

Incremental Few-Shot Semantic Segmentation via Multi-Level Switchable Visual Prompts

Anonymous ICCV submission

Paper ID 7277

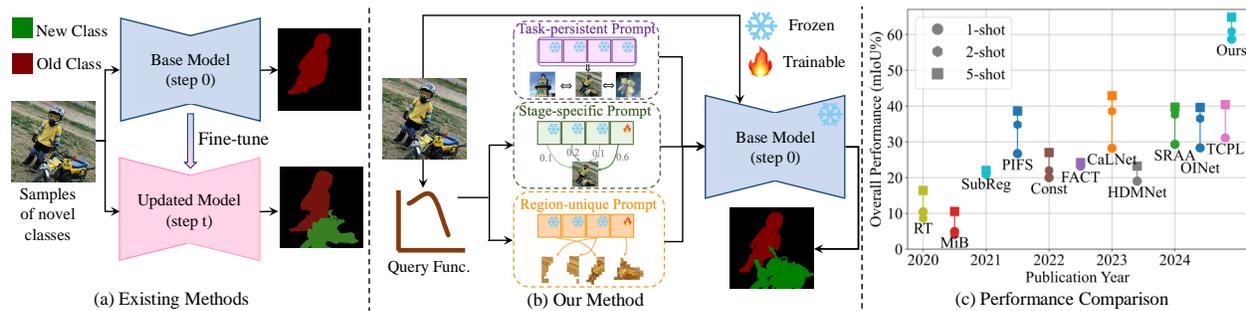


Figure 1. **Comparisons between existing methods and ours.** (a) Existing methods generally fine-tune parameters inherited from old stages and avoid catastrophic forgetting by keeping distribution of old class features. (b) Our method learns novel classes by tailoring multi-granular knowledge to input images achieved by adaptive switching prompts. (c) We compare the overall performance evaluated by Harmonic Mean of mIoU on base and novel classes.

Abstract

001 Existing incremental few-shot semantic segmentation
 002 (IFSS) methods often learn novel classes by fine-tuning pa-
 003 rameters from previous stages. This inevitably reduces the
 004 distinguishability of old class features, leading to catas-
 005 trophic forgetting and overfitting to limited new samples. In
 006 this paper, we propose a novel prompt-based IFSS method
 007 with a visual prompt pool to store and switch multi-granular
 008 knowledge across stages, boosting new class learning ca-
 009 pability. Specifically, we introduce three levels of prompts:
 010 1) Task-persistent prompts: capturing generalizable knowl-
 011 edge shared across stages, such as foreground-background
 012 distributions, to ensure consistent recognition guidance; 2)
 013 Stage-specific prompts: adapting to unique requirements
 014 of each stage by integrating its discriminative knowledge
 015 (e.g., shape difference) with common knowledge from pre-
 016 vious stages; and 3) Region-unique prompts: encoding
 017 category-specific structures (e.g., edges) to accurately guide
 018 the model to retain local details. In particular, we introduce
 019 a prompt switching mechanism that adaptively allocates
 020 the knowledge required for base and new classes, avoid-
 021 ing interference between prompts and preventing catas-
 022 trophic forgetting and reducing increasing computation.
 023 Our method achieves new state-of-the-art performance,

outperforming previous SoTA methods by 30.28% mIoU-N 024
 on VOC and 13.90% mIoU-N on COCO under 1-shot. 025

1. Introduction

026 Incremental few-shot semantic segmentation (IFSS) [1–5] 027
 028 aims to extend segmentation models to novel classes con- 029
 030 tinuously using limited new data without accessing old data. 031
 032 As models expand with few-shot samples, IFSS faces two 033
 034 critical challenges: 1) catastrophic forgetting of old classes 035
 036 and 2) overfitting to limited samples of novel classes. To ad- 037
 038 dress the first, existing methods [1, 6–9] reduce old knowl- 039
 040 edge forgetting by preserving the distribution consistency 041
 042 of old class features, but fine-tuning still harms the abil- 043
 044 ity to distinguish old classes when learning novel classes 045
 046 and cannot solve the dilemma of competition between new 047
 048 and old abilities (see Fig. 1 (a)). Regarding the second is- 049
 050 sue, current approaches facilitate the rapid learning of novel 051
 052 classes by tailoring the prototypes of these classes [1–3] or 053
 054 using surrogate modalities to construct inter-class relation- 055
 056 ships [4, 5]. Unfortunately, due to insufficient model gen- 057
 058 eralization ability and inadequate modal alignment, these 059
 060 methods struggle to effectively differentiate novel classes, 061
 062 resulting in confusion between novel and old classes. 063

046 Recently, prompt-based incremental learning [10–
047 14] has garnered attention in image classification. This
048 paradigm maintains an updatable prompt pool for frozen
049 pre-trained vision transformers, adding new prompts to ac-
050 commodate novel classes and prevent catastrophic forget-
051 ting in succeeding stages. However, its application in se-
052 mantic segmentation remains challenging. On the one hand,
053 embedding multiple background classes into a single con-
054 tinuous feature space leads to a decrease in feature dis-
055 tinguishability and confusion between unseen classes and
056 background features. On the other hand, semantic segmen-
057 tation requires collaborative prompts with multiple gran-
058 ularities to capture both global context and local details, en-
059 abling more meticulous segmentation.

060 In this paper, we innovatively propose an IFSS method
061 based on multi-level switchable prompts (see Fig. 1 (b)).
062 Specifically, we first introduce a prompt-based dense pre-
063 diction framework that leverages well-structured text se-
064 mantics to achieve seamless integration of novel classes
065 with existing ones. First, the framework predefines multiple
066 background semantics to help the model construct clearer
067 background feature boundaries in the feature space, reduc-
068 ing overlap between background and foreground features.
069 Secondly, it introduces a focus decomposition decoder con-
070 sisting of two separators to align text embeddings with
071 pixel features of foreground and background, respectively.
072 Meanwhile, it further enhances alignment by updating vi-
073 sual prompts inserted into the frozen visual encoder, and
074 expands the model by adding new prompts. This design
075 enables the model to process the salient features of novel
076 classes separately and adapt to background changes during
077 incremental learning, effectively alleviating the confusion
078 between novel and background classes.

079 Vanilla visual prompts exhibit several limitations in
080 semantic segmentation, including a lack of fine-grained
081 contextual information, growing computational costs as
082 prompts accumulate, and interference from new prompts
083 causing catastrophic forgetting. To address these chal-
084 lenges, we propose a method of multi-level switchable
085 prompts (MSVP), a prompting strategy to balance knowl-
086 edge retention and adaptability across incremental learn-
087 ing stages. MSVP consists of three levels of prompts: 1)
088 **Task-persistent prompts (TP)**: preserving general knowl-
089 edge shared across stages (e.g., foreground-background dis-
090 tributions); 2) **Stage-specific prompts (SP)**: adapting to the
091 unique requirements of each stage by integrating its dis-
092 criminative knowledge (e.g., shape difference) with com-
093 mon knowledge from previous stages; and 3) **Region-**
094 **unique prompts (RP)**: encoding category-specific struc-
095 tures (e.g., edges) to enhance local detail recovery. TP is
096 frozen after base-stage training, preventing general knowl-
097 edge from being disrupted by new tasks. SP, initialized
098 during base training and continuously expanded, transfers

099 generalizable experience to novel tasks, balancing rigid-
100 ity and plasticity. RP is generated for fine-grained infor-
101 mation aggregation of new classes. These three prompts
102 enable multi-granularity knowledge transfer across stages,
103 enhancing learning ability. To further mitigate interfer-
104 ence and control computational costs, we introduce a flexi-
105 ble prompt-switching mechanism that dynamically tailors
106 prompts for input images. Specifically, a pretrained im-
107 age encoder (e.g., DINOv2 [15]) serves as a query func-
108 tion, generating global and local query features. Global fea-
109 tures aggregate stage-specific prompts through an attention-
110 based mechanism, while local features select region-unique
111 prompts via nearest neighbor matching. This mechanism al-
112 leviates interference between prompts from different stages
113 and keeps a constant number of active prompts, thus avoid-
114 ing catastrophic forgetting and increasing computation.

115 Our method outperforms previous methods by a large
116 margin (see Fig. 1 (c)) proving the effectiveness of the
117 prompt-based dense prediction framework and multi-level
118 switchable visual prompts. In conclusion, our contributions
119 are summarized as:

- 120 1. We propose the first prompt-based IFSS framework,
121 which introduces textual semantics and visual prompts
122 to encode foreground and background classes separately,
123 enabling incremental semantic segmentation.
- 124 2. We propose multi-level switchable visual prompts that
125 customizes multi-granular knowledge tailored to input
126 images, enhancing the model’s ability to learn novel
127 classes while maintaining knowledge of old classes.
- 128 3. Extensive experiments demonstrate the effectiveness of
129 the proposed method. Under the 1-shot condition, it
130 achieves 49.1% mIoU-N on VOC and 25.6% mIoU-N
131 on COCO, setting a new SOTA performance.

132 2. Related Work

133 2.1. Incremental Few-Shot Semantic Segmentation

134 Semantic segmentation [16–21] is a basic computer vi-
135 sion task that involves partitioning an image into mean-
136 ingful segments. Incremental few-shot semantic segmentation
137 aims at continuously learning to segment novel categories
138 via a few given samples, without forgetting knowledge of
139 old categories. To achieve this, PIFS [1] and OINet [3]
140 adopt a distillation training paradigm to avoid forgetting old
141 knowledge and an effective prototype updating strategy of
142 novel categories to learn novel classes. EHNet [2] maintains
143 the old knowledge using the hyperclass representation bank
144 and adaptively updates it to combine novel classes. Instead
145 of distillation or storing old knowledge, CaLNet [4] em-
146 ploys a class-agnostic mask proposal to generate masks for
147 both base and novel categories and integrates language em-
148 bedding into visual features to enrich the representation of
149 a few novel categories. However, the mask proposal mod-

150 ule is prone to overfit base data, causing low recall rates for
 151 novel categories. Different from these methods, we propose
 152 an innovative prompt-based incremental few-shot learning
 153 method, which learns multi-level visual prompts and filters
 154 appropriate prompts for input images, facilitating keeping
 155 old knowledge and learning novel classes.

156 **2.2. Prompt-based Incremental Learning**

157 Prompt-based incremental learning has been studied in image
 158 classification [10–14]. Inspired by VPT [22], these
 159 methods generally freeze the pre-trained parameters and
 160 fine-tune only a set of novel-added learnable prompts at
 161 the incremental stage without a rehearsal buffer to store past
 162 pristine examples for experience replay, which achieves re-
 163 markable performance. L2P [10] is the first method that
 164 introduces this training paradigm, which selects the most re-
 165 levant prompts from a prompt pool in a key-value mechanism.
 166 Instead of merely leveraging task-specific prompts, Dual-
 167 Prompt [11] proposes to introduce general prompts shared
 168 by all tasks, achieving novel SOTA performance. Unlike the
 169 above two methods which learn a pool of key-value pairs
 170 to select learnable prompts, CODA-Prompt [12] introduces
 171 a decomposed prompt that consists of learnable prompt
 172 components that assemble to produce attention-conditioned
 173 prompts and optimizes the model in an end-to-end fashion.
 174 As far as we know, this is the first work that introduces
 175 prompt-based incremental learning methods in IFSS. In par-
 176 ticular, we propose multi-level prompts to meet the needs of
 177 dense prediction tasks for multi-granular contexts.

178 **3. Methods**

179 In this section, we propose a prompt-based IFSS method
 180 that expands and updates multi-granularity switchable
 181 prompts to learn novel classes. It involves a prompt-based
 182 IFSS framework (see Fig. 2) to generate robust class rep-
 183 resentation and enable the model to expand by simply adding
 184 prompts, and an enhanced multi-level prompt generation
 185 method (see Fig. 3) to provide fine-grained knowledge and
 186 switch prompts to expand the model dynamically.

187 **3.1. Prompt-based IFSS Framework**

188 To avoid the interference of background classes on new
 189 class learning, we design an IFSS framework based on
 190 frozen visual-language models, which leverages textual se-
 191 mantics and visual prompts to encode foreground and back-
 192 ground classes separately, thus enabling incremental learn-
 193 ing using prompts. The framework encompasses four key
 194 components: image encoding, text encoding, pixel decod-
 195 ing, and targeted optimization objectives.

196 **Image encoding.** An image is encoded into a sequence
 197 of tokens by a Patch Embedding block as in [23]. Visual
 198 prompts concatenated with image tokens are input into each

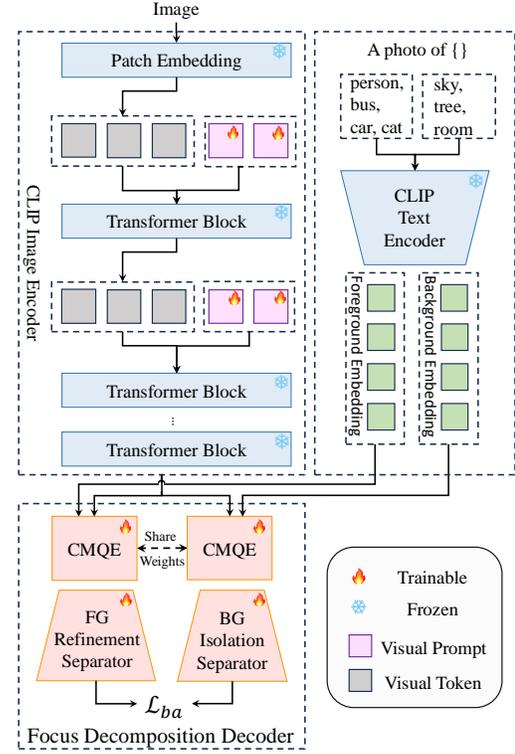


Figure 2. **The proposed IFSS framework.** Image tokens, concatenated with visual prompts, are encoded through successive Transformer blocks. Foreground and background text embeddings, along with image features, are fed into the Focus Decomposition Decoder for final predictions. Two pixel separators are used to distinguish foreground from background and identify specific foreground classes, respectively. The cross-modal query engine (CMQE) generates queries with robust fused modal information.

Transformer block as:

$$[x_{i+1}; p_o] = f_i(x_i, p^{(i)}) \tag{1}$$

where x_i denotes the input tokens output by layer $i - 1$, x_{i+1} denotes the output of layer i , p_o denotes the encoded $p^{(i)}$ which will not be input into next layer. These visual prompts are composed of learnable vectors and are distinct between layers, providing an effective way to inject knowledge into pre-trained models. In the incremental stage, the expanded visual prompts are concatenated with the original ones as supplementary to recognize novel classes, as:

$$[x_{i+1}; p'_o] = f_i(x_i, [p^{(i), t-1}; p_e^{(i)}]) \tag{2}$$

where $p_e^{(i)}$ denotes the expanded prompts of layer i of stage $t - 1$, $p^{(i), t-1}$ denotes the prompts of stage $t - 1$, and $[\cdot; \cdot]$ denotes the concatenation operation. We take this vanilla prompt extension approach as our baseline.

Text encoding. To discriminate background pixels, instead of adding a class of "background", we pre-define a

series of background classes, such as "sky", "room", etc, to pre-isolate background into possible classes. In line with previous works, class names are formatted as "a photo of CLASS_NAME" and encoded using the frozen CLIP Text Encoder to obtain embeddings for both background and foreground classes, denoting as $T_{bg} \in \mathbb{R}^{C_{bg} \times D}$ and $T_{fg} \in \mathbb{R}^{C_{fg} \times D}$, where C_{bg}, C_{fg} represent the number of introduced background classes and foreground classes.

Pixel Decoding. We propose a focus decomposition decoder (FDD) to predict pixel semantics by aligning pixel features with class embeddings derived from the language modality. This decoder is composed of a cross-modal query engine (CMQE), inspired by [24], which generates class queries with robust generalization capabilities, along with two separators: one for identifying foreground and background pixels and the other for classifying the semantics of foreground pixels. The dual-decoder architecture enables the model to separately process salient features of new classes and adapt to background variations during incremental learning, making it easier to seamlessly embed features of novel classes into existing class representations.

CMQE integrates visual and linguistic information to generate robust class-specific queries, enhancing segmentation performance on both base and novel classes. Denoting $T_{fg} = \{t_0, t_1, \dots, t_{C_{fg}}\}$ and $g \in \mathbb{R}^D$, the class-specific query can be denoted as $Q_{fg} = \{\hat{q}_0, \hat{q}_1, \dots, \hat{q}_i, \dots, \hat{q}_{C_{fg}}\}$, where \hat{q}_i is calculated as:

$$\hat{q}_i = \text{MLP}([t_i \odot g; t_i]) \quad (3)$$

where g denotes the global feature, t_i denotes the class embedding of the i -th class, \odot denotes the Hadamard product and MLP is used to align the dimension with pixel features. The background queries Q_{bg} is calculated the same as Q_{fg} .

Subsequently, class-specific queries Q_{bg}, Q_{fg} are forward to BG isolation separator and FG refinement separator to align with pixel features, respectively. We employ a cascaded cross-attention structure to facilitate this alignment, where Q_{bg} and Q_{fg} are the query, and pixel features $\mathbf{P} \in \mathbb{R}^{HW \times d}$ are the keys and values. Finally, we take the scaled dot-product attention in the separator's last block as the final semantic masks, as:

$$\mathbf{M}_{bg} = \frac{\phi_q^{bg}(\hat{\mathbf{Q}}_{bg})\phi_k^{bg}(\mathbf{P})^T}{\sqrt{d_k}}, \mathbf{M}_{fg} = \frac{\phi_q^{fg}(\hat{\mathbf{Q}}_{fg})\phi_k^{fg}(\mathbf{P})^T}{\sqrt{d_k}}, \quad (4)$$

where $\hat{\mathbf{Q}}_{bg}, \hat{\mathbf{Q}}_{fg}$ denote queries aligned by cross-attention blocks before the last block, ϕ denotes the linear projection, d_k is the dimension of the keys, and $\mathbf{M}_{bg} \in \mathbb{R}^{HW \times C_{bg}}, \mathbf{M}_{fg} \in \mathbb{R}^{HW \times C_{fg}}$ are the scores of each class.

Optimization objectives. Since we pre-define multiple background classes while the ground truth includes only a generic 'background' label, we propose a background aggregation loss \mathcal{L}_{ba} , to address this discrepancy. We regard

the maximum scores of all background classes at each pixel as final background scores $\hat{\mathbf{M}}_{bg} \in \mathbb{R}^{HW \times 1}$, and we regard the maximum logits of all foreground classes at each pixel to represent the likelihood of being a foreground pixel $\hat{\mathbf{M}}_{fg} \in \mathbb{R}^{HW \times 1}$, as:

$$\begin{aligned} \hat{\mathbf{M}}_{bg} &= \max_i \mathbf{M}_{bg}^{j,i}, j = 1, 2, \dots, HW \\ \hat{\mathbf{M}}_{fg} &= \max_i \mathbf{M}_{fg}^{j,i}, j = 1, 2, \dots, HW. \end{aligned} \quad (5)$$

An intuitive and vanilla way is to optimize two separators independently, as:

$$\begin{aligned} \mathcal{L}_{van} &= \mathcal{L}_{seg}(y, [1 - \hat{\mathbf{M}}_{fg}; \mathbf{M}_{fg}]) \\ &\quad + \alpha_1 \mathcal{L}_{seg}(\bar{y}, [\hat{\mathbf{M}}_{bg}; 1 - \hat{\mathbf{M}}_{bg}]), \end{aligned} \quad (6)$$

where $y \in \mathbb{R}^{HW \times (C_{fg}+1)}$ denotes the one-hot labels of all classes including background and $\bar{y} \in \mathbb{R}^{HW \times 2}$ denotes the one-hot labels of foreground and background, α_1 is the weight to balance representation of FG and BG, and \mathcal{L}_{seg} denotes the widely used pixel-wise classification loss as in [24] [18], the combination of focal loss [25] and dice loss [26]. To mitigate the potential issue of misaligned optimization directions, we introduce a more flexible loss function that prevents feature confusion among novel, old, and background classes during incremental learning, as:

$$\mathcal{L}_{ba} = \mathcal{L}_{seg}(y, [\hat{\mathbf{M}}_{bg}; \mathbf{M}_{fg}]) + \alpha_2 \mathcal{L}_{seg}(\bar{y}, [\hat{\mathbf{M}}_{bg}; \hat{\mathbf{M}}_{fg}]), \quad (7)$$

which jointly constrains the masks output by the two heads, on the one hand to distinguish the specific categories of foreground pixels as the first term of \mathcal{L}_{ba} , and on the other hand to separate foreground and background pixels to alleviate the problem of class imbalance as the second term. This dual constraint encourages accurate pixel classification and facilitates model expansion.

3.2. Multi-level Switchable Visual Prompts

Simply adding visual prompts during incremental learning has shown improvements (see Sec. 4.4). However, it presents three challenges: 1) computation increase: an increase in visual prompts raises computational complexity, 2) information dilution: stage-wise incremental prompts dilute the information of each stage, leading to knowledge forgetting and a diminished capacity to learn novel classes, and 3) insufficient granularity: a single level of prompts fails to meet the multi-granularity contextual needs essential for semantic segmentation. To address these, we propose to switch appropriate multi-level visual prompts tailored for input images, as shown in Fig. 3. It includes task-persistent prompts, stage-specific prompts, and region-unique prompts, alongside a flexible prompt switching mechanism.

Task-persistent Prompts (TP). According to the theory of Complementary Learning Systems (CLS) [27, 28],

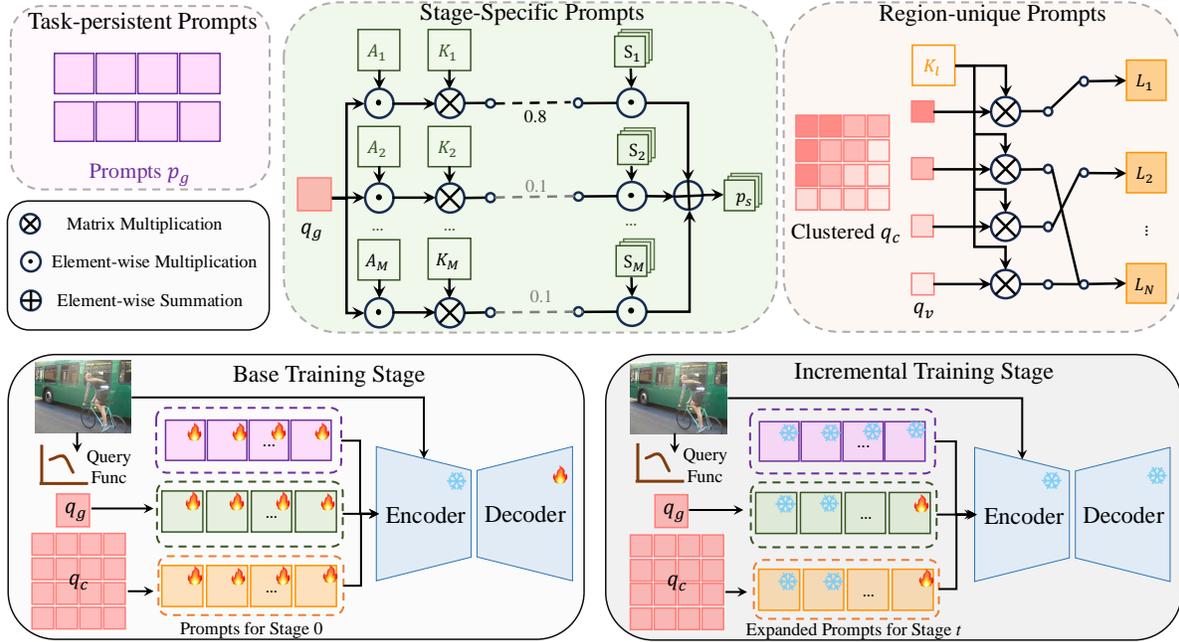


Figure 3. **The pipeline of the proposed multi-level switchable visual prompts.** Images are inputs into the query function to obtain global query features q_g and pixel-wise query features q_c . Stage-specific prompts p_s are generated by an attention-like integration way through q_g . Region-unique prompts p_l are generated by nearest neighbor matching through clustered q_c . Image tokens, concatenated with these selected prompts, are input into pre-trained models to produce predictions. At incremental training stage, the model extends by fine-tuning newly-added stage-specific prompts and region-unique prompts.

312 humans learn continually via the synergy between the hip-
 313 pocampus and the neocortex. The former learns pattern-
 314 separated representation on specific experiences while the
 315 later focuses on learning more general and transferable
 316 representation to enhance the capability of learning future
 317 stages. Inspired by this, we propose task-persistent prompts
 318 to mitigate catastrophic forgetting and capture generaliz-
 319 able representation to facilitate encoding shared seman-
 320 tic structures or relationships for novel classes, including
 321 foreground-background distribution, shared contours and
 322 edges between categories, and so on. Specifically, $p_g^{(i)} \in$
 323 $\mathbb{R}^{L_g \times D}$ are learnable vectors with pre-defined sequence
 324 length L_g and embedding dimension D , which is trained
 325 during base training stage to learn generalizable knowledge
 326 and frozen during incremental stages.

327 **Stage-specific Prompts (SP).** Although TP maintains
 328 invariant core knowledge, they lack flexibility for rapid
 329 adaptation in incremental stages. To address this, we
 330 propose stage-specific prompts that dynamically guide the
 331 model to adapt fine-grained distinctions to learn new stages.
 332 However, limited training samples hinder learning robust
 333 stage-specific representations. To overcome this, we de-
 334 sign a knowledge inheritance mechanism that integrates dis-
 335 criminative information from current stages with preserved
 336 prior knowledge. This is achieved through an attention-
 337 like mechanism that aggregates stage-specific prompts based on

the correlation between knowledge required for inputting
 images and learned knowledge across stages.

Specifically, we take a pre-trained image encoder (e.g.,
 DINOv2 [15]) as query function to obtain the global fea-
 tures $q_g \in \mathbb{R}^{1 \times D}$ to aggregate relevant knowledge stored
 in stage-specific prompts $\mathbf{S}^{(i)} \in \mathbb{R}^{M \times L_s \times D}$, where M
 denotes the number of current training stage and L_s denotes
 the number of learnable vectors for each stage. We then cal-
 culate the correlation γ between the knowledge required by
 the input image and knowledge learned from all stages, as:

$$\gamma = \text{Softmax}(\langle q_g \odot \mathbf{A}^{(i)}, \mathbf{K}^{(i)} \rangle / \tau) \in \mathbb{R}^M \quad (8)$$

where $\langle \cdot \rangle$ denotes cosine similarity, \odot denotes element-
 wise multiplication, τ denotes the temperature coefficient,
 $\mathbf{K}^{(i)} \in \mathbb{R}^{M \times D}$ are keys corresponding to each stage-
 specific prompt. To allow the query to focus on specific pat-
 terns, an attention vector $\mathbf{A}^{(i)} \in \mathbb{R}^{M \times D}$ corresponding to
 each stage-specific prompt is added. For example, a prompt
 designed for recognizing car textures can focus on details
 like headlights while ignoring unrelated features. Addition-
 ally, in contrast to [12], where the weight vector is derived
 directly from the cosine similarity, we employ a normal-
 ized similarity computed with softmax as the weight vector,
 which ensures that stage-specific prompts unrelated to the
 input image are not aggregated into the final prompts. Note
 that all the vectors $\mathbf{S}^{(i)}, \mathbf{A}^{(i)}, \mathbf{K}^{(i)}$ are learnable vectors.

After obtaining γ , stage-specific prompts $p_s^{(i)}$ for the input image are aggregated as:

$$p_s^{(i)} = \gamma \cdot \mathbf{S}^{(i)} = \gamma_1 \mathbf{S}_1^{(i)} + \gamma_2 \mathbf{S}_2^{(i)} + \dots + \gamma_M \mathbf{S}_M^{(i)}, \quad (9)$$

where $p_s^{(i)} \in \mathbb{R}^{L_s \times D}$. Thus, knowledge required by each image is integrated according the similarity between q_g and $\mathbf{K}^{(i)}$. For example, assuming that the input image contains class 'COW', the knowledge required to segment this class is similar to that of class 'SHEEP' learned from stage s . γ between the query and the key corresponding to stage s and current stage can be calculated as 0.8, 0.2. It means most of knowledge required by the input image can be inherited from stage s and combined with discriminative information learned from current stage.

Region-unique Prompts (RP). As a dense prediction task, semantic segmentation requires more fine-grained knowledge of typical structures for specific categories used for storing local details. To this end, we propose to query the best region-unique prompts conditioned on the distance between the keys and local features of the input image.

Firstly, we obtain the pixel-wise features q_c from the query function. Secondly, querying a prompt for pixel-wise features undoubtedly increases the computational complexity and generates redundant noise. Instead, we apply the KMeans algorithm to cluster these features into h categories and let the clustered centroids $q_v \in \mathbb{R}^{h \times D}$ represent the local information of the image. Thirdly, q_v is utilized to query local prompts corresponding to the key $\mathbf{K}_l^{(i)} \in \mathbb{R}^{N \times D}$ closest to q_v , where $N = C_{bg} + C_{fg}$. We have a set of local prompts $\mathbf{L}^{(i)} \in \mathbb{R}^{N \times D}$ for layer i . This process can be formalized as:

$$q_v = \text{KMeans}(q_c, h) \quad (10)$$

$$idx = \text{argmax}_{1 \leq j \leq N} q_v \mathbf{K}_{l,j}^{(i)} \quad (11)$$

$$p_l^{(i)} = \mathbf{L}^{(i)} [idx, :] \quad (12)$$

where $p_l^{(i)} \in \mathbb{R}^{h \times D}$ denotes the selected local prompts when given the query q_c . Thereby, a prompt is assigned to each background and foreground class and is queried using the image's local information via nearest matching.

However, the argmax operator prevents gradient back-propagation, necessitating additional supervision. To enable end-to-end training, we employ a Gumbel_Softmax operation to replace the above procedure, similar to [29], which can be simply formulated as

$$p_l^{(i)} = \text{Gumbel_Softmax}(q_v \mathbf{K}_l^{(i)}) \mathbf{L}^{(i)}. \quad (13)$$

The detailed process is described in Appendix. A.1.

Furthermore, we incorporate attention masks into these visual prompts to prevent prompts specific to individual pixel segments from influencing other segments. When

given pixel-wise query features q_c , the feature similarity between pixels can be calculated as $\mathbf{S}_c \in \mathbb{R}^{HW \times HW}$. The pixels with similarity higher than the threshold ζ are the pixels prompted by region-unique prompts, while the rest are the masked pixels, which can be formulated as:

$$\hat{\mathbf{S}}_c = \begin{cases} 0, & \mathbf{S}_c > \zeta \\ -\infty, & \mathbf{S}_c \leq \zeta \end{cases}, \quad (14)$$

where a higher ζ means that the area prompted by the region is smaller, while the opposite means it is larger. Finally, extracting the mask corresponding to centroids q_v from $\hat{\mathbf{S}}_c$ and inserting it into the attention mask of the self-attention structure can limit the scope of the region-unique prompts.

Incremental training. During incremental training, task-persistent prompts, stage-specific prompts of previous stages and region-unique prompts of old classes remain frozen. We expand $\mathbf{S}^{(i)}$, $\mathbf{A}^{(i)}$, $\mathbf{K}^{(i)}$ to learn knowledge of new stages and expand $\mathbf{K}_l^{(i)}$, $\mathbf{L}^{(i)}$ to learn local details of new classes, and exclusively train the newly expanded components. The detail is formulated in Appendix A.5.

The multi-level switchable prompts generate multi-granular contextual information essential for semantic segmentation, enhancing new-class adaptability. And it addresses two critical limitations of the conventional method by simply adding prompts: 1) preventing information interference and dilution through adaptive prompts selection thus reducing catastrophic forgetting, (See Tab. 3), and 2) alleviating increasing computation by maintaining fixed input sequence length (See Appendix B.6).

4. Experiments

4.1. Datasets

We conduct experiments on Pascal VOC 2012 [40] and COCO [41, 42] as in previous works [1, 2, 4]. VOC contains 20 classes and one background class. In COCO, we use the 80 classes and the residual classes are labeled as background. We consider 15 and 60 of the classes as base and 5 and 20 classes as novel, for VOC and COCO respectively. The protocols start with pretraining on base classes and multiple steps on novel classes in line with [1, 4], i.e., 5 steps of 1 novel class on VOC and 4 steps of 5 novel classes on COCO. We divide the VOC dataset into 4 folds of 5 classes each and the COCO dataset into 4 folds of 20 classes each. We run experiments 5 times, with each experiment considering one fold at a time as the set of novel classes. In each setting, we explore incremental steps using 1, 2, or 5 images. Following the previous methods [1, 2, 4], we evaluate the performance via three metrics based on the mean Intersection over Union (mIoU): mIoU on base classes (mIoU-B), mIoU on novel classes (mIoU-N) and the Harmonic Mean (HM) of the two.

Table 1. Comparison with SOTA methods on VOC. Bold/Underline indicate SoTA/The Second Best.

Method	1-shot			2-shot			5-shot			
	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	
FSC	WI [30]	66.6	16.1	25.9	66.6	19.8	30.5	66.6	21.9	33.0
	DWI [31]	67.2	16.3	26.2	67.5	21.6	32.7	67.6	25.4	36.9
	RT [32]	49.2	5.8	10.4	36.0	4.9	8.6	45.1	10.0	16.4
FSS	AMP [33]	58.6	14.5	23.2	58.4	16.3	25.5	57.1	17.2	26.4
	SPN [34]	49.8	8.1	13.9	56.4	10.4	17.6	61.6	16.3	25.8
IL	LwF [35]	42.1	3.3	6.2	51.6	3.9	7.3	59.8	7.5	13.4
	ILT [36]	43.7	3.3	6.1	52.2	4.4	8.1	59.0	7.9	13.9
	MiB [6]	43.9	2.6	4.9	51.9	2.1	4.0	60.9	5.8	10.5
IFL	SubReg [37]	55.4	13.2	21.3	56.7	12.7	20.8	59.7	13.5	22.0
	Const [38]	58.4	12.1	20.0	61.3	13.4	22.0	62.2	17.2	27.0
	FACT [39]	57.0	14.6	23.2	57.4	15.1	23.9	58.8	15.2	24.2
PIFS [1]	64.1	16.9	26.7	65.2	23.7	34.8	64.5	27.5	38.6	
OINet [3]	66.1	18.0	28.3	66.3	25.2	36.5	66.4	28.2	39.6	
CaLNet [4]	74.2	17.4	28.2	74.4	26.1	38.6	74.7	<u>30.1</u>	42.9	
SRAA [5]	66.4	<u>18.8</u>	<u>29.3</u>	65.1	<u>26.4</u>	37.6	64.3	<u>28.7</u>	39.7	
Ours	<u>73.04</u>	49.08	58.71	<u>73.21</u>	52.19	60.93	<u>73.36</u>	58.13	64.86	

Table 2. Comparison with SOTA methods on COCO. Bold/Underline indicate SoTA/The Second Best.

Method	1-shot			2-shot			5-shot			
	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	
FSC	WI [30]	46.3	8.3	14.0	46.5	9.3	15.4	46.3	10.3	16.8
	DWI [31]	46.2	9.2	15.3	46.5	11.4	18.3	46.6	14.5	22.1
	RT [32]	38.4	5.2	9.1	43.8	10.1	16.4	44.1	16.0	23.5
FSS	AMP [33]	36.6	7.9	13.1	36.0	9.2	14.6	33.2	11.0	16.5
	SPN [34]	40.3	8.7	14.3	41.7	12.5	19.2	41.4	18.2	25.3
IL	LwF [35]	41.0	4.1	7.4	42.7	6.5	11.3	42.3	12.6	19.4
	ILT [36]	43.7	6.2	10.8	47.1	10.0	16.5	45.3	15.3	22.8
	MiB [6]	40.4	3.1	5.8	42.7	5.2	9.3	43.8	11.5	18.2
IFL	SubReg [37]	38.4	8.0	13.2	39.5	10.1	16.0	40.0	10.3	16.4
	Const [38]	39.0	8.2	13.6	40.6	11.4	17.8	41.1	11.3	17.7
	FACT [39]	37.9	8.6	14.0	38.9	11.7	18.0	39.4	12.3	18.7
PIFS [1]	40.4	10.4	16.6	40.1	13.1	19.8	41.1	18.3	25.3	
OINet [3]	41.4	<u>11.7</u>	<u>18.2</u>	41.5	14.4	21.4	41.5	<u>19.7</u>	26.7	
CaLNet [4]	<u>48.4</u>	10.6	17.4	<u>48.5</u>	13.4	21.0	<u>48.6</u>	18.6	26.9	
SRAA [5]	40.7	11.3	17.7	40.5	<u>15.2</u>	<u>22.1</u>	41.0	<u>19.7</u>	26.6	
Ours	48.85	25.60	33.59	48.52	28.05	35.54	48.61	32.38	38.86	

4.2. Implementation Details

In this work, we choose ViT-B [23] as the image encoder, and pretrained ViT-B of DINOv2 [15] as our query function, which can output accurate q_g and q_c to query appropriate prompts. Meanwhile, to make training stable and provide meaningful initial keys for region-unique prompts, we take text embeddings encoded by CLIP text encoder in a manner of CoOP [43] as the keys for region-unique prompts, which is detailed in Appendix. A.2. Stage by stage, the backbone and the query function are frozen, all prompts as well as the decoder are trainable. During incremental training, we freeze all parameters but expand and update the stage-specific prompts and region-unique prompts, which is described in Appendix. A.5. We add orthogonality constraints to parameters of stage-specific prompts to avoid interference between existing and new knowledge and reduce catastrophic forgetting. During base training, the model is trained for 20k iterations on VOC and 80k iterations on COCO. During incremental training, for both VOC and COCO, the model is trained for 400 iterations per step.

4.3. Compare with State-of-the-art methods

We mainly conduct comparison between few-shot classification methods (FSC) [30–32], few-shot semantic segmentation methods (FSS) [33, 34], incremental learning methods (IL) [37–39], and incremental few-shot semantic segmentation methods (IFSS).

Evaluation on VOC. The results of 1-shot, 2-shot and 5-shot experiments are presented in Tab. 1. In general, our method achieves novel SOTA performance on novel classes of 49.08% , 52.19% and 58.13% mIoU for 1-shot, 2-shot, and 5-shot scenarios respectively. Additionally, our method also attains the best overall performance with HM scores of 58.71% and 60.93% and 64.86%. Comparing our methods with other methods, it is evident that FSC meth-

ods excel in retaining knowledge of base classes, achieving competitive performance on these classes, since FSC methods, such as WI [30] and DWI [31], expand classifiers using class prototypes, thereby preventing the corruption of learned knowledge. In contrast, FSS and IL methods perform poorly on both base and novel classes. While meta-learning helps FSS methods adjust to novel classes, they struggle to retain old knowledge. And IL methods require many novel samples, resulting in low performance on few-shot tasks. Our method remarkably outperforms prior SOTA IFSS method SRAA [5] on novel classes by 30.28%, 25.79%, and 29.43% mIoU in 1-shot, 2-shot and 5-shot scenarios respectively.

Evaluation on COCO. The results are presented in Tab. 2. In general, our method achieves a new SOTA performance on novel classes of 25.60% , 28.05% and 32.38% mIoU for 1-shot, 2-shot, and 5-shot scenarios respectively. Although methods like PIFS [1], OINet [3], and CaLNet [4] enhance novel class learning with refined prototypes or textual knowledge, they remain limited in representing pixel features for novel classes. In contrast, our method customizes contextual information per image, effectively improving novel class representation.

4.4. Ablation Studies

Table 3. Ablation on Multi-level Prompts. Table 4. Ablation on the Framework.

SA TP SP RP	1-shot			method	1-shot		
	Base	Novel	HM		Base	Novel	HM
✓	65.80	47.32	55.05	Ours	73.04	49.08	58.71
✓ ✓	72.82	48.45	58.18	- FDD	72.86	44.89	55.55
✓ ✓ ✓	63.26	47.82	54.46	- \mathcal{L}_{ba}	70.53	48.73	57.63
✓ ✓ ✓	71.08	46.70	56.36	- CMQE	70.50	22.20	33.76
✓ ✓ ✓	73.04	49.08	58.71				

516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568

Ablation on different prompts. We simply add visual prompts (SA) of each layer of the proposed framework to learn novel classes, which is regarded as our baseline as in Eq. (2). We perform ablations on SA and the proposed prompts. To ensure fairness, each experiment uses the same total number of prompts. Tab. 3 reports results on the 1-shot VOC benchmark. The baseline achieves 65.80% mIoU-B and 47.32% mIoU-N by expanding with a fixed number of prompts, yielding competitive results. Adding TP and SP improves performance of both base and novel classes. TP conveys universal knowledge across stages, while SP provides discriminative representations of the current stage and inherits relevant knowledge from previous ones, enhancing learning ability. With finer-grained RP, novel class performance increases to 49.08% mIoU-N, highlighting the benefits of multi-granular contextual information. Additionally, using only SP and RP leads to a performance decline, as SP may weaken generalizable knowledge from previous stages when adapting to new ones. Introducing TP preserves this knowledge from the data-rich base stage, improving overall performance. We further prove MSVP shows stronger learning capability for novel classes with more training samples in Appendix B.1.

Ablation on the framework. We validate the importance of each component of the framework by removing them one at a time, as in Tab. 4. It can be observed that replacing FDD with a single head markedly reduces performance, particularly on novel classes. That’s because FDD decompose semantic segmentation into foreground refinement and background isolation, making the model focus on the former and reducing background interference. \mathcal{L}_{ba} outperforms \mathcal{L}_{van} by 2.51% mIoU on base classes, for the reason that \mathcal{L}_{van} optimizes the two separators independently causing misaligned optimization directions. Besides, as CMQE generates queries with high generalization ability, it enhances the model’s capacity to expand, with an improvement of 26.88% mIoU on novel classes.

Ablation on the number of task-persistent prompts L_g . We conduct an experiment on the number of task-persistent prompts as shown in Fig. 4 (a). The graph shows that the performance for both base and novel classes initially increases, peaking at 24 prompts, and then declines as the prompt count increases further. This trend occurs because few prompts provides insufficient shared knowledge, making it harder to prompt the model effectively. Conversely, using too many prompts causes the model to overfit on the base classes, which restricts its flexibility and reduces its capacity to generalize to novel classes.

Ablation on the number of vectors of per stage-specific prompt L_s . We conduct an experiment on the number of vectors of per stage-specific prompt as shown in Fig. 4 (b). Our method achieves its best performance when L_s is 8. This is because when the number of prompts is in-

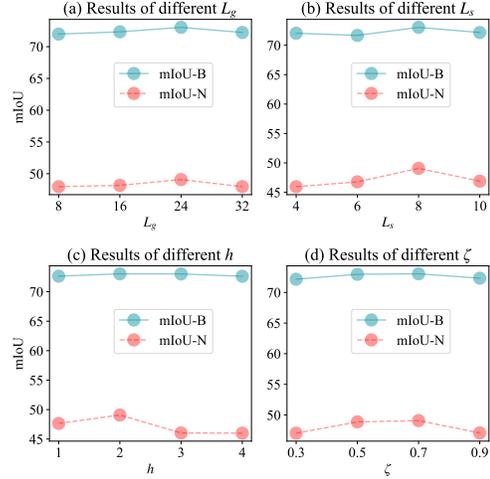


Figure 4. Ablation studies on hyperparameters.

sufficient, stage-specific knowledge is lacking. Conversely, when the number of prompts is excessive, the model is prone to overfitting a small number of novel class samples.

Ablation on the cluster h and the threshold ζ . We conduct two experiments on the hyper-parameters of region-unique prompts, as shown in Fig. 4 (c) and (d). Fig. 4 (c) indicates that the novel class achieves the highest performance when h is 2. Both too few and too many clusters result in decreased performance. Specifically, a low h leads to insufficient granularity in region-unique level prompts, while a high h results in overly fine granularity, which can introduce redundant information. Fig. 4 (d) indicates that when ζ is 0.7, performance on novel classes gets best. A low ζ increases the area of the region-unique prompt, augmenting incorrect pixels, while a high ζ reduces the prompt area, resulting in too few pixels being augmented.

More ablation studies and visualizations are shown in Appendix B and C, including ablation on orthogonality constraints, query function and computation efficiency.

5. Conclusion

In this paper, we propose a novel multi-level prompt-based IFSS method that incorporates a visual prompt pool to store and switch multi-granular knowledge across different stages to enhance the incremental learning. We first design a prompt-based IFSS framework, which leverages textual semantics and visual prompts to encode foreground and background classes separately, enabling incremental semantic segmentation using prompts. Further, we introduce multi-level visual prompts with a switching mechanism to provide the model with multi-granularity contextual information tailored to the image content, thus instructing the model to learn novel classes effectively without forgetting old classes. Extensive experiments on various datasets demonstrate the effectiveness of our method.

604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660

References

[1] Fabio Cermelli, Massimiliano Mancini, Yongqin Xian, Zeynep Akata, and Barbara Caputo. Prototype-based incremental few-shot semantic segmentation. In *Proceedings of the 32nd British Machine Vision Conference*, November 2021. 1, 2, 6, 7

[2] Guangchen Shi, Yirui Wu, Jun Liu, Shaohua Wan, Wenhai Wang, and Tong Lu. Incremental few-shot semantic segmentation via embedding adaptive-update and hyper-class representation. In *Proceedings of the 30th ACM international conference on multimedia*, pages 5547–5556, 2022. 2, 6

[3] Lianlei Shan, Wenzhang Zhou, Wei Li, and Xingyu Ding. Organizing background to explore latent classes for incremental few-shot semantic segmentation. *arXiv preprint arXiv:2405.19568*, 2024. 1, 2, 7

[4] Leo Shan, Wenzhang Zhou, and Grace Zhao. Incremental few shot semantic segmentation via class-agnostic mask proposal and language-driven classifier. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8561–8570, 2023. 1, 2, 6, 7

[5] Yuan Zhou, Xin Chen, Yanrong Guo, Jun Yu, Richang Hong, and Qi Tian. Advancing incremental few-shot semantic segmentation via semantic-guided relation alignment and adaptation. In *International Conference on Multimedia Modeling*, pages 244–257. Springer, 2024. 1, 7

[6] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Buló, 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. 1, 7

[7] 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.

[8] 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.

[9] 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. 1

[10] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022. 2, 3

[11] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European Conference on Computer Vision*, pages 631–648. Springer, 2022. 3

[12] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11909–11919, 2023. 3, 5

[13] Chenxi Liu, Zhenyi Wang, Tianyi Xiong, Ruibo Chen, Yihan Wu, Junfeng Guo, and Heng Huang. Few-shot class incremental learning with attention-aware self-adaptive prompt. *arXiv preprint arXiv:2403.09857*, 2024.

[14] Martin Menabue, Emanuele Frascaroli, Matteo Boschini, Enver Sangineto, Lorenzo Bonicelli, Angelo Porrello, and Simone Calderara. Semantic residual prompts for continual learning. *arXiv preprint arXiv:2403.06870*, 2024. 2, 3

[15] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 5, 7, 3

[16] Kaige Li, Qichuan Geng, Maoxian Wan, Xiaochun Cao, and Zhong Zhou. Context and spatial feature calibration for real-time semantic segmentation. *IEEE Transactions on Image Processing*, 32:5465–5477, 2023. 2

[17] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018.

[18] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. 2022. 4

[19] Kaige Li, Qichuan Geng, Maoxian Wan, Xiaochun Cao, and Zhong Zhou. Context and spatial feature calibration for real-time semantic segmentation. *IEEE Transactions on Image Processing*, 32:5465–5477, 2023.

[20] Kaige Li, Qichuan Geng, and Zhong Zhou. Exploring scale-aware features for real-time semantic segmentation of street scenes. *IEEE Transactions on Intelligent Transportation Systems*, 25(5):3575–3587, 2024.

[21] 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. 2

[22] 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

[23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3, 7

[24] Ziqin Zhou, Yinjie Lei, Bowen Zhang, Lingqiao Liu, and Yifan Liu. Zegclip: Towards adapting clip for zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Con-*

661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717

- ference on Computer Vision and Pattern Recognition, pages 11175–11185, 2023. 4
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 4
- [26] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. 4
- [27] James L McClelland, Bruce L McNaughton, and Randall C O’Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, 102(3):419, 1995. 4
- [28] Dharshan Kumaran, Demis Hassabis, and James L McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in cognitive sciences*, 20(7):512–534, 2016. 4
- [29] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022. 6
- [30] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5822–5830, 2018. 7
- [31] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4367–4375, 2018. 7
- [32] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 266–282. Springer, 2020. 7
- [33] Mennatullah Siam, Boris N Oreshkin, and Martin Jagersand. Amp: Adaptive masked proxies for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5249–5258, 2019. 7
- [34] Yongqin Xian, Subhabrata Choudhury, Yang He, Bernt Schiele, and Zeynep Akata. Semantic projection network for zero-and few-label semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8256–8265, 2019. 7
- [35] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. 7
- [36] 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. 7
- [37] Afra Feyza Akyürek, Ekin Akyürek, Derry Tanti Wijaya, and Jacob Andreas. Subspace regularizers for few-shot class incremental learning. *arXiv preprint arXiv:2110.07059*, 2021. 7
- [38] Michael Hersche, Geethan Karunaratne, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas Rahimi. Constrained few-shot class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9057–9067, 2022. 7
- [39] Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, Liang Ma, Shiliang Pu, and De-Chuan Zhan. Forward compatible few-shot class-incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9046–9056, 2022. 7
- [40] Mark Everingham and John Winn. The pascal visual object classes challenge 2012 (voc2012) development kit. *Pattern Anal. Stat. Model. Comput. Learn., Tech. Rep.*, 2007(1-45):5, 2012. 6
- [41] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocomstuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 6
- [42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 6
- [43] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision (IJCV)*, 2022. 7, 1
- [44] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 3
- [45] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 3

Incremental Few-Shot Semantic Segmentation via Multi-Level Switchable Visual Prompts

Supplementary Material

812 A. Implementation Details

813 In this section, we first elaborate on the process of region-
814 unique prompt switching using `Gumbel_Softmax` and the
815 design of the keys \mathbf{K}_l of region-unique prompts. Secondly,
816 we outline the settings of the proposed method, including
817 the optimizer, learning rate, hyper-parameters, to ensure
818 the reproducibility of the method. Thirdly, we describe
819 the overall data flow and model structure. Finally, we ex-
820 plain the ways to perform incremental learning to expand
821 the model. Due to the double-blind principle, **we will open**
822 **source the code at our github after the review.**

823 A.1. Detailed Region-unique Prompt Switching 824 Process

825 To enable end-to-end training, we employ a
826 `Gumbel_Softmax` operation to switch the appropri-
827 ate region-unique prompts. We first calculate the similarity
828 between q_v and the keys of region-unique prompts, as

$$829 \quad d_{k,j} = \frac{\exp(q_{v,k} \mathbf{K}_{l,j} + \epsilon_k)}{\sum_{m=1}^N \exp(q_{v,k} \mathbf{K}_{l,m} + \epsilon_m)} \quad (\text{A1})$$

830 where $\{\epsilon_k\}$ are i.i.d random samples drawn from the
831 $\text{Gumbel}(0, 1)$ distribution. We compute the region-
832 unique prompt to assign a centroid to by taking the one-hot
833 operation of it `argmax` over all the keys. Since the one-
834 hot assignment operation via `argmax` is not differentiable,
835 we instead use the straight through trick in to compute the
836 assignment matrix as

$$837 \quad \hat{d} = \text{one-hot}(d_{\text{argmax}}) + d - \text{sg}(d) \quad (\text{A2})$$

838 where `sg` is the stop gradient operator. With the straight
839 through trick, \hat{d} has the one-hot value of assignment to a
840 single region-unique prompt, but its gradient is equal to
841 the gradient of d , which makes the whole procedure dif-
842 ferentiable and end-to-end trainable. After assigning q_v to
843 keys of region-unique prompts, we can easily get the region-
844 unique prompts response to q_v by merging all prompts, as:

$$845 \quad p_{l,k}^{(i)} = \frac{\sum_{j=1}^N \hat{d}_{k,j} \mathbf{L}_j^{(i)}}{\sum_{j=1}^N \hat{d}_{k,j}}. \quad (\text{A3})$$

846 This approach effectively solves the problem of gradient
847 backpropagation by transforming the `argmax` process into a
848 discrete variable sampling process.

Table A1. Ablation on the design of keys of region-unique prompts.

Method	1-shot		
	Base	Novel	HM
Rand	72.14	47.53	57.30
CoOP	73.04	49.08	58.71

A.2. The Design of the Keys of Region-unique prompts

851 As described in Sec. 4.2, to ensure stable training and pro-
852 vide meaningful initial keys for region-unique prompts, we
853 leverage text embeddings generated by the CLIP text en-
854 coder following the CoOP[43] approach. Specifically, each
855 region-unique prompt, which corresponds to a particular
856 class, is assigned a key that aids in aligning the local fea-
857 tures of input images. For each key, we use the CLIP text
858 encoder to produce stable and meaningful embeddings by
859 feeding it a prompt $\sigma = [V]_1[V]_2 \dots [V]_n[CLASS]$, where
860 $[V]_i$ represents vectors with the same dimension as word
861 embeddings, n is a hyperparameter specifying the number
862 of context tokens, and $[CLASS]$ denotes the class name's
863 word embedding. By processing the prompt σ through the
864 text encoder, we obtain a key tailored to each region-unique
865 prompt. This method provides an effective initial value for
866 the keys, mitigating the convergence issues often caused
867 by random initialization, especially in few-shot scenarios.
868 Moreover, it incorporates text modality knowledge, reduc-
869 ing the risk of overfitting to the limited samples in few-shot
870 novel classes.

871 We further carry out an experiment to compare the per-
872 formance with and without the keys generated by CoOP
873 on VOC, as shown in Tab. A1. The table shows that with
874 this key generation method, the model's ability to learn new
875 classes is significantly enhanced.

A.3. Settings of the Proposed Method

Table A2. Settings of different datasets.

Dataset	L_g	L_s	h	ζ
VOC	24	8	2	0.7
COCO	40	16	4	0.7

877 During base training, the backbone and the query func-
878 tion are frozen, all prompts as well as the decoder are
879 trainable. We use AdamW as optimizer with $\beta_1 = 0.9$,

880 $\beta_2 = 0.9$, weight decay 0.01, and a polynomial learning
 881 rate policy with a linear learning rate warmup. The model
 882 is trained for 20k iterations on VOC and 80k iterations on
 883 COCO with a learning rate of 2×10^{-4} and batch size of
 884 8. During incremental training, we freeze all parameters
 885 but update the expanded stage-specific prompts and region-
 886 unique prompts. For both VOC and COCO, the model is
 887 trained for 400 iterations per step with a learning rate of
 888 2×10^{-4} without the learning rate warmup. We compute the
 889 results via single-scale full-resolution images without any
 890 post-processing. The settings of model hyper-parameters
 891 on different datasets are shown in Tab. A2.

892 Additionally, the pre-defined background classes are
 893 "sky", "wall", "tree", "wood", "grass",
 894 "road", "sea", "river", "mountain",
 895 "sands", "desk", "bed", "building",
 896 "cloud", "lamp", "door", "window",
 897 "wardrobe", "ceiling", "shelf",
 898 "curtain", "stair", "floor", "hill",
 899 "rail", "fence".

900 **A.4. Detail Data Flow and Model Structure**

901 (1) Prompts Generation: The images are input into pre-
 902 trained query function (DINOv2) to get global and local
 903 query features q_g and q_c which are used to match optimal
 904 prompts, formulated by Eq. (8), Eq. (9) and Eq. (13). (2)
 905 Image and Text Encoding: Image tokens concatenated with
 906 TP, SP and RP are input into each block of CLIP image en-
 907 coder (ViT-B), which outputs the class token g and pixel
 908 features P . The names of each BG and FG classes are input
 909 into CLIP text encoder to get text embeddings T_{bg} and T_{fg} .
 910 (3) Pixel Decoding: CMQE integrates the image global fea-
 911 ture g with text embeddings T_{fg} and T_{bg} to generate class-
 912 specific queries Q_{fg} and Q_{bg} as formulated by Eq. (3). FG
 913 isolation separator and BG refinement separator share the
 914 same structures, composed of three cross-attention blocks.
 915 For each block, pixel-wise features P are input as the keys
 916 and values. Q_{fg} and Q_{bg} are input as the queries of the first
 917 block and the output of each block is input as the queries
 918 of the next block. The last block outputs the final masks
 919 (Eq. (4)).

920 **A.5. Incremental Learning**

921 We elaborate on the ways in which the model is extended
 922 during the incremental stage as follow. During incremen-
 923 tal training stage, we need to expand three components: 1)
 924 text embeddings of novel classes, 2) slots of stage-specific
 925 prompts, 3) slots of region-unique prompts.

926 For expanding text embeddings, we just add novel
 927 classes to foreground text embeddings, which can be de-
 928 notes as $T_{fg}^t \in \mathbb{R}^{(C_{fg}^{t-1} + N_e) \times D}$, where N_e denotes the num-
 929 ber of novel classes, C_{fg}^{t-1} denotes the number of fore-
 930 ground classes of stage $t - 1$.

For expanding slots of stage-specific prompts, we con-
 931 catenate the expanded $\mathbf{A}_e^{(i)} \in \mathbb{R}^{1 \times D}$, $\mathbf{K}_e^{(i)} \in \mathbb{R}^{1 \times D}$, $\mathbf{S}_e^{(i)} \in$
 932 $\mathbb{R}^{1 \times L_s \times D}$ to matrix of stage $t - 1$. Thereby, the stage-
 933 specific prompts integration of stage t can be formulated
 934 as:
 935

$$\begin{aligned} \gamma &= \text{Softmax}(\langle q_g \odot [\mathbf{A}^{(i),t-1}; \mathbf{A}_e^{(i)}], \\ &[\mathbf{K}^{(i),t-1}; \mathbf{K}_e^{(i)}] \rangle > / \tau) \\ p_s^{(i),t} &= \gamma[\mathbf{S}^{(i),t-1}; \mathbf{S}_e^{(i)}] \in \mathbb{R}^{L_s \times D} \end{aligned} \quad (\text{A4})$$

where $\mathbf{A}^{(i),t-1}, \mathbf{K}^{(i),t-1}, \mathbf{S}^{(i),t-1}$ denotes the attention ma-
 936 trix, key matrix and stage-specific prompts of stage $t - 1$ for
 937 layer i , $p_s^{(i),t}$ denotes the integrated stage-specific prompts.
 938

For expanding slots of region-unique prompts, we con-
 939 catenate the expanded $\mathbf{K}_{l_e}^{(i)} \in \mathbb{R}^{N_e \times D}$, $\mathbf{L}_e^{(i)} \in \mathbb{R}^{N_e \times D}$ with
 940 matrix of stage $t - 1$. Thereby, the region-unique prompts
 941 integration of stage t can be formulated as:
 942

$$\begin{aligned} p_l^{(i),t} &= \text{Gumbel.Softmax}(q_v[\mathbf{K}_l^{(i),t-1}; \mathbf{K}_{l_e}^{(i)}]) \\ &[\mathbf{L}^{(i),t-1}; \mathbf{L}_e^{(i)}], \end{aligned} \quad (\text{A5})$$

where N_e denotes the number of novel classes,
 943 $\mathbf{K}_l^{(i),t-1}, \mathbf{L}^{(i),t-1}$ denote the all keys and prompts of
 944 stage $t - 1$, and $p_l^{(i),t}$ denotes the integrated region-unique
 945 prompts.
 946

Note that due to the adoption of a switching mechanism,
 947 the number of visual prompts input to the model remain
 948 consistent at different stages, effectively preventing a sig-
 949 nificant increase in computational complexity.
 950

951 **B. More Ablation Studies**

In this section, we first conduct experiments to demonstrate
 952 the effectiveness of MSVP compared to the baseline, and
 953 prove the effectiveness of orthogonality constraints. Sec-
 954 ondly, we perform an ablation study on the scale and types
 955 of query functions. Thirdly, we prove the computation ef-
 956 ficiency of MSVP. Finally, we carry out ablation experi-
 957 ments on the COCO dataset to further validate the proposed
 958 method.
 959

960 **B.1. Effectiveness of MSVP**

To further demonstrate the effectiveness of the proposed
 961 MSVP, we compare the performance under different novel
 962 shots between the baseline and MSVP. As shown in Fig. A1,
 963 the performance of the baseline model does not increase
 964 as much as the proposed model with the growth of the
 965 shot. This demonstrates that MSVP effectively boosts
 966 the model's capacity to learn novel classes by efficiently
 967 switching prompts. In contrast, the baseline model suffers
 968 from diluted information as more visual prompts are added,
 969 leading to inadequate learning of novel classes. Mean-
 970 while, the baseline's performance on base classes signifi-
 971 cantly outperforms that of the model with MSVP, and
 972

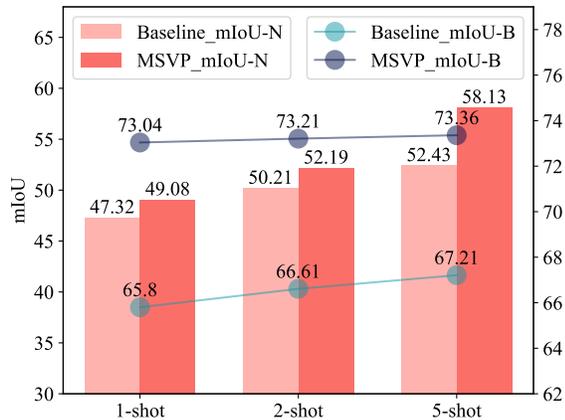


Figure A1. Performance comparison between the baseline and the proposed method on VOC under different shots.

Table A3. Ablation on \mathcal{L}_{orth} .

\mathcal{L}_{orth}	1-shot		
	Base	Novel	HM
wo	67.35	48.62	56.47
w	73.04	49.08	58.71

978 improves with more novel class samples. The proposed
 979 MSVP stores knowledge of different stages in independent
 980 visual prompts and switches them dynamically, thus avoid-
 981 ing knowledge of old stages corrupted.

982 B.2. Effectiveness of Orthogonality Constraints

983 We add orthogonality constraints to parameters of stage-
 984 specific prompts to avoid interference between existing and
 985 new knowledge and reduce catastrophic forgetting. In this
 986 section, we conduct an experiment on this loss to prove its
 987 effectiveness as in Tab. A3. The table shows that adding
 988 this loss function significantly enhances the performance
 989 of the base class by 5.69% mIoU, suggesting that it effec-
 990 tively minimizes interference from new knowledge on exist-
 991 ing knowledge and helps prevent forgetting of previously
 992 learned classes.

993 B.3. Ablation on Query Function

994 We choose pretrained ViT-B of DINOv2 as the query func-
 995 tion to produce high quality global query features and local
 996 query features. In this section, we conduct an ablation study
 997 on the scale of the query function, as shown in Tab. A4.
 998 The model’s performance improves significantly when the
 999 query function transitions from ViT-S to ViT-B. However,
 1000 further scaling from ViT-B to ViT-L results in minimal per-
 1001 formance gains, indicating that the larger model size does
 1002 not substantially enhance effectiveness in our method.

1003 We also perform experiments to evaluate different types
 1004 of query functions, including MAE[44], BEiT[45], and

Table A4. Ablation on the scale of DINOv2.

Model	1-shot		
	Base	Novel	HM
ViT-S	71.62	48.57	57.88
ViT-B	73.04	49.08	58.71
ViT-L	72.68	48.76	58.36

Table A5. Ablation on different query functions.

Pre-train method	1-shot		
	Base	Novel	HM
MAE	50.44	39.00	43.98
BEiT	42.95	32.69	37.12
DINOv2	73.04	49.08	58.71

DINOv2[15], as summarized in Tab. A5. For these exper-
 1005 iments, we use ViT-B with various pre-training methods.
 1006 The results show that DINOv2 achieves the best perfor-
 1007 mance, likely due to its ability to provide both global and
 1008 local features with high generalizability, thereby switch-
 1009 ing appropriate prompts accurately. In contrast, MAE and
 1010 BEiT yield relatively inferior results, which may stem from
 1011 their limitations in effectively capturing image-specific dif-
 1012 ferences across different stages. Consequently, they strug-
 1013 gle to generate tailored prompts that align with the unique
 1014 characteristics of images at various stages.
 1015

1016 B.4. Ablation of Different Prompts on COCO.

1017 We also carry out an experiment of different prompts on
 1018 COCO to prove the effectiveness of our method, as shown
 1019 in Tab. A6. When simply adding prompts as Eq. (2),
 1020 the baseline achieves 44.98% mIoU-B and 23.72% mIoU-
 1021 N, which has outperformed previous SOTA methods by
 1022 a large margin. Replacing the vanilla prompt expand-
 1023 ing strategy with the proposed task-persistent prompts and
 1024 stage-specific prompts, the performance increases by 3.28%
 1025 mIoU-B and 0.93% mIoU-N. That’s because task-persistent
 1026 prompts provide transferable knowledge across stages and
 1027 stage-specific prompts extract relevant knowledge from
 1028 other stages and enhance it with discriminative knowledge
 1029 of the current stage, which offers a flexible way to switch
 1030 knowledge of different stages, thereby achieving better abil-
 1031 ities to keep old knowledge and learn new classes. Fur-
 1032 thermore, with finer-grained region-unique prompts, per-
 1033 formance on novel classes further rises to 25.60%, for the
 1034 reason that region-unique prompts provide the model with
 1035 knowledge of local details of specific classes. The table also
 1036 shows that excluding task-persistent prompts leads to a per-
 1037 formance decrease on both base and novel classes. It proves
 1038 that general knowledge can not only help models maintain
 1039 old abilities but also assist models in learning new abilities.
 1040

Table A6. Ablation of Different Prompts on COCO.

SA	TP	SP	RP	1-shot		
				SA	Novel	HM
✓				44.98	23.72	31.06
	✓	✓		48.26	24.65	32.63
	✓		✓	41.85	24.78	31.13
		✓	✓	44.24	20.54	28.05
	✓	✓	✓	48.85	25.60	33.59

Table A7. Ablation of the Framework on COCO.

Method	1-shot		
	Base	Novel	HM
Ours	48.85	25.60	33.59
- FDD	48.18	23.13	31.25
- \mathcal{L}_{ba}	45.81	21.92	29.65
- CMQE	44.72	16.57	24.18

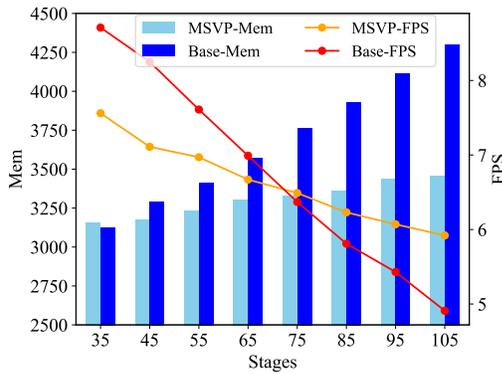


Figure A2. Comparison of MSVP and Base on FPS and Memory.

1040 **B.5. Ablation of the Framework on COCO.**

1041 We also carry out an experiment of the framework on
 1042 COCO to prove the effectiveness of our method, as shown
 1043 in Tab. A6. Replacing FDD with a single decoder results
 1044 in a decrease in performance on novel classes. That’s be-
 1045 cause FDD enables the model to process the salient fea-
 1046 tures of novel classes and adapt to background changes sepa-
 1047 rately, alleviating the confusion between novel classes and
 1048 the background. Additionally, jointly optimizing the two
 1049 separators with \mathcal{L}_{ba} leads to a better performance, which
 1050 mitigates the misaligned optimization directions caused by
 1051 \mathcal{L}_{van} . It can also be concluded from the table that CMQE
 1052 plays an important role in maintaining base knowledge and
 1053 learning novel classes by generating generalizable class em-
 1054 beddings for base classes and novel classes.

1055 **B.6. Computation Efficiency.**

1056 We conduct an experiment to compare FPS and memory use-
 1057 age between MSVP and Baseline as the incremental phase
 1058 increases, on an NVIDIA A40 GPU. MSVP exhibits signifi-

cant advantages in memory efficiency compared to the base-
 line. This improvement stems from our prompt-switching
 mechanism, which effectively maintains a constant input se-
 quence length to the transformer model throughout progres-
 sive training stages, thereby avoiding memory accumula-
 tion. While the baseline shows marginally higher FPS dur-
 ing early incremental stages (<75 stages) due to the intro-
 duced query operation of MSVP, our method demonstrates
 superior computational sustainability as training progresses.
 Notably, the baseline suffers a sharp FPS degradation as
 its growing prompt inventory quadratically increases trans-
 former’s computational complexity ($O(n^2)$). 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070

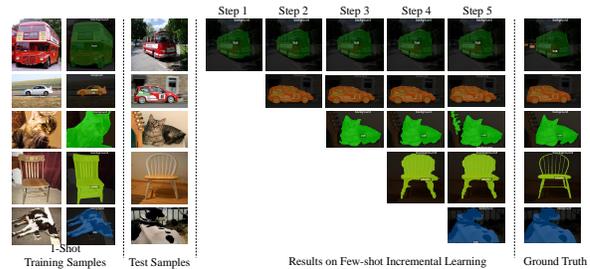


Figure A3. Step-by-step segmentation results of our method. Zoom in for better visualization.

C. Visualization 1071

In Fig. A3, we visualize our step-by-step segmentation re-
 sults for novel classes. The figure shows that our method
 effectively retains old class knowledge, enabling the model,
 even after multiple training rounds, to correctly predict old
 class samples. Additionally, our approach demonstrates
 strong novel class learning capability, as it can generalize
 to other test samples by learning from just one novel class
 sample. The visualization results effectively demonstrate
 the validity of the proposed prompt-based IFSS method. 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079
 1080