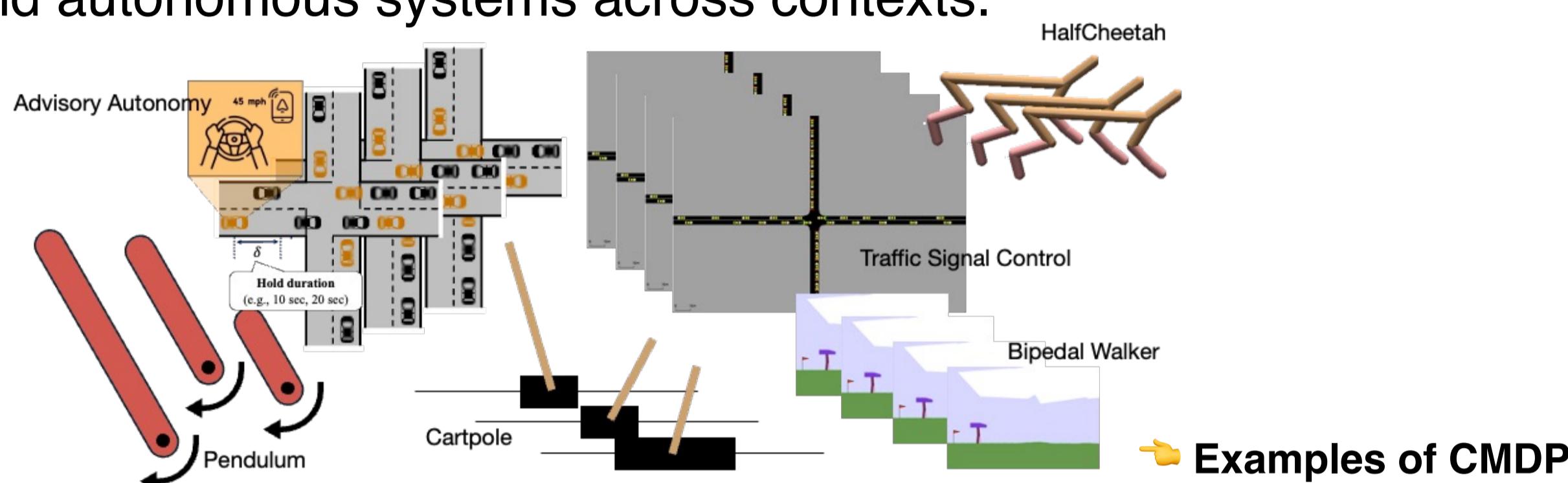


Efficient Source Tasks Selection for Zero-shot Transfer in Contextual Reinforcement Learning

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

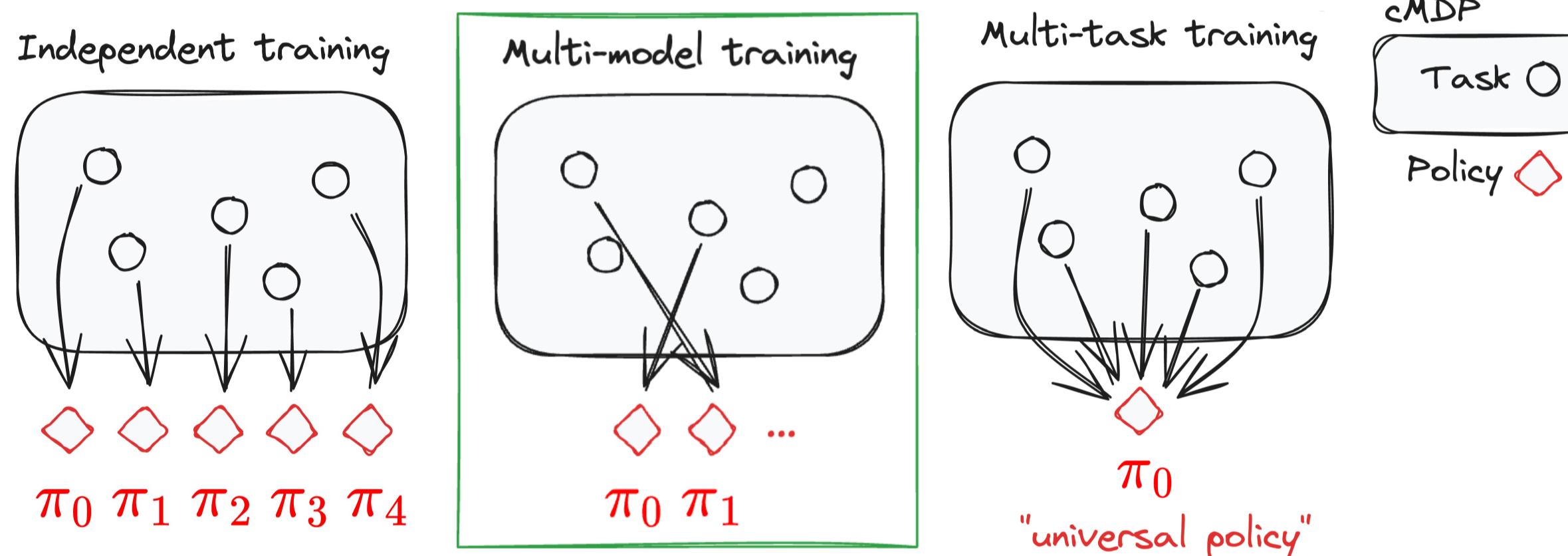
Introduction

- 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)

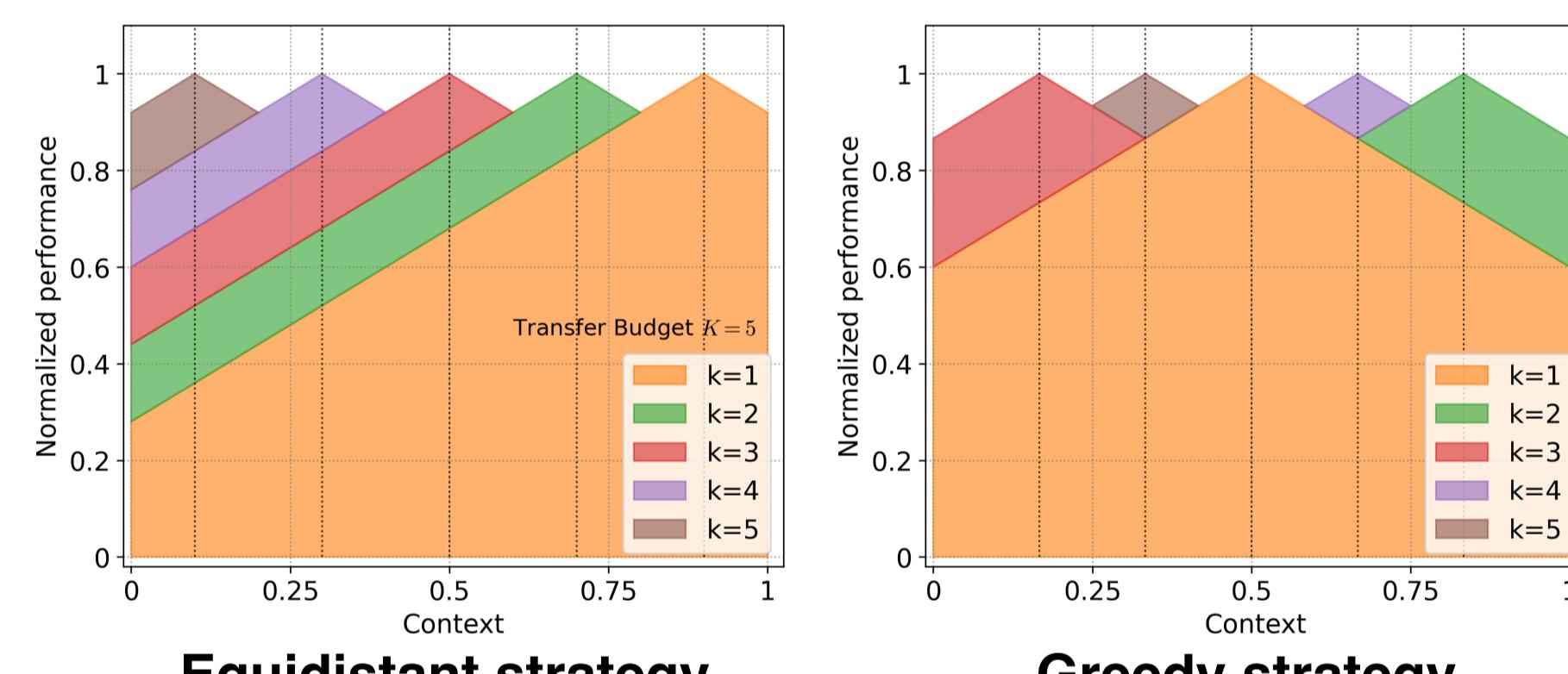
```

Input: CMDPs  $c\mathcal{M}(x)$ , 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:    $\pi_k \leftarrow \text{Train}(\mathcal{M}(x_k))$ 
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  $V$ 

```

Simple methods [1]

Assumptions: Training performance is constant.



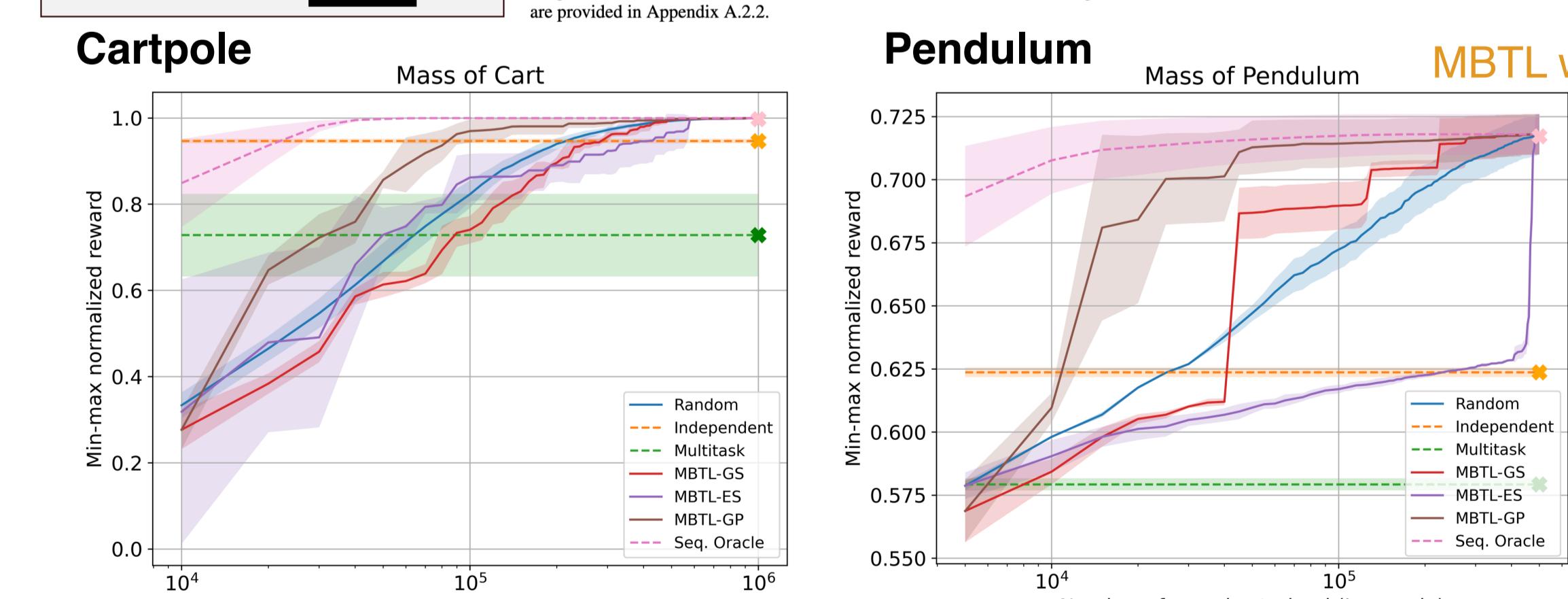
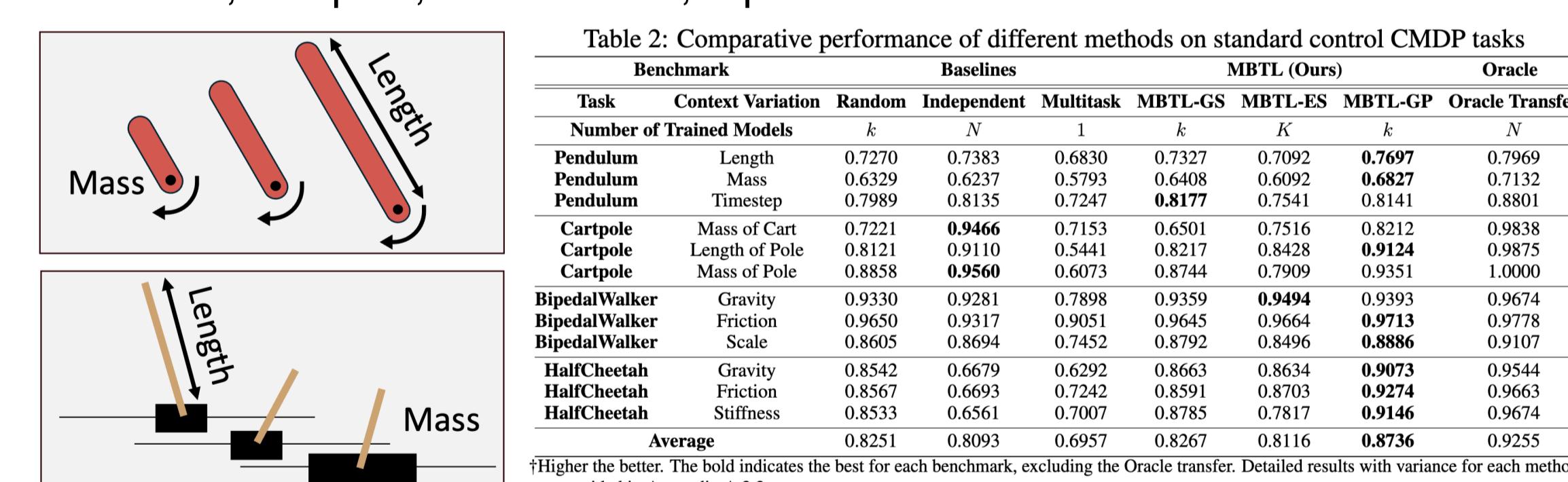
Other Baselines

Independent training, Multitask training, Oracle

Numerical Experiments

Standard control CMDP benchmarks

Pendulum, Cartpole, HalfCheetah, BipedalWalker

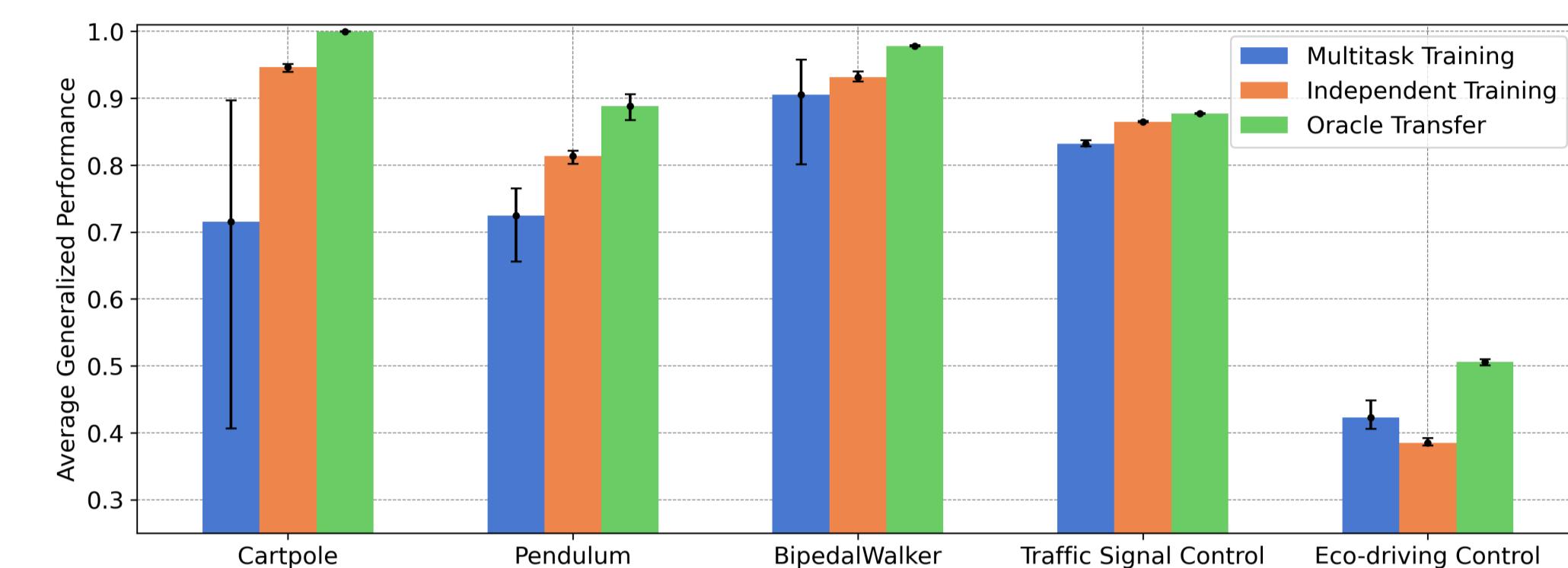


References

- [1] J.-H. Cho, S. Li, J. Kim, and C. Wu, "Temporal Transfer Learning for Traffic Optimization with Coarse-grained Advisory Autonomy." In revision.
[2] J.-H. Cho, V. Jayawardana, S. Li, and C. Wu, "Model-Based Transfer Learning for Contextual Reinforcement Learning." EWRL17, NeurIPS 2024. To Appear.

Zero-shot Transfer

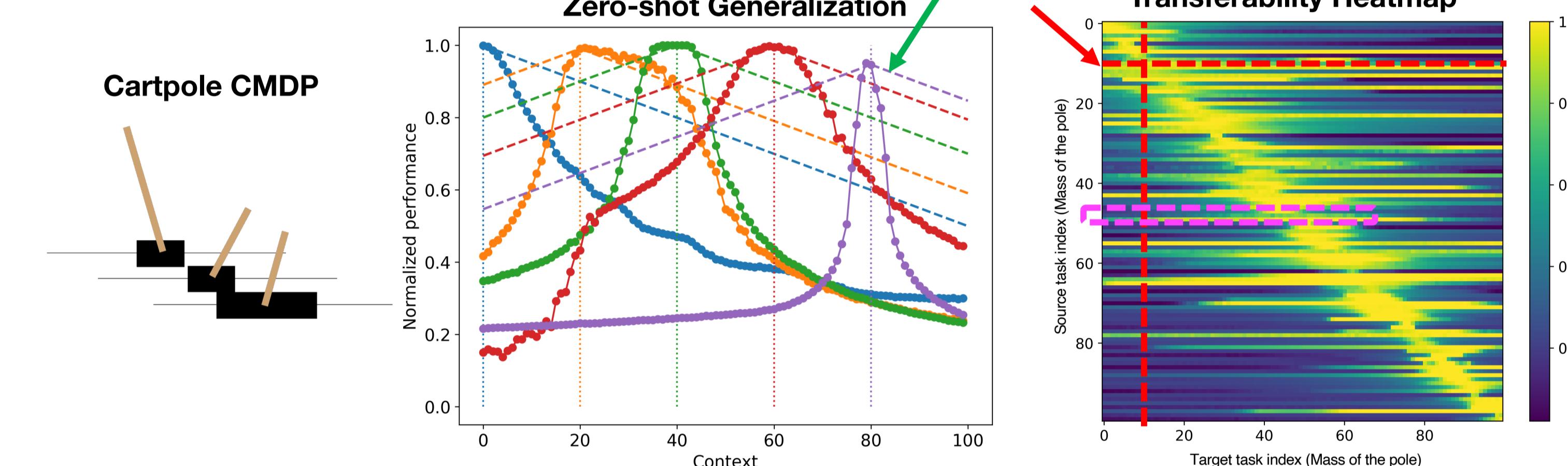
- Transfer Learning:** Utilize already trained policy to solve another task
"Training is expensive, but inference (evaluation) is cheap."
- Zero-shot transfer is **cheap & works remarkably well**



- Challenges:** It's difficult to decide which tasks to train first!
- Algorithmic question:** How to select which task to train?
→ **Source task selection problem**
- Optimize task selection to maximize the performance across contexts

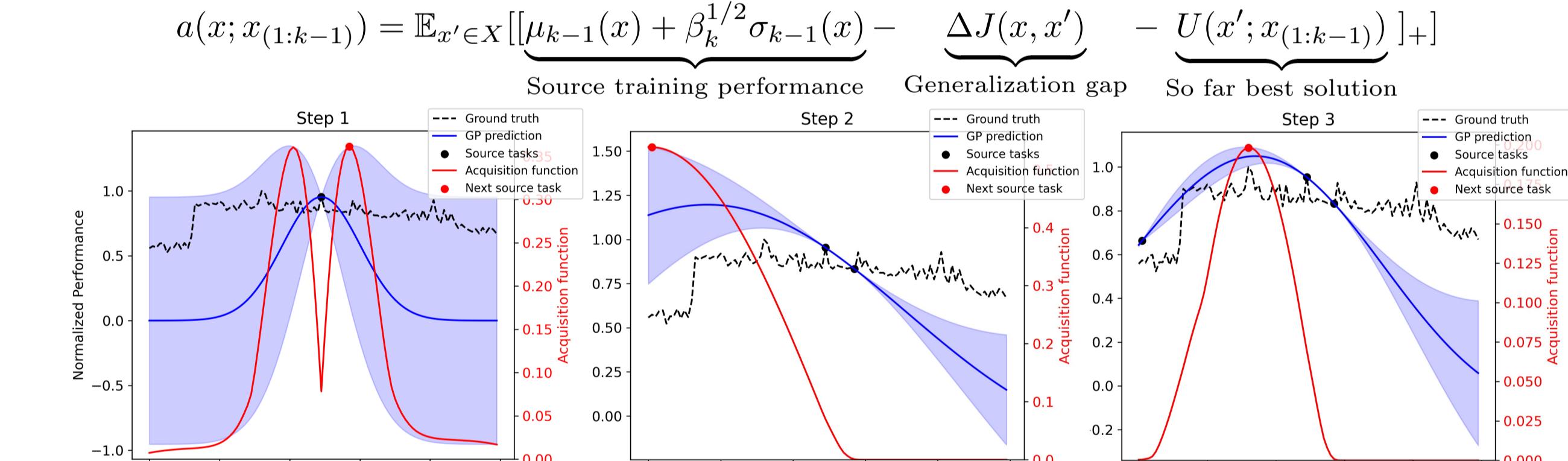
Generalization Gap

- Empirical observation**



MBTL with Bayesian optimization

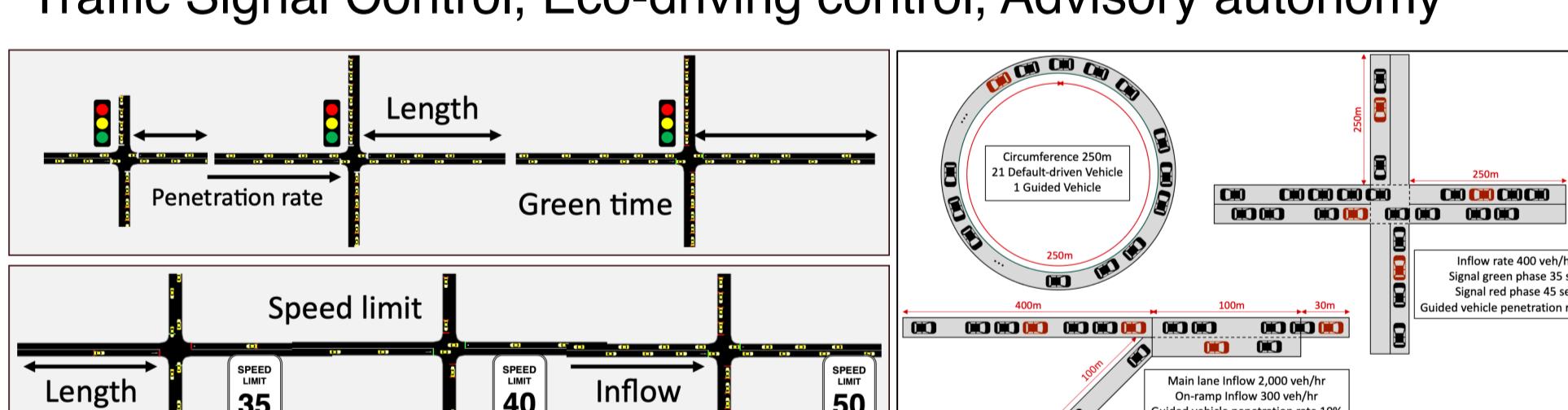
- Challenge:** Training performance is not constant!
- Gaussian Process (GP) to estimate training performance**
- Acquisition function (estimated performance considering gen. gap)



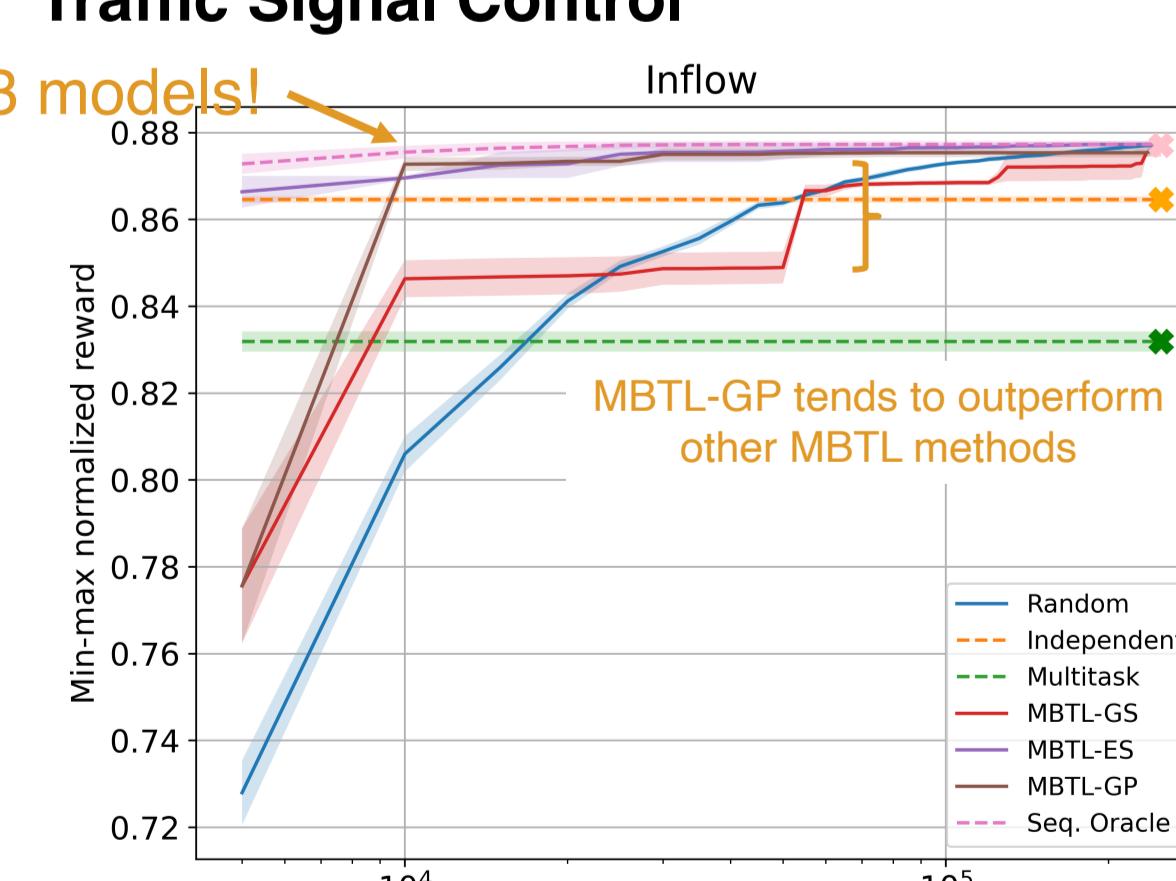
Theorem 3: MBTL-GP achieves **sublinear regret** $R_k \leq O(\sqrt{k \log k})$ with high probability

Transportation CMDP benchmarks

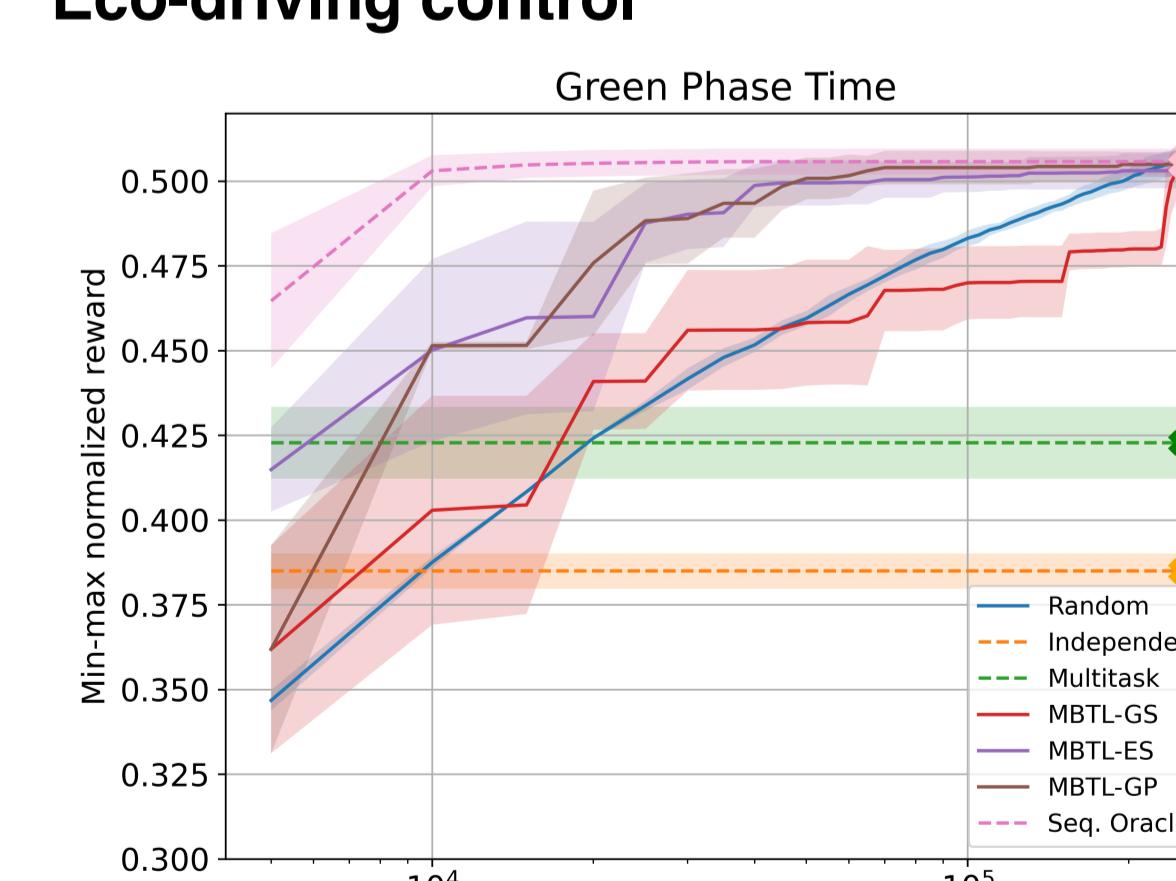
Traffic Signal Control, Eco-driving control, Advisory autonomy



Traffic Signal Control



Eco-driving control



Real-time driving advisory

