

000

001

002

003

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

GAPS: Few-Shot Incremental Semantic Segmentation via Guided Copy-Paste Synthesis

Anonymous CVPR L3D-IVU submission

Paper ID 10

Abstract

Few-shot incremental segmentation is the task of updating a segmentation model, as novel classes are introduced online over time with a small number of training images. Although incremental segmentation methods exist in the literature, they tend to fall short in the few-shot regime and when given partially-annotated training images, where only the novel class is segmented. This paper proposes a data synthesizer, Guided copy-And-Paste Synthesis (GAPS), that improves the performance of few-shot incremental segmentation in a model-agnostic fashion. Despite the great success of copy-paste synthesis in conventional offline visual recognition, we demonstrate substantially degraded performance of its naïve extension in our online scenario, due to newly encountered challenges. To this end, GAPS (i) addresses the partial-annotation problem by leveraging copy-paste to generate fully-labeled data for training, (ii) helps augment the few images of novel objects by introducing a guided sampling process, and (iii) mitigates catastrophic forgetting by employing a diverse memory-replay buffer. Compared to existing state-of-the-art methods, GAPS dramatically boosts the novel IoU of baseline methods on established few-shot incremental segmentation benchmarks by up to 80%. More notably, GAPS maintains good performance in even more impoverished annotation settings, where only single instances of novel objects are annotated.

1. Introduction

Incremental segmentation is an important capability for open-world AI systems. For example, consider a housekeeping robot that has been trained to segment common household objects, but once deployed in a user’s home it encounters a previously unseen type of furniture. For such practical applications, incremental segmentation would be capable of expanding the set of recognized classes to contain the new object. There are a few desired properties of incremental segmentation algorithms to operate under these

scenarios. First of all, the algorithm should be equipped with **few-shot learning capability**, which means that the algorithm can benefit from as few as one image provided by a user rather than requiring hundreds of images annotated offline by professional annotators. Second, providing full segmentation annotation of an image is time-consuming. To avoid causing substantial burdens for untrained users, the algorithm needs to be trainable with **partially-annotated** images where only novel classes are annotated.

A few attempts have been made by recent works [3, 5, 7, 28, 30] on *non-few-shot* incremental segmentation to investigate learning with partially-annotated images, which is termed *semantic background shift* [3]. Background shift describes a challenge unique to incremental semantic segmentation where classes that are not in the current learning step are assigned ‘background’ labels, which prohibits direct end-to-end training. Recent work uses either modified loss [3, 30] or pseudo-labeling [5, 7, 28] as *proxies* to train on partially-annotated images. However, although these proxying methods demonstrate good performance under the non-few-shot settings, they rely on rich annotations and fall short when only a limited amount of data is presented to the model, due to a lack of diversity of data. An even more restrictive setting occurs when users label only a single instance of the novel class, which can dramatically hurt performance of proxy models, due to the training containing non-annotated instances of the novel class (which are treated as negative pixels).

To address the aforementioned challenges, we propose GAPS (Guided copy-And-Paste Synthesis), which improves the training of incremental segmentation models by synthesizing fully-annotated images from partially-annotated examples. It is *model-agnostic*, and can be inserted as a plug-and-play module into different incremental learning algorithms, e.g., standard fine-tuning or PIFS [4]. Copy-paste generates diverse training data to boost performance under few-shot settings, enables the model to learn with partially-annotated images with as few as one annotated novel instance out of many novel instances in an image (e.g., as illustrated at the lower left part of Fig. 1), which

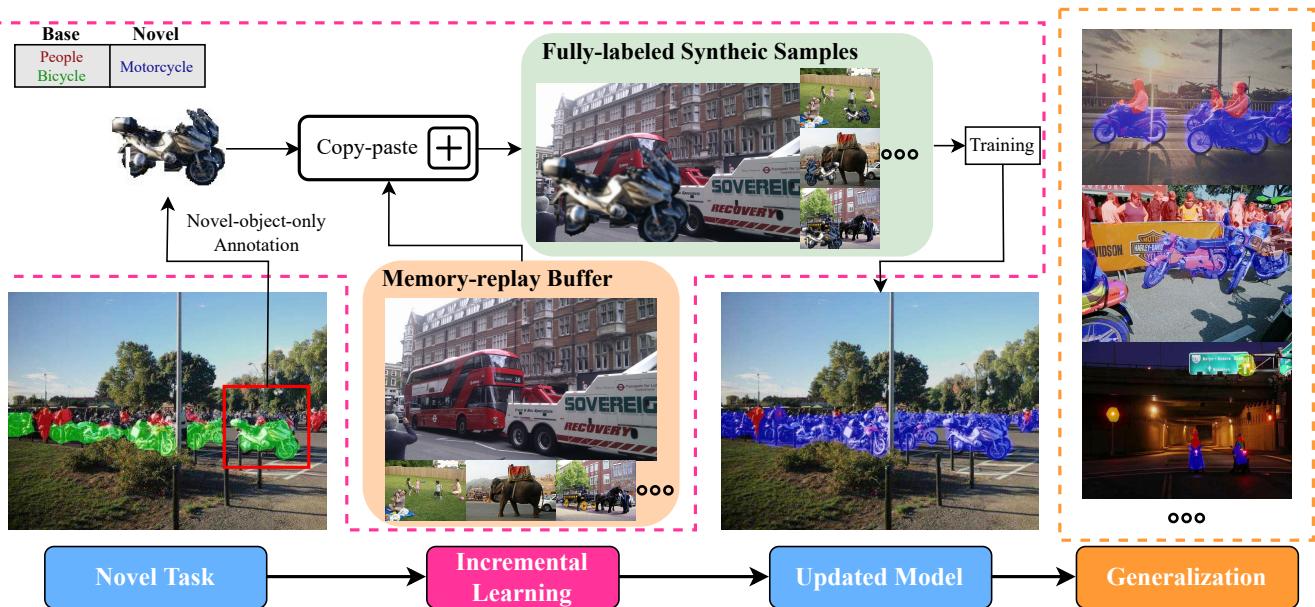
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127

Figure 1. Our proposed method utilizes guided copy-paste augmentation to synthesize diverse training data, using as few as **one single novel instance** for training. For example, the model encounters an image of many **motorcycles**, which is novel to the model. As a result, the model incorrectly assigns learned **bicycle** labels to these pixels and therefore needs to be updated. Our proposed method can adapt to the novel motorcycle class with an annotation of a single motorcycle, which can be efficiently annotated; whereas previous work [3, 4] require time-consuming annotation of all instances of motorcycles or even the entire image. Best view in color.

is a *stricter* setting than semantic background shift [3].

To the best of our knowledge, we are the *first* to introduce copy-paste as a synthesis technique to create a diverse data source for few-shot incremental segmentation. Although copy-paste [12] has been shown to be an effective data augmentation technique for offline visual recognition tasks, we identify new key technical challenges to adapting it to few-shot incremental settings. First, how should the synthesizer pick representative samples from the base dataset to construct a *diverse* pool of fully-annotated base scenes? Second, given the constructed pool of fully-annotated images, how should it select the most *suitable* base images to be pasted on? Third, after an informative image is selected, from what distribution should it sample current and previously learned novel objects to *balance* sample frequency and avoid over-sampling or under-sampling? Our GAPS method differs from a naïve (e.g., uniform random sampling) copy-paste process by a *guided* strategy that considers diversity of the memory-replay buffer, imbalanced class frequencies between base classes and novel classes, and contextual similarity of images.

In summary, our contributions are as follow:

1. We are the first to introduce copy-paste as a synthesis technique to address partially-labeled images for incremental segmentation.
2. To address the gaps between copy-paste under the of-

fine setting as an augmentation technique and under the online setting as a synthesis technique, we design a guided copy-paste process that improves the distribution of synthesized images by enforcing diversity of the memory-replay buffer, exploiting contextual information, and balancing class frequencies.

3. The proposed GAPS technique consistently boosts the performance of a variety of incremental learning algorithms from simple fine-tuning to sophisticated state-of-the-arts under the few-shot setting. Furthermore, we demonstrate the strength of GAPS to cope with a more challenging task setting where only one instance out of many novel instances in an image is annotated, which highlights copy-paste as a better alternative to pseudo-labeling or modified loss for practical incremental segmentation applications.

2. Related Work

Incremental Learning for Semantic Segmentation. It is known that many learning-based models suffer from catastrophic forgetting [19], a phenomenon that causes models to perform significantly worse on old tasks when they are fine-tuned to adapt to new tasks. *Incremental learning* studies how to enable models to adapt to new classes while mitigating catastrophic forgetting without accessing the old dataset or full-scale re-training. This problem has

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

been studied extensively in image classification [1, 2, 16–18, 22, 29]; whilst relatively less work have been done to study incremental learning under the task setting of semantic segmentation [3, 5, 7, 20, 30]. Noticeably, a few attempts have been made by recent work to address the semantic background shift problem proposed by [3] via either pseudo-labeling [5, 7, 28] or modified loss [3, 30] to train on partially-annotated images of novel classes. However, existing work relies on rich annotations and tends to fail when only a limited amount of data is available. In contrast, our work enables incremental segmentation learning with *few data* via a guided copy-paste process, which demonstrates promising performance under the few-shot and more impoverished single-instance setting. Furthermore, GAPS is a *model-agnostic data pre-processor*, which is orthogonal to incremental learning techniques such as regularization [16].

Few-Shot Semantic Segmentation. *Few-shot semantic segmentation* methods predict segmentation masks of novel classes using only a few training examples of the novel class. Many meta-learning-based methods [23, 26, 27, 31] and even specialized datasets [15] have been proposed to address such a problem. However, few-shot semantic segmentation methods produce novel-class-only binary foreground-background segmentation. In comparison, our proposed method works in a more challenging and realistic setting where both base classes and novel classes need to be segmented.

Few-Shot Incremental Segmentation. While there are many works in few-shot incremental image classification [6, 24], relatively fewer works have been done to investigate few-shot incremental segmentation [4, 11, 25]. [25] designs a meta-learning-based classifier that adjusts learned prototypes by modeling interaction between base classes and incoming novel class. Unlike [25], which only performs a single update of weights in the classifier, PIFS [4] apply regularization techniques to allow fine-tuning of the entire network, achieving state-of-the-art result in few-shot incremental semantic segmentation. However, PIFS [4] is fine-tuned on only a small number of samples, which leads to sub-optimal performance due to overfitting. In addition, PIFS requires fully-annotated images as input, which hinders its potential for practical applications.

Copy-Paste Augmentation. Copy-and-paste is an augmentation technique that copies a subset of objects from one image and pastes onto the other image using their segmentation masks. Many works [8, 9, 12] have been done to investigate how copy-and-paste augmentation can help with various visual tasks. Dvornik *et al.* apply copy-paste augmentation in object detection by designing a neural network to consider context and guide copy-paste. However, the context guidance method proposed by Dvornik *et al.* can not be trivially applied to our application since it requires abundant fully-annotated training data. More recently, Ghiasi *et*

al. conduct extensive experiments to demonstrate the effectiveness of simple copy-paste in the instance segmentation problem. We extend the augmentation strategy from [12] and construct an intuitive baseline called Naïve copy-Paste Synthesis (NPS) to adapt it to our online task setting. However, as we will demonstrate in the ablation study, such naïve adaptation gives unsatisfactory performance in our task setting because of *gaps* between the static offline learning and continual online learning. In our work, we propose a series of techniques to guide the copy-paste synthesizer to address these gaps, whose effectiveness is evident from the significant improvement from NPS.

3. Method

Problem Setup. Let $\mathcal{X} \subset \mathbb{R}^{H \times W \times 3}$ be a set of RGB images with size $H \times W$, $\mathcal{C} \subset \mathbb{N}$ be a set of category labels, and $\mathcal{Y}^{\mathcal{C}} \subset \mathbb{R}^{H \times W \times |\mathcal{C}|}$ be a set of label masks (*i.e.*, per-pixel category labels in \mathcal{C}). In semantic segmentation, we aim to learn a model ϕ that maps an image $x \in \mathcal{X}$ to a segmentation mask $y \in \mathcal{Y}^{\mathcal{C}}$. Different from standard semantic segmentation, in few-shot incremental segmentation, \mathcal{C} is expanded over time through two stages. During the *base learning stage*, the model is provided with a base dataset $\mathcal{D}_0 = \{(x_i, y_i) | x_i \in \mathcal{X}, y_i \in \mathcal{Y}^{\mathcal{C}_0}\}$, where \mathcal{C}_0 is a set of classes in the base dataset. \mathcal{D}_0 generally contains many fully-annotated image-mask pairs and is used to train the model $\phi_0 : \mathcal{X} \rightarrow \mathcal{Y}^{\mathcal{C}_0}$ from scratch.

During the *incremental learning stage*, a sequence of tasks $\{D_1, D_2, \dots\}$ with novel categories is presented to the model, where $D_j = \{(x_i, y_i) | x_i \in \mathcal{X}, y_i \in \mathcal{Y}^{\mathcal{C}_j}\}$ and \mathcal{C}_j is a set of classes for task D_j . In few-shot learning, the size of the training sets for the novel tasks is small, *i.e.*, $|D_j| \ll |\mathcal{D}_0|$. After adapting to task D_j , the model is updated as $\phi_j : \mathcal{X} \rightarrow \mathcal{Y}^{\cup_{i=0, \dots, j} \mathcal{C}_i}$. The goal of incremental learning is to optimize the model performance jointly on both previous tasks and the current task. To enforce the partially-annotated image setting, we follow Cermelli *et al.* and assume that only novel classes are annotated, *i.e.*, $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset$ for all $i \neq j$.

Method overview. Fig. 2 illustrates our proposed Guided copy-Paste Synthesis (GAPS) framework for few-shot incremental segmentation. It is a *generic and model-agnostic* data synthesis framework that generates fully-labeled scenes from partially-annotated images of novel objects as a preprocessor to the underlying segmentation model. After the standard base learning stage with base dataset \mathcal{D}_0 , we build a memory-replay buffer $\hat{\mathcal{D}}_0$ using an *diversity-guided exemplar selection strategy* (Section 3.2). During the incremental learning stage, fully-labeled samples are synthesized by copying from the masked novel objects in D_1, \dots, D_j and pasting onto base exemplars from the replay buffer $\hat{\mathcal{D}}_0$. The strategy by which we choose base exemplars and novel segments is *context-guided* (Sec-

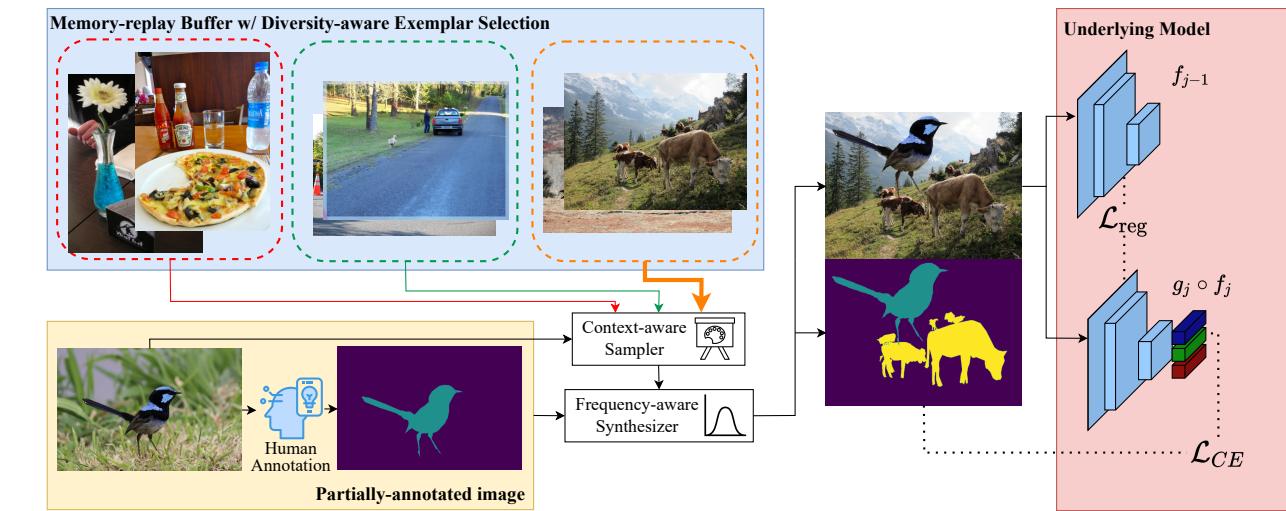


Figure 2. Overview of GAPS. During the incremental learning stage, GAPS takes in as few as one annotated instance of a single image. It is more probable for GAPS to select a scene contextually similar to the provided image from memory-replay buffer \hat{D}_0 . The image is then probabilistically pasted to generate synthetic fully-labeled scenes. Note that GAPS is model-agnostic, and here we use PIFS [4] as an example for the underlying segmentation model to illustrate how GAPS is applied as a pre-processor. Best seen in color.

tion 3.3) and *class-frequency-guided* (Section 3.4).

3.1. Few-Shot Incremental Segmentation Model

In principle, GAPS is model-agnostic, which means that it can work with many incremental segmentation models as a diverse data source to improve their performance. Here we adopt PIFS [4] as the main baseline underlying segmentation model for its state-of-the-art performance on few-shot incremental semantic segmentation and support for end-to-end training. The PIFS segmentation model ϕ is composed of a convolution-based feature extractor f and a per-pixel classification layer g using prototypical representation – g is configured to classify the pixels into n classes, so it is parameterized with prototypes $W = [w_1, w_2, \dots, w_n]$. Intuitively, f maps every pixel in an input image onto the unit hyper-sphere in a high-dimensional representation space. g then generates probability prediction by comparing cosine similarity of feature vectors with learned class prototypes w_i in the representation space and applying softmax of the resulting similarities. Following PIFS, given a previously unseen class $n + 1$ from task $D_{n+1} = \{(x_i, y_i) | x_i \in \mathcal{X}, y_i \in \mathcal{Y}^{C_{n+1}}\}$, instead of randomly initializing the prototype w_{n+1} , we apply the MAP function to estimate the prototype,

$$w_{n+1} = \frac{1}{|D_{n+1}|} \text{MAP}(x_i, y_i, n + 1), i \in D_{n+1} \quad (1)$$

$$= \frac{1}{|D_{n+1}|} \sum_{(x_i, y_i) \in D_{n+1}} \frac{M_{n+1}(y_i, j) \frac{f_j^i(x_i)}{\|f_j^i(x_i)\|}}{M_{n+1}(y_i, j)} \quad (2)$$

where $M_{n+1}(y_i, j)$ is a binary function that returns 1 if

the j -th pixel in mask y_i is class $n + 1$, and 0 otherwise. $f_n^i(x_i)$ denotes the feature vector at the j -th pixel of $f_n(x_i)$.

In addition, we want to note that our re-implementation of PIFS [4] uses \mathcal{L}_2 regularization rather than the prototype distillation loss proposed by Cermelli *et al.* We found experimentally that when a diverse data source is used (i.e., our proposed GAPS), \mathcal{L}_2 regularization works better.

To be more precise, we construct a penalization term \mathcal{L}_{REG} to regularize the output before the classifier. For incremental learning task D_j with image-mask pairs (x, y) , we have

$$\mathcal{L}_{REG} = \|f_j(x) - f_{j-1}(x)\|_2. \quad (3)$$

The final training loss is given by

$$\mathcal{L}(x, y) = \mathcal{L}_{CE}(\phi_j(x), y) + \lambda \mathcal{L}_{REG}, \quad (4)$$

where \mathcal{L}_{CE} is either the standard cross-entropy loss or the modified cross-entropy loss from [3]. λ is a hyper-parameter used to weight the regularization loss. All other components are the same as in [4]. We denote our re-implementation of PIFS with \mathcal{L}_2 regularization loss as PIFS(\mathcal{L}_2).

3.2. Diversity-guided Exemplar Selection with Learned Prototypes

For methods with memory-replaying (e.g., SSUL [5]), GAPS can work directly on top of their constructed buffers with minimal modification. For other methods such as PIFS [4], we propose a diversity-guided exemplar selection process that builds a small yet diverse memory-replay buffer \hat{D}_0 from D_0 to mitigate catastrophic forgetting. Selecting

432 diverse examples that are representative of the base dataset
 433 helps mitigate catastrophic forgetting, as suggested by [22].
 434 Inspired by Bang et al. [1], we select samples distributed
 435 uniformly along a spectrum from easy to hard for diversity.
 436

437 Here, we present an algorithm (Algorithm 1) to construct
 438 \hat{D}_0 by exploiting the Masked Average Pooling (MAP) function
 439 from [4]. Intuitively, we approximate the difficulty
 440 of every sample by their similarity between the estimated
 441 prototype with learned prototypes. Estimated prototypes
 442 that are close to the learned prototype are considered easy
 443 samples and vice versa. After building a list of base sam-
 444 ples sorted by difficulties, we select samples from equally-
 445 spaced intervals to ensure samples of all difficulties are se-
 446 lected for diversity.

447 During the incremental learning stage, we select at most
 448 k samples for each novel class using the same algorithm
 449 to memorize novel classes. To maintain the size of the
 450 memory-replay buffer, we remove old samples from the
 451 memory-replay buffer but keep at least 80% of the samples
 452 to be fully-annotated samples, so that we have diverse base
 453 images for copy-pasting.

454 3.3. Context-guided Sampling

455 We hypothesize that synthesizing novel objects onto con-
 456 textually consistent base images would result in an im-
 457 proved learning process. For example, a TV is more likely
 458 to appear in an apartment rather than in the middle of traf-
 459 fic on streets, and thus a TV object should more likely be
 460 pasted onto an image of another apartment rather than an
 461 outdoor landscape. To guide the copy-pasting process, we
 462 design a context-guided sampling algorithm to select im-
 463 ages from \hat{D}_0 that are contextually similar to the provided
 464 partially-labeled images.

465 One way of estimating pairwise contextual similarities
 466 between two images is to design a mapping $h : \mathcal{X} \rightarrow \mathbb{R}^m$
 467 that maps an image into a metric space, where the metrics
 468 serve as a proxy of the contextual similarity between two
 469 images. Here, in GAPS, we extract the knowledge of the
 470 learned feature extractor. In incremental learning task D_i ,
 471 scene embedding of image i is estimated by,

$$472 h_i = \frac{\text{GAP}(f_{i-1}(h_i))}{\|\text{GAP}(f_{i-1}(h_i))\|_2} \quad (5)$$

473 where $\text{GAP} : \mathbb{R}^{B \times CHW} \rightarrow \mathbb{R}^{B \times C}$ is the commonly
 474 used global average pooling function. Two scene embed-
 475 ding vectors h_i and h_j can then be compared by cosine sim-
 476 ilarity.

477 To find contextually similar base images to each novel
 478 image, we evaluate the cosine similarity of the novel image
 479 to each of the examples in \hat{D}_0 , and construct a contextu-
 480 ally similar subset \mathbb{S} with $|\hat{D}_0|/10$ most contextually sim-
 481 ilar examples (If interested, visualization of query samples
 482 and their contextually-similar counterparts are included in
 483 Fig. 4 in the supplementary material). When there are mul-
 484 tiple novel images, we take a union of selected examples.
 485 To allow other base scenes to be sampled to mitigate catas-
 486 trophic forgetting, we sample from \mathbb{S} with a probability of
 487 α , and sample from \hat{D}_0 with a probability of $1 - \alpha$, where α
 488 is a hyperparameter set to 0.9 in our implementation. Note
 489 that we only need to compute scene embedding once for ev-
 490 ery image in \hat{D}_0 and incoming partially-annotated images.
 491 Hence, the context-guided sampling algorithm poses only
 492 minor computational overhead to GAPS.

493 3.4. Class-frequency-guided Probabilistic Synthesis

494 Now the final question is, given a fully-annotated im-
 495 age x_B and an image of a novel object x_N , how frequent
 496 should we apply copy-paste? There is a trade-off between
 497 oversampling and undersampling. As one extreme, one can
 498 follow [12] and always apply copy-paste augmentation to
 499 paste novel objects onto every base image. However, this
 500 will lead to oversampling of novel categories in the current
 501 task, which we found to hurt the performance of existing
 502 classes. On the other hand, rarely pasting novel instances
 503 would lead to undersampling of the novel class. Therefore,
 504 to guide copy-paste in the online setting, we design a syn-
 505 thesis strategy called vRFS based on RFS (Repeat Factor
 506 Sampling) described by [13] to perform synthesis.

507 To apply vRFS, we first need to compute category-wise
 508 sampling factor r_c for every c as in RFS. If $c \in \mathcal{C}_0$, we
 509 set $r_c = 1$ as since during the construction of \hat{D}_0 we al-
 510 ready consider class balance by class-wise uniformly sam-
 511 pling. If $c \in \mathcal{C}_j$ with $j \geq 1$, we first compute its class
 512 frequency by $f_c = n_{\text{Shot}}/|\hat{D}_0|$, where n_{Shot} denotes

496 Algorithm 1 Construct Memory-replay Buffer

```
497 Require: number of exemplars  $n$ 
498    $k \leftarrow \text{FLOOR}(n/|Y_0|)$  // Sample per class
499   for  $c$  from 1 to  $|Y_0|$  do
500      $S_c \leftarrow \{(x_i, y_i) \in D_0, c \in y_i\}$ 
501     for  $(x_i, y_i) \in S_c$  do
502        $p_i \leftarrow \text{MAP}(x_i, y_i, c)$  // Pred. Proto.
503        $s_i \leftarrow \text{COSINESIMILARITY}(p_{ic}, w_c)$ 
504     end for
505     Sort  $S_c$  by similarity score  $s_i$ 
506      $ES_c \leftarrow \{\}$  // final exemplar set of class  $c$ 
507     for  $j = 1, 2, \dots, k$  do
508        $L_{idx} \leftarrow j \cdot |S_c|/k$ 
509        $U_{idx} \leftarrow \text{MIN}(\text{Lower} + |S_c|/k, |S_c|)$ 
510        $(x, y) \leftarrow \text{SAMPLE}(S_c[L_{idx} : U_{idx}])$ 
511        $ES_c \leftarrow ES_c \cup (x, y)$ 
512     end for
513   end for
514    $\hat{D}_0 \leftarrow \text{UNIFORMSAMPLE}(\bigcup_{i=1, \dots, |Y_0|} ES_i, n)$ 
```

540 the number of images in D_j with at least one pixel of c .
 541 Then, the category-wise sampling factor for c is given by
 542 $r_c = \text{MAX}(1, \sqrt{t/f_c})$. Note that in [13], t is chosen as a hyperparameter to be tuned.
 543 However, we empirically found that setting t to be the multiplicative inverse of total number
 544 of classes, or $t = 1/|\cup_{1,\dots,j} Y_j|$, is enough to yield stable results across different datasets and under different few-shot settings.
 545 This eliminates the need to search a hyperparameter for different settings and make our proposed method
 546 more robust towards different task settings.
 547

548 During the synthesis process, we first randomly select a novel class c_N from C_j , and another class c_o from
 549 $\cup_{1,\dots,j} C_j \setminus \{c_N\}$. We first decide if c_o should be pasted onto x_B . To apply vRFS resampling, we hallucinate two
 550 *virtual samples*: in the first sample where copy-paste would not be applied, the image-level sampling factor is given by
 551 1. In the second sample where copy-paste synthesis were
 552 to be performed, we would obtain a sample with image-level sampling factor of $r_i = \text{MAX}_{c \in i} r_c = r_{c_o}$. Thus,
 553 the probability to synthesize class c_o onto x_B is given by
 554 $r_{c_o}/(1 + r_{c_o})$. We then repeat the process for the novel
 555 class c_N . Note that vRFS synthesis is applied twice for every class, resulting in up to two pasted instances of c_N in
 556 the final image.
 557

558 4. Experiments

559 4.1. Datasets

560 We follow literature in few-shot segmentation and few-shot incremental segmentation [4, 21, 23, 25] and evaluate
 561 our model on the PASCAL-5ⁱ dataset [23] and the COCO-20ⁱ dataset [21]. PASCAL-5ⁱ is artificially built from the
 562 PASCAL VOC 2012 Semantic Segmentation dataset [10] with additional annotations from the SBD [14] dataset.
 563 The original VOC segmentation dataset provides segmentation annotations for 20 object categories. The PASCAL-
 564 5ⁱ dataset manually splits the original dataset into 4 folds for cross-validation. For each fold, 5 categories are selected
 565 as novel categories, while the remaining 15 categories are regarded as base categories. In our experiments, images
 566 containing at least one pixel of the novel categories are excluded from the base dataset. The construction of the
 567 COCO-20ⁱ dataset handles the 80 thing classes in COCO in a similar manner, where the dataset is split into 4 folds
 568 and each fold contains 20 categories. The rest of the process to construct the base dataset and the novel dataset in
 569 COCO-20ⁱ is same as the PASCAL-5ⁱ dataset.
 570

580 4.2. Evaluation Protocols

581 In the base learning stage, the model is trained using
 582 the entire base dataset. In incremental learning stages, sequences of tasks are presented to the model. We use the
 583 same evaluation protocol as proposed in [4] for fair com-

584 parisons, where 5 incremental learning tasks are used for
 585 PASCAL-5ⁱ and each task contains 1 class from the novel
 586 split. On the COCO-20ⁱ dataset, there are 4 incremental
 587 learning tasks, and each task contains 5 classes from the
 588 novel split.
 589

590 We evaluate the performance of the model on the entire validation set of the corresponding dataset after every
 591 step. For fair comparisons with our main baseline PIFS [4], we average results across different steps and exclude completely
 592 unseen classes from evaluation of current step. We use three different metrics to evaluate the performance of
 593 the model: mean Intersection-over-Union (mIoU) over base
 594 categories, mIoU over novel categories, and harmonic mean
 595 of the base mIoU and the novel mIoU. Unless otherwise
 596 noted, the numbers are computed by averaging results over
 597 splits in a cross-validating fashion.
 598

599 To average out randomness due to few training samples,
 600 we also average results over multiple runs with different set
 601 of few-shot training samples. For experiments on splits on
 602 PASCAL-5ⁱ, we found that averaging results over 10 runs
 603 with randomly sampled few-shot novel images yields stable
 604 results. For COCO-20ⁱ, we found that averaging results
 605 over 5 runs is enough to yield stable results.
 606

607 4.3. Main Results

608 In Table 1, we evaluate various incremental segmentation
 609 methods on the PASCAL-5ⁱ dataset and the COCO-20ⁱ
 610 dataset, and combine them with GAPS where appropriate.
 611

612 **Baselines.** There are two main baselines we are
 613 comparing to. The first one is SSUL [5], which is the state-
 614 of-the-art method in non-few-shot incremental segmentation.
 615 The second one is PIFS [4], for it is the state-of-
 616 the-art method in few-shot incremental semantic segmen-
 617 tation. We also report the performance of two variants of
 618 PIFS: one is our re-implementation PIFS(\mathcal{L}_2) described
 619 in Sec. 3.1. The other variant is PIFS(\mathcal{L}_2)+MEM, which
 620 uses a memory-replay buffer of the same size of GAPS. To
 621 handle partially-labeled images, we follow SSUL and per-
 622 form pseudo-labeling on partially annotated samples before
 623 adding them to memory for PIFS(\mathcal{L}_2)+MEM. In addition,
 624 we also evaluate simple fine-tuning and MiB [3].
 625

626 **GAPS consistently increases performance under few-
 627 shot settings.** Methods combined with our proposed data
 628 source, GAPS, consistently outperform their un-augmented
 629 counterpart on both the base and novel categories’ perfor-
 630 mance. It is worth noting that GAPS substantially boosts
 631 the performance of methods that originally require fully-
 632 annotated training images (i.e., fine-tuning and PIFS), de-
 633 spite using only partially-annotated images now. Even for
 634 methods that do not carry out end-to-end training and up-
 635 date only the classifier (i.e., SSUL), GAPS still steadily in-
 636 creases performance on novel categories. Compared to our
 637 implemented variant PIFS(\mathcal{L}_2)+MEM, our method demon-
 638 strates better performance on novel categories. In addition,
 639 GAPS consistently outperforms PIFS(\mathcal{L}_2)+MEM on novel
 640 categories. This indicates that GAPS is able to learn more
 641 information from partially-annotated images and transfer
 642 them to novel categories. The performance of GAPS is
 643 comparable to that of SSUL on novel categories, which
 644 indicates that GAPS is able to learn from partially-annotated
 645 images and achieve comparable performance to SSUL.
 646

648	METHOD	BASE	NOVEL	HM	BASE	NOVEL	HM	702
								703
PASCAL-5 ⁱ	1-SHOT	PASCAL-5 ⁱ	5-SHOT					
MIB [3]	43.9	2.6	4.9	60.9	5.8	10.5		704
FINETUNE*	47.2	3.9	7.2	58.7	7.7	13.6		705
FINETUNE+GAPS	64.2(+17.0)	16.2(+12.3)	25.9(+18.7)	66.8(+8.1)	38.1(+30.4)	48.5(+34.9)		706
SSUL [5]	73.9	16.4	26.8	74.8	27.8	40.5		707
SSUL+GAPS	74.0(+0.1)	19.9(+3.5)	31.3(+4.5)	74.9(+0.1)	30.0(+2.2)	42.8(+2.3)		708
PIFS* [4]	64.1	16.9	26.7	64.5	27.5	38.6		709
PIFS(\mathcal{L}_2) ^{*1}	64.6	19.7	30.2	57.7	24.5	34.4		710
PIFS(\mathcal{L}_2)+MEM	68.1	17.4	27.8	69.3	39.7	50.5		711
PIFS(\mathcal{L}_2)+GAPS	66.8(+2.2)	23.6(+3.9)	34.9(+4.7)	68.2(+10.5)	43.9(+19.4)	53.4(+19.0)		712
COCO-20 ⁱ	1-SHOT	COCO-20 ⁱ	5-SHOT					713
MIB [3]	40.4	3.1	5.8	43.8	11.5	18.2		714
FINETUNE*	38.5	4.8	8.5	39.5	11.5	17.8		715
FINETUNE+GAPS	44.5(+6.0)	11.0(+6.2)	17.7(+9.5)	46.4(+6.9)	24.9(+13.4)	32.4(+14.6)		716
SSUL [5]	51.0	6.3	11.3	51.6	15.0	23.2		717
SSUL+GAPS	50.8(-0.2)	11.0(+4.7)	18.1(+6.8)	51.9(+0.3)	17.1(+2.1)	25.7(+2.5)		718
PIFS* [4]	40.4	10.4	16.5	41.1	18.3	25.3		719
PIFS(\mathcal{L}_2) ^{*1}	45.7	10.3	16.8	46.2	20.2	28.1		720
PIFS(\mathcal{L}_2)+MEM	47.8	11.2	18.1	46.8	22.0	29.9		721
PIFS(\mathcal{L}_2)+GAPS	46.8(+1.1)	12.7(+2.4)	20.0(+3.2)	49.1(+2.9)	25.8(+5.6)	33.8(+5.7)		722

Table 1. Methods augmented with our proposed GAPS consistently outperform their un-augmented counterparts in terms of IoU across different few-shot settings on COCO-20ⁱ and PASCAL-5ⁱ. Methods noted with * are privileged and use fully-annotated images, others use images with novel-class-only partial annotation. ¹: our re-implementation using \mathcal{L}_2 regularization. Highest results are colored red and the second highest results are colored blue. HM stands for harmonic mean. (Best view in color).

strates considerable relative improvement on novel categories but performs slightly worse on base categories due to the introduction of pseudo-labeling in PIFS(\mathcal{L}_2)+MEM, which has been shown in previous work [5] to have regularization effects on base classes.

4.4. Ablation Study

In Table 2, we ablate guidance designs in GAPS to illustrate how different types of guidance contribute to the final incremental learning performance than naïve copy-paste synthesis. Though GAPS is a synthesis method that applies to many base learning algorithms, due to the highest harmonic mean of PIFS(\mathcal{L}_2)+GAPS on all settings, here we use PIFS(\mathcal{L}_2)+GAPS for the ablation study.

Our diversity-guided exemplar selection method consistently increases performance on base categories, which suggests that it is capable of choosing diverse samples to construct a representative memory-replay buffer and mitigate catastrophic forgetting to improve performance on base classes after sequential adaptations.

Context-guided sampling steadily improves performance on novel classes, which is consistent with findings in previous work [8] that background context is an important factor to consider in copy-paste synthesis. (Example visualization is available in Fig. 4 in the supplementary material).

Frequency-guided probabilistic synthesis boosts results on novel classes. On the other hand, it does not influence the performance of base categories in a statistically significant manner. We take a closer look at step-wise performance and found that the reason is due to unguided copy-paste’s oversampling of novel classes that are being adapted, and forgetting of classes learned in the previous incremental learning stage and not in the memory-replay buffer.

4.5. More Challenging Single-Instance Experiment

Though the semantic background shift proposed by [3] relaxes the requirement to provide full segmentation annotations, it still requires *all novel instances* in images to be annotated, which can be time-consuming to obtain in cluttered scenes and hinder potential applications (e.g., the cluttered motorcycle image in Figure 1). Here we consider a more challenging task setting, which we term *single-instance incremental learning*. Namely, for training images provided in incremental learning stages, if there are multiple instances of a novel class in the image, we assume that only *one instance* will be annotated for the model.

To simulate this setting, we use the instance-level segmentation annotation provided by the COCO dataset to enforce only annotation of one novel instance in every image is available to the model. Since state-of-the-art incre-

756	MEM	COPY-PASTE	F-GUIDE	D-GUIDE	C-GUIDE	BASE	NOVEL	HM	810
757	—	—*	—	—	—	46.2(± 0.3)	20.2(± 0.7)	28.1(± 0.3)	811
758	✓	—*	—	—	—	49.3(± 0.2)	19.4(± 0.7)	27.9(± 0.3)	812
759	✓	✓	—	—	—	47.0(± 0.2)	19.8(± 0.6)	27.8(± 0.3)	813
760	✓	✓	✓	—	—	47.2(± 0.2)	25.2(± 0.6)	32.9(± 0.3)	814
761	✓	✓	✓	✓	—	48.2(± 0.2)	25.0(± 0.7)	32.9(± 0.3)	815
762	✓	✓	✓	✓	✓	49.1 (± 0.2)	25.8 (± 0.6)	33.8 (± 0.3)	816
763									817

Table 2. Ablation study of components in GAPS on PIFS(\mathcal{L}_2) on the COCO-20ⁱ dataset under 5-shot setting. Note that when only combined with the memory-replay buffer, the base IoU is higher because model has access to additional full annotations. When diversity guidance (D-guide) is disabled, \tilde{D}_0 consists of random examples from the base dataset, resulting in worse base performance. When context guidance (C-guide) is disabled, a base image is uniformly sampled. When frequency guidance (F-guide) is disabled, a novel instance is sampled uniformly and is always pasted onto the base image. 95% confidence intervals over 20 trials are reported assuming that trial results are normally distributed. *: privileged. use fully-annotated masks when copy-paste is turned off.

770	METHOD	BASE	NOVEL	HM	BASE	NOVEL	HM	824
					ALL INSTANCES		SINGLE-INSTANCE ONLY	
773	PIFS(\mathcal{L}_2) [†]	46.2	20.2	28.1	46.1 (-0.2%)	17.6 (-12.9%)	25.4 (-9.6%)	825
774	PIFS(\mathcal{L}_2)+GAPS	49.1	25.8	33.8	49.2 (+0.2%)	25.1 (-2.7%)	33.2 (-1.8%)	826

Table 3. Performance of pseudo-labeling methods and GAPS under the more challenging single-instance learning setting on COCO-20ⁱ 5-shot. Only 1 novel instance out of potentially many instances in individual training images is annotated. The pseudo-labeling baseline, PIFS(\mathcal{L}_2)[†], yields substantially worse performance; whereas PIFS(\mathcal{L}_2)+GAPS has only minor performance decreases.

mental segmentation approaches use pseudo-labeling [5], we design a method PIFS(\mathcal{L}_2)[†], which simulates combining PIFS(\mathcal{L}_2) with pseudo-labeling to cope with partially-annotated sample. Here we allow PIFS(\mathcal{L}_2)[†] to be privileged and have access to additional information – the annotation of other non-novel background pixels – to simulate an oracle pseudo-labeling model which perfectly segments learned classes but recognize unseen novel classes as background.

The results are given in Table 3. We can observe that the pseudo-labeling baseline, PIFS(\mathcal{L}_2)[†], yields substantially worse performance when the model receives single-instance annotations despite having privileged access. We reason this is due to noisy labels generated by the pseudo-labeling process, where novel instances are incorrectly labeled as background. On the contrary, PIFS(\mathcal{L}_2)+GAPS shows only a minor performance decrease with single instances. This highlights the potential of copy-paste synthesis as an alternative to the existing pseudo-labeling paradigm to cope with the more realistic single-instance setting and robustness against false negative annotations.

More results and visualization. Due to space limits, we kindly refer readers to the supplementary material for more quantitative results that justify our design choices such as memory-replay buffer strategies and vRFS over other simple baselines. In addition, visualized qualitative results of sample segmentation and contextually-similar set construction can also be found in the attached supplementary material.

5. Conclusion and Discussion

In this paper, we demonstrate how the judicious use of copy-paste dramatically boosts the performance of incremental segmentation methods under the few-shot setting and enables learning with partially-annotated images. Our proposed GAPS technique selects representative exemplars in the memory-replay buffer and addresses the problems of class imbalance and contextual mismatch in synthesis.

In future work, we are interested in further application of copy-paste as a synthesis technique to cope with the background shifting problem for incremental segmentation. We believe that copy-paste can serve as a promising alternative to pseudo-labeling and modified loss to enable learning on partially-annotated images. We also believe that further optimizing exemplar selection and sampling strategies can lead to better guidance and lead to even better performance. Finally, the ability to learn with as few as one annotated instance in an image raises several intriguing possibilities. For example, integrating our work with learning-based interactive segmentation will enable human operators to continually and adaptively teach novel classes and correct failed predictions. This workflow has many interesting applications such as robot teleoperation where sparse annotations are preferable. Learning with weaker annotations, like bounding boxes or single clicks, and even self-supervision, is also an interesting direction to explore.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, 2021. 3, 5
- [2] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *ECCV*, 2018. 3
- [3] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Bulo, Elisa Ricci, and Barbara Caputo. Modeling the background for incremental learning in semantic segmentation. In *CVPR*, 2020. 1, 2, 3, 4, 6, 7
- [4] Fabio Cermelli, Massimiliano Mancini, Yongqin Xian, Zeynep Akata, and Barbara Caputo. Prototype-based incremental few-shot semantic segmentation. *arXiv preprint arXiv:2012.01415*, 2021. 1, 2, 3, 4, 5, 6, 7
- [5] Sungmin Cha, YoungJoon Yoo, Taesup Moon, et al. Ssl: Semantic segmentation with unknown label for exemplar-based class-incremental learning. *NeurIPS*, 2021. 1, 3, 4, 6, 7, 8
- [6] Ali Cheraghian, Shafin Rahman, Sameera Ramasinghe, Pengfei Fang, Christian Simon, Lars Petersson, and Mehrtash Harandi. Synthesized feature based few-shot class-incremental learning on a mixture of subspaces. In *ICCV*, 2021. 3
- [7] Arthur Douillard, Yifu Chen, Arnaud Dapogny, and Matthieu Cord. Plop: Learning without forgetting for continual semantic segmentation. In *CVPR*, 2021. 1, 3
- [8] Nikita Dvornik, Julien Mairal, and Cordelia Schmid. Modeling visual context is key to augmenting object detection datasets. In *ECCV*, 2018. 3, 7
- [9] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *ICCV*, 2017. 3
- [10] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 6
- [11] Dan Andrei Ganea, Bas Boom, and Ronald Poppe. Incremental few-shot instance segmentation. In *CVPR*, 2021. 3
- [12] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *CVPR*, 2021. 2, 3, 5
- [13] Agrim Gupta, Piotr Dollar, and Ross Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 5, 6
- [14] Bharath Hariharan, Pablo Arbelaez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *ICCV*, 2011. 6
- [15] Xiang Li, Tianhan Wei, Yau Pun Chen, Yu-Wing Tai, and Chi-Keung Tang. Fss-1000: A 1000-class dataset for few-shot segmentation. In *CVPR*, 2020. 3
- [16] Zhizhong Li and Derek Hoiem. Learning without forgetting. *TPAMI*, 2017. 3
- [17] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *NeurIPS*, 2017. 3

- [18] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *ECCV*, 2018. 3
- [19] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*. Elsevier, 1989. 2
- [20] Umberto Michieli and Pietro Zanuttigh. Incremental learning techniques for semantic segmentation. In *ICCVW*, 2019. 3
- [21] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *ICCV*, 2019. 6
- [22] Sylvestre-Alvise Rebiffé, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017. 3, 5
- [23] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. 2017. 3, 6
- [24] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong. Few-shot class-incremental learning. In *CVPR*, 2020. 3
- [25] Zhuotao Tian, Xin Lai, Li Jiang, Michelle Shu, Hengshuang Zhao, and Jiaya Jia. Generalized few-shot semantic segmentation. *arXiv preprint arXiv:2010.05210*, 2020. 3, 6
- [26] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. *TPAMI*, 2020. 3
- [27] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *ICCV*. 2019. 3
- [28] Shipeng Yan, Jiale Zhou, Jiangwei Xie, Songyang Zhang, and Xuming He. An em framework for online incremental learning of semantic segmentation. In *Proceedings of the 29th ACM International Conference on Multimedia*, 2021. 1, 3
- [29] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. In *ICLR*, 2018. 3
- [30] Chang-Bin Zhang, Jia-Wen Xiao, Xialei Liu, Ying-Cong Chen, and Ming-Ming Cheng. Representation compensation networks for continual semantic segmentation. In *CVPR*, 2022. 1, 3
- [31] Xiaolin Zhang, Yunchao Wei, Yi Yang, and Thomas S Huang. Sg-one: Similarity guidance network for one-shot semantic segmentation. *IEEE transactions on cybernetics*, 2020. 3