

RESPONSE-BASED DISTILLATION FOR INCREMENTAL OBJECT DETECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Traditional object detection are ill-equipped for incremental learning. However,
2 fine-tuning directly on a well-trained detection model with only new data will
3 leads to catastrophic forgetting. Knowledge distillation is a straightforward way
4 to mitigate catastrophic forgetting. In Incremental Object Detection (IOD), previ-
5 ous work mainly focuses on feature-level knowledge distillation, but the different
6 response of detector has not been fully explored yet. In this paper, we propose
7 a fully response-based incremental distillation method focusing on learning re-
8 sponse from detection bounding boxes and classification predictions. Firstly, our
9 method transferring category knowledge while equipping student model with the
10 ability to retain localization knowledge during incremental learning. In addition,
11 we further evaluate the qualities of all locations and provides valuable response
12 by adaptive pseudo-label selection (APS) strategies. Finally, we elucidate that
13 knowledge from different responses should be assigned with different importance
14 during incremental distillation. Extensive experiments conducted on MS COCO
15 demonstrate significant advantages of our method, which substantially narrow the
16 performance gap towards full training.

17 1 INTRODUCTION

18 In the natural world, the visual system of creatures could constantly acquire, integrate and optimize
19 knowledge. Learning mode is inherently incremental for them. In contrast, currently, the classic
20 training paradigm of the object detection model (Tian et al., 2019; Li et al., 2021b) does not have
21 such capability. Supervised object detection paradigm relies on accessing pre-defined labeled data.
22 This learning paradigm implicit assumes data distribution is fixed or stationary, while data from
23 real world is represented by continuous and dynamic data flow, whose distribution is non-stationary.
24 When the model continuously obtains knowledge from non-stationary data distribution, new knowl-
25 edge would interfere with the old one, triggering catastrophic forgetting (Goodfellow et al., 2015;
26 McCloskey & Cohen, 1989).

27 A straightforward way in incremental object detection is based on knowledge distillation (Hinton
28 et al., 2015). Peng et al. (2021) stressed that the Tower layers could reduce catastrophic forgetting
29 significantly. They implemented incremental learning on an anchor-free detector and selectively per-
30 formed distillation on non-regression outputs. In knowledge distillation for object detection where
31 incremental learning was not introduced, previous work extracted knowledge from the combined
32 distillation of different components. For example, Chen et al. (2017) and Sun et al. (2020) dis-
33 tilled all components of the detector. Nevertheless, the nature of these methods are designed using
34 feature-based knowledge distillation, fully response-based method (Gou et al., 2021) has not been
35 explored in incremental object detection yet. Besides, since different components in the detection
36 make different contributions to incremental distillation, an elaborate design for different responses
37 is essential.

38 This paper focused on a practical and challenging problem concerning incremental object detection:
39 *how to learn response from detecting bounding boxes and classification predictions*. Responses in
40 object detection contain logits together with the offset of bounding box (Gou et al., 2021). Firstly,
41 since the number of ground truth on each new image is uncertain, one of the foremost considerations
42 is that validate the object of all samples, determining which object is positive or negative and which
43 ground truth each object should regress towards. A troublesome issue is that the output of the

44 regression branch may be substantially different from that of the ground truth. Furthermore, the
 45 localization knowledge of each edge in the detection bounding boxes is also response that should be
 46 taken seriously. To sum up, we use the response on the location where teacher detector generates
 47 high-quality predictions as the ground truth to guide the student detector following the behavior of
 48 teacher on the old object. In this case, it is of great significance to use the old detector to provide
 49 valuable incremental information from detection bounding boxes and classification predictions.

50 To tackle the above problems, this paper rethinks response-based knowledge distillation method,
 51 finding that distillation at proper locations is crucial in facilitating incremental object detection.
 52 We believe that student detector can acquire high-quality knowledge from the teacher detector’s
 53 high-quality predictions. Driven by this inspiration, we proposed an incremental distillation scheme
 54 that learns specific responses from the classification head and regression head respectively. Unlike
 55 previous work, we introduce incremental localization distillation (Zheng et al., 2021) in regression
 56 response to equip student detector with the ability to learn location ambiguity during incremental
 57 learning. Besides, we propose adaptive pseudo-label selection (APS) strategies to automatically
 58 select distillation nodes based on statistical characteristics from different responses, which evaluates
 59 the qualities of all locations and provides valuable response. We alleviate catastrophic forgetting
 60 greatly and significantly narrow the gap with full training by distilling the response alone. Extensive
 61 experiments on the MS COCO dataset support our analysis and conclusion.

62 The main contributions of this work can be summarized,

- 63 1. To the best of our knowledge, this paper is first work to explore the fully response-based
 64 distillation method in incremental object detection.
- 65 2. We propose a novel distillation scheme elaborate designed for incremental detection focus-
 66 ing on detection bounding boxes and classification predictions.
- 67 3. We propose adaptive pseudo-label selection strategies to automatically select distillation
 68 nodes based on statistical characteristics from the different responses.

69 2 RELATED WORK

70 **Incremental Learning.** Catastrophic forgetting is the core challenge for incremental learning. In-
 71 cremental learning based on parameter constraints is a candidate solution for such problem, which
 72 protects the old knowledge by introducing an additional parameter-related regularization term to
 73 modify the gradient. EWC (Kirkpatrick et al., 2016) and MAS (Aljundi et al., 2018) are two typical
 74 representatives of such method. Another solution is incremental learning based on knowledge dis-
 75 tillation, as well as the topic of the study. This kind of method mainly projects old knowledge by
 76 transferring knowledge in old tasks to new tasks through knowledge distillation. LwF (Li & Hoiem,
 77 2018) is the first algorithm that introduces the concept of knowledge distillation into incremental
 78 learning, in the purpose of making predictions of the new model on new tasks similar to that of
 79 the old model and thereby protecting the old knowledge in the form of knowledge transfer. How-
 80 ever, it would cause knowledge confusion when the correlation between new and old tasks is low.
 81 iCaRL (Rebuffi et al., 2017) algorithm uses knowledge distillation to avoid excessive deterioration
 82 of knowledge in the network, while BiC (Wu et al., 2019) algorithm added a bias correction layer
 83 after the FC layer to offset the category bias of new data when using the distillation loss.

84 **Incremental Object Detection.** Compared with incremental classification, achievements on incre-
 85 mental object detection is much less. Meanwhile, the high complexity of the detection task also
 86 adds the difficulty of incremental object detection. Shmelkov et al. (2017) proposed to apply LwF
 87 to Fast RCNN detector (Girshick, 2015), which is the first work on incremental object detection.
 88 Thereafter, some researchers move this area forward. Peng et al. (2021) proposed SID approach
 89 for incremental object detection on anchor-free detector and conducted experiments on FCOS (Tian
 90 et al., 2019) and CenterNet (Zhou et al., 2019). Li et al. (2021a) studied object detection based
 91 on class-incremental learning on Faster RCNN detector with emphasis given to few-shot scenarios,
 92 which is also the focus of ONCE algorithm (Perez-Rua et al., 2020). Li et al. (2019) designed an
 93 incremental object detection system with RetinaNet detector (Lin et al., 2020) under the scenario
 94 of edge device. the latest work, Joseph et al. (2021) introduced the concept of incremental learning
 95 when defining the problems of Open World Object Detection (OWOD).

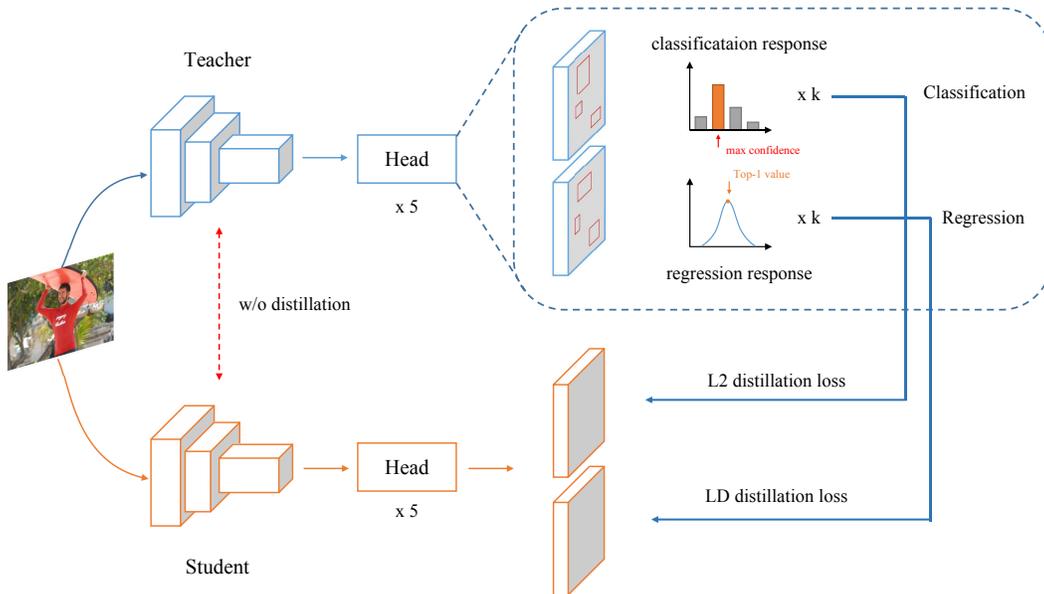


Figure 1: Overview of response-based incremental distillation.

96 **Knowledge Distillation for Object Detection.** Knowledge distillation (Bucila et al., 2006) is an
 97 effective way to transfer knowledge between models. Widely applied in image classification tasks
 98 in previous researches, knowledge distillation is now used in object detection tasks more and more
 99 frequently. Chen et al. (2017) implemented distillation in all components on Faster RCNN detector
 100 (including backbone, proposals in RPN, and head). To imitate the high-level feature response of
 101 the teacher model with the student model, Wang et al. (2019) proposed a distillation method based
 102 on fine-grained feature imitation. By synthesizing category-conditioned objects through inverse map-
 103 ping, Chawla et al. (2021) proposed a data-free knowledge distillation technology applicable for
 104 object detection, but the method would trigger dream-image. Guo et al. (2021) believing that fore-
 105 ground and background both play a unique role in object detection, proposed an object detection
 106 distillation method that could decouple foreground and background. Zheng et al. (2021) proposed
 107 a localization distillation method introducing knowledge distillation into the regression branch of
 108 the object detector, so as to enable the student network to solve the localization ambiguity in object
 109 detection as the teacher network. However, existing object detection distillation framework does not
 110 pay enough attention to the significant role of the head. In this study, we found the head has its
 111 particularly significant.

112 3 METHOD

113 3.1 OVERALL STRUCTURE

114 In general, a one-stage object detector is composed of three components: (i.) backbone for feature
 115 extraction; (ii.) neck for fusion of multi-level features; (iii.) head for classification and regression.
 116 The purpose of incremental distillation is to transfer old knowledge to the student detector, and this
 117 knowledge could be the features of the intermediate layer in the backbone or neck or the soft predic-
 118 tions in the head. Here, we incrementally learn a strong and efficient student object detector by the
 119 distillation of incremental knowledge from responses of the different heads. The overall incremental
 120 detection framework is shown in Figure 1. Firstly, knowledge distillation is applied to learn incre-
 121 mental response from the classification head and regression head of the teacher detector. Secondly,
 122 incremental localization distillation loss is also applied to enhance the localization information ex-
 123 traction ability of the student detector. Notably, the adaptive pseudo-label selection strategies are
 124 proposed to gain more meaningful incremental responses from the teacher detector, that is, selec-

125 tive calculation of the distillation loss from the pseudo label provided by the teacher detector. The
 126 overall learning target of the student detector is therefore defined as,

$$\mathcal{L}_{total} = \mathcal{L}_{model} + \lambda_1 \mathcal{L}_{dist.cls}(\mathcal{C}_S, \mathcal{C}_T) + \lambda_2 \mathcal{L}_{dist.bbox}(\mathcal{R}_S, \mathcal{R}_T) \quad (1)$$

127 where λ is the parameters that balances the weights of different loss terms. The loss term \mathcal{L}_{model}
 128 is standard loss function used in GFocal (Li et al., 2020) to train object detector for the new object
 129 class. The second loss term $\mathcal{L}_{dist.cls}$ is the L2 incremental distillation loss for classification branch.
 130 The third loss term $\mathcal{L}_{dist.bbox}$ is the incremental localization distillation loss for regression branch.
 131 In the above, we set $\lambda_1 = \lambda_2 = 1$.

132 3.2 DISTILLATION AT CLASSIFICATION-BASED RESPONSE

133 The soft predictions from the classification head contains the knowledge of various categories dis-
 134 covered by the teacher model. Through the learning of soft prediction, the student model can inherit
 135 hidden knowledge, which is intuitive for classification tasks. Let \mathcal{T} be the teacher model, we use
 136 SoftMax function to transform logits \mathcal{Z}_T in final score output, responding probability distribution
 137 \mathcal{P}_T is defined as,

$$\mathcal{P}_T = \text{SoftMax} \left(\frac{\mathcal{Z}_T}{t} \right) \quad (2)$$

138 Similarly, we define \mathcal{P}_S for the student model \mathcal{S} ,

$$\mathcal{P}_S = \text{SoftMax} \left(\frac{\mathcal{Z}_S}{t} \right) \quad (3)$$

139 where t is temperature to soften the probability distribution for \mathcal{P}_T and \mathcal{P}_S .

140 Previous works usually directly use all the prediction responses in the classification head and treat
 141 each position equally. If there is any inappropriate balance, the response generated by the back-
 142 ground category may overwhelm the response generated by the foreground category, thereby inter-
 143 fering with the retention of old knowledge. To tackle this problem, the L2 incremental distillation
 144 loss for the classification-based response is as follows,

$$\mathcal{L}_{dist.cls}(\mathcal{C}_S, \mathcal{C}_T) = \sum_{i=1}^m (\mathcal{P}_T^i - \mathcal{P}_S^i)^2 \quad (4)$$

145 where \mathcal{P}_T^i is the category response of the frozen teacher detector from m pseudo object classes
 146 using the new data, and \mathcal{P}_S^i is the category response of the student detector for the old object classes.
 147 By distilling the selected response, the student detector inherits the knowledge of the positive object
 148 category to a greater extent.

149 3.3 DISTILLATION AT REGRESSION-BASED RESPONSE

150 The bounding box response from the regression branch is also quite important for incremental detec-
 151 tion. Contrary to the discrete class information, there is a possibility that the output of the regression
 152 branch may provide a regression direction that contradicts the ground truth. That’s because, even if
 153 the image does not contain any objects of the old category, the regression branch will still predict
 154 the bounding box, although the confidence is relatively low. That poses a challenge for learning the
 155 knowledge of the old model to correctly predict the bounding box of the old object. On the other
 156 hand, in previous works, only the bounding box of a relatively high-confidence object was learned
 157 as the knowledge of the teacher detector, ignoring the localization information.

158 Benefit from the general distribution of bounding box \mathcal{B} from GFocal detector, each edge of \mathcal{B} can
 159 be represented by probability distribution through SoftMax function (Zheng et al., 2021). Further,
 160 the probability matrix of bounding box \mathcal{B} is defined as,

$$\mathcal{B} = [p_t, p_b, p_l, p_r] \in \mathbb{R}^{n \times 4} \quad (5)$$

161 Therefore, we can extract the incremental localization knowledge of bounding box \mathcal{B} from teacher
162 detector \mathcal{T} and transfer it to student detector \mathcal{S} by using KL-Divergence loss,

$$\mathcal{L}_{LD}^e = \mathcal{L}_{KL}(\mathcal{P}_{\mathcal{S}^j}, \mathcal{P}_{\mathcal{T}^j}) \quad (6)$$

163 Finally, incremental localization distillation loss for the regression-based response is defined as,

$$\mathcal{L}_{dist.bbox}(\mathcal{R}_{\mathcal{S}}, \mathcal{R}_{\mathcal{T}}) = \sum_{j=1}^J \sum_{e \in \mathcal{B}} \mathcal{L}_{LD}^e \quad (7)$$

164 where $\mathcal{R}_{\mathcal{T}^j}$ is the regression response of the frozen teacher detector from J pseudo bounding box
165 using the new object, and $\mathcal{R}_{\mathcal{S}^j}$ is the regression response of the student detector for the old bounding
166 box. Compared to only use the bounding box in previous works, incremental localization distillation
167 can provide extra localization response.

168 3.4 ADAPTIVE PSEUDO-LABEL SELECTION

169 When an incremental object detector is trained, the gap of knowledge between the teacher detector
170 and the student detector is obvious. For a new sample, it’s preferable for the teacher detector to
171 provide the high-quality knowledge, as the student detector will benefit from positive response. To
172 this end, a basic problem related to incremental object detection has been thoroughly studied: *how*
173 *to select distillation nodes as positive response*. Traditional selection strategies depend on sensitive
174 hyper-parameters such as setting confidence and Top-K. Those empirical practices in which rules
175 are fixed have such consequences that too small thresholds lead to the ignoring of some objects
176 while too large ones probably result in the introduction of negative response.

177 To solve this problem, the adaptive pseudo-label selection (APS) strategy is proposed. Algorithm
178 1 describes how the proposed strategy works for an input image. We obtain positive response from
179 the category and bounding box as distillation nodes respectively.

180 **Classification head.** The statistical characteristics of the category information are utilized to deter-
181 mine the response of classification, as described in L-2 to L-12. We first calculate the classification
182 confidence of each position. After that, we calculate the mean μ_C and standard deviation σ_C in
183 L-6 and L-7. With these statistical, the threshold τ_C is obtained in L-8. Finally, we select these
184 candidates whose confidence are greater than the threshold τ_C in L-9 to L-12.

185 **Regression head.** The statistical characteristics of the distribution information are utilized to deter-
186 mine the response of regression, as described in L-14 to L-23. For the GFocal detector, the author
187 points out that a certain and unambiguous bounding box, whose distribution is usually sharp. There-
188 fore, the Top-1 value is usually very large if the distribution is sharp. Based on these statistical
189 characteristics, the top-1 is used to measure the confidence of the bounding box. We first calculate
190 the Top-1 of each distribution. After that, we calculate the mean μ_B and the standard deviation σ_B
191 of all Top-1 in L-17 and L-18. Then, the threshold τ_B is obtained in L-19. Finally, we select these
192 candidates whose confidence are greater than the threshold τ_B in L-20 to L-23.

193 The proposed APS strategy has the following advantages: 1. guaranteeing fair selection of pseudo
194 labels of different objects. 2. using statistical characteristics of different branches to adaptively
195 select pseudo labels to provide the incremental response.

196 4 EXPERIMENTS AND DISCUSSION

197 In this section, we perform experiments on several incremental scenarios on the MS COCO dataset
198 using baseline detector GFocal to validate our method. Then, we perform ablation studies to prove
199 the effectiveness of each component of our method. Finally, we discuss a question: What are the
200 bottlenecks in our method?

Algorithm 1 Adaptive Pseudo-label Selection (APS)

Input: Unlabeled image I , image-level labels c, b , teacher detector θ'
Output: Sampled pseudo-label sets C', B'

- 1: Inference I with θ' yields the classification score C and predicted distribution B
- 2:
- 3: Classification branch:
- 4: **for** $k = 1$ to C **do**
- 5: $G_C \leftarrow confidence(C_k)$
- 6: Compute $\mu_C = mean(G_C)$
- 7: Compute $\sigma_C = std(G_C)$
- 8: Compute threshold $\tau_C = \mu_C + \sigma_C$
- 9: **for** each candidate c in C **do**
- 10: **if** $G_{C_k} \geq \tau_C$ **then**
- 11: Add candidate c to C'
- 12: **return** C'
- 13:
- 14: Regression branch:
- 15: **for** $k = 1$ to B **do**
- 16: $G_B \leftarrow Max(B_k)$
- 17: Compute $\mu_B = mean(G_B)$
- 18: Compute $\sigma_B = std(G_B)$
- 19: Compute threshold $\tau_B = \mu_B + \sigma_B$
- 20: **for** each candidate b in B **do**
- 21: **if** $G_{S_b} \geq \tau_B$ **then**
- 22: Add candidate b to B'
- 23: **return** B'

201 **Implementation Details.** We build our method on top of the GFocal detector using their public
202 implementations. The teacher and student detectors defined in our experiments are standard GFocal
203 architectures. For the GFocal detector, ResNet-50 is used as its backbone, FPN (Lin et al., 2017) is
204 used as its neck. We trained our detector to follow the same parameters described in their paper. All
205 the experiments are performed on 8 NVIDIA Tesla V100 GPU, with batch size of 8.

206 **Datasets and Evaluation Metric.** MS COCO 2017 (Chen et al., 2015) is a challenging benchmark
207 in object detection which contains 80 object classes. For experiments on the COCO dataset, we use
208 train and validation set for training and test set for testing. The standard COCO protocols are used
209 as an evaluation metric, i.e. $AP, AP_{50}, AP_{75}, AP_S, AP_M$ and AP_L .

210 **Experiment Setup for MS COCO.** The detector is trained by 12 epochs (1x mode) for each incre-
211 mental step for the MS COCO dataset. The setting is consistent for all the detectors in the different
212 scenarios. We set up experiments in the following scenarios:

- 213 • **40 + 40:** we train a base detector with the first 40 classes and then the last 40 classes are
214 learned incrementally as new object classes.
- 215 • **75 + 5:** we train a base detector with the first 75 classes and then the last 5 classes are
216 learned incrementally as new object classes.
- 217 • **Last 40 + First 40:** we specially train a base detector with the last 40 classes and then the
218 first 40 classes are learned incrementally as new object classes.

219 4.1 OVERALL PERFORMANCE

220 We reported the incremental results under the first 40 classes + last 40 classes scenario in Table 1. In
221 this scenario, we observed that if the old detector and new data were directly used to conduct fine-
222 tuning process, then the AP dropped to 17.8% as compared to the 40.2% in full data training. This
223 is because the fine-tuning made the detector’s memory of old object classes close to 0, resulting in
224 catastrophic forgetting (ref to Figure 2(b)). Our method far outperformed fine-tuning across various
225 IoUs evaluation criteria from 0.5 to 0.95. The experimental results show that when IoU is 0.5, 0.75

Table 1: Incremental results based on GFocal detector on COCO benchmark under first 40 classes + last 40 classes. (“ Δ ” represents an improvement over Catastrophic Forgetting. “ ∇ ” represents the gap with Upper Bound.)

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Upper Bound	40.2	58.3	43.6	23.2	44.1	52.2
Catastrophic Forgetting	17.8	25.9	19.3	8.3	19.2	24.6
LwF (Li & Hoiem, 2018)	17.2($\Delta - 0.6/\nabla 23.0$)	25.4	18.6	7.9	18.4	24.3
RILOD (Li et al., 2019)	29.9($\Delta 12.1/\nabla 10.3$)	45.0	32.0	15.8	33.0	40.5
SID (Peng et al., 2021)	34.0($\Delta 16.2/\nabla 6.2$)	51.4	36.3	18.4	38.4	44.9
Ours	36.9 ($\Delta 19.1/\nabla 3.2$)	54.5	39.6	21.3	40.4	47.5

Table 2: Incremental results based on GFocal detector on COCO benchmark under last 40 classes + first 40 classes. (“ Δ ” represents an improvement over Catastrophic Forgetting. “ ∇ ” represents the gap with Upper Bound.)

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Upper Bound	40.2	58.3	43.6	23.2	44.1	52.2
Catastrophic Forgetting	22.6	32.7	24.2	15.1	25.0	27.6
LwF (Li & Hoiem, 2018)	20.5($\Delta - 2.1/\nabla 19.7$)	29.9	22.1	13.0	22.5	25.3
RILOD (Li et al., 2019)	34.1($\Delta 11.5/\nabla 6.1$)	51.1	36.8	19.1	38.0	43.9
SID (Peng et al., 2021)	33.5 ($\Delta 10.9/\nabla 6.7$)	50.9	36.3	19.0	37.7	43.0
Ours	37.5 ($\Delta 14.9/\nabla 2.8$)	55.1	40.4	21.3	41.1	48.2

226 and 0.95, the AP improves by 19.1%, 28.6% and 20.3%, respectively. This indicates that our method
 227 can well address catastrophic forgetting. Notably, even compared with the full data training where
 228 the entire dataset was used, our method only had a gap of 3.2%. This indicates that the student
 229 detector maintained a good memory of the old objects while learning new objects. To put it more
 230 intuitively, we visualized the incremental results of all object classes, as shown in Figure 2. The blue
 231 column denotes the AP of the first 40 classes, while the orange column denotes the AP of the last
 232 40 classes. As can be seen, our method has produced significant outcomes. In Figure 3, we further
 233 visualized the AP of all objects of the first 40 classes and the last 40 classes.

234 Considering the long-tail problem of the COCO dataset, we particularly configured an incremental
 235 experiment under the last 40 classes + first 40 classes scenario. In this scenario, the first 40 classes
 236 object contain more memories that should be retained, which means that more incremental responses
 237 can be obtained. As can be seen from Table 2, the incremental performance of our method has been
 238 further improved, with the gap against full data training reduced to 2.8% and the improvement on
 239 catastrophic forgetting increased to 14.9%. This also validates our inference that the method we
 240 propose benefits from more incremental responses.

241 In addition, we also compared our method with LwF, RILOD, and SID. Both Table 1 and Table 2
 242 show that although LwF works well in incremental classification, it is even lower AP than directly
 243 fine-tuning in detection tasks. To a fair comparison with RILOD and SID, we replicated them
 244 based on GFocal detector. For RILOD, we completely followed their method. For SID, we used
 245 the component with the greatest improvement proposed by the authors. Both tables show that the
 246 improvement of our method to catastrophic forgetting is outstanding.

247 4.2 ABLATION STUDY

248 As shown in Table 3, we validated the effectiveness of different components of the proposed method
 249 on the COCO benchmark to highlight our improvement in performance. “all cls + all reg” denotes
 250 that responses from both the classification branch and regression branch are treated equally in the
 251 incremental distillation, which is also our baseline performance. “all cls” denotes that only classi-
 252 fication responses in the incremental distillation process are treated equally. “all reg” denotes that
 253 only regression responses in the incremental distillation process are treated equally. “cls + APS”

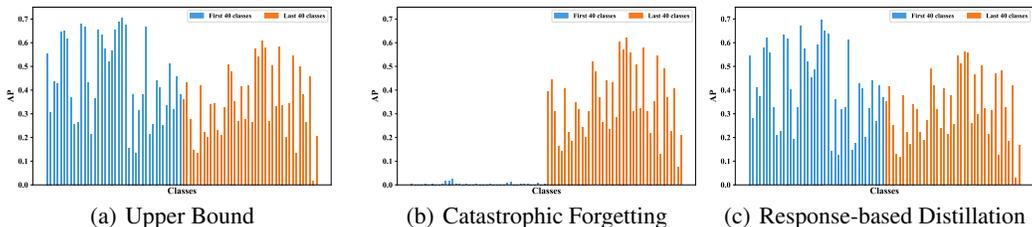


Figure 2: AP of Per-class among different learning schemes. (a) Detector is trained with all data.(b) Student detector is finetuned with new classes.(c) Student detector is distilled via different response.

254 denotes that the APS strategy is employed to conduct incremental distillation over classification re-
 255 sponses, as shown in Equation 4. “cls + reg +APS” denotes that responses based on regression are
 256 also used, as shown in Equation 7. In Table 3, separately distillation all responses from classifica-
 257 tion and regression, obtained 23.8% and 13.0% of AP. When only all responses from the regression
 258 branch are used, AP is even lower than the fine-tuning performance, which also supports our as-
 259 sumption stated in the introduction section. Comparatively, the direct incremental distillation of
 260 all responses from classification and regression branches obtains 31.5% of AP. By utilizing APS
 261 to decouple classification responses, the student detector obtained higher results. Our decoupling
 262 proposal can improve the result from 31.5% of AP to 33.2%. The incremental distillation process
 263 further utilized the APS strategy to decouple regression responses, obtaining 36.9% of AP on the
 264 COCO benchmark, a 5.4% improvement compared with the baseline performance. All these results
 265 clearly point to the advantageous performance of our method.

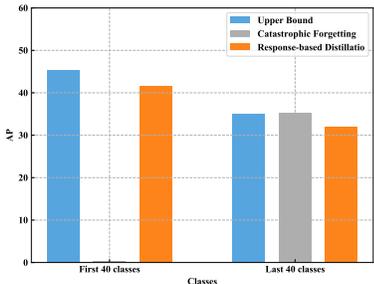


Figure 3: First 40 classes vs. Last 40 classes.

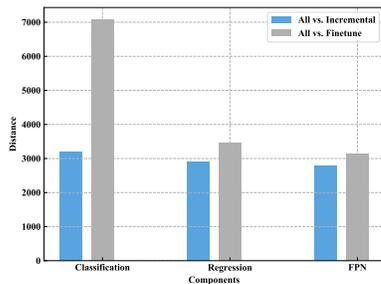


Figure 4: L2 distance analysis.

266 4.3 DISCUSSION

267 In this section, we present further insights into response-based incremental distillation. We reveal the
 268 contribution of different components for distillation detection and discuss the impact of incremental
 269 response in the head.

270 **Distance between different components.** We calculate the feature distance between different com-
 271 ponents to illustrate why response-based distillation can attain higher performance compared to
 272 other components. We randomly choose 10 images from COCO minival and calculate the L2 dis-
 273 tance of features in different components of different training strategies. As shown in Figure 4,
 274 “All” denotes that the detector with full data training; ‘Incremental’ denotes that the detector with
 275 incremental data training; “Finetune” denotes that the detector with finetuning training. Distilling
 276 student detector via classification-based and regression-based incremental response in the head can
 277 substantially narrow the distances with upper bound. However, neither the L2 distance between “All
 278 vs. Incremental” and “All vs. Finetune” improves significantly in the FPN representing the feature-
 279 based distillation. This also supports our assumption that different response from the head has its
 280 particularly significant, especially classification response.

281 **Incremental response helps both learning and generalization.** We notice that the incremental
 282 response from the head can provide an effective guidance to avoid catastrophic forgetting problems.

Table 3: Ablation study based on GFocal detector using the COCO benchmark under first 40 classes + last 40 classes. ("Δ" represents an improvement over Catastrophic Forgetting. "∇" represents the gap with Upper Bound.)

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Upper Bound	40.2	58.3	43.6	23.2	44.1	52.2
Catastrophic Forgetting	17.8	25.9	19.3	8.3	19.2	24.6
KD:all cls + all reg	31.5(Δ13.7/∇8.7)	48.3	33.4	17.7	35.3	41.3
KD:all cls	23.8(Δ10.1/∇16.4)	36.6	24.9	11.8	27.2	32.9
KD:all reg	13.0(Δ - 4.8/∇27.2)	21.1	13.4	5.0	14.7	18.6
KD:cls + APS	33.2(Δ15.4/∇7.0)	51.2	35.2	18.5	37.8	43.8
KD:cls + reg + APS	36.9(Δ19.1/∇3.2)	54.5	39.6	21.3	40.4	47.5

Table 4: Incremental results based on GFocal detector on COCO benchmark under first 75 classes + last 5 classes. ("Δ" represents an improvement over Catastrophic Forgetting.)

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Catastrophic Forgetting	3.8	5.9	3.8	1.9	5.3	6.5
All response	32.5 (Δ28.7)	48.9	34.7	18.3	35.9	41.1
Adaptive response	28.3 (Δ24.5)	42.4	30.3	15.0	31.5	37.0

283 Thereby, the student detector achieves noticeable improvement in different scenarios. In Table 4,
 284 our method can still learn new object classes without forgetting old ones even with a little data. But,
 285 due to the insufficient incremental response provided in the +5 classes scenario, our method did
 286 not achieve a more competitive AP. However, our method still contributes to generalization. In this
 287 case, we can degrade the adaptive response to all responses in exchange for a better compromise.
 288 Comparatively, when sufficient incremental responses emerge, our method is easy to achieve (near)
 289 perfect AP.

290 5 CONCLUSION

291 In this paper, we design an entirely response-based incremental object detection paradigm. This
 292 method uses only the detection head to achieve incremental detection, which significantly alleviates
 293 catastrophic forgetting. We innovatively learn responses from detection bounding boxes and classi-
 294 fication predictions, and specifically introduce incremental localization distillation in the regression
 295 response. Second, the adaptive selection technique is designed to provide a fair incremental response
 296 in the different heads. Extensive experiments validate the effectiveness of our method. Finally, our
 297 empirical analysis reveals the contribution of different responses and components in incremental
 298 detection, which could provide insights to further advancement in the field.

299 REFERENCES

- 300 Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars.
 301 Memory aware synapses: Learning what (not) to forget. In Vittorio Ferrari, Martial Hebert,
 302 Cristian Sminchisescu, and Yair Weiss (eds.), *Computer Vision - ECCV 2018 - 15th European*
 303 *Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part III*, volume 11207 of
 304 *Lecture Notes in Computer Science*, pp. 144–161. Springer, 2018.
- 305 Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Tina Eliassi-
 306 Rad, Lyle H. Ungar, Mark Craven, and Dimitrios Gunopulos (eds.), *Proceedings of the Twelfth*
 307 *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadel-*
 308 *phia, PA, USA, August 20-23, 2006*, pp. 535–541. ACM, 2006.

- 309 Akshay Chawla, Hongxu Yin, Pavlo Molchanov, and Jose M. Alvarez. Data-free knowledge dis-
310 tillation for object detection. In *IEEE Winter Conference on Applications of Computer Vision*,
311 *WACV 2021, Waikoloa, HI, USA, January 3-8, 2021*, pp. 3288–3297. IEEE, 2021.
- 312 Guobin Chen, Wongun Choi, Xiang Yu, Tony X. Han, and Manmohan Chandraker. Learning effi-
313 cient object detection models with knowledge distillation. In Isabelle Guyon, Ulrike von Luxburg,
314 Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.),
315 *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Informa-*
316 *tion Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 742–751, 2017.
- 317 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and
318 C. Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server. *CoRR*,
319 abs/1504.00325, 2015.
- 320 Ross Girshick. Fast r-cnn, 2015.
- 321 Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical
322 investigation of catastrophic forgetting in gradient-based neural networks, 2015.
- 323 Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation: A
324 survey. *Int. J. Comput. Vis.*, 129(6):1789–1819, 2021.
- 325 Jianyuan Guo, Kai Han, Yunhe Wang, Han Wu, Xinghao Chen, Chunjing Xu, and Chang Xu. Dis-
326 tilling object detectors via decoupled features. In *IEEE Conference on Computer Vision and*
327 *Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 2154–2164. Computer Vision
328 Foundation / IEEE, 2021.
- 329 Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network.
330 *CoRR*, abs/1503.02531, 2015.
- 331 K. J. Joseph, Salman H. Khan, Fahad Shahbaz Khan, and Vineeth N. Balasubramanian. Towards
332 open world object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*,
333 *CVPR 2021, virtual, June 19-25, 2021*, pp. 5830–5840. Computer Vision Foundation / IEEE,
334 2021.
- 335 James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, An-
336 dreei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis
337 Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic for-
338 getting in neural networks. *CoRR*, abs/1612.00796, 2016.
- 339 Dawei Li, Serafettin Tasci, Shalini Ghosh, Jingwen Zhu, Junting Zhang, and Larry P. Heck. RILOD:
340 near real-time incremental learning for object detection at the edge. In Songqing Chen, Ryo-
341 kichi Onishi, Ganesh Ananthanarayanan, and Qun Li (eds.), *Proceedings of the 4th ACM/IEEE*
342 *Symposium on Edge Computing, SEC 2019, Arlington, Virginia, USA, November 7-9, 2019*, pp.
343 113–126. ACM, 2019.
- 344 Pengyang Li, Yanan Li, and Donghui Wang. Class-incremental few-shot object detection, 2021a.
- 345 Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang.
346 Generalized focal loss: Learning qualified and distributed bounding boxes for dense object de-
347 tection. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and
348 Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Con-*
349 *ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020,*
350 *virtual, 2020*.
- 351 Xiang Li, Wenhai Wang, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss V2:
352 learning reliable localization quality estimation for dense object detection. In *IEEE Conference*
353 *on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 11632–
354 11641. Computer Vision Foundation / IEEE, 2021b.
- 355 Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE Trans. Pattern Anal. Mach.*
356 *Intell.*, 40(12):2935–2947, 2018.

- 357 Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Be-
358 longie. Feature pyramid networks for object detection. In *2017 IEEE Conference on Computer
359 Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 936–944.
360 IEEE Computer Society, 2017.
- 361 Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
362 object detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(2):318–327, 2020.
- 363 M. McCloskey and N. J. Cohen. Catastrophic interference in connectionist networks: The sequential
364 learning problem. *Psychology of Learning and Motivation*, 24:109–165, 1989.
- 365 Can Peng, Kun Zhao, Sam Maksoud, Meng Li, and Brian C. Lovell. SID: incremental learning
366 for anchor-free object detection via selective and inter-related distillation. *Comput. Vis. Image
367 Underst.*, 210:103229, 2021.
- 368 Juan-Manuel Perez-Rua, Xiatian Zhu, Timothy Hospedales, and Tao Xiang. Incremental few-shot
369 object detection, 2020.
- 370 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl:
371 Incremental classifier and representation learning. In *2017 IEEE Conference on Computer Vision
372 and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 5533–5542.
373 IEEE Computer Society, 2017.
- 374 Konstantin Shmelkov, Cordelia Schmid, and Karteek Alahari. Incremental learning of object de-
375 tectors without catastrophic forgetting. In *IEEE International Conference on Computer Vision,
376 ICCV 2017, Venice, Italy, October 22-29, 2017*, pp. 3420–3429. IEEE Computer Society, 2017.
- 377 Ruoyu Sun, Fuhui Tang, Xiaopeng Zhang, Hongkai Xiong, and Qi Tian. Distilling object detectors
378 with task adaptive regularization. *CoRR*, abs/2006.13108, 2020.
- 379 Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: fully convolutional one-stage object
380 detection. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul,
381 Korea (South), October 27 - November 2, 2019*, pp. 9626–9635. IEEE, 2019.
- 382 Tao Wang, Li Yuan, Xiaopeng Zhang, and Jiashi Feng. Distilling object detectors with fine-grained
383 feature imitation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019,
384 Long Beach, CA, USA, June 16-20, 2019*, pp. 4933–4942. Computer Vision Foundation / IEEE,
385 2019.
- 386 Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu.
387 Large scale incremental learning. In *IEEE Conference on Computer Vision and Pattern Recog-
388 nition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 374–382. Computer Vision
389 Foundation / IEEE, 2019.
- 390 Zhaohui Zheng, Rongguang Ye, Ping Wang, Jun Wang, Dongwei Ren, and Wangmeng Zuo. Local-
391 ization distillation for object detection, 2021.
- 392 Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. In *arXiv preprint
393 arXiv:1904.07850*, 2019.