
On Overcompression in Continual Semantic Segmentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Class-Incremental Semantic Segmentation (CISS) is an emerging challenge of
2 Continual Learning (CL) in Computer Vision. In addition to the well-known issue
3 of catastrophic forgetting, CISS suffers from the semantic drift of the background
4 class, further increasing forgetting. Existing attempts aim to solve this using pseudo-
5 labelling, knowledge distillation or model freezing. We argue and demonstrate
6 that frozen or rigid models suffer from poor expressibility due to overcompression.
7 We improve on these methods by focusing on the offline training process and
8 the expressiveness of the learnt representations. Beyond the characterisation and
9 demonstration of this issue in terms of the Information Bottleneck principle, we
10 show the benefit of two practical measures: (1) using shared but wider convolution
11 modules before final classifiers to improve scaling for new, continual tasks; (2)
12 introducing dropout into the encoder-decoder architecture to improve regularisation
13 and decrease the overcompression of information in the representation space. We
14 improve the IoU on the 15-1 and 10-1 scenarios by over 2% and 3% respectively
15 while maintaining a smaller memory and MAdds footprint. Last, we propose a
16 new benchmark setting that lies closer to the nature of lifelong learning to drive the
17 development of more realistic and valuable architectures in the future.

18 1 Introduction

19 Continual Learning (CL) aims to address the shortcoming of standard supervised learning that
20 requires large quantities of labelled data. It pursues a more natural, human-like ability to quickly and
21 continuously learn from small exposures to training data. The learner experiences a stream of tasks
22 whose relatedness is not known beforehand [4] and has the potential to quickly learn a new task by
23 re-using past experiences. Clearly, we have to prevent the model from naively storing all experiences.
24 This constraint enforces the need to efficiently compress the knowledge to a modest size and limit the
25 computation used to recreate representation for a task at hand. Recent work on Continual Learning
26 focuses on the numerous ways to combat forgetting using replay methods, continually changing the
27 models' architecture to adapt to new knowledge or distillation methods.

28 Class-Incremental Semantic Segmentation (CISS) is a recently recognised problem that, on top of
29 forgetting, deals with the unique issue of *background shift* [2] of the unknown pixels (see Section
30 2.3.1). Existing work focuses on the efficient transfer of existing knowledge between the continual
31 steps while acquiring enough information to work with the tasks at hand. The prevalent approach
32 is using rigid encoders that aim to maintain initial features throughout training while focusing on
33 attempts to learn new classes. As a result, the most successful approaches tend to fair well with initial
34 classes while struggling with the ones learned in the online fashion.

35 Most problems, however, tend to also have an offline phase where the model is initialised and learns
36 the initial data and tasks. Mirzadeh et al. [27] show that wide convolutional models outperform

37 standard ResNet-based encoders [17] in Continual Learning, proving that the choice of the initial
 38 architecture is fundamental to combating catastrophic forgetting. Further inspired by the Information
 39 Bottleneck principle [35] we claim that we should put just as much effort into the training of the
 40 *offline* phase as we put into the *online* phase for best results.

41 1.1 Deep neural networks and overcompression

42 When developing models for CL, it is important to develop from pre-trained models that have many
 43 diverse features. Adding additional filters to an existing model is resource intensive, and provides
 44 training challenges. Therefore, it is important to start with an extensive feature set that we can adapt to
 45 the different scenarios as they appear. However, in initial offline training on the limited data, the many
 46 unimportant features are pulled closer to existing relevant filters, and such excessive compression
 47 means that we might remove features that are crucial for learning a prediction of one of the future
 48 classes.

49 The Information Bottleneck principle (IB) [35] can
 50 partially explain and support the above difficulty of
 51 training Continual Learning models with gradient-
 52 based optimisation. The IB interpretation of learn-
 53 ing suggests that hidden layers in neural networks
 54 learn to maximally compress the characteristics in
 55 the input data that are irrelevant to the task at hand,
 56 and where the higher level layers learn to compress
 57 first and the the lower level layers as training pro-
 58 gresses. On the other hand, the uncertainty about
 59 the future requirements posed by CL requires us
 60 to maintain as many high-quality features as possi-
 61 ble. Thus, it is in our best interest to maintain
 62 diverse features while limiting overcompression in
 63 Continual Learning.

64 Although the importance of the Information Bottleneck
 65 theory has been questioned by some papers [31],
 66 it serves as an interesting point of view on the
 67 learning phases of neural networks and can be
 68 noticed empirically. We know that wider models
 69 better manage catastrophic forgetting [27], proving
 70 that passing information through the bottleneck,
 71 which inherently imposes compression, can have
 adverse effects on future tasks. Therefore, we use
 two simple approaches to verify if the reduction of
 overcompression can give us attainable gains in CL,
 in particular in CISS:

- 72 1. Increase the width of the last layer to reduce the bottleneck effect and increase the amount
 73 of information for each subsequent continual task;
- 74 2. Introduce dropout layer between the pre-trained encoder and decoder to introduce uncertainty
 75 into the model, encouraging a more robust representation that is later reused by the future
 76 steps.

77 By developing Dropout Continual Semantic Segmentation (DCSS), which incorporates only these two
 78 simple additions while improving on previous works by a sizeable margin (Figure 1), we demonstrate
 79 that improving the representation provided by the initial network has a significant impact on overall
 80 results.

81 2 Related work

82 2.1 Catastrophic forgetting

83 Catastrophic forgetting [30, 24] is a phenomenon observed primarily in Continual Learning where
 84 the earlier learned concepts are forgotten while incorporating the more recent samples. Forgetting
 85 appears when previously-learned representations are degraded by more recent exposures, a typical
 86 case in all SGD-based algorithms. Solutions proposed to address this issue can be grouped into the
 87 following categories:

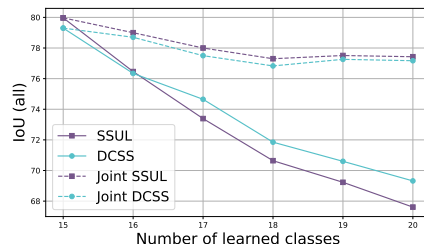


Figure 1: **IoU and Joint upper-bounds comparison** of the proposed DCSS model and SSUL [3] in the 15-1 scenario. DCSS achieves better online scores despite being worse in offline training, suggesting that we can improve CL performance by increasing the expressiveness of the learnt representations.

88 **Rehearsal learning** In rehearsal learning, we utilise the fact that, although we have lost access to
 89 the original training data from the past, we are usually allowed to maintain the data in a different
 90 form. [3] store exemplar images in a small memory buffer that can be replayed to simulate previous
 91 data distribution. Compressed features in the form of embeddings [15, 20] can be used as an efficient
 92 form of memory that can also have advantages in terms of privacy. Finally, we can use generative
 93 strategies [32] that can efficiently store and reconstruct past experiences to enable the replay of raw
 94 images. [21] uses a brain-inspired dual-memory system where the new memories are consolidated
 95 from recent memory to a long-term memory using a generative model, similar to mechanisms that
 96 occur during sleep.

97 **Adaptive architecture** Other approaches focus on the adaptability of the model architectures. To
 98 integrate new classes and tasks, we can extend the architecture using feature extractors with domain-
 99 specific trainable layers [24, 39]. Model freezing can help with performance degradation during
 100 fine-tuning of the new models for new tasks. It is also possible to adapt the networks without explicitly
 101 adding new modules. [12, 11] dynamically re-arranges existing sub-networks, each specialised in
 102 one specific task, to account for new knowledge. Moreover, continually changing data distributions
 103 can be accounted for by explicitly correcting the classifier drift [1, 38].

104 **Knowledge distillation** Knowledge distillation methods consider the use of a *teacher-student*
 105 approach. A knowledgeable teacher model trained on past inaccessible experiences can inform the
 106 current model of the past. Therefore, distillation aims at constraining the model to prevent forgetting
 107 as it changes to adopt new data. There are several ways to constrain the model, with the most salient
 108 methods being applied to the weights [22], gradients [4, 25] or output probabilities [24, 29].

109 2.2 Information Bottleneck principle

110 Tishby et al. in the Information Bottleneck principle [35] claim that modern Deep Neural Networks
 111 undergo two phases of learning. During the generalisation phase, the model learns an internal feature
 112 representation to extract high-level information, used for final prediction. In the compression phase,
 113 unused and noisy features that prevent robust adaptation to the data are stripped away and removed
 114 from the feature space, ultimately improving useful features’ signal quality. Tishby et al. measured
 115 the mutual information $I(X; T)$ and $I(T; Y)$, which quantify the hidden layer’s information about
 116 the input and the output, where X is the input, Y is the output and T is the internal representation.
 117 They showed that the amount of information about the output $I(T; Y)$ steadily increases until a
 118 high number of epochs is experienced, and the model starts overfitting. The information about the
 119 input $I(X; T)$ raises initially but quickly starts decreasing. This decrease can be considered as the
 120 compression phase because we remove information from the input that does not contribute much to
 121 the output. Compression accelerates rapidly when the number of training epochs is high.

122 2.3 Class-incremental semantic segmentation

123 2.3.1 Problem definition and notation

124 We follow the definition of van de Ven et al. [36] who use three distinct scenarios for Continual
 125 Learning: Task-Incremental, Domain-Incremental and Class-Incremental learning. This work focuses
 126 on the Class-Incremental setting initially formulated by Cermelli et al. [2] where we split the set
 127 of semantic segmentation classes into $t = 1, \dots, T$ incremental tasks, with $t = 1$ being the offline
 128 step. The model learns each disjoint set of classes incrementally while trying not to forget the
 129 previous steps. At each step a new task t arrives with a training dataset D_t that consists of pairs $(\mathbf{x}^t$
 130 $, \mathbf{y}^t)$, where $\mathbf{x}^t \in \mathcal{X}^N$ denotes an input image of N pixels, and $\mathbf{y}^t \in (\mathcal{C}^t \cup \{c_b^t\})^N$ denotes the
 131 corresponding ground-truth pixel labels, where \mathcal{C}^t is the set of classes that are seen in task t and c_b^t is
 132 the background class for that task.

133 At test time, the semantic segmentation model f_θ^t is required to predict whether a pixel belongs to
 134 a class learned so far, $c \in \mathcal{C}^{1:t-1} \cup \mathcal{C}^t$, or the true background class c_b . However, the labels of an
 135 image from a task t only contains classes from \mathcal{C}^t , not classes seen in past or future tasks. Thus,
 136 during training, the background label c_b^t is also assigned to the pixels of potential objects that belong
 137 to past classes $\mathcal{C}^{1:t-1}$ and the future classes $\mathcal{C}^{t+1:T}$ (Equation 1).

$$c_b^t = c_b \cup \underbrace{\mathcal{C}^{1:t-1}}_{past} \cup \underbrace{\mathcal{C}^{t+1:T}}_{future} \quad (1)$$

138 This counter-intuitive behaviour of the background class’ labels causes the problem of *background*
 139 *shift* [2], on top of the catastrophic forgetting found in CL. Each pixel with background label can be a
 140 class from the past, future or the actual background. This introduces unwanted noise into the model
 141 because the same object can be labelled differently, depending on the step t . Once the learning of task
 142 t by the model f_θ^t is done, the prediction for pixel i of an input image \mathbf{x} at test time is obtained by

$$\tilde{\mathbf{y}}_i^t = \arg \min_{c \in \mathcal{Y}^{1:t}} f_{\theta,c}^t(\mathbf{x}_i^t) \quad (2)$$

143 for all the classes learned so far in $\mathcal{C}^{1:t}$. The performance is measured by the Intersection-over-Union
 144 (IoU) metric for the three splits of $\mathcal{C}^{1:t}$. Cermelli et al. [2] introduce two different settings for
 145 Class-Incremental Semantic Segmentation: *Disjoint* and *Overlapped*. Both cases consider training
 146 examples where only the background class c_b^t and classes in \mathcal{C}^t are used as labels. In *Disjoint*, only
 147 the old and present classes can appear, while in *Overlapped* pixels can belong to any old, current and
 148 future classes. Thus, following the previous works, we also focus only on the *Overlapped* setting as it
 149 is more challenging and closer to real conditions.

150 2.3.2 Existing work

151 Class-Incremental Semantic Segmentation is a relatively new problem, with few works attempting to
 152 address it. Cermelli et al. in MiB [2] first defined the issue of *background shift*, a problem unique to
 153 Continual Semantic Segmentation. To alleviate the issue, they proposed to adapt the loss function
 154 such that it takes into consideration potential previous classes $\mathcal{C}^{1:t-1}$ that could be included in c_b^t .
 155 They sum the output probabilities of $\mathcal{C}^{1:t}$ and c_b to match the available label c_b^t that contains them
 156 when computing the cross-entropy loss. Conversely, they sum the probabilities of \mathcal{C}^t and c_b^{t-1} for
 157 knowledge distillation to match the output of f_θ^{t-1} . The summation of class probabilities to modify
 158 the cross-entropy and distillation losses is a novel approach in semantic segmentation, but it makes it
 159 hard to learn the underlying probability distribution for the classes at each pixel.

160 PLOP [9] extended the knowledge distillation on output probabilities from [2] with Localised Pooled
 161 Distillation [9], a form of heavy attention transfer that prevents from forgetting at the cost of decreased
 162 flexibility. Additionally for a task t , the target labels \mathbf{y}_i^t are augmented with generated *pseudo-labels*
 163 [23] with a model f_θ^{t-1} from a previous step. This yields the target semantic maps $\tilde{\mathbf{y}}^t$ containing
 164 background class c_b^t , pseudo labels $\hat{\mathbf{y}}^{t-1}$ and current labels \mathbf{y}^t .

165 SSUL [3] finds success with further increase of model’s rigidity using model freezing. Instead of
 166 fine-tuning the whole model at task t , all parameters from f_θ^{t-1} are being frozen in f_θ^t , preventing
 167 any changes to the prediction of classes $\mathcal{C}^{1:t}$. Moreover, SSUL uses saliency maps generated by a
 168 separate, off-the-shelf salient object detector [19] to predict a region of interest from the background
 169 which helps to differentiate the difference between true background and background that may contain
 170 a past or future class, which they label using the *unknown class*. Similarly to pseudo-labels, the
 171 saliency map is added to the labels \mathbf{y}^t and predicted by the model during training. The unknown class
 172 c_u and the background class c_b predictions are combined at test time, representing the unclassified
 173 pixels. Taking advantage of the foreground prediction, SSUL performs weight transfer from the
 174 unknown class to the new classifier $\phi_{c_u}^{t-1} \rightarrow \phi_c^t$ at the beginning of each continual step.

175 3 DCSS model design

176 3.1 Wider classifiers with shared representation

177 Overcompression can hurt the learning of future tasks. Our goal is to maintain as many features as
 178 possible for future use. One simple way to achieve this is to use wider networks that are no longer
 179 a bottleneck to the information signal. To verify our hypothesis we decide to increase the number
 180 of channels in the last layer of the frozen, offline representation. In this way, we should have more
 181 information available for the continual steps.

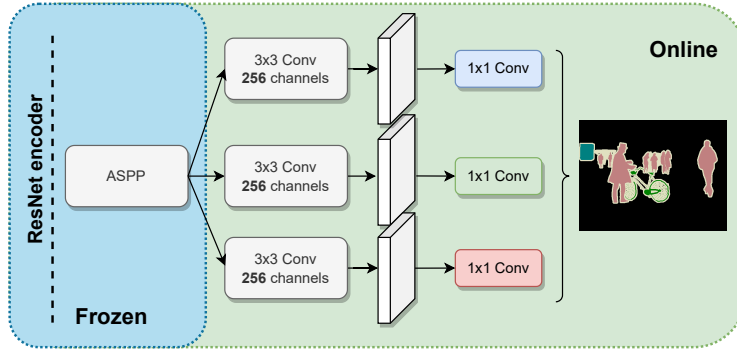


Figure 2: **High-level architecture of SSUL** [3]. Each step has its own 3×3 convolution block that is trained in the online phase.

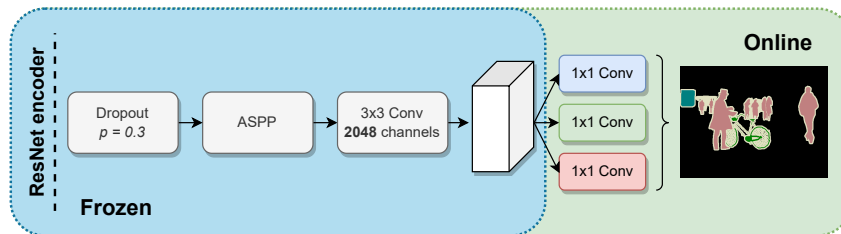


Figure 3: **High-level architecture of our DCSS** model. One HEAD module is trained in the offline step and frozen. Dropout layer is placed right after the pre-trained encoder to increase uncertainty in the learnt decoder.

182 In the DeepLabV3 architecture [5] used for this project the final 1×1 classifier is preceded by the
 183 ASPP module and one 3×3 convolution. The last 3×3 convolution is important to aggregate and
 184 mix multi-scale features produced by the ASPP module. We will refer to this 3×3 convolution layer
 185 as HEAD for readability reasons. We increase the size of HEAD output channels from 256, as in the
 186 original DeepLabV3, to 2048 channels. Hence, a more detailed feature space is available for the
 187 classifiers to exploit. Recent papers tackling CISS, including MiB, PLOP and SSUL, replicate the
 188 HEAD module for each step. Therefore, each set of classes \mathcal{C}^t has its own HEAD and final classifier,
 189 as seen in Figure 2. However, the HEAD layer is a costly operation that has over 0.5M parameters
 190 with 256 channels and over 4.5M parameters when using 2048 channels. Therefore, to prevent
 191 the large parameter accumulation we use only a single HEAD layer, frozen after offline training
 192 and add only 1×1 convolutions at online steps $t > 1$. Additionally, we use depthwise-separable
 193 convolution in the shared HEAD to produce the same activation map using a cheaper operation [8]
 194 with fewer parameters and MAdds. In our work, we have used an intermediate width multiplier of 4
 195 and 4 convolution groups. Standard shared HEAD has almost twice the parameters of the depthwise-
 196 separable implementation (4.71M vs 2.68M). The high-level architecture comparison can be found in
 197 Figures 2 and 3.

198 The reason why at step $t = 1$ the DCSS is larger than standard DeepLabV3, despite both being trained
 199 in an offline fashion, is the requirement of having a separate classifier for the background class c_b and
 200 the unknown class c_u , which means that we already have three HEAD modules at $t = 1$, including the
 201 one for classes in $\mathcal{C}^{t=1}$. In contrast, increasing the size of the HEAD to produce 2048 channels with
 202 separable convolution as in DCSS, compared to 256 channels in SSUL, yields a significantly smaller
 203 model in the long run (Figure 4) while achieving better results than SSUL.

204 3.1.1 Dropout as a source of uncertainty

205 Dropout [34] randomly sets activations of a layer to 0 with probability p , effectively disabling
 206 their contribution to the calculation of the output and introducing uncertainty that should hinder
 207 overcompression. Since the placement of the Dropout is essential, we follow Spilsbury et al. [33] and
 208 place the Dropout layer in between the encoder and decoder modules. The idea behind this decision

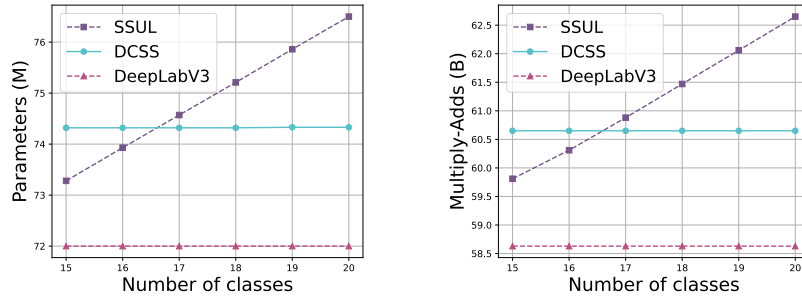


Figure 4: **Resource comparison** in the **15-1** scenario. SSUL adds significantly more trainable parameters at each step, surpassing DCSS size and complexity in just 2 steps.

209 is that we want to prevent interference in the pre-trained encoders while maximising the effect on the
 210 decoder.

211 We test two popular Dropout approaches in CNNs: standard Dropout and channel Dropout, where the
 212 whole convolutional channel can be dropped with probability p . Moreover, we use ScheduledDropout
 213 [33] that linearly increases the dropout probability from 0 to p during training. During the exploration
 214 phase (generalisation), smaller values of p help with feature accumulation. Heavier regularisation is
 215 more useful during the exploitation phase (compression), thus requiring higher values of p . In the end,
 216 we obtain similar benefits while improving learning convergence. Last, we suggest the removal of the
 217 Dropout after offline training. Limiting the access to features in online training will only restrain the
 218 model from using them while offering few benefits in generalisation since we freeze the encoder and
 219 are unlikely to learn new features.

220 4 Experiments

221 **Implementation details** The code for this project has been based on the implementations of PLOP
 222 [9] papers and SSUL [3] that use the DeepLabV3 architecture [5]. To compare the results of DCSS to
 223 previous work [2, 9, 3], we also use ResNet101 while keeping in mind that a smaller encoder would
 224 have a better performance-to-cost ratio. A standard learning rate of 0.01 has been used, with the
 225 value decreased to 0.001 for the backbone when using pre-trained weights [5]. We use *poly*¹ learning
 226 rate scheduler, as proposed by the DeepLabV3 authors. The same data augmentation has been used
 227 as in SSUL. We use a batch size of 24 for all experiments compared to the size of 32 used by SSUL,
 228 as we found more success with smaller values, especially for the continual steps. We use binary
 229 cross-entropy loss with sigmoid activation function, as in SSUL. Pseudo-labels are added to labels in
 230 the continual steps, with an entropy-based threshold of $\tau = 0.9$. In CISS problem we train each step
 231 for 50 epochs. The experiments were run on 4 RTX 2080Ti GPUs using the `torch.distributed`
 232 library [7] to work around the limited GPU memory. Each experiment, including the offline step,
 233 takes approximately 2-4 hours, depending on the scenario. In the case where the offline step can be
 234 reused, the training usually is decreased to 1-2 hours.

235 **PASCAL VOC 2012** The experiments are evaluated on the PASCAL VOC 2012 [10] semantic
 236 segmentation dataset with 20 foreground object classes and one common background class. The
 237 dataset is augmented with the extra contour annotations [14] of the PASCAL VOC 2011 dataset.
 238 Following [9], we split the training set and used 80% for training and 20% for validation, with
 239 the testing set left untouched for model evaluation. To help distinguish the background containing
 240 potential future classes from the true background, SSUL [3] proposed the use of an off-the-shelf
 241 salient object detector [19] to predict a region of interest. Thus, we follow their implementation and
 242 extend the labels of the PASCAL VOC training set to include the additional foreground class.

¹The scheduler has been called “poly” by the authors of DeepLabV3, even though the learning rate is not a polynomial [37].

Table 1: **Main results for the CISS problem.** DCSS previous works in new classes (middle columns) in all scenarios and achieve improvements in combined scores in the more difficult scenarios with multiple tasks.

Method	VOC 10-1 (11 tasks)			VOC 15-1 (6 tasks)			VOC 5-3 (6 tasks)			VOC 19-1 (2 tasks)			VOC 15-5 (2 tasks)		
	0-10	11-20	all	0-15	16-20	all	0-5	6-20	all	0-19	20	all	0-15	15-20	all
ILT [26]	7.15	3.67	5.50	8.75	7.99	8.56	22.51	31.66	29.04	67.75	10.88	65.05	67.08	39.23	60.45
MiB [2]	12.25	13.09	12.65	34.22	13.50	29.29	57.10	42.56	46.71	71.43	23.59	69.15	76.37	49.97	70.08
PLOP [9]	44.03	15.51	30.45	65.12	21.11	54.64	17.48	19.16	18.68	75.35	37.35	73.54	75.73	51.71	70.09
SSUL [3]	71.31	45.98	59.25	77.31	36.59	67.61	71.17	45.38	52.75	77.73	29.68	75.44	77.82	50.10	71.22
DCSS (ours)	73.34	50.20	62.32	77.66	42.69	69.33	68.10	48.83	54.34	77.22	36.85	75.30	77.49	51.49	71.30
Joint (V3)	78.41	76.35	77.43	79.77	72.35	77.43	76.91	77.63	77.43	77.51	77.04	77.43	79.77	72.35	77.43
Joint (DCSS)	77.80	76.58	77.22	78.77	72.25	77.22	76.29	77.61	77.22	77.30	75.60	77.22	78.77	72.25	77.22

243 **Evaluation metrics** The difficulty of the Class-Incremental challenge also depends on the number
 244 of steps and the number of new classes. Thus, our model is evaluated on five different scenarios
 245 generated from the PASCAL VOC, each with a varying level of difficulty: **10-1** (11 tasks), **15-1**
 246 (**6 tasks**), **5-3** (6 tasks), **19-1** (2 tasks) and **15-5** (2 tasks). The numbers in each scenario define the
 247 number of classes introduced at each step. For example, **5-3** (6 tasks) means learning 5 base classes
 248 at the offline step $t = 1$, followed by 5 incremental steps $t \in \{2, \dots, 6\}$ introducing 3 new classes at
 249 each step, yielding 6 training steps covering all 20 classes for PASCAL VOC. We report the results as
 250 mean intersection-over-union (IoU) for three categories of classes: old (offline), new (online) and all
 251 classes combined (offline + online). This split helps to distinguish the difference between forgetting
 252 and issue with obtaining the knowledge for the new tasks.

253 4.1 Experimental results

254 In Table 1 we observe that DCSS consistently outperforms other models while having a simpler and
 255 more extendable architecture that is also easier to train (Figure 3). In the most popular **15-1** scenario,
 256 DCSS achieves a 1.72% improvement over SSUL. The biggest overall gain of over 3% can be found
 257 in **10-1** scenario that has a larger number of continual steps. Most importantly, in all scenarios we
 258 have improved the mean score of the new, continual classes, up to 6% in the case of **15-1**. This gain
 259 can be a surprising result, considering that we are only training a single 2048-dimensional vector at
 260 each step. We conclude that the current counter-forgetting measures used in CISS effectively prevent
 261 models from adding new features to the representation and, thus, it is enough to simply learn a linear
 262 mapping of *existing* features during the continual phase.

263 To compare our results to the ones produced by
 264 previous works, we have to constrain ourselves to
 265 a model of a similar learning capacity. Figure 1
 266 proves that we increase Continual Learning performance
 267 despite the decreased overall performance of DCSS
 268 in offline learning, further signifying the importance
 269 of our contributions to online learning. What is most
 270 important, all of this is achieved with a smaller and
 271 easier to train model. In **15-1** DCSS has over 2M
 272 parameters less than SSUL (Figure 4), while in the
 273 extreme case of **10-1** the model is actually smaller
 274 by over 5M parameters and the difference grows
 275 linearly with each additional step.

276 Figure 5 shows improvement over the range of the
 277 Dropout probabilities p compared to no Dropout
 278 ($p = 0$). Simply by adding the Dropout layer during
 279 the offline training we have increased the performance
 280 during the online steps, yielding a more accurate model
 281 on average. We have found more success with the
 282 standard dropout (Dropout1d) rather than the channel
 dropout (Dropout2d), despite literature stating the reverse
 for CNNs [33, 18]. The results follow our intuition,
 which suggests that dropping the whole channels
 removes a feature across the whole image, whereas

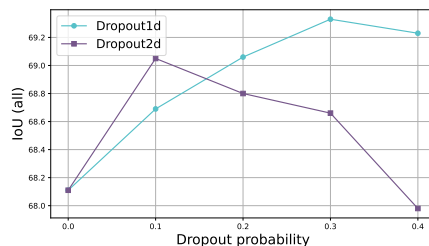


Figure 5: **IoU comparison of Dropout in 15-1 scenario.** We notice that a small amount of uncertainty produced by the Dropout can increase the IoU of the model.

Table 2: **Experimental results on the proposed RCISS protocol.** Both SSUL and DCSS decrease their performance in new tasks, although the scores lie surprisingly close to the original scores. This suggests that more training cannot extract missing features from the frozen model.

Method	Epochs	VOC 10-1 (11 tasks)			VOC 15-1 (6 tasks)			VOC 5-3 (6 tasks)			VOC 15-5 (2 tasks)		
		0-10	11-20	all	0-15	16-20	all	0-5	6-20	all	0-15	15-20	all
SSUL [3]	50	71.31	45.98	59.25	77.31	36.59	67.61	71.17	45.38	52.75	77.82	50.10	71.22
DCSS	50	73.34	50.20	62.32	77.66	42.69	69.33	68.10	48.83	54.34	77.49	51.49	71.30
SSUL [3]	1	73.23	38.10	56.50	77.34	32.44	66.65	69.90	27.52	39.63	76.68	47.64	69.76
DCSS	1	73.78	43.26	59.24	77.16	37.46	67.71	68.44	38.84	47.30	77.46	45.60	69.88

283 standard dropout can work inter-channel, which is vital in pixel-level tasks. Removing whole
 284 channels can prevent the propagation of features, while removing specific neurons can decrease the
 285 reliance on more salient parts of the image and promote more holistic attention.

286 4.1.1 New protocols and evaluation

287 We have seen the performance of DCSS on the tasks proposed initially by [2]. DCSS offers a better
 288 IoU score while being a generally cheaper model to train in the long run. We argue, however, that the
 289 evaluation protocol has an unrealistic assumption about the ability to store the continual training data
 290 and use it multiple times over the step t . This is an unnatural setting that promotes irrelevant models.
 291 Continual Learning models should aim to efficiently use the available data and aim for an approach
 292 similar to few-shot learning, where only a few examples for each class are available. Therefore, we
 293 also evaluate our DCSS model on the proposed Restricted Class-Incremental Semantic Segmentation
 294 protocol (RCISS), where we train each continual step for only one epoch instead of the 50 as in
 295 SSUL. This scenario lies closer to the *lifelong learning* setting by Chaudhry et al. [4].

296 In the experiments on RCISS we remove learning rate scheduling and use a fix learning rate of 0.01.
 297 Moreover, we reduce the batch size in the continual phase from 24 to 4, effectively increasing the
 298 number of backpropagation steps in the hope that, while being noisy, it will better utilise the limited
 299 exposure to data. Table 2 shows the result of our experiments on four scenarios in RCISS, similar to
 300 previous scenarios in standard CISS. We notice that both models struggle more with learning multiple
 301 classes at once (**5-3**), although DCSS manages to outperform SSUL in the continual classes by a
 302 considerable margin. However, we notice that the results for 1 epoch are still similar to the ones with
 303 full 50 epochs of training. We find this result interesting as it signifies the difficulty of extracting
 304 performance in continual learning even with extensive training.

305 5 Discussion and conclusions

306 We showed that the current approaches to CISS problem put too much importance on the continual
 307 phase of learning. The small difference between RCISS and CISS shows that there is a limit to the
 308 amount of information that can be obtained by simply learning new tasks with a continual manner
 309 while having significant constraints on the model flexibility to prevent catastrophic forgetting. The
 310 balance between rigidity and flexibility is hard to achieve, and most recent works found success with
 311 the focus on the former. With heavily constrained or frozen representations we are unable to adapt
 312 to new tasks or classes. Thus, assuming that the potential cost of adding features is much higher
 313 than removing them, we take a step back to consider what can be done to maintain potentially useful
 314 features for the future that would be pruned otherwise. We show that with a few simple changes to
 315 the offline training protocol we can propagate more information to future learning steps, lessening
 316 the need for complex online methods.

317 We adapted the idea of wide continual networks [27] for model freezing. Wider architecture of the
 318 decoder can help to maintain more features for CL; surprisingly, even though we are essentially
 319 learning a linear mapping of the embedding from the frozen model we improve the IoU in CISS. This
 320 further suggests that having a rich representation outweighs most attempts at learning new features
 321 using additional trainable layers while having the benefit of being easily scalable and far simpler to
 322 train. The simple addition of a single Dropout layer further improved the IoU, surpassing all current
 323 approaches based on the DeepLabV3 model.

324 Despite that, our work aims to show not the benefits of the proposed methods but rather the general
325 lack of consideration of overcompression in Continual Learning. Information Bottleneck principle
326 suggests that deep learning models have a generalisation and compression phase. Although useful
327 in offline learning, compression is counterproductive in Continual Learning where we don't know
328 the future requirements. We believe that most short-term gains in CL could be achieved by focusing
329 on training with future's uncertainty in mind. Further gains could be achieved with a different
330 approach to learning that will increase the generalisation. Recent trends tend to remove the reliance
331 on supervised training for a more contextualised signal using self-supervised learning in the form of
332 contrastive models [16, 6] or multi-modal representation learning [13], where a more holistic type of
333 learning must appear to capture the whole context. Altogether, these approaches should allow for
334 encoders that are less task-specific, providing the required flexibility in Continual Learning.

335 5.1 Future work

336 For subsequent work on the topic of CISS, we think a good idea would be to look at expanding our
337 findings to different architectures. For example, SATS [28], used a Transformer-based encoder and
338 self-attention transfer to achieve state-of-the-art results on CISS. So, while we expect our encoder-
339 agnostic findings to maintain their validity for other architectures, it would be interesting to explore
340 the topic of overcompression of Transformer models in the context of CL. Additionally, the research
341 community is in dire need of dedicated Continual Semantic Segmentation datasets, which need to
342 emphasize the issue of the background shift problem and introduce temporal and spatial locality of
343 the data, which should spur new interest in this topic.

344 Societal impacts

345 This paper touches upon the topic of Continual Learning which greatly increases the move towards
346 more private and efficient training. Our contribution emphasizes the need for optimal offline training
347 procedures that have immense impact on the future utility during continual phase, improving the ease
348 of training and, in turn, decreasing the required compute and data. We do not see any foreseeable
349 negative impacts of this work while noticing the positive impacts of more efficient CL models that
350 should reduce the computational burden in the short term and allow for more intelligent agents in the
351 long term.

352 References

- 353 [1] Eden Belouadah and Adrian Popescu. ScaLL: Classifier Weights Scaling for Class Incremental
354 Learning, 2020. arXiv:2001.05755.
- 355 [2] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Bulò, Elisa Ricci, and Barbara Ca-
356 puto. Modeling the Background for Incremental Learning in Semantic Segmentation, 2020.
357 arXiv:2002.00718.
- 358 [3] Sungmin Cha, Beomyoung Kim, Youngjoon Yoo, and Taesup Moon. SSUL: Semantic
359 Segmentation with Unknown Label for Exemplar-based Class-Incremental Learning, 2021.
360 arXiv:2106.11562.
- 361 [4] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient
362 Lifelong Learning with A-GEM. In *ICLR*, 2019.
- 363 [5] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking Atrous
364 Convolution for Semantic Image Segmentation, 2017. arXiv:1706.05587.
- 365 [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework
366 for Contrastive Learning of Visual Representations, 2020. arXiv:2002.05709.
- 367 [7] R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A Matlab-like Environment for Machine
368 Learning. In *BigLearn, NIPS Workshop*, 2011.
- 369 [8] Elliot J. Crowley, Gavin Gray, and Amos Storkey. Moonshine: Distilling with Cheap Convolu-
370 tions, 2017. arXiv:1711.02613.

- 371 [9] Arthur Douillard, Yifu Chen, Arnaud Dapogny, and Matthieu Cord. PLOP: Learning without
372 Forgetting for Continual Semantic Segmentation, 2021. arXiv:2011.11390.
- 373 [10] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman.
374 The pascal visual object classes challenge: A retrospective. *International Journal of Computer
375 Vision*, 111(1):98–136, January 2015.
- 376 [11] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu,
377 Alexander Pritzel, and Daan Wierstra. PathNet: Evolution Channels Gradient Descent in Super
378 Neural Networks. 2017. arXiv:1701.08734.
- 379 [12] Siavash Golkar, Michael Kagan, and Kyunghyun Cho. Continual Learning via Neural Pruning.
380 03 2019. arXiv:1903.04476.
- 381 [13] Wenzhong Guo, Jianwen Wang, and Shiping Wang. Deep Multimodal Representation Learning:
382 A Survey. *IEEE Access*, 7:63373–63394, 2019.
- 383 [14] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse
384 detectors. In *2011 International Conference on Computer Vision*, 2011.
- 385 [15] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. RE-
386 MIND Your Neural Network to Prevent Catastrophic Forgetting. In *Proceedings of the European
387 Conference on Computer Vision (ECCV)*, 2020.
- 388 [16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for
389 Unsupervised Visual Representation Learning, 2019. arXiv:1911.05722.
- 390 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image
391 Recognition, 2015. arXiv:1512.03385.
- 392 [18] Mingjie He, Jie Zhang, Shiguang Shan, Xiao Liu, Zhongqin Wu, and Xilin Chen. Locality-
393 Aware Channel-Wise Dropout for Occluded Face Recognition. *IEEE Transactions on Image
394 Processing*, 31:788–798, 2022.
- 395 [19] Qibin Hou, Ming-Ming Cheng, Xiaowei Hu, Ali Borji, Zhuowen Tu, and Philip H. S. Torr.
396 Deeply Supervised Salient Object Detection with Short Connections. *IEEE Transactions on
397 Pattern Analysis and Machine Intelligence*, 41(4):815–828, apr 2019.
- 398 [20] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-Efficient
399 Incremental Learning Through Feature Adaptation. 2020. arXiv:2004.00713.
- 400 [21] Ronald Kemker and Christopher Kanan. FearNet: Brain-Inspired Model for Incremental
401 Learning. In *International Conference on Learning Representations*, 2018.
- 402 [22] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, An-
403 drei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis
404 Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic
405 forgetting in neural networks, 2016. arXiv:1612.00796.
- 406 [23] Dong-Hyun Lee. Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method
407 for Deep Neural Networks. *ICML 2013 Workshop : Challenges in Representation Learning
408 (WREPL)*, 07 2013.
- 409 [24] Zhizhong Li and Derek Hoiem. Learning without Forgetting. *IEEE Transactions on Pattern
410 Analysis and Machine Intelligence*, 40(12):2935–2947, 2018.
- 411 [25] David Lopez-Paz and Marc’ aurelio Ranzato. Gradient Episodic Memory for continual learning.
412 June 2017. arXiv:1706.08840.
- 413 [26] Umberto Michieli and Pietro Zanuttigh. Incremental Learning Techniques for Semantic Seg-
414 mentation. In *International Conference on Computer Vision (ICCV), Workshop on Transferring
415 and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2019.

- 416 [27] Seyed Iman Mirzadeh, Arslan Chaudhry, Dong Yin, Timothy Nguyen, Razvan Pascanu,
417 Dilan Gorur, and Mehrdad Farajtabar. Architecture Matters in Continual Learning, 2022.
418 arXiv:2202.00275.
- 419 [28] Yiqiao Qiu, Yixing Shen, Zhuohao Sun, Yanchong Zheng, Xiaobin Chang, Weishi Zheng, and
420 Ruixuan Wang. SATS: Self-Attention Transfer for Continual Semantic Segmentation, 2022.
421 arXiv:2203.07667.
- 422 [29] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert.
423 iCaRL: Incremental Classifier and Representation Learning, 2016. arXiv:1611.07725.
- 424 [30] Anthony V. Robins. Catastrophic Forgetting, Rehearsal and Pseudorehearsal. *Connect. Sci.*,
425 7:123–146, 1995.
- 426 [31] Andrew Michael Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky,
427 Brendan Daniel Tracey, and David Daniel Cox. On the Information Bottleneck Theory of Deep
428 Learning. In *International Conference on Learning Representations*, 2018.
- 429 [32] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual Learning with Deep
430 Generative Replay. In *Proceedings of the 31st International Conference on Neural Information
431 Processing Systems*, 2017.
- 432 [33] Thomas Spilisbury and Paavo Camps. Don't ignore Dropout in Fully Convolutional Networks,
433 2019. arXiv:1908.09162.
- 434 [34] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
435 Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine
436 Learning Research*, 15(56):1929–1958, 2014.
- 437 [35] Naftali Tishby and Noga Zaslavsky. Deep Learning and the Information Bottleneck Principle,
438 2015. arXiv:1503.02406.
- 439 [36] Gido M. van de Ven and Andreas S. Tolias. Three scenarios for continual learning, 2019.
440 arXiv:1904.07734.
- 441 [37] Yanzhao Wu, Ling Liu, Juhyun Bae, Ka-Ho Chow, Arun Iyengar, Calton Pu, Wenqi Wei, Lei
442 Yu, and Qi Zhang. Demystifying Learning Rate Policies for High Accuracy Training of Deep
443 Neural Networks, 2019. arXiv:1908.06477.
- 444 [38] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu.
445 Large Scale Incremental Learning. In *Proceedings of the IEEE Conference on Computer Vision
446 and Pattern Recognition*, pages 374–382, 2019.
- 447 [39] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong Learning with
448 Dynamically Expandable Networks. In *International Conference on Learning Representations*,
449 2018.

450 Checklist

- 451 1. For all authors...
- 452 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
453 contributions and scope? [Yes]
- 454 (b) Did you describe the limitations of your work? [Yes]
- 455 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 456 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
457 them? [Yes]
- 458 2. If you are including theoretical results...
- 459 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 460 (b) Did you include complete proofs of all theoretical results? [N/A]
- 461 3. If you ran experiments...

- 462 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
463 mental results (either in the supplemental material or as a URL)? [Yes]
- 464 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
465 were chosen)? [Yes]
- 466 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
467 iments multiple times)? [No] Error bars are not reported because we did not have
468 enough compute resources. The results report an average of 3 random seeds for each
469 experiment.
- 470 (d) Did you include the total amount of compute and the type of resources used (e.g., type
471 of GPUs, internal cluster, or cloud provider)? [Yes]
- 472 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 473 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 474 (b) Did you mention the license of the assets? [Yes] MIT License
- 475 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 476 (d) Did you discuss whether and how consent was obtained from people whose data you're
477 using/curating? [No]
- 478 (e) Did you discuss whether the data you are using/curating contains personally identifiable
479 information or offensive content? [No]
- 480 5. If you used crowdsourcing or conducted research with human subjects...
- 481 (a) Did you include the full text of instructions given to participants and screenshots, if
482 applicable? [N/A]
- 483 (b) Did you describe any potential participant risks, with links to Institutional Review
484 Board (IRB) approvals, if applicable? [N/A]
- 485 (c) Did you include the estimated hourly wage paid to participants and the total amount
486 spent on participant compensation? [N/A]