

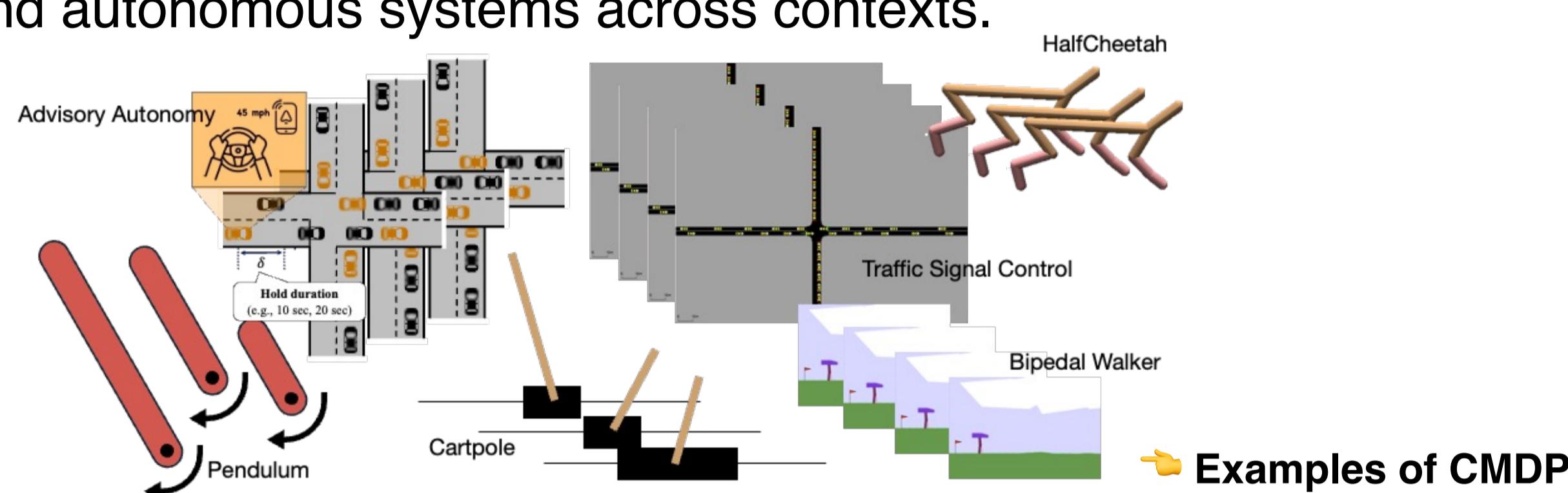
Model-Based Transfer Learning for Contextual Reinforcement Learning

Jung-Hoon Cho, Vindula Jayawardana, Sirui Li, Cathy Wu
Massachusetts Institute of Technology



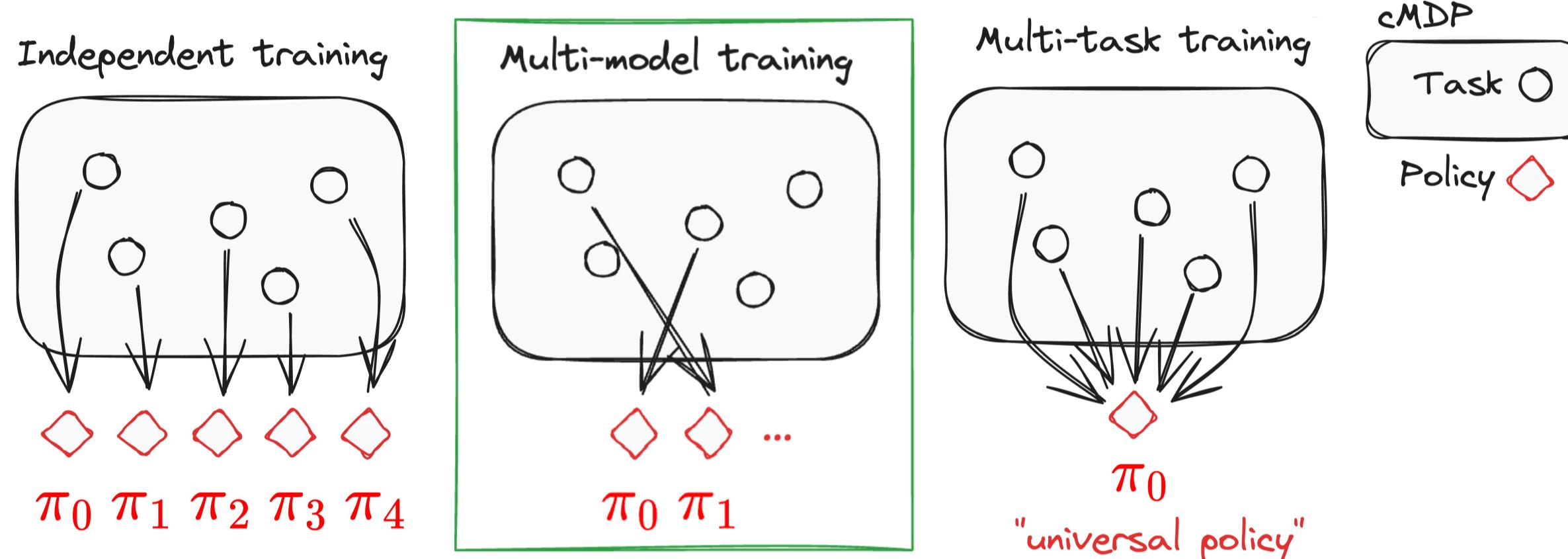
Introduction

- Despite deep reinforcement learning (DRL) remarkable success across domains, DRL models often exhibit **brittleness when exposed to small variations in task settings**.
- Contextual Markov Decision Process (CMDP)** extends traditional MDPs with additional contextual information, leading to multiple, closely related MDPs.
- Motivation:** Aim to improve generalization for practical applications like robotics and autonomous systems across contexts.



Model-Based Transfer Learning

- Resource-intensive traditional DRL approaches** (e.g. Independent training or multi-task training)



MBTL Algorithm

Strategically select source tasks for multi-model training

How? Simple & explicit modeling of generalization gap (model-based)

Algorithm 1 Model-based Transfer Learning (MBTL)

```

Input: CMDPs  $\mathcal{M}(x)$ , Task  $x \in X$ , Transfer budget  $K$ 
Output:  $\pi$  and  $V$ 
Initialize:  $J, V = 0 \forall x \in X, \pi = \{\}, k = 1$ 
1: while  $k \leq K$  do
2:    $x_k \leftarrow \text{NextTask}(J_{k-1}, V_{k-1})$ 
3:    $\pi_{x_k} \leftarrow \text{Train}(\mathcal{M}(x_k))$ 
4:    $\pi \leftarrow \pi \cup \{\pi_k\}$ 
5:   Update  $J(x), V(x)$ 
6:    $k \leftarrow k + 1$ 
7: end while
8: return  $\pi$  and  $V$ 

```

Normalized performance

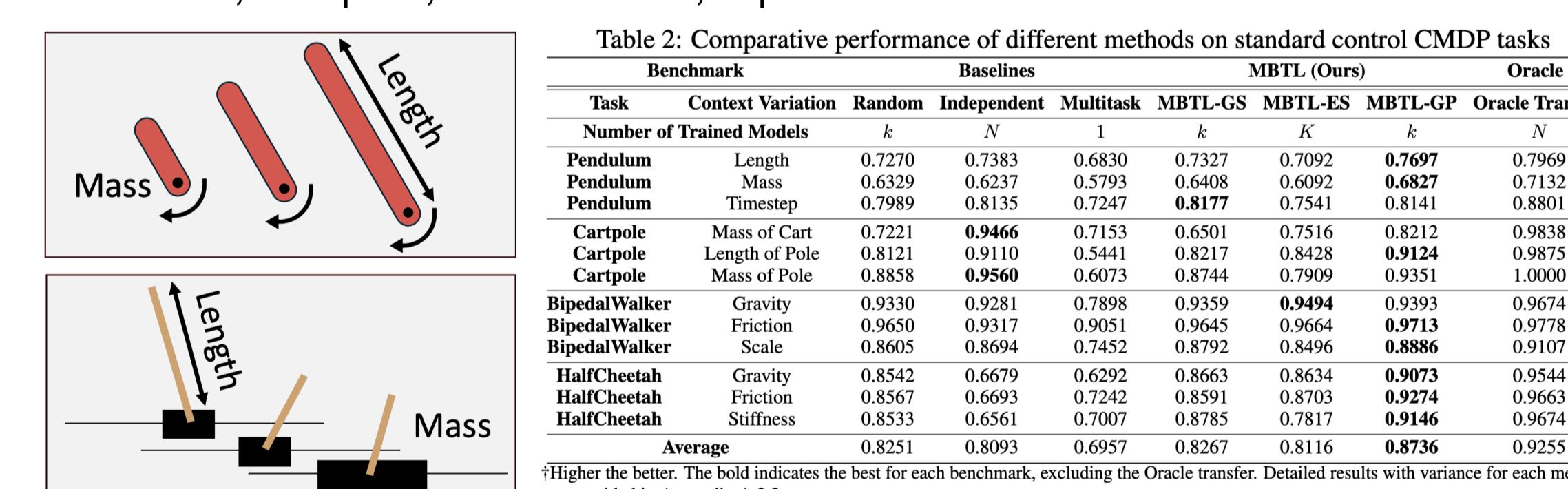
Simple methods [1]



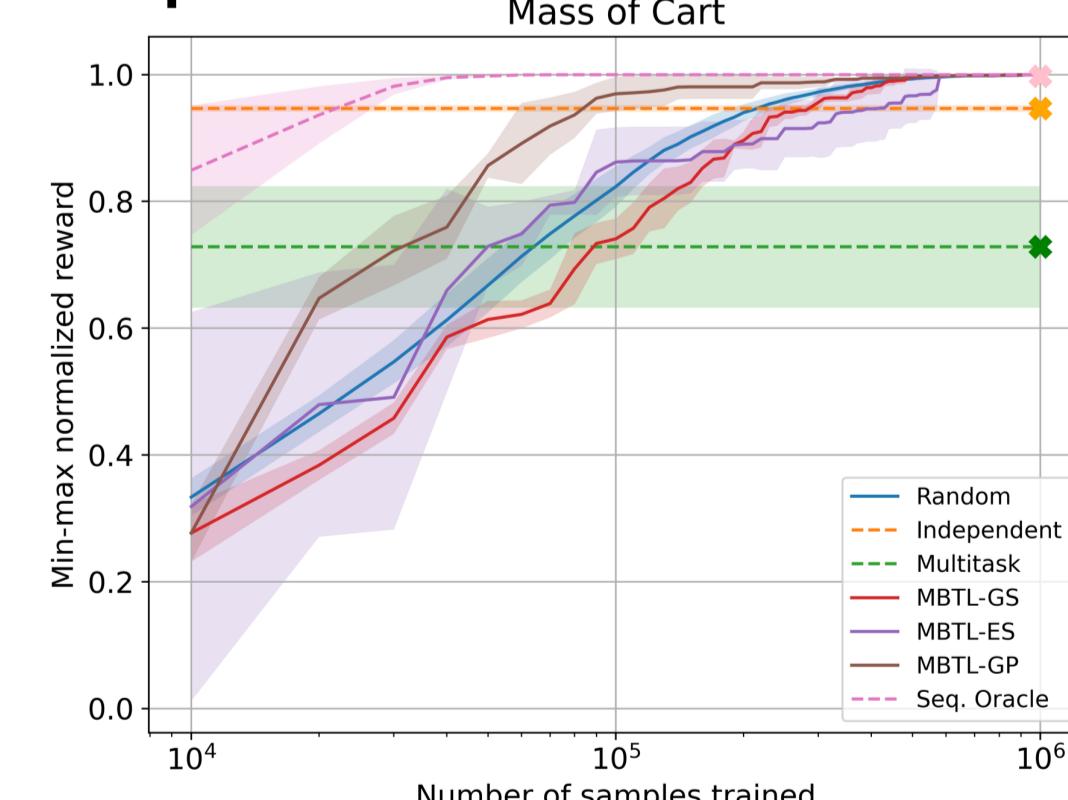
Numerical Experiments

Standard control CMDP benchmarks

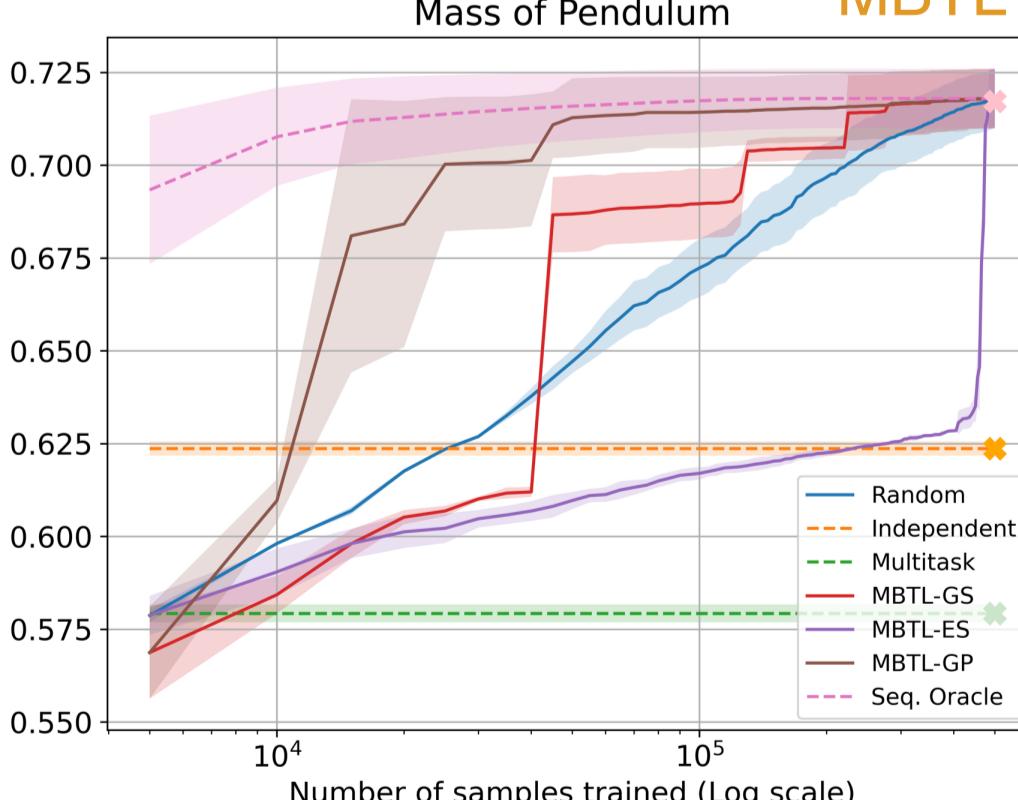
Pendulum, Cartpole, HalfCheetah, BipedalWalker



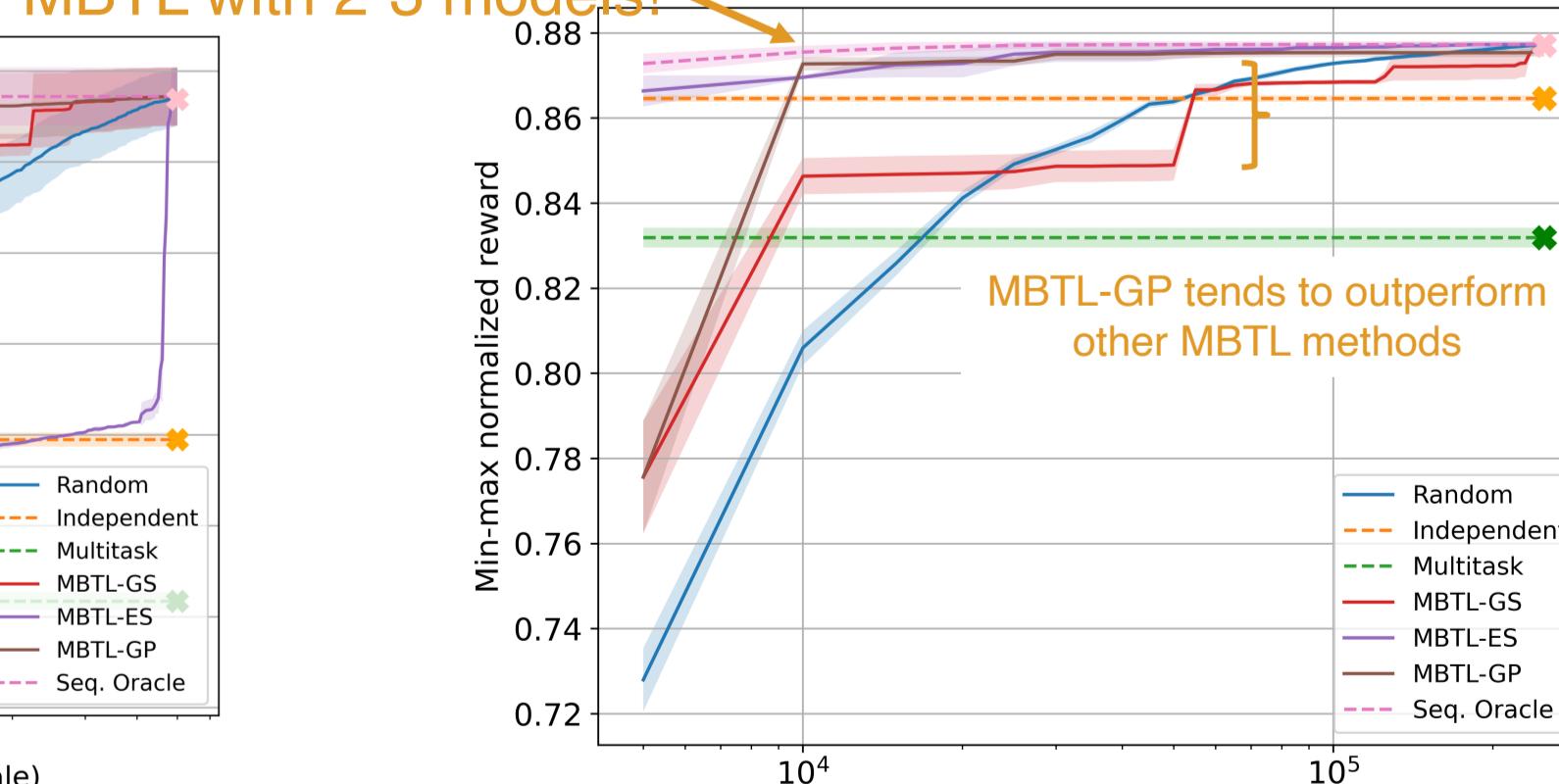
Cartpole



Pendulum

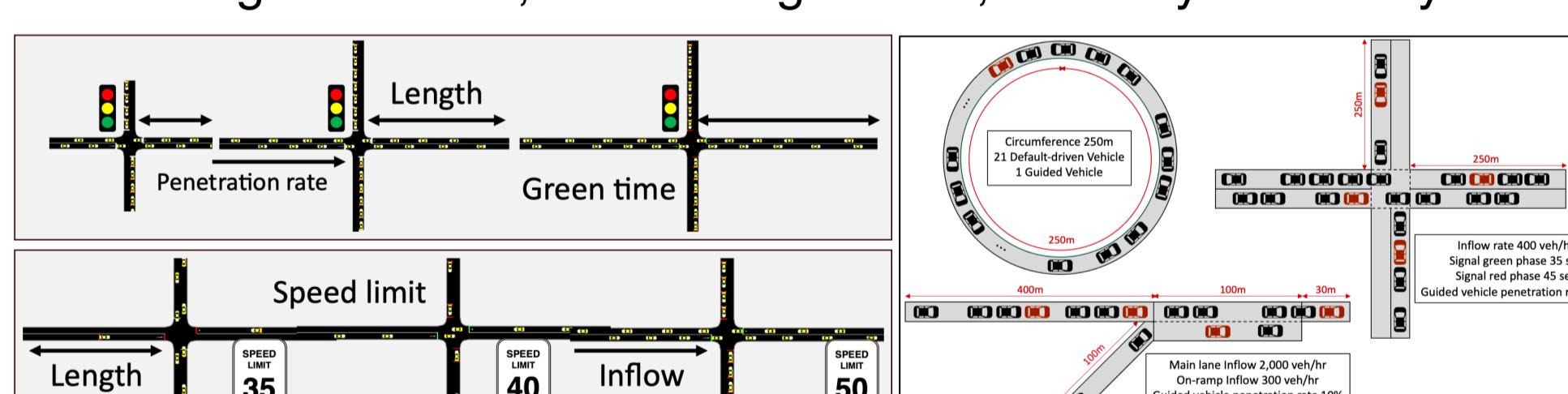


MBTL with 2-3 models!

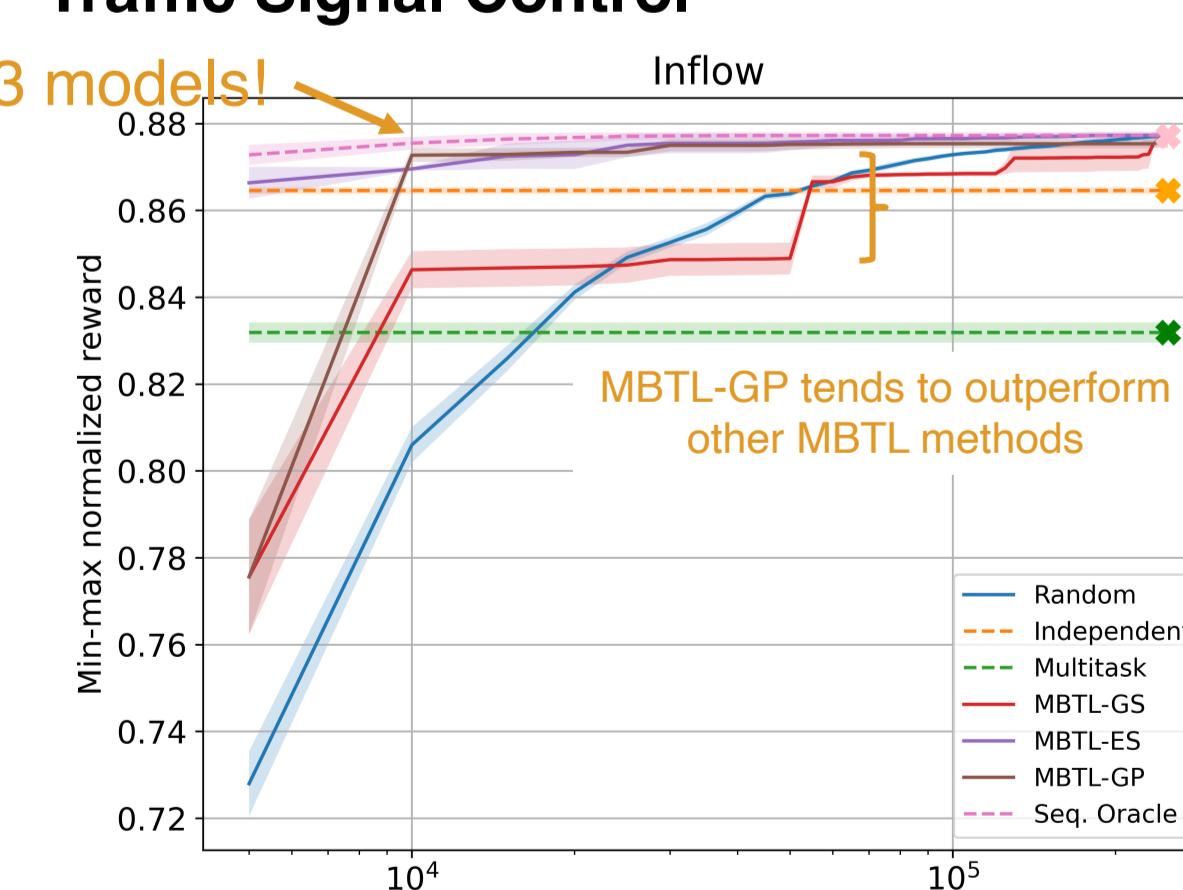


Transportation CMDP benchmarks

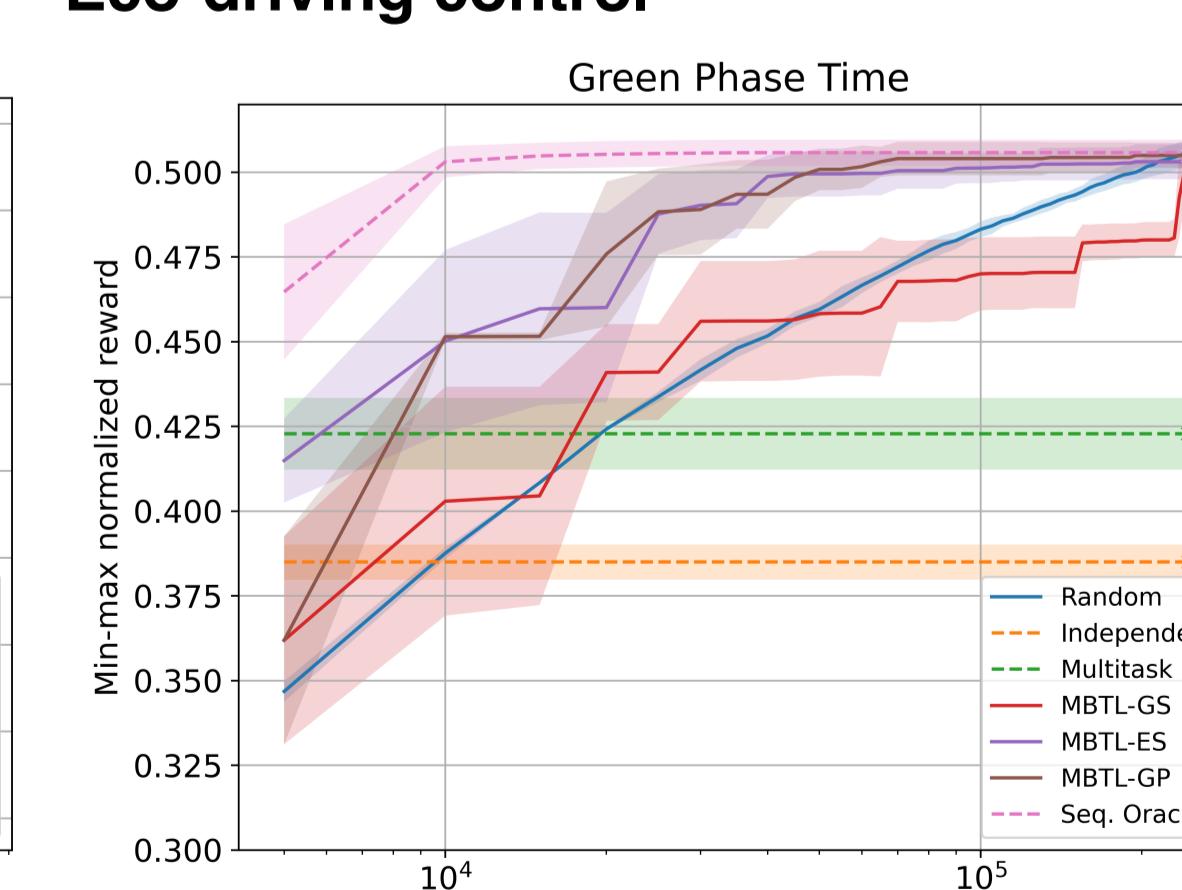
Traffic Signal Control, Eco-driving control, Advisory autonomy



Traffic Signal Control



Eco-driving control



Real-time driving advisory

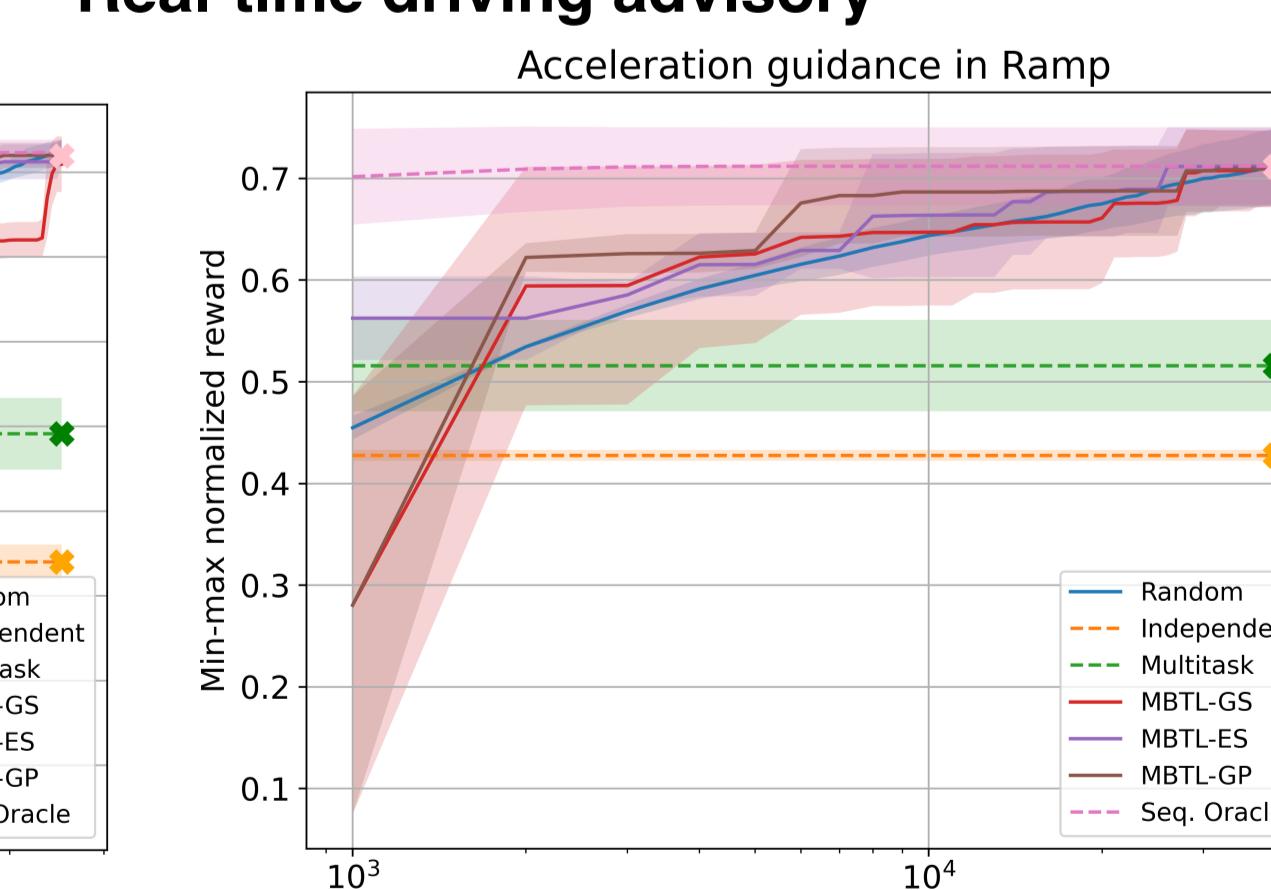


Table 1: Comparative performance of different methods on traffic CMDP tasks

Benchmark Context Variation Random Baseline Multitask MBTL-GS MBTL-ES MBTL-GP Oracle Transfer

Benchmark	Context Variation	Random	Independent	Multitask	MBTL-GS	MBTL-ES	MBTL-GP	Oracle Transfer
Task		k	N	k	N	k	N	
Traffic Signal	Road Length	0.9249	0.9409	0.8242	0.9278	0.9213	0.9371	0.9409
Traffic Signal	Inflow	0.8457	0.8646	0.8319	0.8495	0.8700	0.8673	0.8768
Traffic Signal	Speed Limit	0.8821	0.8857	0.6083	0.8862	0.8854	0.8876	
Eco-Driving	Penetration Rate	0.9509	0.5260	0.9445	0.5827	0.9394	0.6223	0.6060
Eco-Driving	Inflow	0.9774	0.4700	0.9249	0.5927	0.9350	0.5745	0.5526
Eco-Driving	Green Phase	0.4406	0.3850	0.4228	0.4431	0.4557	0.4700	0.5027
AA-Ring-Acc	Hold Duration	0.8924	0.8362	0.9209	0.8776	0.9057	0.9242	0.9552
AA-Ring-Vel	Hold Duration	0.9785	0.9589	0.9720	0.9807	0.9772	0.9816	0.9711
AA-Ramp-Acc	Hold Duration	0.6050	0.4276	0.5158	0.6143	0.5956	0.6318	0.6182
AA-Ramp-Vel	Hold Duration	0.6090	0.5474	0.5304	0.5967	0.6787	0.7082	0.7666
Average		0.7112	0.6778	0.6017	0.7220	0.7354	0.7559	0.7844

*Higher the better. The bold indicates the best for each benchmark, excluding the Oracle transfer. Detailed results with variance for each method are provided in Appendix A.2.

†AA: Advisory autonomy tasks, Ring: Single lane ring, Ramp: Highway ramp, Acc: Acceleration guidance, Vel: Speed guidance.

References

- [1] J.-H. Cho, S. Li, J. Kim, and C. Wu, "Temporal Transfer Learning for Traffic Optimization with Coarse-grained Advisory Autonomy." In revision.

Model-Based Transfer Learning (MBTL) effectively selects training data for solving contextual MDPs by modeling the generalization gap with respect to the context.
(20-35x more sample efficient than independent & multi-task baselines)