

# Generalizing Cooperative Eco-driving via Multi-residual Task Learning

Vindula Jayawardana<sup>†</sup>, Sirui Li<sup>†</sup>, Cathy Wu<sup>†</sup>, Yashar Farid<sup>‡</sup>, Kentaro Oguchi<sup>‡</sup>

<sup>†</sup> MIT <sup>‡</sup> Toyota Motor North America

Correspondence: vindula@mit.edu

## Motivation

- Real-world autonomous driving contends with a multitude of diverse traffic scenarios that are challenging for conventional model-based planning algorithms.
- Model-free deep reinforcement learning (DRL) on the other hand presents a promising avenue to devise control algorithms, but learning DRL controllers that generalize to multiple traffic scenarios is still a challenge.
- In tackling this challenge, we introduce **Multi-residual Task Learning (MRTL)**, a generic learning framework based on multi-task learning that, for a set of task scenarios, decomposes the control into nominal components that are effectively solved by conventional control methods and residual terms which are solved using DRL.

## Problem Formulation

- We study the algorithmic generalization of DRL algorithms across a family of MDPs (scenarios) that originate from a single task.
- Formally, consider a contextual Markov Decision Process (cMDP)  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{C}, p_c, r_c, \rho_c, \gamma \rangle$  which extends Markov Decision Processes (MDP) with a context space  $\mathcal{C}$  (scenarios), and the action space  $\mathcal{A}$  and state space  $\mathcal{S}$  remain unchanged. The transition  $p_c$ , rewards  $r_c$ , and initial state distribution  $\rho_c$  are changed based on the context  $c \in \mathcal{C}$ .
- We seek to find policy  $\pi$  that solve a given cMDP by solving the problem of algorithmic generalization within that task (i.e., finding a policy that performs well in the cMDP overall).

$$\pi^*(s) = \operatorname{argmax}_{\pi} \mathbb{E} \left[ \sum_{c \in \mathcal{C}} \sum_{t=0}^H \gamma^t r_c(s_t, a_t) | s_0^c, \pi \right]$$

## Method

- We introduce **Multi-Residual Task Learning (MRTL)**, a unified learning approach that harnesses the synergy between multi-task learning and residual reinforcement learning.
- We aim to learn the MRTL policy  $\pi(s|c): \mathcal{S} \times \mathcal{C} \rightarrow \mathcal{A}$  by learning a residual function  $f_{\theta}(s|c): \mathcal{S} \times \mathcal{C} \rightarrow \mathcal{A}$  on top of a given nominal policy  $\pi_n(s|c): \mathcal{S} \times \mathcal{C} \rightarrow \mathcal{A}$  such that,

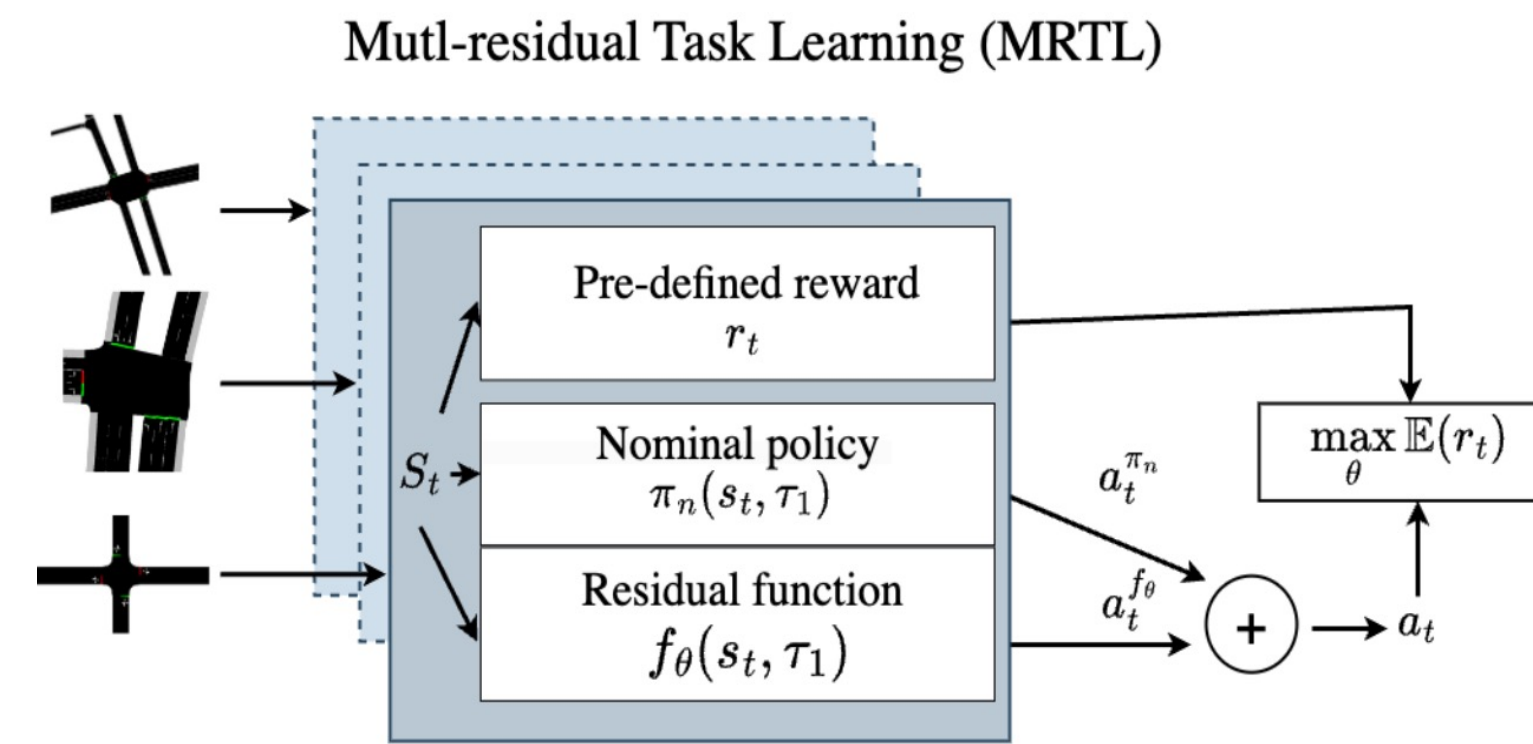
$$\pi(s|c) = \pi_n(s|c) + f_{\theta}(s|c)$$

MRTL policy

Nominal policy

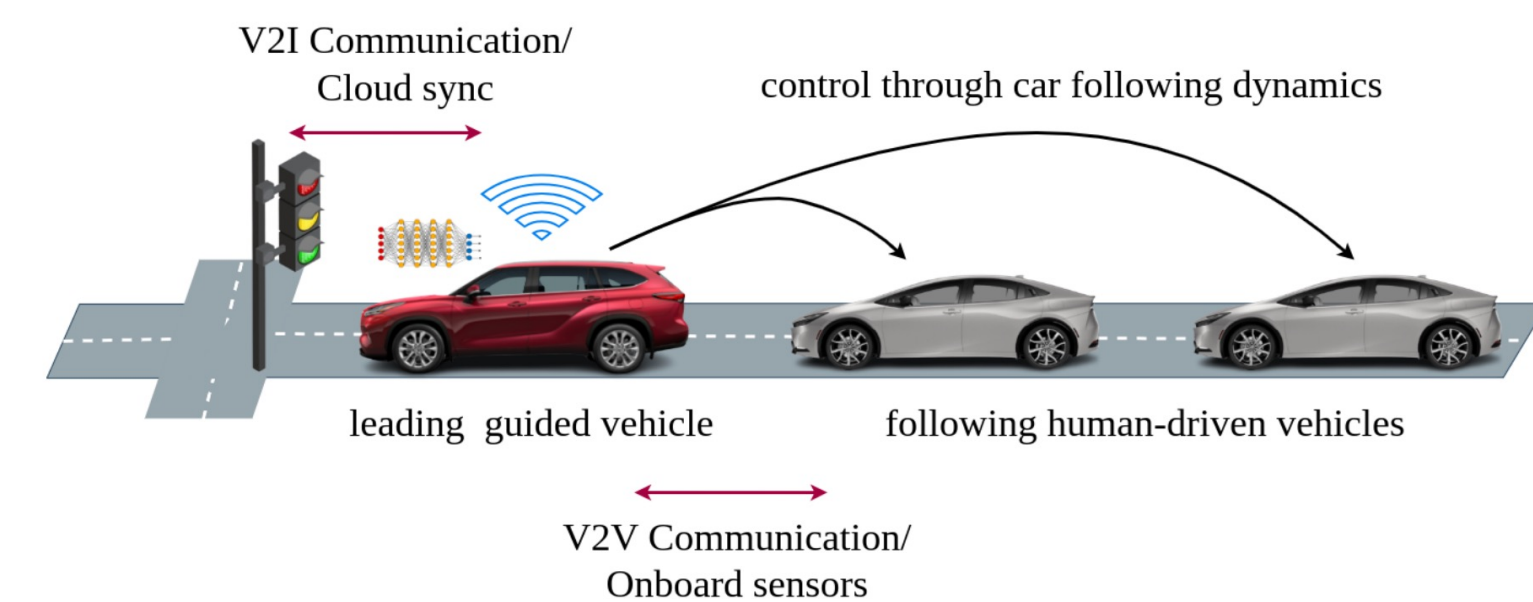
Residual function

- The gradient of the  $\pi$  does not depend on the  $\pi_n$ . This enables flexibility with nominal policy choice.
- Intuition:** If the nominal policy is nearly perfect, the residual term can be viewed as a corrective term. If not, nominal policy provide useful hints to guide the exploration of DRL training.



## Evaluations

- We apply MRTL to cooperative multi-agent eco-driving at signalized intersections.



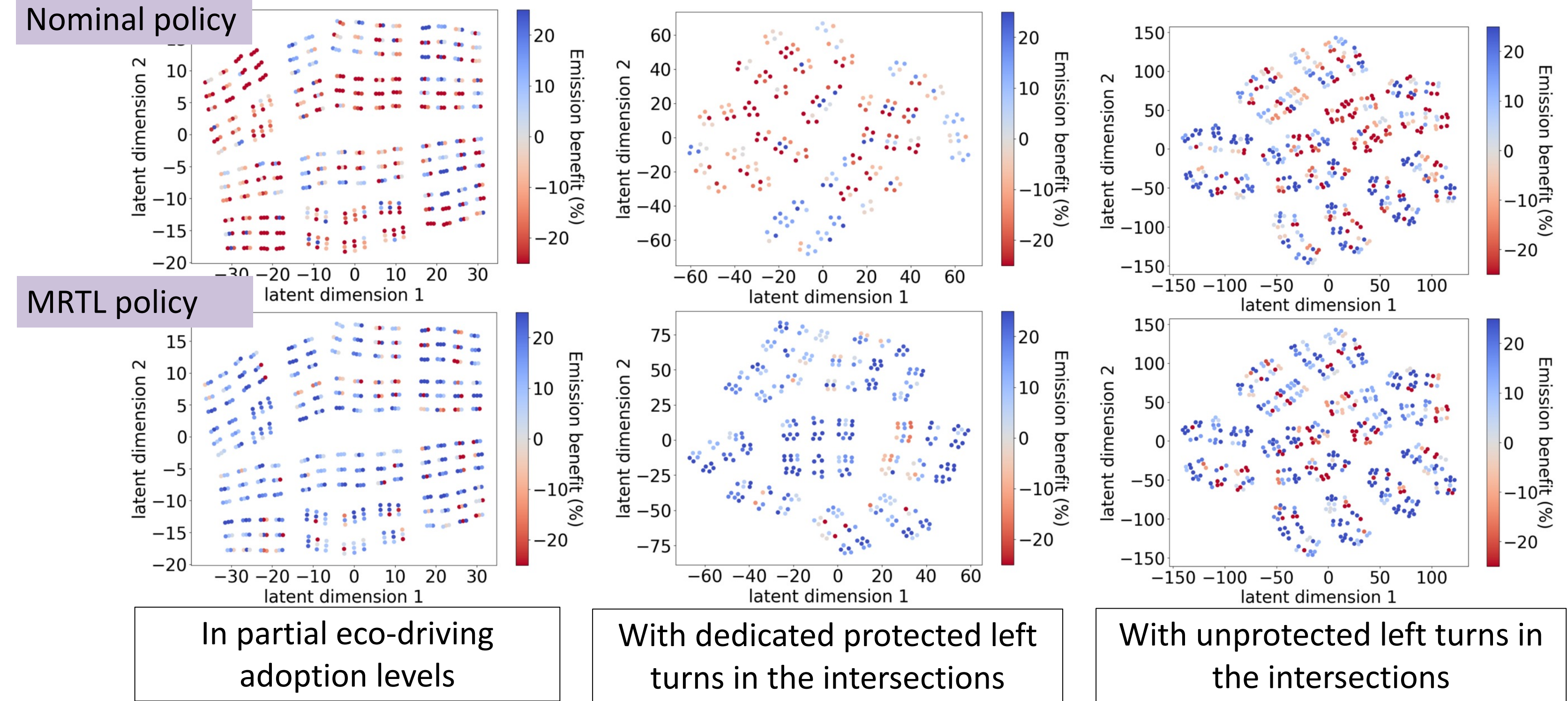
- Goal:** reduce fleet-wide emissions while having less impact on travel time.
- Setting:** 600 signalized intersections synthetically generated to match high-level real world intersection statistics. Both 20% and 100% eco-driving adoption levels were tested.
- Nominal policy:** A model-based heuristic (GLOSA algorithm)
- Baselines:** Multi-task learning from scratch (MTL) and the nominal policy alone (NP)

## Performance comparison across 600 signalized intersections

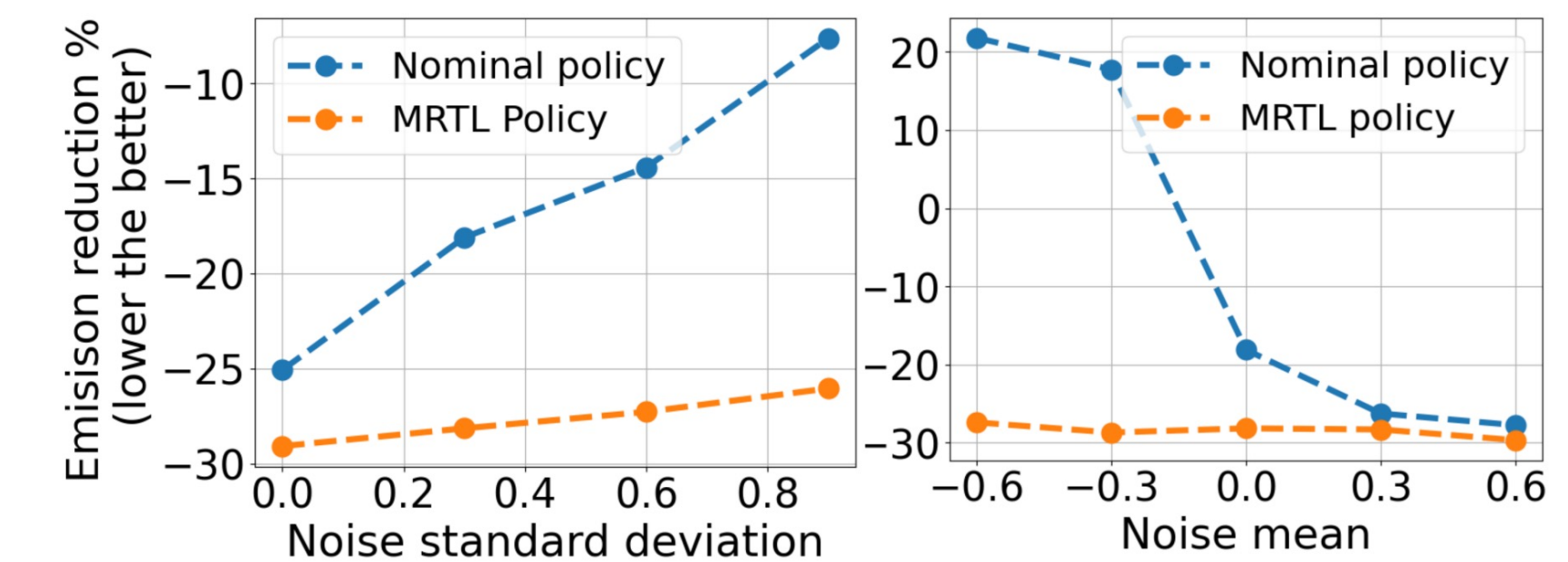
Method	20% penetration			100% penetration		
	Emission ↓	Speed ↑	Throughput ↑	Emission ↓	Speed ↑	Throughput ↑
MTL	64.08%	-27.70%	-34.70%	95.86%	-30.87%	-68.11%
NP	13.13%	-21.11%	-30.07%	-25.09%	11.72%	-3.90%
MRTL (Ours)	<b>-13.95%</b>	<b>12.35%</b>	<b>7.95%</b>	<b>-29.09%</b>	<b>17.10%</b>	<b>5.72%</b>

## Visualization of t-SNE plots illustrating emission benefits in assessing the efficacy of MRTL policy in mitigating nominal policy limitations

- Each dot represents a signalized intersection approach and the colors denote the emission benefit levels.
- Here, higher the emission benefits the better the results.



## Robustness of MRTL to control noise (left) and bias noise (right)



## Takeaway

- Combining conventional control with residual terms learned through DRL is a promising approach to achieve algorithmic generalization in solving contextual Markov decision processes.