

1 We thank the reviewers for their time and feedback.

2 **R1: Score 6 (confidence 3)**

3 **Q1: "Not novel enough".** We disagree. The content of Sec. 3.2 and 3.3 are entirely new. Both R2 and R3 agree.

4 **Q2: "What is the fundamental difference between converting whole network vs only the last layer"?** In the
5 beginning of training, weights are changing frequently from task to task, therefore their uncertainty is high. Using only
6 the last layer ignores the uncertainty in the rest of the weights. This could hurt performance a lot in the beginning. We
7 will add more explanation in the paper, and include a small illustrative experiment in the Appendix.

8 **Q3: "What role does the ... regularization term play ... compared with FRCL"?** This term ensures that the current
9 DNN outputs (mean) after task t remain close to the DNN after task $t - 1$, which is desirable to avoid forgetting. In
10 contrast, FRCL only optimise the 'variance information' to be close, which may not be as effective (see line 44).

11 **Q4: "Is it possible to do task detection?"** Yes, it is possible to do this in the same way as FRCL paper.

12 **Q5: "Does sampling according to Λ really makes sense?"** Yes, this is closely related to Kernel methods, such as
13 kernel ridge-regression and SVM, where this strategy enables us to pick boundary/high-leverage points. This also leads
14 to an intuitive and computationally cheap method (both R2 and R3 agree).

15 **Q6: "FROMP does not outperform FRORP"... Is random sampling enough?** It is not enough, and a careful
16 selection is required when memory size is small. See Split CIFAR results in Figures 3b and 3c. On MNIST, FRORP
17 is similar to FROMP because the task is very simple, and very little is gained from a careful selection of memorable
18 examples. We will add results for 11 tasks on Split CIFAR where we expect to see even more of a difference.

19 **Q7: The complexity of computing Λ .** It costs $O(N)$. We will clarify this in the paper.

21 **R2: Score 7 (confidence 3)**

22 **Q8: Reliability under distribution shift:** This should not be an issue as long as the memory is large enough.

23 **Q9: "How does this compare to simple gradient based sample ranking schemes (e.g. Deep Batch Active learn-
24 ing)".** Very good point. Our method turns out to be similar to 'Conf' selection, but we arrived at it with a different
25 approach. We will add this reference and add a discussion in the paper.

26 **Q10: Include experiment with 10s of tasks.** Yes, we will add this. Also see our answer to Q14.

28 **R3: Score 8 (confidence 4)**

29 **Q11: Add Forward and Backward transfer to all results.** Yes, we will add it.

30 **Q12: "Add clarifying assumptions earlier to help categorise this work".** Yes, we will add it (around line 73).

31 **Q13: Broader review of recent work in Continual Learning.** We agree and will improve it. Due to page limit, our
32 Related Work section is short. For camera-ready, we can devote around half a page to include many other related works.

33 **Q14: Ready to raise score if a challenging and convincing CL result is added.** Thanks for the suggestion. Since
34 our main contribution is to propose a new method, we focused on standard benchmarks. We will add the Split CIFAR
35 experiment with 11 tasks (previous papers used only 6 tasks), and the "predictive uncertainty for change-point detection"
36 in the Appendix. Adding many more challenging experiments takes more time and space, and we will explore scalability
37 (e.g., large classes and datasets) and generality (e.g. in RL and Bandits) of the method in a separate future work.

38 **Q15: Behaviour as number of memorable points increases.** Good question. The performance saturates as expected.
39 It is not clear how to find the number of examples beforehand, but one could easily design an online, greedy scheme to
40 determine this number.

42 **R4: Score 5 (confidence 4)**

43 **Q16: "Limited contributions".** We disagree. The content of Sec. 3.2 and 3.3 are entirely new. Both R2 and R3 agree.

44 **Q17: Discuss more related approaches other than regularization ones.** Thanks. We will fix this (including adding
45 all the references mentioned, and a comparison with GEM). Also, see Q13.

46 **Q18: "Is optimization constraints violated?"** This is a misunderstanding. There are no constraints in our approach.

47 **Q19: Gap between FRORP and FROMP.** This is a misunderstanding. Please see our answer to Q6.