

Leveraging Hidden Positives for Unsupervised Semantic Segmentation

Hyun Seok Seong, WonJun Moon, SuBeen Lee, Jae-Pil Heo*
Sungkyunkwan University

{gustjrd195, wjun0830, leesb7426, jaepilheo}@skku.edu

Abstract

Dramatic demand for manpower to label pixel-level annotations triggered the advent of unsupervised semantic segmentation. Although the recent work employing the vision transformer (ViT) backbone shows exceptional performance, there is still a lack of consideration for task-specific training guidance and local semantic consistency. To tackle these issues, we leverage contrastive learning by excavating hidden positives to learn rich semantic relationships and ensure semantic consistency in local regions. Specifically, we first discover two types of global hidden positives, task-agnostic and task-specific ones for each anchor based on the feature similarities defined by a fixed pre-trained backbone and a segmentation head-in-training, respectively. A gradual increase in the contribution of the latter induces the model to capture task-specific semantic features. In addition, we introduce a gradient propagation strategy to learn semantic consistency between adjacent patches, under the inherent premise that nearby patches are highly likely to possess the same semantics. Specifically, we add the loss propagating to local hidden positives, semantically similar nearby patches, in proportion to the predefined similarity scores. With these training schemes, our proposed method achieves new state-of-the-art (SOTA) results in COCO-stuff, Cityscapes, and Potsdam-3 datasets. Our code is available at: <https://github.com/hynnsk/HP>.

1. Introduction

Semantic segmentation is a major task for scene understanding and plays a crucial role in many applications including medical imaging and autonomous driving [5, 11, 28, 34, 38, 41]. However, existing supervised approaches demand large-scale pixel-level annotations which require huge labeling costs. It has triggered the advent of weakly-supervised [23, 33, 35, 36] and unsupervised semantic segmentation [8, 14, 19, 37] which are to learn without expensive pixel-level annotations.

*Corresponding author

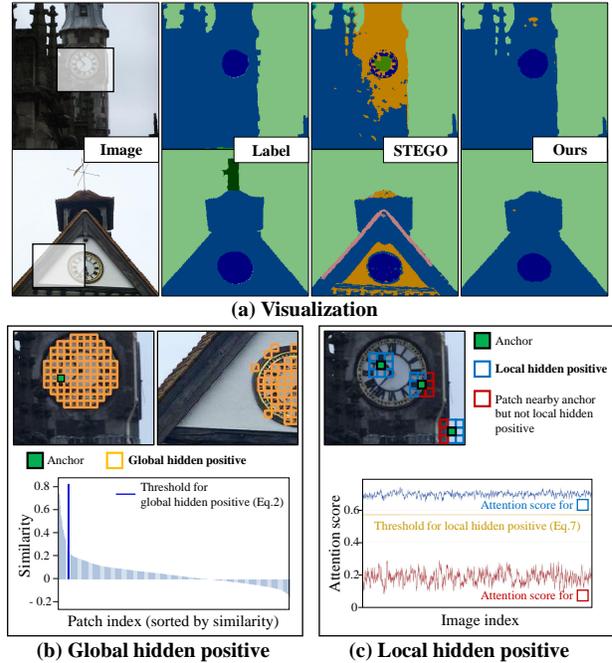


Figure 1. Assuming a mini-batch comprising two images shown in (a), we describe two types of hidden positives, to leverage for contrastive learning. (a) With two types of hidden positives introduced in (b) and (c), we provide an example of how our training scheme provides more precise and consistent semantics. (b) (top) Semantically analogous patches throughout the mini-batch are selected as global hidden positives. (bottom) Data-driven criterion per anchor is designed for reliable positive collection. With the criterion, selected positives are illustrated in (b top). (c) (top) We define local hidden positives for each anchor to be the adjacent patches with high semantic consistency, i.e., blue boxes. (bottom) Average attention scores for adjacent patches from the pretrained transformer architecture. The blue line represents the attention score for local hidden positives while the red line for patches neighboring anchors but having low semantic consistency.

Particularly, unsupervised semantic segmentation is one of the most challenging tasks, since it needs to capture pixel-level semantics from unlabeled data. In this context, clustering-based approaches have been proposed to learn

semantic-preserving clusters by attracting the augmented views in the pixel-level [8, 19]. They implemented the intuition of contrastive learning [3, 6, 7, 13, 15, 25] by ensuring the augmented pairs yield symmetric cluster assignments. More recently, as discovering pixel-level semantics from scratch is challenging, STEGO [14] broke down the problem into learning the representation and learning the segmentation head. With the learned patch-level representation from the seminal work in unsupervised learning [4, 39], they train the segmentation head with a distillation strategy. Although they have made great advancements, we point out their limitations in that they rely solely on a fixed backbone that is not specifically trained for the segmentation task and overlook the importance of semantic consistency along the adjacency that could be a crucial clue for segmentation.

To take these into consideration, we leverage contrastive learning based on the mined hidden positives to ensure contextual consistency along the patches with analogous semantics, particularly the nearby patches, as described in Fig. 1. Specifically, we elaborately select the pseudo-positive samples (i.e., global hidden positive, GHP) for contrastive learning to learn semantic consistency. Also, to ensure local consistency, we propagate the loss gradient to the adjacent patches (i.e., local hidden positive, LHP) in proportion to their equivalency. First, the GHP selection process is designed with two types of data reference pools, task-agnostic and task-specific, to collect the semantically consistent patch features throughout the mini-batch per anchor. For instance, the task-agnostic data reference pool is composed of features extracted by the unsupervised pretrained backbone. On the other hand, the task-specific reference pool is constructed with the features from the segmentation head-in-training to complement task relevance. Based on the two reference pools, two sets of GHP are selected each with generalized and task-specific perspectives. Second, to implement the property of locality and prevent the semantics from fluctuating, we propagate the loss gradient to adjacent patches (i.e., LHP) in proportion to the similarity scores built within the pretrained backbone. This enables the model to learn the relevance of the local context that nearby patches often belong to the same instance.

Our main contributions are summarized as:

- We propose a novel method to discover semantically similar pairs, called global hidden positives, to explicitly learn the semantic relationship among patches for unsupervised semantic segmentation.
- We utilize the task-specific features from a model-in-training and validate the effectiveness of progressive increase of their contribution.
- A gradient propagation to nearby similar patches, local hidden positives, is developed to learn local semantic consistency which is the nature of segmentation.

- Our approach outperforms existing state-of-the-art methods across extensive experiments.

2. Related Work

2.1. Unsupervised Semantic Segmentation

Semantic segmentation has been extensively studied for its wide applicability [5, 11, 29, 34, 38, 41], but collecting pixel-level annotations requires expensive costs. Therefore, many studies [8, 14, 17, 19, 30, 37] attempted to address semantic segmentation without any supervision. Earlier techniques tried to learn semantic correspondence at the pixel level. IIC [19] maximizes the mutual information between the features of two differently augmented images, and PiCIE [8] learns photometric and geometric invariances as an inductive bias. Yet, their training process highly depends on data augmentation, and learning semantic consistency without any prior knowledge is challenging. Therefore, recent methods [14, 37] adopted the ViT model trained in a self-supervised manner, i.e., DINO [4], as a backbone architecture. For instance, TransFGU [37] relocates the high-level semantic features from DINO into low-level pixel-wise features by generating pixel-wise pseudo labels. On the other hand, STEGO [14] utilizes knowledge distillation that learns correspondences between features extracted from DINO. Although STEGO shows a dramatic performance improvement compared to the prior works, it heavily relies on the pretrained backbone and overlooks the property of local consistency that the adjacent pixels are likely to belong to the same category. On the other hand, our training is driven by both the task-agnostic and task-specific pseudo-positive features, and the gradients are conditionally propagated to the neighboring patches, thereby ensuring task-specificity and locality.

2.2. Contrastive Learning

Self-supervised learning methods [2, 3, 6, 7, 12, 13, 15] aim to learn general representations without any annotations. Thanks to their good representation capability, they have been employed to yield remarkable performances in various downstream tasks [8, 16, 20, 26, 40]. Among them, contrastive learning approaches [3, 6, 7, 13, 15, 25] have shown unrivaled performances. In general, they learn representation by attracting a self-augmented set and repulsing other images [3, 6]. Other variants use additional memory [15] or only exploit the positive set for the attraction [7, 13]. This training scheme has also been utilized for unsupervised semantic segmentation [8, 19]. However, the aforementioned methods only considered augmented pairs for the positive which makes them very sensitive to the quality of the augmentation techniques. Our work differs in that our key idea is to collect and exploit reliable pseudo-positives throughout the mini-batch as described in Fig. 1 (b).

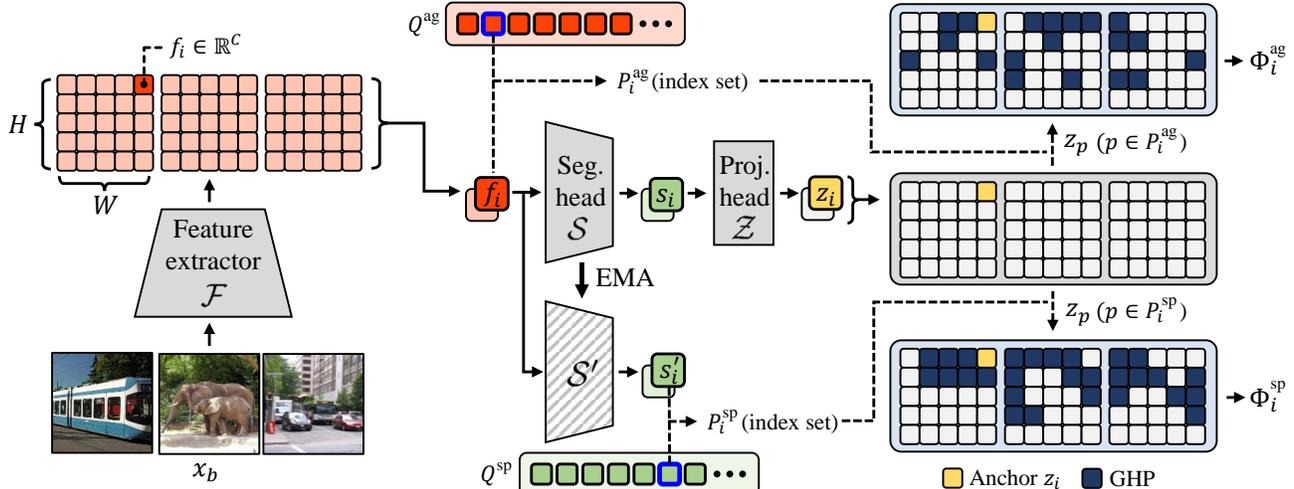


Figure 2. Illustration of the global hidden positive (GHP) selection process. Our GHP can be divided into two sub-sets: task-agnostic and task-specific. An index set of task-agnostic GHP P_i^{ag} comprises the indices of positives discovered within the task-agnostic reference pool Q^{ag} . Note that, Q^{ag} is composed of randomly sampled features extracted by the feature extractor \mathcal{F} . Once the anchor feature f_i is projected to z_i , other patches in the mini-batch are gathered as positives if their similarity with the anchor feature exceeds the similarity between the anchor and the most similar feature in Q^{ag} . On the other hand, task-specific GHP is discovered in a similar manner but with task-specific reference pool Q^{sp} which keeps being updated with the features from the momentum segmentation head S' . Whereas the task-agnostic GHP set solely contributes to the initial training, the task-specific GHP set gradually replaces the portion of the task-agnostic set until the end of training.

3. Method

In this section, we introduce the pseudo-positive selection strategy to discover hidden positives with analogous semantics in Sec. 3.2, the training objective with discovered positives in Sec. 3.3, and the gradient propagation scheme to preserve the property of locality in Sec. 3.4.

3.1. Preliminary

In unsupervised semantic segmentation, the model utilizes unlabeled image set $X = \{x_b\}_{b=1}^B$ where B is the number of training data in the mini-batch. Given an image x_b processed to the feature extractor \mathcal{F} , we have $H \cdot W$ features of $f_i \in \mathbb{R}^C$, where $i \in [1, \dots, H \cdot W]$. Subsequently, the segmentation head \mathcal{S} maps a patch feature f_i to the corresponding segmentation feature $s_i \in \mathbb{R}^K$. And then, the projection head \mathcal{Z} produces a projected vector $z_i \in \mathbb{R}^K$ to formulate a contrastive loss function. In the inference stage, we use the segmentation feature s_i .

Based on the projected vector z_i for the i -th patch, let j be the index of the augmented patch of i -th one. Then, the conventional self-supervised contrastive loss [6] for i -th patch in unsupervised semantic segmentation can be defined as follows:

$$L_i^{\text{self}} = -\log \frac{\exp(\text{sim}(z_i, z_j/\tau))}{\sum_{a \in A} \exp(\text{sim}(z_i, z_a)/\tau)}, \quad (1)$$

where A indicates a set of all indexes except i , τ denotes the

scalar temperature parameter, and $\text{sim}(\cdot, \cdot)$ is cosine similarity between two vectors.

3.2. Global Hidden Positives

Learning mutual information with augmented pixels only provides insufficient training signal in unsupervised semantic segmentation [14]. Therefore, it is important to discover hidden pseudo-positives to tailor the contrastive loss for unsupervised segmentation. To discover the hidden positives at the initial stage, we utilize the self-supervised pretrained backbone [4] as the task-agnostic criterion. Then, we gradually increase the contribution of hidden positives found in a task-specific way for the training. Fig. 2 provides the overview of the global hidden positive (GHP) selection process.

Initially, the pretrained backbone is utilized to construct a task-agnostic reference pool to assess whether other features in the mini-batch are semantically-alike for each anchor feature. Specifically, task-agnostic reference pool, $Q^{\text{ag}} = \{q_m\}_{m=1}^M$, is composed of M randomly sampled features that are extracted by the unsupervised pretrained backbone \mathcal{F} . Note that, we only sample a single patch feature per image to ensure the semantic randomness of the reference pool. This task-agnostic reference pool is fixed as the pretrained backbone is frozen throughout the training.

Once the reference pool is gathered, for each patch feature f_i , we define an anchor-dependent similarity criterion

c_i to collect positives, as the distance to the closest feature within the reference pool Q^{ag} by the cosine similarity:

$$c_i = \max_{q_m \in Q^{\text{ag}}} \text{sim}(q_m, f_i). \quad (2)$$

For each anchor feature f_i , we basically treat the other feature in the mini-batch f_j as positive if the similarity between f_i and f_j is greater than c_i . Still, although one patch feature might be the positive sample for the other, it may not hold mutually. This is because the criterion c_i is anchor-dependent. To endow consistency in training, we make the GHP selection symmetric to prevent the relation between two patches from being ambiguous. Therefore, index set of GHP P_i^{ag} for each i -th anchor feature f_i is defined as follows:

$$P_i^{\text{ag}} = \{j \mid \text{sim}(f_i, f_j) > c_i \vee \text{sim}(f_i, f_j) > c_j\}, \quad (3)$$

where j indicates the index for different patch features in the mini-batch. Accordingly, such a distribution-aware reference pool allows the discovery of globally analogous features in consideration of each anchor.

However, although the reference pool built upon the features from an unsupervised pretrained network can serve as an appropriate basis for positivity, it may be insufficient since it lacks task-specificity. We argue that features from the segmentation head are more task-specific than those from the pretrained backbone. Therefore, along with the GHP selected by P^{ag} , we construct additional task-specific GHP utilizing the features from the segmentation head.

Specifically, an index set of task-specific GHP P_i^{sp} is formed similarly to Eq. 3 by comparing the features $s' = S'(f)$ and task-specific reference pool Q^{sp} , where S' indicates the momentum segmentation head, and the task-specific reference pool Q^{sp} comprises s' . Also, along with the update of the segmentation head, during training, the reference pool is periodically renewed. Formally, with c'_i , calculated by substituting Q^{ag} and f_i in Eq. 2 with Q^{sp} and s'_i , respectively, the P_i^{sp} is expressed as follows:

$$P_i^{\text{sp}} = \{j \mid \text{sim}(s'_i, s'_j) > c'_i \vee \text{sim}(s'_i, s'_j) > c'_j\}. \quad (4)$$

Note that, the usage of the momentum segmentation head is for the stability of the reference pool [22, 42].

3.3. Objective Function

To formulate a contrastive objective with mined GHP in Sec. 3.2, we also need negative features. As we collected the positives throughout the mini-batch, the naive implementation of contrastive learning would utilize all features except the selected positives in the mini-batch as the negatives. However, since an immoderate increase in the size of the negative set may disturb the model training [32], we form a negative set N_i by randomly choosing $\rho\%$ of the remaining patches for each i -th anchor. Note that, the separate

index sets of negative samples N_i^{ag} and N_i^{sp} are defined for each P_i^{ag} and P_i^{sp} , correspondingly. Also, unlike Eq. 1, our contrastive loss for each i -th anchor is more like a supervised objective [20] since we are given multiple positives:

$$L^{\text{cont}}(z_i, P, N) = \frac{-1}{|P|} \sum_{p \in P} \log \frac{\exp(\text{sim}(z_i, z_p)/\tau)}{\sum_{n \in (N \cup P)} \exp(\text{sim}(z_i, z_n)/\tau)}, \quad (5)$$

where z_i , P , and N are the projected anchor vector $\mathcal{Z}(\mathcal{S}(f_i))$, positive index set, and negative index set, respectively. For simplicity, we use Φ_i^{ag} and Φ_i^{sp} to denote the objective functions with task-agnostic GHP P_i^{ag} and task-specific GHP P_i^{sp} for each i -th anchor as follows:

$$\begin{aligned} \Phi_i^{\text{ag}} &= L^{\text{cont}}(z_i, P_i^{\text{ag}}, N_i^{\text{ag}}) \\ \Phi_i^{\text{sp}} &= L^{\text{cont}}(z_i, P_i^{\text{sp}}, N_i^{\text{sp}}). \end{aligned} \quad (6)$$

3.4. Gradient Propagation to Local Hidden Positives

Besides considering the semantically analogous features globally, it is a common hypothesis that nearby pixels are highly likely to belong to the same semantic class. To this end, we consider the property of locality by propagating the loss gradient to the surrounding features of the anchor. Still, the propagation should be cautiously designed since semantic labels of the adjacent patches are not given; semantic consistency between adjacent patches mostly holds, but sometimes does not (i.e., at object boundaries). Thus, to decide the semantically consistent patches nearby, we utilize the attention scores from the unsupervised pretrained ViT backbone \mathcal{F} .

In detail, we first define the index set I_i^{surr} for surrounding patches of the i -th anchor including i -th anchor itself. Also, given the spatial attention score for i -th anchor $\tilde{T}_i \in \mathbb{R}^{H \cdot W}$ from the last self-attention layer in the backbone \mathcal{F} , we use the average value of \tilde{T}_i as a threshold to select an index set of LHP I_i^{local} among I_i^{surr} :

$$I_i^{\text{local}} = \{j \mid j \in I_i^{\text{surr}} \wedge t_j > \text{Avg}(\tilde{T}_i)\}, \quad (7)$$

where t_j is j -th element of \tilde{T}_i and $\text{Avg}(\cdot)$ denotes the averaging function. The visualization of the attention scores for the adjacent patches is shown in Fig. 1 (b). By utilizing I_i^{local} , we obtain the surrounding positive features G_i and the corresponding attention score T_i for LHP as follows:

$$\begin{aligned} G_i &= \{s_j \mid j \in I_i^{\text{local}}\} \\ T_i &= \{t_j \mid j \in I_i^{\text{local}}\}, \end{aligned} \quad (8)$$

where s_j is the feature of j -th patch extracted by the segmentation head \mathcal{S} .

To propagate the loss gradient to the LHP, we mix the patch features in G_i in proportion to the corresponding attention scores in T_i . Formally, the mixed patch composed

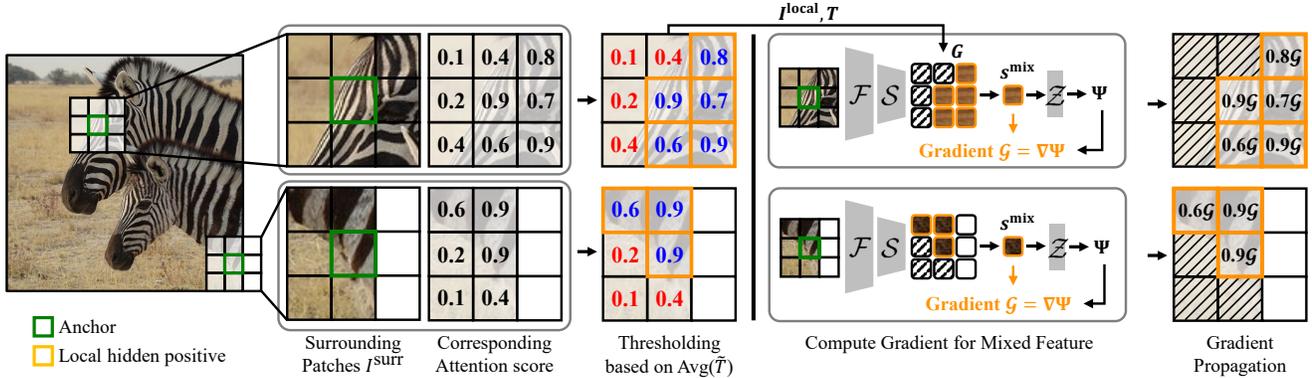


Figure 3. Illustration of our gradient propagation strategy to preserve local semantic consistency. For each anchor, with its surrounding patches I^{surr} and corresponding attention scores from feature extractor \mathcal{F} , local hidden positives (LHP) I^{local} are appointed based on the threshold $\text{Avg}(\bar{T})$ (Eq. 7). In a forward pass, the features of LHP G (Eq. 8) are mixed by weighted average according to the attention scores T to compute the objective function Ψ . In this way, the loss gradient propagates toward the LHP in proportion to T in the backward pass.

of LHP s_i^{mix} is expressed as:

$$s_i^{\text{mix}} = \frac{1}{|I_i^{\text{local}}|} \sum_{j \in I_i^{\text{local}}} \sigma g_j t_j, \quad (9)$$

where σ is a scalar value to scale the attention score, and g_j and t_j indicate the j -th element of G_i and T_i , respectively. Then, we define the objective functions Ψ_i^{ag} and Ψ_i^{sp} to learn the locality by inserting the projected mixed vector $z_i^{\text{mix}} = \mathcal{Z}(s_i^{\text{mix}})$ into Eq. 5 as follows:

$$\begin{aligned} \Psi_i^{\text{ag}} &= L^{\text{cont}}(z_i^{\text{mix}}, P_i^{\text{ag}}, N_i^{\text{ag}}) \\ \Psi_i^{\text{sp}} &= L^{\text{cont}}(z_i^{\text{mix}}, P_i^{\text{sp}}, N_i^{\text{sp}}). \end{aligned} \quad (10)$$

Since these functions are calculated by utilizing the mixed vector z_i^{mix} , the loss gradients are propagated to all features composing the G_i , as described in Fig. 3. Therefore, the semantically-alike surrounding vectors in G_i are updated in the same direction, thereby retaining semantic consistency within the neighboring patches.

Overall, by combining all loss formulations with consistency regularizer R_i that minimizes Euclidean distance between the projected vectors of two differently augmented patches, the final loss function is defined as follows:

$$L_i = (\Phi_i^{\text{ag}} + \Psi_i^{\text{ag}}) + \lambda(\Phi_i^{\text{sp}} + \Psi_i^{\text{sp}}) + \alpha R_i, \quad (11)$$

where λ and α control the contribution of each loss. For instance, λ gradually increases from 0 to 1 during the training and α remains constant at 0.05 throughout the training. Note that, the sample with zero positive is excluded from training although it rarely exists.

4. Experiments

4.1. Datasets and Experimental Settings

Datasets. We utilize COCO-stuff [1], Cityscapes [9], and Potsdam-3 datasets following the existing works [8, 14,

19]. COCO-stuff is a large-scale scene understanding dataset that consists of dense pixel-level annotations and Cityscapes is a more recently publicized dataset having street scenes across 50 different cities. Potsdam-3 dataset contains satellite images. Following the baselines [8, 14, 19], we choose the 27 classes for COCO-stuff and Cityscapes datasets, and 3 classes for Potsdam-3 dataset.

Evaluation Protocols and Metrics. To evaluate our approach, we conduct two testing methods; clustering and linear probe [14]. Clustering is to measure how well the semantic-preserving clusters are formed. Once the unlabeled clusters are computed with the extracted representations, clusters are matched with the ground truth class labels using the Hungarian matching algorithm. On the other hand, the linear probe is a popular method to evaluate the quality of representations learned in unsupervised manners. Specifically, an additional linear layer is learned with representations extracted from the frozen unsupervised model for evaluation and ground truth labels. Inferring the representations of test data extracted by the frozen model with the learned linear layer, we can measure the quality of representations. With two types of protocols, the performance is measured by two common metrics; accuracy (Acc.) and mean Intersection Over Union (mIoU).

Implementation Details. For fair comparisons with our baselines [14, 37], we follow them to mainly use DINO pretrained ViT models as a backbone network \mathcal{F} for the COCO-stuff dataset. In addition, we also test with the advanced backbone, SelfPatch [39], on the COCO-stuff dataset. The segmentation head \mathcal{S} is constructed with a two-layer RELU MLP as STEGO, and the projection head \mathcal{Z} is composed of a linear layer equipped with a normalization layer. The embedding dimension K is set to 512 for ViT-

Method	Backbone	Unsupervised		Linear	
		Acc.	mIoU	Acc.	mIoU
DC [2]	R18+FPN	19.9	-	-	-
MDC [2]	R18+FPN	32.2	9.8	48.6	13.3
IIC [19]	R18+FPN	21.8	6.7	44.5	8.4
PiCIE [8]	R18+FPN	48.1	13.8	54.2	13.9
PiCIE+H [8]	R18+FPN	50.0	14.4	54.8	14.8
DINO [4]	ViT-S/8	28.7	11.3	68.6	33.9
+ TransFGU [37]	ViT-S/8	52.7	17.5	-	-
+ STEGO [14]	ViT-S/8	48.3	24.5	74.4	38.3
+ HP (Ours)	ViT-S/8	57.2	24.6	75.6	42.7
DINO [4]	ViT-S/16	22.0	8.0	50.3	18.1
+ STEGO [14]	ViT-S/16	52.5	23.7	70.6	34.5
+ HP (Ours)	ViT-S/16	54.5	24.3	74.1	39.1
SelfPatch [39]	ViT-S/16	35.1	12.3	64.4	28.5
+ STEGO [14]	ViT-S/16	52.4	22.2	72.2	36.0
+ HP (Ours)	ViT-S/16	56.1	23.2	74.9	41.3

Table 1. Experimental results on COCO-stuff dataset with various backbones and pretrained models.

Method	Backbone	Unsupervised		Linear	
		Acc.	mIoU	Acc.	mIoU
MDC [2]	R18+FPN	40.7	7.1	-	-
IIC [19]	R18+FPN	47.9	6.4	-	-
PiCIE [8]	R18+FPN	65.5	12.3	-	-
DINO [4]	ViT-S/8	34.5	10.9	84.6	22.8
+ TransFGU [37]	ViT-S/8	77.9	16.8	-	-
+ HP (Ours)	ViT-S/8	80.1	18.4	91.2	30.6
DINO [4]	ViT-B/8	43.6	11.8	84.2	23.0
+ STEGO [14]	ViT-B/8	73.2	21.0	90.3	26.8
+ HP (Ours)	ViT-B/8	79.5	18.4	90.9	33.0

Table 2. Experimental results on Cityscapes dataset.

S/8 and ViT-B/8 models, and 256 for ViT-S/16. We train the model for 3, 20, and 10 epochs for COCO-stuff, Cityscapes, and Potsdam-3 datasets, respectively, based on the AdamW optimizer with a learning rate of 0.0005 and weight decay of 0.1. The task-specific reference pool Q^{sp} is renewed every 100 iterations throughout the training. The percentage of negative samples usage ρ is set to 2. In the last stage, we add a feature refinement step utilizing Conditional Random Field [21] as did in STEGO. Evaluation metrics, i.e., clustering and linear probe, are optimized with the Adam optimizer each with learning rates of 0.005 and 0.001.

4.2. Experimental Results

We compare our proposed method against the prior techniques for the unsupervised segmentation [2, 8, 14, 19, 37]. Most of the results in the result tables are brought from the literature [8, 14]. In Tab. 1, it is observed that self-supervised models, i.e., DINO and SelfPatch, are already good segmentation predictors with the linear probe, which makes them a new baseline over the prior works for unsu-

Method	Backbone	Unsup. Acc.
Random CNN [19]	VGG11	38.2
K-Means [27]	VGG11	45.7
SIFT [24]	VGG11	38.2
ContextPrediction [10]	VGG11	49.6
CC [18]	VGG11	63.9
DeepCluster [2]	VGG11	41.7
IIC [19]	VGG11	65.1
DINO [4]	ViT-B/8	53.0
+ STEGO [14]	ViT-B/8	77.0
+ HP (Ours)	ViT-B/8	82.4

Table 3. Experimental results on Potsdam-3 dataset.

pervised segmentation. Furthermore, we utilize two pre-trained backbones with two kinds of architectures to compare with STEGO, in detail. As reported, our proposed model provides consistent performance improvements over the previous SOTA model in almost all cases on the COCO-stuff dataset.

Results on Cityscapes also show a similar tendency. As shown in Tab. 2, ours outperform previous methods except for the mIoU when clustering is used for the evaluation. For instance, we achieve 8.6% and 23% improvements in cluster accuracy and linear mIoU over STEGO, respectively. For the slight decrease in cluster mIoU with ViT-B/8 architecture, we argue that it is insignificant since the linear probe better describes the quality of representations. Specifically, clustering evaluation highly depends on the purpose of the dataset so that it is sensitive to the degree of class-specificity as different body parts can be either classified as human or as independent body parts. However, whereas it is more appropriate to detect each body part independently for unsupervised learning as ours do in Fig. 5, the annotations for general datasets, e.g., COCO-stuff and cityscape, treat these body parts as a human class. Such circumstances make the clustering evaluation vulnerable to the degree of the class hierarchy. In contrast, the linear probe projects these features to close proximity if the given label space considers them as a human body. Thus, we believe the linear probe is a more appropriate measure of representation quality.

Also, we compare the cluster accuracy in Tab. 3 on Potsdam-3 dataset. We have achieved a 7% boost over STEGO, confirming that ours also performs well even in a completely new domain. Likewise, our superior results verify the effectiveness of our process of discovering global- and local-hidden positive patches.

Qualitative Results. In addition to the quantitative results, we report qualitative results in Fig. 4. In comparison to STEGO, we observe that our results include fewer mis-predicted pixels throughout the images, while the strength



Figure 4. Qualitative comparison results of Ours and STEGO on the COCO-stuff dataset with DINO pretrained ViT-S/8 backbone.

of preserving semantic locality is also validated. For instance, whereas we find that the predicted label of the wheel is only partially correct in the 4th column, our results are consistent along the neighboring pixels. These results also demonstrate the superiority of our method.

5. Ablation Study and Further Analysis

In this section, we provide ablation studies and an analysis of our model. Particularly, we explore the contributions of main components and test with varying hyperparameters. Most experiments for ablation studies are conducted on the COCO-stuff dataset using the DINO pretrained ViT-S/8 model, except for Sec. 5.2, which utilized the ViT-S/16.

5.1. Importance of the Main Components.

Tab. 4 reports the performances when an individual component or various combinations of them are not utilized. We found that task-specific GHP and LHP are essential in improving the performance of the unsupervised segmentation task. Compared to (d) where both the task-specific GHP and LHP are not used, the use of them each leads to 6.5% (b) and 11.6% (c) improvements. Also when used together, they bring 16% of performance boosts (a). Furthermore, the importance of preserving symmetricity in selecting GHP can be found by comparing (a) to (e) as its usage boosts 5.9%, and consistency regularizer enhances performance by 2.5% by comparing (a) to (f). Lastly, (g) shows the performance of naively implemented contrastive learning (Eq. 1) with photometric perturbations used in PicIE [8]. This also confirms the strengths of our proposed method as (a) enhances (g) by 51.3%.

5.2. Alternatives to Gradient Propagation Strategy

There can be alternative ways to meet our goal of reflecting semantic consistency between adjacent patches. As an alternative method of gradient propagation, we simply apply the identical loss to the surrounding patches pro-

	GHP		LHP	SA	Reg	Unsupervised	
	TA	TS				Acc.	mIoU
(a)	✓	✓	✓	✓	✓	57.2	24.6
(b)	✓	✓		✓	✓	52.5	23.1
(c)	✓		✓	✓	✓	55.0	19.1
(d)	✓			✓	✓	49.3	20.1
(e)	✓	✓	✓		✓	54.0	23.6
(f)	✓	✓	✓	✓		55.8	24.5
(g)						37.8	10.4

Table 4. Ablation study for each component. GHP, LHP, TA, TS, SA, and Reg denote Global Hidden Positive, Local Hidden Positive, task-agnostic, task-specific, symmetrical assignments, and consistency regularizer, respectively.

Method	Unsupervised		Linear	
	Acc.	mIoU	Acc.	mIoU
DINO + GradProp	54.5	24.3	74.1	39.1
DINO + LossProp	54.7	23.2	74.3	40.5
SelfPatch + GradProp	56.1	23.2	74.9	41.3
SelfPatch + LossProp	54.5	22.2	75.1	41.4

Table 5. Comparison results between gradient propagation and loss propagation strategies.

portionally to their attention score (i.e., loss propagation). The results in Tab. 5 show that the loss propagation strategy performs comparably to the gradient propagation strategy. Nonetheless, we observed that the loss propagation approach incurs higher computational costs ($1.2\times$ memory and $3\times$ time). Likewise, our method is implemented considering both effectiveness and efficiency.

5.3. Visualization

Discovered GHP. To demonstrate the effectiveness of our GHP selection process, we visualize the selected GHP sets for different anchors in a single image in Fig. 5. First, we observe that the corresponding reference point in the second column is semantically correlated with the anchor patch.



Figure 5. Discovered patches by our GHP selection process. From the left to right columns, red boxes indicate the anchors, the closest patches within the task-agnostic reference pool (reference), and GHP sets chosen in the mini-batch, respectively.

Also, the obtained GHP sets verify the appropriateness of the use of reference points as each anchor’s criterion as they are capable of precisely discovering the semantically-similar positives. For example, we find that all the anchors, reference points, and GHP sets have the same semantic labels in the first and the second rows (solid (mountain) and ground (snow)). More intriguing results are in the third and fourth rows where we find that the GHP selection process distinguishes the body parts in a much more fine-grained manner than the given annotation, i.e., person. These results imply that the designed GHP selection is well-designed and capable of capturing detailed semantic contexts.

5.4. Robustness to Hyperparameters

In this subsection, we enumerate our use of hyperparameters in Tab. 6 and conduct an ablation study. Overall, our proposed training scheme is robust to hyperparameters as Fig. 6 supports the claim. Below, we illustrate the influences of each parameter. Despite the slight drop in performance from time to time, it is insignificant since the change in performance is very marginal and results are consistent.

The Number of Data in the Reference Pool. Fig. 6 (a) shows the performances with varying numbers of data M in the reference pool. Generally, the reference pool is not very vulnerable to its size, unless the capacity is either too small or too large. When the reference pool is skinny, it may not be sufficient to represent all kinds of semantics present in the dataset. In other words, the small reference pool may induce a biased criterion for each anchor which may incur biased training. On the other hand, when the reference pool is too large, a tight threshold (c_i) could be derived, which can interfere with gathering GHP.

Dataset	M	τ
COCO-stuff	2048	0.8
Cityscapes	2048	0.6
Potsdam-3	1024	0.4

Table 6. Hyperparameter used for each dataset.

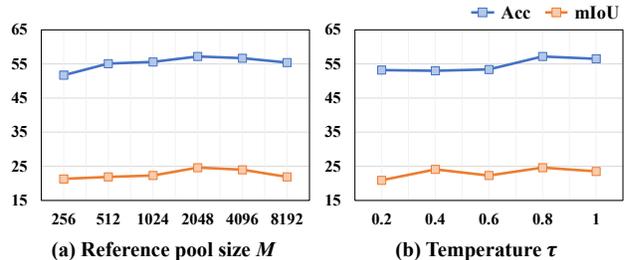


Figure 6. Ablation studies on hyperparameters.

Temperature Parameter. The temperature parameter τ is the scaling parameter to manipulate the sensitivity of the contrastive loss. If τ gets bigger, the objective function becomes robust to the difference between the similarity of the positive and negative samples. On the other hand, when τ gets lower, the embedding distribution is likely to be more uniform [31]. Although the results do not fluctuate much, we observe that training uniformly distributed embedding space leads to a slight performance drop in Fig. 6 (b).

6. Conclusion

In this paper, we introduced a novel unsupervised semantic segmentation method by discovering and leveraging two types of hidden positives, global hidden positive (GHP) and local hidden positive (LHP), to learn rich semantic information with local consistency. First, anchor-dependent GHP comprises task-agnostic and task-specific positive sets which are used to tailor the contrastive learning for the unsupervised semantic segmentation task. Whereas the task-agnostic features are collected to guide the initial training, task-specific features are progressively engaged to learn the task-specific semantics information. Moreover, under the inherent premise that the adjacent patches are likely to be semantically similar, we propagate the loss gradient to the surrounding patches in proportion to their attention scores. This encourages the semantically similar peripheral patches to have the same objective as the anchor, resulting in semantic consistency between adjacent patches unless they belong to different objects. Finally, our proposed method achieves new state-of-the-art results in various datasets.

Acknowledgements. This work was supported in part by MSIT/IITP (No. 2022-0-00680, 2019-0-00421, 2020-0-01821, 2021-0-02068), and MSIT&KNPA/KIPoT (Police Lab 2.0, No. 210121M06).

References

- [1] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocosuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 5
- [2] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 2, 6
- [3] 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*, 33:9912–9924, 2020. 2
- [4] 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. 2, 3, 6
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *ICLR*, 2015. 1, 2
- [6] 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. 2, 3
- [7] 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. 2
- [8] Jang Hyun Cho, Utkarsh Mall, Kavita Bala, and Bharath Hariharan. Picie: Unsupervised semantic segmentation using invariance and equivariance in clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16794–16804, 2021. 1, 2, 5, 6, 7
- [9] 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. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 5
- [10] 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. 6
- [11] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3146–3154, 2019. 1, 2
- [12] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018. 2
- [13] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 2
- [14] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T. Freeman. Unsupervised semantic segmentation by distilling feature correspondences. In *International Conference on Learning Representations*, 2022. 1, 2, 3, 5, 6
- [15] 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. 2
- [16] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32, 2019. 2
- [17] Jyh-Jing Hwang, Stella X Yu, Jianbo Shi, Maxwell D Collins, Tien-Ju Yang, Xiao Zhang, and Liang-Chieh Chen. Segsort: Segmentation by discriminative sorting of segments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7334–7344, 2019. 2
- [18] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H Adelson. Learning visual groups from co-occurrences in space and time. *arXiv preprint arXiv:1511.06811*, 2015. 6
- [19] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9865–9874, 2019. 1, 2, 5, 6
- [20] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. 2, 4
- [21] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24, 2011. 6
- [22] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *ICLR*, 2017. 4
- [23] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5267–5276, 2019. 1
- [24] David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1150–1157. Ieee, 1999. 6
- [25] 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. 2

- [26] WonJun Moon, Ji-Hwan Kim, and Jae-Pil Heo. Tailoring self-supervision for supervised learning. In *European Conference on Computer Vision*, pages 346–364. Springer, 2022. [2](#)
- [27] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. [6](#)
- [28] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [1](#)
- [29] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021. [2](#)
- [30] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Unsupervised semantic segmentation by contrasting object mask proposals. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10052–10062, 2021. [2](#)
- [31] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2495–2504, 2021. [8](#)
- [32] Shujun Wang, Lequan Yu, Caizi Li, Chi-Wing Fu, and Pheng-Ann Heng. Learning from extrinsic and intrinsic supervisions for domain generalization. In *European Conference on Computer Vision*, pages 159–176. Springer, 2020. [4](#)
- [33] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020. [1](#)
- [34] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. [1](#), [2](#)
- [35] Lian Xu, Wanli Ouyang, Mohammed Bennamoun, Farid Boussaid, and Dan Xu. Multi-class token transformer for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4310–4319, 2022. [1](#)
- [36] Yazhou Yao, Tao Chen, Guo-Sen Xie, Chuanyi Zhang, Fumin Shen, Qi Wu, Zhenmin Tang, and Jian Zhang. Non-salient region object mining for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2623–2632, 2021. [1](#)
- [37] Zhaoyuan Yin, Pichao Wang, Fan Wang, Xianzhe Xu, Hanling Zhang, Hao Li, and Rong Jin. Transfgu: a top-down approach to fine-grained unsupervised semantic segmentation. In *European Conference on Computer Vision*, pages 73–89. Springer, 2022. [1](#), [2](#), [5](#), [6](#)
- [38] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *European conference on computer vision*, pages 173–190. Springer, 2020. [1](#), [2](#)
- [39] Sukmin Yun, Hankook Lee, Jaehyung Kim, and Jinwoo Shin. Patch-level representation learning for self-supervised vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8354–8363, 2022. [2](#), [5](#), [6](#)
- [40] Xiangyun Zhao, Raviteja Vemulapalli, Philip Andrew Mansfield, Boqing Gong, Bradley Green, Lior Shapira, and Ying Wu. Contrastive learning for label efficient semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10623–10633, 2021. [2](#)
- [41] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. [1](#), [2](#)
- [42] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Sessd: Self-ensembling single-stage object detector from point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14494–14503, 2021. [4](#)