

Rethinking Query-based Transformer for Continual Image Segmentation

Yuchen Zhu^{1*} Cheng Shi^{1*} Dingyou Wang¹ Jiajin Tang¹ Zhengxuan Wei¹
Yu Wu³ Guanbin Li² Sibei Yang^{2†}

¹School of Information Science and Technology, Shanghai Tech University

²Sun Yat-sen University

³Wuhan University

Abstract

Class-incremental/Continual image segmentation (CIS) aims to train an image segmenter in stages, where the set of available categories differs at each stage. To leverage the built-in objectness of query-based transformers, which mitigates catastrophic forgetting of mask proposals, current methods often decouple mask generation from the continual learning process. This study, however, identifies two key issues with decoupled frameworks: loss of plasticity and heavy reliance on input data order. To address these, we conduct an in-depth investigation of the built-in objectness and find that highly aggregated image features provide a shortcut for queries to generate masks through simple feature alignment. Based on this, we propose SimCIS, a simple yet powerful baseline for CIS. Its core idea is to directly select image features for query assignment, ensuring "perfect alignment" to preserve objectness, while simultaneously allowing queries to select new classes to promote plasticity. To further combat catastrophic forgetting of categories, we introduce cross-stage consistency in selection and an innovative "visual query"-based replay mechanism. Experiments demonstrate that SimCIS consistently outperforms state-of-the-art methods across various segmentation tasks, settings, splits, and input data orders. All models and codes will be made publicly available at https://github.com/SooLab/SimCIS.

1. Introduction

Continual learning empowers models to progressively acquire, learn, and assimilate new knowledge from an everevolving environment. It serves as a fundamental task in image classification [4, 8, 18, 20, 25, 31, 40, 42, 48, 49, 57, 58, 60, 64, 71, 72, 75] where models are required to recognize new classes (**plasticity**) and preserve old class knowledge (avoid **catastrophic forgetting**). Extending be-

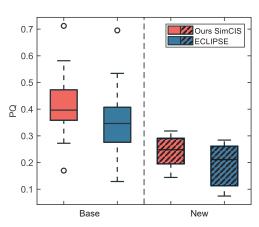


Figure 1. **Boxplots** of PQ metric for our SimCIS and previous SOTA [38] on ADE20K. We train each model on randomly shuffled continual data input orders and report average PQ for base and novel classes. We observe that recent query-based transformers suffer from a loss of plasticity (low average PQ) and heavy reliance on the input data order (high variance).

yond classification, continual image segmentation adapts this to the image segmentation, unlocking a myriad of practical applications [50, 53]. However, it also confronts more challenges: 1) Additional catastrophic forgetting of mask prediction, beyond that of class prediction; 2) Background semantic shift occurs when the current foreground becomes background in subsequent stages, driven by the need for image segmentation to predict the background class and the constraint of only having class annotations from current stage. Recently, query-based transformers [10, 17, 34, 35, 55, 56, 59, 63, 79] are introduced into continual image segmentation, as their built-in objectness has been shown to mitigate catastrophic forgetting in mask generation. Leveraging this built-in objectness, many studies [7, 26, 38, 74] decouple mask segmentation from the continual learning process by freezing the parameters associated with mask proposal generation. However, we observe two notable yet suboptimal behaviors in the aforementioned methods.

^{*}Equal contribution

[†]Corresponding author

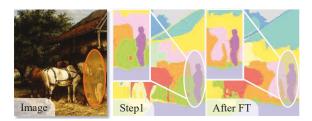


Figure 2. **Clustering results** from feature map. Pixel feature provides sufficient semantic priors (Person) even after finetuning.

- The advantage of objectness diminishes and even has a detrimental effect on plasticity as the task sequence shortens. In the shortest two-task setting, they typically achieve performance comparable to or even slightly lower than the baseline.
- The built-in objectness is fragile and lacks robustness, showing heavy dependence on the split and order of input data. As shown in Fig 1, in ten random trials, the worst trial shows a significant performance drop on new classes compared to the default setting.

Therefore, in this work, we aim to understand the built-in objectness and achieve consistent improvements (especially on plasticity) across different task lengths and varying data input orders. This is crucial, as it is impractical to assume fixed task lengths and data sequences in real-world scenarios. The conclusion from a series of investigations is:

- * The built-in objectness diminishes over training stages due to the query's failure to align with the semantic priors of the feature map. As shown in Fig 3 (left), since semantic priors vary at different stages due to background semantic shift, causing the updated learnable query to gradually misalign with the pixel feature from old classes in previous stages, even after the decoder's post-alignment (observed in 1).

Inspired by ① and ②, to ensure objectness is preserved throughout the continual learning stages, we propose a lazy Query Pre-Alignment (QPA) method, where query features are selected from specific locations in the image feature map, rather than being learned from scratch, to "perfectly" pre-align query feature with semantic priors. Specif-

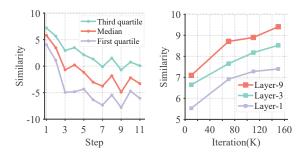


Figure 3. **Similarity** between queries and feature map changes across decoder layers and training stages (right). The query gradually misaligns with the pixel feature (left).

ically, based on the current stage's semantic classes, we select the most semantically significant locations in the image feature, preserving objectness at each stage. However, objectness is still lost across stages due to varying semantic classes in different stages.

To overcome cross-stage selection issues, a naive solution involves distillation on the feature map or query features between stages. However, in turn, while it preserves old priors from previous stages, it re-introduces incorrect priors for current stages (where old priors label current semantics as background), leading to a loss of plasticity. Fortunately, thanks to our query pre-alignment method, we can easily maintain old classes by keeping queries corresponding to old class positions, while enabling the selection of remaining queries for new classes in the current stage. Thus, we propose a **Consistent Selection Loss (CSL)** to ensure that, for the same image, the most semantically significant locations selected in the previous stage are revisited in the current stage.

With QPA and CSL, objectness in the query-based transformer is fully utilized to generate mask proposals. However, for class prediction, catastrophic forgetting may still occur. Previous methods typically rely on image replay to mitigate catastrophic forgetting. In contrast, thanks to our query pre-alignment, our query inherently contains category semantics. By storing the query feature, we can simulate specific semantics without requiring the actual image to contain the corresponding category. Therefore, we propose a novel **Virtual Query (VQ)** strategy to replay the virtual queries corresponding to previous classes in the decoder layer to avoid catastrophic forgetting. Compared to conventional image replay methods, our approach reduces storage requirements by 10x, is independent of input data order, and preserves dataset privacy.

In summary, our contributions are multi-fold:

- We provide a thorough analysis of the built-in objectness, revealing the reasons behind its emergence and demise.
- By addressing the root cause, we can successfully leverage built-in objectness to mitigate catastrophic forgetting

- and background semantic shift through the introduction of three simple yet novel modules—QPA, CSL, and VQ.
- Our model, SimCIS, consistently and significantly outperforms state-of-the-art results on ADE20K in both continual panoptic and semantic segmentation.
- We introduce new dataset splits to evaluate the model's robustness to input order in continual learning. SimCIS shows superior robustness over state-of-the-art methods, thanks to the effective utilization of built-in objectness.

2. Related Work

Continual Learning is a longstanding field which possesses significant importance in addressing dynamic environments, enhancing model adaptability, and improving resource efficiency. The objective of continuous learning is to enable the model to efficiently acquire and adapt to new tasks and data, while retaining previously learned knowledge as it encounters additional information. The greatest challenge of continual learning is catastrophic forgetting [24, 49, 65]. The early research are categorized into three primary types: those that rely on regularization constraints [8, 9, 19, 20, 39, 41], those employing replay techniques [46, 49, 60], and those based on dynamic structures [22, 42, 43, 61, 71, 75]. Regularization-based methods aim to reduce the interference of new tasks on old knowledge by constraining the learning process of the model, ensuring that the model parameters remain closely aligned with previously learned representations when updated due to task changes. Replay-based methods employ strategies to store, replay [4, 33, 48, 67], or generate [46, 60, 66] samples from old tasks to mitigate catastrophic forgetting. Those methods based on dynamic structure [42, 43, 52] allocate distinct subsets of parameters to various subtasks by facilitating the expansion of their network architecture.

Universal Image Segmentation. Before MaskFormer proposed, traditional segmentation methods developed specialized architectures and models for each task to achieve top performance [2, 12–15, 28, 30, 36, 62, 69, 73, 78]. MaskFormer [16] is the first unified segmentation architecture to achieve state-of-the-art performance across three image segmentation tasks. Mask2Former [17] improves MaskFormer by adapting multi-scale features and introducing mask attention mechanism and achieve better performance. Follow its success in segmentation, we use Mask2Former as our baseline aims to extend its capability into the field of continual learning.

Continual Segmentation is the application of continual learning within the field of image segmentation. The challenge of continual segmentation tasks lies in the ability to identify new categories while generating high-quality masks for each category. This dual requirement underscores the complexity of maintaining accurate segmentation performance while adapting to an evolving set of class labels.

Methods for continual segmentation are also categorized into three types as previously mentioned: regularizationbased [5, 6, 21, 44, 45, 47, 54, 68, 74, 77], replay-based [7, 11, 23, 76, 81], and dynamic structure-based [1, 26, 27, 38, 70]. Among these methods, those query-based architectures demonstrate notable performance. CoMFormer [6] is the first query-based method in the field of continuous panoptic segmentation, employing distillation and pseudo label to combat catastrophic forgetting. CoMasTRe [26] is inspired by the methods of CoMFormer and, while maintaining the use of distillation loss, decouples mask and class predictions in continuous segmentation tasks. ECLIPSE [38] adapts the strategy of VPT [37], freezing the majority of model parameters and providing a set of trainable queries for fine-tuning across different tasks. BalConpas [11] attempts to combat catastrophic forgetting by employing a method that combines feature-based distillation and a replay sample set, aiming to learn new classes without negatively impacting previously acquired knowledge.

3. Preliminary

3.1. Problem Setting

Following the same continual learning setting in [6], we train our model over T steps. At each step t, the model \mathcal{M}^t has access only to a subset $\mathcal{D}^t = \{x^t, y^t\}$ of the entire dataset $\mathcal{D}^{1:T}$, where $x^t \in \mathbb{R}^{C \times H \times W}$ denotes the image at the current step and y^t represents the corresponding annotations (where it can only contain annotations for classes \mathcal{C}^t). This setup, where each stage involves learning different classes, makes the model highly susceptible to catastrophic forgetting as it tends to lose previously acquired knowledge at each training step. Meanwhile, as the same image may appear across different learning steps with entirely different annotations, we also face the issue of so-called background shift [5]. Given these challenges, our objective is to design a model $\mathcal M$ such that, at any stage t, the model $\mathcal M^t$ not only effectively learns from $\mathcal D^t$ but also preserve the previous class knowledge from $\mathcal D^{1:t-1}$.

3.2. Mask2Former

We leverage Mask2former [17] as our meta-architecture for image segmentation. Mask2Former is a transformer-based model, which predicts a set of binary masks instead of perpixel classification, for universal segmentation tasks. It primarily consists of three components: 1) An image encoder as backbone f_{backbone} to extract image embeddings. 2) A pixel decoder f_{pixel} to embed image embeddings to multiscale pixel features, which we denote as F:

$$F = \{ \mathcal{F}_{(l,h,w)} \mid \forall (l,h,w) \in \Omega \}, \ \mathcal{F} \in \mathbb{R}^{D \times H_l \times W_l}, \quad (1)$$

where l denotes the multi-scale layer, D represents the hidden dimension, $\mathcal{F}_{(l,h,w)}$ refers to the feature point at po-

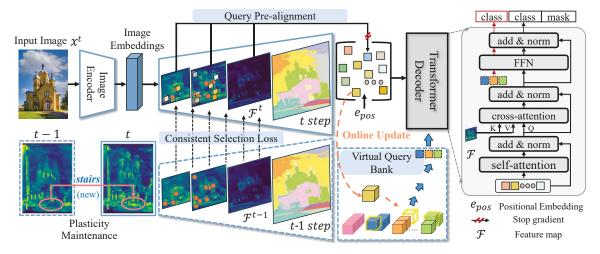


Figure 4. The Overall Architecture of our SimCIS: a lazy Query Pre-Alignment (Sec 4.1) with a Consistent Selection loss (Sec 4.2) to ensure built-in objectness inner and across stages, and Virtual Query (Sec 4.3) to avoid catastrophic forgetting in class prediction.

sition (h,w) on the l-th layer and Ω represents the spatial set of multi-scale features. 3) A transformer decoder f_{decoder} takes N learnable queries $Q_N = \{q_1,q_2,\ldots,q_N\} \in \mathbb{R}^{N \times D}$ with positional encodings $e_{pos} \in \mathbb{R}^{N \times D}$ to first conduct cross-attention and then self-attention with \mathcal{F} as follows:

$$Q'_{N} = FFN(SA(CA(Q_{N} + e_{nos}, F))), \tag{2}$$

where CA(,) denotes the cross-attention, SA(·) represents self-attention, and Q_N' denotes the updated query feature. The final prediction for each query is $Z_N = \{(c_i, m_i)\}_{i=1}^N$, where $c_i \in \mathbb{R}^C$ and $m_i \in \mathbb{R}^{H \times W}$ represent the predicted class and mask for q_i , respectively.

4. Method

In this section, we introduce the overall architecture of our proposed SimCIS model for continual image segmentation. As shown in Fig 4, SimCIS contains three modules: 1) Lazy Query Pre-alignment (Sec 4.1), 2) Consistent Selection Loss (Sec 4.2) and Virtual Query (Sec 4.3).

4.1. Lazy Query Pre-alignment

To preserve the objectness across continual learning stages, we propose to pre-align the object query Q_N with semantic priors in the pixel feature $\mathcal{F}_{(l,h,w)}$ by directly initializing query feature with the most semantically significant pixel feature. To determine the semantic score of each pixel feature, we learn a prototype for each category and select pixel features as initial features by calculating the similarity between the pixel feature and each prototype.

Specifically, for each training step t, we maintain a set of trainable prototypes $\{p^i | i \in C^t\}, p^i \in \mathbb{R}^D$ for each

class in C^t . By concatenating the prototypes of the past step, \mathcal{P}^{t-1} , with those of the current classes, we obtain the current prototype set \mathcal{P}^t as follows,

$$\mathcal{P}^t = \operatorname{concat}(\mathcal{P}^{t-1}, \{ p^i \mid i \in C^t \}). \tag{3}$$

Then for each feature point on F, we compute its similarity with \mathcal{P}^t to select the best feature points. The selection process is as follows:

$$\mathcal{I}^{t} = \operatorname{topK}\left(\left\{\max S(\mathcal{F}_{(l,h,w)}^{t}, \mathcal{P}^{t}) \mid \forall (l,h,w) \in \Omega\right\}, N\right),$$
(4)

$$Q_N = \mathcal{E}_{m=t}^{n=t} = \left\{ \mathcal{F}_i^{m=t} \mid i \in \mathcal{I}^{n=t} \right\},\tag{5}$$

where $\mathcal{I}=\{(l_i,h_i,w_i)\}_{i=0}^N\in\Omega$ represents the spatial positions of the selected feature points, \mathcal{E}_m^n represents the feature points from \mathcal{F}^m selected by \mathcal{I}^n and S(,) denotes the similarity calculation by dot product. The topK(X,Y) function returns the indices of the Y largest values in X and X is the number of object query Q_N . We select X feature points with the highest similarity with the prototype to initialize Q_N . To supervise our selection process, we use a classification loss during training and update \mathcal{P}^t through backpropagation [51]. Additionally, we apply stop gradient on Q_N to ensure that the information in Y is not disrupted during training, keeping the objectness information stable across different stages.

4.2. Consistent Selection Loss

To ensure selection $\mathcal I$ is stable for the same image across stages, we propose a consistent selection loss. Specifically, when training our model $\mathcal M^t$ at current stage, we can easily obtain feature points $\mathcal E_t^{t-1} = \{\mathcal F_i^t \mid i \in \mathcal I^{t-1}\}$. Then,

to maintain consistency in object selection across different steps, we calculate the similarity between selected feature points with \mathcal{P}^{t-1} , after that, we use the Kullback-Leibler (KL) divergence loss [32] to compute the loss:

$$L_{csl} = \frac{1}{|\mathcal{I}^{t-1}|} \sum_{i=1}^{|\mathcal{I}^{t-1}|} S(\mathcal{E}_{t-1}^{t-1}, \mathcal{P}^{t-1}) \log \frac{S(\mathcal{E}_{t-1}^{t-1}, \mathcal{P}^{t-1})}{S(\mathcal{E}_{t}^{t-1}, \mathcal{P}^{t-1})}.$$
(6

In this way, we successfully maintain the most semantically significant locations from the previous stage, ensuring that the selection of Q_N remains stable across stages.

4.3. Virtual Query

To overcome catastrophic forgetting in class prediction, we propose the virtual query to bypass the limitations of previous methods that rely on data order. Virtual Query replays the previous query feature in the decoder layer to simulate semantics. Specifically, our innovative virtual query strategy can be divided into three steps: Firstly, we use the results of bipartite matching to select object queries and build our VQ bank. Then we analyze the pseudo-distribution to focus on rare categories in the current stage. Finally, we sample VQs in the new stage according to the pseudo-distribution and concatenate them into the object query Q_N for input into the decoder.

(1) Query Storage. During training, we maintain a queue of length h for each class, forming our virtual query bank

$$\mathcal{B}_{\text{vq}} = \{b_1^h, b_2^h, \dots, b_{|c^{1:T}|}^h\},\tag{7}$$

where b_i^h represents a queue of length h for class i where b_i^h is the queue for class i. Queries matched through bipartite matching [3] from the decoder's final layer output, Z_N (defined in Sec 3.2), are stored in the appropriate class queues based on their bipartite matching results with ground truth u.

$$\begin{cases}
\mathcal{I}_{b} = \operatorname{Bipartite}(Z_{N}, \boldsymbol{y}), \\
\mathcal{B}_{\operatorname{vq}} \leftarrow \underset{\forall i = (i_{q}, i_{y}) \in \mathcal{I}_{b}}{\operatorname{Enqueue}}(Q_{N}(i_{q}), b_{\hat{y}^{(i_{y})}}),
\end{cases} (8)$$

where N denotes the number of queries. The set \mathcal{I}_b consists of tuples, where each tuple $i=(i_q,i_y)$ represents the correspondence between query and ground truth. Here, i_q denotes the query index, and i_y denotes the ground truth index. \hat{y}^i represents the class label of the i^{th} ground truth.

(2) Pseudo-Distribution Statistics. In each continual learning step, the category distribution of images changes at each stage. To ensure the decoder retains the category information for all old classes, we use the pre-trained last-stage model \mathcal{M}^{t-1} 's outputs on current stage's dataset D^t to simulate the distribution of real classes which helps mitigate the forgetting of rare classes in the current stage. We use this pseudo-distribution statistics by calculating

$$\omega = \left\{ \left(\left(\sum_{i=1}^{m} \sigma_i \right) / \sigma_j \right)^{\frac{1}{2}} \right\}_{i=1}^{m}, \tag{9}$$

where σ_i is the pseudo number of class i in the current stage and $m=|c^{1:t-1}|$ represents the number of categories from the previous stages.

(3) **VQ Utilization.** Based on the pseudo-distribution statistics, in each iteration, we sample j virtual queries $Q_j = \{vq_1, \ldots, vq_j\}$ for each batch based on ω . These queries are then concatenated with Q_N as

$$Q_{N+j} = \{q_1, \dots, q_N, vq_1, \dots, vq_j\},$$
(10)

and fed into the decoder. As shown in Fig 4, within the decoder, we design a skip attention strategy for the VQs. Specifically, since the objects represented by the VQs do not appear in the image, to prevent the VQs from influencing Q_N during the self-attention and cross-attention processes, we allow the VQs to bypass the attention layers and directly affect the FFN layers as follows:

$$Q'_{N+j} = FFN(\operatorname{concat}[\operatorname{CA}(\operatorname{SA}(Q_N + e_{pos}, F)), Q_j]). \tag{11}$$

Finally, the virtual query only computes L_{class} to address the model's category forgetting.

5. Experiments

5.1. Experimental Setup

Dataset and Evaluation Metric. Following previous works [6, 11, 38], we compare our SimCIS with other approaches using the ADE20K dataset [80] to evaluate its effectiveness. The images in the dataset include annotations for 150 classes, which are ranked by their total pixel ratios in the whole dataset. Among these 150 classes, 50 amorphous background classes are labeled as "stuff" classes, while 100 discrete object classes are labeled as "thing" classes. Following[6], we use Panoptic Quality (PQ) as the performance metric for continual panoptic segmentation and mean Inter-over-Union (mIoU) for continual semantic segmentation. After incremental learning steps, we report results for base classes (\mathcal{C}^1), new classes ($\mathcal{C}^{2:T}$), all classes ($\mathcal{C}^{1:T}$), and an average of all visible classes at each step (avg), respectively.

Continual Learning Protocol. Following existing continual segmentation methods [5–7, 11, 21, 26, 38], we evaluate our method on different continual learning settings. In particular, our incremental learning tasks are represented in the form of A-B, where A denotes the number of base classes partitioned from the dataset, and B denotes the number of new classes. For both continual panoptic (CPS) and semantic segmentation (CSS), we conduct tasks of 100 - 5, 100 -

Method		100-5 (11 tasks)			100-10 (6 tasks)			100-50 (2 tasks)				
iviethou	1-100	101-150	all	avg	1-100	101-150	all	avg	1-100	101-150	all	avg
FT	0.0	2.2	0.7	4.7	0.0	4.8	1.6	8.9	0.0	32.4	10.8	26.8
MiB [5]	2.3	0.0	1.5	13.4	6.8	0.2	4.6	19.1	23.3	14.9	20.5	31.7
PLOP [21]	31.1	11.9	24.7	31.3	37.7	23.3	32.9	37.8	42.4	23.7	36.2	39.5
SSUL [7]	30.2	7.9	22.8	27.9	31.6	11.9	25.0	30.3	35.9	18.1	30.0	33.8
CoMFormer [6]	34.4	15.9	28.2	34.0	36.0	17.1	29.7	35.3	41.1	27.7	36.7	38.8
BalConpas [11]	36.1	20.3	30.8	35.8	40.7	22.8	<u>34.7</u>	38.8	42.8	<u>25.7</u>	<u>37.1</u>	40.0
ECLIPSE [38]	41.1	16.6	32.9	-	<u>41.4</u>	18.8	33.9	-	41.7	23.5	35.6	-
Our SimCIS	42.1	21.9	35.4	38.7	42.2	30.1	38.1	40.5	44.7	30.8	40.0	42.7
joint	43.6	34.2	40.4	-	43.6	34.2	40.4	-	43.6	34.2	40.4	-

Table 1. **Continual Panoptic Segmentation** results on ADE20K dataset in PQ. All methods use the same network of Mask2Former [17] with ResNet-50 [29] backbone. *joint* means an oracle setting training all classes offline at once.

Method	50-	10 (11 task	s)	50-20 (6 tasks)			50-50 (3 tasks)		
Method	1-50	51-150	all	1-50	51-150	all	1-50	51-150	all
FT	0.0	1.7	1.1	0.0	4.4	2.9	0.0	12.0	8.1
MiB [5]	34.9	7.7	16.8	38.8	10.9	20.2	42.4	15.5	24.4
PLOP [21]	39.9	15.0	23.3	43.9	16.2	25.4	45.8	18.7	27.7
CoMFormer [6]	38.5	15.6	23.2	42.7	17.2	25.7	45.0	19.3	27.9
ECLIPSE [38]	45.9	17.3	26.8	46.4	19.6	28.6	46.0	20.7	29.2
BalConpas [11]	44.6	24.8	31.4	49.2	28.2	35.2	51.2	26.5	34.7
Our SimCIS	48.8	30.0	36.3	51.6	31.9	38.5	52.1	30.7	37.9
joint	51.1	35.1	40.4	51.1	35.1	40.4	51.1	35.1	40.4

Table 2. Continual Panoptic Segmentation results on ADE20K dataset in PQ. All methods use Mask2Former [17] with ResNet-50 [29].

10, and 100 - 50. Additionally, we conduct tasks of 50 - 10, 50 - 20, and 50 - 50 for panoptic segmentation.

Implementation Details. We adapt an pre-trained ResNet-50 [29] backbone for CPS and an pre-trained ResNet-101 for CSS. Following previous work [11], the input image resolution for the CPS tasks is set to 640×640 , while for the CSS tasks, it is set to 512×512 . For the number of virtual queries N, it be set up to 80. For more detailes, please refer to the Appendix.

5.2. Quantitative Results

Tab 1, Tab 2 and Tab 3 present the performance of Sim-CIS and other approaches on the continual panoptic segmentation and semantic segmentation benchmark. In these tables, "FT" refers to fine-tuning the base model without employing continual learning methods, while "joint" indicates training the base model using all available data. They represent the lower and upper-performance bounds for continual learning methods, respectively.

Continual Panoptic Segmentation. Tab 1 and Tab 2 present the performance of SimCIS and other approaches under different continual panoptic segmentation settings. (1) Compared to regularization-based methods MiB [5], PLOP [21], and CoMFormer [6], SimCIS achieves superior

results on both new and base classes. Notably, compared to CoMFormer, the best-performing among them, SimCIS improves PQ by +6.0% on new classes and +7.7% on base classes in the 100 - 5 task, maintaining a consistent lead in the 100 - 10 and 100 - 50 tasks. Especially in the 100 - 10 task, it surpasses CoMFormer by +6.2% PQ on base and +13.0% PQ on new classes. When using 50 base classes, SimCIS significantly outperforms these methods, demonstrating its superiority. (2) Compared with the method also using built-in objectness, SimCIS achieves better performance on new classes without freezing the model parameters. In the 100 - 5, 100 - 10, and 100 - 50 tasks, SimCIS outperforms ECLIPSE [38] by +5.3% PQ, +11.3% PQ, and +7.6% PQ, respectively. In the tasks with 50 classes as base classes, SimCIS outperforms ECLIPSE [38] by over +10% PQ on new classes, demonstrating the stability of our approach. (3) BalConpas [11] is a continual learning method based on the Mask2Former [17] architecture. In the 100 - 10 and 100 - 50 tasks, SimCIS outperforms BalConpas [11] by more than +5.0% PQ on new classes. In the longer step sequence of the 100 - 5 task, SimCIS surpasses BalConpas [11] by +6.0% PQ on base classes. In the 50 -20 and 50 - 50 tasks, SimCIS maintains strong performance, averaging +4% PQ higher than BalConpas [11] on new

Model	100-5 (11 tasks)			100-10 (6 tasks)			100-50 (2 tasks)					
Model	1-100	101-150	all	avg	1-100	101-150	all	avg	1-100	101-150	all	avg
FT	0.0	0.3	0.1	5.6	0.0	0.1	0.0	9.1	0.0	3.2	1.1	26.3
MiB [5]	36.0	5.7	26.0	-	31.8	14.1	25.9	-	37.9	27.9	34.6	-
PLOP [21]	39.1	7.8	28.8	35.3	40.5	14.1	31.6	36.6	41.9	14.9	32.9	37.4
SSUL [7]	42.9	17.8	34.6	-	42.9	17.7	34.5	-	42.8	17.5	34.4	-
EWF [68]	41.4	13.4	32.1	-	41.5	16.3	33.2	-	41.2	21.3	34.6	-
CoMFormer [6]	39.5	13.6	30.9	36.5	40.6	15.6	32.3	37.4	39.5	26.2	38.4	41.2
ECLIPSE [38]	43.3	16.3	34.2	-	43.4	17.4	34.6	-	45.0	21.7	37.1	-
BalConpas [11]	42.1	17.2	33.8	41.3	47.3	24.2	38.6	43.6	49.9	30.1	43.3	47.4
CoMasTRe [26]	40.8	15.8	32.6	38.6	42.3	18.4	34.4	38.4	45.7	26.0	39.2	41.6
Our SimCIS	46.7	22.8	38.7	47.4	49.7	27.4	42.3	49.2	54.9	36.0	48.6	52.0
Joint	57.1	39.1	51.2	-	57.1	39.1	51.2	-	57.1	39.1	51.2	-

Table 3. Continual Semantic Segmentation results on the ADE20K dataset, measured by mIoU.

Psd	QPA	CSL	VQ	Panopt 1-100	ic 100-5 (1 101-150	1 tasks) all	Seman 1-100	tic 100-5 (1 101-150	1 tasks) all
✓ ✓ ✓ ✓	\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	✓ ✓	√ ✓	31.6 30.7 35.7 35.1 42.1	21.3 22.3 24.0 23.3 21.9	28.2 27.9 31.8 31.2 35.4	15.6 37.4 43.2 42.5 46.7	8.5 16.7 17.0 19.5 22.8	13.2 30.5 34.5 34.8 38.7

Table 4. **Ablation Study on Proposed Components.** Psd: pseudo label, QPA: lazy query pre-alignment, CSL: consistent selection loss, and VQ: virtual query.

classes. In the longer step sequence of the 50 - 10 task, Sim-CIS exceeds BalConpas [11] by +4.2% PQ on base classes. It is noteworthy that in the 100 - 50 task, SimCIS almost matches the performance of the "joint", with base classes performance even exceeding that of the "joint".

Continual Semantic Segmentation. As shown in Tab 3, we further compare SimCIS with state-of-the-art works in continual semantic segmentation. (1) Across three tasks, SimCIS surpasses prior approaches by at least +4% mIoU on base classes. For new classes, it outperforms SSUL [7] by +5.0% and +9.7% mIoU in the 100 - 5 and 100 - 10 tasks, respectively. In the 100 - 50 task, SimCIS surpasses MiB [5], which achieves 27.9% mIoU, by +8.1% mIoU. (2) Among Mask2Former [17]-based methods, SimCIS also achieves the best results. In the 100 - 5 task, it outperforms ECLIPSE [38] on base classes by +3.4% mIoU and Bal-Conpas [11] on new classes by +5.6% mIoU. In the 100 - 10 task, SimCIS achieves the performance of new classes exceeding all other architectures by at least +3.0% mIoU while maintaining high performance on base classes.

5.3. Qualitative Comparison.

Comparison with Previous SOTAs. We compare SimCIS with BalConpas [11] in the 100 - 5 continual panoptic segmentation task of the ADE20K dataset, and the visual results are illustrated in Fig 5. In the first, second, and fifth examples, BalConpas [11] encounters forgetting on base

classes such as path, bus, and building. Additionally, in the third example, BalConpas incorrectly classifies the microwave and bag as cabinet and box, respectively. Benefiting from the VQ, our SimCIS has a significant advantage in preserving class information, allowing it to perform well in these examples. Furthermore, BalConpas [11] fails to provide segmentation masks for the bus and refrigerator instances in the second and third examples. In contrast, our proposed the keep built-in objectness strategy effectively preserves object information within the encoder, enabling SimCIS to accurately segment object instances.

Comparison in Different Steps. To further illustrate the effectiveness of our method, we select certain visual examples from the continual learning steps of the 100 - 5 task. In the two examples shown in Fig 6, our method is able to correct errors during the continual learning steps, such as the microwave and bag in the first image, as well as the sink, vase, and stair in the second image. SimCIS refines itself during the continual learning process, ultimately achieving accurate classification and segmentation of object instances based on our proposed flexible VQ.

5.4. Ablation Study

In this section, we report the results of the ablation experiments to validate the effectiveness of each component and configuration in our SimCIS. We select the 100 - 5 task in CPS and CSS to report the performance of SimCIS.

Main Components. As shown in Tab 4, each component contributes to the overall performance. We take Mask2Former [17] with pseudo label as our baseline performance. The second row of the table shows the performance of QPA with an increase of +18.2% mIoU on base classes and an increase of +8.2% mIoU on new classes. With the help of CSL (the third row), the CSL strategy achieves increases of +8.2% PQ and +5.8% mIoU for base classes, respectively.

Effectiveness of VQ. As shown in Tab 5, compared to the

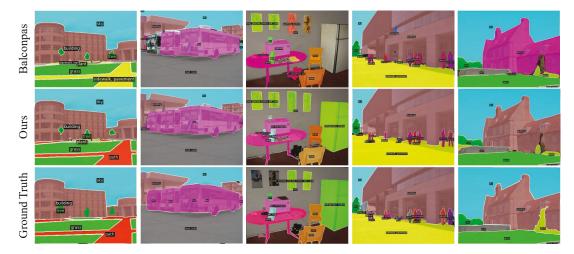


Figure 5. **Qualitative comparisons** between SimCIS and BalConpas [11] on the ADE20K 100-5 continual panoptic segmentation scenario. Our SimCIS demonstrates significant results, highlighting the effectiveness of our strategies.

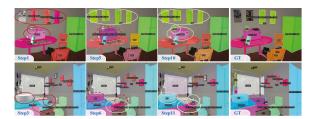


Figure 6. Qualitative examples in continual learning.

Reply	Num	Disk	100-5	(11 tasks)
Type	Samples	Memory	base	all
Image	0 (*20)	0.0MB	35.7	31.8
	75 (*20)	3.4MB	38.9	33.4
	150 (*20)	6.1MB	38.9	34.0
	300 (*20)	11.8MB	38.5	33.7
	600 (*20)	21.9MB	39.2	34.3
Virtual Query	0 (*150)	0.0MB	35.7	31.8
	20 (*150)	1.5MB	40.6	34.6
	40 (*150)	3.0MB	40.4	34.1
	80 (*150)	5.9MB	42.1	35.4
	160 (*150)	12.0MB	40.9	34.2

 $\label{thm:conditional} \mbox{Table 5. Effect of Replay Type and Storage Requirements}.$

conventional image replay method, our VQ strategy demonstrates significant improvements in both storage efficiency and performance. Firstly, when using 300 samples for the image replay and 80 samples for VQ, we achieve a +1.4% increase in PQ across all classes while using almost the same disk memory. When comparing the optimal cases for both storage methods, our VQ strategy outperforms the conventional image replay method by +1.1% PQ, while utiliz-

Method	100		
Method	1-100	101-150	all
BalConpas [11]	38.9(39.4)	27.8(26.8)	35.2
ECLIPSE [38]	32.7(32.1)	22.3(23.8)	29.3
Ours	40.3(40.2)	25.4(25.7)	35.3
Joint	(43.6)	(34.2)	(40.4)

Table 6. Continual Panoptic Segmentation with random order. We also report the performance evaluated in the original class order in (\cdot) . For detailed experiments, please refer to the Appendix.

ing only 27% of the storage space.

Robust to Input Data Order. As shown in Tab 6, our model has great robustness in random data order. We have a +0.1% PQ increase compared to BalConpas and a +6.0% PQ increase against ECLIPSE across all classes.

6. Conclusion

In this work, we present a novel class-incremental image segmentation (CIS) method called SimCIS, which addresses the challenges of catastrophic forgetting and background shift. We first explore the emergence and diminishing of built-in objectness in query-based transformers and then propose two novel modules: lazy query pre-alignment and consistent selection loss, to ensure both intra-stage and cross-stage built-in objectness. Additionally, we introduce virtual queries to mitigate catastrophic forgetting in class prediction. Comparisons with previous state-of-the-art CIS methods and our ablation study demonstrate the superiority of each individual component in our model, highlighting its effectiveness in overcoming the challenges of incremental learning. Acknowledgment: This work was supported by the National Natural Science Foundation of China (No.62206174).

References

- Donghyeon Baek, Youngmin Oh, Sanghoon Lee, Junghyup Lee, and Bumsub Ham. Decomposed knowledge distillation for class-incremental semantic segmentation. Advances in Neural Information Processing Systems, 35:10380–10392, 2022. 3
- [2] 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. 3
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In *European confer*ence on computer vision, pages 213–229. Springer, 2020. 5
- [4] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European conference* on computer vision (ECCV), pages 233–248, 2018. 1, 3
- [5] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Bulo, Elisa Ricci, and Barbara Caputo. Modeling the background for incremental learning in semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9233–9242, 2020. 3, 5, 6, 7
- [6] Fabio Cermelli, Matthieu Cord, and Arthur Douillard. Comformer: Continual learning in semantic and panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3010–3020, 2023. 3, 5, 6, 7
- [7] Sungmin Cha, YoungJoon Yoo, Taesup Moon, et al. Ssul: Semantic segmentation with unknown label for exemplarbased class-incremental learning. Advances in neural information processing systems, 34:10919–10930, 2021. 1, 3, 5, 6, 7
- [8] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European conference on computer vision* (ECCV), pages 532–547, 2018. 1, 3
- [9] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with agem. arXiv preprint arXiv:1812.00420, 2018. 3
- [10] Chaoqi Chen, Yushuang Wu, Qiyuan Dai, Hong-Yu Zhou, Mutian Xu, Sibei Yang, Xiaoguang Han, and Yizhou Yu. A survey on graph neural networks and graph transformers in computer vision: A task-oriented perspective. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2024.
- [11] Jinpeng Chen, Runmin Cong, Yuxuan Luo, Horace Ho Shing Ip, and Sam Kwong. Strike a balance in continual panoptic segmentation, 2024. 3, 5, 6, 7, 8
- [12] Liang-Chieh Chen. Semantic image segmentation with deep convolutional nets and fully connected crfs. arXiv preprint arXiv:1412.7062, 2014. 3
- [13] Liang-Chieh Chen. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint* arXiv:1706.05587, 2017.

- [14] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern* analysis and machine intelligence, 40(4):834–848, 2017.
- [15] Bowen Cheng, Liang-Chieh Chen, Yunchao Wei, Yukun Zhu, Zilong Huang, Jinjun Xiong, Thomas S Huang, Wen-Mei Hwu, and Honghui Shi. Spgnet: Semantic prediction guidance for scene parsing. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5218–5228, 2019. 3
- [16] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Perpixel classification is not all you need for semantic segmentation. Advances in Neural Information Processing Systems, 34:17864–17875, 2021. 3
- [17] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022. 1, 3, 6, 7
- [18] Qiyuan Dai and Sibei Yang. Curriculum point prompting for weakly-supervised referring image segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13711–13722, 2024. 1
- [19] Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyan Wu, and Rama Chellappa. Learning without memorizing. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 5138–5146, 2019. 3
- [20] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In Computer vision–ECCV 2020: 16th European conference, Glasgow, UK, August 23–28, 2020, proceedings, part XX 16, pages 86–102. Springer, 2020. 1, 3
- [21] Arthur Douillard, Yifu Chen, Arnaud Dapogny, and Matthieu Cord. Plop: Learning without forgetting for continual semantic segmentation. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pages 4040–4050, 2021. 3, 5, 6, 7
- [22] Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9285–9295, 2022. 3
- [23] Mostafa ElAraby, Ali Harakeh, and Liam Paull. Bacs: Background aware continual semantic segmentation. arXiv preprint arXiv:2404.13148, 2024. 3
- [24] Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.
- [25] Weifeng Ge, Sibei Yang, and Yizhou Yu. Multi-evidence filtering and fusion for multi-label classification, object detection and semantic segmentation based on weakly supervised learning. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 1277–1286, 2018. 1

- [26] Yizheng Gong, Siyue Yu, Xiaoyang Wang, and Jimin Xiao. Continual segmentation with disentangled objectness learning and class recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 3848–3857, 2024. 1, 3, 5, 7
- [27] Dipam Goswami, René Schuster, Joost van de Weijer, and Didier Stricker. Attribution-aware weight transfer: A warmstart initialization for class-incremental semantic segmentation. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 3195–3204, 2023. 3
- [28] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VII 13, pages 297–312. Springer, 2014. 3
- [29] 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. 6
- [30] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 3
- [31] Xiang He, Sibei Yang, Guanbin Li, Haofeng Li, Huiyou Chang, and Yizhou Yu. Non-local context encoder: Robust biomedical image segmentation against adversarial attacks. In *Proceedings of the AAAI Conference on Artificial Intelli*gence, pages 8417–8424, 2019. 1
- [32] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. stat, 1050:9, 2015.
- [33] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *Proceedings of the IEEE/CVF conference* on computer vision and pattern recognition, pages 831–839, 2019. 3
- [34] Hanzhuo Huang, Yufan Feng, Cheng Shi, Lan Xu, Jingyi Yu, and Sibei Yang. Free-bloom: Zero-shot text-to-video generator with llm director and ldm animator. Advances in Neural Information Processing Systems, 36:26135–26158, 2023. 1
- [35] Hanzhuo Huang, Yuan Liu, Ge Zheng, Jiepeng Wang, Zhiyang Dou, and Sibei Yang. Mvtokenflow: High-quality 4d content generation using multiview token flow. arXiv preprint arXiv:2502.11697, 2025. 1
- [36] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 603–612, 2019. 3
- [37] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer, 2022. 3
- [38] Beomyoung Kim, Joonsang Yu, and Sung Ju Hwang. Eclipse: Efficient continual learning in panoptic segmentation with visual prompt tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 3346–3356, 2024. 1, 3, 5, 6, 7, 8

- [39] Zhizhong Li and Derek Hoiem. Learning without forgetting. IEEE transactions on pattern analysis and machine intelligence, 40(12):2935–2947, 2017. 3
- [40] Liang Lin, Pengxiang Yan, Xiaoqian Xu, Sibei Yang, Kun Zeng, and Guanbin Li. Structured attention network for referring image segmentation. *IEEE Transactions on Multime*dia, 24:1922–1932, 2021.
- [41] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. Advances in neural information processing systems, 30, 2017. 3
- [42] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pages 7765–7773, 2018. 1, 3
- [43] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggy-back: Adapting a single network to multiple tasks by learning to mask weights. In *Proceedings of the European conference on computer vision (ECCV)*, pages 67–82, 2018. 3
- [44] Umberto Michieli and Pietro Zanuttigh. Incremental learning techniques for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 0–0, 2019. 3
- [45] Umberto Michieli and Pietro Zanuttigh. Continual semantic segmentation via repulsion-attraction of sparse and disentangled latent representations. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 1114–1124, 2021. 3
- [46] Oleksiy Ostapenko, Mihai Puscas, Tassilo Klein, Patrick Jahnichen, and Moin Nabi. Learning to remember: A synaptic plasticity driven framework for continual learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 11321–11329, 2019. 3
- [47] Minh Hieu Phan, Son Lam Phung, Long Tran-Thanh, Abdesselam Bouzerdoum, et al. Class similarity weighted knowledge distillation for continual semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 16866–16875, 2022. 3
- [48] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE con*ference on Computer Vision and Pattern Recognition, pages 2001–2010, 2017. 1, 3
- [49] Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 7(2):123–146, 1995.
- [50] David A Ross, Jongwoo Lim, Ruei-Sung Lin, and Ming-Hsuan Yang. Incremental learning for robust visual tracking. *International journal of computer vision*, 77:125–141, 2008.
- [51] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986. 4
- [52] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint arXiv:1606.04671, 2016. 3

- [53] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint arXiv:1606.04671, 2016.
- [54] Chao Shang, Hongliang Li, Fanman Meng, Qingbo Wu, Heqian Qiu, and Lanxiao Wang. Incrementer: Transformer for class-incremental semantic segmentation with knowledge distillation focusing on old class. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7214–7224, 2023. 3
- [55] Cheng Shi and Sibei Yang. Edadet: Open-vocabulary object detection using early dense alignment. In *Proceedings of* the *IEEE/CVF international conference on computer vision*, pages 15724–15734, 2023. 1
- [56] Cheng Shi and Sibei Yang. Logoprompt: Synthetic text images can be good visual prompts for vision-language models. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 2932–2941, 2023. 1
- [57] Cheng Shi and Sibei Yang. The devil is in the object boundary: Towards annotation-free instance segmentation using foundation models. arXiv preprint arXiv:2404.11957, 2024.
- [58] Cheng Shi, Yulin Zhang, Bin Yang, Jiajin Tang, Yuexin Ma, and Sibei Yang. Part2object: Hierarchical unsupervised 3d instance segmentation. In European Conference on Computer Vision, pages 1–18. Springer, 2024. 1
- [59] Cheng Shi, Yuchen Zhu, and Sibei Yang. Plain-det: A plain multi-dataset object detector. In European Conference on Computer Vision, pages 210–226. Springer, 2024. 1
- [60] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. Advances in neural information processing systems, 30, 2017. 1, 3
- [61] Pravendra Singh, Vinay Kumar Verma, Pratik Mazumder, Lawrence Carin, and Piyush Rai. Calibrating cnns for lifelong learning. Advances in Neural Information Processing Systems, 33:15579–15590, 2020. 3
- [62] 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. 3
- [63] Jiajin Tang, Ge Zheng, Cheng Shi, and Sibei Yang. Contrastive grouping with transformer for referring image segmentation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 23570–23580, 2023. 1
- [64] Jiajin Tang, Ge Zheng, and Sibei Yang. Temporal collection and distribution for referring video object segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 15466–15476, 2023. 1
- [65] Sebastian Thrun. Lifelong learning algorithms. In *Learning to learn*, pages 181–209. Springer, 1998. 3
- [66] Chenshen Wu, Luis Herranz, Xialei Liu, Joost Van De Weijer, Bogdan Raducanu, et al. Memory replay gans: Learning to generate new categories without forgetting. Advances in neural information processing systems, 31, 2018. 3
- [67] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF conference*

- on computer vision and pattern recognition, pages 374–382, 2019. 3
- [68] Jia-Wen Xiao, Chang-Bin Zhang, Jiekang Feng, Xialei Liu, Joost van de Weijer, and Ming-Ming Cheng. Endpoints weight fusion for class incremental semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 7204–7213, 2023. 3, 7
- [69] 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. 3
- [70] Zhengyuan Xie, Haiquan Lu, Jia-wen Xiao, Enguang Wang, Le Zhang, and Xialei Liu. Early preparation pays off: New classifier pre-tuning for class incremental semantic segmentation. arXiv preprint arXiv:2407.14142, 2024. 3
- [71] Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3014–3023, 2021. 1, 3
- [72] Sibei Yang, Meng Xia, Guanbin Li, Hong-Yu Zhou, and Yizhou Yu. Bottom-up shift and reasoning for referring image segmentation. In *Proceedings of the IEEE/CVF Con*ference on Computer Vision and Pattern Recognition, pages 11266–11275, 2021.
- [73] Yuhui Yuan, Lang Huang, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. Ocnet: Object context network for scene parsing. arXiv preprint arXiv:1809.00916, 2018. 3
- [74] Chang-Bin Zhang, Jia-Wen Xiao, Xialei Liu, Ying-Cong Chen, and Ming-Ming Cheng. Representation compensation networks for continual semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7053–7064, 2022. 1, 3
- [75] Gengwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learner with classifier alignment for continual learning on a pre-trained model. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 19148–19158, 2023. 1, 3
- [76] Zekang Zhang, Guangyu Gao, Zhiyuan Fang, Jianbo Jiao, and Yunchao Wei. Mining unseen classes via regional objectness: A simple baseline for incremental segmentation. Advances in neural information processing systems, 35:24340– 24353, 2022. 3
- [77] Zekang Zhang, Guangyu Gao, Jianbo Jiao, Chi Harold Liu, and Yunchao Wei. Coinseg: Contrast inter-and intra-class representations for incremental segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 843–853, 2023. 3
- [78] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 2881–2890, 2017. 3
- [79] Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. Advances in Neural Information Processing Systems, 36:5168–5191, 2023. 1

- [80] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on* computer vision and pattern recognition, pages 633–641, 2017. 5
- [81] Lanyun Zhu, Tianrun Chen, Jianxiong Yin, Simon See, and Jun Liu. Continual semantic segmentation with automatic memory sample selection. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pages 3082–3092, 2023. 3