

1 We thank all reviewers for their insightful and constructive comments.

2 **Response to R#2 Q2.1:...**The main weakness is novelty ... the novel contributions are in data scheduler and
3 **negative sample strategy, which in my opinion do not constitute a contribution worth publication at NeurIPS...**

4 While R#2 well noticed our new techniques, our contribution of *bringing a simple and effective line back to the sight*
5 *for unsupervised learning* may be overlooked. It breaks the inherent belief that parametric instance classification has
6 intrinsic limitation [32]. Such knowledge is new to the community and is better valued independent of tech contribution.
7 As for the two novel techniques, we respect the reviewer’s personal opinion, but we would greatly appreciate if the
8 reviewer could also take the following facts into consideration: 1) “a simple baseline” (by R#1, R#3, R#4) which
9 does not “require special handling to avoid data leakage and can be easily implemented” (by R#4); 2) simple and
10 effective solutions to address *crucial* issues in parametric instance classification (by R#1, R#3, R#4): the sliding window
11 scheduler addresses the extremely infrequent instance visiting issue and the *corrected updating technique* enables
12 applying the PIC framework to unlimited data scale.

13 **Q2.2: ...Another weakness of this paper is on the empirical side...when SimCLR is trained for 1000 epochs it**
14 **outperforms the proposed method...** SimCLR performs worse than our approach in longer training settings (69.3%
15 vs 70.8% in Table 6). Nevertheless, the main goal of this work is not to achieve higher accuracy number, but more
16 importantly to recall another line of unsupervised learning, which could increase the diversity of research.

17 **Q2.3: ...notation details not properly introduced...** We would check it carefully and improve the writing accordingly.

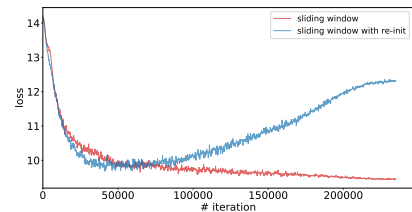
18 **Q2.4: ...missing references on cross-level relationship modeling...** The listed papers perform unsupervised learning
19 via pairwise similarity learning. We will add discussion of these papers in the revision.

20 **Response to R#3 Q3.1: ... make it clear in introduction that the proposed solution is a "fix" of [12] ...** Thanks for
21 the suggestion. We will clarify it in the introduction like “our paper is basically a revisit of [12]” (L184-188).

22 **Q3.2: ...does the improvement in training comes only from being able to "learn well" a single instance class**
23 **before moving to another one? ...** Yes, the improvement in training should come from being able to "learn well" a
24 single instance class because the opposite effect could not help training.

25 **Q3.3: How about the opposite effect like forgetting this instance class [1*]**

26 ... To answer this question, we consider an experiment by re-initializing all
27 the instance class weights which have not been seen for a long time (12.5%
28 in each window). The figure illustrates its loss curve compared to that of the
29 original sliding-window based training. At early steps, explicitly forgetting the
30 instances performs similarly well with the original sliding-window based
31 method, indicating the forgetting issue does not affect optimization. This



32 is probably because the large learning rate enables learning classification weights well even when given random
33 initialization for the long-period-not-visited instances. At later steps, the standard sliding-window has steadily reduced
34 loss while explicitly forgetting the instances results in much poorer performance. This indicates that the weights for
35 long-period-not-visited instances may still effect well, although their weights are not updated for a long period, probably
36 because of the small learning rates. On the whole, the forgetting of the *minority* may not be a serious issue, and the
37 benefit of well learning the *majority* weights overcomes such disadvantages. We will add discussion in the revision.

38 **Q3.4: ...it would be much more valuable if we could see a comparison of visualizations and statistics (Fig. 3)**
39 **with other methods such as MoCo and SimCLR...** Thanks for the suggestion. We conducted similar experiments for
40 MoCo and SimCLR and their behavior is similar as PIC. Actually, the goal of this section is not to compare PIC with
41 MoCo/SimCLR, but to understand PIC itself, where the conclusion may be generalized to other instance discrimination
42 based methods. Note we also perform experiments distinct to PIC, shown in Appendix Sec. F and Fig. 4, which build a
43 smooth transition between PIC and supervised method.

44 **Q3.5: ...how to count epochs for sliding window sampler...** As stated in the footnote of page 5, we use the term
45 “epoch” to indicate the equivalent training length with that of epoch-based scheduler, to simplify the description.

46 **Response to R#4 Q4.1: ...results on larger backbone (e.g., deeper or wider)...** ResNet-50 (2×) with 200-epoch
47 pre-training achieves 71.2% top-1 accuracy, which is 3.9% better than that of the standard ResNet-50 (1×, 67.3%).

48 **Q4.2: ...if projection head is dropped...** Yes, we follow [5,15,22] to drop the projection head in downstream tasks.

49 **Q4.3: ...comment on why cosine softmax brings such a significant improvement in the current setting...** The only
50 difference between cosine softmax loss and standard softmax loss is that the standard softmax loss accounts vector
51 *magnitudes* for similarity computation, in addition to the *angles* between two vectors used by cosine softmax. We
52 think the significant accuracy drop by standard softmax loss is due to the significantly worse generalization ability of
53 vector *magnitudes* than *angles*. In fact, similar behavior has been widely observed in numerous applications of metric
54 learning, including face recognition [9, 30], person re-ID (“Deep Cosine Metric Learning for Person Re-Identification ,
55 WACV-18”) and few-shot learning (“MatchingNet-NeurIPS16”, “Improving Generalization via Scalable Neighborhood
56 Component Analysis, ECCV-18”, “A Closer Look at Few-shot Classification, ICLR-20”). To our knowledge, strict
57 proof of such behavior is still an open question.