



Figure 1: Distribution of embedded negative representatives using t-SNE; The selected hard negative proposals from support images (3-shot, i.e. 3 buses) of bus. [Zoom In]

Shot	1	2	3	5	10
Novel Set 1	26.1	32.9	34.4	38.6	41.3
Novel Set 2	17.2	22.1	23.4	28.3	35.8
Novel Set 3	27.5	31.1	31.5	34.4	37.2

Table 1: Performance of RepMet on PASCAL VOC novel sets.

2 We thank all reviewers for their constructive and valuable comments.

3 **(R1) Baseline: training a simple object detector on base classes and fine-tuning its head on novel classes.**

4 A1: Indeed, the proposed baseline has already been implemented in RepMet[17] (see its Page 6 and Table 3 for
5 “baseline-FT”). This baseline on ImageNet-LOC performs much inferior than RepMet/NP-RepMet. We will add this.

6 **(R1) The backbone and pre-training regime of the network should be controlled for comparison.**

7 A2: We agree ! We mentioned this in Sec. 4.2: we use ResNet-101 as the backbone (L213); its weights initialization
8 follows [17] for ImageNet-LOC and [16] for PASCAL VOC (L215); weights of other modules (e.g. FPN, RPN)
9 are randomly initialized (L216). That is, we pre-train the backbone on COCO to fairly compare with RepMet on
10 ImageNet-LOC; pre-train it on ImageNet to align with [16, 36,38] on PASCAL VOC. We are also very careful to follow
11 train/test details (class splits (L211-212), support sets (L224-226), inference scheme (L311-319)) in [17] and [16].

12 **(R2) [16,36,38] do not rely on representatives and hence do not suffer from the problem addressed in this work.**

13 A3: The key observation of this work is that hard negative information within support images is not carefully exploited
14 in previous works [16,17,36,38]. They simply extracts positive proposals w.r.t. ground truth from support images while
15 negatives proposals containing e.g. partial objects or ambiguous surroundings are not considered (L37). These negatives
16 however are important false positives in object detection (L40-43). Although we build our work on RepMet, our idea of
17 restoring negative information should be essential and beneficial to many few-shot detection works. Judging from the
18 results, NP-RepMet drastically outperforms other few-shot methods [16, 36, 38] on PASCAL VOC (e.g. up to 18%).

19 **(R2) It is less intuitive that negative representatives can be used in a query image with different background.**

20 A4: First, we emphasize that hard negative proposals mined in this work are mainly focused on proposals containing
21 partial, occluded object or entire object with massive surroundings. To realize this, we propose to specifically choose
22 them via IoU thresh (0.2, 0.3) and Cluster-Min. Fig. 1 illustrates the selected hard negatives from the support set of bus.
23 We also use t-SNE to draw the embedded negative representatives from the support sets of multiple classes: different
24 classes are clearly distinguished from each other; given a query, its hard negative proposals (e.g. partial object) for
25 certain class can be easily filtered out by comparing to those negative representatives from the support set. We provide
26 ablation study in Sec 4.3 and Table 2 to show that these negative representatives contribute substantially to NP-RepMet.

27 **(R2) Discussion on whether the pre-trained data include samples of the novel classes.**

28 A5: Our pre-training follows previous works (see R1-A2). Class overlap between the pre-trained data and novel data
29 may exist, noting that ImageNet contains 1000 classes. Yet, we argue 1) the class distribution on pretrained data is very
30 different to that on novel data; 2) detection-related modules in NP-RepMet are randomly initialized; 3) as long as base
31 class data are enough, we do not think using a pretrained model would cast a strong impact on the final performance.

32 **(R2) More refs.** A6: Thanks! We will carefully discuss them! Particularly, the NeurIPS paper indeed has very different
33 train/test settings (class splits, support sets, queries, etc.) to ours and [16,38,36], so the results are not comparable.

34 **(R3) The proposed method is not restricted to few-shot regime, it can be used to train a standard detector**

35 A7: Yes! This can be validated from results in Table 1, 4 in the paper and Table 1 in the supp, where we train the
36 detector on the train set of base classes and evaluate it on the test set of base classes to show it maintains a good
37 accuracy; some comparison to standard detectors can also be found in Table 3 in [38]; notice our NP-RepMet is inferred
38 for all the classes together on PASCAL VOC (L315-319), like a standard detector. On the other hand, if R3 meant
39 fine-tuning NP-RepMet with novel classes instead of meta-testing as in the paper, yes, we can further fine-tune it and
40 our results can be improved to 71.3, 78.6, and 81.7 for 1, 5, and 10-shot on ImageNet-LOC unseen (novel) classes.

41 **(R3) Performance of RepMet on PASCAL VOC.**

42 A8: RepMet did not report results on PASCAL VOC novel sets. We can reproduce it under the same setting with ours,
43 results are in Table 1: RepMet performs clearly inferior to our NP-RepMet (Table 3 in the paper).

44 **(R4) The authors did not report the results of [16] for 3-shot base classes.**

45 A9: It is because [16] did not report it in their paper. One might find it being reported in [38] as 64.8 which is lower
46 than our 66.6; but we also found that [38] reports the 10-shot result of [16] as 63.6 instead of 69.7. Overall, the base
47 class results among NP-RepMet and [16,38] are competitive to each other and all maintain a good accuracy.

48 **(R4) IoU hyper-parameter.** A10: Taking IoU>0.7 for positive samples is a common practice for many detection
49 works, e.g. Faster R-CNN and RepMet; we simply follow this. We experimented with IoU>0.5, the result is 68.3 (v.s.
50 68.5 for IoU>0.7, very close) on the 1-shot case of ImageNet-LOC. As for negative sample selection, we perform
51 ablation study for different IoU in Table 2 (right bottom), this hyper-parameter is quite robust within certain range.

52 **(R4) No comparison to other hard negative mining methods, though they have not applied to few-shot detection.**

53 A11: One representative hard negative mining strategy OHEM [30] is indeed used in RepMet (L101) where class
54 representatives from different classes are considered as negatives to each other. We actually adopt this in our loss
55 function (Eq 1 and 2). We bootstrap the classifier with negative information both within and across images and classes.
56 More importantly, we show that the idea of mining hard negatives within the same image of positives is essential to
57 few-shot detection. We will add more discussion about hard negative mining in literature in the revised version.