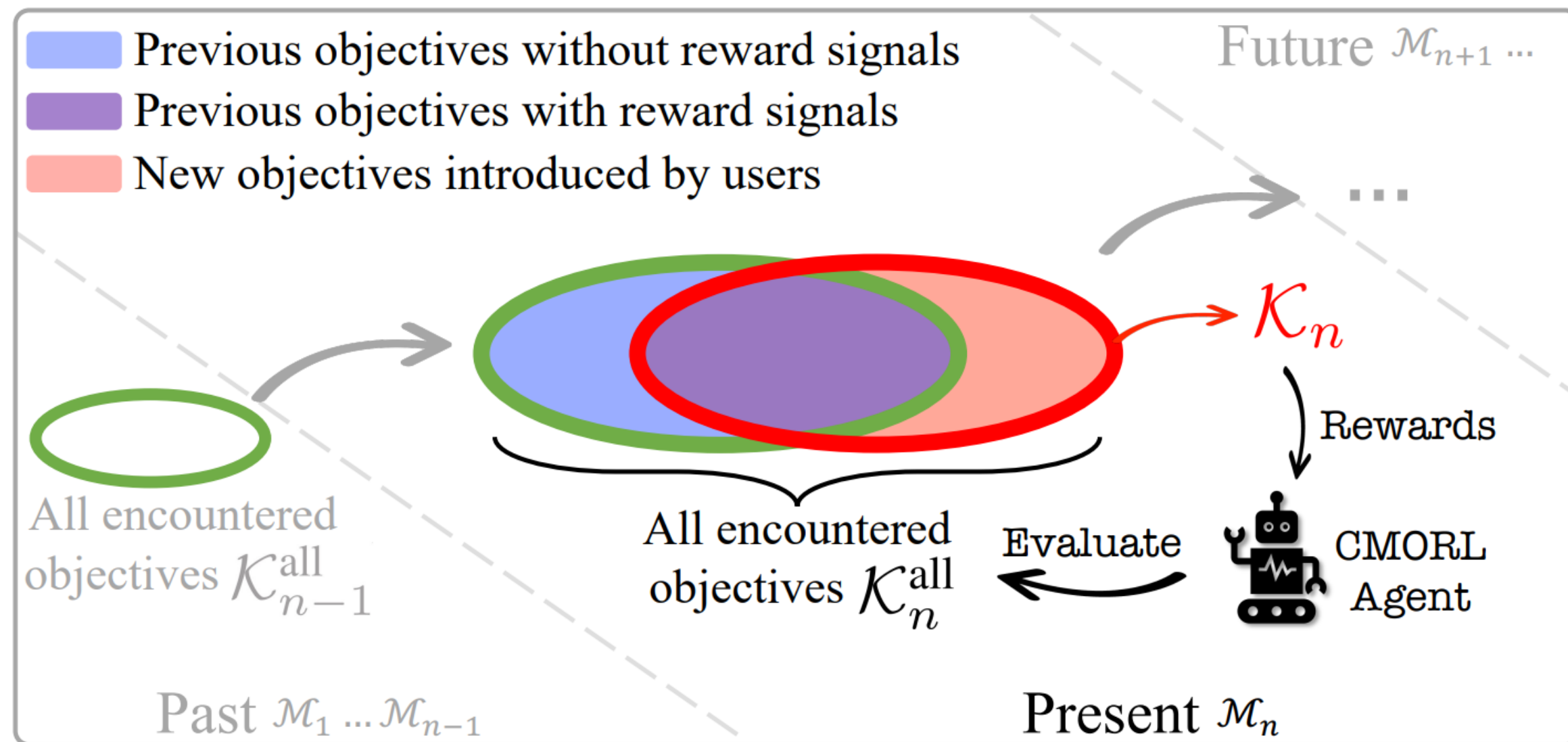


Introduction

MORL methods focus on learning to complete and trade-off a **fixed** set of objectives. But the objective set may **change over time** in some real-world scenarios.

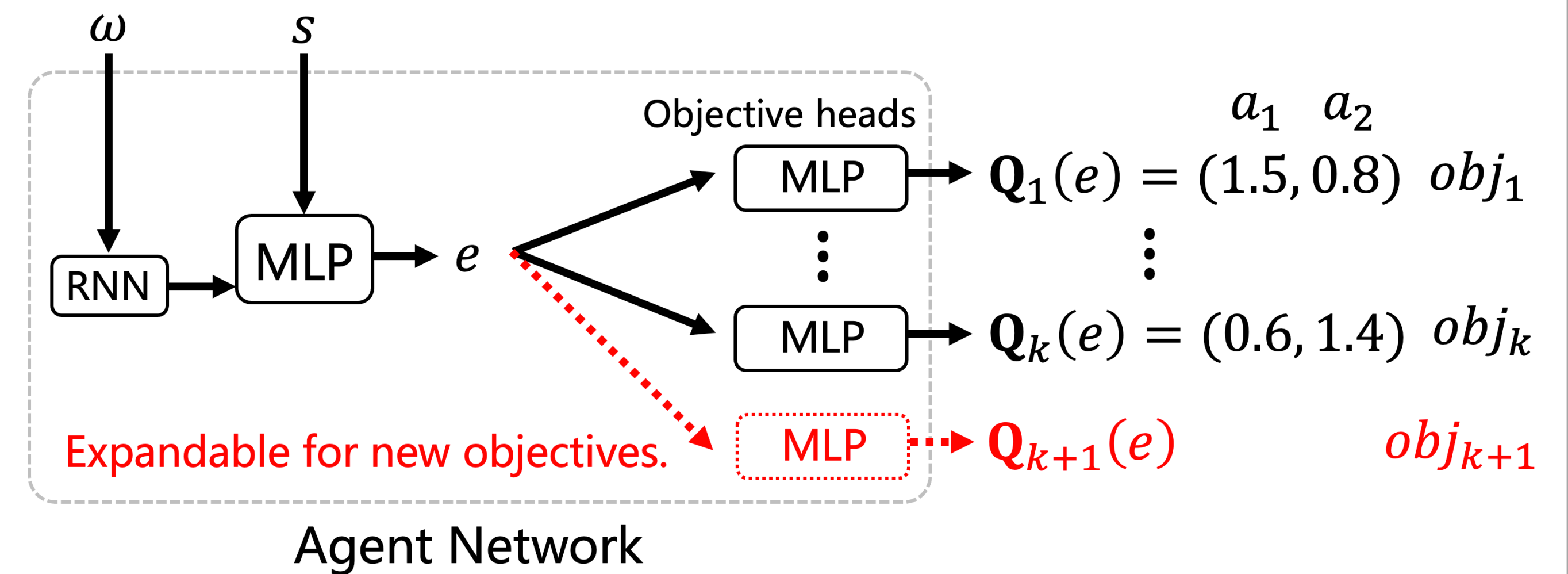


To bridge this gap, we formalize the **continual MORL** problem and try to solve it with our method **Continual Multi-Objective Reinforcement Learning via Reward Model Rehearsal (CORE3)**.

Dynamic Network for Evolving Objectives

- For new objectives, expand the agent network with multi-head architecture.
- No need to train the entire network from scratch.
- Results: **Learn to complete new objectives faster.**

Method



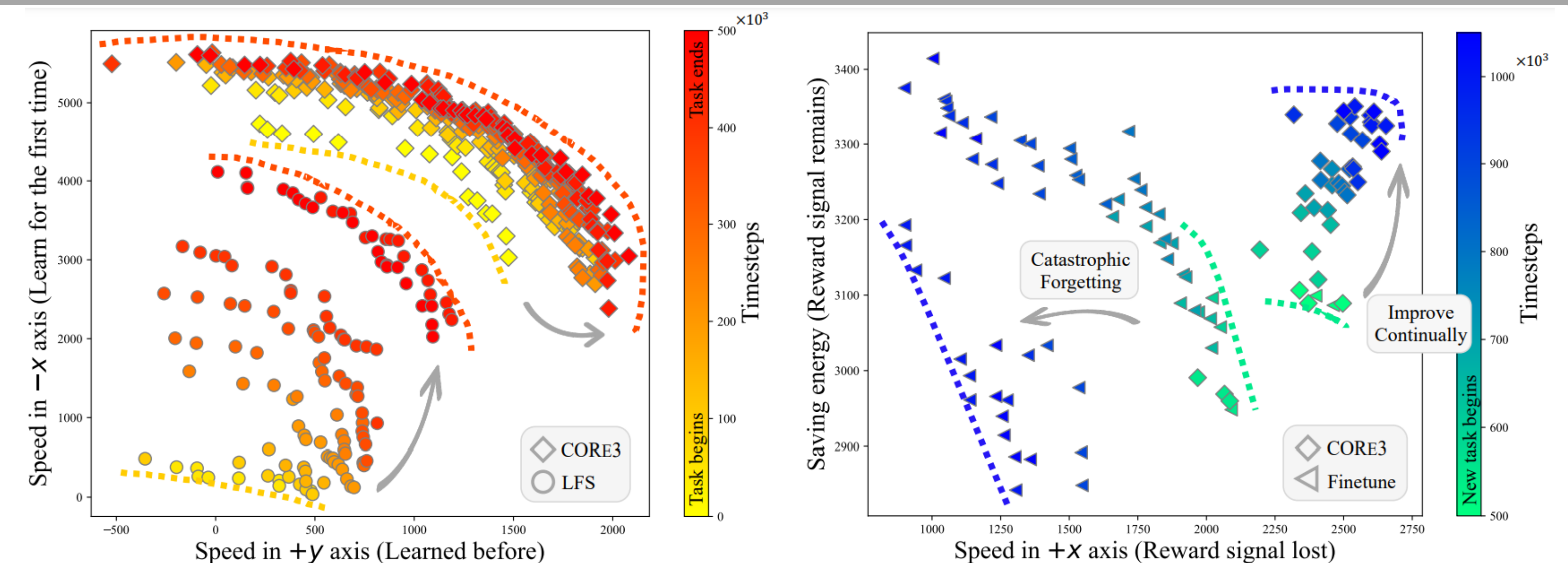
Multi-Objective Reward Model Rehearsal

- For old objectives lacking reward signals, we learn a multi-objective reward model $\hat{\mathbf{R}}(s, a) = (\hat{R}_1(s, a), \dots, \hat{R}_M(s, a))$ to recover these signals, and use them to train the agent, i.e., reward model rehearsal.
- Performance guarantee: For all preferences, performance loss $\leq \delta / (1 - \gamma)$, where δ is the model prediction error.
- Results: **Maintain the ability to complete old objectives lacking reward signals.**

Experiments

| Envs | Metrics | Finetune | CLEAR | EWC | CORE3 (Oracle) | CORE3 (Ours) |
|------------------------|---------------------|------------------|------------------------------------|-----------------------------------|------------------|-----------------------------------|
| FTN | Hv (\uparrow) | 1.00 \pm 0.07 | 0.85 \pm 0.07 | 1.04 \pm 0.07 | 1.10 \pm 0.04 | 1.06 \pm 0.00 |
| | SP (\downarrow) | 1.00 \pm 0.14 | 0.73 \pm 0.02 | 0.58 \pm 0.06 | 0.72 \pm 0.01 | 0.82 \pm 0.02 |
| | BWT (\uparrow) | -0.12 \pm 0.02 | 0.00 \pm 0.01 | -0.14 \pm 0.01 | 0.03 \pm 0.00 | 0.02 \pm 0.01 |
| Grid | Hv (\uparrow) | 1.00 \pm 0.82 | 4.21 \pm 0.66 | 2.26 \pm 1.34 | 4.60 \pm 0.26 | 4.54 \pm 0.60 |
| | SP (\downarrow) | 1.00 \pm 0.58 | 0.42 \pm 0.11 | 0.47 \pm 0.10 | 0.37 \pm 0.02 | 0.44 \pm 0.10 |
| | BWT (\uparrow) | -0.44 \pm 0.07 | -0.14 \pm 0.17 | -0.58 \pm 0.05 | -0.07 \pm 0.12 | -0.17 \pm 0.15 |
| Ant | Hv (\uparrow) | 1.00 \pm 0.42 | 1.78 \pm 0.37 | 1.03 \pm 0.60 | 2.30 \pm 0.18 | 1.89 \pm 0.35 |
| | SP (\downarrow) | 1.00 \pm 0.56 | 0.33 \pm 0.06 | 1.43 \pm 1.78 | 1.06 \pm 0.57 | 0.73 \pm 0.72 |
| | BWT (\uparrow) | -0.27 \pm 0.02 | 0.01 \pm 0.18 | -0.35 \pm 0.1 | 0.34 \pm 0.18 | 0.11 \pm 0.16 |
| Hopper | Hv (\uparrow) | 1.00 \pm 1.22 | 2.03 \pm 3.49 | 1.01 \pm 1.43 | 7.84 \pm 0.29 | 8.23 \pm 1.42 |
| | SP (\downarrow) | 1.00 \pm 0.76 | 7.59 \pm 2.98 | 4.76 \pm 3.54 | 0.35 \pm 0.11 | 0.21 \pm 0.10 |
| | BWT (\uparrow) | -0.40 \pm 0.07 | -0.38 \pm 0.22 | -0.37 \pm 0.07 | 0.11 \pm 0.09 | 0.06 \pm 0.05 |
| API (%) (\uparrow) | Hv | \ | 121.72 | 33.61 | 296.03 | 293.18 |
| | SP | \ | -126.81 | -80.85 | 37.61 | 45.37 |
| | BWT | \ | 69.09 | -17.87 | 145.36 | 109.09 |

CORE3 learns to complete and trade-off all encountered objectives the best.



(a) CORE3 v.s. LFS on learning a new task.

(b) CORE3 v.s. Finetune on an ended task.

CORE3 learns new objectives faster and better remembers old objectives.