

1 We thank the reviewers for the valuable feedback. There
 2 were two common concerns: lack of a complex benchmark
 3 and unclear terminology in the experiments. We address
 4 these first and then follow with responses to each reviewer.

5 **Other Benchmarks:** We didn't run experiments on a more
 6 complex dataset because even on split-omniglot, existing
 7 continual learning methods perform extremely poorly. It
 8 is fair question if the conclusions extend to other datasets.
 9 We therefore ran our method on Mini-Imagenet and report
 10 the results in Figure 1. We incrementally learn a 20-way
 11 classifier using 30 samples per class for both train and test.
 12 The results support the same conclusions. Note that we go
 13 over the training trajectory *only once*, one sample at a time.

14 We will include these results to provide further evidence for MRCL and that the strategy scales to more complex settings.

15 **To clarify the experiment protocol** we have decided to use updated terminology in the paper. Our method is divided into two
 16 phases : (1) meta-training phase and (2) meta-testing phase. The meta-training phase involves optimizing the OML objective
 17 for learning a representation and the meta-testing phase involves training the TLN on a highly correlated trajectory *in a single*
 18 *pass*. There is no overlap between data used in meta-training and meta-testing. All results are reported in the meta-testing
 19 phase; we also do not use IID sampling or multiple epochs for MRCL in any of the reported results. The first figure of
 20 Figure 4 is very meaningful because it shows MRCL—which learns incrementally in class order—is almost as effective as the
 21 Oracle—which learns using IID data. This highlights how much a representation trained for online updating can help mitigate
 22 interference.

23 **Reviewer 1: Writing and clarity:** Thank you for pointing out issues with the writing and L73; we will fix those and move
 24 the algorithm to the main paper. In Appendix L379-380, we meant that a meta-learned initialization alone can not solve the
 25 interference problem; it is important to transform the input into a representation with non-interfering solution manifolds.

26 **Improvements: ... increase your score? (1)** The intuition behind MRCL is that instead of using sparsity as a proxy for good
 27 representations for continual learning, we directly measure interference caused by highly correlated updates over a finite
 28 horizon, and minimize this interference to learn a representation. We assume that a representation that minimizes interference
 29 for k correlated updates would also reduce interference in the long run. In the incremental sine experiment, k is actually equal
 30 to length of the complete trajectory whereas for Omniglot, k is much smaller than the complete trajectory. Empirical results
 31 support that in both cases, OML can recover a good non-interfering representation. One explanation is that as long as k is large
 32 enough to cause measurable interference, minimizing OML will result in a good representation.

33 **(3):** It's not clear how to compare our method with the three suggested approaches. One of the three approaches, TADAM,
 34 is specific to few-shot-learning - a different problem setting than ours. The remaining two approaches improve on gradient
 35 based meta-learning in general and **are complementary to our work** i.e. they can be combined with our objective function to
 36 potentially further improve the results. To better clarify this, we will extend the related work section of our paper to explain
 37 why a comparison with these approaches is tangential to our contributions.

38 **Reviewer 2: Quality and Clarity: (1)** Yes. The model is trained on the meta-training set using iid sampling, and online
 39 learning is done on meta-test set in a single pass. **(2)** Random batch is sampled from the entire meta-training set. Meta-training
 40 involves revisiting the data, whereas training during meta-testing involves a single pass through data. **(3)** We agree that since
 41 class label is the same as class id, class id is implicitly available. However, our method does not exploit it in anyway, and learns
 42 a single classifier over all classes across task ids. We will nonetheless fix the inaccuracy in our claim. **(4)** The first figure in
 43 Figure 4 is the **training error during meta-testing** and does not involve IID sampling or multiple epochs. It measures degree
 44 of forgetting without taking into account the generalization error, and does in-fact perform very close to Oracle.

45 **Limitations: (2)** We fully agree that a fixed representation can not solve continual learning. We addressed this limitation in
 46 L273-L275 by suggesting one strategy which can be used to continuously update the representation. For the purpose of this
 47 paper, we focused on demonstrating that an effective representation can greatly reduce interference. We are currently extending
 48 this work using this proposed strategy, with a slowly changing representation updated using the OML objective.

49 **Reviewer 3: Please... sparse SR-NN method works:** SR-NN regularizes the activations across a mini-batch to be instance
 50 sparse where a feature is $x\%$ instance sparse if it is non-zero for $x\%$ of examples in a batch/mini-batch of data. We will add a
 51 more detailed description of SR-NN, and a precise definition of instance sparsity in the appendix.

52 **Equation (3)** The expectation is taken with respect to all possible length k trajectories, starting from $X_t = x$: over all the
 53 random variables $\{(X_{t+i}, Y_{t+i})\}$. The outer integral is an expectation over X_t , according to distribution μ .

54 **Diagram and pseudocode** We will move the pseudocode and diagram to the main paper.

55 **Could ... How do these related to Figure 4?** The results on the left of Table 1 (One class per task, 50 tasks) correspond
 56 to $x=50$ in Figure 4. Online + MRCL correspond to the MRCL line at $x=50$ whereas Online + Pretraining corresponds to
 57 Pretraining line at $x=50$.

58 **...EWC are surprisingly low .. why?** This is a great question. There are two reasons. (1) EWC tends to do extremely poorly
 59 on incremental classification tasks. (2) It does poorly when using a single pass through the trajectory, because the model
 60 does not necessarily converge on a task in a single pass. Our results are consistent with those reported by Riemer *et.al* 2019,
 61 Chaudhry *et.al* 2019 and others.

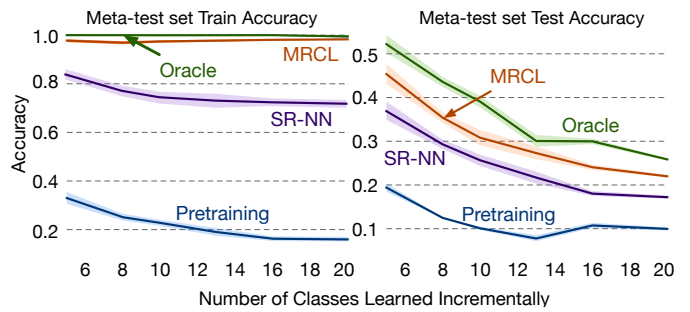


Figure 1: Reproducing the classification results on Mini-Imagenet.