

Unsupervised Universal Image Segmentation

Dantong Niu^{*†} Xudong Wang^{*†} Xinyang Han^{*} Long Lian Roei Herzig Trevor Darrell
Berkeley AI Research, UC Berkeley

Code: <https://github.com/u2seg/U2Seg>

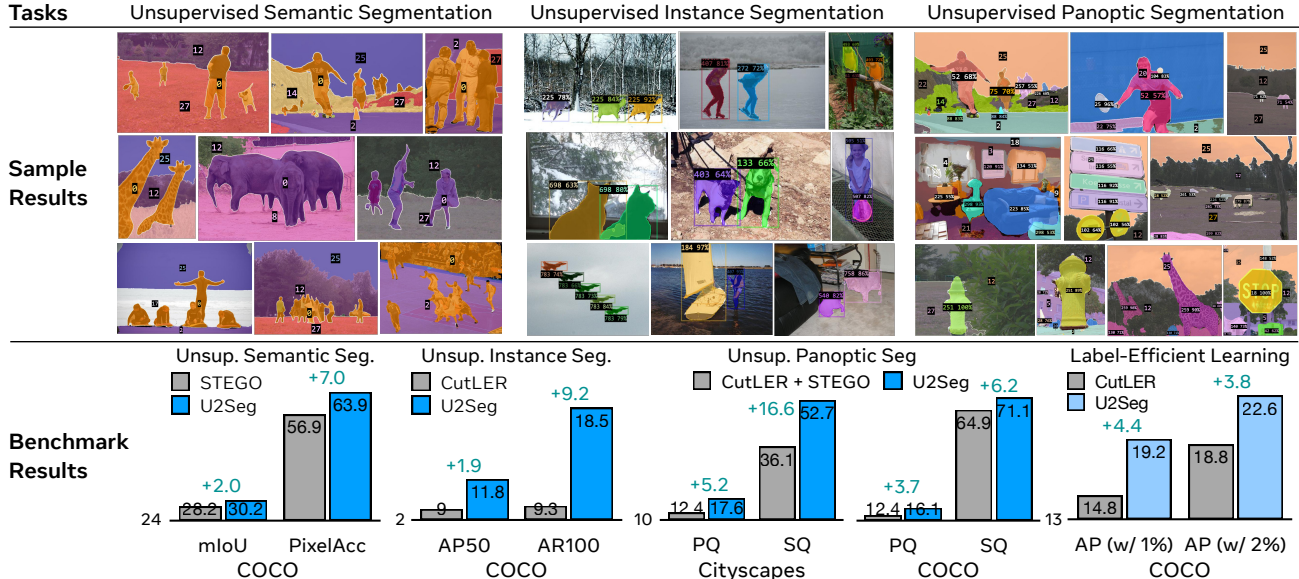


Figure 1. We present **U2Seg**, a unified framework for Unsupervised Universal image Segmentation that consistently outperforms previous state-of-the-art methods designed for individual tasks: CutLER [60] for unsupervised instance segmentation, STEGO [24] for unsupervised semantic segmentation, and the naive combination of CutLER and STEGO for unsupervised panoptic segmentation. We visualize instance segmentation results with “semantic label” + confidence score and semantic predictions with “semantic label”. Zoom in for the best view.

Abstract

Several unsupervised image segmentation approaches have been proposed which eliminate the need for dense manually-annotated segmentation masks; current models separately handle either semantic segmentation (e.g., STEGO) or class-agnostic instance segmentation (e.g., CutLER), but not both (i.e., panoptic segmentation). We propose an Unsupervised Universal Segmentation model (U2Seg) adept at performing various image segmentation tasks—instance, semantic and panoptic—using a novel unified framework. U2Seg generates pseudo semantic labels for these segmentation tasks via leveraging self-supervised models followed by clustering; each cluster represents different semantic and/or instance membership of pixels. We then self-train the model on these pseudo semantic labels, yielding substantial performance gains over specialized

methods tailored to each task: a +2.6 AP^{box} boost (vs. CutLER) in unsupervised instance segmentation on COCO and a +7.0 PixelAcc increase (vs. STEGO) in unsupervised semantic segmentation on COCOStuff. Moreover, our method sets up a new baseline for unsupervised panoptic segmentation, which has not been previously explored. U2Seg is also a strong pretrained model for few-shot segmentation, surpassing CutLER by +5.0 AP^{mask} when trained on a low-data regime, e.g., only 1% COCO labels. We hope our simple yet effective method can inspire more research on unsupervised universal image segmentation.

1. Introduction

The field of image segmentation has witnessed significant advancements in the recent years [4, 5, 12, 20, 26, 35, 38, 40]. Nonetheless, the effectiveness of these segmentation methods heavily depends on the availability of extensive densely human-labeled data for training these models,

^{*}Equal Contribution.

[†]Project Lead.

which is both labor-intensive and costly and thus less scalable. In this paper, our objective is to explore the extent to which unsupervised image segmentation can be achieved without relying on any human-generated labels.

Several recent works such as CutLER [60] and STEGO [24] have emerged as promising approaches for unsupervised image segmentation. CutLER leverages the property of the self-supervised model DINO [8] to ‘discover’ objects without supervision, and learns a state-of-the-art localization model on pseudo instance segmentation masks produced by MaskCut [60] (based on Normalize Cuts [45]). Similarly leveraging DINO [8], STEGO [24] introduces a novel framework that distills unsupervised features into discrete semantic labels. This is achieved using a contrastive loss that encourages pixel features to form compact clusters while preserving their relationships across the corpora [24]. However, these methods have limitations:

- The output of *unsupervised instance segmentation* methods such as CutLER [60] comprises *class-agnostic* segments for “things”, ignoring the “stuff” categories that represent pixel semantics. Moreover, CutLER often treats several *overlapping instances* as one instance, especially when these instances belong to the same semantic class.
- On the other hand, *unsupervised semantic segmentation* methods such as STEGO [24] focus on the segmentation of semantically coherent regions, lacking the capability to distinguish between individual instances.
- *Unsupervised panoptic segmentation* has not been addressed. Supervised panoptic segmentation methods [12, 29, 31] predict both “stuff” and “things” classes simultaneously; to the best of our knowledge there has not been work on unsupervised panoptic segmentation heretofore.

To address these limitations, we propose **U2Seg**, a novel *Unsupervised Universal image Segmentation* model. U2Seg offers comprehensive scene understanding—instance, semantic and panoptic—without relying on human annotations, segmenting semantically meaningful regions in the image as well as identifying and differentiating between individual instances within those regions.

U2Seg is comprised of three steps. First, we create high-quality, discrete semantic labels for instance masks obtained from MaskCut and DINO, by clustering semantically similar instance masks into distinct fine-grained clusters, as described in Sec. 3.2. Next, we amalgamate the semantically pseudo-labeled “things” pixels (from the first step) with “stuff” pixels (from STEGO) to produce pseudo semantic labels for each pixel in the image. Lastly, a universal image segmentation model is trained using these pseudo-labels, resulting in a model capable of simultaneously predicting pixel-level (*i.e.*, semantic segmentation and class-agnostic instance segmentation) and instance-level semantic labels, detailed in Sec. 3.3.

Despite the inherent noise in these pseudo-labels, self-

training the model with them yields substantial performance gains over specialized methods tailored to each task: U2Seg achieves a **+2.6 AP^{box}** boost (*vs.* CutLER) in unsupervised instance segmentation on COCO and a **+7.0 PixelAcc** increase (*vs.* STEGO) in unsupervised semantic segmentation on COCOStuff. Moreover, our method sets up a new baseline for unsupervised panoptic segmentation. We also find that the multi-task learning framework and learning unsupervised segmentor with semantic labels enable our model to generate a more discriminative feature space, which makes it a superior representation for downstream supervised detection and segmentation tasks. When trained on a low-data regime, such as 1% COCO labels, U2Seg surpasses CutLER by **+5.0 AP^{mask}**.

Contributions. Our main contribution is the first universal unsupervised image segmentation model that can tackle unsupervised semantic-aware instance, semantic and panoptic segmentation tasks using a unified framework. We establish a suite of benchmarks on unsupervised semantic-aware instance segmentation and panoptic segmentation, areas previously unexplored. Despite using a single framework, we demonstrate that U2Seg surpasses previous methods specialized for each task across all experimented benchmarks (instance, semantic, panoptic, *etc.*) and datasets (COCO, Cityscapes, UVO, VOC, *etc.*).

2. Related Work

Self-supervised Representation Learning focuses on feature learning from a large amount of unlabeled data without using human-made labels. *Contrastive Learning-Based Methods* [9, 27, 43, 62] learn representation by comparing similar instances or different versions of a single instance while separating dissimilar ones. *Similarity-Based Self-Supervised Learning* [10, 23] mainly reduces differences between different augmented versions of the same instance. *Clustering-Based Feature Learning* [2, 7, 55, 63, 65] finds natural data groups in the hidden space. *Masked Autoencoders* [3, 18, 28] learn by masking and then reconstructing masked parts of the image.

Unsupervised Object Detection and Instance Segmentation. DINO [8] shows that self-supervised learning (SSL) Vision Transformers (ViT) [19] can reveal hidden semantic segmentation in images, which is not obvious in supervised counterparts [8, 66]. Extending this, LOST [46] and TokenCut [61] use DINO’s patch features to identify main objects in images. FreeSOLO [58] performs unsupervised class-agnostic instance segmentation by creating coarse masks first, which are later improved through self-training. MaskDistill [50] uses a self-supervised DINO to get initial masks from an affinity graph but only allows one mask per image during distillation, limiting multi-object detection. Meanwhile, CutLER [59] introduces the MaskCut method, which aims to identify multiple instances in a sin-

gle image. Yet, MaskCut frequently consolidates overlapping instances into a single segment and lacks the capability to assign semantic labels to each instance.

Unsupervised Semantic Segmentation. IIC [30] maximizes mutual information for clustering, while PiCIE [14] uses invariance to photometric effects and equivariance to geometric transformations for segmentation. MaskContrast [49] learns unsupervised semantic segmentation by contrasting features within saliency masks. STEGO [24] refines pretrained SSL visual features to distill correspondence information embedded within these features, thereby fostering discrete semantic clusters.

Universal Segmentation has been introduced to deliver instance, semantic and panoptic segmentation tasks using a unified architecture [6, 11–13, 29, 33, 34, 37, 52, 64]. In this work, we propose U2Seg to tackle this challenging task without relying on human-annotated data.

Unsupervised Image Classification methods mainly focus on providing a semantic label for each query image that can be mapped to ground truth classes by hungarian matching. SCAN [48] proposes a three-stage pipeline that includes representation learning, deep clustering, and self-labeling. NNM [16] enhances SCAN by incorporating local and global nearest neighbor matching. RUC [44] further improves SCAN using a robust loss as training objective. However, these approaches only provide one classification prediction per image, whereas our method provides classification per-instance for instance segmentation and per-pixel for semantic segmentation.

3. Unsupervised Universal Segmentation

3.1. Preliminaries

We first explain the previous Unsupervised Instance Segmentation method CutLER [60], and Unsupervised Semantic Segmentation method STEGO [24].

CutLER [60] exploits self-supervised learning models like DINO [8] to ‘discover’ objects and train a state-of-the-art detection and segmentation model using a cut-and-learn pipeline. It first uses MaskCut to extract multiple initial masks from DINO [8] features. MaskCut first generates a patch-wise affinity matrix $W_{ij} = \frac{K_i K_j}{\|K_i\|_2 \|K_j\|_2}$ using the “key” features K_i for patch i from DINO’s last attention layer. Subsequently, the cut-based clustering method Normalized Cut [45] is employed on the affinity matrix by finding the eigenvector x that corresponds to the second smallest eigenvalue. A foreground instance mask M^s is derived through bi-partitioning of the vector x , enabling segmentation of individual objects in the image. For multi-instance segmentation, MaskCut iteratively refines the affinity matrix by masking out already segmented objects, allowing for

subsequent extractions

$$W_{ij}^t = \frac{(K_i \prod_{s=1}^t M_{ij}^s)(K_j \prod_{s=1}^t M_{ij}^s)}{\|K_i\|_2 \|K_j\|_2} \quad (1)$$

and repeating above steps by N times. CutLER then refines detection and segmentation through a loss-dropping strategy and iterative self-training.

STEGO [24] harnesses the semantically rich feature correlations produced by unsupervised methods like DINO [8] for segmentation. It trains a segmentation head to refine these correlations within an image, with its K-Nearest Neighbors (KNNs), and across randomly chosen images. Specifically, STEGO distills DINO’s unsupervised features into distinct semantic labels by optimizing a correspondence loss. This loss function measures the feature correspondences F^{SC} between image feature pairs generated by DINO and the feature correspondence S_{hwij} derived from a trainable, lightweight segmentation head [24]:

$$L_{\text{corr}}(x, y, b) = - \sum_{hwij} (F_{hwij}^{SC} - b) \max(S_{hwij}, 0) \quad (2)$$

3.2. Unsupervised Instance Segmentation

Although CutLER [60] provides high-quality instance segmentation masks without human annotations, the predicted masks are class-agnostic, and thus do not include semantic labels for each instance. Our method addresses this issue by grouping the detected instances with a clustering method. In this way, instances assigned to the same cluster are associated with identical or closely related semantic information, while instances residing in separate clusters exhibit semantic dissimilarity.

Pseudo Semantic Labels. To train a detection and instance segmentation model, we vector quantize the model targets (pseudo semantic labels) by clustering the instance-level features of the entire dataset, under constraints derived from self-supervision. Specifically, our approach starts with the generation of instance segmentation masks using MaskCut [60]. Subsequently, we utilize the efficient K -Means clustering method as implemented in USL [57] to cluster all segmentation masks into semantically meaningful clusters.

We employ K -Means clustering to partition n instances into $C (\leq n)$ clusters, where each cluster is represented by its centroid c [22, 41]. Each instance is assigned to the cluster with the nearest centroid. Formally, we conduct a C -way node partitioning, denoted as $\mathcal{S} = S_1, S_2, \dots, S_C$, that minimizes the within-cluster sum of squares [36]:

$$\min_{\mathcal{S}} \sum_{i=1}^C \sum_{V \in S_i} |V - c_i|^2 = \min_{\mathcal{S}} \sum_{i=1}^C |S_i| \text{Var}(S_i) \quad (3)$$

This optimization process is carried out iteratively using the EM algorithm [42], starting from selecting random samples as initial centroids. As a result, this process assigns

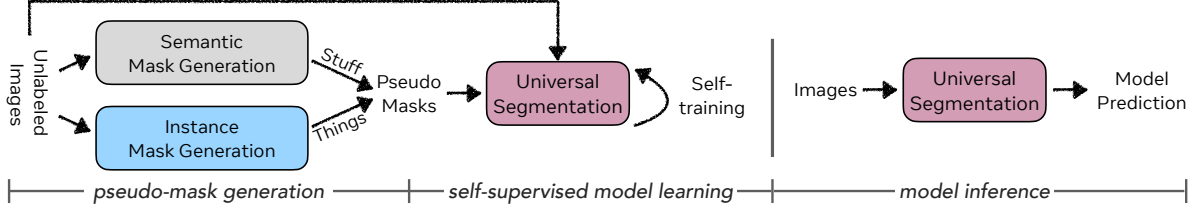


Figure 2. Overview of the training and inference pipeline for the proposed Unsupervised Universal Segmentation model (U2Seg) adept at performing various image segmentation tasks—instance, semantic and panoptic—using a novel unified framework.

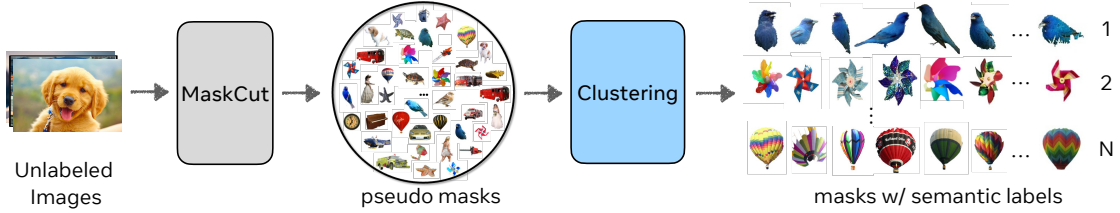


Figure 3. Pipeline overview for generating masks and their semantically meaningful pseudo labels in semantic-aware instance segmentation. We first use MaskCut to generate class-agnostic instance masks, which are then grouped into semantically meaningful clusters. These pseudo semantic labels are used for training a semantic-aware instance segmentor.

pseudo semantic labels, denoted as y_i , to each instance i , with y_i falling within the range of $[1, C]$.

The resulting semantic labels serve multiple purposes: **1) Semantic-aware copy-paste augmentation**, which significantly improves CutLER’s capability to differentiate overlapping instances, especially when they share similar semantic information. **2) Training instance segmentation models**: They serve as pseudo ground-truth labels for training a non-agnostic instance segmentor.

Semantic-aware Copy-Paste Augmentation. In cluttered natural scenes, previous unsupervised instance segmentation model often fail to distinguish instances from the same semantic class. This results in multiple instances being captured in the same mask. To distinguish multiple overlapping objects and small objects in existing unsupervised detectors, we employ semantic-aware copy-paste augmentation, which includes several steps:

1) We begin by randomly selecting two instances, denoted as I_1 and I_2 , both belonging to the same pseudo-category (or group/cluster). **2)** One of these instances undergoes a transformation function \mathcal{T} , which randomly resizes and shifts the associated pseudo-masks. **3)** The resized instance is then pasted onto another image, creating synthetic overlapping cases using the following equation:

$$I_3 = I_1 \cdot (1 - \mathcal{T}(M_c)) + I_2 \cdot \mathcal{T}(M_c) \quad (4)$$

where \cdot denotes element-wise multiplication.

Learning Unsupervised Instance Segmentor. Traditionally, unsupervised segmentation community focused primarily on class-agnostic instance segmentation [58, 60, 61], whose outputs lack class labels. However, by incorporating clustering information obtained from pseudo-labels on Im-

ageNet, as discussed above, our method allows the model to predict not only the location and segmentation of an object but also its pseudo semantic labels.

As observed by [60], “ground-truth” masks may miss instances. However, a standard detection loss penalizes predicted regions r_i that do not overlap with the “ground-truth”. Therefore, following [60], we drop the loss for each predicted region r_i that has a maximum overlap of τ^{IoU} with any of the ‘ground-truth’ instances: $\mathcal{L}_{\text{drop}}(r_i) = \mathbb{1}(\text{IoU}_i^{\text{max}} > \tau^{\text{IoU}}) \mathcal{L}_{\text{vanilla}}(r_i)$, where $\text{IoU}_i^{\text{max}}$ denotes the maximum IoU with all ‘ground-truth’ for r_i and $\mathcal{L}_{\text{vanilla}}$ is the vanilla loss function of detectors. $\mathcal{L}_{\text{drop}}$ encourages the exploration of image regions missed in the “ground-truth”.

3.3. Unsupervised Universal Image Segmentation

Pseudo Labels for Panoptic Segmentation. For each pixel (i, j) in the image, we vector quantize pixels with different semantics or instance membership, generating pseudo semantic labels for panoptic segmentation. We assign each pixel a semantic label based on “stuff” or “things” identity. This results in an instance label $(I(i, j))$ for “things” or a semantic label $(S(i, j))$ for “stuff”. The critical challenge in this process is distinguishing between pixels associated with “things” (countable, often foreground) and “stuff” (uncountable, usually background) [1].

To resolve this problem, our method unfolds in three steps: **1) Semantic Labeling for “Things”**: Utilizing the class-agnostic instance segmentation capabilities of CutLER [60], we first identify “things” within an image, generating class-agnostic instance masks. These masks then undergo deep clustering to attribute a semantic label $I_C(i, j)$ to each instance, detailed in Sec. 3.2. **2) Semantic Labeling**

Task → Datasets → Metric →	Agn Instance Seg.		Instance Seg.				UVO		Semantic Seg.		Panoptic Seg.					
	COCO		COCO		VOC		UVO		COCO		COCO			Cityscapes		
	AP ^{box}	AP ₅₀ ^{box}	AP ₅₀ ^{box}	AR ₁₀₀ ^{box}	AP ₅₀ ^{box}	AR ₁₀₀ ^{box}	AP ₅₀ ^{box}	AR ₁₀₀ ^{box}	PixelAcc	mIoU	PQ	SQ	RQ	PQ	SQ	RQ
FreeSOLO [58]	9.6	4.2	-	-	-	-	-	-	-	-	-	-	-	-	-	-
TokenCut [61]	5.8	3.2	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CutLER [59]	21.9	12.3	-	-	-	-	-	-	-	-	-	-	-	-	-	-
DINO [8]	-	-	-	-	-	-	-	-	30.5	9.6	-	-	-	-	-	-
PiCIE + H [14]	-	-	-	-	-	-	-	-	48.1	13.8	-	-	-	-	-	-
STEGO [24]	-	-	-	-	-	-	-	-	56.9	28.2	-	-	-	-	-	-
CutLER+	-	-	9.0	10.3	26.8	27.2	10.6	11.8	-	-	-	-	-	-	-	-
CutLER+STEGO	-	-	-	-	-	-	-	-	-	-	12.4	64.9	15.5	12.4	36.1	15.2
U2Seg	22.8	13.0	11.8	21.5	31.0	48.1	10.8	25.0	63.9	30.2	16.1	71.1	19.9	17.6	52.7	21.7
vs. prev. SOTA	+0.9	+0.7	+2.8	+11.2	+4.2	+20.9	+0.2	+13.2	+7.0	+2.0	+3.7	+6.2	+4.4	+5.2	+16.6	+6.5

Table 1. With a unified framework, U2Seg outperforms previous state-of-the-art methods tailored for individual tasks across various datasets, including CutLER for unsupervised instance segmentation, STEGO for unsupervised semantic segmentation, and CutLER+STEGO for unsupervised panoptic segmentation. ‘‘Agn Instance Seg’’ denotes class-agnostic instance segmentation.

for ‘‘Stuff’’: For ‘‘stuff’’ pixels, we deploy the unsupervised semantic segmentation model STEGO [24], which distills DINO’s unsupervised features into discrete semantic labels, as outlined in 3.1. This step assigns a ‘‘stuff’’ semantic label to all pixels, including those of ‘‘Things’’ identified earlier. 3) *Integrating Labels for ‘‘Things’’ and ‘‘Stuff’’*. We determine a pixel’s classification as ‘‘things’’ or ‘‘stuff’’ using the following logic:

$$I(i, j) = \begin{cases} I_C(i, j), & \text{if } I_C(i, j) \neq 0 \\ S_S(i, j), & \text{if } I_C(i, j) = 0 \text{ \& } S_S(i, j) \neq 0 \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

This process merges the semantic labels, assigning priority to ‘‘things’’ labels over ‘‘stuff’’ where applicable. We then train a universal segmentation model on these pseudo-labels for instance, semantic and panoptic segmentation tasks.

Learning Unsupervised Universal Image Segmentor. After we obtain the pseudo labels for panoptic segmentation, following [33], we construct an unsupervised universal image segmentation model, that has two branches: instance segmentation branch and semantic segmentation branch, to address corresponding segmentation tasks. The model is trained jointly for both branches, employing the following loss function: $\mathcal{L} = \lambda_i(\mathcal{L}_c + \mathcal{L}_b + \mathcal{L}_m) + \lambda_s \mathcal{L}_s$, where \mathcal{L}_c represents the classification loss, \mathcal{L}_b is the detection loss, \mathcal{L}_m is the segmentation loss, and \mathcal{L}_s signifies the semantic loss. The \mathcal{L}_s is computed as a per-pixel cross-entropy loss between the predicted and ground-truth labels. The hyperparameters λ_i and λ_s balance these two parts.

4. Experiments and Results

4.1. Experimental Setup

Training Data. Our model is trained on 1.3M unlabeled images from ImageNet [17] and is evaluated directly across various benchmarks, unless otherwise noted. For unsupervised semantic segmentation comparisons with

STEGO [24], we additionally fine-tune our model using MSCOCO’s unlabeled images, following STEGO [24].

Test Data. For unsupervised instance segmentation, we test our model on COCO val2017, PASCAL VOC val2012 [21] and UVO val [53]. For unsupervised panoptic segmentation, we evaluate our model on COCO val2017 and Cityscapes val [15].

Evaluation Metrics. We use AP, AP₅₀, AP₇₅ and AR₁₀₀ to evaluate the unsupervised instance segmentation; PixelAcc and mIoU for unsupervised semantic segmentation; PQ, RQ, SQ for unsupervised universal image segmentation. After predicting the instance with its semantic labels, we use Hungarian matching to map the semantic labels to class names in the real dataset (details in B). It evaluates the consistency of the predicted semantic segments with the ground truth labels, remaining unaffected by any permutations in the predicted class labels.

Implementation Details. Following [32], we employ Panoptic Cascade Mask R-CNN [4, 32] with a ResNet50 backbone [25]. Following CutLER’s training recipe [60], our model, initialized with DINO pre-trained weights, is trained on unlabeled ImageNet for two epochs. It starts with an initial learning rate of 0.01, which then decreases after the first epoch to 5×10^{-5} , with a batch size of 16 for all models. For unsupervised panoptic segmentation, we maintain the same training schedule as unsupervised instance segmentation for zero-shot evaluation. In non-zero-shot scenarios, the models undergo training on a combination of unlabeled COCO and ImageNet datasets, beginning with a learning rate of 0.01 over 90k steps.

4.2. Unsupervised Universal Image Segmentation

To the best of our knowledge, U2Seg represents the first effort in addressing unsupervised *semantic-aware* instance, semantic and panoptic segmentation, all unified under a single framework. Due to the absence of benchmarks for unsupervised semantic-aware instance segmentation and panop-

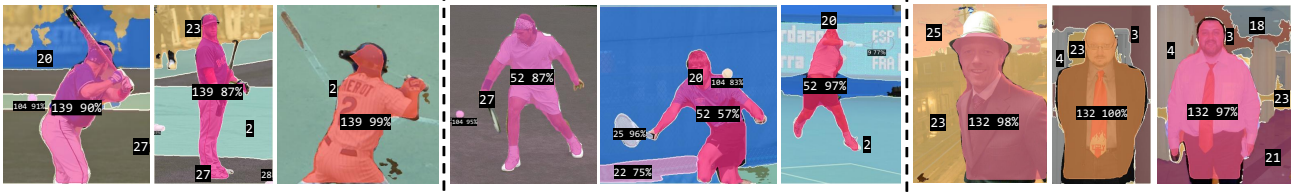


Figure 4. **Universal image segmentation visualization** in COCO val2017. We show the results with predicted cluster ID directly from the model, with athletes playing hockey (left columns) as “139”, playing badminton (middle columns) as “52” and the gentlemen (right columns) as “132”. After Hungarian matching, these IDs are automatically matched to the category “person” for subsequent quantitative evaluation.

Metric	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	AR ₁₀₀ ^{box}	AP ^{mask}	AP ₅₀ ^{mask}	AP ₇₅ ^{mask}	AR ₁₀₀ ^{mask}
CutLER+	5.9	9.0	6.1	10.3	5.3	8.6	5.5	9.3
U2Seg	7.3	11.8	7.5	21.5	6.4	11.2	6.4	18.5
Δ	+1.4	+2.8	+1.4	+11.2	+1.1	+2.6	+0.9	+9.2

Table 2. The results for **zero-shot unsupervised object detection and instance segmentation** on COCO val2017. The model is trained on ImageNet with a cluster number of 800. We compare it with CutLER+, a combination of CutLER and offline clustering.

tic segmentation, we establish comprehensive benchmarks and baselines for both tasks.

In Tab. 1, we demonstrate that U2Seg, utilizing a unified framework, significantly outperforms all previous approaches across various benchmarks and datasets. For **class-agnostic unsupervised instance segmentation**, our method achieves an increase of **+0.9** in AP^{box} compared to CutLER [60]. This improvement is largely attributed to our novel semantic-aware copy-paste augmentation, as detailed in Sec. 3.2. For **unsupervised semantic-aware instance segmentation**, we benchmark against the advanced baseline CutLER+, derived from CutLER, and record a substantial gain of over 11.2% in AR. A more comprehensive analysis of these results is provided in Sec. 4.3. For **unsupervised semantic segmentation**, our approach surpasses the state-of-the-art STEGO with impressive margins of **+7.0** in PixelAcc and **+2.0** in mIoU. Lastly, for **unsupervised panoptic segmentation**, we compare against the strong baseline of CutLER+STEGO, a hybrid of CutLER+ and STEGO, and observe performance gains of over 6.2% in SQ on MSCOCO and a notable 16.6% improvement in SQ on Cityscapes. Further comparisons and discussions on this task are elaborated in Sec. 4.4.

4.3. Unsupervised Instance Segmentation

We performed extensive experiments for zero-shot unsupervised instance segmentation. Given that prior methods [51, 58, 60, 61] are limited to class-agnostic instance segmentation, we developed CutLER+, a strong baseline for unsupervised semantic-aware instance segmentation, building upon the current state-of-the-art CutLER [60]. CutLER+ operates in two steps: it first uses the pre-trained

Methods	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	AR ₁₀₀ ^{box}
CutLER+	17.1	26.8	18.1	27.2
U2Seg	19.0	31.0	19.5	48.1
Δ	+1.9	+4.2	+1.4	+20.9

Table 3. The results for **zero-shot unsupervised object detection on PASCAL VOC val2012**. The model is trained on ImageNet with a cluster number of 800. We compare it with CutLER+, a combination of CutLER and offline clustering.

Metric	AP ^{box}	AP ₅₀ ^{box}	AR ₁₀₀ ^{box}	AP ^{mask}	AP ₅₀ ^{mask}	AR ₁₀₀ ^{mask}
CutLER+	6.3	10.6	11.8	6.0	9.0	10.4
U2Seg	6.8	10.8	25.0	6.2	9.5	21.0
Δ	+0.5	+0.2	+13.2	+0.2	+0.5	+10.6

Table 4. The results for **zero-shot unsupervised object detection and instance segmentation on UVO val1**. The model is trained on ImageNet with a cluster number of 800. We compare with CutLER+, a combination of CutLER and offline clustering.

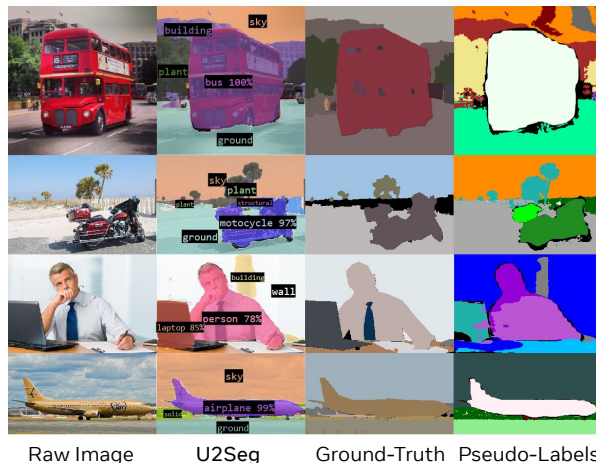


Figure 5. Visualizations of U2Seg’s **unsupervised Panoptic segmentation** results on COCO val2017 (after Hungarian matching). The pseudo label is the naive combination of previous state-of-the-art instance segmentation, *i.e.* CutLER [60], and semantic segmentation, *i.e.*, STEGO [24], results.

CutLER to generate class-agnostic instance masks, and subsequently assigns semantic labels to all instance masks through offline clustering.

Tab. 2 demonstrates that U2Seg markedly improves

Methods	Pretrain	PQ	SQ	RQ
<i>zero-shot methods</i>				
U2Seg	IN	15.7	46.6	19.8
<i>non zero-shot methods</i>				
CutLER+STEGO	COCO	12.4	36.1	15.2
U2Seg	COCO	15.4	51.5	19.0
U2Seg	COCO+IN	17.6	52.7	21.7
Δ		+5.2	+16.6	+6.5

Table 5. **Unsupervised Panoptic image segmentation** on Cityscapes val. We show PQ, SQ and RQ on zero-shot and non-zero shot settings with the cluster number of 800. We compare with CutLER+STEGO, a combination of CutLER+ and STEGO.

Methods	Pretrain	PQ	SQ	RQ
<i>zero-shot methods</i>				
U2Seg	IN	11.1	60.1	13.7
<i>non zero-shot methods</i>				
CutLER+STEGO	COCO	12.4	64.9	15.5
U2Seg	COCO	15.3	66.5	19.1
U2Seg	COCO+IN	16.1	71.1	19.9
Δ		+3.7	+6.2	+4.4

Table 6. **Unsupervised Panoptic image segmentation on COCO val2017**. We show PQ, SQ and RQ on zero-shot and non-zero shot settings. We use CutLER+STEGO, a combination of CutLER+ and STEGO, as a strong baseline.

performance in both unsupervised object detection and instance segmentation on MSCOCO, delivering a **+2.8** boost in AP_{50}^{box} and a **+2.6** rise in AP_{50}^{mask} over CutLER+. Additionally, our method sees a substantial increase of approximately **+10.0** in AR_{100} . Results on PASCAL VOC val2012 and UVO val are detailed in Tab. 3 and Tab. 4, respectively. Notably, we achieve gains exceeding **+20%** in AR for PASCAL VOC and **+10%** for UVO.

4.4. Unsupervised Panoptic Segmentation

For unsupervised panoptic/universal image segmentation, our experiments span two scenarios. In the zero-shot setting, the model is trained exclusively on unlabeled ImageNet images. For non-zero-shot (in-domain) scenarios, we train on unlabeled COCO images or a mix of COCO and ImageNet. With no existing benchmarks for unsupervised panoptic segmentation, we establish a new baseline by integrating the state-of-the-art unsupervised semantic segmentation from STEGO [24] with semantic-aware instance segmentation from CutLER+ (discussed in Sec. 4.3), which are then merged to create panoptic/universal segmentation outcomes, referred to as CutLER+STEGO.

Tab. 6 presents the PQ, SQ, and RQ scores of U2Seg on COCO val2017. U2Seg surpasses the strong baseline CutLER+STEGO with a **+3.5** improvement in PQ and an increase of over **+4.0** in RQ. Qualitative results of U2Seg’s performance is provided in Fig. 4, with the predicted semantic labels visualized. The qualitative results suggest that an

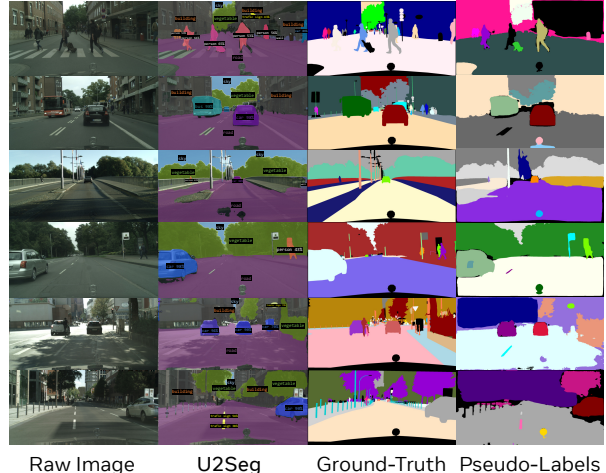


Figure 6. Qualitative results of U2Seg’s **Panoptic image segmentation** results on Cityscapes val (after Hungarian matching).

over-clustering strategy in pseudo-label generation, *e.g.* setting the number of clusters to 300 or 800, leads to highly granular model predictions. For instance, as in Fig. 4, the model distinctly categorizes hockey players as “139”, badminton players as “52”, and gentlemen in suits as “132”, showcasing its potent discriminative capabilities.

To quantitatively measure the quality of segmentation masks and their corresponding semantic labels, we use Hungarian matching (detailed in Appendix B) to align semantic labels with the category names from the test dataset; for instance, all three sub-clusters depicted in Fig. 4 are assigned to the “person” category. The qualitative outcomes post-Hungarian matching are shown in Fig. 5, where our model demonstrates superior panoptic segmentation mask quality. For instance, while the baseline tends to segment parts separately (as seen with the man’s head and torso being treated as separate entities in the third row), our model correctly identifies them as parts of a single object. This level of recognition is also evident with the “trunks of the motorcycle” example in the second row. For additional results, please see Appendix D. We also present results of the more challenging dataset Cityscapes in Tab. 5 and Fig. 6.

4.5. Efficient Learning

Specifically, for object detection and instance segmentation, we employ our unsupervised instance segmentation model, with cluster count set to 300, to initialize the model weights. We adopt the recipe from [56, 59] for model fine-tuning across various annotation splits. For label-efficient panoptic segmentation, we fine-tune the model initialized with our zero-shot unsupervised framework on the same data splits.

The results are depicted in Fig. 7, where our model’s instance segmentation performance is benchmarked against MoCo-V2, DETReg, and CutLER. Our model consistently

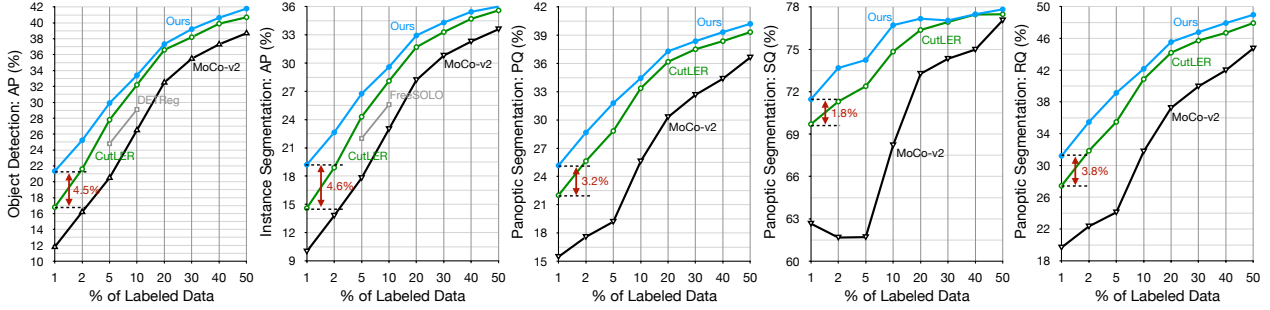


Figure 7. We evaluate the **label-efficient learning** performance on 3 different tasks: object detection (the left), instance segmentation (the second left) and panoptic image segmentation (the last three).

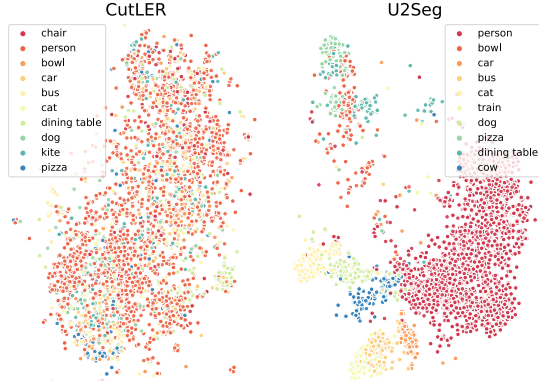


Figure 8. U2Seg learns features that are more discriminative than those learned by CutLER. The **t-SNE** [47] visualization of the features from the model’s FC layer. We color-code each dot based on its ground-truth category.

surpasses the state-of-the-art with consistent gains in both AP^{box} and AP^{mask} . In scenarios with panoptic image segmentation as the downstream task, we contrast our results with MoCo-V2 and CutLER in terms of PQ, SQ, and RQ metrics. The results illustrate a remarkable improvement, effectively doubling the performance boost from MoCo-V2 to CutLER, especially in few-shot contexts with limited annotations (1% or 2% labeled samples). This highlights the practical value of our approach in real-world unsupervised learning applications, where annotations are often scarce.

We attribute the performance gains primarily to the discriminative features our model learns, as in Fig. 8, obtaining effective model initialization for few-shot learning.

4.6. Ablation Studies

In this section, we conduct ablation study on U2Seg.

Numbers of clusters. The choice of cluster quantity significantly affects the model’s representation granularity. Our ablation study on various cluster counts, as detailed in Tab. 7, reveals their impact on model performance. Over-clustering generally leads to a finer level of detail, prompting the model to learn more discriminative features.

Hungarian matching. As our trained model could pre-

# Cluster	COCO		UVO		VOC	
	AP_{50}^{box}	AR_{100}^{box}	AP_{50}^{box}	AR_{100}^{box}	AP_{50}^{box}	AR_{100}^{box}
300	9.3	20.1	9.8	22.6	29.6	45.7
800	11.8	21.5	10.8	25.0	31.0	48.0
2911	13.3	22.1	15.1	25.8	31.6	48.3

Table 7. Over-clustering can improve the model performance. We show results on different datasets for the unsupervised object detection using **different cluster numbers**.

conf	# matched	AP_{50}^{box} AR_{100}^{box}		IoU	# matched	AP_{50}^{box} AR_{100}^{box}	
		AP_{50}^{box}	AR_{100}^{box}			AP_{50}^{box}	AR_{100}^{box}
0.9	109	10.9	13.1	0.9	295	10.8	19.7
0.7	225	11.6	18.0	0.8	348	11.4	20.7
0.6	282	11.8	19.7	0.4	414	11.5	21.6
0.4	389	11.8	21.5	0.2	450	11.5	21.1
0.2	513	11.3	21.8	0.0	494	9.2	17.7
0.0	718	8.6	18.4	0.6	389	11.8	21.5

(a) **Conf’s effect on accuracy.**

(b) **IoU’s effect on accuracy.**

Table 8. **Impact of Confidence and IoU** on Hungarian Matching Performance: The left table illustrates the outcomes at a fixed IoU of 0.6 while varying the confidence scores. Conversely, the right table displays the results with a constant confidence of 0.4, altering the IoU values. The cluster number is 800.

dict the instance with corresponding semantic labels, we are able to go further beyond unsupervised class-agnostic instance segmentation. To quantitatively evaluate the performance, Hungarian matching is employed to match the predicted semantic labels to the ground-truth dataset categories. See Appendix B for details of the adopted Hungarian matching used in our evaluation. As shown in Tab. 8, the two parameters conf threshold and IoU threshold also affect the precision and recall.

5. Summary

We present **U2Seg**, a novel **Unsupervised Universal Image Segmentation** model, adept at performing unsupervised instance, semantic, and panoptic segmentation tasks within a unified framework. Evaluated on extensive benchmarks, U2Seg consistently outperforms previous state-of-the-art methods designed for individual tasks. Additionally, U2Seg achieves the new state-of-the-art for label-efficient panoptic

segmentation and instance segmentation. We anticipate that U2Seg, free from the constraints of human annotations, will demonstrate enhanced performance when scaled up with more training data, representing an important direction for our future research.

6. Acknowledgement

Trevor Darrell, Dantong Niu and Xudong Wang were funded by DoD including DARPA LwLL and the Berkeley AI Research (BAIR) Commons.

Appendix Materials

A. Datasets used for Evaluation

We provide information about the datasets used in this work as shown in Tab. A1

COCO. The COCO dataset, introduced by [39], is used for object detection and instance segmentation. It has 115,000 training images, 5,000 validation images, and a separate batch of 123,000 unannotated images. We test our unsupervised instance segmentation on the COCO val2017 set with zero-shot setting. We report results using standard COCO metrics, including average precision and recall for detection and segmentation. Also, for unsupervised universal image segmentation, we test the performance on COCO val2017. We report results using panoptic segmentation COCO metrics.

PASCAL VOC. The PASCAL VOC dataset [21] is a widely-used benchmark for object detection. We test our model using the trainval07 split and adopt COCO-style evaluation metrics.

UVO. The UVO dataset [54] is designed for video object detection and instance segmentation. We test our unsupervised instance segmentation on the UVO val split, which includes 256 videos with each one annotated at 30 fps. We remove the extra 5 non-COCO categories which are marked as “other” in their official annotations. For evaluation, we employ COCO-style metrics.

Cityscapes. Cityscapes is a dataset dedicated to semantic urban scene understanding, focusing primarily on semantic segmentation of urban scenes. In our research, we tested our unsupervised universal image segmentation on the Cityscapes val splits, using COCO-style panoptic evaluation metrics.

B. Hungarian Matching for Unsupervised Segmentation Evaluation

In unsupervised object detection and instance segmentation, category IDs are predicted without referencing any predefined labels. For convenience, we differentiate the predicted category ID of U2Seg as “cluster ID” while keep the ground

truth category ID as “category ID” in the following analysis. To evaluate the segmentation performance, particularly concerning category accuracy, an optimal correspondence between the cluster ID and the ground truth category ID is essential. We leverage a multi-to-one Hungarian matching for evaluation of U2Seg.

Hungarian Matching. Given a set of predicted bounding boxes, masks associated with predicted cluster IDs and the corresponding ground truth, the objective is to find the best match from “cluster ID” to “category ID”. To do this, we first use the predicted confidence score *conf* as a threshold to filter the predicted instance, removing the ones with low confidence. Then, for each predicted instance with its cluster ID, we calculate the IoU of the predicted bounding box or mask with all ground truth instances, then select the one whose IoU is bigger than the predefined threshold, regarding it as the ground truth category ID for this cluster ID. After we get these cluster ID and ground truth category ID pairs, we form a histogram for each kind of cluster ID based on its overlap with all kinds of ground truth category ID. The ground truth category ID that appears most frequently in this histogram becomes the mapping for this cluster ID. This process may result in multiple predicted cluster IDs being mapped to the same ground truth category ID, leading to a multi-to-one matching scenario.

In our experiment, the confidence score threshold *conf* to filter the predicted instance and the IoU threshold to match predicted instance with its ground truth instance are both hyperparameters, some ablations can be found in Sec. 4.6.

Evaluation Implications. The multi-to-one Hungarian matching method provides a systematic and efficient way to assess the performance of unsupervised segmentation models. By mapping predicted cluster ID to their most likely ground truth counterparts, the method ensures that the evaluation reflects the true categorization capability of the model. This, in turn, allows for a fair and consistent comparison across different unsupervised segmentation techniques.

C. Unsupervised Instance Segmentation

In this section, we provide complete results for the unsupervised instance segmentation of U2Seg. The results are presented over various datasets and classes to furnish a comprehensive evaluation of our model’s capability.

Tab. A2 and Tab. A3 display the results for unsupervised object detection and instance segmentation on different datasets. One trend can be observed across the different datasets: as the number of the predicted cluster ID increases (*e.g.*, moving from 300 to 2911), there is a consistent increase for most of the metrics. This trend can be succinctly attributed to the intrinsic properties of the multi-to-one Hungarian matching approach (we also show the parameter IoU and Conf used for Hungarian matching). With

Datasets	Domain	Testing Data	#Images	Instance Segmentation Label
COCO [39]	natural images	val2017 split	5,000	✓
UVO [54]	video frames	val split	21,235	✓
PASCAL VOC [21]	natural images	trainval07 split	9,963	✗
Cityscapes [15]	urban scenes	val split	500	✓

Table A1. **Summary of datasets** used for evaluation.

Datasets	# cluster	IoU	Conf	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	AP _S ^{box}	AP _M ^{box}	AP _L ^{box}	AR ₁ ^{box}	AR ₁₀ ^{box}	AR ₁₀₀ ^{box}
UVO	2911	0.6	0.1	9.7	15.1	9.3	0.6	5.2	14.4	18.0	25.3	25.8
	800	0.4	0.1	6.8	10.8	7.2	0.6	2.9	10.2	17.2	24.5	25.0
	300	0.7	0.1	6.5	9.8	6.5	0.8	2.6	9.2	16.0	22.2	22.6
VOC	2911	0.5	0.2	19.2	31.6	19.7	1.0	6.4	26.6	28.6	44.9	48.3
	800	0.8	0.2	19.0	31.0	19.5	0.6	4.8	26.6	28.8	45.2	48.1
	300	0.8	0.4	18.4	29.6	18.8	0.3	3.8	26.0	27.1	41.0	42.8
COCO	2911	0.5	0.3	8.2	13.3	8.4	1.4	7.0	18.2	14.1	21.4	22.1
	800	0.6	0.4	7.3	11.8	7.5	1.2	5.8	15.8	13.3	20.8	21.5
	300	0.6	0.3	5.7	9.3	5.9	0.5	4.6	12.9	11.9	19.5	20.1

Table A2. **Complete results for unsupervised object detection.** We show results on UVO val, PASCAL VOC val2012 and COCO val2017, with corresponding clustering numbers. The IoU and Conf are the Hungarian matching parameter we use for evaluation.

Datasets	# cluster	IoU	Conf	AP ^{mask}	AP ₅₀ ^{mask}	AP ₇₅ ^{mask}	AP _S ^{mask}	AP _M ^{mask}	AP _L ^{mask}	AR ₁ ^{mask}	AR ₁₀ ^{mask}	AR ₁₀₀ ^{mask}
UVO	2911	0.6	0.1	8.8	13.9	8.4	0.5	6.4	14.4	16.0	21.7	22.1
	800	0.4	0.1	6.2	9.5	6.0	0.5	2.1	9.8	15.7	20.6	21.0
	300	0.7	0.1	6.1	9.5	5.8	0.7	1.0	8.8	14.1	19.2	19.4
COCO	2911	0.5	0.3	7.3	12.4	7.4	0.8	4.9	17.9	12.8	18.7	19.2
	800	0.6	0.4	6.4	11.2	6.4	0.7	3.7	15.0	11.9	18.0	18.5
	300	0.6	0.3	4.9	8.6	5.0	0.3	2.6	11.8	10.7	16.9	17.3

Table A3. Complete results for **unsupervised instance segmentation.** We show results on UVO val and COCO val2017, with corresponding clustering numbers. The IoU and Conf is the Hungarian matching parameter we use for evaluation.

Datasets	Pretrain	# Cluster	PQ	PQ St	PQ Th	SQ	SQ Th	SQ St	RQ	RQ Th	RQ St
COCO	IN	300	11.1	9.5	19.3	60.1	60.3	59.0	13.7	11.6	25.0
	IN	800	11.9	10.5	19.6	65.9	67.4	58.2	14.8	12.8	25.3
	COCO	300	15.3	14.2	21.6	66.5	67.2	62.4	19.1	17.5	27.5
	COCO	800	15.5	14.6	20.5	69.7	71.1	62.6	19.1	17.8	26.1
	IN+COCO	300	15.5	14.4	21.2	67.1	67.7	64.3	19.2	17.8	26.9
	IN+COCO	800	16.1	15.1	21.2	71.1	72.5	63.8	19.9	18.6	26.8
Cityscapes	IN	300	15.3	4.1	23.4	48.8	54.7	44.6	19.5	5.4	29.7
	IN	800	15.7	4.3	24.0	46.6	47.5	45.9	19.8	5.5	30.2
	COCO	300	18.4	7.8	26.1	47.4	47.3	47.4	22.6	9.8	31.9
	COCO	800	15.4	5.8	22.3	51.5	62.9	43.2	19.0	7.5	27.4
	IN+COCO	300	16.5	6.2	24.1	44.1	45.2	43.3	20.5	7.9	29.7
	IN+COCO	800	17.6	8.4	24.2	52.7	67.5	42.0	21.7	10.5	29.9

Table A4. Complete results for **unsupervised universal image segmentation.** We show results for different models pretrained on various dataset and test on COCO val2017, Cityscapes val, with corresponding cluster numbers.

an increase of the cluster numbers, the Hungarian matching has a broader set of predictions to associate with a single label. This inherently increases the chances of having at least one correct prediction for the given label, making the matching process more amenable. In essence, larger cluster

numbers afford easier matching, thereby boosting the evaluation metrics.

Furthermore, the qualitative results are shown in Fig. A1, with the samples selected in COCO val2017 and PASCAL VOC val2012. After Hungarian matching, we are

Model	AP^{box}	AP_{50}^{box}	AP^{mask}	AP_{50}^{mask}
CutLER+	5.9	9.0	5.3	8.6
Panoptic	6.1	9.8	5.8	9.0
Instance	7.3	11.8	6.4	11.2

Table A5. **Limitation of U2Seg.** We show the zero-shot unsupervised instance segmentation results on COCO val2017. CutLER+ is evaluated on the combination of CutLER and offline clustering, Panoptic is trained on both “*stuff*” and “*things*” pseudo labels, Instance is trained solely on “*things*” labels.

able to get the real categories of the predicted instances.

D. Unsupervised Universal Image Segmentation

Our model’s performance for unsupervised universal image segmentation closely mirrors the trends observed in instance segmentation. Specifically, as the number of the predicted clusters increases, the performance of the panoptic segmentation also improves. Detailed universal segmentation results are shown in Tab. A4 and Fig. A2.

E. Limitation

The primary goal of our research is to develop a comprehensive model capable of excelling in all areas of unsupervised segmentation. As shown in Tab. A5, in terms of the individual sub-task, the universal model exhibits a slight underperformance compared to its counterpart model trained with task-specific annotations. This suggests that U2Seg is adaptable to various tasks, yet it requires task-specific training to achieve the best outcomes for a specific sub-task. Looking ahead, we aim to develop a more versatile model that can be trained once to effectively handle multiple tasks.



Figure A2. **Unsupervised universal image segmentation** visualizations of COCO val2017 (after Hungarian matching).

References

- [1] Edward H Adelson. On seeing stuff: the perception of materials by humans and machines. In *Human Vision and Electronic Imaging*, 2001.
- [2] YM Asano, C Rupprecht, and A Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *International Conference on Learning Representations*, 2019.
- [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *International Conference on Learning Representations*, 2021.
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33, 2020.
- [8] 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.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [10] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.
- [11] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020.
- [12] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. 2021.
- [13] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. 2022.
- [14] Jang Hyun Cho, U. Mall, K. Bala, and Bharath Hariharan. Picie: Unsupervised semantic segmentation using invariance and equivariance in clustering. *ArXiv*, abs/2103.17070, 2021.
- [15] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. *CoRR*, abs/1604.01685, 2016.
- [16] Zhiyuan Dang, Cheng Deng, Xu Yang, Kun Wei, and Heng Huang. Nearest neighbor matching for deep clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13693–13702, 2021.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [18] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015.
- [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [20] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6569–6578, 2019.
- [21] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [22] Edward W Forgy. Cluster analysis of multivariate data: efficiency versus interpretability of classifications. *biometrics*, 21:768–769, 1965.
- [23] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- [24] 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.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017.
- [27] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- [28] 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.
- [29] Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. OneFormer: One Transformer to Rule Universal Image Segmentation. 2023.

- [30] Xu Ji, João F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *ICCV*, 2019.
- [31] Alexander Kirillov, Kaiming He, Ross B. Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9396–9405, 2018.
- [32] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6399–6408, 2019.
- [33] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, 2019.
- [34] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019.
- [35] 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. *arXiv preprint arXiv:2304.02643*, 2023.
- [36] Hans-Peter Kriegel, Erich Schubert, and Arthur Zimek. The (black) art of runtime evaluation: Are we comparing algorithms or implementations? *Knowledge and Information Systems*, 52(2):341–378, 2017.
- [37] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Yukang Chen, Lu Qi, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation with point-based supervision. *arXiv preprint arXiv:2108.07682*, 2021.
- [38] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European Conference on Computer Vision*, pages 280–296. Springer, 2022.
- [39] 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 *European conference on computer vision*, pages 740–755. Springer, 2014.
- [40] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- [41] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.
- [42] Geoffrey J McLachlan and Thriyambakam Krishnan. *The EM algorithm and extensions*. John Wiley & Sons, 2007.
- [43] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020.
- [44] Sungwon Park, Sungwon Han, Sundong Kim, Danu Kim, Sungkyu Park, Seunghoon Hong, and Meeyoung Cha. Improving unsupervised image clustering with robust learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12278–12287, 2021.
- [45] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- [46] 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.
- [47] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [48] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *European conference on computer vision*, pages 268–285. Springer, 2020.
- [49] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Unsupervised semantic segmentation by contrasting object mask proposals. *arXiv preprint arXiv:2102.06191*, 2021.
- [50] Wouter Van Gansbeke, Simon Vandenhende, and Luc Van Gool. Discovering object masks with transformers for unsupervised semantic segmentation. *arXiv preprint arXiv:2206.06363*, 2022.
- [51] Van Huy Vo, Elena Sizikova, Cordelia Schmid, Patrick Pérez, and Jean Ponce. Large-scale unsupervised object discovery. *Advances in Neural Information Processing Systems*, 34:16764–16778, 2021.
- [52] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. MaX-DeepLab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021.
- [53] Weiyao Wang, Matt Feiszli, Heng Wang, and Du Tran. Unidentified video objects: A benchmark for dense, open-world segmentation. *CoRR*, abs/2104.04691, 2021.
- [54] Weiyao Wang, Matt Feiszli, Heng Wang, and Du Tran. Unidentified video objects: A benchmark for dense, open-world segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10776–10785, 2021.
- [55] Xudong Wang, Ziwei Liu, and Stella X Yu. Unsupervised feature learning by cross-level instance-group discrimination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12586–12595, 2021.
- [56] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [57] Xudong Wang, Long Lian, and Stella X Yu. Unsupervised selective labeling for more effective semi-supervised learning. In *European Conference on Computer Vision*, pages 427–445. Springer, 2022.
- [58] 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.
- [59] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. *arXiv preprint arXiv:2301.11320*, 2023.
- [60] 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.
- [61] 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. *arXiv preprint arXiv:2209.00383*, 2022.
- [62] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.
- [63] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, pages 478–487. PMLR, 2016.
- [64] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. Upsnet: A unified panoptic segmentation network. In *CVPR*, 2019.
- [65] Chengxu Zhuang, Alex Lin Zhai, Daniel Yamins, , et al. Local aggregation for unsupervised learning of visual embeddings. In *ICCV*, 2019.
- [66] Adrian Ziegler and Yuki M Asano. Self-supervised learning of object parts for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14502–14511, 2022.