

1 We thank the reviewers for their work, comments, relevant questions, and generally positive feedback, addressed in turn:

2 **[R1] Sprechmann et al [S]**, thank you for this reference we will add it. Key differences: [S] do retrieval at inference
3 time (vs training), our method is *not* retrieving nearest neighbors, but uses the more sophisticated MIR criteria. The
4 experiments in [S] are more limited as the retrieval is done in pixel space ([S] Sec. 4.1 paragraph 3)

5 **case of dissimilar tasks** We are not certain we understood this criticism correctly. For clarification, the only nearest
6 neighbor lookup is done in the hybrid method. We employ it to find the training datapoint that are the closest (in latent
7 space) to the ones retrieved via MIR optimization. If the MIR criteria retrieves diverse latent codes, then the nearest
8 neighbor lookup will find diverse samples as well. In Permuted MNIST, samples have very different appearance and we
9 were still able to improve over compared methods.

10 **diversity** We use a diversity penalty (L113-115) in Generative MIR. L246-250 and Figure 5 further studies the effects
11 of the diversity penalty. In ER-MIR, diversity is enforced via sampling prior to applying the criterion (L102-104).

12 **more challenging data** we note these datasets are also used in the related work on this challenging online continual
13 learning with shared classifier setting. We now extend our ER-MIR experiments to Mini-ImageNet split. We train on 17
14 tasks as in [6] but with shared-head. Over 20 runs we obtain an accuracy of $26.4\% \pm 0.6$ vs ER accuracy of $25.5\% \pm 0.7$
15 and a substantial gain in forgetting ($19.1\% \pm 0.8$ vs $23.5\% \pm 1.2$ for ER baseline).

16 **[R2] additional compute of a separate virtual gradient step. comparison with baselines with twice as many**
17 **gradient steps. overhead/memory of this approach** For Generative MIR the number of online updates is a hyper-
18 parameter we searched for both GEN-ER baseline and our GEN-MIR (see Appendix B.2), thus in both cases adding
19 more gradient steps would not help. Indeed in the online setting there is not a clear correlation between more iterations
20 on the incoming data (and buffer) and performance, as more iterations leads to more forgetting and overfitting. We
21 add an ablation on Split MNIST where we count the virtual update as an iteration i.e. the random baseline is allowed
22 2x more real updates \rightarrow 2 iterations: GEN 57.4% / GEN-MIR 82.1% , 10 iterations: GEN 76.8% / GEN-MIR 83.3% ,
23 100 iterations GEN 70.7% / GEN-MIR: 69.7% . Note that the regular GEN is much more sensitive to using very few
24 iterations. For ER-MIR experiments only 1 update is done in our experiments except in Table 3, observe there that even
25 with 5 gradient steps for ER, ER-MIR with 1 is still superior. We emphasize our work's aim was to determine if the
26 non-uniform sampling strategy works versus uniform, which we have shown it generally is, a critical first step to future
27 work that can find more efficient approximation of the criteria. Notably existing continual learning methods can indeed
28 be computationally very expensive compared to standard training methods. For example GEM can be 10x slower than
29 ER and other methods[6]. Our un-optimized ER-MIR implementation (for CIFAR-10) is approximately 3x slower in
30 wall clock time than regular ER. In terms of memory consumption it is the same as ER with equivalent buffer.

31 **Why is hybrid approach in the appendix?** The primary results of the hybrid experiments are indeed in the main
32 paper (see Fig 6 and L281-285). Due to space constraints we put the (sizable) algorithm block of the hybrid method as
33 well as the ablation study in the appendix.

34 **"Performance of baselines seems poorer than related work"** Our experiments focus on the online setting (e.g.
35 [3,6,26])with a shared head classifier (see [3,9]). The accuracy for DGR you refer to is for offline and multi-headed
36 evaluation aka at test time the classification only chooses between 2 categories versus 10 and is thus a very different
37 setting. Note that our baselines are similar to those reported in [3] which also considers online and shared classifier.
38 Details of all setups are given in experiments and appendix, code is included in supplementary material and will be
39 further extended in release.

40 **More external comparisons (suggested VCL).** GEM and ER are the state of the arts in the online continual learning
41 case. For EWC [1] and [2] report extremely poor performance in this setting compared GEM/ER. We have added an
42 additional comparison to VCL. Using the official VCL code, and the same experimental setup as ours (online, single
43 head, 1k samples per task, 50 buffer slots per class) we ran VCL on Split MNIST and Permuted MNIST. VCL with
44 buffer gives (87.2% , 65.9%) and without (73.9% , 64.4%), on Split and Permuted MNIST respectively. These results fall
45 below the MIR performance (87.6% , 80.1%) by 0.4% and 14.2% . We also tried VCL with 2 gradients steps, however
46 this did not help performance. Note the buffer of VCL is used after training with an additional offline training steps
47 on it. Before any prediction step during learning, training on the buffer has to be performed which violates the online
48 continual learning setting we consider in this paper. Moreover, VCL is orthogonal to MIR, and both can be used jointly.

49 **Hyperparameter sweeps for baseline, checklist** We use a validation set as described in Sec 4 and Appendix for
50 baseline and our method. In terms of learning rates ranges are similar to those used in other works [5,7]. For compute
51 we utilized a single GPU in all experiments. We reiterate that all methods were given the same number of trial runs.

52 **[R3] Theory** Unfortunately there is very little existing tools for theoretical understanding of the continual learning
53 methods (and especially in the typical non-convex setting) much less the online counterparts to provide a basis for
54 analyzing MIR.

55 **Clarity** We will correct the ambiguous notation/text you note.

56 **Longer sequence** See response to R1 for new comparisons on Mini-ImageNet, we note the online and shared-head
57 setting is very challenging and longer sequences lead to extreme forgetting even with strong baselines like ER.