# UNSAMV2:
# Self-Supervised Learning Enables Segment Anything at Any Granularity

Junwei Yu    Trevor Darrell    XuDong Wang*
UC Berkeley

Project Page: https://yujunwei04.github.io/UnSAMv2-Project-Page/



**Figure 1.** What is an object? The notion has long been debated: should it follow a model's learned semantics or an annotator's subjective judgment? UNSAMV2 takes a third path, granting users full flexibility to define objectness through promptable segmentation with a single point and a continuous, controllable granularity score. Built on SAM-2 [35], UNSAMV2 introduces a granularity-aware, self-supervised training pipeline based on divide-and-conquer pseudo-labels [47]. Trained on just 6000 unlabeled images, it segments anything from fine-grained parts to holistic objects, achieving state-of-the-art performance across interactive, whole-image, and video segmentation tasks.

## Abstract

*The Segment Anything Model (SAM) family has become a widely adopted vision foundation model, but its ability to control segmentation granularity remains limited. Users often need to refine results manually — by adding more prompts or selecting from pre-generated masks — to achieve the desired level of detail. This process can be ambiguous, as the same prompt may correspond to several plausible masks, and collecting dense annotations across all granularities is prohibitively expensive, making supervised solutions infeasible. To address this limitation, we introduce UNSAMV2 which enables segment anything at any granularity without human annotations. UNSAMV2 extends the divide-and-conquer strategy of UnSAM by discovering abundant mask-granularity pairs and introducing a novel granularity control embedding that enables precise, continuous control over segmentation scale. Remarkably, with only 6K unlabeled images and 0.02% additional parameters, UNSAMV2 substantially enhances SAM-2, achieving segment anything at any granularity across interactive, whole-image, and video segmentation tasks. Evaluated on over 11 benchmarks, UNSAMV2 improves $NoC_{90}$ (5.69 → 4.75), 1-IoU (58.0 → 73.1), and $AR_{1000}$ (49.6 → 68.3), showing that small amounts of unlabeled data with*

*Corresponding author

*granularity-aware self-supervised learning method can unlock the potential of vision foundation models.*

## 1. Introduction

What is an object? This question has long been debated in both vision and cognition. Should an object be defined by a model's learned semantics or by an annotator's subjective judgment? We propose a third perspective: granting users full flexibility to define objectness through a single point prompt and a continuous granularity score.
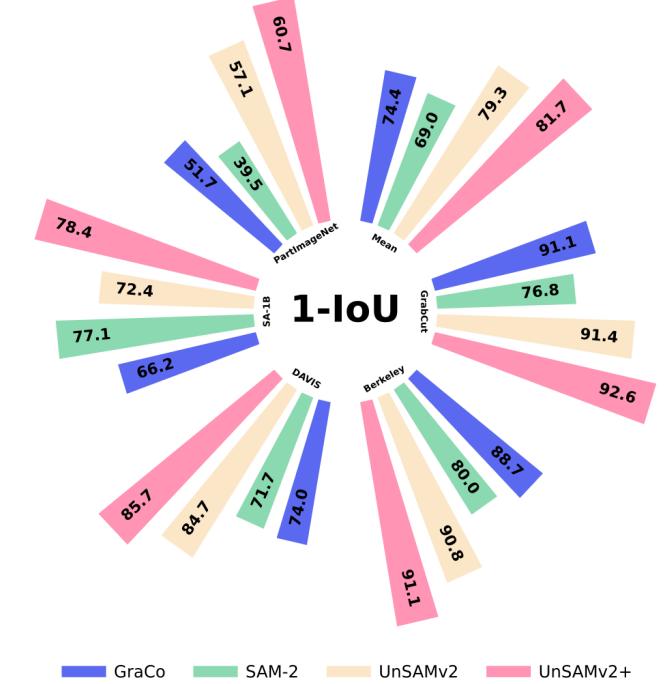
Imagine clicking on a single point in an image and smoothly adjusting a slider that controls segmentation granularity. At low granularity values, the model reveals fine-grained parts with precise boundaries; as the granularity increases, it gradually merges regions into larger and more semantically coherent entities.

In essence, this transforms SAM's [24] three discrete mask hypotheses into a continuous granularity axis capable of producing multiple masks per point, each corresponding to a user-defined granularity. Such controllable segmentation enables users to flexibly define what constitutes an "object" for their specific task, whether it involves part-level analysis, instance grouping, or large-scale region editing.

The emergence of the Segment Anything family [24, 35] has positioned SAM as a vision foundation model, significantly advancing diverse tasks such as video object tracking [51, 56], multi-scale perception [18, 48], and compositional understanding [6, 9, 45]. However, SAM and its successors rely heavily on *supervised learning* from SA-1B dataset [24]. This pipeline directly ties the model's notion of an "object" to human annotation bias and limits its output to three discrete mask hypotheses per prompt. As a result, segmentation in SAM is *defined by supervision* rather than *discovered from data*. This design assumes a shallow hierarchy of objectness, while real-world scenes exhibit complex, nested part–whole structures that cannot be captured through fixed human labels. Although SA-1B covers a broad range of object sizes, it lacks explicit correspondences between instance- and part-level masks, making it difficult for supervised models to learn how granularity should vary continuously.

*We argue that segmentation granularity should be learned through unsupervised learning, which allows models to infer object hierarchies directly from image statistics instead of depending on predefined labels.* To this end, we introduce **UNSAMV2**, a self-supervised framework that enables segmentation at any granularity without human supervised labels. UNSAMV2 learns a continuous representation of granularity that bridges the gap between parts and wholes, giving users full control over segmentation masks.

UNSAMV2 builds upon UnSAM [47], which introduced an unsupervised divide-and-conquer strategy for discover-



**Figure 2. UNSAMV2 achieves state-of-the-art performance across interactive segmentation benchmarks.** Across multiple datasets, UNSAMV2 consistently outperforms SAM-2 and prior methods, turning segmentation into a controllable and interpretable process rather than a fixed prediction.

ing hierarchical masks. While UnSAM focused on unsupervised hierarchy construction, **UNSAMV2** extends this idea to *granularity-controllable segmentation*. In the *divide* stage, we use the normalized-cut method MaskCut [44] to extract instance-level masks. In the *conquer* stage, we recursively merge similar pixels within each instance to discover finer parts, forming hierarchical pseudo-labels that encode relative scale. From these hierarchies, we compute a continuous granularity scalar for each mask, representing its position along the part–whole continuum.

We then augment SAM-2 [35] with a lightweight granularity encoder and a granularity-aware mask token. Given a point prompt and a granularity scalar $g$, UNSAMV2 predicts the mask corresponding to the desired granularity, turning segmentation into a controllable function of scale. Training only the lightweight SAM-2 decoder for four hours with 2 A100 GPUs on 6,000 unlabeled images (0.02% additional parameters) enables smooth interpolation between coarse and fine structures and reveals the latent hierarchy within SAM's feature space.

Across interactive, whole-image, and video segmentation benchmarks, UNSAMV2 consistently surpasses SAM-2 [35] and previous state-of-the-art promptable segmentation methods. Evaluated on more than 11 widely used datasets, including SA-1B, COCO, and PartImageNet, UN-

SAMv2 improves $NoC_{90}$ from 5.69 to 4.75, 1-IoU from 58.0 to 73.1, and $AR_{1000}$ from 49.6 to 68.3, all achieved using only unlabeled data. The resulting model empowers users to define their own notion of objectness and to explore segmentation as a continuous, controllable process rather than a static prediction.
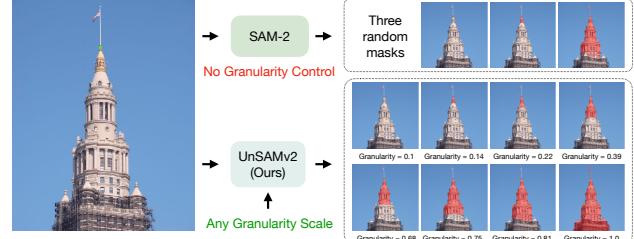
**Contributions.** *(i)* We propose UNSAMV2, a granularity-controllable segmentation framework that enables continuous control of mask granularity from a single point prompt and scalar input. *(ii)* We develop an unsupervised granularity discovery pipeline that learns hierarchical instance–part structures and assigns each mask a continuous scale, applicable to various promptable segmentation models. *(iii)* Trained on only 6,000 unlabeled images, UNSAMV2 achieves state-of-the-art results across interactive, whole-image, and video segmentation benchmarks.

## 2. Related Work

**Multi-Granularity Segmentation.** Segment Anything project [24, 35] has greatly advanced segmentation performance by leveraging large-scale human-annotated data and extensive compute. Extensions such as Semantic-SAM [25] improve fine-grained predictions through a multiple-choice learning design [12, 26]. However, these approaches constrain point-prompt predictions to a fixed number of candidate masks, forcing users to manually select from limited outputs or give additional prompts. This restriction highlights the need for explicit control over mask granularity. Recent work [28, 42] has begun to tackle such ambiguity. GARField [23] and SAMPart3D [54, 55] address scale ambiguity in 3D scene decomposition via absolute scale conditioning, while GraCo [58] achieves granularity-controllable interactive segmentation by extending SimpleClick [29] with discrete granularity inputs.

In contrast, our UNSAMV2 tackles mask ambiguity in a fully self-supervised manner by treating granularity as a continuous, relative concept. We enable granularity-aware segmentation within the widely adopted SAM framework without requiring manual annotation.

**Self-Supervised Learning and Unsupervised Segmentation.** Self-supervised learning (SSL) methods such as MAE [16], JEPA [2], and DINO [5, 32, 40] demonstrate that large-scale pretraining can endow vision transformers with strong semantics-aware representations, benefiting a wide range of downstream tasks [10, 19, 20, 50, 52, 53, 57]. In parallel, unsupervised segmentation has gained lots of attention [1, 13, 14, 22, 31, 39, 43, 49, 60]. CutLER [44], as a recent foundational work of unsupervised image segmentation, greatly advanced unsupervised instance segmentation by introducing MaskCut, a normalized-cuts–based strategy [37] that iteratively extracts multiple objects from images. VideoCutLER [46] extended this framework to video through a cut–synthesize–learn pipeline. CutS3D [38] in-



**Figure 3. From ambiguity to control.** Without granularity input, SAM-2 yields up to three masks per point, requiring users to manually choose one. UNSAMV2 resolves this ambiguity by introducing a continuous granularity variable, allowing users to obtain the intended object at any scale with a single prompt. This simple addition turns segmentation from a discrete guess into a continuous, controllable reasoning process.

troduces the concept of projecting 2D image into 3D space via ZoeDepth [3] to enhance unsupervised segmentation performance on overlapping objects. SOHES [4] adopts a bottom-up merging scheme, grouping pixels based on cosine similarity to progressively discover objects. More recently, UnSAM [47] introduced a divide-and-conquer paradigm to generate hierarchical pseudo labels.

Building on these efforts, UnSAMv2 leverages the hierarchical mask discovery perspective in the divide-and-conquer [47] pipeline and extends it to assign an explicit granularity scale for each pseudo mask. This enables unsupervised segmentation at arbitrary levels of detail, realizing the goal of segment anything at any granularity.

## 3. Method

We present **UNSAMV2**, a self-supervised framework that enables segmentation at arbitrary levels of granularity without human annotations. Unlike the supervised SAM [24] pipeline, which depends on human-labeled object masks, UNSAMV2 learns granularity directly from image statistics through a hierarchy-aware divide-and-conquer process. This enables segmentation granularity to be continuously controlled by a single scalar input, rather than restricted to a fixed number of discrete mask tokens.

We first review prior unsupervised segmentation methods (Sec. 3.1) and the limitations of supervised training paradigms (Sec. 3.2). We then present our granularity-aware divide-and-conquer algorithm that automatically constructs mask–granularity pairs from unlabeled data (Sec. 3.3). Next, we describe the architectural design that empowers any promptable segmentation model to interpret and control segmentation granularity (Sec. 3.4). We then briefly discuss the difference with prior supervised learning works (Sec. 3.5). Finally, we present UNSAMV2+, a lightly supervised variant that integrates SA-1B annotations to further refine the granularity learning process (Sec. 3.6).

3

## 3.1. Preliminaries

**UnSAM, MaskCut, and SOHES.** UnSAM [47] introduced a divide-and-conquer strategy to generate pseudo masks without supervision. In the *divide* stage, a cut-based segmentation method MaskCut [44] is applied to obtain instance/semantic-level masks. Then, inspired by bottom-up hierarchical segmentation, *e.g.*, SOHES [4], the *conquer* stage iteratively merges similar pixels under a sequence of thresholds, constructing part–whole hierarchies. Formally, for a local image region $I_{\text{local}}$, patch-level features $K = \text{DINO}(I_{\text{local}})$ are extracted using DINO [40]. Neighboring patches are merged according to the cosine similarity of their DINO features against thresholds $\theta_1, \ldots, \theta_l$, yielding part-level masks nested inside instances. This process produces a discrete but rich hierarchy of mask granularity.

**Segment Anything family.** Segment Anything models (SAM and SAM-2) [24, 35] have advanced promptable segmentation with a scalable encoder–decoder design. **(1) Image encoder:** a ViT that maps an input image to multi-scale dense embeddings while preserving spatial structure. **(2) Prompt encoder:** embeddings for user inputs such as points, boxes, or masks that condition the segmentation process and guide attention to regions of interest. **(3) Mask decoder:** a lightweight transformer that fuses image and prompt features to predict segmentation masks.

Despite these strengths, the training pipeline is fully *supervised* on SA-1B [24], which ties the notion of objectness to human-labeled masks. Moreover, the decoder employs three fixed mask tokens (small, medium, large), producing at most three hypotheses per prompt. This discretization limits controllable granularity and discourages hierarchical reasoning about parts and wholes, motivating an unsupervised formulation that learns granularity from data rather than from fixed labels.

## 3.2. Limitations of SAM's Supervised Paradigm

**Lack of granularity control.** When a single point corresponds to multiple plausible objects (*e.g.*, a part versus the whole), SAM generates up to three discrete masks and requires manual selection by the user. Without an explicit granularity variable, the model cannot traverse scales continuously—fine details and coarse structures remain disconnected. This limitation not only reduces efficiency in interactive segmentation but also prevents smooth, interpretable control over the level of detail.

**Lack of hierarchical reasoning.** Supervised training on human-labeled masks encourages SAM to learn a flat object representation, where parts and instance/semantic-level segments are treated as isolated entities rather than components within a hierarchy. As a consequence, SAM lacks structural awareness and fails to capture relationships across



**Figure 4. Granularity distribution of discovered masks.** Our divide-and-conquer pipeline produces a rich, left-tailed hierarchy of pseudo-masks, dominated by fine-grained structures. Despite this imbalance, UNSAMV2 learns stable semantics across all scales. Hierarchical perception can emerge from unlabeled data!

scales. It struggles to segment scenes at intermediate levels of detail or to uncover the nested hierarchical structure of visual scenes (Fig. 3). This limitation underscores the necessity of *unsupervised learning*, which can recover hierarchical dependencies directly from image statistics rather than relying on human annotations.

## 3.3. Granularity-Aware Divide-and-Conquer

We extend UnSAM's divide-and-conquer framework to automatically construct dense mask–granularity pairs from unlabeled images. The process consists of four stages.

**Stage 1: Instance discovery via N-Cuts.** We employ the normalized-cut-based method CutLER [44] to generate segmentation masks $\mathcal{M} = \{m_1, \ldots, m_n\}$ at varying levels of granularity. We filter out noisy CutLER mask outputs with a confidence threshold $\tau_{\text{conf}}$,
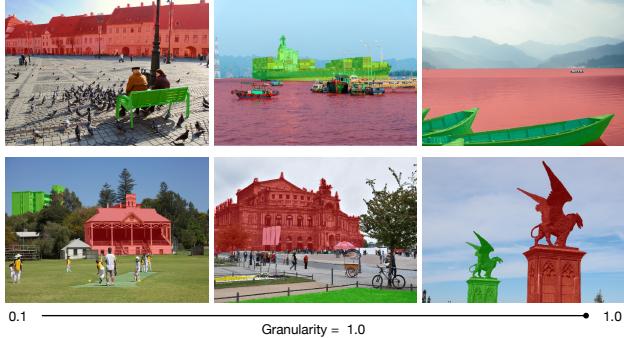
$$\mathcal{M}_{\text{high}} = \{m_i \mid \text{conf}(m_i) \geq \tau_{\text{conf}}\}, \forall i \in [n]. \quad (1)$$

**Stage 2: Instance–part relationship discovery.** We identify instance-level masks $\mathcal{M}_{\text{inst}} \subseteq \mathcal{M}_{\text{high}}$ based on two criteria: (i) their area-to-image-area ratio exceeds the threshold $\tau_{\text{area}}$, and (ii) they dominate overlapping masks, which means for any $m_i \in \mathcal{M}_{\text{inst}}$ and $m_j \in \mathcal{M}_{\text{high}}$, we have

$$\text{IoU}(m_i, m_j) \geq \tau_{\text{overlap}} \Rightarrow \text{Area}(m_i) \geq \text{Area}(m_j). \quad (2)$$

The remaining masks form $\mathcal{M}_{\text{rest}} = \mathcal{M}_{\text{high}} \setminus \mathcal{M}_{\text{inst}}$. Each $m_r \in \mathcal{M}_{\text{rest}}$ is assigned as a part of $m_i$ if $\text{IoU}(m_r, m_i) > \tau_{\text{overlap}}$; otherwise, it is discarded. Therefore, for every instance mask $m_i$, we have a set of part-level masks $\mathcal{M}_{i,\text{part}}$.

**Stage 3: Fine-grained mask discovery.** To enrich granularity, we further apply patch merging to each $m_i \in \mathcal{M}_{\text{inst}}$, decomposing it into finer structures. The resulting fine-grained masks $\mathcal{M}_{i,\text{conquer}}$ are merged with $\mathcal{M}_{i,\text{part}}$ through Non-Maximum Suppression (NMS), denoted as $\mathcal{M}_{i,\text{final}}$.

**Figure 5. Granularity as a relative notion.** At a fixed granularity value, mask sizes vary widely across scenes, showing that UNSAMV2 learns granularity relationally, consistent with human perception of parts and wholes rather than simply associating it with absolute size.

This process expands granularity richness in a hierarchical perspective while maintaining mask quality.

**Stage 4: Continuous granularity assignment.** After obtaining the instance–part hierarchies, we assign each mask $m_i \in \mathcal{M}_{i,\text{final}}$ a continuous granularity score $g_i \in [0.1, 1.0]$ according to its relative area within the hierarchy:

$$g_i = \left( \frac{\sqrt{A_i} - \sqrt{A_{\min}}}{\sqrt{A_{\max}} - \sqrt{A_{\min}}} \right) \cdot 0.9 + 0.1, \qquad (3)$$

where $A_i$ denotes the area of mask $m_i$. As illustrated in Fig. 4, this generates a dense hierarchy of masks spanning fine to coarse scales. We observe that most masks have $g_i < 0.4$, indicating an abundance of part-level structures; yet UNSAMV2 learns stable semantics across all scales (Fig. 5), suggesting robust hierarchical understanding.

### 3.4. UNSAMV2 Architecture

**Granularity encoding.** As shown in Fig. 6, we introduce a lightweight granularity encoding module that converts a scalar $g$ into a high-dimensional Fourier embedding [41], denoted as $\phi(g) \in \mathbb{R}^{d_{\text{Fourier}}}$. Then, we project $\phi(g)$ into the decoder feature space using a three-layer MLP: $E_g = \text{MLP}(\phi(g)) \in \mathbb{R}^{d_{\text{decoder}}}$. This granularity embedding is concatenated with the point prompt features $E_p$: $E_{\text{prompt}} = \text{Concat}(E_p, E_g)$. This design allows the mask decoder to jointly interpret spatial prompts and granularity cues.

**Granularity-aware mask decoding.** We replace SAM's three fixed-size mask tokens with a single learnable *granularity-aware mask token*. This token attends to both prompt embeddings and image features through self- and cross-attention, producing a mask that corresponds to the target granularity. The design is model-agnostic and can be integrated with any promptable segmentation framework.

**Parameter efficiency.** The additional encoding and token modules introduce less than 0.02% extra parameters while providing continuous control over mask detail.

### 3.5. Distinctions with prior methods.

Prior methods such as GraCo [58] and GARField [23] treat granularity as a discrete variable or directly associate it with absolute mask size. However, discretizing granularity imposes artificial boundaries between adjacent levels, preventing smooth transitions across scales, while absolute-size definitions fail to account for contextual differences, *e.g.*, a "small" mask could correspond to an entire object in one image but only a part in another. In contrast, UN-SAMV2 formulates granularity as a ***continuous, relative*** measure within the instance–part hierarchy, aligning more closely with human perception, where objects and parts are interpreted in relation to one another rather than by their absolute scales. This formulation enables UNSAMV2 to capture fine-grained hierarchical structure and traverse segmentation levels seamlessly.
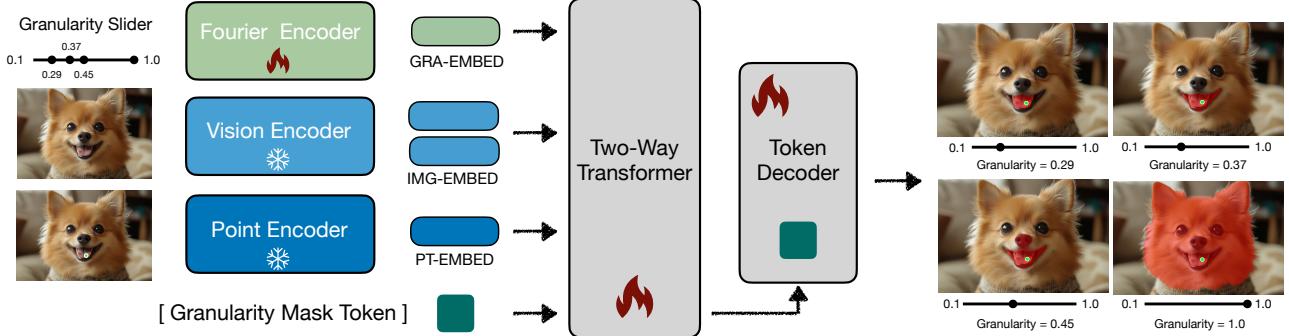
### 3.6. UNSAMV2+: Light-Supervised Variant

Although the unsupervised pipeline generates abundant mask–granularity pairs, pseudo labels can be noisy. To further stabilize learning, we introduce **UNSAMV2+**, a lightly supervised variant that incorporates SA-1B ground-truth masks into the divide stage: $\mathcal{M}_{\text{UNSAMV2+}} = \mathcal{M}_{\text{CutLER}} \cup \mathcal{M}_{\text{SA-1B}}$. The combined masks are then processed through the same conquer and granularity-assignment stages as in UNSAMV2. This hybrid supervision effectively balances human-curated and self-discovered hierarchies, producing cleaner masks while preserving the scalability and flexibility of unsupervised granularity learning.
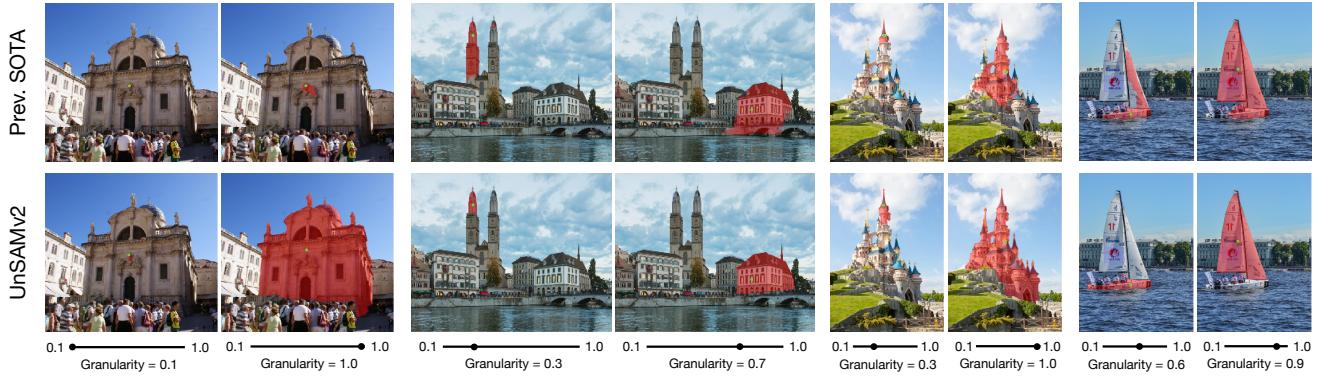
## 4. Experiments

### 4.1. Model Training Settings

**Pseudo mask-granularity data generation.** For the divide stage of our unsupervised data pipeline, we set $\tau_{\text{conf}}$ to 0.3 and $\tau_{\text{overlap}}$ to 0.8. In the conquer stage, we leverage DINOv3 [40] vision transformer to extract feature embedding and set the iterative merging thresholds as $\theta = [0.9, 0.8, 0.7, 0.6, 0.5]$. We run the unsupervised pipeline to generate mask-granularity pairs on 6,000 images from SA-1B. **UNSAMV2 training.** We adopt SAM-2-small as the base model, and with only 8 A-100 GPU hours, we finetune the base model for 5 epochs with 6,000 unlabeled images. During this process, we freeze the vision encoder and train the granularity encoding module, granularity mask token, and the two-way transformer block in the mask decoder with LoRA [17]. See more experiment details in Sec. A2.

5

**Figure 6. Architecture of UNSAMV2.** Built on SAM-2, UNSAMV2 introduces a Fourier-based granularity encoder and a granularity-aware mask token to enable segmentation at arbitrary granularity. A scalar granularity input $g \in [0.1, 1]$ is mapped to a high-dimensional embedding via Fourier transformation and an MLP, then injected into the transformer alongside the sparse point prompt embedding and dense image embedding. The granularity-aware mask token attends to image, point, and granularity embeddings, and is finally decoded by a token decoder into a mask at the requested granularity.



**Figure 7. Qualitative comparison with the previous state-of-the-art method** [58]. Each scene shows results at different target granularity values. Prior methods often break one object into parts at high granularity or include extra regions at low granularity. In contrast, UNSAMV2 produces clear and consistent masks with smooth transitions across scales.

## 4.2. Evaluation Datasets and Metrics

**Interactive Image Segmentation.** Following SimpleClick [29] and GraCo [58], we evaluate UNSAMV2's capability to segment objects at various granularity levels with the number of clicks (NoC) required to reach a certain Intersection over Union (IoU) threshold to assess the model's efficiency in segmenting the target object. Specifically, we choose NoC$_{80}$ and NoC$_{90}$ which record the average number of clicks needed to achieve 80% and 90% IoU between predicted mask and ground-truth mask. In addition, we measure the IoU between the predicted mask and the ground-truth mask with just one click, denoted as 1-IoU in the table. For object levels, we conduct evaluations on 5 commonly used datasets [15, 24, 30, 33, 36], and for part level, we evaluate models on [7, 8]. Note that for the fairness of comparison, we only compare datasets across models that are not in-distribution with their training data, e.g., SAM-2 [35] is trained on SA-1B, SimpleClick [29] and GraCo [58] are trained with SBD and PascalPart. We pro-

vide detailed descriptions of evaluation datasets in Sec. A1.
**Whole-Image Segmentation.** We evaluate our model's performance in identifying all possible masks for given images in a zero-shot setting across 5 datasets that span a broad range of granularity [11, 24, 27, 34, 59]. Importantly, each benchmark annotates only specific semantic classes and typically emphasizes certain levels of the part–instance hierarchy, whereas our method predicts masks at all levels and for any class. Consequently, COCO Average Precision (AP) does not faithfully capture open-world performance. Following prior work [4, 44, 47], we report Average Recall (AR$_{1000}$) for comparison across methods.

## 4.3. Experimental Results

**Interactive Segmentation.** Remarkably, UNSAMV2 surpasses SAM-2 across all datasets in zero-shot evaluation settings as summarized in Table 1. Finetuned with only 6,000 images with unsupervised pseudo-labels, UNSAMV2 demonstrates superior performance in segmenting objects at various granularity levels. It indicates that

| Datasets → | Averaged | | | GrabCut | | | Berkeley | | | SBD | | | PascalPart | | | PartImageNet | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | NoC$_{80\downarrow}$ | NoC$_{90\downarrow}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{85}$ | 1-IoU | NoC$_{80}$ | NoC$_{85}$ | 1-IoU |
| SAM-2 [35] | 4.02 | 5.69 | 58.0 | 1.48 | 1.54 | 76.8 | 1.29 | 1.72 | 80.0 | 2.85 | 6.73 | 64.6 | 9.37 | 11.80 | 29.0 | 5.09 | 6.64 | 39.5 |
| UNSAMV2* | **3.41** | **4.75** | **73.1** | **1.16** | **1.30** | **91.4** | **1.10** | **1.60** | **90.8** | **2.75** | **5.97** | **74.8** | **7.84** | **9.59** | **51.5** | **4.21** | **5.29** | **57.1** |
| *vs. sup. SAM-2* | -0.61 | -0.94 | +15.1 | -0.32 | -0.24 | +14.6 | -0.19 | -0.12 | +10.8 | -0.10 | -0.76 | +10.2 | -1.53 | -2.21 | +22.5 | -0.88 | -1.35 | +17.6 |
| UNSAMV2+* | 3.30 | 4.56 | 74.7 | 1.22 | 1.26 | 92.6 | 1.12 | 1.37 | 91.1 | 2.55 | 5.72 | 77.3 | 7.87 | 9.67 | 52.0 | 3.75 | 4.77 | 60.7 |

**Table 1. Comparison with SAM-2 on interactive segmentation benchmarks.** Trained on 6,000 images with purely unsupervised pseudo-labels, UNSAMV2 significantly outperforms SAM-2. We report promptable segmentation performance in terms of NoC and 1-IoU across five benchmarks. * Following [58], we select the optimal granularity from 0.1 to 1.0 in steps of 0.1 and report averaged results.

| Datasets → | Averaged | | | GrabCut | | | Berkeley | | | DAVIS | | | SA-1B | | | PartImageNet | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | NoC$_{80\downarrow}$ | NoC$_{90\downarrow}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{90}$ | 1-IoU | NoC$_{80}$ | NoC$_{85}$ | 1-IoU |
| SimpleClick [29] | 3.32 | 4.87 | 60.2 | 1.32 | 1.48 | 89.6 | 1.26 | 2.44 | 84.6 | 2.88 | 5.38 | 75.6 | 4.80 | 7.27 | 21.9 | 6.34 | 7.76 | 29.2 |
| GraCo [58]* | 2.35 | 3.42 | 74.4 | 1.24 | 1.32 | 91.1 | 1.17 | 1.63 | 88.7 | 2.84 | 5.09 | 74.0 | 2.41 | 4.01 | 66.2 | 4.11 | 5.03 | 51.7 |
| SAM-2 [35] | 2.44 | 3.63 | 69.0 | 1.48 | 1.54 | 76.8 | 1.29 | 1.72 | 80.0 | 2.51 | 4.50 | 71.7 | **1.82** | 3.76 | 77.1 | 5.09 | 6.64 | 39.5 |
| UNSAMV2* | 2.28 | 3.40 | 79.3 | **1.16** | 1.30 | 91.4 | **1.10** | 1.60 | 90.8 | 2.75 | 4.66 | 84.7 | 2.18 | 4.16 | 72.4 | 4.21 | 5.29 | 57.1 |
| UNSAMV2+* | **2.07** | **3.10** | **81.7** | 1.22 | **1.26** | **92.6** | 1.12 | **1.37** | **91.1** | **2.36** | **4.40** | **85.7** | 1.87 | **3.69** | **78.4** | **3.75** | **4.77** | **60.7** |
| *vs. prev. SOTA* | -0.28 | -0.32 | +7.30 | -0.02 | -0.06 | +1.50 | -0.05 | -0.26 | +2.40 | -0.48 | -0.69 | +11.70 | -0.54 | -0.32 | +12.20 | -0.36 | -0.26 | +9.00 |

**Table 2. State-of-the-art performance on interactive segmentation at various granularity levels.** Trained with only 6,000 images using a combination of supervised and unsupervised labels, UNSAMV2+ demonstrates superior segmentation quality and indicates how effectively unsupervised methods can complement supervised data. Results of [29, 58] are reproduced with official code and checkpoints. *: Following [58], we select optimal granularity from 0.1 to 1.0 with a step of 0.1 and report average results.

| Methods | Avg. | COCO | LVIS | ADE | Entity | SA-1B |
|---|---|---|---|---|---|---|
| SAM [24] | 49.6 | 49.6 | 46.1 | 45.8 | 45.9 | 60.8 |
| UnSAM [47] | 39.2 | 40.5 | 37.7 | 35.7 | 39.6 | 41.9 |
| UnSAM+ [47] | 52.6 | 52.2 | 50.8 | 45.3 | 49.8 | 64.8 |
| UNSAMV2 | 68.3 | 74.5 | 63.8 | 68.4 | 64.3 | 70.6 |
| UNSAMV2+ | **74.1** | **79.1** | **70.3** | **72.7** | **70.3** | **77.9** |
| *vs. prev. SOTA* | +21.5 | +26.9 | +19.5 | +27.4 | +20.5 | +13.1 |

**Table 3. State-of-the-art performance on whole image segmentation.** UNSAMV2 outperform baseline methods [24, 47] on evaluation datasets that contain instances over a wide range of granularity levels. The evaluation metric is $AR_{1000}$. We copy results of SAM and UnSAM from [47]. For UNSAMV2, we aggregate masks generated at granularity levels ranging from 0.1 to 1.0 in increments of 0.1 and filter out low-confidence masks.

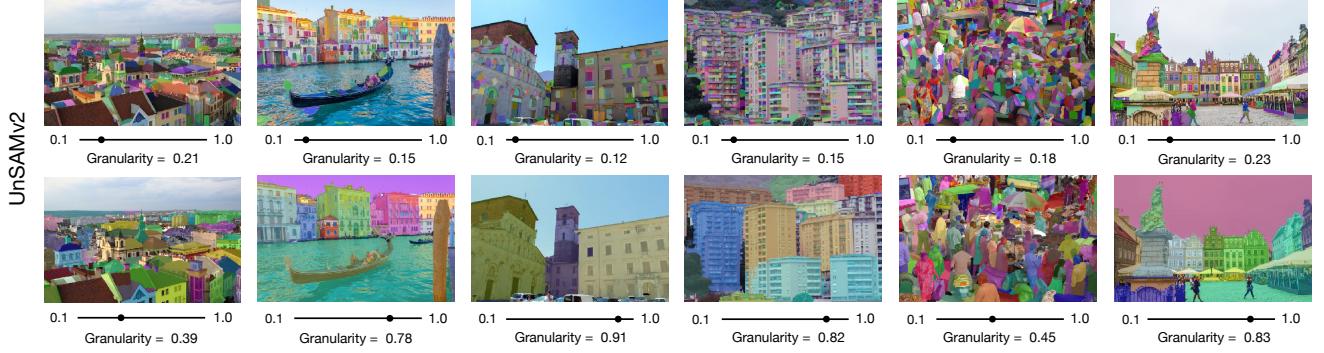| Methods | # Images | PascalPart | | | PartImageNet | | |
|---|---|---|---|---|---|---|---|
| | | NoC$_{80\downarrow}$ | NoC$_{85\downarrow}$ | 1-IoU | NoC$_{80\downarrow}$ | NoC$_{85\downarrow}$ | 1-IoU |
| SAM-2 | – | 9.37 | 11.8 | 29.0 | 5.09 | 6.64 | 39.5 |
| UnSAMv2 | 1K | 7.82 | 9.41 | 49.1 | 4.48 | 5.50 | 55.5 |
| UnSAMv2 | 3K | 7.85 | 9.55 | 51.0 | 4.27 | 5.32 | 55.7 |
| UnSAMv2 | 6K | 7.84 | 9.59 | 51.5 | 4.21 | 5.29 | 57.1 |

**Table 4. Ablations for training data size**. UNSAMV2 starts to demonstrate a decent understanding of granularity scale even trained with only 1,000 images with unsupervised pseudo-labels.

our model architecture and unsupervised data pipeline effectively guide the base model to understand the meaning of granularity scaler. On average, UNSAMV2 surpasses SAM-2 by 15.2% in NoC$_{80}$, 16.5% in NoC$_{90}$, and 26.0% in 1-IoU. With UNSAMV2, user could segment their desired objects accurately at any granularity level with fewer prompt points, which greatly enhances the flexibility of object segmentation and prepares for downstream tasks. In addition, as shown in Table 2, UNSAMV2+, trained with a combination of supervised and unsupervised labels on 6,000 images achieves state-of-the-art performance on all metrics across multiple evaluation datasets. Qualitative comparisons with previous SOTA [58] are shown in Fig. 7.

**Whole-Image Segmentation.** Apart from SOTA results in point-based interactive segmentation, UNSAMV2 achieves superior performance on whole-image segmentation across datasets with instances of abundant granularity levels, out-

performing SAM [24] by 37.7% and UnSAM [47] by 29.8% in $AR_{1000}$. The results are shown in Table 3. With a granularity scalar as input, UNSAMV2 can surface instances over a wide range of detail. As shown in Fig. 8, users simply set the desired granularity to obtain all candidate masks in images at that level, creating new possibilities on how to integrate segmentation models into vision task pipelines.
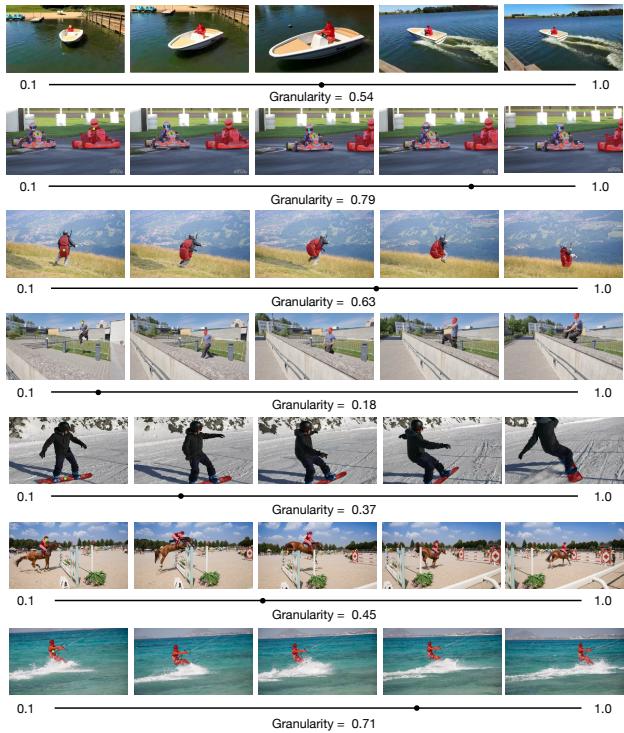
**Video Segmentation Results.** Despite being trained solely on images, UNSAMV2 demonstrates promising capabilities on interactive video segmentation. Although we keep SAM-2's memory module frozen during training, UNSAMV2 still achieves competitive results on video data, demonstrating that the granularity embeddings and mask token propagate across frames effectively. This further shows that our self-supervised pipeline enables granularity to be assimilated into pretrained model's reasoning flow. We show qualitative video segmentation results in Fig. 9.

## 5. Ablations

**Efficiency of UNSAMV2 training.** In Table 4, we ablate the training data size used for UNSAMV2. We find that
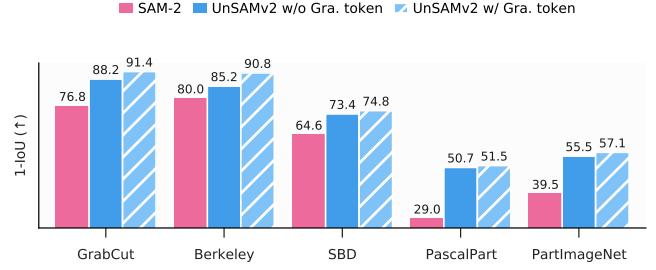
**Figure 8. Visualizations for whole-image segmentation.** Low granularity reveals fine parts, while higher values recover whole objects. UNSAMV2 offers controllable, scalable whole-image segmentation capability, even for scenes with many densely packed entities.



**Figure 9. Granularity generalizes to video.** We prompt UN-SAMV2 with a point and granularity value on the first frame, then propagate the mask to later frames. Even though trained only on images, UNSAMV2 maintains consistent masks over time, showing strong temporal coherence and transferability.



**Figure 10. Ablation of granularity-aware mask token.** Directly finetuning the original SAM-2 mask tokens leads to limited gains, suggesting they already encode strong mask priors. Adding our granularity-aware token enables efficient learning of granularity.

| Methods | Sup. Data | Unsup. Data | Berkeley 1-IoU | DAVIS 1-IoU | PascalPart 1-IoU | PtIn 1-IoU |
|---------|-----------|-------------|----------------|-------------|------------------|------------|
| SAM-2 | – | – | 80.0 | 71.7 | 29.0 | 39.5 |
| UnSAMv2 | ✓ | ✗ | 89.7 | 80.8 | 42.5 | 54.3 |
| UnSAMv2 | ✗ | ✓ | 90.8 | 84.7 | 51.5 | 57.1 |
| UnSAMv2+ | ✓ | ✓ | 91.1 | 85.7 | 52.0 | 60.7 |

**Table 5. Ablation for purely supervised granularity training.** Only training with SA-1B ground-truth labels results in mediocre performance compared to data from our unsupervised pipeline. This indicates how we crucial our unsupervised pseudo-labels are to maximize UNSAMV2's capability in learning granularity.

granularity training is highly sample-efficient: UNSAMV2 already shows a solid grasp of granularity with only 1,000 images. We attribute this efficiency to two factors. First, we derive the granularity scale hierarchically, mirroring how humans perceive scale—by relative size within a hierarchy rather than absolute size. Second, during training we update only $0.1\%$ of the parameters and keep the rest frozen, which preserves pretrained model's segmentation ability. Overall, our procedure effectively teaches pretrained model the concept of mask granularity: with just a few example images, UNSAMV2 learns a latent representation of granularity rather than merely memorizing instances.

**Granularity-aware mask token.** We observe that training UNSAMV2 directly with the original mask tokens and token decoding module in SAM-2 leads to unsatisfying performance (Fig. 10). This phenomenon has been noticed in HQ-SAM [21]. It indicates SAM-2's original mask tokens have already shown strong-prior knowledge on what constitutes as an object. Teaching these pretrained tokens to understand the meaning of granularity is difficult. Thus, we introduce a new mask token and its MLP to do mask decoding that is only trained with mask data accompanied with gran-

| rank → | none | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|
| $NoC_{80\downarrow}$ | 4.47 | 4.41 | 4.21 | 4.40 | 4.38 |

| $N$ → | 3 | 5 | 7 |
|---|---|---|---|
| $NoC_{80\downarrow}$ | 4.21 | 4.31 | 4.23 |

| $N$ → | 20 | 30 | 50 |
|---|---|---|---|
| $NoC_{80\downarrow}$ | 4.41 | 4.21 | 4.25 |

| $d_{fourier}$ → | 32 | 64 | 128 | 256 |
|---|---|---|---|---|
| $NoC_{80\downarrow}$ | 4.44 | 4.30 | 4.21 | 4.37 |

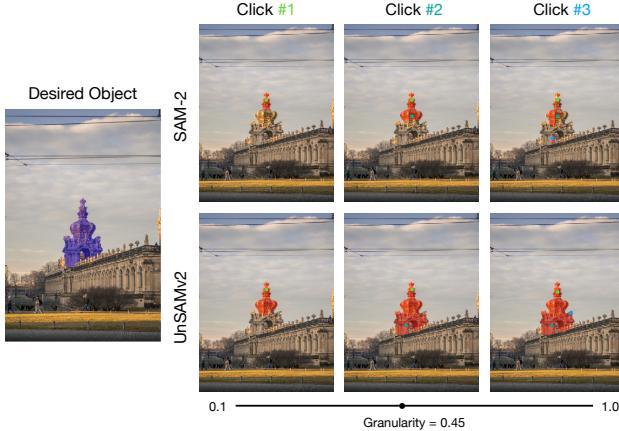**(a) LoRA rank.**  **(b) # points per mask.**  **(c) # masks per image.**  **(d) Fourier dimension for granularity.**

Table 6. **Ablations for hyperparameter and design choices** used for training UNSAMV2. We report UNSAMV2's interactive segmentation performance on validation set of PartImageNet. **(a)** We vary the rank of LoRA weights in mask decoder. **(b)** We vary the number of correction points sampled per mask. **(c)** We vary the number of masks randomly sampled per iteration when training UNSAMV2 models. **(d)** We study different choices of dimension to encode granularity scaler in Fourier feature space. Default settings are highlighted in teal.



Figure 11. **Fewer prompts, more control.** SAM-2 often needs multiple clicks to isolate the target object. With a single granularity value, UNSAMV2 finds the correct mask efficiently, and it can also work with multi-point prompts for finer control.

ularity scales. By conducting self-attention with the granularity embedding encoded by Fourier module and cross-attention with image embeddings, the newly introduced token learn the representation of masks and their corresponding granularity simultaneously. The newly introduced token is a crucial component of UNSAMV2 and a key component to achieve the goal – segment anything at any granularity.

**Unsupervised pseudo-labels** are pivotal for granularity training, as shown in Table 5. We observe that when trained solely on human-labeled SA-1B annotations, UN-SAMV2 performs unsatisfactorily compared with training that also incorporates pseudo-labels from our divide-and-conquer pipeline. This points to an inherent issue with supervised datasets: they are heavily biased by human labelers' notions of what constitutes an object. By contrast, our unsupervised approach focuses on intrinsic relationships among patches, enabling coverage of both instances and parts in a coherent hierarchy.

**Design choices in UnSAMv2 training.** In Table 6, we present the ablation studies on design choices for UN-SAMV2 model architecture and training procedure. We study the use of LoRA in the mask decoder in Table 6a. The results show that, with LoRA, UNSAMV2 learns the concept of granularity efficiently while retaining the strong segmentation capability of SAM-2. Next, we ablate the num-

ber of correction points sampled per mask in one training step in Table 6b. With the granularity scaler, our model can identify the target mask with few clicks, so we use 3 correction points per mask for efficiency. In Table 6c, we study the number of masks per image in one training iteration. UNSAMV2 benefits from a moderate number of masks per step, which best supports learning granularity while maintaining training stability. Finally, we study the effect of the Fourier feature dimension used to encode the granularity input in Table 6d. We observe that a moderate dimension best matches the granularity representation and enables UNSAMV2 to distinguish granularity levels smoothly in a continuous manner.

# 6. Conclusion

We presented UNSAMV2, a self-supervised framework that equips pretrained segmentation model to segment anything at any granularity. By deriving continuous granularity scales from unlabeled data, UNSAMV2 learns to traverse part–whole hierarchies and control segmentation with a single scalar. Trained on only 6K unlabeled images, it achieves state-of-the-art results across interactive, whole-image, and video segmentation. Our results highlight that self-supervised learning can unlock latent hierarchical structure in vision foundation models, transforming segmentation from discrete prediction into continuous, controllable reasoning.

# References

[1] Shahaf Arica, Or Rubin, Sapir Gershov, and Shlomi Laufer. Cuvler: Enhanced unsupervised object discoveries through exhaustive self-supervised transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23105–23114, 2024. 3

[2] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15619–15629, 2023. 3

[3] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 3

[4] Shengcao Cao, Jiuxiang Gu, Jason Kuen, Hao Tan, Ruiyi Zhang, Handong Zhao, Ani Nenkova, Liang-Yan Gui, Tong Sun, and Yu-Xiong Wang. Sohes: Self-supervised open-world hierarchical entity segmentation. *arXiv preprint arXiv:2404.12386*, 2024. 3, 4, 6

[5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3

[6] Danhui Chen, Ziquan Liu, Chuxi Yang, Dan Wang, Yan Yan, Yi Xu, and Xiangyang Ji. Conformalsam: Unlocking the potential of foundational segmentation models in semi-supervised semantic segmentation with conformal prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 24045–24055, 2025. 2

[7] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1971–1978, 2014. 6, 13

[8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 6, 13

[9] Carl Doersch, Pauline Luc, Yi Yang, Dilara Gokay, Skanda Koppula, Ankush Gupta, Joseph Heyward, Ignacio Rocco, Ross Goroshin, Joao Carreira, et al. Bootstap: Bootstrapped training for tracking-any-point. In *Proceedings of the Asian Conference on Computer Vision*, pages 3257–3274, 2024. 2

[10] Zhirui Gao, Renjiao Yi, Yuhang Huang, Wei Chen, Chenyang Zhu, and Kai Xu. Partgs: Learning part-aware 3d representations by fusing 2d gaussians and superquadrics. *arXiv preprint arXiv:2408.10789*, 2024. 3

[11] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019. 6, 13

[12] Abner Guzman-Rivera, Dhruv Batra, and Pushmeet Kohli. Multiple choice learning: Learning to produce multiple structured outputs. *Advances in neural information processing systems*, 25, 2012. 3

[13] Oliver Hahn, Christoph Reich, Nikita Araslanov, Daniel Cremers, Christian Rupprecht, and Stefan Roth. Scene-centric unsupervised panoptic segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24485–24495, 2025. 3

[14] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T Freeman. Unsupervised semantic segmentation by distilling feature correspondences. *arXiv preprint arXiv:2203.08414*, 2022. 3

[15] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 international conference on computer vision*, pages 991–998. IEEE, 2011. 6, 13

[16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 3

[17] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 5

[18] Nan Huang, Wenzhao Zheng, Chenfeng Xu, Kurt Keutzer, Shanghang Zhang, Angjoo Kanazawa, and Qianqian Wang. Segment any motion in videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3406–3416, 2025. 2

[19] Hanwen Jiang, Hao Tan, Peng Wang, Haian Jin, Yue Zhao, Sai Bi, Kai Zhang, Fujun Luan, Kalyan Sunkavalli, Qixing Huang, et al. Rayzer: A self-supervised large view synthesis model. *arXiv preprint arXiv:2505.00702*, 2025. 3

[20] Markus Karmann and Onay Urfalioglu. Repurposing stable diffusion attention for training-free unsupervised interactive segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24518–24528, 2025. 3

[21] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36: 29914–29934, 2023. 8

[22] Chanyoung Kim, Woojung Han, Dayun Ju, and Seong Jae Hwang. Eagle: Eigen aggregation learning for object-centric unsupervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3523–3533, 2024. 3

[23] Chung Min Kim, Mingxuan Wu, Justin Kerr, Ken Goldberg, Matthew Tancik, and Angjoo Kanazawa. Garfield: Group anything with radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21530–21539, 2024. 3, 5

[24] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 2, 3, 4, 6, 7, 13

[25] Feng Li, Hao Zhang, Peize Sun, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, Lei Zhang, and Jianfeng Gao. Segment and recognize anything at any granularity. In *European Conference on Computer Vision*, pages 467–484. Springer, 2024. 3

[26] Zhuwen Li, Qifeng Chen, and Vladlen Koltun. Interactive image segmentation with latent diversity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3

[27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 6, 13

[28] Zheng Lin, Nan Zhou, Chen-Xi Du, Deng-Ping Fan, and Shi-Min Hu. Refcut: Interactive segmentation with reference guidance. *arXiv preprint arXiv:2503.17820*, 2025. 3

[29] Qin Liu, Zhenlin Xu, Gedas Bertasius, and Marc Niethammer. Simpleclick: Interactive image segmentation with simple vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22290–22300, 2023. 3, 6, 7

[30] Kevin McGuinness and Noel E O'connor. A comparative evaluation of interactive segmentation algorithms. *Pattern Recognition*, 43(2):434–444, 2010. 6, 13

[31] Dantong Niu, Xudong Wang, Xinyang Han, Long Lian, Roei Herzig, and Trevor Darrell. Unsupervised universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22744–22754, 2024. 3

[32] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 3

[33] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016. 6, 13

[34] Lu Qi, Jason Kuen, Yi Wang, Jiuxiang Gu, Hengshuang Zhao, Philip Torr, Zhe Lin, and Jiaya Jia. Open world entity segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 6, 13

[35] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 1, 2, 3, 4, 6, 7, 13

[36] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. " grabcut" interactive foreground extraction using iterated graph cuts. *ACM transactions on graphics (TOG)*, 23(3): 309–314, 2004. 6, 13

[37] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000. 3

[38] Leon Sick, Dominik Engel, Sebastian Hartwig, Pedro Hermosilla, and Timo Ropinski. Cuts3d: Cutting semantics in 3d for 2d unsupervised instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21265–21275, 2025. 3

[39] Oriane Siméoni, Gilles Puy, Huy V Vo, Simon Roburin, Spyros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud Marlet, and Jean Ponce. Localizing objects with self-supervised transformers and no labels. *arXiv preprint arXiv:2109.14279*, 2021. 3

[40] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 3, 4, 5, 13

[41] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in neural information processing systems*, 33:7537–7547, 2020. 5

[42] Bin Wang, Anwesa Choudhuri, Meng Zheng, Zhongpai Gao, Benjamin Planche, Andong Deng, Qin Liu, Terrence Chen, Ulas Bagci, and Ziyan Wu. Order-aware interactive segmentation. *arXiv preprint arXiv:2410.12214*, 2024. 3

[43] Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz, Anima Anandkumar, Chunhua Shen, and Jose M Alvarez. Freesolo: Learning to segment objects without annotations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14176–14186, 2022. 3

[44] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3124–3134, 2023. 2, 3, 4, 6

[45] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023. 2

[46] Xudong Wang, Ishan Misra, Ziyun Zeng, Rohit Girdhar, and Trevor Darrell. Videocutler: Surprisingly simple unsupervised video instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22755–22764, 2024. 3

[47] XuDong Wang, Jingfeng Yang, and Trevor Darrell. Segment anything without supervision. *Advances in Neural Information Processing Systems*, 37:138731–138755, 2024. 1, 2, 3, 4, 6, 7, 13

[48] Xuehao Wang, Zhan Zhuang, Feiyang Ye, and Yu Zhang. Mtsam: Multi-task fine-tuning for segment anything model. In *The Thirteenth International Conference on Learning Representations*, 2025. 2

[49] Yangtao Wang, Xi Shen, Yuan Yuan, Yuming Du, Maomao Li, Shell Xu Hu, James L Crowley, and Dominique Vaufreydaz. Tokencut: Segmenting objects in images and videos with self-supervised transformer and normalized cut. *IEEE transactions on pattern analysis and machine intelligence*, 45(12):15790–15801, 2023. 3

[50] Xiaoyang Wu, Daniel DeTone, Duncan Frost, Tianwei Shen, Chris Xie, Nan Yang, Jakob Engel, Richard Newcombe, Hengshuang Zhao, and Julian Straub. Sonata: Self-supervised learning of reliable point representations. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 22193–22204, 2025. 3

[51] Jinyu Yang, Mingqi Gao, Zhe Li, Shang Gao, Fangjing Wang, and Feng Zheng. Track anything: Segment anything meets videos. *arXiv preprint arXiv:2304.11968*, 2023. 2

[52] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024. 3

[53] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv:2406.09414*, 2024. 3

11

[54] Yunhan Yang, Yukun Huang, Yuan-Chen Guo, Liangjun Lu, Xiaoyang Wu, Edmund Y Lam, Yan-Pei Cao, and Xihui Liu. Sampart3d: Segment any part in 3d objects. *arXiv preprint arXiv:2411.07184*, 2024. 3

[55] Yunhan Yang, Yufan Zhou, Yuan-Chen Guo, Zi-Xin Zou, Yukun Huang, Ying-Tian Liu, Hao Xu, Ding Liang, Yan-Pei Cao, and Xihui Liu. Omnipart: Part-aware 3d generation with semantic decoupling and structural cohesion. *arXiv preprint arXiv:2507.06165*, 2025. 3

[56] Mingqiao Ye, Seoung Wug Oh, Lei Ke, and Joon-Young Lee. Entitysam: Segment everything in video. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24234–24243, 2025. 2

[57] Yujia Zhang, Xiaoyang Wu, Yixing Lao, Chengyao Wang, Zhuotao Tian, Naiyan Wang, and Hengshuang Zhao. Concerto: Joint 2d-3d self-supervised learning emerges spatial representations. *arXiv preprint arXiv:2510.23607*, 2025. 3

[58] Yian Zhao, Kehan Li, Zesen Cheng, Pengchong Qiao, Xiawu Zheng, Rongrong Ji, Chang Liu, Li Yuan, and Jie Chen. Graco: Granularity-controllable interactive segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3501–3510, 2024. 3, 5, 6, 7, 14

[59] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. 6, 13

[60] Tianfei Zhou, Jianwu Li, Shunzhou Wang, Ran Tao, and Jianbing Shen. Matnet: Motion-attentive transition network for zero-shot video object segmentation. *IEEE transactions on image processing*, 29:8326–8338, 2020. 3

# UNSAMV2:
# Self-Supervised Learning Enables Segment Anything at Any Granularity

## Supplementary Material

## A1. Evaluation Datasets

**Interactive Segmentation.** Specifically, we adopt the following 7 datasets as benchmarks.

- **GrabCut** [36] consists of 50 images, each containing a single instance.
- **Berkeley** [30] includes 96 images with 100 instances, some of which are more challenging for segmentation.
- **DAVIS** [33] contains 50 high-quality videos and we use 345 frames for evaluation.
- **SA-1B** [24] contains 11 million high-resolution images (average size 3300×4950 pixels) and 1.1 billion high-quality segmentation masks. For evaluation, we randomly choose 1,000 images that are not included in UNSAMV2 training data as our evaluation set.
- **SBD** [15] includes 8,498 training images (with 20,172 instances) and 2,857 validation images (with 6,671 instances).
- **PascalPart** [7] provides part-level annotations for 20 classes from PascalVOC, totaling 193 part categories. We evaluate models on the whole validation set.
- **PartImageNet** [8] organizes 158 ImageNet classes into 11 super-categories and defines 40 distinct part categories. We utilize its validation set to evaluate model's performance on segmenting part-level objects, which consists of 1,206 images and 4,553 annotated parts.

**Whole-Image Segmentation.** We adopt the following 5 datasets to evaluate UNSAMV2's capability on discovering all instances in images.

- **COCO** (Common Objects in Context) [27] is a widely utilized object detection and segmentation dataset. It consists of 115,000 labeled training images and 5,000 labeled validation images. We evaluate our model on COCO Val2017 with 5000 validation images in a zero-shot manner. We use averaged recall ($AR_{1000}$) as the metrics for the whole-image segmentation task.
- **SA-1B** [24] consists of 11 million high-resolution images and 1.1 billion segmentation masks. Following interactive segmentation, we randomly choose 1,000 images that are not included in UNSAMV2 training procedure as evaluation set.
- **LVIS** (Large Vocabulary Instance Segmentation) [11] has 164,000 images with over 1,200 categories and 2 million high-quality instance-level segmentation masks. It covers a large number of object categories. We evaluate UN-SAMV2 using its 5000 validation images.
- **EntitySeg** [34] is an open-world, class-agnostic dataset

with 33,277 images, averaging 18.1 annotated entities per image. We conduct zero-shot evaluation on 1,314 low-resolution images in the validation set.
- **ADE20K** [59] contains 25,574 training and 2,000 testing images covering 365 scenes, emphasizing semantic-level segmentation. It provides labels for 150 semantic categories and 707,868 objects drawn from 3,688 categories. We evaluate zero-shot whole-image segmentation performance on the 2,000-image test split.
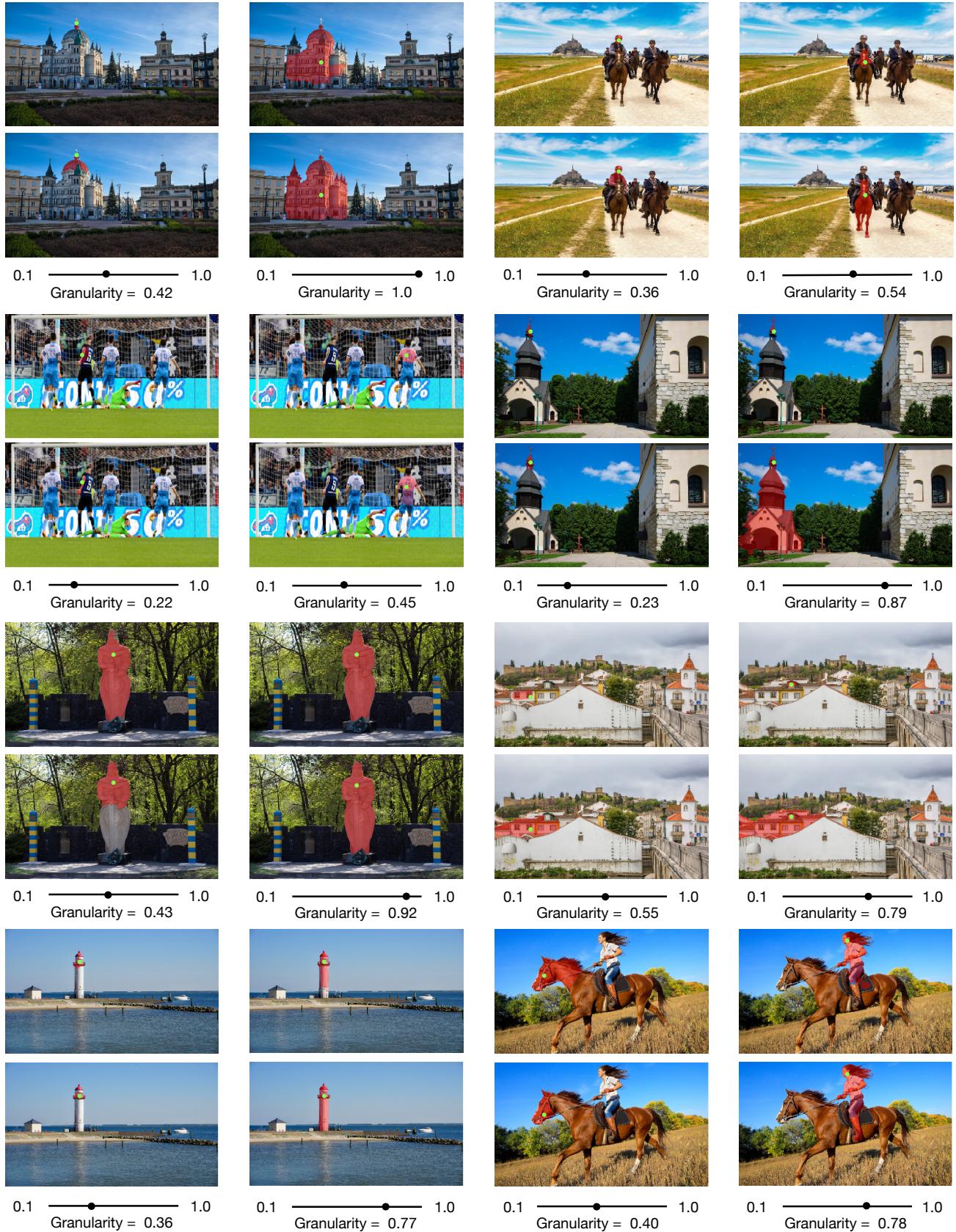
## A2. Training Details.

**Pseudo mask-granularity data generation.** For the divide stage of our unsupervised data pipeline, we set $\tau_{\mathrm{conf}}$ to 0.3 and $\tau_{\mathrm{overlap}}$ to 0.8. In the conquer stage, we leverage DINOv3 [40] ViT-B/16 backbone to extract feature embedding from the last layer of vision transformer and merge adjacent patches together based on predefined cosine similarity thresholds $\theta = [0.9, 0.8, 0.7, 0.6, 0.5]$. We run the pipeline to generate mask-granularity pairs on 6,000 images from SA-1B in a fully unsupervised manner. On average, we have 112 pseudo-labels on each image. Note that unlike UnSAM [47], UNSAMV2 is designed to learn instance–part relationships rather than only instance representations, so our granularity-aware divide-and-conquer pipeline intentionally produces fewer pseudo-labels than UnSAM, which produces 448 pseudo-labels per image.

**UNSAMV2 Training.** We finetune SAM-2-small model for 5 epochs with pseudo-labels on 6,000 images. During finetuning, we freeze the heavy-weight Hiera image encoder and only train the granularity encoding module, granularity mask token, and the two-way transformer block in the mask decoder. For the granularity encoder, we first adopt Fourier transformation to granularity scalar with dimension $d_{\mathrm{Fourier}} = 128$ and followed by a 3-layer MLP. Following SAM-2 [35], UNSAMV2 adopts a combination of focal loss and dice loss with a ratio of 20:1. The batch size is set to 4, with the learning rate initialized at 1e-4. We apply LoRA technology to all projection layers of the transformer, setting the LoRA rank to 8. All experiments are conducted on either 2 A-100 or 4 RTX 3090 GPUs.
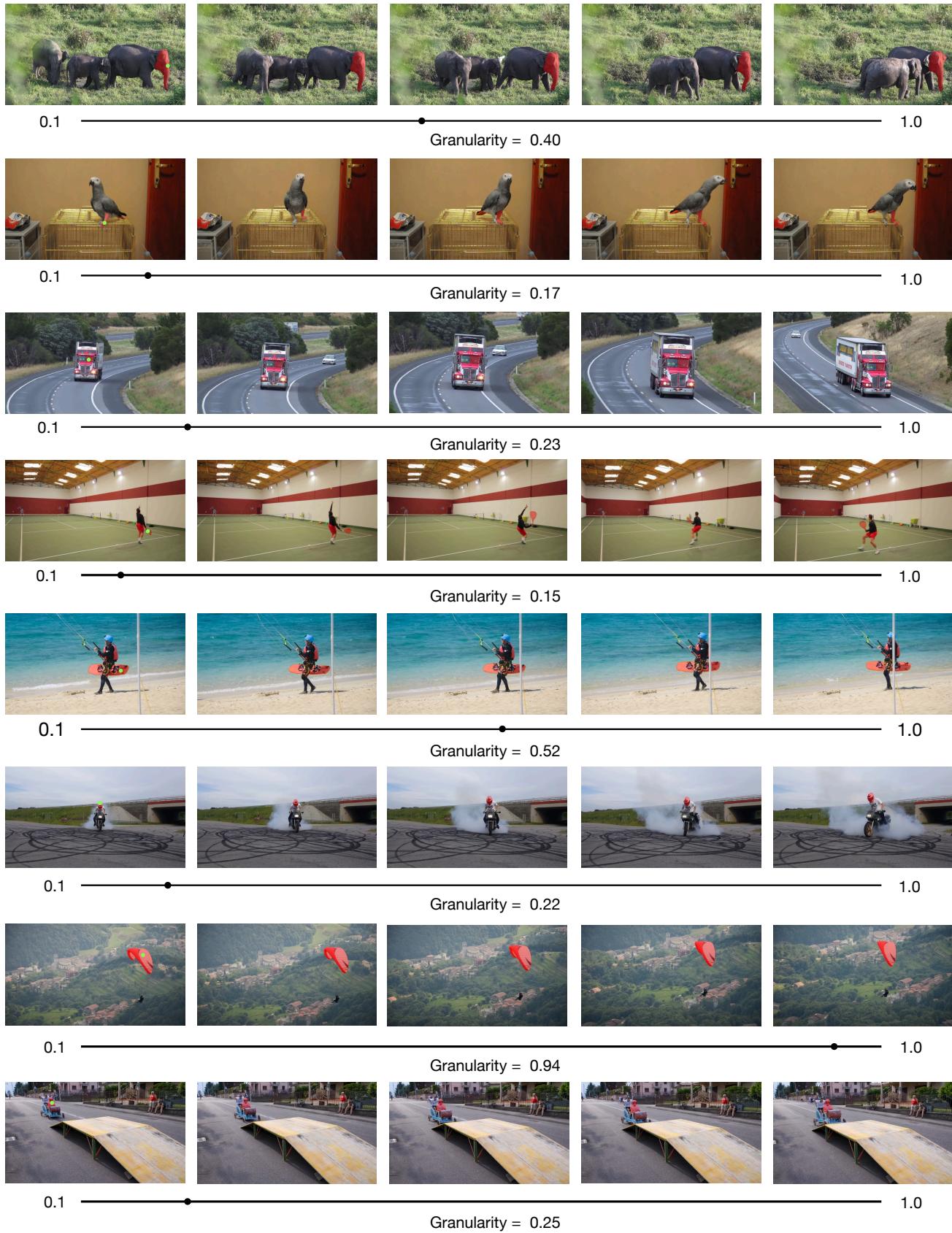
## A3. More Visualizations

We present more UNSAMV2's qualitative results on interactive image segmentation in Fig. A1, whole image segmentation in Fig. A2, and video segmentation in Fig. A3.

**Figure A1.** Interactive segmentation results on SA-1B. The top row is previous SOTA GraCo [58] and the bottom row is UNSAMv2.

Granularity = 0.21

Granularity = 0.15

Granularity = 0.12

Granularity = 0.14

Granularity = 0.39

Granularity = 0.78

Granularity = 0.91

Granularity = 0.82

Granularity = 0.13

Granularity = 0.17

Granularity = 0.15

Granularity = 0.16

Granularity = 0.96

Granularity = 0.91

Granularity = 0.72

Granularity = 0.31

**Figure A2.** Whole image segmentation on SA-1B. From top to bottom are raw images, segmentation by SAM-2 and UNSAMV2.

15

0.1 ———————————————————————————— 1.0

Granularity = 0.40

0.1 ———————————————————————————— 1.0

Granularity = 0.17

0.1 ———————————————————————————— 1.0

Granularity = 0.23

0.1 ———————————————————————————— 1.0

Granularity = 0.15

0.1 ———————————————————————————— 1.0

Granularity = 0.52

0.1 ———————————————————————————— 1.0

Granularity = 0.22

0.1 ———————————————————————————— 1.0

Granularity = 0.94

0.1 ———————————————————————————— 1.0

Granularity = 0.25

**Figure A3.** UNSAMV2's interactive video segmentation results on YoutubeVIS dataset at various granularity levels.