

DDA 6201 Online Decision-Making Lecture 12

Application 2: Value-Based RL



Tongxin Li

School of Data Science

The Chinese University of Hong Kong (Shenzhen)

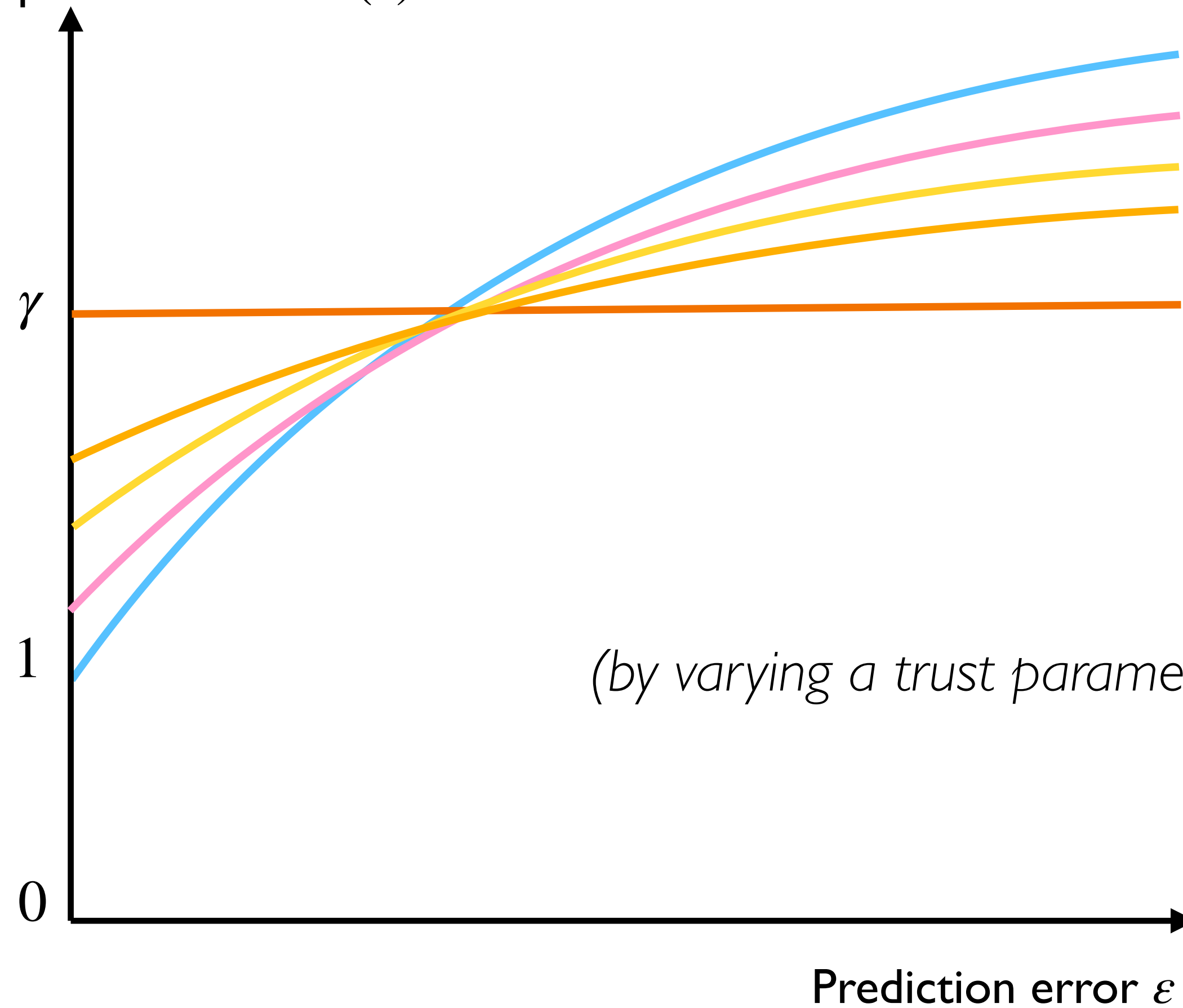
Revisit: Learning-Augmented Algorithms

Performance Benchmark

e.g. Competitive ratio $CR(\varepsilon)$

Meta-algorithms

Consistency vs Robustness Trade-off



Consistent ML Algorithm (Good when ε is small)

Intermediate Regimes

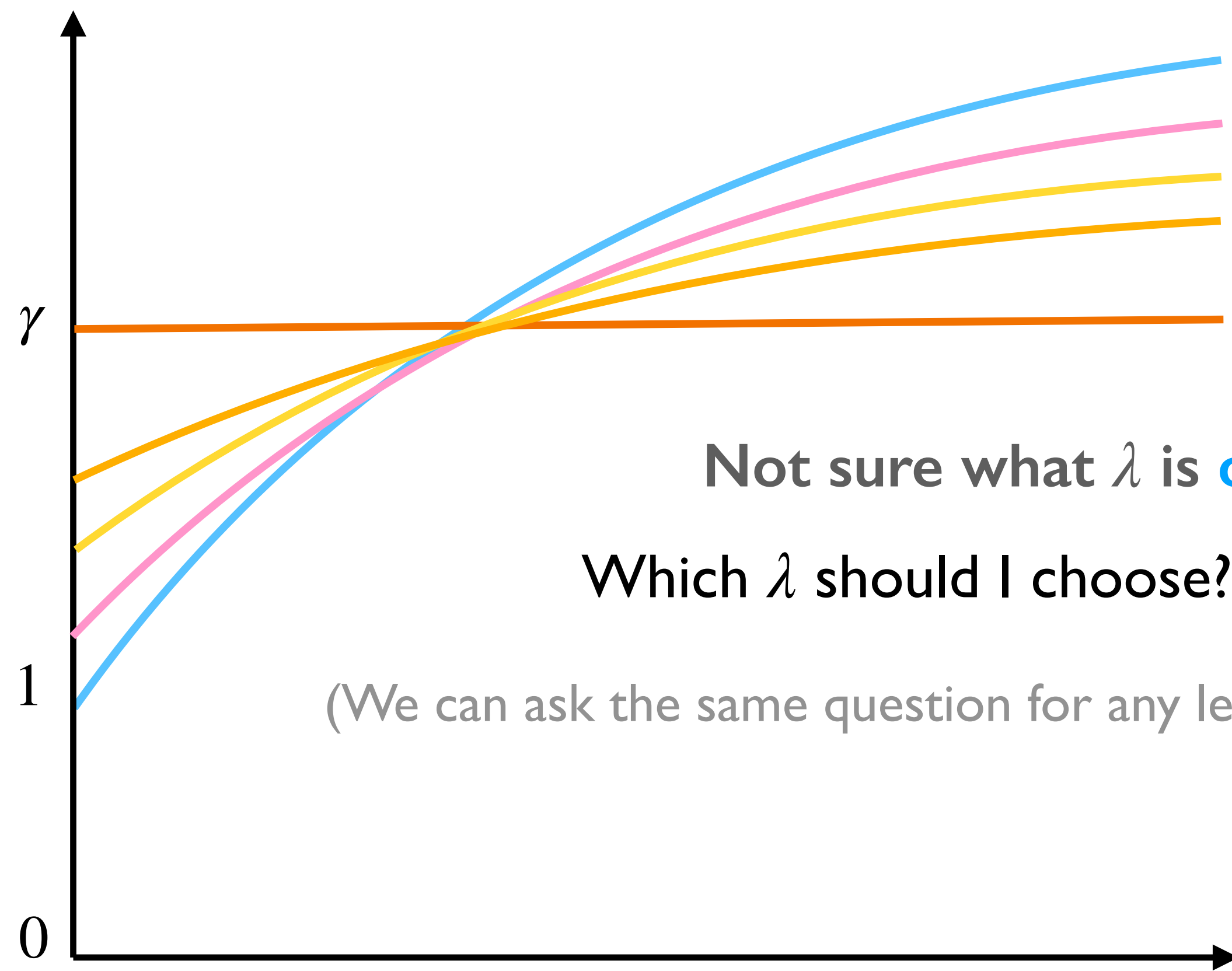
Robust Classic Algorithm (Good when ε is large)

First Limitation

General Goal of Learning-Augmented Algorithms

Consistency vs Robustness Trade-off

Competitive ratio $CR(\varepsilon)$



$\lambda = 1$

$\lambda = 0.7$

$\lambda = 0.5$

$\lambda = 0.2$

$\lambda = 0$

Not sure what λ is **optimal** ...

Which λ should I choose? (ε is unknown)

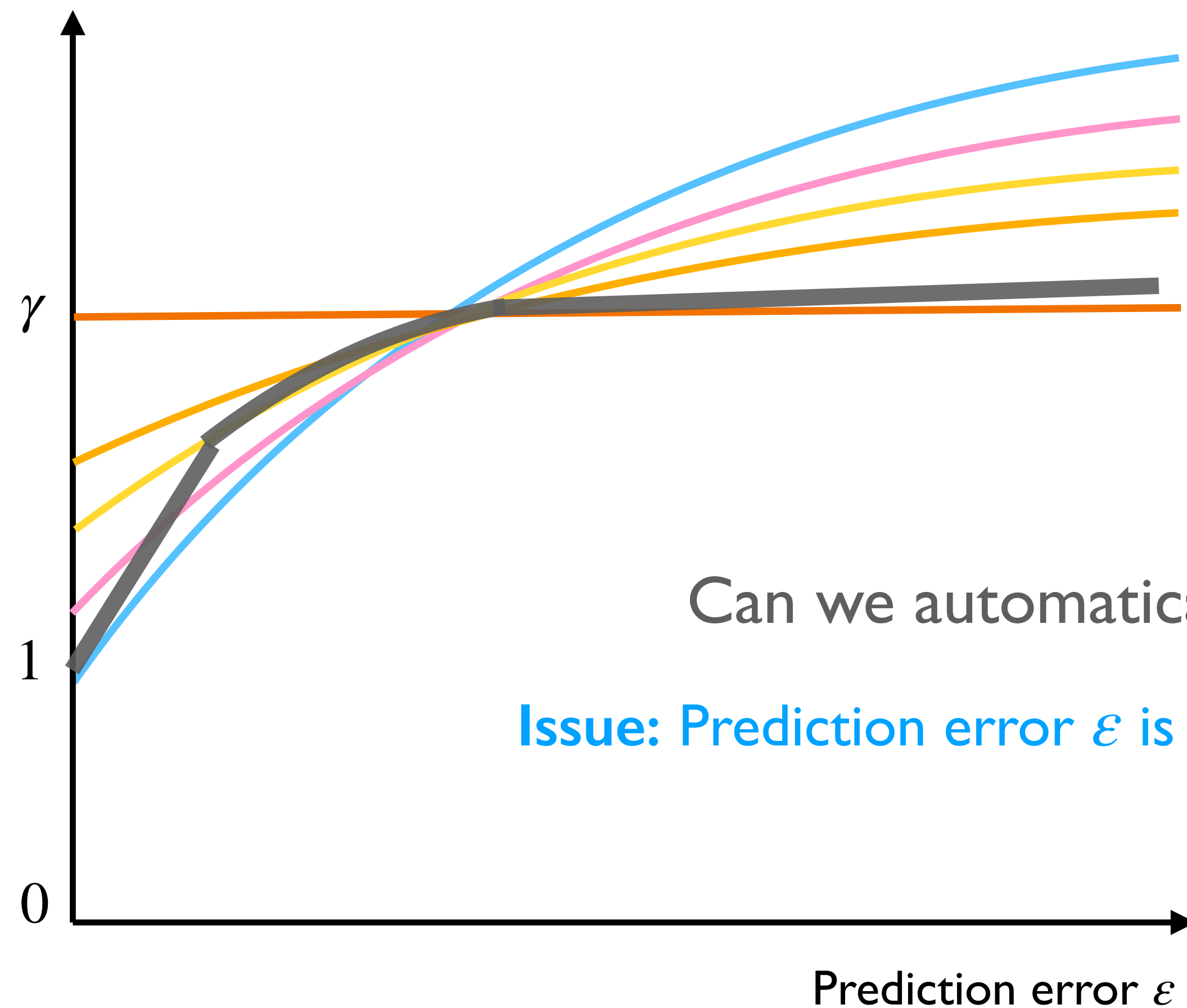
(We can ask the same question for any learning-augmented online algorithms)

Prediction error ε

First Limitation

Goal: Find an online algorithm with good Competitive Ratio **CR** regardless of prediction error ε

Competitive ratio $CR(\varepsilon)$



$$\lambda = 1$$

$$\lambda = 0.7$$

$$\lambda = 0.5$$

$$\lambda = 0.2$$

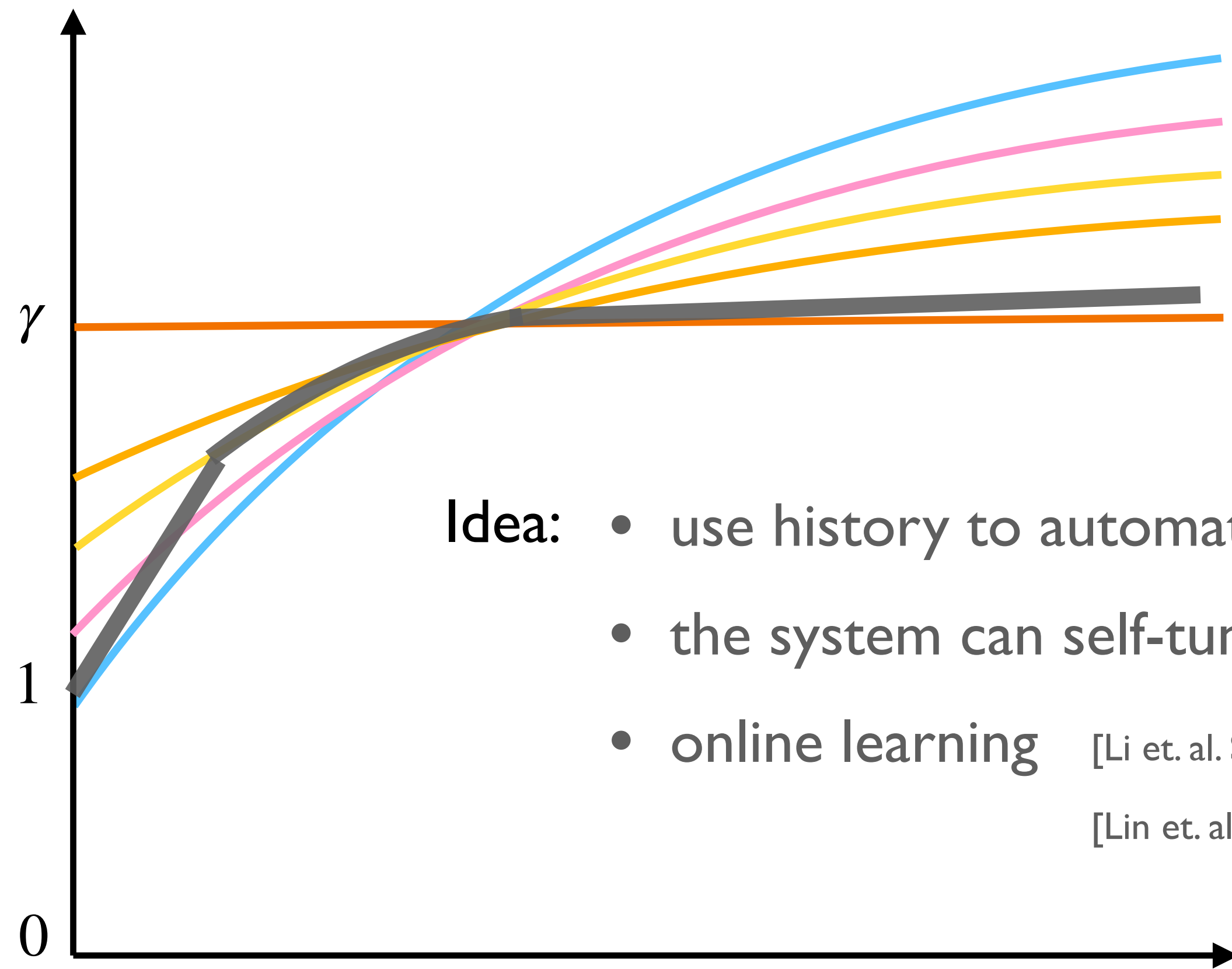
$$\lambda = 0$$

One Solution: Online Learning

General Goal of Learning-Augmented Algorithms

Consistency vs Robustness Trade-off

Competitive ratio $CR(\varepsilon)$



Idea: • use history to automatically select λ

• the system can self-tune

• online learning [Li et. al. SIGMETRICS 2022] [Khodak et. al. NeurIPS 2022]

[Lin et. al. Preprint 2023] ... [Li et. al. NeurIPS 2024]

A Real-World Problem



Image: Paired Power



The UC San Diego/EVgo project

EV Charging with Uncertainties

Use RL for scheduling?

- Tons of existing policies

Question: Do they work well in practice?

Question: If so, why is it hard to see them being used?

Ideally, they work well, but ...

Main Sources of Uncertainties → Data

Electricity price varies Solar generation Charging behavior



The UC San Diego/EVgo project

Large-Scale Adaptive Charging Network

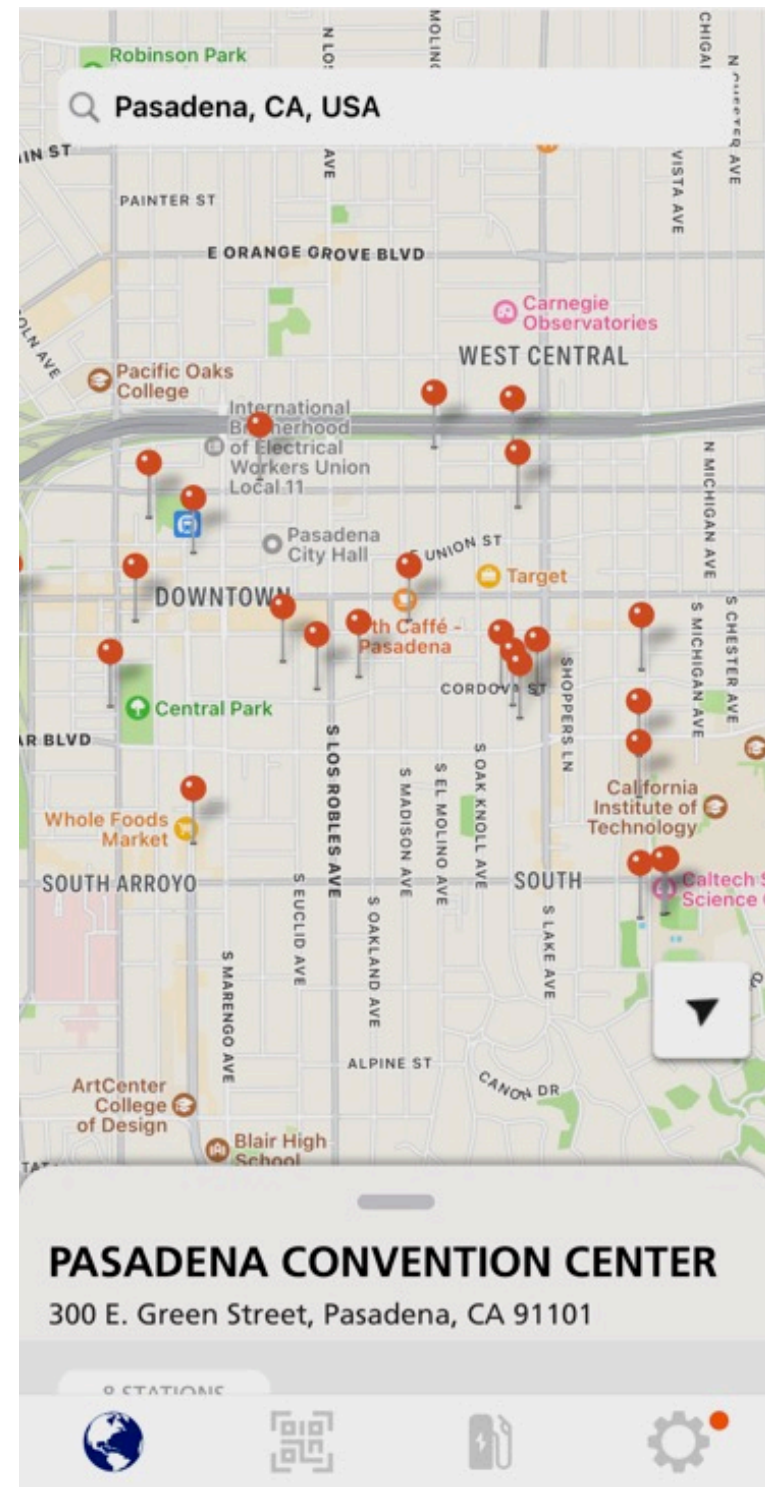


Adaptive Charging Network (ACN@Caltech)

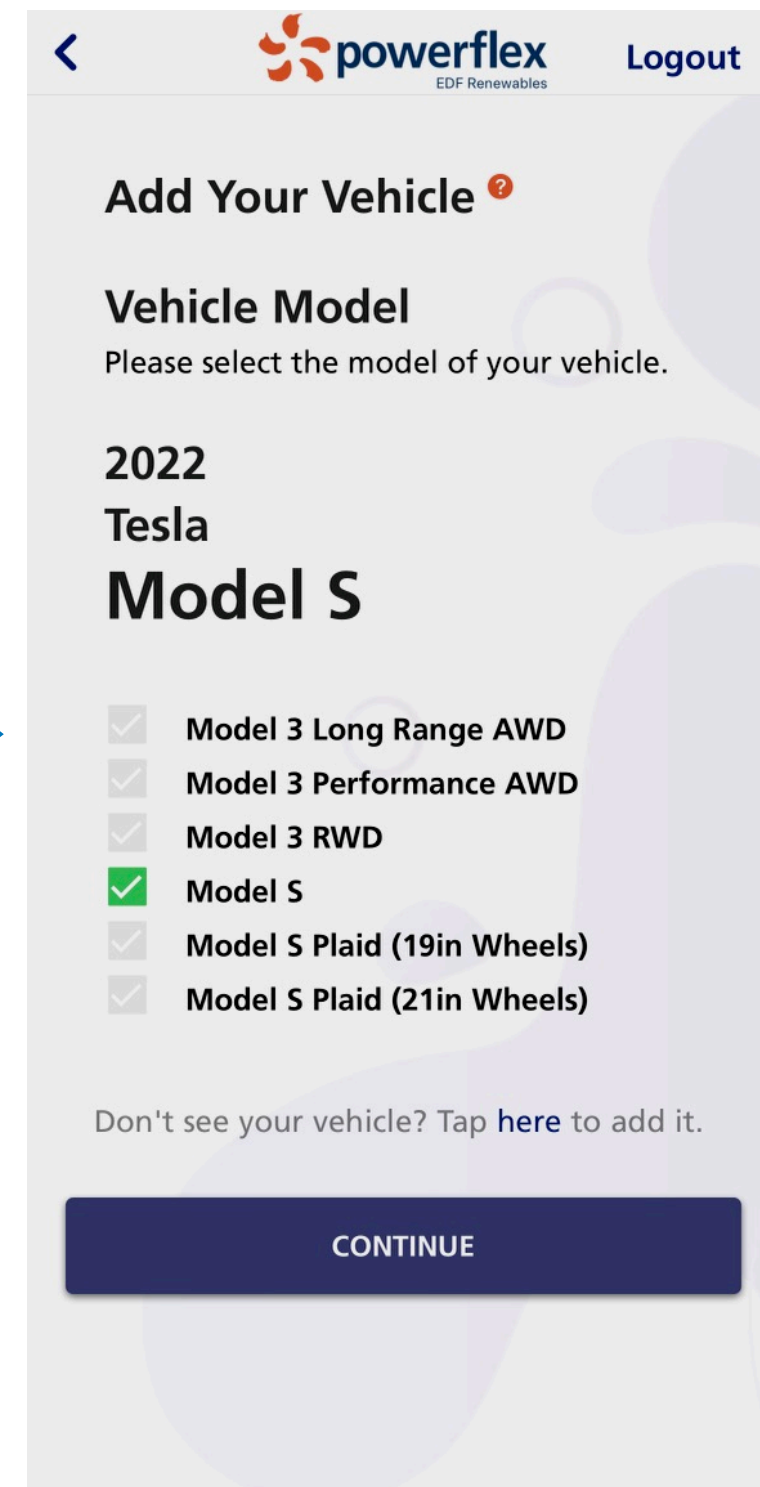
- A Parking lot with 54 chargers
- How to schedule EV charging is challenging



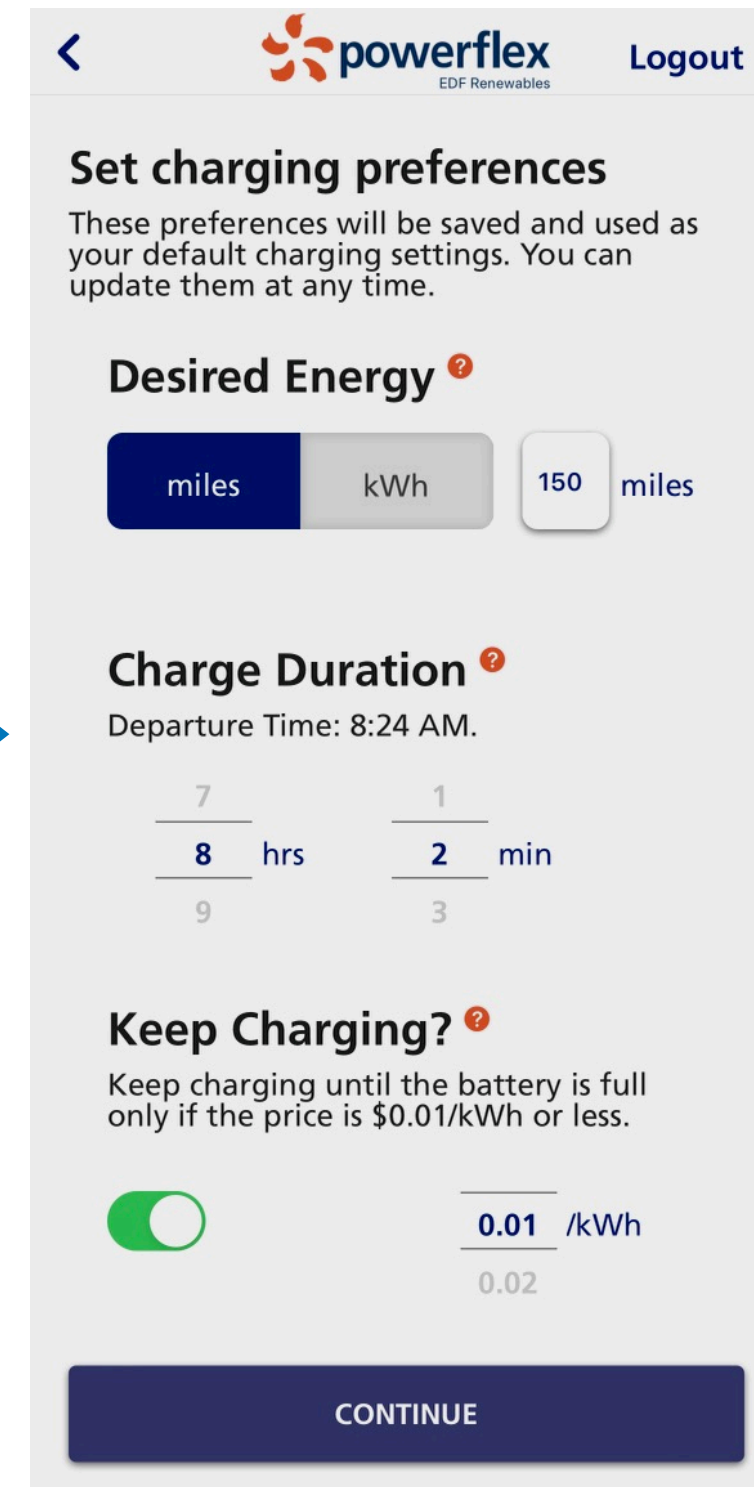
Large-Scale Workplace EV Charging



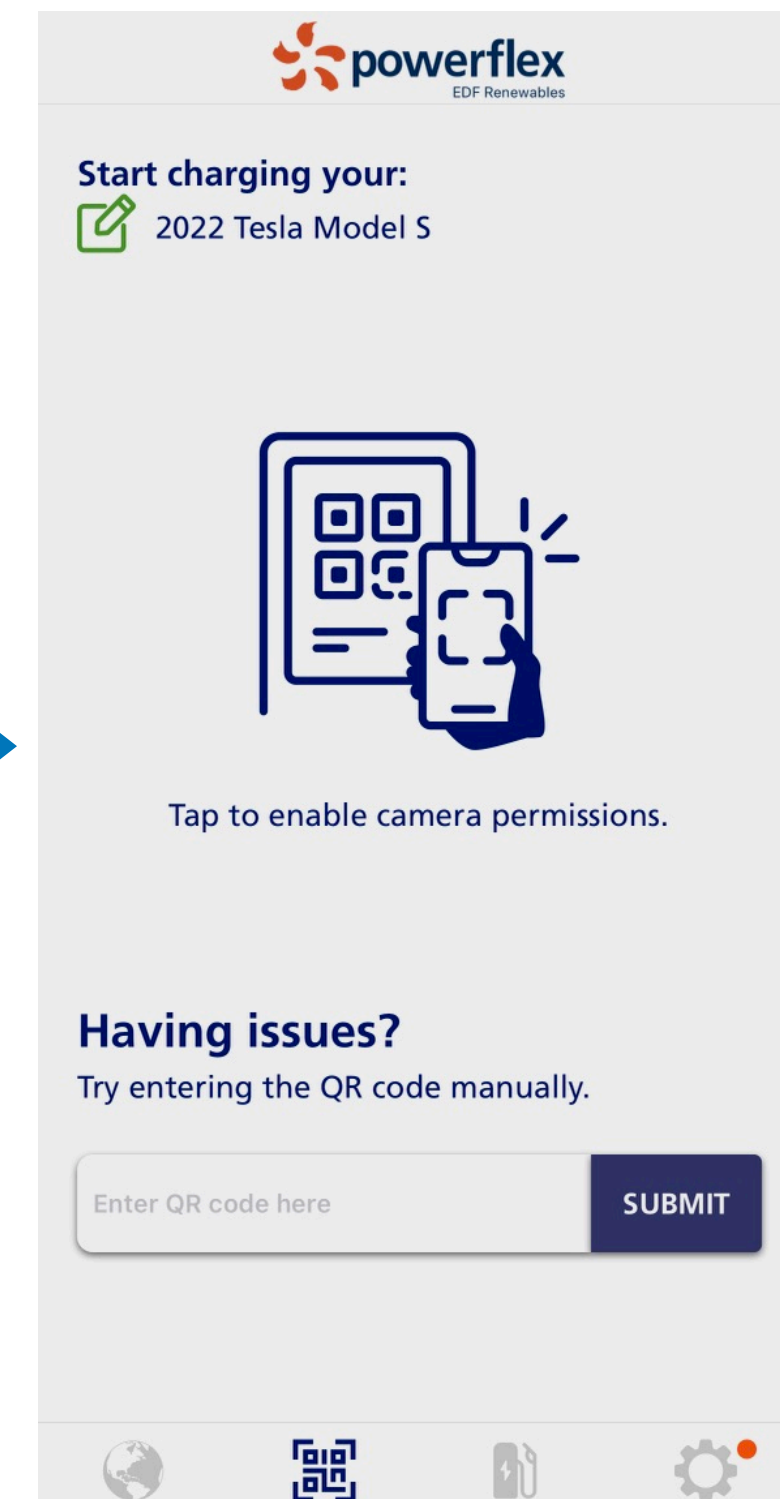
Fina a Location



Choose Model



Select Energy



Scan and Charge

Large-Scale Workplace EV Charging

Classic Scheduling Algorithms

- Least laxity first (LLF)
 - Earliest deadline first (EDF)
 - Model predictive control (MPC)
- (Currently used in Caltech ACNs)

...

Set charging preferences
These preferences will be saved and used as your default charging settings. You can update them at any time.

Desired Energy [?]

miles kWh miles

Charge Duration [?]
Departure Time: 8:24 AM.

hrs min

Keep Charging? [?]
Keep charging until the battery is full only if the price is \$0.01/kWh or less.

/kWh

Large-Scale Workplace EV Charging

Classic Scheduling Algorithms

- Least laxity first (LLF)
 - Earliest deadline first (EDF)
 - Model predictive control (MPC)
- (Currently used in Caltech ACNs)

...

Set charging preferences

These preferences will be saved and used as your default charging settings. You can update them at any time.

Desired Energy [?]

miles kWh miles

Charge Duration [?]

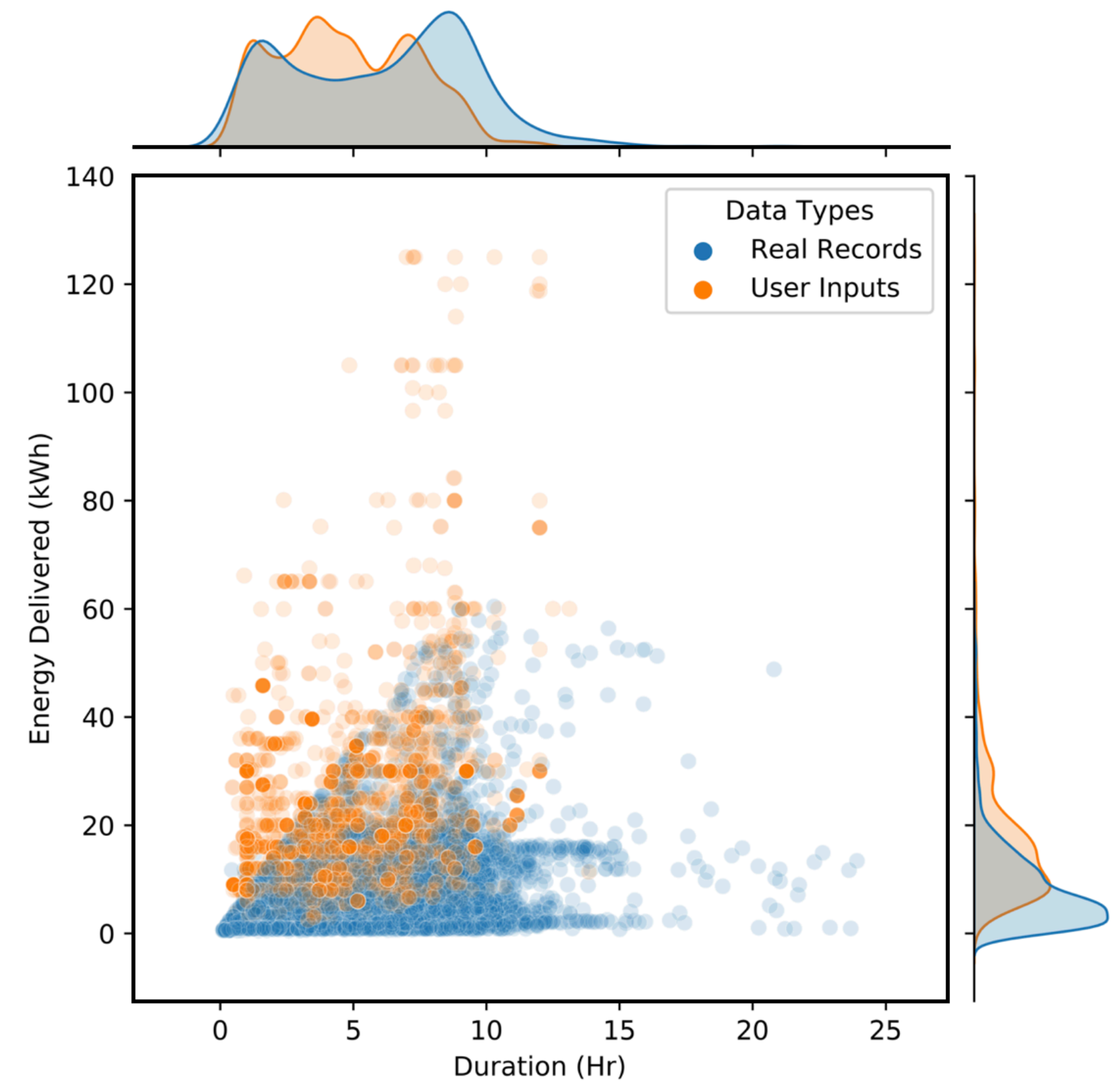
Departure Time: 8:24 AM.

hrs min

Keep Charging? [?]

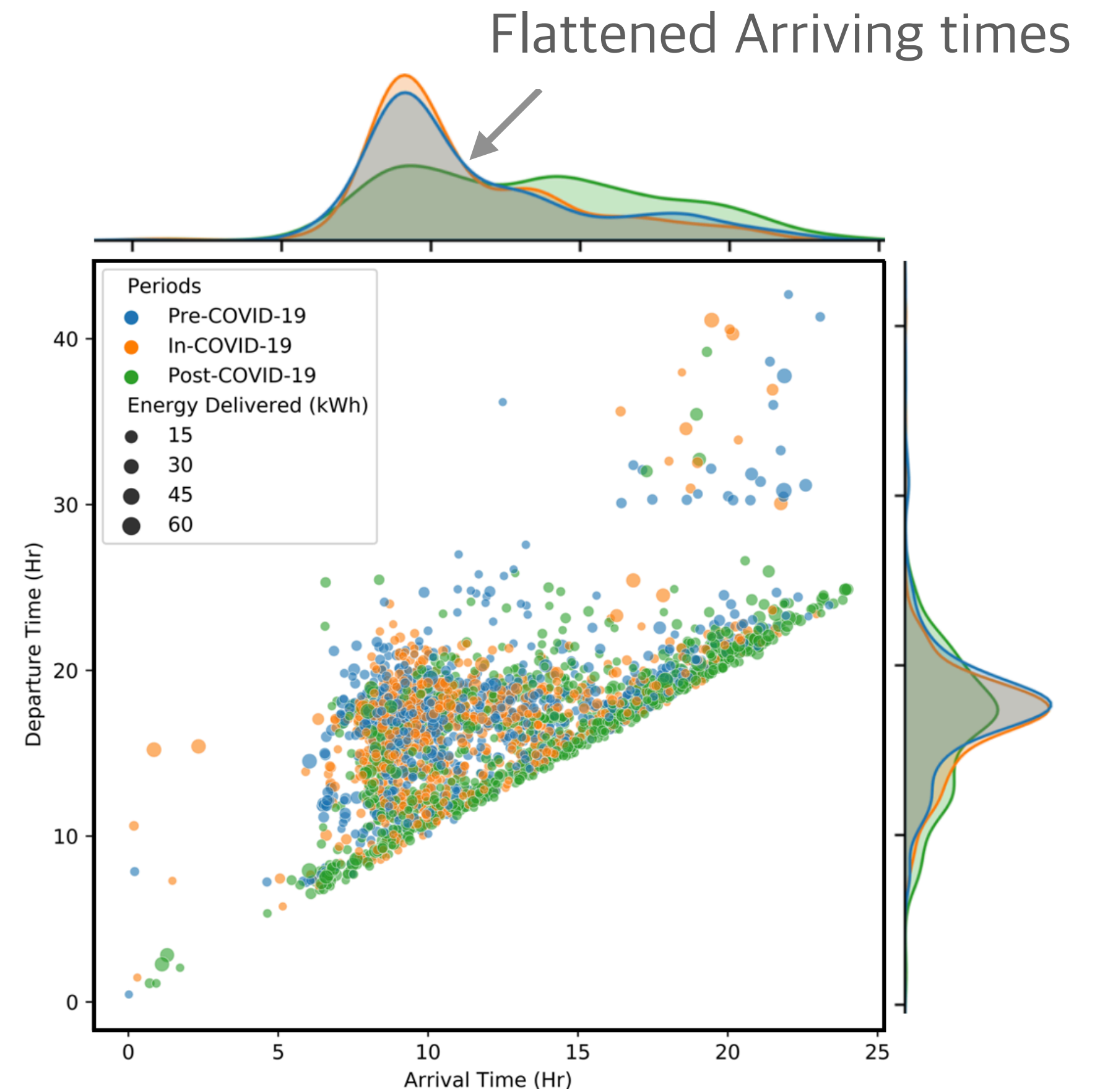
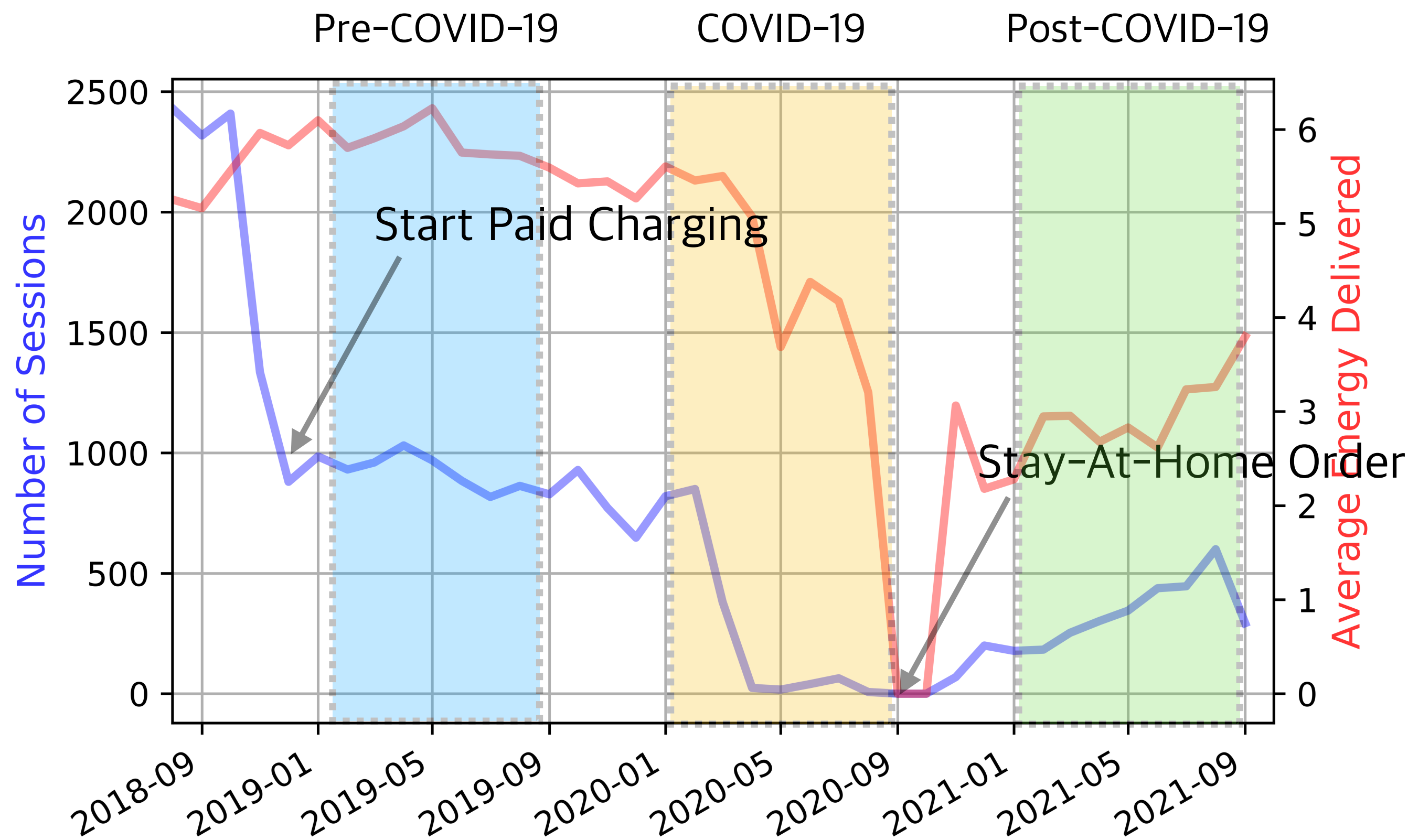
Keep charging until the battery is full only if the price is \$0.01/kWh or less.

/kWh



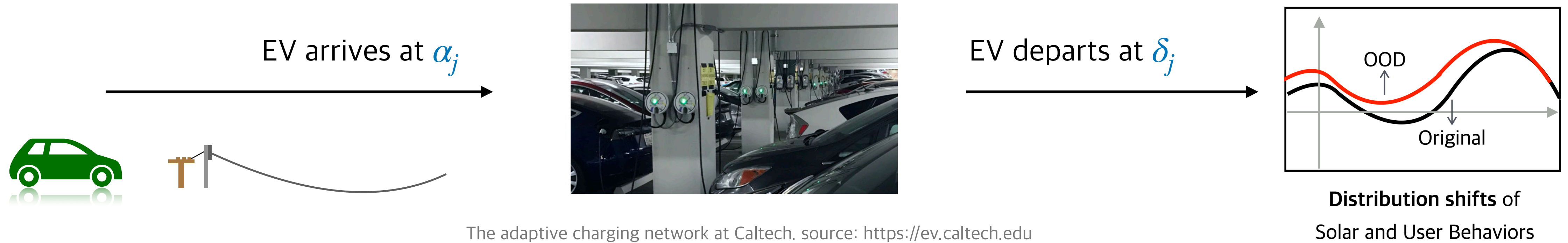
COVID-19 Caused Dataset Shift

Statistical distributions shifted (Data from the real Caltech system) [e-Energy '19]



RL policies trained on out-of-distribution data can perform poorly

Combing MPC and RL Scheduling



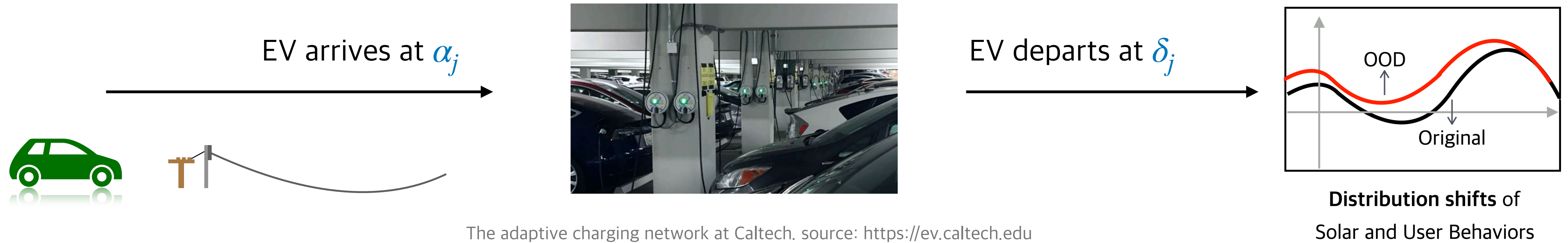
$$(e_{t+1} \| b_{t+1}) = s_{t+1} = g_{\mathcal{S}} \left[\underbrace{A_t s_t + B_t g_{\mathcal{A}}(a_t)}_{\text{Battery Dynamics}} + \underbrace{\ell'_t - \Delta h'_t}_{\text{(Uncertain) Behavior/Solar Perturbations}} \right], t \geq 0$$

Battery Dynamics

(Uncertain) Behavior/Solar Perturbations

$(\alpha_j, \delta_j, \kappa_j, i)$ is a charging session: At time α_j , EV j arrives at charger i , with an EV battery capacity κ_j , and departs at time δ_j

Combing MPC and RL Scheduling

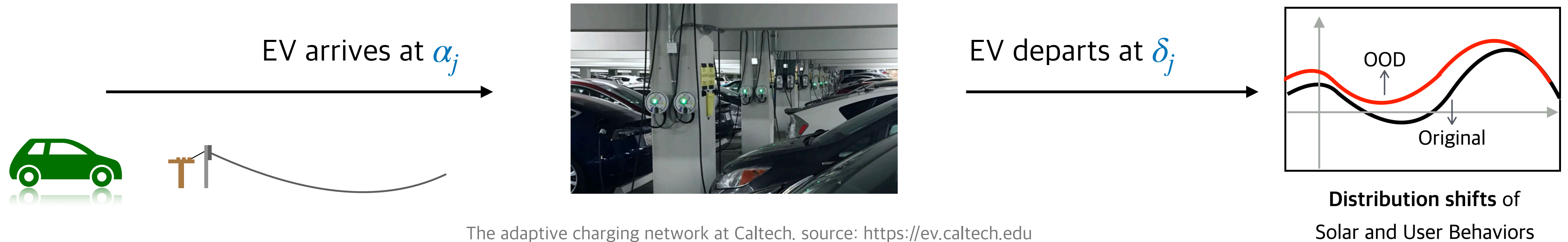


$$(e_{t+1} \| b_{t+1}) = s_{t+1} = g_{\mathcal{S}} \left[\underbrace{A_t s_t + B_t g_{\mathcal{A}}(a_t)}_{\text{Battery Dynamics}} + \underbrace{\ell'_t - \Delta h'_t}_{\text{(Uncertain) Behavior/Solar Perturbations}} \right], t \geq 0$$

Battery SoC Charging Rates Change of Charging Rates Human behaviors arrival/departure Solar generation

$(\alpha_j, \delta_j, \kappa_j, i)$ is a charging session: At time α_j , EV j arrives at charger i , with an EV battery capacity κ_j , and departs at time δ_j

Combing MPC and RL Scheduling



$$(e_{t+1} \| b_{t+1}) = s_{t+1} = g_{\mathcal{S}} \left[\underbrace{A_t s_t + B_t g_{\mathcal{A}}(a_t)}_{\text{Battery Dynamics}} + \underbrace{\ell'_t - \Delta h'_t}_{\text{(Uncertain) Behavior/Solar Perturbations}} \right], t \geq 0$$

Battery Dynamics

(Uncertain) Behavior/Solar Perturbations

Projections $g_{\mathcal{S}}$ and $g_{\mathcal{A}}$ capture network constraints, such as line constraints by the Kirchhoff's Current Law:

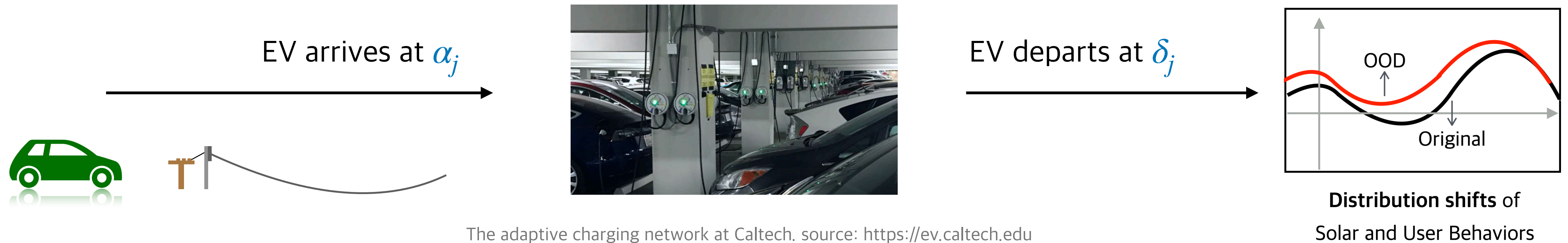
$$\left| \sum_{i=1}^N D_{ij} b_t(i) e^{j\phi_i} \right| \leq \gamma, \quad \forall t \in \mathcal{T}\$$$

Formed by circuit analysis

Phase angle of current phasor

Current magnitude limit

Combing MPC and RL Scheduling



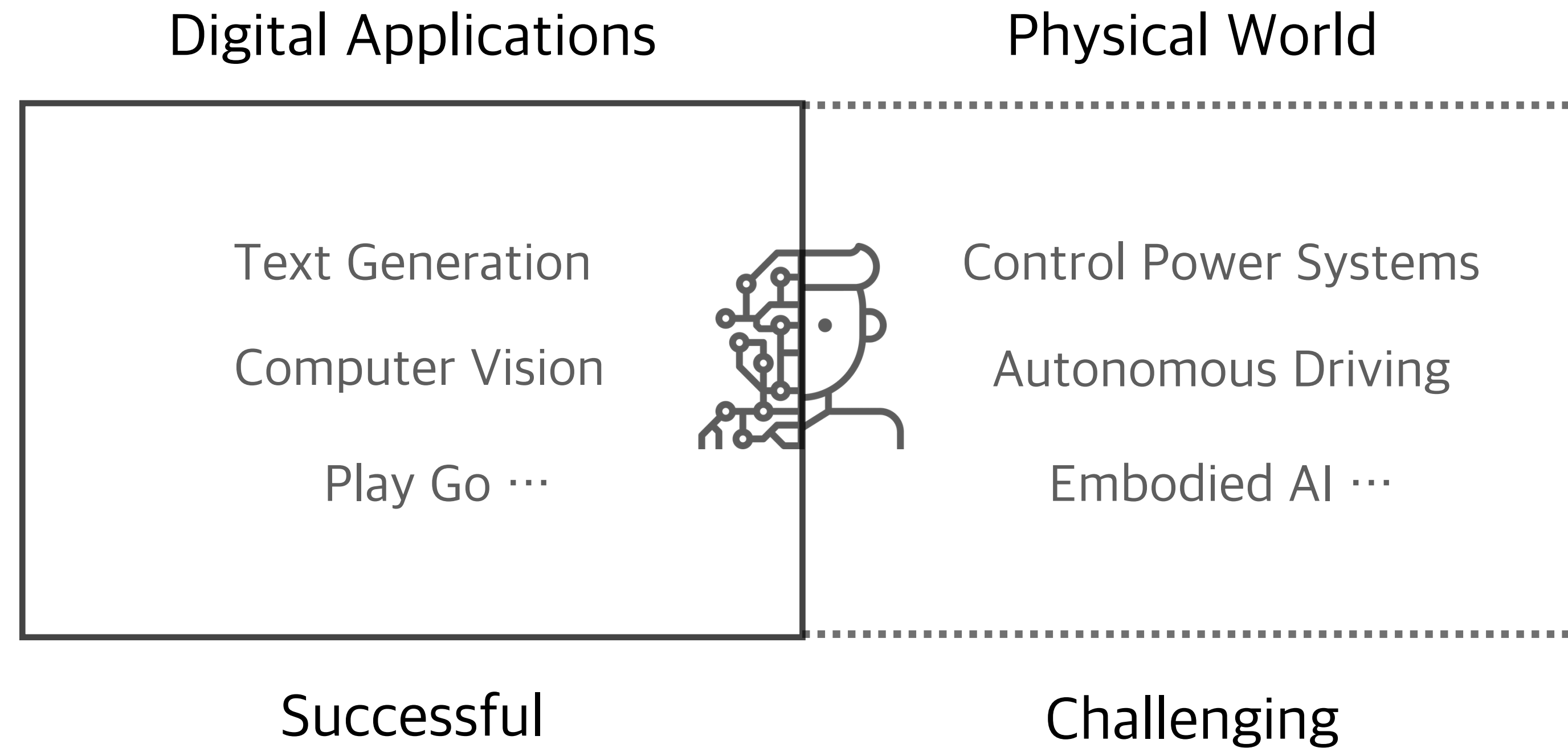
$$(e_{t+1} \| b_{t+1}) = s_{t+1} = g_{\mathcal{S}} \left[\underbrace{A_t s_t + B_t g_{\mathcal{A}}(a_t)}_{\text{Battery Dynamics}} + \underbrace{\ell'_t - \Delta h'_t}_{\text{(Uncertain) Behavior/Solar Perturbations}} \right], t \geq 0$$

Battery Dynamics

(Uncertain) Behavior/Solar Perturbations

- **Robustness** Classic algorithm (MPC) depends on battery dynamics and user inputs
- **Consistency** RL policy can better learn uncertain residuals when they are not out-of-distribution
- This is a general paradigm in many real-world applications

ML in Real-World Decision-Making ...



Machine-learned policies have the advantage of utilizing data

On average near-optimal performance

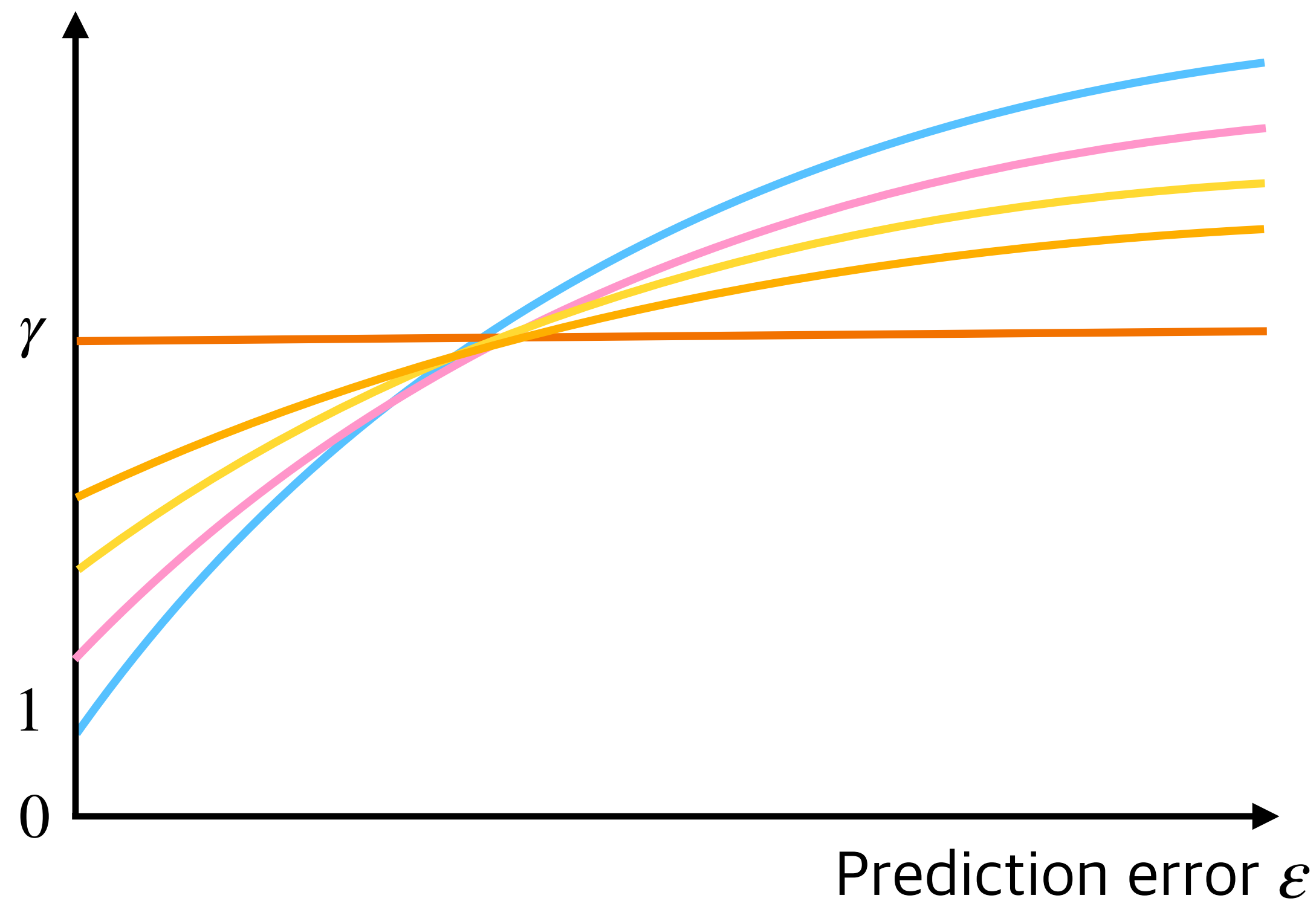
Existing well-established **classic methods** that are **hard** to be replaced entirely

Worst-case guarantee

Consistency vs Robustness Trade-off

On average good performance vs Worst-case guarantee

Performance $f(\varepsilon)$



Consistent ML Policy (Good when ε is small)

Intermediate Regimes

Robust Classic Policy (Good when ε is large)

$f(0)$ -consistent

$\sup_{\varepsilon \geq 0} f(\varepsilon)$ -robust

Today's Topic: Value-Based RL

Nonlinear Model is Harder if the ML Agent is a Black-Box

System Model	Classic Agent	ML Agent	Remarks	Tradeoffs
Linear Dynamics	LQR	MPC+Perturbation Predictions	Convex Combination	Consistency vs Robustness
NonLinear Dynamics	LQR	Black-Box RL	Switching	Consistency vs Stability
MDP	Robust Policy	?	?	?

Moving to general MDP ...

- Linear combination of two stabilizing controllers can be unstable
- Many learning-augmented online algorithms consider black-box predictions or advice
- What if we move beyond black-box advice?

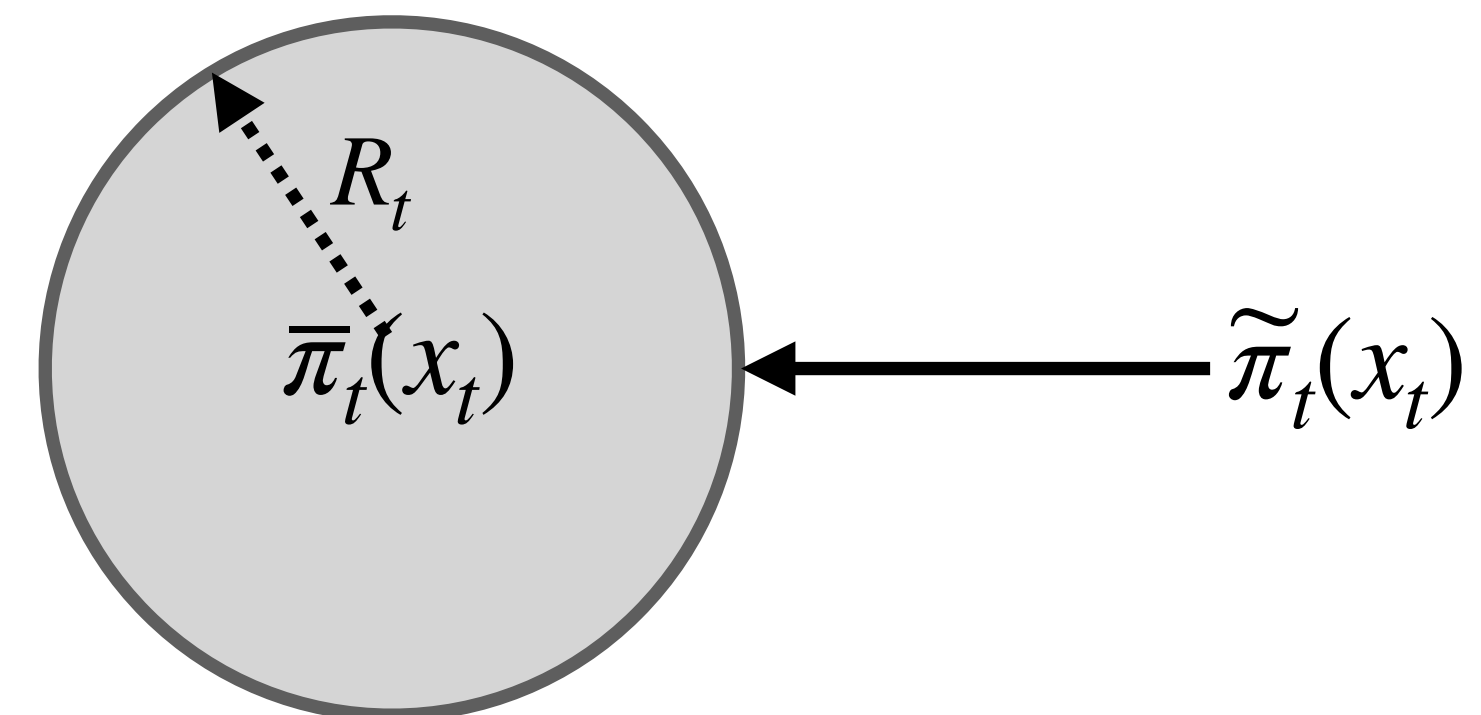
In General MDP ...

- Linear combination of two stabilizing controllers can be unstable
- Many learning-augmented online algorithms consider black-box predictions or advice
- What if we move beyond black-box advice?
- Need to consider more structural information, i.e., grey-box agents

Value-Based RL $\tilde{\pi} : X \rightarrow U$

$$\tilde{u}_t = \operatorname{arginf}_{u \in U} \tilde{Q}_t(u, x_t)$$

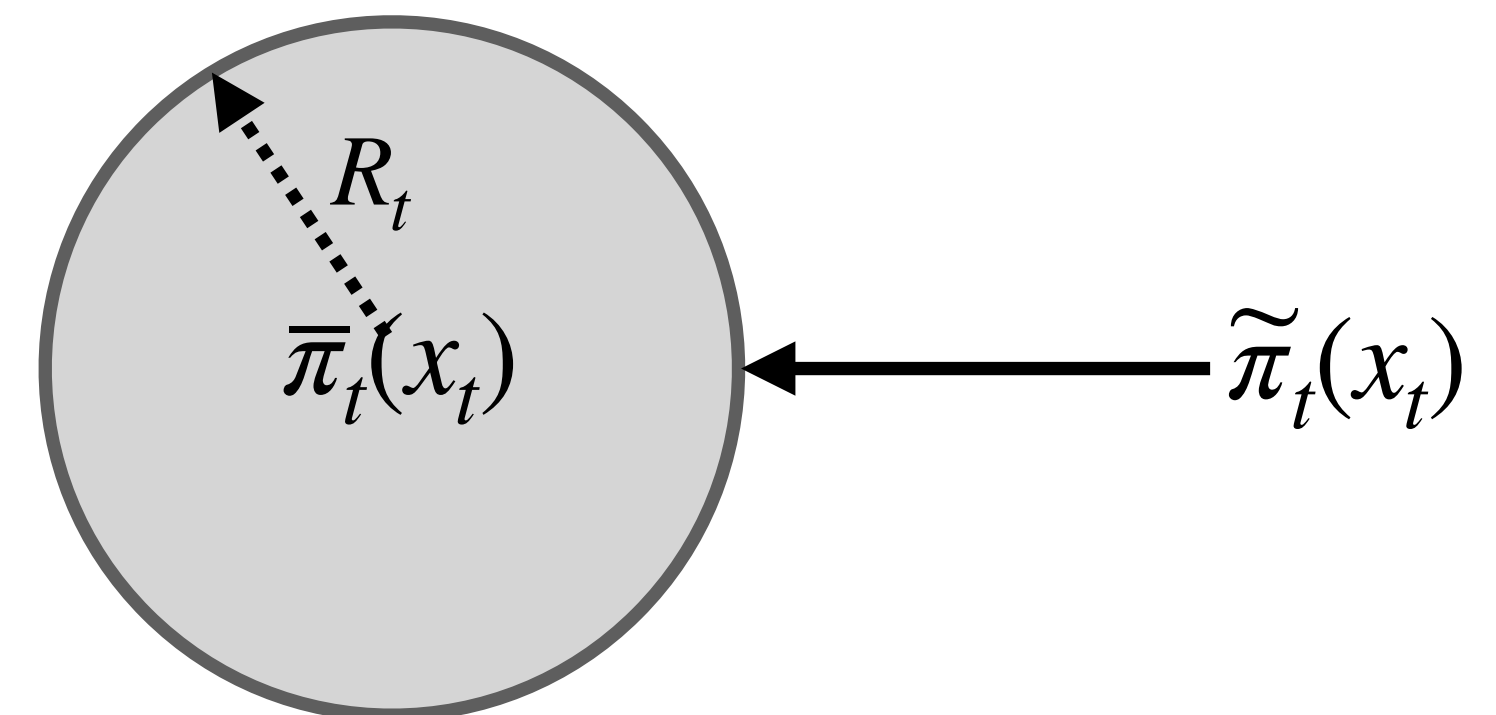
Q-value functions contain useful information



In General MDP ...

- Linear combination of two stabilizing controllers can be unstable
- Many learning-augmented online algorithms consider black-box predictions or advice
- What if we move beyond black-box advice?
- Need to consider more structural information, i.e., grey-box agents

How to select R_t ?



In General MDP ...

Classic Agent

$$\bar{\pi} : X \rightarrow U$$

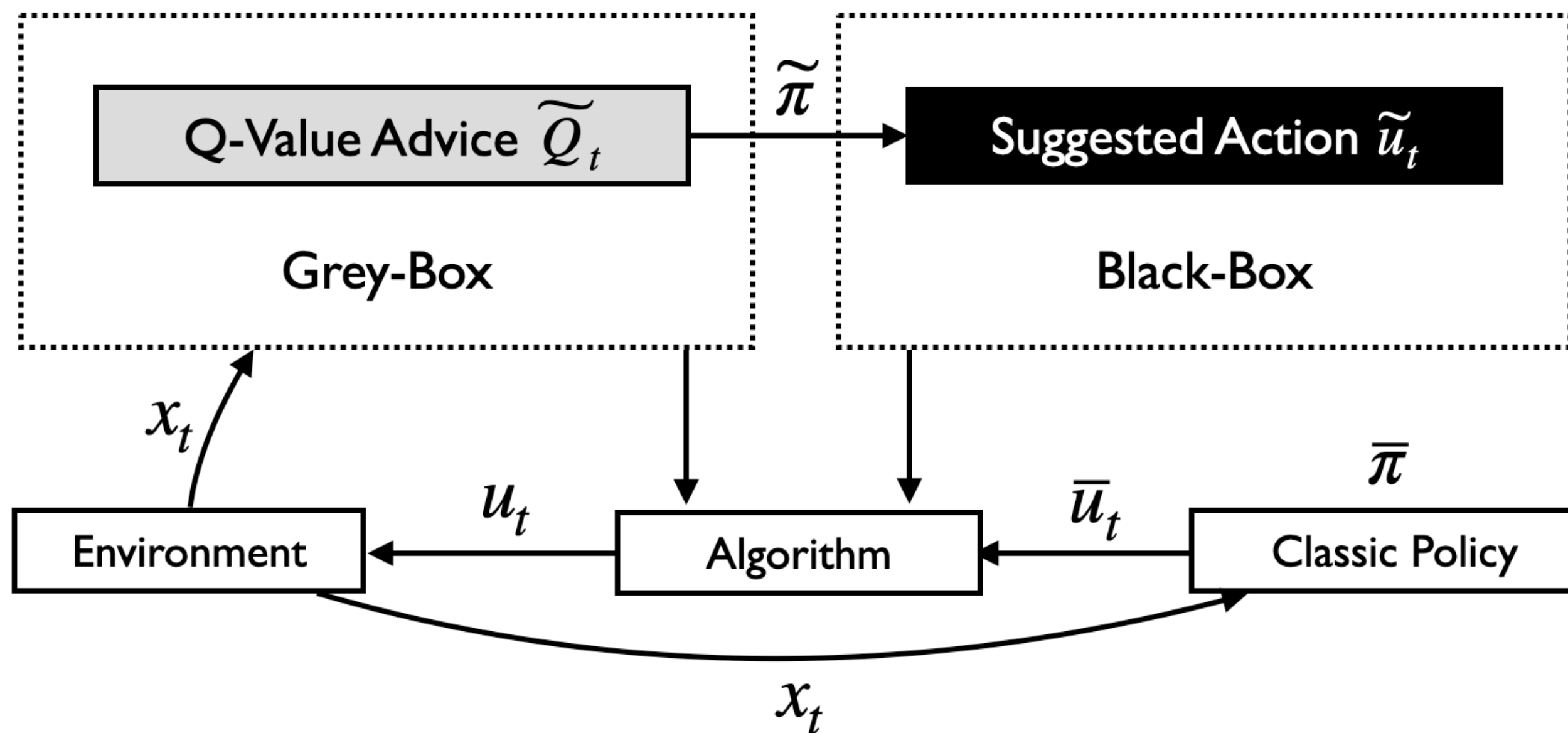
State Space: X

Action Space: U

ML Agent

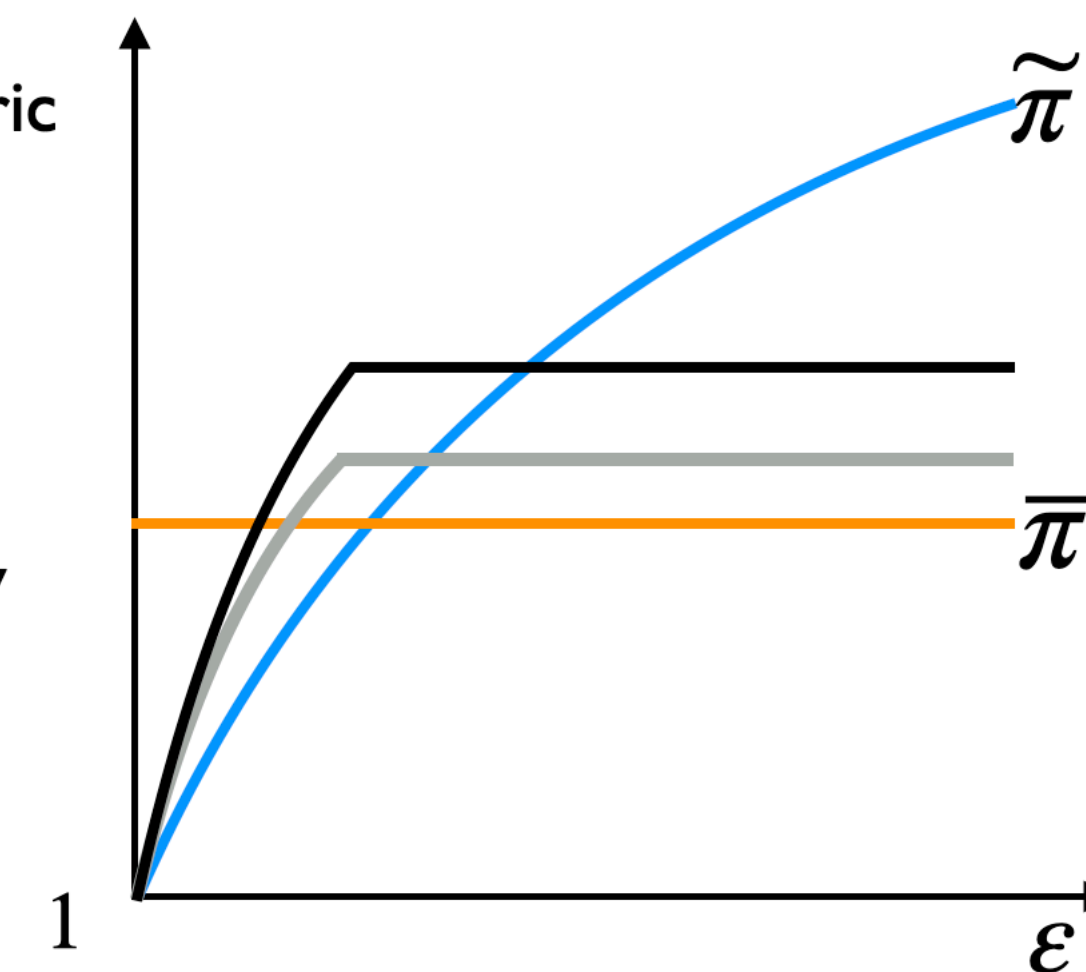
$$\tilde{\pi} : X \rightarrow U$$

Value-Based RL $\tilde{u}_t = \operatorname{arginf}_{u \in U} \tilde{Q}_t(u, x_t)$



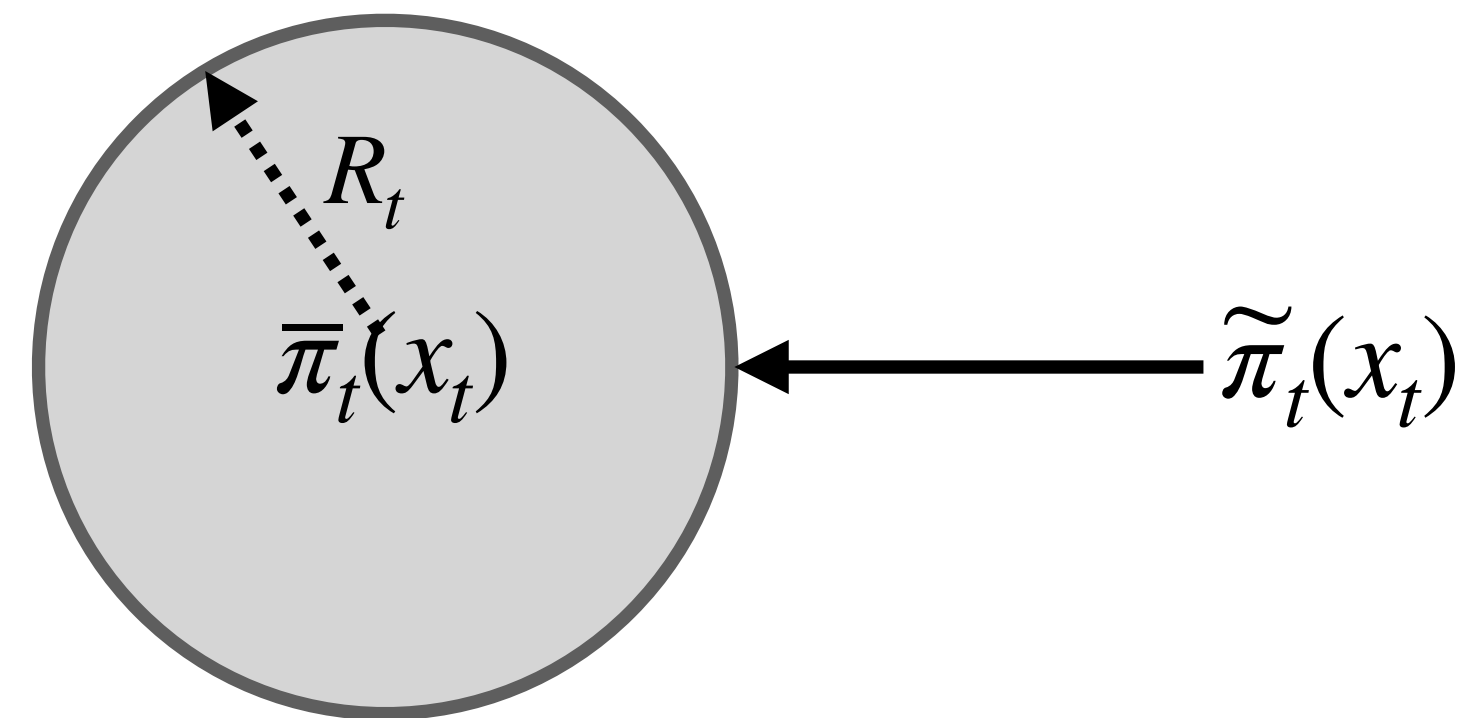
Performance Metric

- Grey-Box
- Black-Box
- Classic Policy
- ML Policy



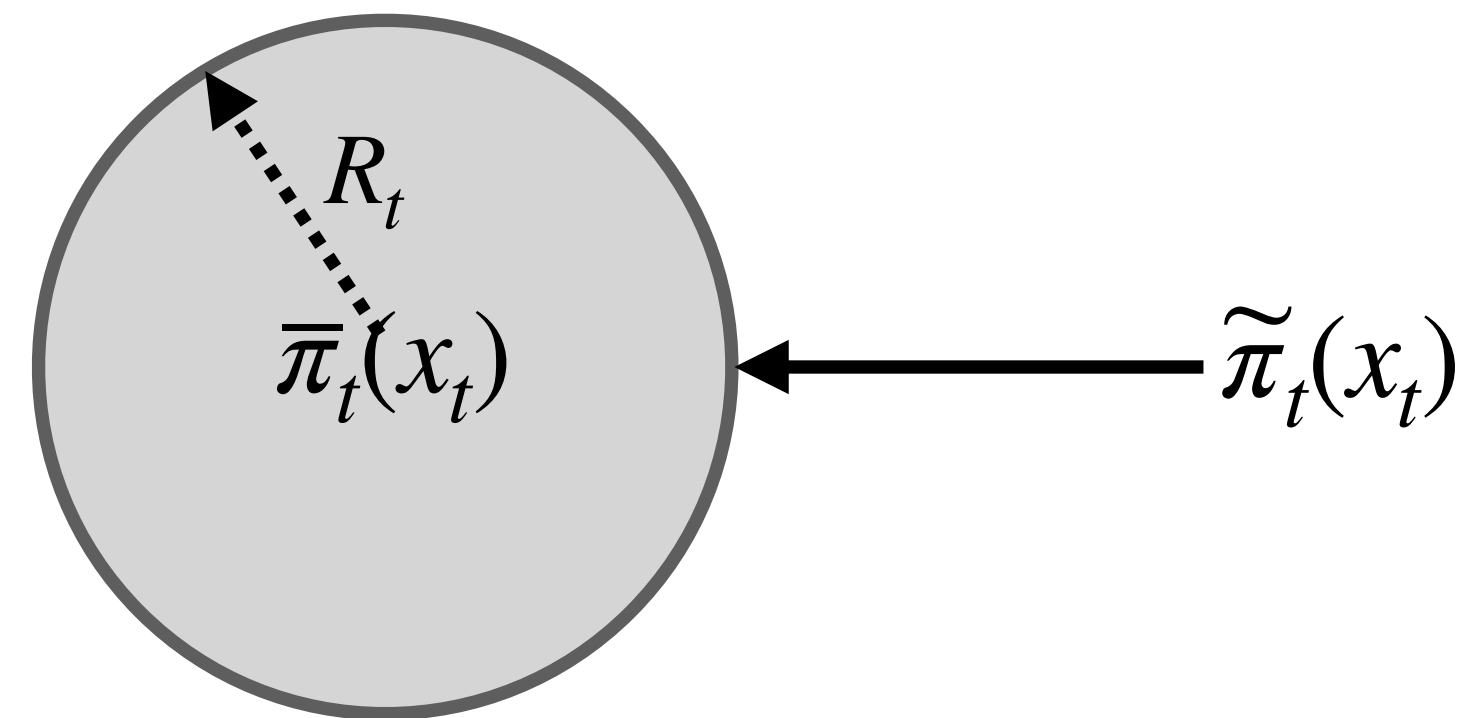
In the General MDP setting, can **Q-value advice** provide a better **consistency** vs **robustness** tradeoff?

Idea: Use Temporal Difference (TD) Error



Temporal Difference (TD)-Error: $TD_t = c_{t-1} + \mathbb{P}_{t-1} \tilde{V}_t - \tilde{Q}_{t-1}$

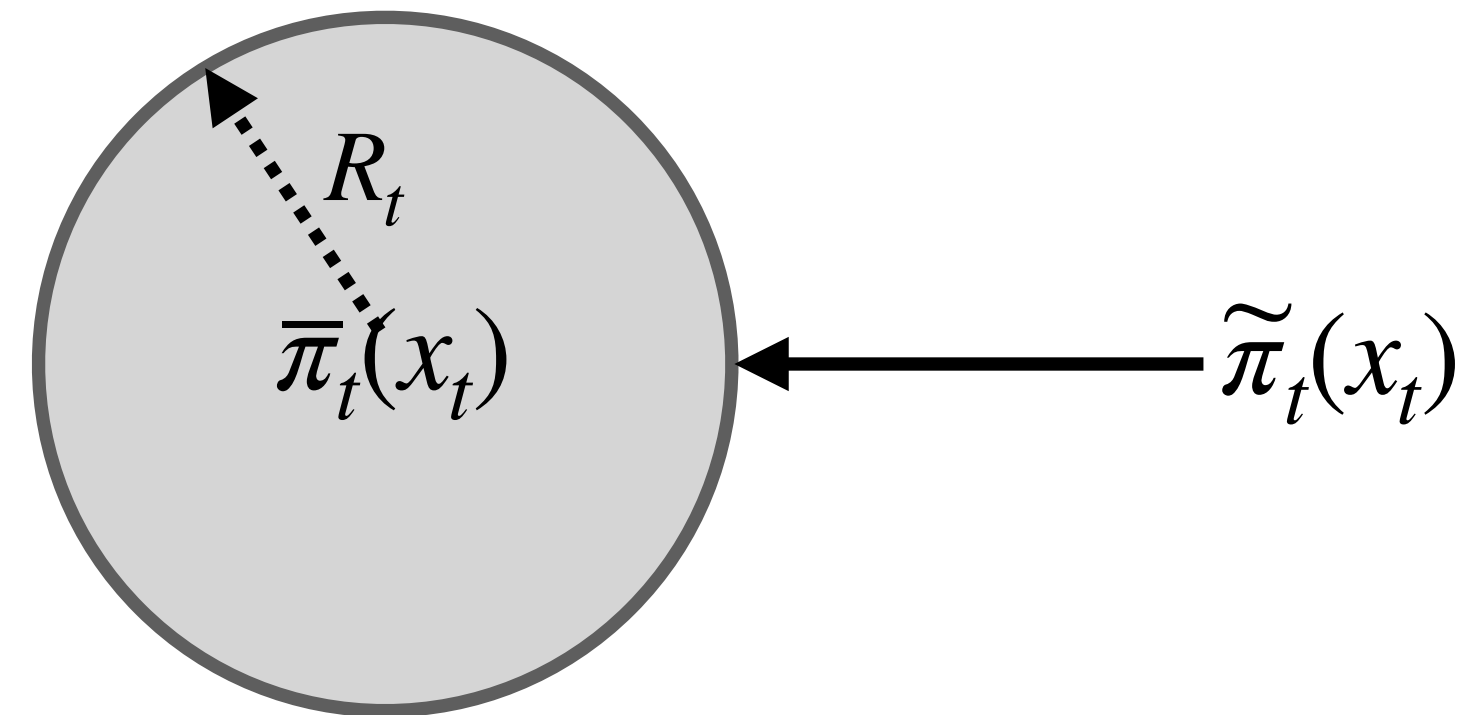
Idea: Use Temporal Difference (TD) Error



Temporal Difference (TD)-Error: $TD_t = c_{t-1} + \mathbb{P}_{t-1} \tilde{V}_t - \tilde{Q}_{t-1}$

Hard to compute since we don't know \mathbb{P}

Idea: Use Temporal Difference (TD) Error



Temporal Difference (TD)-Error: $\text{TD}_t = c_{t-1} + \mathbb{P}_{t-1} \tilde{V}_t - \tilde{Q}_{t-1}$

Approximate TD-Error: $\delta_t(x_t, x_{t-1}, u_{t-1}) := c_{t-1}(x_{t-1}, u_{t-1}) + \inf_{v \in \mathcal{U}} \tilde{Q}_t(x_t, v) - \tilde{Q}_{t-1}(x_{t-1}, u_{t-1})$

$$R_t := \left[\underbrace{\|\tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t)\|_{\mathcal{U}}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+ \quad \beta \text{ is a hyper-parameter}$$

Lipschitz constant of costs/rewards

Robust Baseline $\bar{\pi}$

Not all classic policies can be used ...

We need to regulate the behaviors of the **classic policies** so they become baselines (to guarantee worst-case performance)

Robust Baseline $\bar{\pi}$

Definition (*r*-locally *p*-Wasserstein robustness)

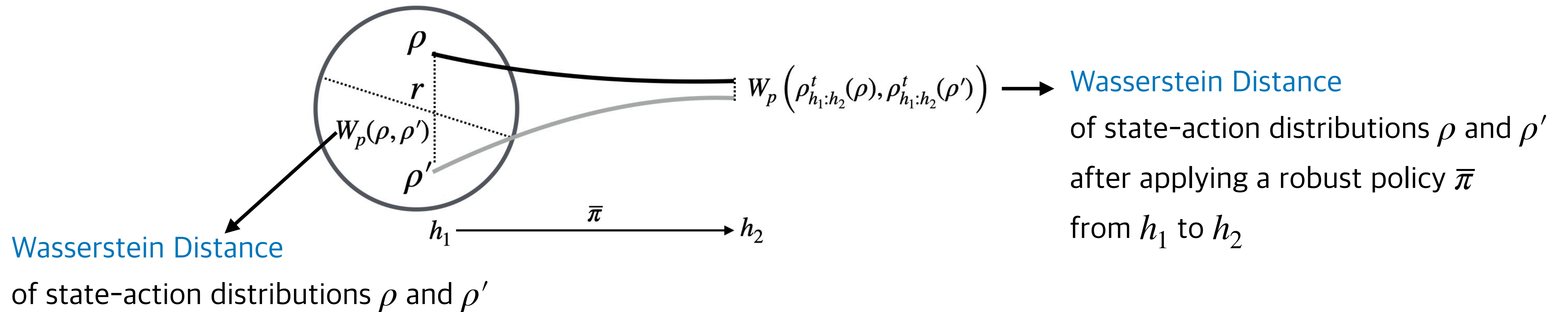
A policy $\bar{\pi} = (\pi_t : t \in [T])$ is *r*-locally *p*-Wasserstein-robust if for any $0 \leq t_1 \leq t_2 < T$ and state-action distributions ρ, ρ' such that $W_p(\rho, \rho') \leq r$, for some radius $r > 0$,

$$W_p \left(\rho_{t_1:t_2}(\rho), \rho_{t_1:t_2}(\rho') \right) \leq s(t_2 - t_1) W_p(\rho, \rho')$$

for some function $s : [T] \rightarrow \mathbb{R}_+$ such that $\sum_{t \in [T]} s(t) \leq C_s$ for some constants $C_s > 0$.

Robust Policy

A general class of **robust** classic policies



Many practical instances:

- Discrete MDP: Any Policy that Induced a Regular Markov Chain
- Time-varying LQR: MPC with Robust Predictions
- Extends a contraction property in [Lin 2022]

PROjection Pursuit Policy (PROP)

Algorithm PROjection Pursuit (PROP)

Initialize: $\tilde{\pi} = (\tilde{\pi}_t : t \in [T])$ and $\bar{\pi} = (\bar{\pi}_t : t \in [T])$

for $t = 0, \dots, T - 1$

Get R_t using approximate TD-error

Take $u_t = \text{Proj}_{\bar{U}_t}(\tilde{u}_t)$ where $\bar{U}_t := \left\{ u \in \mathcal{U} : \|u - \bar{\pi}_t(x_t)\|_{\mathcal{U}} \leq R_t \right\}$

Sample next state $x_{t+1} \sim \mathbb{P}_t(x_t, u_t)$

OOD-Aware EV Charging

Algorithm OOD-Aware EV Charging (OOD-Charging)

Initialize: $\tilde{\pi} = (\tilde{\pi}_t : t \in [T])$ and $\bar{\pi} = (\bar{\pi}_t : t \in [T])$

for $t = 0, \dots, T - 1$

↳ NN that updates every t

↳ MPC Procedure with user inputs and estimated state

Approximately Wasserstein robust

Receive user inputs

Get R_t using approximate TD-error

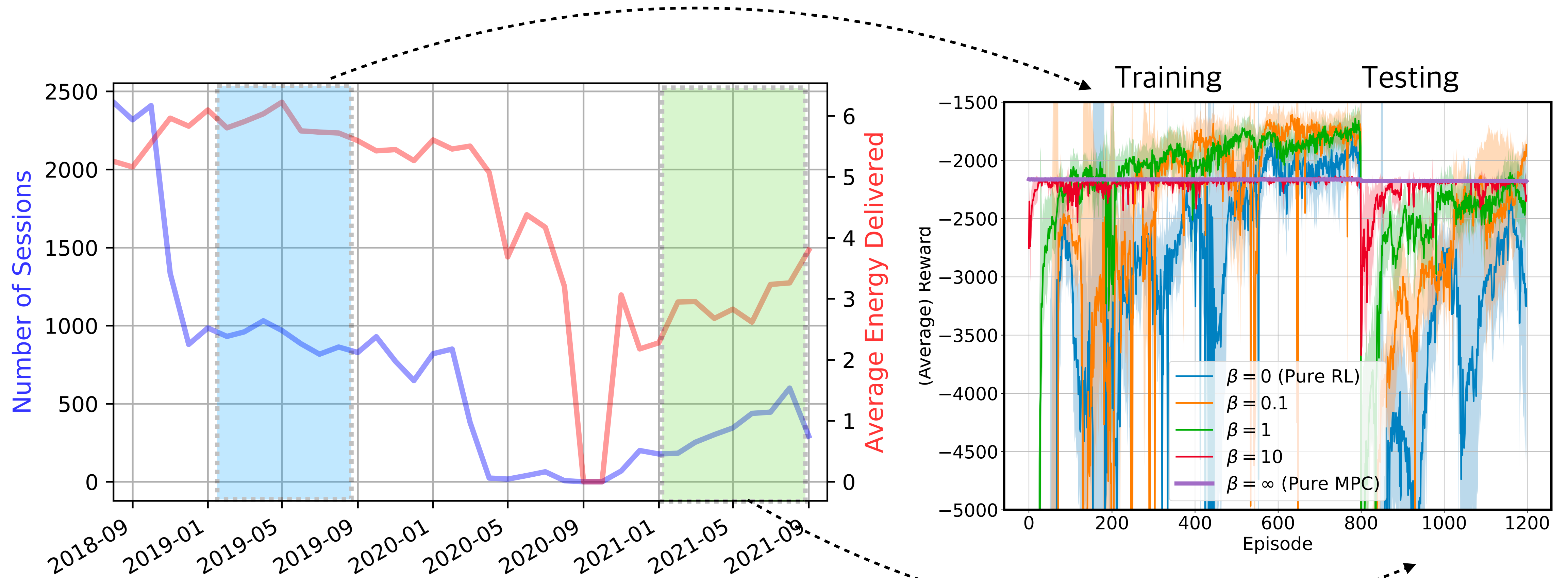
Take $u_t = \text{Proj}_{\bar{U}_t}(\tilde{u}_t)$ where $\bar{U}_t := \left\{ u \in \mathcal{U} : \|u - \bar{\pi}_t(x_t)\|_{\mathcal{U}} \leq R_t \right\}$

Sample next state $x_{t+1} \sim \mathbb{P}_t(x_t, u_t)$

Estimate previous state \tilde{x}_t

Update replay buffer and retrain $\tilde{\pi}$

Out-of-Distribution EV Charging

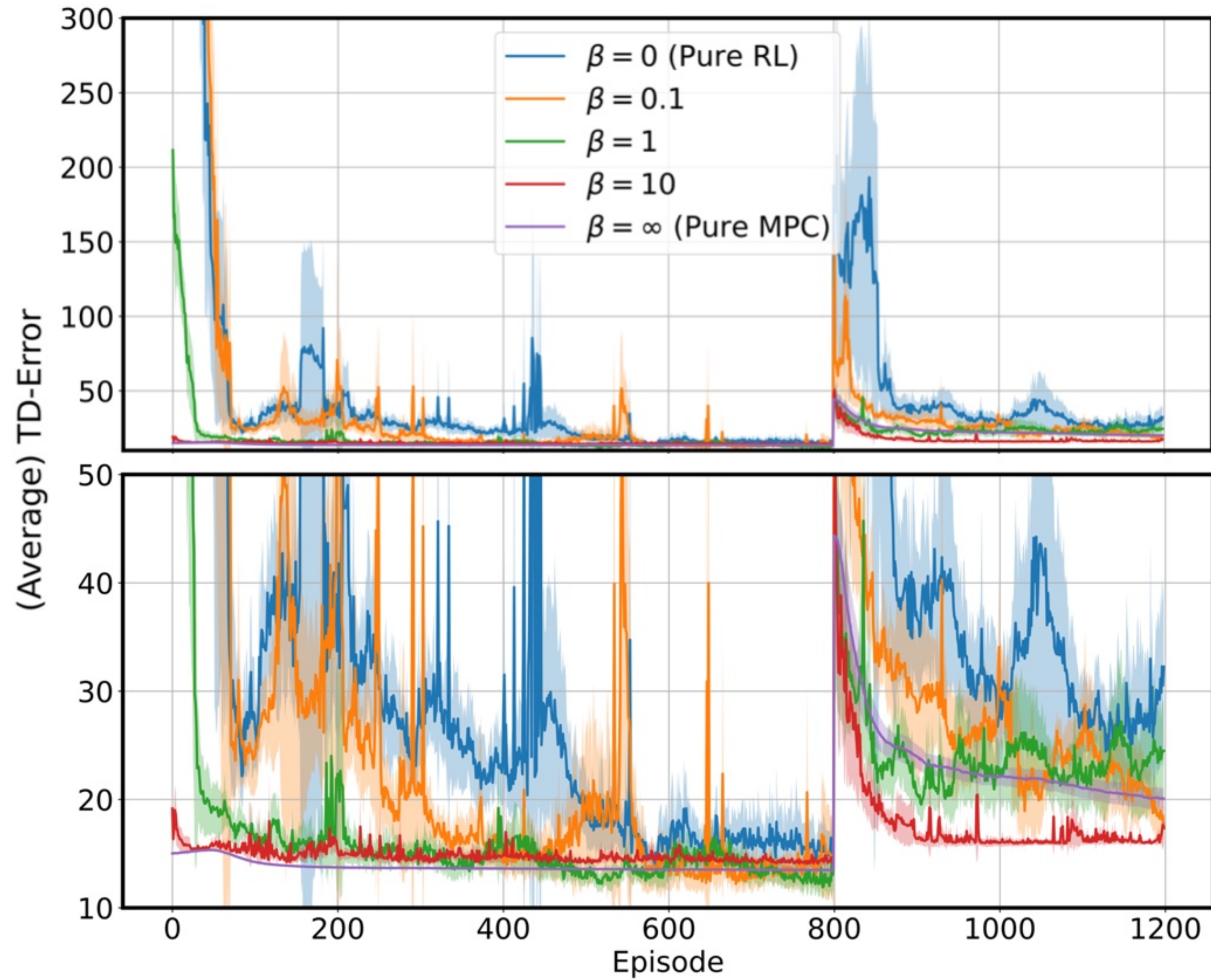


Pre-COVID19

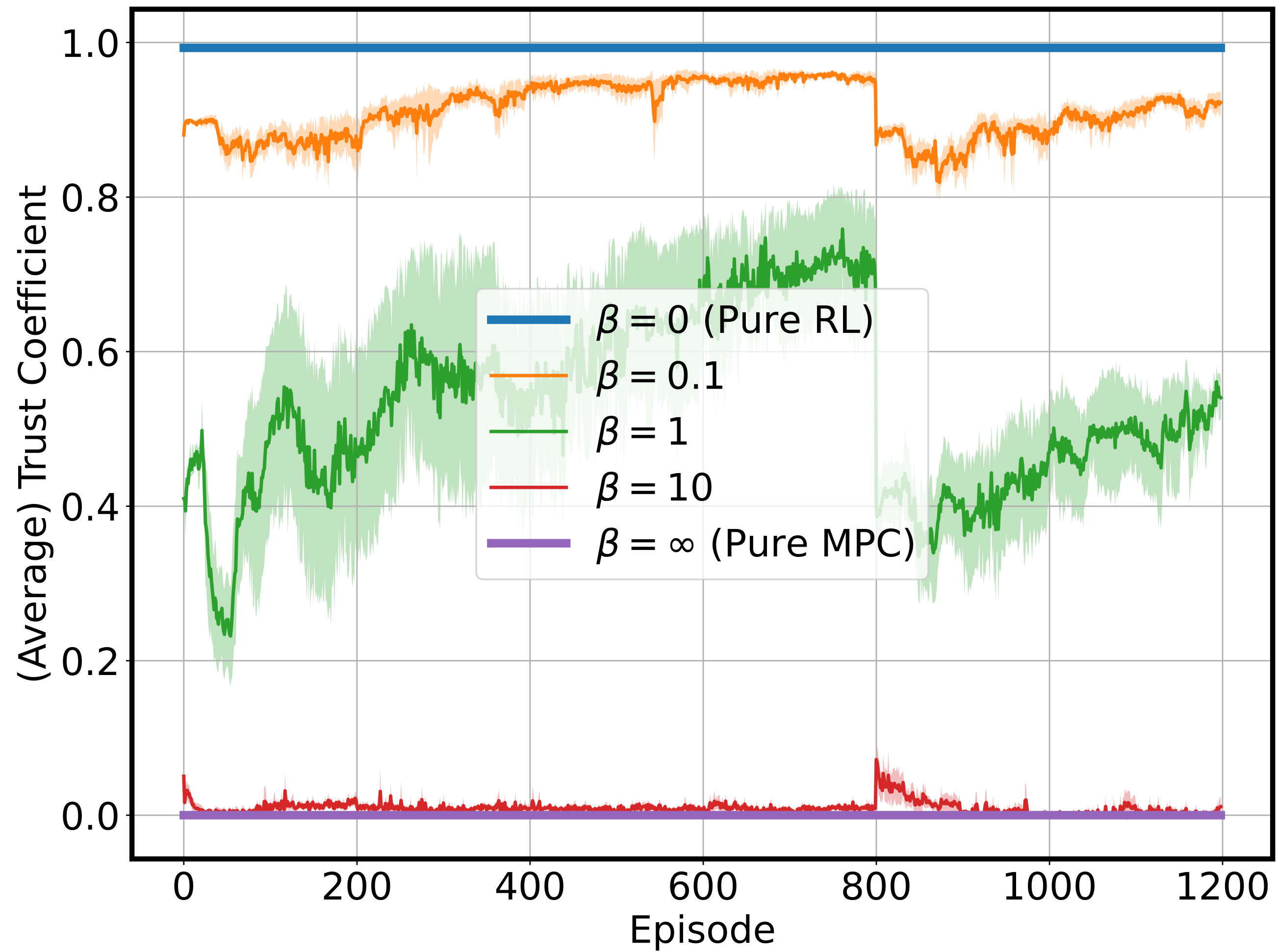
ACN-Data

Post-COVID19

Out-of-Distribution EV Charging



Out-of-Distribution EV Charging



Trust Coefficient

$$\lambda(R_t) = \min \{ 1, R_t / \|\tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t)\|_2 \}$$

Theoretical Guarantees

Consistency and Robustness

k -Consistency: Ratio of Expectations (RoE) satisfies $\text{RoE}(\varepsilon) \leq k$ for $\varepsilon = 0$

l -Robustness: Ratio of Expectations (RoE) satisfies $\text{RoE}(\varepsilon) \leq k$ for any ε

$$\text{where } \varepsilon := \sum_{t \in [T]} \left(\|\widetilde{Q}_t - Q_t^*\|_\infty + \left\| \inf_{v \in \mathcal{U}} \widetilde{Q}_t - \inf_{v \in \mathcal{U}} Q_t^* \right\|_\infty \right)$$

(can be generalize to:)

$$\varepsilon(p, \rho) := \sum_{t \in [T]} \left(\|\widetilde{Q}_t - Q_t^*\|_{p, \rho_t} + \left\| \inf_{v \in \mathcal{U}} \widetilde{Q}_t - \inf_{v \in \mathcal{U}} Q_t^* \right\|_{p, \phi_t} \right)$$

Black-Box Impossibility

Theorem (Informal)

There exists an algorithm with a black-box agent that is $(1 + \mathcal{O}((1 - \lambda)\gamma))$ -consistent and $(\text{ROB} + \mathcal{O}(\lambda\gamma))$ -robust where $0 \leq \lambda \leq 1$ is a hyper-parameter.

(ROB is a ratio of expectation upper bound for the robust baseline)

Theorem (Informal) **Impossibility**

Any algorithm with a black-box agent cannot be both $(1 + o((1 - \lambda)\gamma))$ -consistent and $(\text{ROB} + o(\lambda\gamma))$ -robust for any $0 \leq \lambda \leq 1$.

Proof Highlights

Theorem (Informal) **Impossibility**

Any algorithm with a black-box agent cannot be both $(1 + o((1 - \lambda)\gamma))$ -consistent and $(\text{ROB} + o(\lambda\gamma))$ -robust for any $0 \leq \lambda \leq 1$.

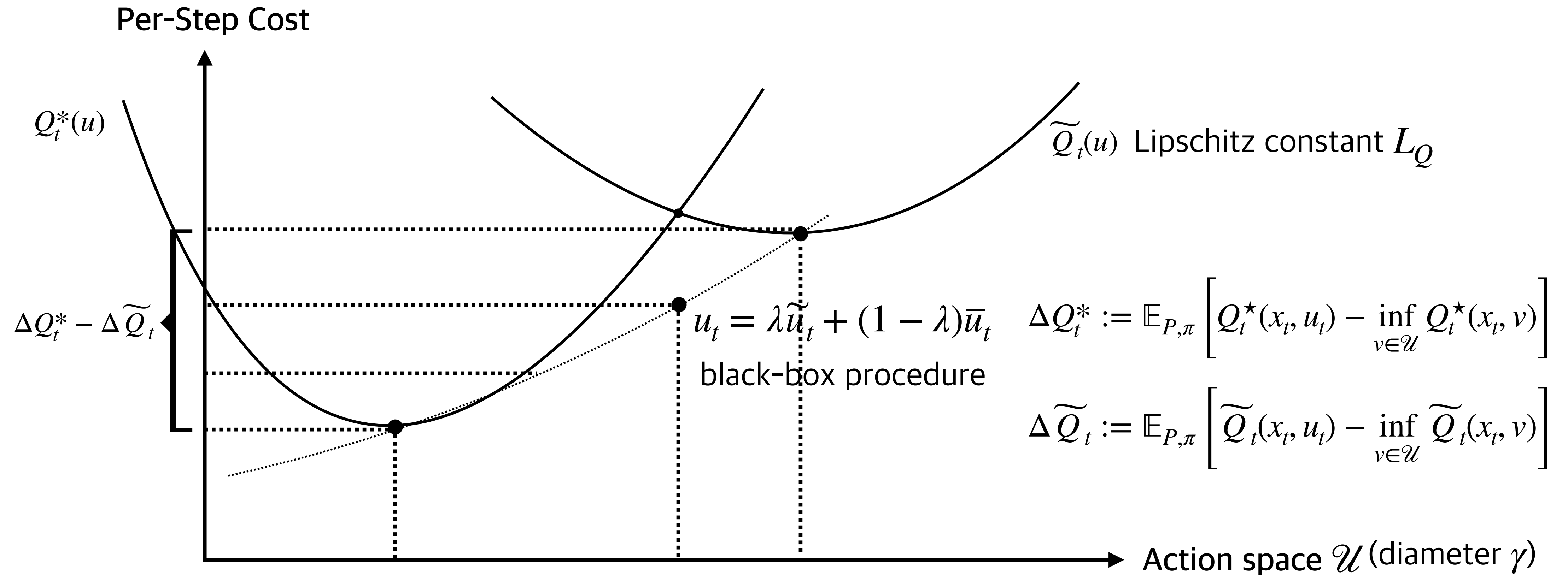
Proof Idea:

Construct a special case (satisfying all model assumptions) with decoupled and identical cost at each t



Then argue with fixed λ , can separate Q^* and \tilde{Q} so that a lower bound can be derived

Proof Highlights

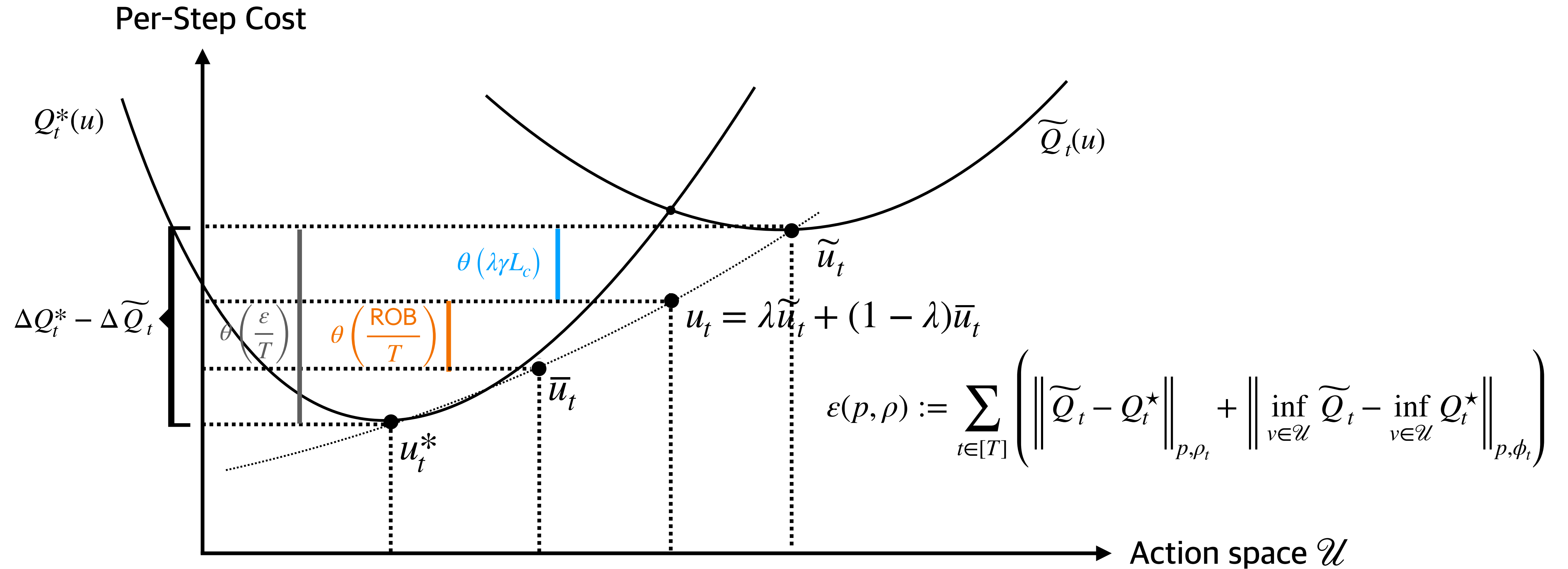


Dynamic regret:

$$\text{DR (PROP(Black-Box))} = \sum_{t \in [T]} \mathbb{E}_{P, \pi} \left[Q_t^*(x_t, u_t) - \inf_{v \in \mathcal{U}} Q_t^*(x_t, v) \right]$$

$$\text{DR (PROP(Black-Box))} \geq \sum_{t \in [T]} \left(\Delta Q_t^*(P, \pi) - \Delta \tilde{Q}_t(P, \pi) + (1 - \lambda)L_Q\gamma \right)$$

Proof Highlights



$$\frac{\sum_{t \in [T]} \left(\Delta Q_t^*(P, \pi) - \Delta \tilde{Q}_t(P, \pi) \right)}{\text{OPT}} \geq \Omega \left(\text{ROB} + \frac{\lambda\gamma L_c}{\text{OPT}} T \right) \implies \text{RoE}(\text{PROP}(\text{Black-Box})) = 1 + \Omega \left((1 - \lambda)L_Q\gamma + \min\{\varepsilon, \lambda\gamma L_c + \text{ROB}\} \right)$$

(cannot be both $(1 + o((1 - \lambda)\gamma))$ -consistent and $(\text{ROB} + o(\lambda\gamma))$ -robust for any $0 \leq \lambda \leq 1$)

Grey-Box Setting

Theorem (Informal)

PROP with a grey-box agent that is 1-consistent and $(\text{ROB} + o(1))$ -robust for some $\beta > 0$.

Take-Aways: Grey-box information can grant nontrivial improvements on the consistency and robustness tradeoff

Out-of-Distribution EV Charging

Constant

Theorem (Informal)

OOD-Charging with a grey-box agent that is $1 + \mathcal{O}(\bar{W})$ -consistent and $(\text{MPC} + \mathcal{O}(\bar{W}) + o(1))$ -robust for some $\beta > 0$.

$$\text{MPC} \leq \frac{2\xi C^2(1 + C^2)(1 + \bar{A}^2 + \bar{B}^2)}{\mu(1 - \bar{\lambda})^2} \quad \bar{\lambda} := \left(\frac{\bar{\sigma} - \underline{\sigma}}{\bar{\sigma} + \underline{\sigma}} \right)^{\frac{1}{2}} \quad C := \frac{4(\xi + 1 + \bar{A} + \bar{B})}{\underline{\sigma}^2 \cdot \lambda}$$

$$\underline{\sigma} := \min\{\mu, 1\}(\bar{A} + \bar{B} + 1) \left(\xi / (2\mu\xi + \mu\sigma^2) \right)^{\frac{1}{2}}$$

$$\bar{\sigma} := \sqrt{2}(\xi + \bar{A} + \bar{B} + 1)$$

$$\text{State estimation error: } \|\ell'_t - \Delta h'_t\| \leq \bar{W}$$

$$\text{Standard assumptions: } \|A_t\| \leq \bar{A} \quad \|B_t\| \leq \bar{B} \quad \mu I_n \leq Q_t \leq \xi I_n, \quad \mu I_m \leq R_t \leq \xi I_m, \quad \mu I_n \leq P \leq \xi I_n$$

Proof Highlights

Theorem (Informal)

PROP with a grey-box agent that is 1-consistent and $(\text{ROB} + o(1))$ -robust for some $\beta > 0$.

Proof Idea: A general bound on DR, therefore RoE:

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

p-Wasserstein robustness

$\|\tilde{u}_t - \bar{u}_t\|$

Assume $c_t > 0$ for all t , we can bound RoE

Here, $\mu_t := \zeta_t^V - \zeta_t^Q$ $\zeta_t^Q(x_t, u_t) := \widetilde{Q}_t(x_t, u_t) - Q_t^*(x_t, u_t)$ Q-value error

$\zeta_t^V(x_t) := \inf_{v \in \mathcal{U}} \widetilde{Q}_t(x_t, v) - \inf_{v \in \mathcal{U}} Q_t^*(x_t, v)$ V-value error

Proof Highlights

Fix a choice of projection radii $(R_t : t \in [T])$.

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} \left[(R_t)^p \right]^{1/p}}_{\text{Robustness Bound}} \right\}$$

↑
Applying a projection lemma

↑

Applying the Kantorovich-Rubinstein duality theorem

Applying the Wasserstein robustness definition

$$\begin{aligned} J(\pi) - J(\bar{\pi}) &= \sum_{t \in [T]} \mathbb{E}_{(x, u) \sim \rho_t} [c_t(x, u)] - \mathbb{E}_{(x, x) \sim \bar{\rho}_t} [c_t(x, u)] \\ &\leq L_C \sum_{t \in [T]} \sum_{\tau=0}^{t-1} s(\tau) \mathbb{E}_{P, \pi} \left[(R_{t-\tau})^p \right]^{1/p} \\ &\leq L_C C_s \sum_{t \in [T]} \mathbb{E}_{P, \pi} \left[(R_t)^p \right]^{1/p} \end{aligned}$$

Proof Highlights

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P,\pi} [\mu_t] + L_Q \mathbb{E}_{P,\pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P,\pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Consistency: Let $\varepsilon = 0$, $\mathbb{E}_{P,\pi} [\mu_t] = 0$ and the consistency bound becomes

$$\mathbb{E}_{P,\pi} [\eta_t - R_t] \leq \mathbb{E}_{P,\pi} \left[\frac{\beta}{L_Q} \sum_{s=1}^t \delta_s \right] = 0 \quad \text{Applying the radius update rule:}$$

$$R_t := \left[\underbrace{\left\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \right\|_{\mathcal{U}}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

Proof Highlights

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Robustness:

$$R_t := \left[\underbrace{\left\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \right\|_{\mathcal{U}}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

TD-Error:

$$\text{TD}_t = c_{t-1} + \mathbb{P}_{t-1} \tilde{V}_t - \tilde{Q}_{t-1}$$

Approximate TD-Error:

$$\delta_t(x_t, x_{t-1}, u_{t-1}) := c_{t-1}(x_{t-1}, u_{t-1}) + \inf_{v \in \mathcal{U}} \tilde{Q}_t(x_t, v) - \tilde{Q}_{t-1}(x_{t-1}, u_{t-1})$$

Key observation:

$$\mu_t - \delta_t = \zeta_{t-1}^Q - \zeta_t^Q \quad (\zeta_{-1}^Q = 0)$$

$$\implies \sum_{s=0}^t (\mu_s - \delta_s) = \sum_{s=0}^t (\zeta_{s-1}^Q - \zeta_s^Q) = \zeta_t^Q$$

Proof Highlights

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Robustness:

$$R_t := \left[\underbrace{\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \|_{\mathcal{U}}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

$\exists \Delta = o(T)$ such that $|\zeta_t^Q| \leq \Delta$ for all $t \in [T]$ (model assumption)

Consider two cases:

Case I $\sum_{t \in [T]} \mu_t \leq \Delta$ Automatically obtain $\text{RoE}(\varepsilon) \leq \text{ROB} + o(1)$ by the consistency bound

Case II $\sum_{t \in [T]} \mu_t > \Delta$ (Cont.)

Proof Highlights

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Robustness:

$$R_t := \left[\underbrace{\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \|_{\cup}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

Case II $\sum_{t \in [T]} \mu_t > \Delta$

$$\implies \sum_{t \in [T]} \delta_t > 0 \quad (\text{Applying } \sum_{s=0}^t (\mu_s - \delta_s) = \sum_{s=0}^t (\zeta_{s-1}^Q - \zeta_s^Q) = \zeta_t^Q)$$

\implies There exists $\beta > 0$ such that $R_t = 0$

(The action space is compact, discrepancy η_t is bounded)

Proof Highlights

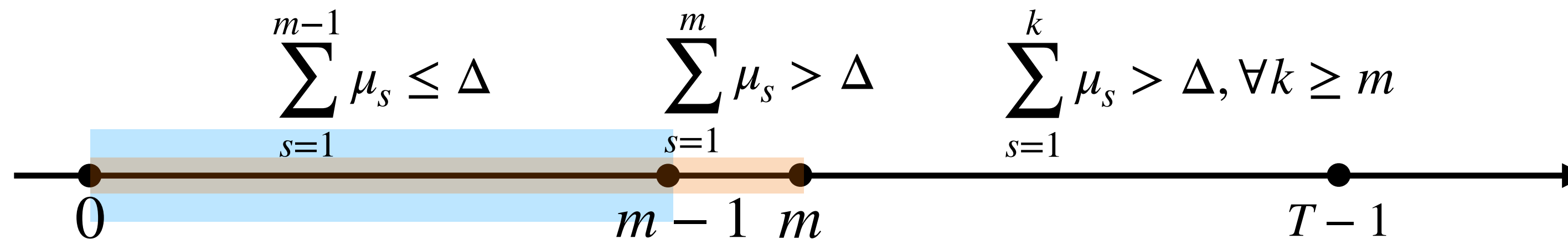
$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Robustness:

$$R_t := \left[\underbrace{\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \|_{\cup}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

Case II $\sum_{t \in [T]} \mu_t > \Delta$

Let m be the latest time index such that $\sum_{s=1}^m \mu_s > \Delta$ and $\sum_{s=1}^{m-1} \mu_s \leq \Delta$



Proof Highlights

$$\text{DR}(\text{PROP}) \leq \sum_{t \in [T]} \min \left\{ \underbrace{\mathbb{E}_{P, \pi} [\mu_t] + L_Q \mathbb{E}_{P, \pi} (\eta_t(x_t) - R_t)}_{\text{Consistency Bound}}, \underbrace{\varphi_t + L_C C_s \mathbb{E}_{P, \pi} [(R_t)^p]^{1/p}}_{\text{Robustness Bound}} \right\}$$

Robustness:

$$R_t := \left[\underbrace{\| \tilde{\pi}_t(x_t) - \bar{\pi}_t(x_t) \|_{\cup}}_{\text{Decision Discrepancy } \eta_t} - \frac{\beta}{L_Q} \sum_{s=1}^t \underbrace{\delta_s(x_s, x_{s-1}, u_{s-1})}_{\text{Approximate TD-Error}} \right]^+$$

Case II $\sum_{t \in [T]} \mu_t > \Delta$ (Cont.)

Per-step robustness bound

$$\implies \text{DR}(\text{PROP}(\text{Grey-Box})) \leq \sum_{s=1}^{m-1} \mu_s + \sum_{s=m}^{T-1} \varphi_s + O(\eta_m) \leq \sum_{s=1}^{m-1} \mu_s + \sum_{s=m}^{T-1} \varphi_s + O(\gamma)$$

(Applying the Wasserstein robustness definition)

$$\implies \text{RoE}(\varepsilon) \leq \text{ROB} + o(1)$$

Comparison

	Noah & Moitra	This Work
MDP Model	Finite Action/State Spaces Episodic setting	Finite or Continuous Single-trajectory setting
Assumption on \tilde{Q}	Approximate distillation for stronger results	Lipschitz continuous $\tilde{Q} - Q^* = o(T)$
Robust Policies	N/A	Wasserstein Robustness
Main Results	Regret bound	Consistency-robustness Tradeoff

Summary: Learning-Augmented Decision-Making

Systems	Untrusted AI		Performance Guarantees	
Nonlinear Dynamics	Black-box Policy	[Li et. al. OJCSYS 2023]	Stability + CR	Scenario II : Battery Scheduling
Anytime-Constrained MDP	Black-box Policy	[Yang et. al. NeurIPS 2023]	Regret	
MDP	Grey-Box Policy	[Li et. al. NeurIPS 2023]	Ratio of Expectations	
Two-Controller System	Black-box Advice	[Li et. al. e-Energy 2020] [Li et. al. TSG 2021] [Li et. al. SIGMETRICS 2021]	Feasibility + CR	Scenario III : Remand Response
Linear Quadratic System	Perturbation Predictions	[Li et. al. SIGMETRICS 2022]	CR	Scenario I : Linear Control
Linear Quadratic System	Disentangled Perturbations		Ongoing	
General Sum Games	Black-box Bayesian probability		Ongoing	

References

- (1) **Tongxin Li**, Ruixiao Yang, Guannan Qu, Guanya Shi, Chenkai Yu, Adam Wierman, and Steven Low.
“Robustness and Consistency in Linear Quadratic Control with Untrusted Predictions”
SIGMETRICS Performance Evaluation Review 50.1 (2022): 107-108.
- (2) Jianyi Yang, Pengfei Li, **Tongxin Li**, Adam Wierman, and Shaolei Ren.
“Anytime-competitive reinforcement learning with policy prior”
Advances in Neural Information Processing Systems 36 (2024).
- (3) Yeja Liu, Jianyi Yang, Pengfei Li, **Tongxin Li**, and Shaolei Ren.
“Building Socially-Equitable Public Models”
Advances in Neural Information Processing Systems 36 (2024).
- (4) **Tongxin Li**, Yiheng Lin, Shaolei Ren, Adam Wierman.
“Beyond Black-Box Advice: Learning-Augmented Algorithms for MDP with Value-Based Policies”
Advances in Neural Information Processing Systems 36 (2024).
- (5) **Tongxin Li**, Chenxi Sun.
“Out-of-Distribution-Aware Electric Vehicle Charging”
IEEE Transactions on Transportation Electrification, 2024
- (6) **Tongxin Li**, Hao Liu, Yisong Yue.
“Disentangling Linear Quadratic Control with Untrusted ML Predictions”
Preprint, 2024