementation. This step ensures
132 our pipeline is consistent with the original benchmark and provides a reliable starting point for
133 subsequent experiments.

134 **5 Experiment A: cross-model generalization**

135 We compare six pretrained ViTs in terms of their ability to generalize across viewpoints. The goal is
136 to assess how different pretraining strategies affect performance when only limited reference views
137 are available in the memory.

138 **5.1 Experimental setup**

139 We evaluate how well the models generalize to unseen viewpoints when trained on selected angular
140 bins from the MVImgNet dataset. The reference and validation splits vary with the difficulty level,
141 as defined in Table 1. Each model is evaluated using mIoU across all 16 classes (background plus
142 15 object categories) and all validation images. This setup enables a controlled comparison of each
143 model’s ability to interpolate between observed viewpoints and extrapolate to unseen ones.

²The 15 categories are ordered as follows: stove, sofa, microwave, bed, toy cat, toy cow, toy dragon, coat
rack, guitar stand, ceiling lamp, toilet, sink, strings, broccoli and durian.

Table 1: **Reference and validation bin splits.** Difficulty levels are defined by the choice of reference bins (stored in the memory bank) and validation bins (unseen views for testing generalization).

Difficulty	Reference bins	Validation bins
Easy	$0^\circ, 30^\circ, 60^\circ, 90^\circ$	$15^\circ, 45^\circ, 75^\circ$
Medium	$0^\circ, 45^\circ, 90^\circ$	$15^\circ, 30^\circ, 60^\circ, 75^\circ$
Hard	$0^\circ, 90^\circ$	$15^\circ, 30^\circ, 45^\circ, 60^\circ, 75^\circ$
Extreme	0°	$15^\circ, 30^\circ, 45^\circ, 60^\circ, 75^\circ, 90^\circ$

144 5.2 Results

145 We evaluated the models across four difficulty levels (see Table 1). As shown in Table 2, DINO
 146 achieved the highest mIoU in the first three difficulty levels, indicating strong generalization when
 147 trained on a wider range of bins. Under the Extreme difficulty, limited to a single reference bin
 148 (0°), DINOv2 surpassed all other models, suggesting greater robustness in extrapolating to unseen
 149 viewpoints.

150 Quantitative results illustrating these trends are shown in Figure 2, with per-category results in sub-
 151 section D.2. Across all classes, DINO and DINOv2 produced clearer segmentations under more
 152 difficult conditions than the other models. For example, the class-specific breakdown on *ceiling*
 153 *lamp* (Figure 14) highlights DINOv2’s robustness in unseen viewpoint bins.

154 Overall, CLIP consistently followed close behind DINO and DINOv2, while C-RADIOv2, TIPS, and
 155 SigLIP2 lagged significantly, especially under limited reference view conditions.

Table 2: **mIoU across difficulty levels.** Scores for each model are reported across increasing difficulty levels. DINO achieves the best performance under easier setups, whereas DINOv2 performs best when fewer reference views are available.

Model	Easy	Medium	Hard	Extreme
CLIP ViT-B/16	0.755 ± 0.130	0.748 ± 0.134	0.734 ± 0.140	0.701 ± 0.149
DINO ViT-B/16	0.782 ± 0.132	0.774 ± 0.137	0.748 ± 0.153	0.686 ± 0.171
DINOv2 ViT-B/14	0.763 ± 0.136	0.758 ± 0.139	0.748 ± 0.143	0.728 ± 0.154
C-RADIOv2 ViT-B/16-CPE	0.653 ± 0.129	0.636 ± 0.132	0.592 ± 0.141	0.506 ± 0.150
SigLIP 2 B/16-512	0.564 ± 0.150	0.551 ± 0.153	0.530 ± 0.157	0.481 ± 0.152
TIPS ViT-B/14-HR	0.667 ± 0.137	0.647 ± 0.146	0.588 ± 0.162	0.462 ± 0.169

156 5.3 Discussion

157 Although none of the models were explicitly trained for viewpoint understanding, the results show
 158 clear differences in their ability to generalize across views. DINO and DINOv2 perform best,
 159 indicating that their features capture object shape more reliably and remain consistent under view-
 160 point changes. Their self-supervised training, which aligns features from different augmentations
 161 of the same image, promotes consistency between local and global cues and leads to structured
 162 representations where similar object parts share similar embeddings.

163 DINO achieved the best results under the Easy difficulty, where multiple views of the object were
 164 available, while DINOv2 performed best under the Extreme difficulty with only a single reference
 165 bin. This contrast suggests that DINO’s simpler self-distillation loss favors interpolation across views,
 166 whereas DINOv2’s larger-scale self-supervised pretraining supports stronger extrapolation to unseen
 167 angles.

168 CLIP was consistently ranked third, reflecting some ability to encode visual–semantic regularities
 169 despite being trained for image–text alignment rather than geometric consistency. By comparison,
 170 SigLIP2, TIPS, and C-RADIOv2 performed substantially worse, indicating that objectives centered
 171 on semantic matching or mixed supervision may weaken part-level consistency needed for viewpoint-
 172 robust segmentation.

Experiment A Results — All classes

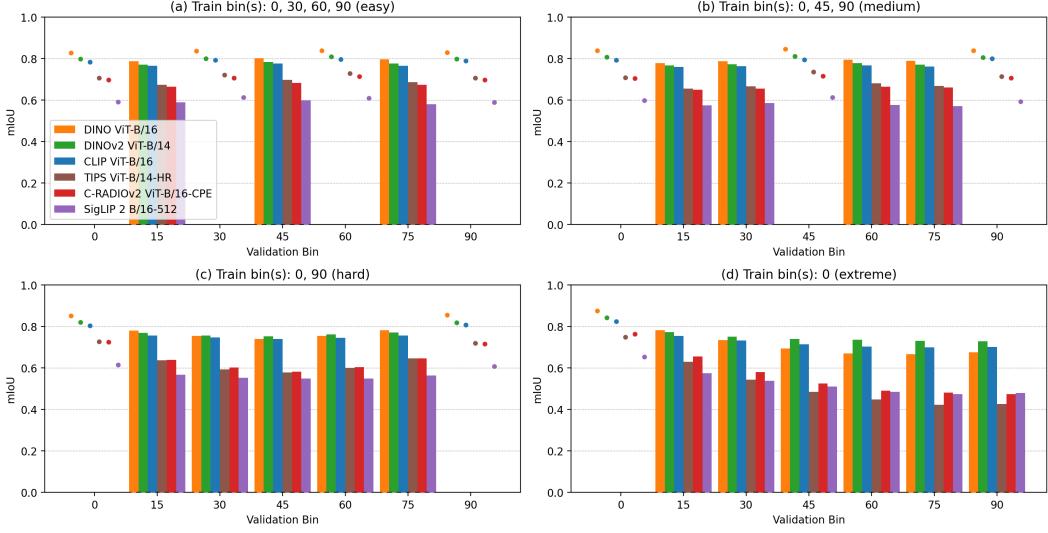


Figure 2: **Segmentation across viewpoint bins.** Each subplot (a)–(d) corresponds to a different difficulty. Bars show performance on validation bins, while dots show performance on the reference bins stored in the memory bank. Performance declines as validation angles move further from the reference bins, with DINO and DINoV2 degrading more slowly than multimodal encoders such as CLIP and SigLIP2.

173 6 Experiment B: breaking point analysis

174 We analyze whether and when each model experiences a sudden failure in viewpoint generalization.
 175 The goal is to identify the angular range within which performance remains stable under the Extreme
 176 setting defined in Table 1.

177 6.1 Experimental setup

178 We perform a breaking point analysis under the Extreme setting from Experiment A, where only a
 179 single reference bin (0°) is available in the memory and models are evaluated across the remaining
 180 bins. We define the *breaking point* as the earliest validation bin where a model’s normalized
 181 performance drops significantly compared to the previous bin. Specifically, we compute

$$\Delta_i = \text{norm mIoU}_i - \text{norm mIoU}_{i-1} \quad (2)$$

$$\text{norm mIoU}_i = \frac{\text{mIoU}_i}{\text{mIoU}_{0^\circ}}, \quad (3)$$

182 where the normalization by the 0° bin ensures that performance drops reflect viewpoint sensitivity
 183 rather than overall model scale.

184 A breaking point is recorded at bin i if $\Delta_i \leq -0.1$, indicating a relative drop of more than 10%. This
 185 analysis highlights each model’s resistance to viewpoint shifts and identifies the angular range within
 186 which performance remains stable. We also report the normalized mIoU degradation curves as the
 187 viewpoint angle increases. The baseline that is used is an evaluation of the same setting, but with the
 188 validation bins being part of the memory.

189 6.2 Results

190 As illustrated in Figure 3, most models show a gradual decline in normalized mIoU as the validation
 191 angle increases. C-RADIOv2 and TIPS are the models that reach a breaking point, both at the 30°

192 bin (see Table 3). TIPS shows the steepest drop, with a relative decrease of -0.125 between 15°
 193 and 30° . By contrast, DINO, DINoV2, CLIP, and SigLIP2 degrade smoothly without reaching the
 194 breaking point.

Table 3: **Breaking points.** For each model in Experiment B, we report the validation bin where a breaking point occurs together with the corresponding normalized mIoU drop ($\Delta_i \leq -0.1$) relative to the previous bin.

Model	Breaking Point Bin	Δ (Drop)
CLIP ViT-B/16	None	–
DINO ViT-B/16	None	–
DINOv2 ViT-B/14	None	–
C-RADIOv2 ViT-B/16-CPE	30	-0.1049
SigLIP2 B/16-512	None	–
TIPS ViT-B/14-HR	30	-0.1250

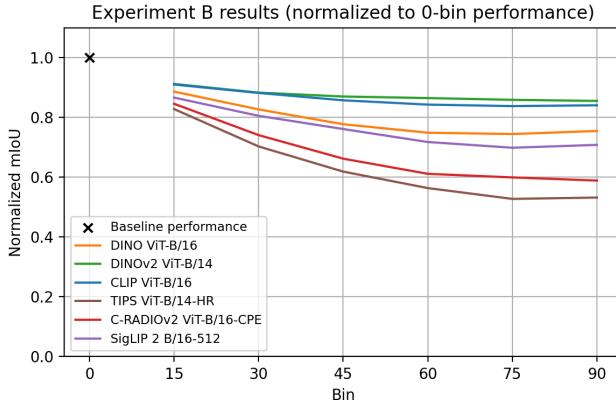


Figure 3: **Normalized mIoU.** We show normalized mIoU curves (Δ_i) per model in Experiment B, measured relative to the 0-bin performance.

195 6.3 Discussion

196 The breaking point analysis highlights differences in robustness across models and shows how their
 197 pretraining objectives affect stability under viewpoint shifts. DINoV2 is the most stable, maintaining
 198 gradual degradation without sudden drops. DINO shows a similar trend, though with slightly weaker
 199 overall performance. CLIP and SigLIP2 also decrease steadily, suggesting that they preserve some
 200 resilience despite being trained for image–text alignment rather than spatial structure.

201 In contrast, TIPS and C-RADIOv2 exhibit clear breakdowns, indicating that their learned features
 202 are less robust when faced with unseen viewpoints. For TIPS, this instability likely stems from its
 203 focus on semantic alignment rather than spatial consistency, while C-RADIOv2’s mixed supervision
 204 appears to weaken the reliability of part-level representations.

205 Relative to the baseline, where validation bins are also included in memory, these findings confirm
 206 that the sudden decline comes from missing reference views rather than instability in the evaluation
 207 pipeline. Overall, the results show that self-supervised encoders (DINO, DINoV2) produce more
 208 geometry-aware features under strong viewpoint shifts, while multimodal or distilled approaches are
 209 more brittle. DINoV2’s stability makes it the most reliable choice for applications that must handle
 210 abrupt perspective changes.

211 7 Experiment C: memory size robustness

212 In this experiment, we analyze how memory bank size affects viewpoint generalization, comparing
 213 the performance gains across difficulty levels against the added computational cost.

Table 4: **Memory size comparison.** Viewpoint-based segmentation performance (mIoU) is reported for three memory sizes: (a) 320,000, (b) 640,000, and (c) 1,024,000. Bold values mark the best score per difficulty-memory configuration, and each cell shows mean \pm standard deviation.

Model	Easy	Medium	Hard	Extreme
(a) Memory = 320,000				
CLIP ViT-B/16	0.729 ± 0.138	0.725 ± 0.141	0.710 ± 0.149	0.681 ± 0.156
DINO ViT-B/16	0.741 ± 0.147	0.736 ± 0.150	0.712 ± 0.163	0.656 ± 0.178
DINOv2 ViT-B/14	0.738 ± 0.145	0.736 ± 0.146	0.726 ± 0.152	0.709 ± 0.158
C-RADIOv2 ViT-B/16-CPE	0.588 ± 0.136	0.579 ± 0.139	0.539 ± 0.146	0.467 ± 0.149
SigLIP 2 B/16-512	0.506 ± 0.154	0.501 ± 0.156	0.483 ± 0.157	0.443 ± 0.149
TIPS ViT-B/14-HR	0.600 ± 0.161	0.590 ± 0.166	0.539 ± 0.178	0.437 ± 0.175
(b) Memory = 640,000				
CLIP ViT-B/16	0.745 ± 0.132	0.740 ± 0.136	0.727 ± 0.141	0.694 ± 0.150
DINO ViT-B/16	0.768 ± 0.137	0.761 ± 0.141	0.736 ± 0.156	0.676 ± 0.173
DINOv2 ViT-B/14	0.754 ± 0.138	0.750 ± 0.142	0.740 ± 0.147	0.721 ± 0.155
C-RADIOv2 ViT-B/16-CPE	0.629 ± 0.132	0.615 ± 0.135	0.572 ± 0.142	0.491 ± 0.150
SigLIP 2 B/16-512	0.541 ± 0.152	0.531 ± 0.154	0.511 ± 0.157	0.466 ± 0.149
TIPS ViT-B/14-HR	0.644 ± 0.144	0.628 ± 0.153	0.571 ± 0.167	0.453 ± 0.170
(c) Memory = 1,024,000				
CLIP ViT-B/16	0.755 ± 0.130	0.748 ± 0.134	0.734 ± 0.140	0.701 ± 0.149
DINO ViT-B/16	0.782 ± 0.132	0.774 ± 0.137	0.748 ± 0.153	0.686 ± 0.171
DINOv2 ViT-B/14	0.763 ± 0.136	0.758 ± 0.139	0.748 ± 0.143	0.728 ± 0.154
C-RADIOv2 ViT-B/16-CPE	0.653 ± 0.129	0.636 ± 0.132	0.592 ± 0.141	0.506 ± 0.150
SigLIP 2 B/16-512	0.564 ± 0.150	0.551 ± 0.153	0.530 ± 0.157	0.481 ± 0.152
TIPS ViT-B/14-HR	0.667 ± 0.137	0.647 ± 0.146	0.588 ± 0.162	0.462 ± 0.169

7.1 Experimental setup

We follow the Hummingbird paper [1] and vary the number of entries in the memory bank. Three configurations are evaluated: 320k, 640k, and 1,024k. For each configuration we report segmentation accuracy and measure computational cost, including runtime and resource usage.

7.2 Results

The results in Table 4 show that increasing memory generally improves mIoU, with the largest gains observed under the Easy difficulty. Doubling memory from 320k to 640k increases performance by about 0.030 mIoU on Easy, while a further increase to 1,024k yields a smaller improvement of 0.017 (see Table 9).

DINO and DINOv2 achieve the highest absolute performance across all memory sizes. However, models such as C-RADIOv2, SigLIP2, and TIPS show the greatest relative gains. Each improves by more than 0.048 when moving from 320k to 1,024k, with most of that improvement coming from the first doubling.

The runtime analysis in Table 10 shows that larger memory banks increase computational cost. The highest configuration adds about two hours per job. Total memory usage remains similar (189–214GB), but CPU efficiency declines as memory grows: for example, C-RADIOv2 drops from 16.35% at 320k to 12.52% at 1,024k.

7.3 Discussion

These results confirm that larger memory banks improve generalization, but the benefits vary across models. Less robust encoders such as C-RADIOv2, SigLIP2, and TIPS gain the most from additional memory, while DINO and DINOv2 already perform strongly even with smaller banks. The diminishing gains beyond 640k and the higher computational cost suggest that the intermediate configuration offers the best trade-off. This shows that memory scaling can compensate for weaker representations, but stable self-supervised features remain more effective overall.

238 **8 Qualitative analysis of segmentations**

239 In addition to quantitative evaluations, we conducted a qualitative inspection of predicted segmentation
240 masks to better understand model behaviour under the large-angle changes of the Extreme difficulty.

241 **Localized shape recovery.** As shown in Figure 27, the model successfully segments the dinosaur
242 figure from an unfamiliar viewpoint. While the prediction (center) is not a perfect pixel-wise match
243 to the ground truth (left), the overall contour and structure of the object are preserved. The prediction
244 predicts that that the object extends further (e.g., tail tip and scattered patches), and the core body
245 is segmented as a bigger area (encompassing more of the bony plates embedded in the back of the
246 dinosaur).

247 **Prediction versus ground truth overlays.** In Figure 28, we visualize the predicted mask (center)
248 versus the ground truth (left) when overlaid on the input image. These overlays reveal that the model
249 tends to underestimate object boundaries in highly self-occluded regions or under harsh lighting (e.g.,
250 dark shadows like those under the dinosaur’s core). Nonetheless, the predictions largely align with
251 visually salient contours of the object, indicating strong local feature retrieval. When looking at the
252 overlay of the ground truth over the prediction (right), we notice that the prediction actually covers
253 more of the core and parts of the dinosaur’s tail, an improvement as compared to the ground truth.

254 **Prediction versus ground truth.** We also visualize mask comparisons using color-coded difference
255 maps (Figure 28 (right) or Figure 29). These maps highlight regions of over- and under-segmentation.
256 Notably, misalignments are often located near object boundaries or thin structures (e.g., legs, tails),
257 suggesting that memory retrieval may struggle with fine-grained spatial resolution under Extreme
258 difficulty views. In addition, dark shadows are often included as parts of the segmentation mask.
259 In some cases, the ground truth masks themselves are of poor quality, as illustrated in Figures 8
260 and 20–25 for the bed class. This leads to decreased scores even when the prediction masks are
261 visually more accurate. Such cases highlight that annotation quality can limit the reliability of
262 quantitative evaluation, since metrics are tied to imperfect labels rather than perceptual correctness.

263 **Failures.** In Figure 29, we show a failure case on a different *toy dragon* instance. Here, the predicted
264 segmentation extends beyond the actual object boundary, indicating that the model retrieved visually
265 similar but geometrically inconsistent patches. Segmentation failures such as this one are causing
266 certain models (e.g., C-RADIOv2, TIPS) to exhibit performance drops beyond 30° angular deviation.

267 **9 Conclusion**

268 We studied the ability of frozen ViT encoders to generalize across unseen camera viewpoints in
269 an in-context segmentation setting, using the Hummingbird architecture and MVIImgNet dataset.
270 Our benchmark covered three aspects: controlled evaluation under view shifts, robustness through
271 breaking point analysis, and the effect of memory bank size.

272 Our findings show that self-supervised encoders, particularly DINO and DINOv2, provide more
273 stable geometry-aware features under viewpoint changes, while multimodal and distilled approaches
274 are more brittle. Larger memory banks help weaker models but offer diminishing returns for stronger
275 encoders, raising questions about the cost–benefit trade-off.

276 Taken together, these results highlight the strengths and limitations of current ViTs for multi-view
277 perception. While modern self-supervised training yields robust features, none of the models are
278 immune to viewpoint-induced degradation. Addressing this gap will likely require objectives or
279 architectures that enforce 3D consistency. Future work could include extending this evaluation to
280 multi-object scenes, wider angular ranges, and compound rotations, providing a broader testbed for
281 building truly viewpoint-robust representations.

282 **References**

- 283 [1] Ivana Balažević, David Steiner, Nikhil Parthasarathy, Relja Arandjelović, and Olivier J. Hé-
284 naff. Towards in-context scene understanding, 2023. URL <https://arxiv.org/abs/2306.01667>.

- 286 [2] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
 287 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson,
 288 Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel,
 289 Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano
 290 Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren
 291 Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto,
 292 Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas
 293 Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling,
 294 Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi,
 295 Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa
 296 Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric
 297 Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman,
 298 Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr,
 299 Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi
 300 Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack
 301 Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan
 302 Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang,
 303 William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga,
 304 Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia
 305 Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models,
 306 2022. URL <https://arxiv.org/abs/2108.07258>.
- 307 [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,
 308 and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021. URL
 309 <https://arxiv.org/abs/2104.14294>.
- 310 [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 311 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
 312 Jakoba Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image
 313 Recognition at Scale. In *International Conference on Learning Representations*, 2020.
- 314 [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 315 recognition, 2015. URL <https://arxiv.org/abs/1512.03385>.
- 316 [6] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
 317 toencoders are scalable vision learners, 2021. URL <https://arxiv.org/abs/2111.06377>.
- 318 [7] Greg Heinrich, Mike Ranzinger, Hongxu, Yin, Yao Lu, Jan Kautz, Andrew Tao, Bryan Catan-
 319 zaro, and Pavlo Molchanov. Radiov2.5: Improved baselines for agglomerative vision foundation
 320 models, 2025. URL <https://arxiv.org/abs/2412.07679>.
- 321 [8] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
 322 Guo. Swin transformer: Hierarchical vision transformer using shifted windows, 2021. URL
 323 <https://arxiv.org/abs/2103.14030>.
- 324 [9] Kevis-Kokitsi Maninis, Kaifeng Chen, Soham Ghosh, Arjun Karpur, Koert Chen, Ye Xia, Bingyi
 325 Cao, Daniel Salz, Guangxing Han, Jan Dlabal, Dan Gnanapragasam, Mojtaba Seyedhosseini,
 326 Howard Zhou, and André Araujo. TIPS: Text-Image Pretraining with Spatial Awareness. In
 327 *ICLR*, 2025.
- 328 [10] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi,
 329 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis, 2020. URL
 330 <https://arxiv.org/abs/2003.08934>.
- 331 [11] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Flo-
 332 rence Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat,
 333 Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao,
 334 Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,
 335 Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman,
 336 Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, An-
 337 drew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis
 338 Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester

- Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiro, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ibai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotstet, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Toootchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- [12] Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. URL <https://arxiv.org/abs/2304.07193>.
- [13] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- [14] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Am-radio: Agglomerative vision foundation model – reduce all domains into one, 2024. URL <https://arxiv.org/abs/2312.06709>.
- [15] Shouwei Ruan, Yinpeng Dong, Hang Su, Jianteng Peng, Ning Chen, and Xingxing Wei. Towards viewpoint-invariant visual recognition via adversarial training, 2023. URL <https://arxiv.org/abs/2307.10235>.

- 395 [16] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In
396 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 397 [17] Johannes Lutz Schönberger, True Price, Torsten Sattler, Jan-Michael Frahm, and Marc Pollefeys.
398 A vote-and-verify strategy for fast spatial verification in image retrieval. In *Asian Conference*
399 *on Computer Vision (ACCV)*, 2016.
- 400 [18] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise
401 view selection for unstructured multi-view stereo. In *European Conference on Computer Vision*
402 (*ECCV*), 2016.
- 403 [19] Ofir Shifman and Yair Weiss. Lost in translation: Modern neural networks still struggle with
404 small realistic image transformations, 2024. URL <https://arxiv.org/abs/2404.07153>.
- 405 [20] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-
406 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
407 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
408 language models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 409 [21] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alab-
410 dulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier
411 Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-
412 language encoders with improved semantic understanding, localization, and dense features,
413 2025. URL <https://arxiv.org/abs/2502.14786>.
- 414 [22] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-
415 context learning as implicit bayesian inference, 2022. URL <https://arxiv.org/abs/2111.02080>.
- 416 [23] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
417 Gao, Chengan Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
418 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
419 Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
420 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
421 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
422 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
423 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
424 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 425 [24] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng
426 Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, Guanying Chen, Shuguang Cui,
427 and Xiaoguang Han. Mvimgnet: A large-scale dataset of multi-view images, 2023. URL
428 <https://arxiv.org/abs/2303.06042>.

430 **A MVIImgNet dataset**

431 The statistics for the reorganized subset of MVIImgNet [24] we used are shown in Table 5. The
 432 selected subset includes 15 object classes, each represented across 7 standardized angle bins:
 433 $[0^\circ, 15^\circ, 30^\circ, 45^\circ, 60^\circ, 75^\circ, 90^\circ]$. Figure 1 shows an instance of each angle for each object class.

Table 5: **Angle selection accuracy per class.** Mean and standard deviation of angular error (in degrees) are reported for each object class after selecting views closest to the predefined angle bins.

Class Number	Category	Std. Error	Mean Error	# Images in each Angle bin
7	Stove	1.28	0.01	197
8	Sofa	1.15	-0.04	91
19	Microwave	1.44	-0.04	120
46	Bed	1.17	-0.00	23
57	Toy Cat	1.96	-0.01	783
60	Toy Cow	1.96	-0.02	735
70	Toy Dragon	1.62	0.01	627
99	Coat Rack	1.41	-0.02	97
113	Ceiling Lamp	1.64	-0.03	154
125	Toilet	1.31	0.02	58
126	Sink	1.20	-0.12	30
152	Strings	1.25	0.03	192
166	Broccoli	2.04	-0.03	210
196	Durian	1.65	0.03	758
100	Guitar Stand	1.53	0.04	218

434 **B Model specifics**

435 The following table details the ViT-based encoders used in our experiments, including their patch
 436 configurations and feature dimensionality.

Table 6: **Vision Transformer configurations.** We show the architectural and input settings of the models evaluated in our experiments.

Model	Architecture	Input Size	Patch Size	Patch Count	Feature Dim.	Batch Size
DINO	ViT-B/16	512	16	1024	768	4
DINOv2	ViT-B/14	504	14	1296	768	4
OpenAI CLIP	ViT-B/16	512	16	1024	768	4
C-RADIOv2	ViT-B/16-CPE	512	16	1024	768	4
SigLIP2	ViT-B/16-512	512	16	1024	768	4
TIPS	ViT-B/14-HR	504	14	1296	768	4

437 **C Licenses for assets**

438 All datasets and models used in this work are released under open licenses. We list them here for
 439 clarity:

Table 7: **Licenses for datasets and models.** All assets are used in compliance with their respective licenses.

Asset	License
MVImgNet	CC BY-NC 4.0 (code), dataset Terms of Use
Hummingbird	MIT License
CLIP	MIT
DINO	Apache 2.0
DINOv2	Apache 2.0
SigLIP2	Apache 2.0
C-RADIOv2	NVIDIA Open Model License
TIPS	Apache 2.0 (code), CC BY 4.0 (docs)

440 **D Additional results**

441 **D.1 Reproduction**

Table 8: **Reproduction results.** We report accuracy (%) at different evaluation scales. “Reported” indicates values taken from prior published results, while “Reproduced” refers to our reproduction with a batch size (BS). Bold values highlight the best result for each model-scale configuration, with reproduced results differing by less than 3%.

Model	Source	1024×10^2	1024×10^3	1024×10^4
ViT-S/16	Reproduced, BS: 256	37.5	45.0	—
	Reproduced, BS: 64	37.9	45.1	49.3
	Reported, BS: 64	37.2	43.1	46.6
ViT-B/16	Reproduced, BS: 64	48.0	54.7	—
	Reproduced, BS: 32	47.8	54.6	—
	Reproduced, BS: 8	—	—	57.9
	Reported, BS: 64	44.9	50.8	55.7
ViT-S/14	Reproduced, BS: 64	69.6	75.1	77.0
	Reported, BS: 64	70.2	74.9	77.0
ViT-B/14	Reproduced, BS: 64	68.0	74.0	—
	Reproduced, BS: 8	—	—	76.6
	Reported, BS: 64	69.1	74.6	76.9
ViT-L/14	Reproduced, BS: 64	64.1	71.2	—
	Reproduced, BS: 8	—	—	74.4
	Reported, BS: 64	64.6	71.7	74.8
ViT-G/14	Reproduced, BS: 32	62.4	—	—
	Reproduced, BS: 16	—	70.1	—
	Reproduced, BS: 8	—	—	73.3
	Reported, BS: 64	62.3	69.9	73.6

442 D.2 Experiment A

443 In Experiment A, we assess model generalization across view-angle difficulties using per-class
 444 performance curves. The background class performs best, while bed, coat rack, guitar stand, and sink
 445 perform worst. For beds, the drop is explained by poor ground-truth annotations that likely reduce
 446 model scores.

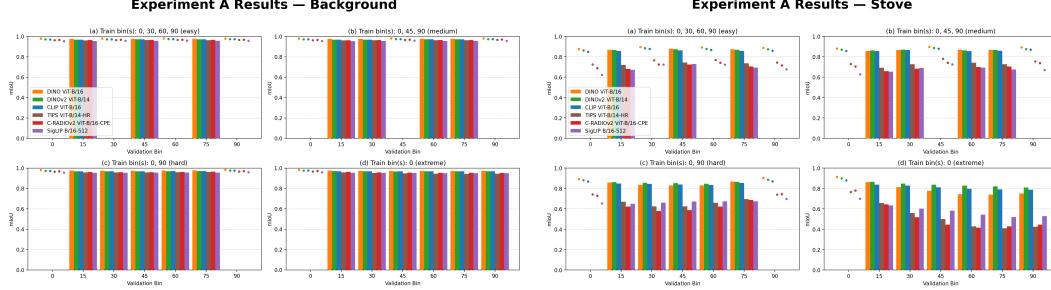


Figure 4: **Experiment A: background.** View-angle generalization for the class *background*, with four training difficulty settings: (a) easy; (b) medium; (c) hard; and (d) extreme.

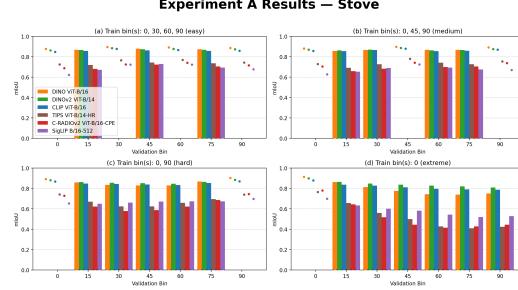


Figure 5: **Experiment A: stove.** SigLIP2 does relatively better than the average performance in Figure 2

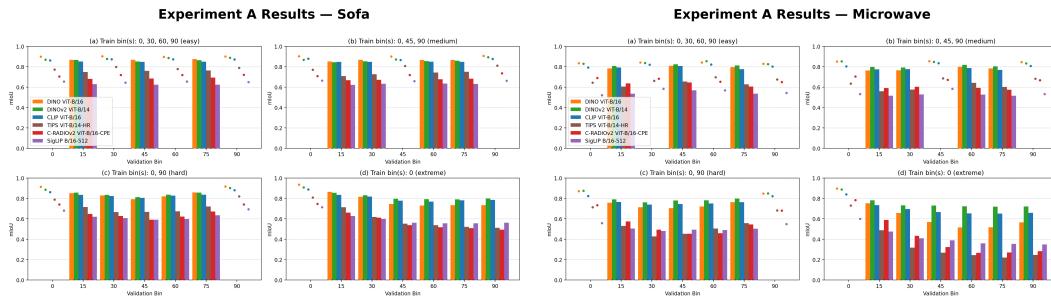


Figure 6: **Sofa.** Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

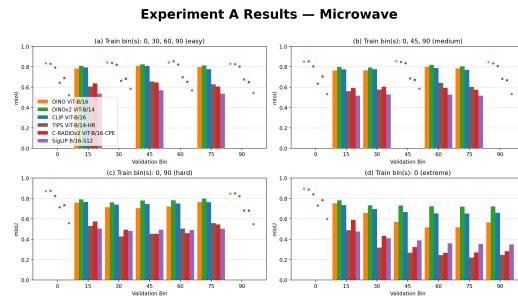


Figure 7: **Microwave.** TIPS appears to struggle when segmenting the microwave. Performs poorer than in Figure 2

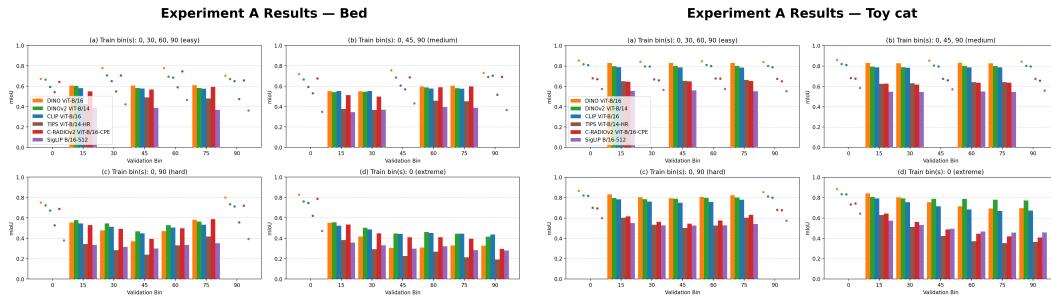


Figure 8: **Bed.** Beds have bad ground truths (seen in Figures 20–25 of Appendix D.2.1), resulting in a drop in overall performance

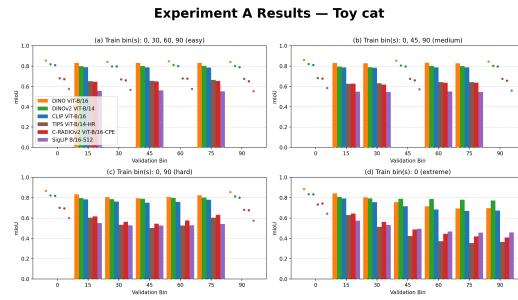


Figure 9: **Cat.** Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Toy cow

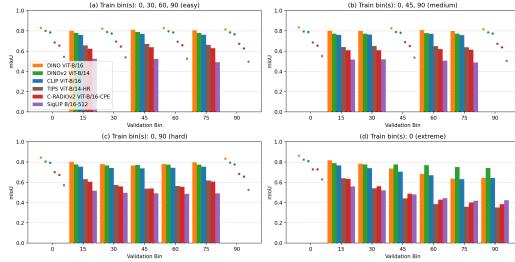


Figure 10: Cow. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Toy dragon

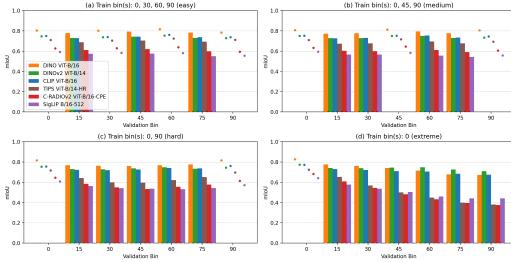


Figure 11: Toy dragon. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Coat rack

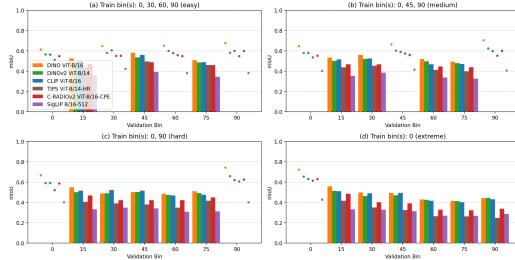


Figure 12: Coat rack. Performance degradation is observed, but its cause remains unclear and is a point for future investigation.

Experiment A Results — Guitar stand

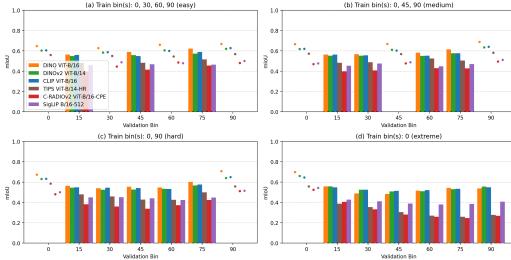


Figure 13: Guitar stand. Performance degradation is observed, but its cause remains unclear and is a point for future investigation.

Experiment A Results — Ceiling lamp

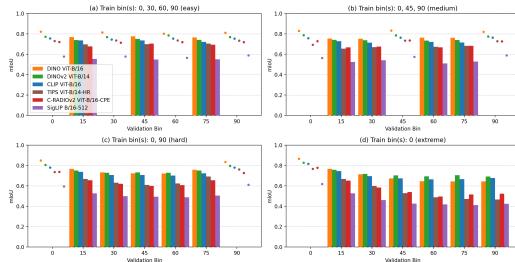


Figure 14: Ceiling lamp. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Toilet

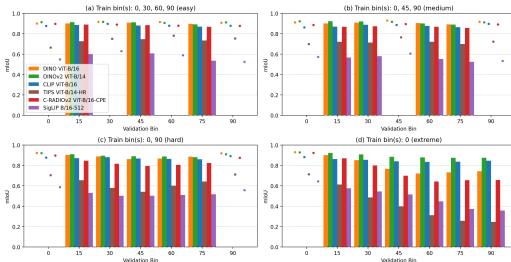


Figure 15: Toilet. C-Radio does significantly better on toilets. Resulting in a greater performance compared to Figure 2

Experiment A Results — Sink

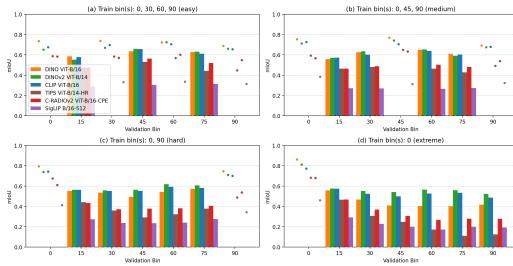


Figure 16: **Sink**. Performance degradation is observed, but its cause remains unclear and is a point for future investigation.

Experiment A Results — Strings

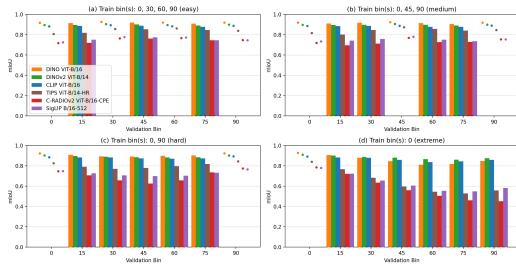


Figure 17: **Strings**. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Broccoli

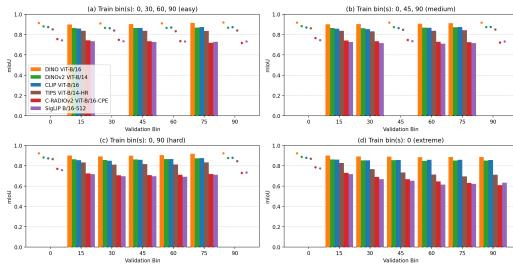


Figure 18: **Broccoli**. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

Experiment A Results — Durian

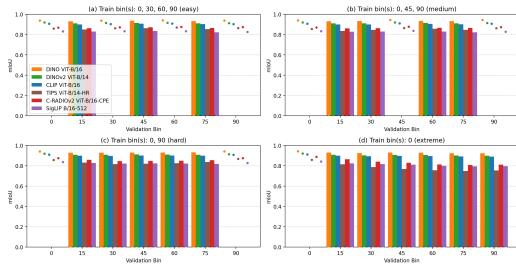


Figure 19: **Durian**. Performance trends mirror the general degradation observed in Figure 2, with accuracy dropping consistently across difficulty settings.

447 **D.2.1 Bed category ground truth analysis**

448 We investigated the poor performance of the 'bed' class (as seen in Figure 8) and identified inconsistent
449 or incomplete ground truth annotations as a likely cause.



Figure 20: **Bed image example 1.** Original image used for annotation.



Figure 21: **Incomplete annotation.** Ground truth excludes parts of the bed (side railing, footboard) and ignores surrounding beds.



Figure 22: **Bed image example 2.** Original image used for annotation.

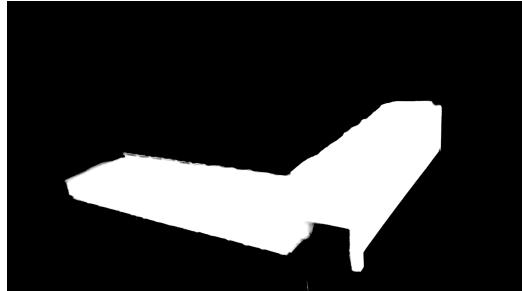


Figure 23: **Erroneous annotation.** The ground truth excludes large parts of the bed (headboard, mattress) and incorrectly includes the shadow between two beds.



Figure 24: **Bed image example 3.** Original image used for annotation.



Figure 25: **Annotation error.** The ground truth omits the surrounding beds, incorrectly includes the gap between the main bed and the bed to its right, and only partially marks the headboard (side only).

450 **D.3 Experiment B**

451 As shown in Figure 3 and Figure 26, most models exhibit a gradual performance decline as the
452 distance of the validation bin from the training bin (0°) increases. Notably, only C-RADIOv2 (ViT-
453 B/16-CPE) and TIPS (ViT-B/14-HR) exhibit sharp performance drops with increasing validation
454 angles, indicating limited generalization capacity under extreme viewpoint shifts.

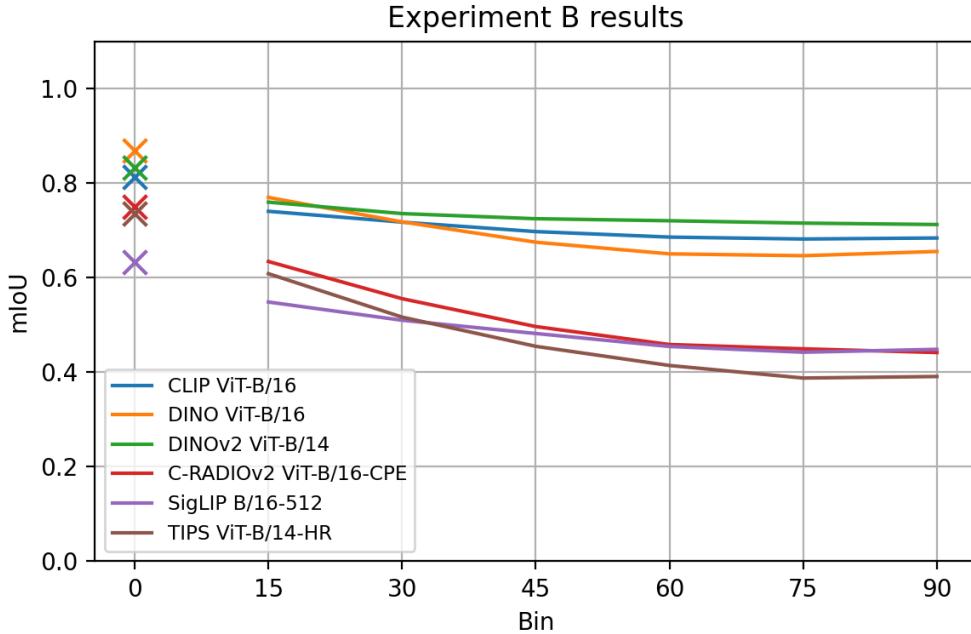


Figure 26: **Absolute mIoU.** We report absolute mIoU scores per validation bin. Trends follow the normalized plot, with DINOv2 showing the most consistent generalization across bins.

455 **D.4 Experiment C**

456 In Experiment C, we evaluate the impact of memory bank size on model performance. Table 9
 457 reports the absolute mIoU gains per difficulty level when increasing the memory bank for entries (a)
 458 from 320,000 to 640,000, (b) from 640,000 to 1,024,000, and (c) from 320,000 to 1,024,000. This
 459 highlights how memory sensitivity varies across model architectures. On the other hand, Table 10
 460 outlines the computational trade-offs. We report the runtime and memory usage associated with
 461 each model under both memory configurations. This provides context for the computational cost
 462 underlying the performance results.

Table 9: **Performance gains from memory.** We report absolute mIoU improvements when increasing memory from (a) 320,000 to 640,000, (b) 640,000 to 1,024,000, and (c) 320,000 to 1,024,000. Values are computed as differences from Table 4, with the final column showing the average gain across all four difficulty levels.

Model	Gain (Easy)	Gain (Medium)	Gain (Hard)	Gain (Extreme)	Average Gain
(a) Memory = 320,000 → 640,000					
CLIP ViT-B/16	0.016	0.015	0.017	0.013	0.01525
DINO ViT-B/16	0.027	0.025	0.024	0.020	0.02400
DINOv2 ViT-B/14	0.016	0.014	0.014	0.012	0.01400
C-RADIOv2 ViT-B/16-CPE	0.041	0.036	0.033	0.024	0.03350
SigLIP 2 B/16-512	0.035	0.030	0.028	0.023	0.02900
TIPS ViT-B/14-HR	0.044	0.038	0.032	0.016	0.03250
Average per task	0.030	0.026	0.0247	0.018	
(b) Memory = 640,000 → 1,024,000					
CLIP ViT-B/16	0.010	0.008	0.007	0.007	0.00800
DINO ViT-B/16	0.014	0.013	0.012	0.010	0.01225
DINOv2 ViT-B/14	0.009	0.008	0.008	0.007	0.00800
C-RADIOv2 ViT-B/16-CPE	0.024	0.021	0.020	0.015	0.02000
SigLIP 2 B/16-512	0.023	0.020	0.019	0.015	0.01925
TIPS ViT-B/14-HR	0.023	0.019	0.017	0.009	0.01700
Average per task	0.017	0.015	0.0138	0.0105	
(c) Memory = 320,000 → 1,024,000					
CLIP ViT-B/16	0.026	0.023	0.024	0.020	0.02325
DINO ViT-B/16	0.041	0.038	0.036	0.030	0.03625
DINOv2 ViT-B/14	0.025	0.022	0.022	0.019	0.02200
C-RADIOv2 ViT-B/16-CPE	0.065	0.057	0.053	0.039	0.05350
SigLIP 2 B/16-512	0.058	0.050	0.047	0.038	0.04825
TIPS ViT-B/14-HR	0.067	0.057	0.049	0.025	0.04950
Average per task	0.047	0.041	0.0385	0.0285	

Table 10: **System performance under memory scaling.** We report wall-clock time, memory usage, CPU utilization, CPU efficiency, and memory efficiency for viewpoint-based segmentation with memory sizes (a) 320,000, (b) 640,000, and (c) 1,024,000. Some entries are marked as *(running)* or show 0.00 usage values; these correspond to cases where monitoring logs did not update correctly, even though the jobs completed successfully. Bold values highlight the largest drops in CPU efficiency for each model.

Model	Wall-clock Time	CPU Used	CPU Eff.	Memory Used	Mem. Eff.
(a) Memory = 320,000					
CLIP ViT-B/16	07:32:14	3-15:33:05	16.13%	198.14 GB	41.28%
DINO ViT-B/16	07:32:35	3-12:20:38	15.53%	214.11 GB	44.61%
DINOv2 ViT-B/14	08:39:50 (running)	00:00:00	0.00%	0.00 GB	0.00%
C-RADIOv2 ViT-B/16-CPE	07:27:39	3-15:50:23	16.35%	199.82 GB	41.63%
SigLIP 2 B/16-512	07:39:28	3-19:58:28	16.68%	189.89 GB	39.56%
TIPS ViT-B/14-HR	08:33:54 (running)	00:00:00	0.00%	0.00 GB	0.00%
(b) Memory = 640,000					
CLIP ViT-B/16	08:21:00	3-10:39:36	13.75%	177.58 GB	37.00%
DINO ViT-B/16	08:28:26 (running)	00:00:00	0.00%	0.00 GB	0.00%
DINOv2 ViT-B/14	09:43:52 (running)	00:00:00	0.00%	0.00 GB	0.00%
C-RADIOv2 ViT-B/16-CPE	08:18:48	3-11:02:33	13.87%	203.63 GB	42.42%
SigLIP 2 B/16-512	08:26:16	3-16:01:55	14.49%	181.39 GB	37.79%
TIPS ViT-B/14-HR	09:33:10 (running)	3-09:16:27	11.82%	209.87 GB	43.72%
(c) Memory = 1,024,000					
CLIP ViT-B/16	09:20:31	3-07:03:04	11.75%	207.41 GB	43.21%
DINO ViT-B/16	09:29:42	3-07:07:59	11.58%	213.23 GB	44.42%
DINOv2 ViT-B/14	11:10:37	3-11:36:03	10.39%	209.23 GB	43.59%
C-RADIOv2 ViT-B/16-CPE	09:28:53	3-13:29:38	12.52%	200.34 GB	41.74%
SigLIP 2 B/16-512	09:36:48 (running)	00:00:00	0.00%	0.00 GB	0.00%
TIPS ViT-B/14-HR	11:04:09	3-11:57:18	10.53%	211.38 GB	44.04%

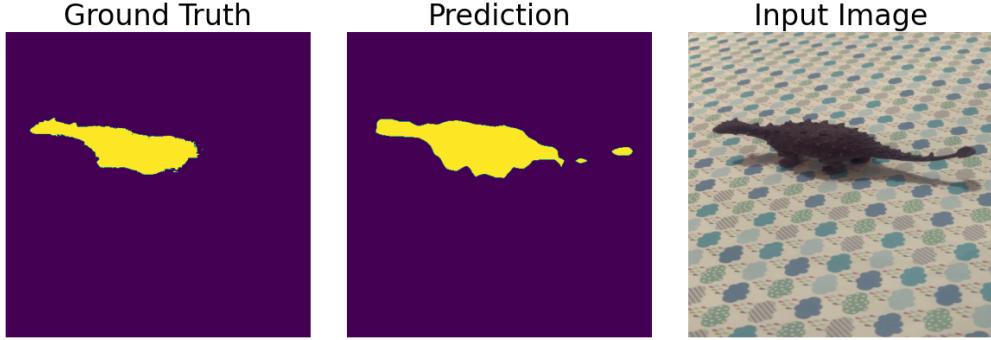


Figure 27: **Qualitative segmentation.** Results shown using DINO. Left: input image. Center: predicted mask. Right: ground truth mask. Interestingly, the prediction of the *toy dragon* aligns more closely with visible object boundaries than the ground truth, which appears coarser and less accurate at fine details (e.g., tail, horns, stomach occlusion).

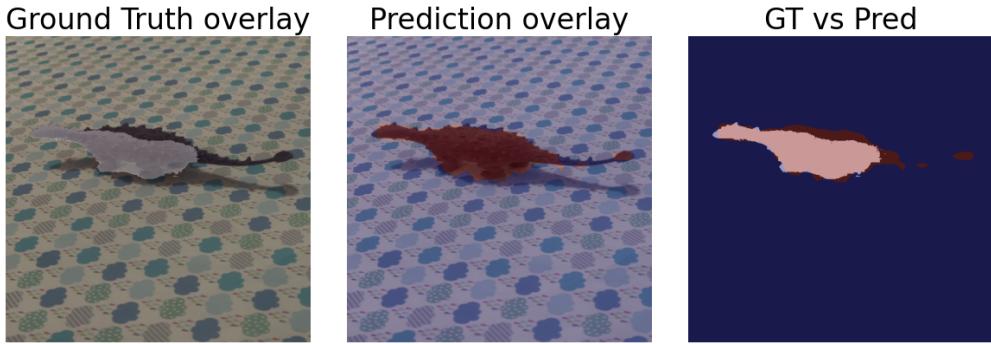


Figure 28: **Overlay comparison of ground truth and prediction.** The same instance as in Figure 27 is shown for better. Left: ground truth overlaid on the input image. Center: prediction overlaid on the input image. Right: difference visualization, with the ground truth overlaid on the predicted mask.

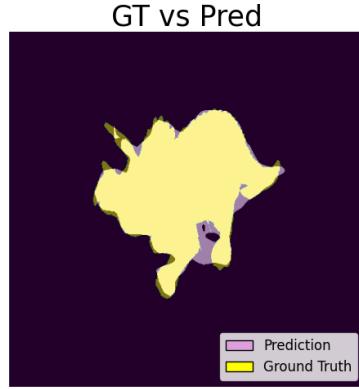


Figure 29: **Another overlay comparison of ground truth and prediction.** Shown using DINO. Yellow: ground truth mask. Purple: predicted mask. We note that discrepancies appear at fine-grained boundaries (e.g., tail, horns, nose) and in occluded or shadowed regions (e.g., under the stomach).

464 **NeurIPS Paper Checklist**

465 **1. Claims**

466 Question: Do the main claims made in the abstract and introduction accurately reflect the
467 paper's contributions and scope?

468 Answer: [Yes]

469 Guidelines:

- 470 • The answer NA means that the abstract and introduction do not include the claims
471 made in the paper.
- 472 • The abstract and/or introduction should clearly state the claims made, including the
473 contributions made in the paper and important assumptions and limitations. A No or
474 NA answer to this question will not be perceived well by the reviewers.
- 475 • The claims made should match theoretical and experimental results, and reflect how
476 much the results can be expected to generalize to other settings.
- 477 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
478 are not attained by the paper.

479 **2. Limitations**

480 Question: Does the paper discuss the limitations of the work performed by the authors?

481 Answer: [Yes]

482 Guidelines:

- 483 • The answer NA means that the paper has no limitation while the answer No means that
484 the paper has limitations, but those are not discussed in the paper.
- 485 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 486 • The paper should point out any strong assumptions and how robust the results are to
487 violations of these assumptions (e.g., independence assumptions, noiseless settings,
488 model well-specification, asymptotic approximations only holding locally). The authors
489 should reflect on how these assumptions might be violated in practice and what the
490 implications would be.
- 491 • The authors should reflect on the scope of the claims made, e.g., if the approach was
492 only tested on a few datasets or with a few runs. In general, empirical results often
493 depend on implicit assumptions, which should be articulated.
- 494 • The authors should reflect on the factors that influence the performance of the approach.
495 For example, a facial recognition algorithm may perform poorly when image resolution
496 is low or images are taken in low lighting. Or a speech-to-text system might not be
497 used reliably to provide closed captions for online lectures because it fails to handle
498 technical jargon.
- 499 • The authors should discuss the computational efficiency of the proposed algorithms
500 and how they scale with dataset size.
- 501 • If applicable, the authors should discuss possible limitations of their approach to
502 address problems of privacy and fairness.
- 503 • While the authors might fear that complete honesty about limitations might be used by
504 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
505 limitations that aren't acknowledged in the paper. The authors should use their best
506 judgment and recognize that individual actions in favor of transparency play an impor-
507 tant role in developing norms that preserve the integrity of the community. Reviewers
508 will be specifically instructed to not penalize honesty concerning limitations.

509 **3. Theory assumptions and proofs**

510 Question: For each theoretical result, does the paper provide the full set of assumptions and
511 a complete (and correct) proof?

512 Answer: [NA]

513 Guidelines:

- 514 • The answer NA means that the paper does not include theoretical results.

- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Guidelines:

- The answer NA means that paper does not include experiments requiring code.

- 568 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
 569
 570 • While we encourage the release of code and data, we understand that this might not be
 571 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
 572 including code, unless this is central to the contribution (e.g., for a new open-source
 573 benchmark).
 574 • The instructions should contain the exact command and environment needed to run to
 575 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
 576
 577 • The authors should provide instructions on data access and preparation, including how
 578 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
 579
 580 • The authors should provide scripts to reproduce all experimental results for the new
 581 proposed method and baselines. If only a subset of experiments are reproducible, they
 582 should state which ones are omitted from the script and why.
 583
 584 • At submission time, to preserve anonymity, the authors should release anonymized
 585 versions (if applicable).
 586 • Providing as much information as possible in supplemental material (appended to the
 587 paper) is recommended, but including URLs to data and code is permitted.

588 **6. Experimental setting/details**

589 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
 590 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
 591 results?

592 Answer: [\[Yes\]](#)

593 Guidelines:

- 594 • The answer NA means that the paper does not include experiments.
 595 • The experimental setting should be presented in the core of the paper to a level of detail
 596 that is necessary to appreciate the results and make sense of them.
 597 • The full details can be provided either with the code, in appendix, or as supplemental
 598 material.

599 **7. Experiment statistical significance**

600 Question: Does the paper report error bars suitably and correctly defined or other appropriate
 601 information about the statistical significance of the experiments?

602 Answer: [\[Yes\]](#)

603 Guidelines:

- 604 • The answer NA means that the paper does not include experiments.
 605 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
 606 dence intervals, or statistical significance tests, at least for the experiments that support
 607 the main claims of the paper.
 608 • The factors of variability that the error bars are capturing should be clearly stated (for
 609 example, train/test split, initialization, random drawing of some parameter, or overall
 610 run with given experimental conditions).
 611 • The method for calculating the error bars should be explained (closed form formula,
 612 call to a library function, bootstrap, etc.)
 613 • The assumptions made should be given (e.g., Normally distributed errors).
 614 • It should be clear whether the error bar is the standard deviation or the standard error
 615 of the mean.
 616 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 617 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 618 of Normality of errors is not verified.
 619 • For asymmetric distributions, the authors should be careful not to show in tables or
 620 figures symmetric error bars that would yield results that are out of range (e.g. negative
 621 error rates).

- 620 • If error bars are reported in tables or plots, The authors should explain in the text how
621 they were calculated and reference the corresponding figures or tables in the text.

622 **8. Experiments compute resources**

623 Question: For each experiment, does the paper provide sufficient information on the com-
624 puter resources (type of compute workers, memory, time of execution) needed to reproduce
625 the experiments?

626 Answer: [Yes]

627 Guidelines:

- 628 • The answer NA means that the paper does not include experiments.
629 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
630 or cloud provider, including relevant memory and storage.
631 • The paper should provide the amount of compute required for each of the individual
632 experimental runs as well as estimate the total compute.
633 • The paper should disclose whether the full research project required more compute
634 than the experiments reported in the paper (e.g., preliminary or failed experiments that
635 didn't make it into the paper).

636 **9. Code of ethics**

637 Question: Does the research conducted in the paper conform, in every respect, with the
638 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

639 Answer: [Yes]

640 Guidelines:

- 641 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
642 • If the authors answer No, they should explain the special circumstances that require a
643 deviation from the Code of Ethics.
644 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
645 eration due to laws or regulations in their jurisdiction).

646 **10. Broader impacts**

647 Question: Does the paper discuss both potential positive societal impacts and negative
648 societal impacts of the work performed?

649 Answer: [NA]

650 Justification: Our work focuses on methodological evaluation of viewpoint generalization in
651 vision encoders. As such, it does not, on its own, have direct positive or negative societal
652 impacts. Any downstream societal considerations would depend on the specific deployment
653 contexts of such encoders, which are outside the scope of this work.

654 Guidelines:

- 655 • The answer NA means that there is no societal impact of the work performed.
656 • If the authors answer NA or No, they should explain why their work has no societal
657 impact or why the paper does not address societal impact.
658 • Examples of negative societal impacts include potential malicious or unintended uses
659 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
660 (e.g., deployment of technologies that could make decisions that unfairly impact specific
661 groups), privacy considerations, and security considerations.
662 • The conference expects that many papers will be foundational research and not tied
663 to particular applications, let alone deployments. However, if there is a direct path to
664 any negative applications, the authors should point it out. For example, it is legitimate
665 to point out that an improvement in the quality of generative models could be used to
666 generate deepfakes for disinformation. On the other hand, it is not needed to point out
667 that a generic algorithm for optimizing neural networks could enable people to train
668 models that generate Deepfakes faster.
669 • The authors should consider possible harms that could arise when the technology is
670 being used as intended and functioning correctly, harms that could arise when the
671 technology is being used as intended but gives incorrect results, and harms following
672 from (intentional or unintentional) misuse of the technology.

- 673 • If there are negative societal impacts, the authors could also discuss possible mitigation
674 strategies (e.g., gated release of models, providing defenses in addition to attacks,
675 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
676 feedback over time, improving the efficiency and accessibility of ML).

677 **11. Safeguards**

678 Question: Does the paper describe safeguards that have been put in place for responsible
679 release of data or models that have a high risk for misuse (e.g., pretrained language models,
680 image generators, or scraped datasets)?

681 Answer: [NA]

682 Guidelines:

- 683 • The answer NA means that the paper poses no such risks.
684 • Released models that have a high risk for misuse or dual-use should be released with
685 necessary safeguards to allow for controlled use of the model, for example by requiring
686 that users adhere to usage guidelines or restrictions to access the model or implementing
687 safety filters.
688 • Datasets that have been scraped from the Internet could pose safety risks. The authors
689 should describe how they avoided releasing unsafe images.
690 • We recognize that providing effective safeguards is challenging, and many papers do
691 not require this, but we encourage authors to take this into account and make a best
692 faith effort.

693 **12. Licenses for existing assets**

694 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
695 the paper, properly credited and are the license and terms of use explicitly mentioned and
696 properly respected?

697 Answer: [Yes]

698 Guidelines:

- 699 • The answer NA means that the paper does not use existing assets.
700 • The authors should cite the original paper that produced the code package or dataset.
701 • The authors should state which version of the asset is used and, if possible, include a
702 URL.
703 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
704 • For scraped data from a particular source (e.g., website), the copyright and terms of
705 service of that source should be provided.
706 • If assets are released, the license, copyright information, and terms of use in the
707 package should be provided. For popular datasets, paperswithcode.com/datasets
708 has curated licenses for some datasets. Their licensing guide can help determine the
709 license of a dataset.
710 • For existing datasets that are re-packaged, both the original license and the license of
711 the derived asset (if it has changed) should be provided.
712 • If this information is not available online, the authors are encouraged to reach out to
713 the asset's creators.

714 **13. New assets**

715 Question: Are new assets introduced in the paper well documented and is the documentation
716 provided alongside the assets?

717 Answer: [Yes]

718 Guidelines:

- 719 • The answer NA means that the paper does not release new assets.
720 • Researchers should communicate the details of the dataset/code/model as part of their
721 submissions via structured templates. This includes details about training, license,
722 limitations, etc.
723 • The paper should discuss whether and how consent was obtained from people whose
724 asset is used.

- 725 • At submission time, remember to anonymize your assets (if applicable). You can either
726 create an anonymized URL or include an anonymized zip file.

727 **14. Crowdsourcing and research with human subjects**

728 Question: For crowdsourcing experiments and research with human subjects, does the paper
729 include the full text of instructions given to participants and screenshots, if applicable, as
730 well as details about compensation (if any)?

731 Answer: [NA]

732 Guidelines:

- 733 • The answer NA means that the paper does not involve crowdsourcing nor research with
734 human subjects.
735 • Including this information in the supplemental material is fine, but if the main contribu-
736 tion of the paper involves human subjects, then as much detail as possible should be
737 included in the main paper.
738 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
739 or other labor should be paid at least the minimum wage in the country of the data
740 collector.

741 **15. Institutional review board (IRB) approvals or equivalent for research with human
742 subjects**

743 Question: Does the paper describe potential risks incurred by study participants, whether
744 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
745 approvals (or an equivalent approval/review based on the requirements of your country or
746 institution) were obtained?

747 Answer: [NA]

748 Guidelines:

- 749 • The answer NA means that the paper does not involve crowdsourcing nor research with
750 human subjects.
751 • Depending on the country in which research is conducted, IRB approval (or equivalent)
752 may be required for any human subjects research. If you obtained IRB approval, you
753 should clearly state this in the paper.
754 • We recognize that the procedures for this may vary significantly between institutions
755 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
756 guidelines for their institution.
757 • For initial submissions, do not include any information that would break anonymity (if
758 applicable), such as the institution conducting the review.

759 **16. Declaration of LLM usage**

760 Question: Does the paper describe the usage of LLMs if it is an important, original, or
761 non-standard component of the core methods in this research? Note that if the LLM is used
762 only for writing, editing, or formatting purposes and does not impact the core methodology,
763 scientific rigorousness, or originality of the research, declaration is not required.

764 Answer: [NA]

765 Guidelines:

- 766 • The answer NA means that the core method development in this research does not
767 involve LLMs as any important, original, or non-standard components.
768 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
769 for what should or should not be described.