

1 Thank you for noting the novelty of LPG-FTW, the fairness of evaluation, and the presentation and strength of results.

2 **R1:** Thank you for the very positive comments and detailed suggestions (**requested non-stationary results below**).

3 **3) I.i.d. assumption and non-stationarity:** We assume tasks are drawn i.i.d. from a task distribution. In our evaluation,
4 we sample tasks at random, and vary their order across trials (line 249). LPG-FTW supports interleaved MTL, but we
5 evaluate a true lifelong setting, without revisiting tasks. Following your suggestion, we ran LPG-FTW on domains
6 violating this i.i.d. assumption, with gravity increasing linearly in $[0.5, 1.5]g$. LPG-FTW, although not specifically
7 designed to handle non-stationarity, performs equivalently to the i.i.d. setting on Half-Cheetah & Hopper (Fig. 1). We
8 did no parameter tuning due to time; we think Walker2D needs tuning to handle the violated assumption.

9 **3) Storing all Hessians:** We don't store H 's in practice after updating A, b , since LPG-FTW only needs H if revisiting
10 tasks, which we don't do. If tasks are revisited (e.g., interleaved MTL), it would require storing H at a cost of $O(d^2T)$.

11 **8) Poor performance of experience replay (ER):** as stated in Sec. 2, ER in PG algorithms (via importance sampling) is
12 unstable. If the policy moves far from earlier tasks, replay stops helping, which is aggravated as training progresses.

13 **R3:** Thank you. We address as many comments as space allows, and will update our draft with cites and clarifications.

14 **3&6) Previous work:** LPG-FTW improves PG-ELLA in 3 main ways. 1) PG-ELLA trains single-task (STL) policies
15 on each task separately and then finds $L, s^{(t)}$ via dictionary learning on the STL policies (line 60), while LPG-FTW
16 learns directly in the factored space (finding $s^{(t)}$ given the current L). 2) PG-ELLA has no explicit initialization, while
17 LPG-FTW uses Algo. 2, ensuring that columns of L are diverse. 3) PG-ELLA assumes STL finds optimal policies,
18 while LPG-FTW adds a linear term to the cost to address non-optimality (line 129). This enables LPG-FTW to work on
19 far more complex tasks than PG-ELLA (Meta-World vs cartpole). We will emphasize these connections in Sec. 4.

20 **3) Chosen baselines:** We chose EWC and PG-ELLA as representatives of single- & multi-model classes. Progressive nets
21 scale poorly in lifelong settings. P&C suffers from the same limitations as EWC due to the single-model assumption,
22 and so it wouldn't scale well to the highly diverse tasks we study in Meta-World.

23 **4) Correctness of PG-ELLA evaluation:** We consulted with one author of PG-ELLA to validate that our evaluation was
24 correct. In PG-ELLA, the lifelong-learned policies are used only as a warm start for subsequent learning. They first
25 train STL policies, then run dictionary learning on the pre-trained policies to find an L matrix, and finally use the $Ls^{(t)}$
26 policies to start a second STL process—for evaluation. Their main result (Fig. 2 in PG-ELLA) shows that initializing
27 STL from a lifelong-learned policy accelerates training. In contrast, our evaluation is entirely lifelong learning: the
28 agents are evaluated as they train on each task sequentially. There, PG-ELLA *does not* leverage information across tasks
29 (Fig. 1 and 3-top in our paper). Once all tasks are trained via STL, PG-ELLA runs dictionary learning and we evaluate
30 the policy in bars 3 and 4 of Fig. 2 and 3-bottom. Bar 4 corresponds to the starting point of the evaluation in PG-ELLA.

31 **5) Clarity of notation:** LPG-FTW keeps one dictionary shared across tasks: the multi-task cost in Eq. 2 is optimized
32 incrementally with auxiliary matrices A, b . L_t denotes the L learned up to task t . We only keep one A matrix for all
33 tasks (A_t is an auxiliary variable). \hat{t} indexes over previous tasks, and t is the current task. α is the policy obtained for a
34 task after its training process (using the L up to the previous task). We will emphasize these clarifications in the text.

35 **8) Additional experience initializing L ?:** There is none: initialization replaces the standard training for tasks 1– k .

36 **8) Why is STL=PG-ELLA in 5.1 but not 5.2?:** Figs. 1 and 3-top show the training process, for which PG-ELLA uses
37 STL. Figs. 2 and 3-bottom show subsequent steps (bars 3 and 4) where PG-ELLA combines information from all tasks.

38 **8) Comparable capacity?:** All methods used policies of the same size, but overall capacity naturally varied across
39 methods. We followed your suggestion of increasing the capacity of EWC to match that of LPG-FTW in Meta-World,
40 and found that results didn't change significantly (Fig. 2).

41 **7) Reproducibility:** The clarifications above, along with the full code we provided, should make our results reproducible.

42 **8) How does training progress?:** Number of iterations is summarized Appendix B. Each task is seen by the agent only
43 once, and there is no possibility of going back for further experience at any point during the process.

44 **8) EWC performance drop after training all tasks:** EWC suffers from catastrophic forgetting, unlike our approach.

45 **R4:** Thank you for the very positive comments. For connections to PG-ELLA, see **R3**.

46 **8.1) Connections between charts:** 'Tune' matches the end of the curves, considering only per-task training but not how
47 that affects previous tasks. 'Update' and 'Final' come later: 'Final' assesses performance after all tasks are trained.

48 **8.2) Does STL performance imply tasks are similar?:** STL *ignores* all other tasks when training on one task, so task
49 similarity does not play a role. STL performs well because STL performance was used in the design of the benchmarks.

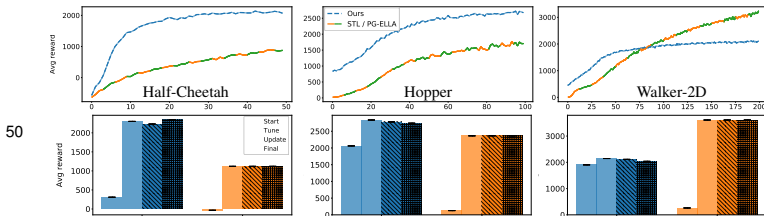


Fig. 1: Non-stationary results requested by **R1**. STL performance (for reference) measured on original i.i.d. task distribution.

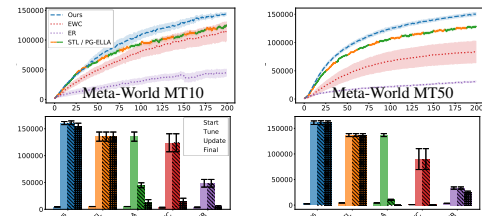


Fig. 2: Higher EWC capacity, requested by **R3**.