

1 We thank the reviewers for their time, feedback and highly encouraging comments. **It was acknowledged that our**
2 **algorithm is intuitive and principled (R4), achieves significantly better results (R1, R3), is clearly presented and**
3 **situated (by all), and is novel and relevant to the community (by all).** We will incorporate suggested improvements
4 to the paper regarding punctuation, notation, algorithm box and typos. We address here the remaining concerns:

5 **R2: Discrepancy in baselines’ numbers, lacking experimental information:**

6 We appreciate the reviewer bringing up this point since it is important for the
7 reader to understand how our comparisons are made, which we will make clearer.
8 Regarding the disparity, the cited works use a higher resource allocation than
9 ours: a large replay buffer of size 5120 [Ref4]/5000 [Ref2 (results from AGEM)]
10 for GEM, while our MNIST [Rotation/Permutation] experiments have a buffer
11 containing only 200 samples. For EWC, the network used is 20 times larger
12 [Ref2], there are only 10 tasks [Ref1, Ref2] (we have 20 tasks) and the setting is
13 not single-pass [Ref1] (they train each task for 20 epochs). It should be noted that
14 the low memory regime is where the performance trends of many CL methods
15 are most pronounced and meaningful. Given unlimited memory/compute, all
16 the methods perform comparably to ER training (as also noted by MER). We will
17 expand the experimental description in the main paper to highlight these details,
18 which we have currently outlined in Appendix F and G due to space constraints.

19 **R2: RHS \neq LHS in Eq. 1:** While we believe the equation is technically correct, we acknowledge that it might be
20 confusing since we have clubbed the implied data arguments into (X^n, Y^n) for brevity. We will separate these for each
21 of the task-specific loss functions in the camera ready for correctness.

22 **R3: Timing comparison:** This is a good point, and we will include a plot for our *Multi-Pass* experiments on CIFAR,
23 showing the total running time for La-MAML, GEM, AGEM, iCARL and ER, as shown here in Fig.1.

24 **R3, All: Multi-headed Problem Setting:** We thank R3 for raising this point and make a correction to our experimental
25 description: While our real-world vision experiments are multi-headed, our MNIST experiments are in the single-headed
26 domain-incremental setup (since as mentioned in Section 5.1, the output space for all tasks is the same set of 10 classes
27 while the common transformation to the digits varies with each task). The paper thus contains **both task-aware and**
28 **task-agnostic** experiments since our algorithm’s working is task-agnostic. We omit Class-IL settings since they have
29 many of their challenges arising primarily from the bias imbalance in the classification layer. Many class-IL works
30 specifically focus on this issue since it is tricky to isolate it to study the general forgetting problem in CL (R3 Ref. [7]).

31 **R3: Choice of backbone:** We have tried to use an architecture (3-4 conv layers + 2 FCs) that is commonly used
32 in meta-learning works for its simplicity^{1 2 3}, and have used it to run all our baselines. In the CL setting with
33 ever-increasing tasks, any model will eventually be under-parameterised. As long as the model performs decently in the
34 *i.i.d* setting and there is a gap between the *i.i.d*-trained model and other CL methods, it should be valid to use it to study
35 the CL problem. If the reviewer recommends, we will add a sensitivity analysis for network sizes to the Appendix.

36 **R3: Relevance of LLL setting:** We agree that the issue of how setup constraints are commonly chosen in CL works
37 is worthy of debate. However, we reiterate that the LLL setting is challenging yet realistic in many cases where it
38 is not feasible to store all the within-task data points, and is also studied in many prior works like AGEM, MER. It
39 should also be noted that we take multiple gradient steps (*glances*) over each sample in the LLL setup (described in
40 Appendices E,F), thus making enough updates to the parameters.

41 **R3, R1: Lookahead plot:** We had hoped to show through Fig.3 in the paper, how the performance varies as more tasks
42 are added to the problem (as asked by R1). Note that the average accuracy stays roughly the same across an increasing
43 number of tasks. We shall remove this figure if it is not considered informative by the reviewers.

44 **R1: RL experiments:** We agree and think it would be particularly interesting to test our algorithm for Model-Based
45 RL, where models learnt online from a temporal data stream should undergo considerable *forgetting*. However, given
46 the ambiguities and lack of benchmarks for properly defining a *continual* setup in RL, we are pursuing it as an extension
47 and it is out of scope for this work.

48 **R1: Lookahead search:** We added the following: *"In optimisation literature, lookahead search usually evaluates the*
49 *fitness of proposed parameter updates based on an auxiliary criterion evaluated after hypothetically applying them.*
50 *These proposals are then modified based on evaluated fitness to make an actual update. In our case, the LRs act as*
51 *the modifiers of the update, and their values result from the evaluation of two criterions: the losses on old and new tasks".*

52

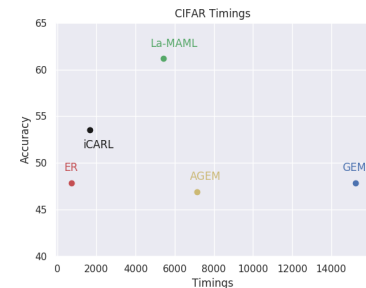


Figure 1: Accuracy vs Timing comparison on CIFAR (time in seconds)

¹Chelsea Finn, Aravind Rajeswaran, Sham Kakade, Sergey Levine. Online Meta-Learning: Proceedings of the 36th International Conference on Machine Learning, PMLR 97:1920-1930, 2019.

²Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing networks: Adaptive selection of non-linear functions for multi-task learning. In International Conference on Learning Representations, 2018.

³Yu, T., Kumar, S., Gupta, A., Levine, S., Hausman, K., and Finn, C. Gradient surgery for multi-task learning. arXiv preprint arXiv:2001.06782, 2020.