

Safe Reinforcement Learning

Minseok Seo¹,

¹Mobility Intelligence and Control Laboratory (MIC Lab)
CCS Graduate School of Mobility
Korea Advanced Institute of Science and Technology (KAIST)



September 5, 2025

1 Introduction

- Motivation
- Reinforcement Learning (RL)

2 Constrained Reinforcement Learning

- Constrained Reinforcement Learning (Constrained RL)
- Constrained Policy Optimization Problem
- State-wise Constrained Policy Optimization

3 Conclusion

- Summary
- Future Work

1

Introduction

- Motivation
- Reinforcement Learning (RL)

2

Constrained Reinforcement Learning

-
-

3

Conclusion

-



Figure: Waymo and Tesla

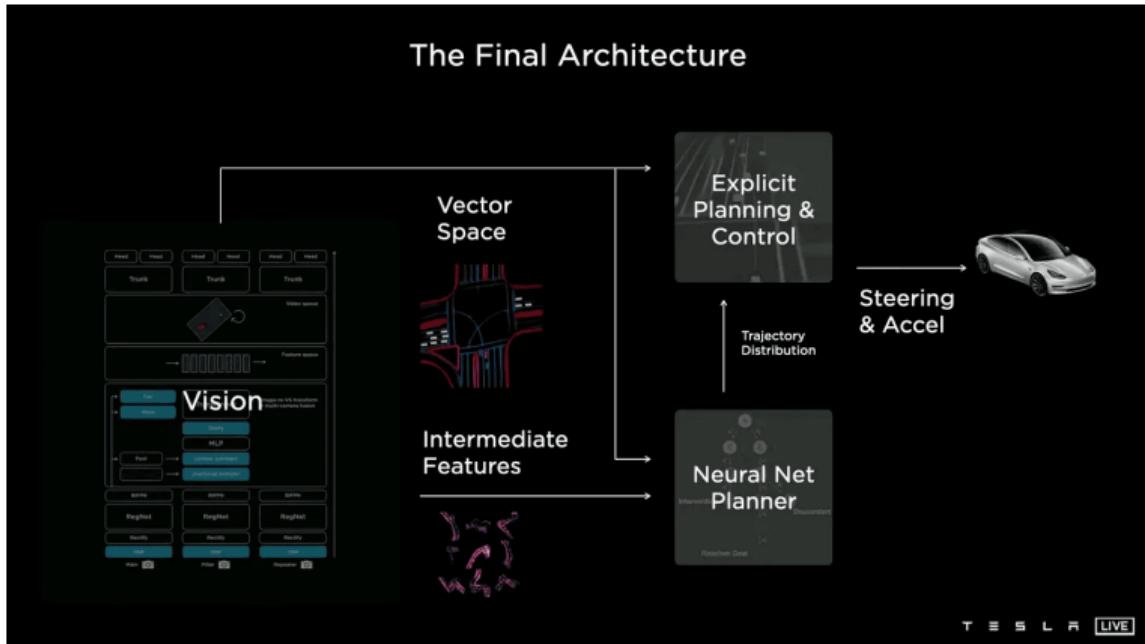


Figure: Tesla's Architecture in 2021 (source: AI Day 2021)

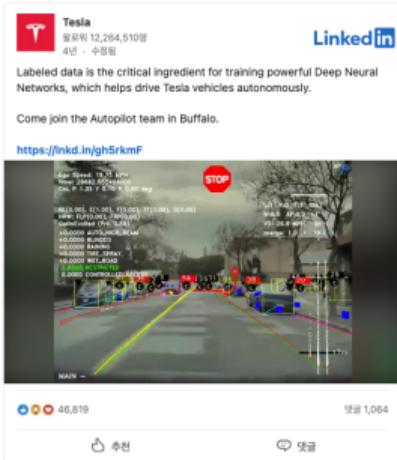


Figure: Tesla's recruitment post for data labeling positions

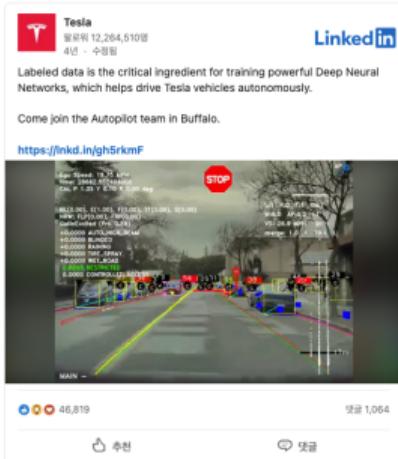


Figure: Tesla's recruitment post for data labeling positions

- Supervised learning requires a large amount of labeled data.

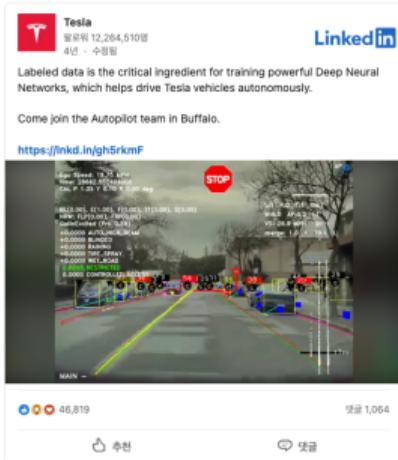


Figure: Tesla's recruitment post for data labeling positions

- Supervised learning requires a large amount of labeled data.
- Since it is created by humans, it is expensive.

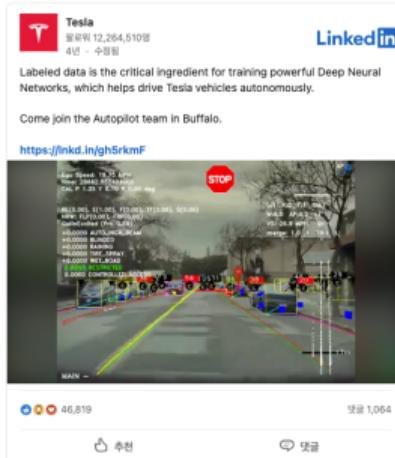


Figure: Tesla's recruitment post for data labeling positions

- Supervised learning requires a large amount of labeled data.
- Since it is created by humans, it is expensive.
- The performance of supervised learning is limited by human-labeled data.

Alpha Go [1]



Deep Q-Network [2]

Proximal Policy Optimization [3]

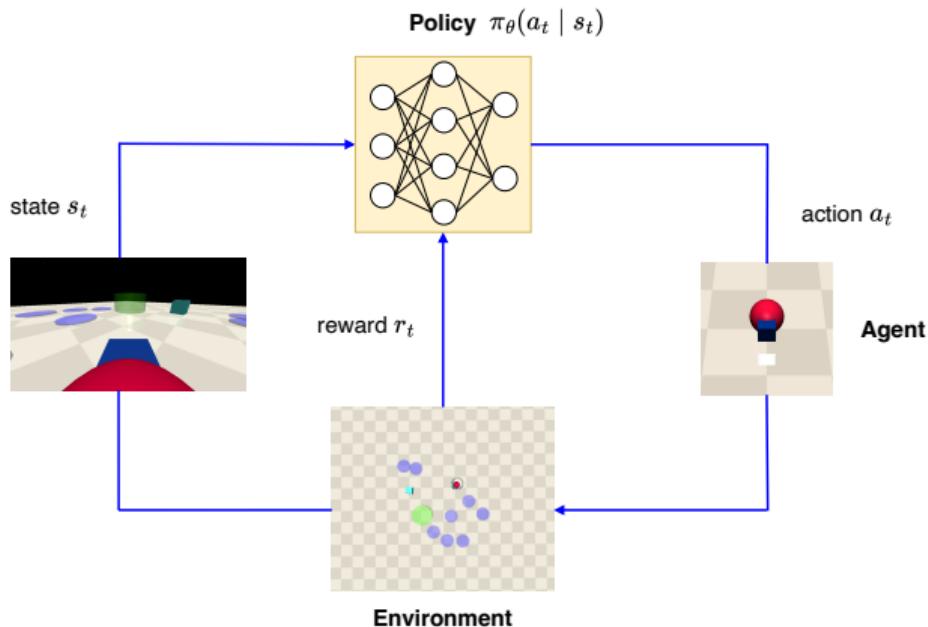


Figure: Overview of the reinforcement learning framework

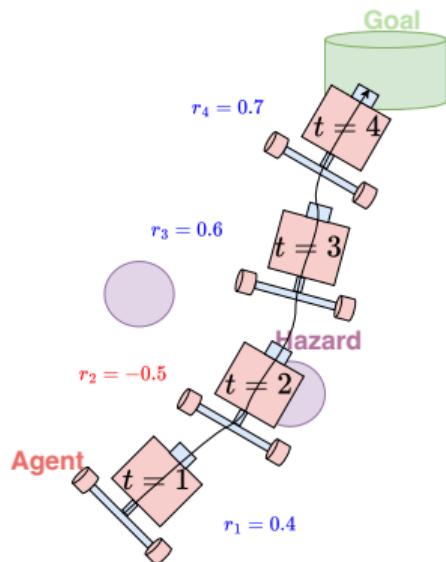


Figure: Illustration of sequential decision making: the agent receives a reward after each action.

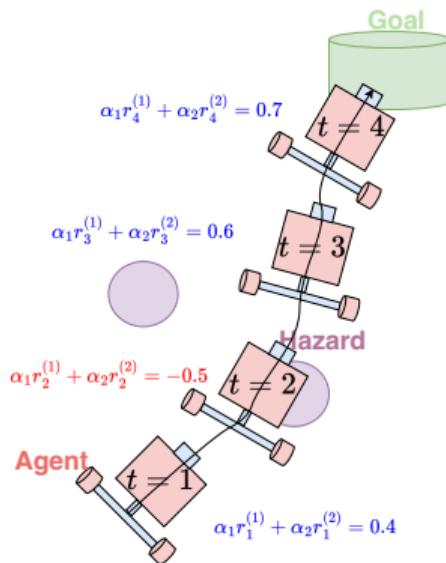


Figure: Illustration of reward engineering: shaping the reward signal to guide the agent's learning process.

- Policy parameterized by θ , denoted as $\pi_\theta(a|s)$
- Goal: find the optimal policy π_θ^* that maximizes the expected cumulative reward

$$\begin{aligned}\theta^* &= \arg \max_{\theta} J(\theta) \\ J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right]\end{aligned}\tag{1}$$

1. Needs **simulation**, not labeled data

1. Needs **simulation**, not labeled data
2. Works in simulation, but may fail in the real world
 - Even after training, the agent may still take **unsafe actions**

1. Needs **simulation**, not labeled data
2. Works in simulation, but may fail in the real world
 - Even after training, the agent may still take **unsafe actions**
3. **Reward engineering** can induce desired behaviors (multi-objective), but it is time-consuming and difficult

1. Needs **simulation**, not labeled data
2. Works in simulation, but may fail in the real world
 - Even after training, the agent may still take **unsafe actions**
3. **Reward engineering** can induce desired behaviors (multi-objective), but it is time-consuming and difficult

My research

Challenges 2 & 3: learning safe policies without relying on reward engineering

1

Introduction

2

Constrained Reinforcement Learning

3

Conclusion

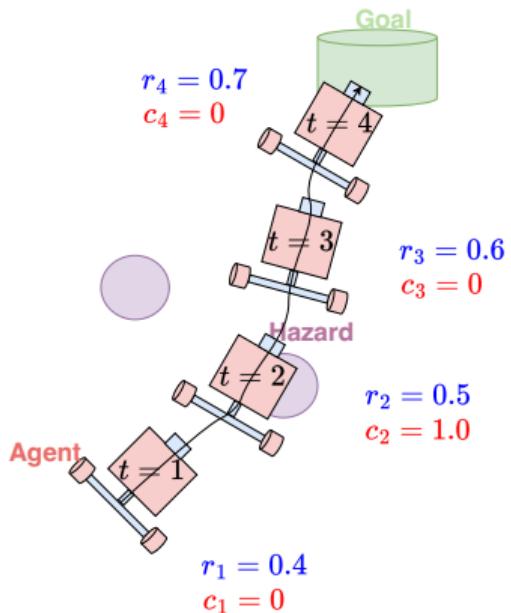


Figure: Constrained RL example — agent trajectory with rewards (blue) and costs (red)

Constrained Policy Optimization Problem

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T c_t \right] \leq d \end{aligned} \tag{2}$$

$$\theta^* = \arg \max_{\theta} J(\theta)$$
$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T c_t \right] \leq d \quad (2)$$

- Find a policy that maximizes rewards while satisfying constraints.

Constrained Policy Optimization Problem

$$\theta^* = \arg \max_{\theta} J(\theta)$$
$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T c_t \right] \leq d \quad (2)$$

- Find a policy that maximizes rewards while satisfying constraints.
- Directly solving constrained optimization is difficult.

$$\theta^* = \arg \max_{\theta} J(\theta)$$
$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T c_t \right] \leq d \quad (2)$$

- Find a policy that maximizes rewards while satisfying constraints.
- Directly solving constrained optimization is difficult.
- By applying Lagrangian relaxation, we can convert it to an unconstrained problem.

$$\theta^* = \arg \max_{\theta} \mathcal{L}(\theta, \lambda)$$

$$\mathcal{L}(\theta, \lambda) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] - \lambda \left(\mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T c_t \right] - d \right) \quad (3)$$

$$\lambda \leftarrow \left[\lambda + \beta (\hat{J}_c - d) \right]_+ \quad (4)$$

Penalty (λ) increases when constraints are violated \rightarrow policy is encouraged to satisfy them.

In CMDPs [4], constraints are imposed on the trajectory level.

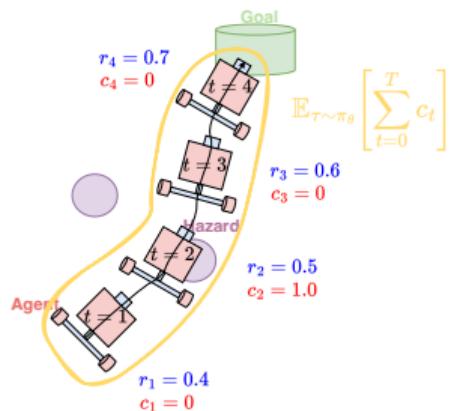


Figure: CMDP example — agent trajectory with rewards (blue) and costs (red)

In CMDPs [4], constraints are imposed on the trajectory level.

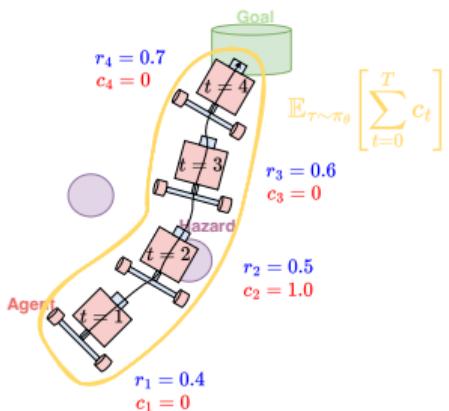


Figure: CMDP example — agent trajectory with rewards (blue) and costs (red)

Wouldn't imposing constraints at the state level allow for more precise constraint enforcement?

In state-wise constrained MDPs, constraints are imposed on each state.

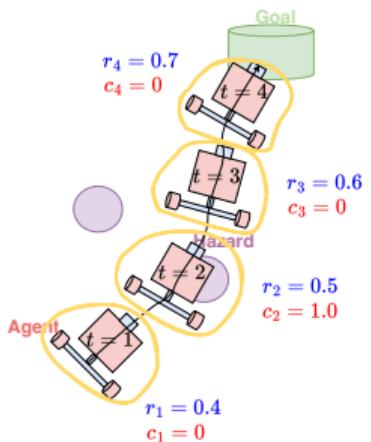


Figure: State-wise constrained MDP example — agent trajectory with rewards (blue) and costs (red)

$$\begin{aligned}\pi^* &= \arg \max_{\pi_\theta} J(\theta) \\ J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} [c(s, a)] \leq w, \quad \forall s \in S\end{aligned}\tag{5}$$

$$\lambda(s) \leftarrow \lambda(s) + \beta(\hat{J}_c - w) \tag{6}$$

$$\pi^* = \arg \max_{\pi_\theta} J(\theta) \quad (5)$$

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T r_t \right] \text{ subject to } \mathbb{E}_{\tau \sim \pi_\theta} [c(s, a)] \leq w, \quad \forall s \in S$$

$$\lambda(s) \leftarrow \lambda(s) + \beta(\hat{J}_c - w) \quad (6)$$

The Lagrange multiplier is replaced by a neural network's output instead of a scalar.

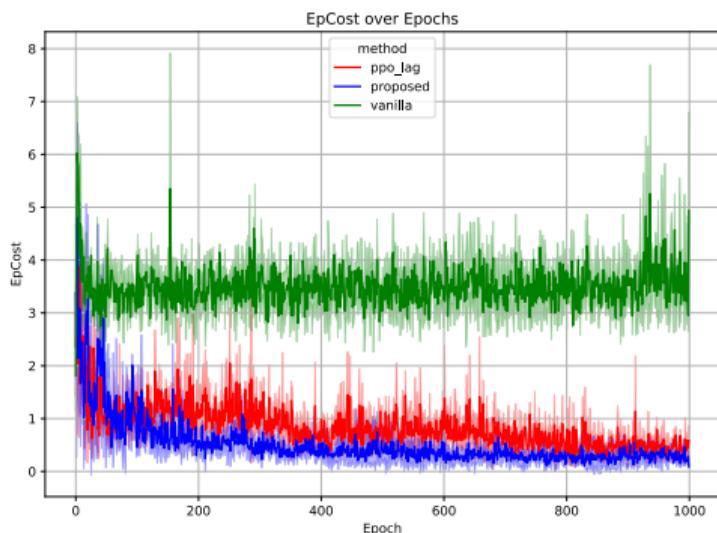
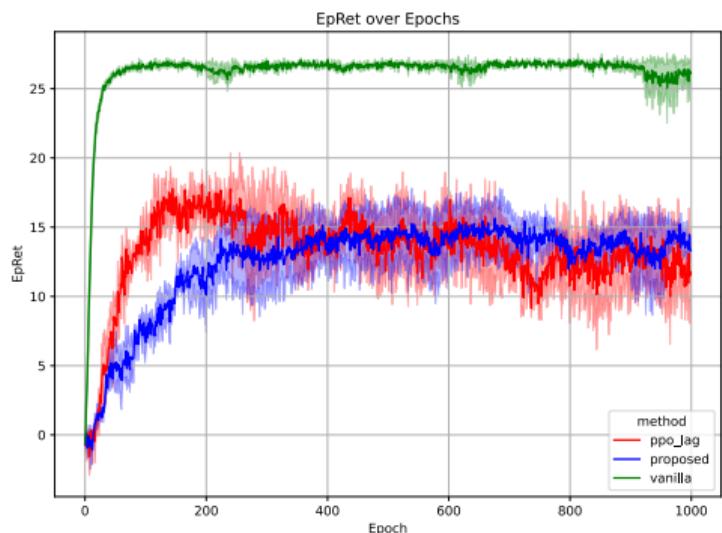


Figure: Performance comparison on Safety Gym tasks

PPO

PPO Lagrangian

Proposed Method

Limitations of the Proposed Approach

- Sensitivity to Lagrange multiplier (init., learning rate)

Limitations of the Proposed Approach

- Sensitivity to Lagrange multiplier (init., learning rate)
- Constraint threshold setting is difficult
 - Too strict → overly conservative policy
 - Too lenient → constraints not enforced

Limitations of the Proposed Approach

- Sensitivity to Lagrange multiplier (init., learning rate)
- Constraint threshold setting is difficult
 - Too strict → overly conservative policy
 - Too lenient → constraints not enforced
- Needs sufficient violations during training

1

Introduction

2

Constrained Reinforcement Learning

3

Conclusion

Summary

Future Work

- Major companies focus on supervised learning for autonomous driving

Summary

- Major companies focus on supervised learning for autonomous driving
- Supervised learning is costly and data-dependent

- Major companies focus on supervised learning for autonomous driving
- Supervised learning is costly and data-dependent
- Reinforcement learning can learn policies in simulation and transfer to the real world

- Major companies focus on supervised learning for autonomous driving
- Supervised learning is costly and data-dependent
- Reinforcement learning can learn policies in simulation and transfer to the real world
- Risk of unsafe behavior due to the Sim2Real gap

- Major companies focus on supervised learning for autonomous driving
- Supervised learning is costly and data-dependent
- Reinforcement learning can learn policies in simulation and transfer to the real world
- Risk of unsafe behavior due to the Sim2Real gap
- Research direction: state-wise constrained policy learning

- Major companies focus on supervised learning for autonomous driving
- Supervised learning is costly and data-dependent
- Reinforcement learning can learn policies in simulation and transfer to the real world
- Risk of unsafe behavior due to the Sim2Real gap
- Research direction: state-wise constrained policy learning

But there are still challenges to solve..

- Extension to Model-Based RL

Thank you for your attention!

- [1] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [3] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [4] E. Altman, *Constrained Markov decision processes*. Routledge, 2021.