

Faster RL by Freezing Slow States

Daniel Jiang (Meta, Applied RL)

joint work with Yijia Wang (Pitt)

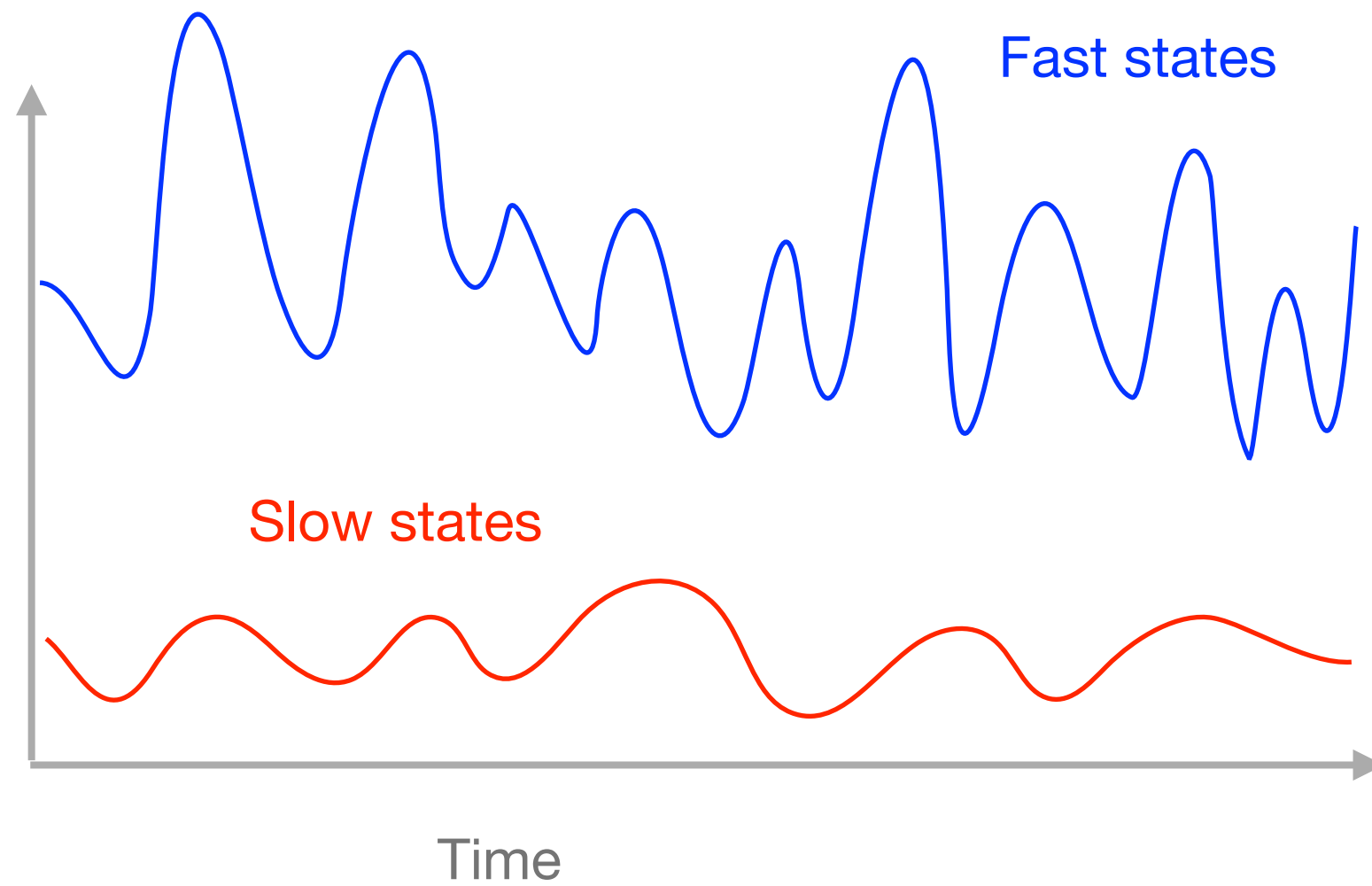
International Symposium on Mathematical Programming

Montreal

July 23, 2024

1. Motivation via example applications

Fast and slow states

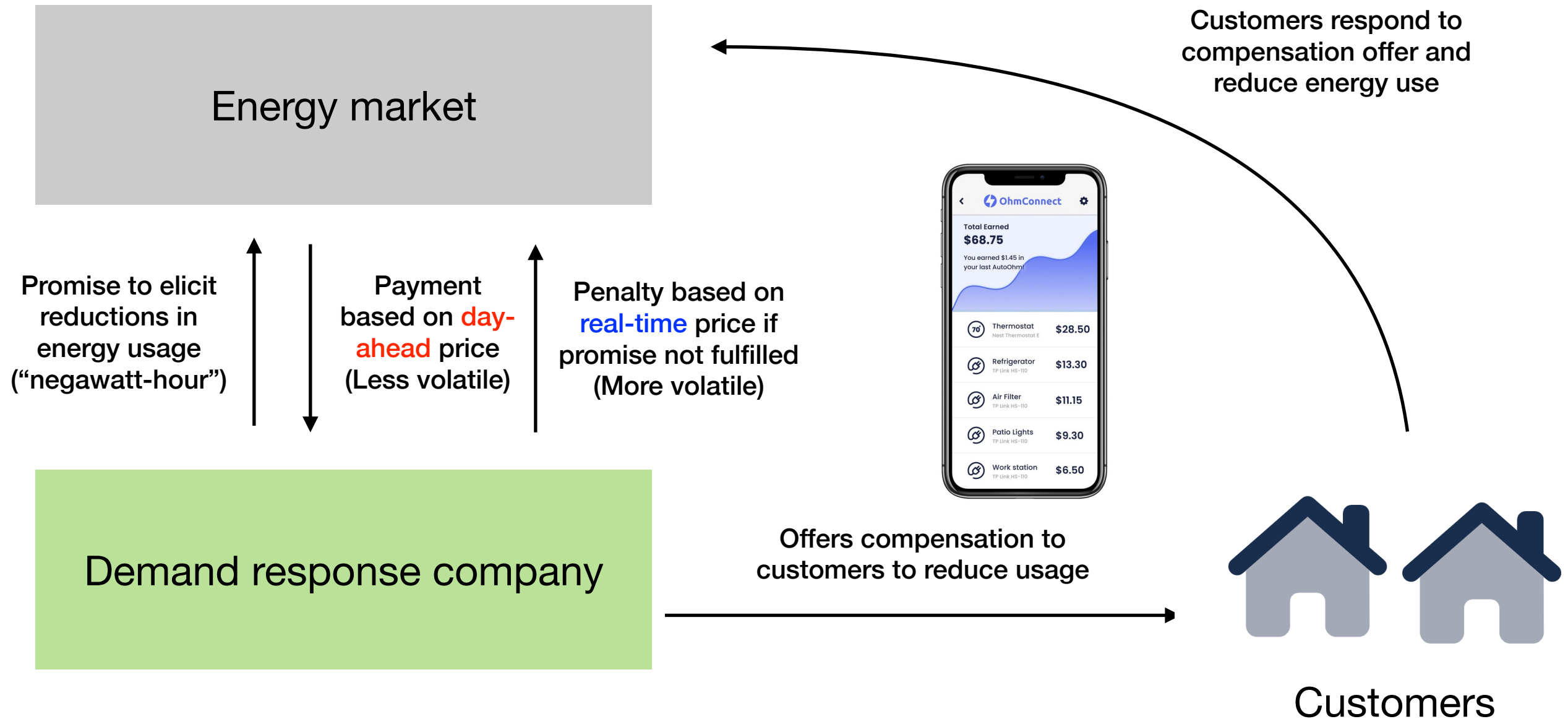


Recommendation systems

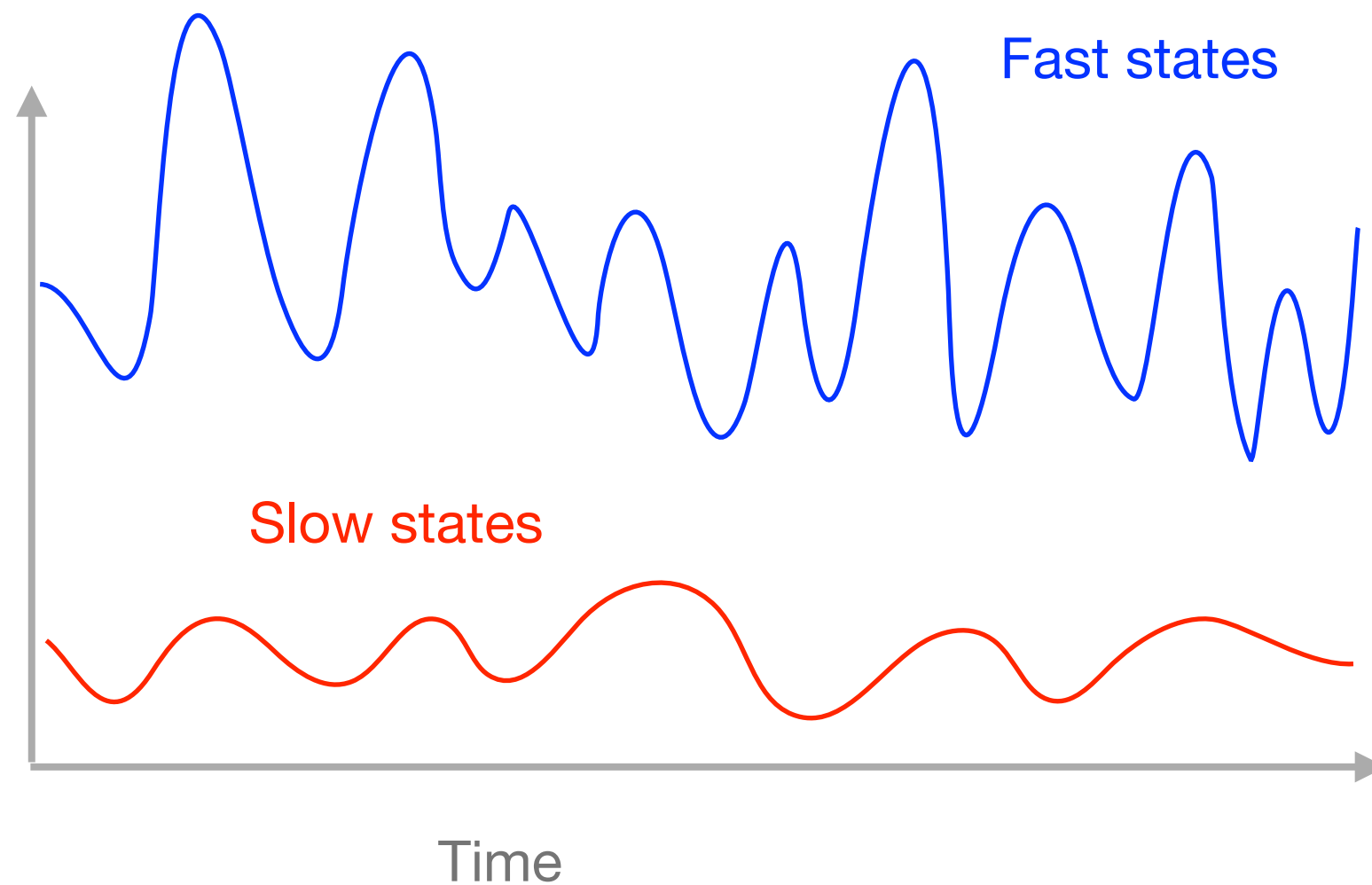
- Consider a recommendation system setting where:
 - Users return (i.e., log on) to the platform more often if they are seeing **content** that is *interesting* (Sumida & Zhou, 2023)
 - Users return to the platform more often if they have a *diverse* **recent content history**.
 - Users **interests** can shift over time as a function of the content they see.



Demand response (shifting electricity demand)



What do they have in common?



Fast states from examples:

- Real-time prices
- Recent content history

Shorter timescales

Slow states from examples:

- Day-ahead prices
- Underlying user interests

Longer timescales


Current practice when modeling a new problem

- While modeling an MDP, additional state variables is expensive:
 - Each iteration of value iteration $\mathcal{O}(S^2A)$
- **What do practitioners do (anecdotally)?**
 - If a state is deemed a “slow state” (contexts, environmental variables, etc), they might be *ignored/omitted*
 - e.g. assume costs are deterministic, demand is stationary, weather doesn’t change
- **This work:** A *compromise* between computational tractability and fully ignoring the slow state
 - We propose to *periodically* ignore slow states
 - We give evidence and argue that completely omitting slow states from the decision model is often not a viable heuristic



2. Fast-slow Markov decision processes

Fast-slow Markov decision processes

- A γ -discounted, infinite horizon MDP:
 - States $s \in \mathcal{S}$
 - Actions $a \in \mathcal{A}$
 - Rewards $r(s, a) \in [0, r_{\max}]$
 - Transition function
 - $s_{t+1} = f(s_t, a_t, w_{t+1}), w_{t+1} \in \mathcal{W}$
- Fast-slow MDP:
 - States $s = (x, y) \in \mathcal{S} = (\mathcal{X} \times \mathcal{Y})$ 
 - Actions $a \in \mathcal{A}$
 - Rewards $r(s, a) \in [0, r_{\max}]$
 - Transition function
 - $x_{t+1} = f_{\mathcal{X}}(s_t, a_t, w_{t+1})$
 - $y_{t+1} = f_{\mathcal{Y}}(s_t, a_t, w_{t+1})$

Main assumption (“fast-slow property”):

$$\|y - f_{\mathcal{Y}}(x, y, a, w)\|_2 \leq d_{\mathcal{Y}} \quad \text{and} \quad \|x - f_{\mathcal{X}}(x, y, a, w)\|_2 \leq \alpha d_{\mathcal{Y}}.$$

Lipschitz assumptions (let $U^{\star}(s)$ be the optimal value function):

$$|r(s, a) - r(s', a')| \leq L_r \|(s, a) - (s', a')\|_2,$$

$$\|f(s, a, w) - f(s', a', w)\|_2 \leq L_f \|(s, a) - (s', a')\|_2,$$

$$\|U^{\star}(s) - U^{\star}(s')\|_2 \leq L_U \|s - s'\|_2.$$

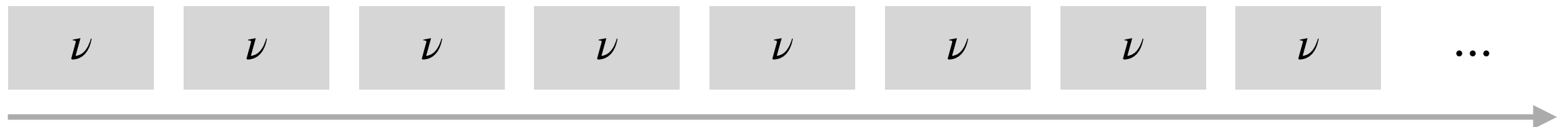
3. Hierarchical reformulation

Hierarchical reformulation (of any MDP)

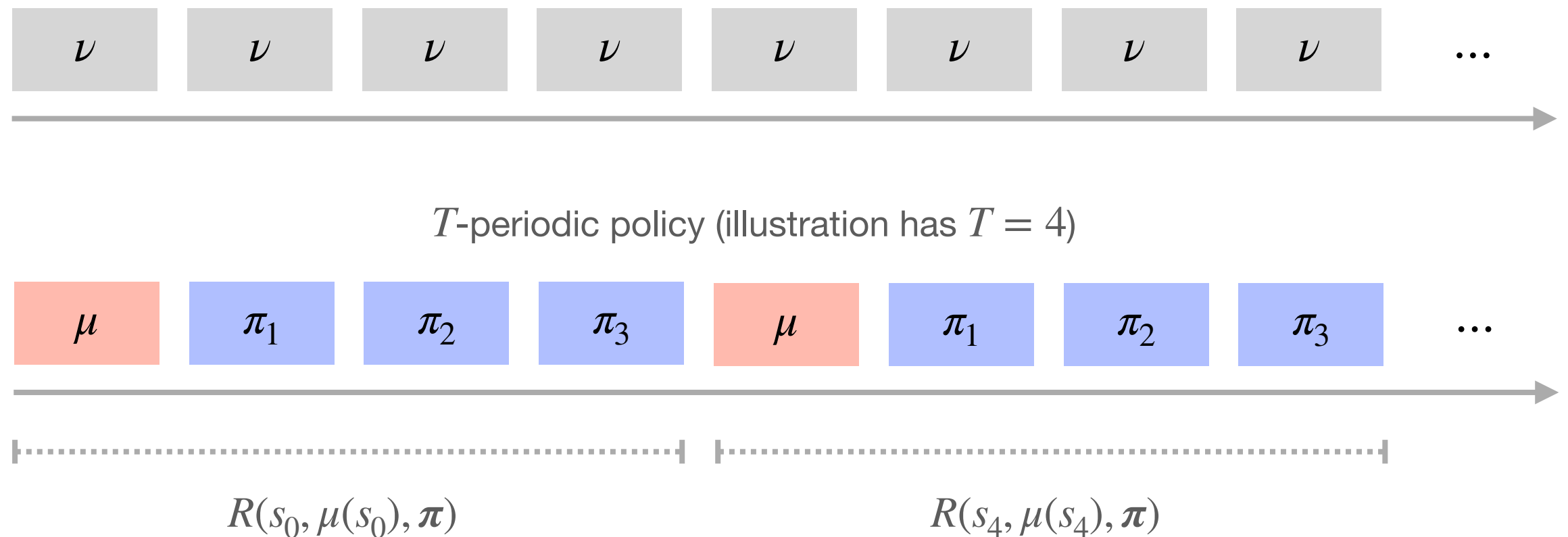
- A hierarchical reformulation is at the basis of our proposed approach
- Consider an MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{W}, f, r, \gamma \rangle$
- Let $\nu : \mathcal{S} \rightarrow \mathcal{A}$ be a stationary policy
- The value function is

$$U^\nu(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \nu) \mid s_0 = s \right] = r(s, \nu) + \gamma \mathbb{E} [U^\nu(s')]$$

- The policy can be thought of as (ν, ν, \dots) :



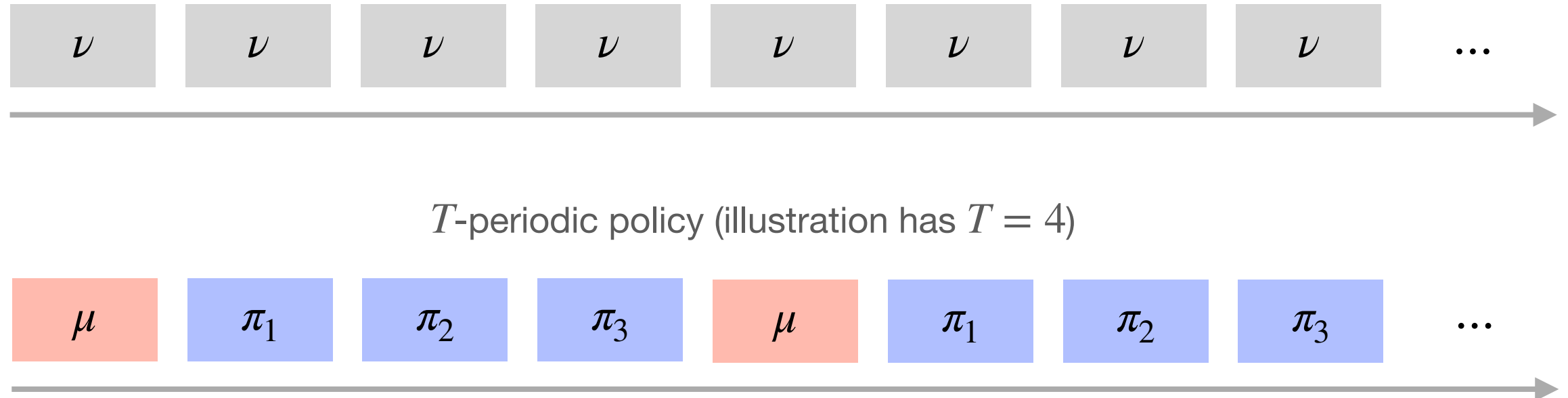
Hierarchical reformulation (of any MDP)



- Given a T -periodic policy $(\mu, \pi) = (\mu, \pi_1, \dots, \pi_{T-1})$, T -horizon reward is

$$R(s_0, \mu(s_0), \pi) = r(s_0, \mu) + \sum_{t=1}^{T-1} \gamma^t r(s_t, \pi_t)$$

Hierarchical reformulation (of any MDP)



- Bellman equations of the base model and its hierarchical reformulation are:

$$U^*(s_0) = \max_a r(s, a) + \gamma \mathbb{E}[U^*(s_1)]$$

$$\bar{U}^*(s_0) = \max_{(\mu, \pi)} \mathbb{E}[R(s_0, \mu(s_0), \pi) + \gamma^T \bar{U}^*(s_T)]$$

How can we take advantage of this?

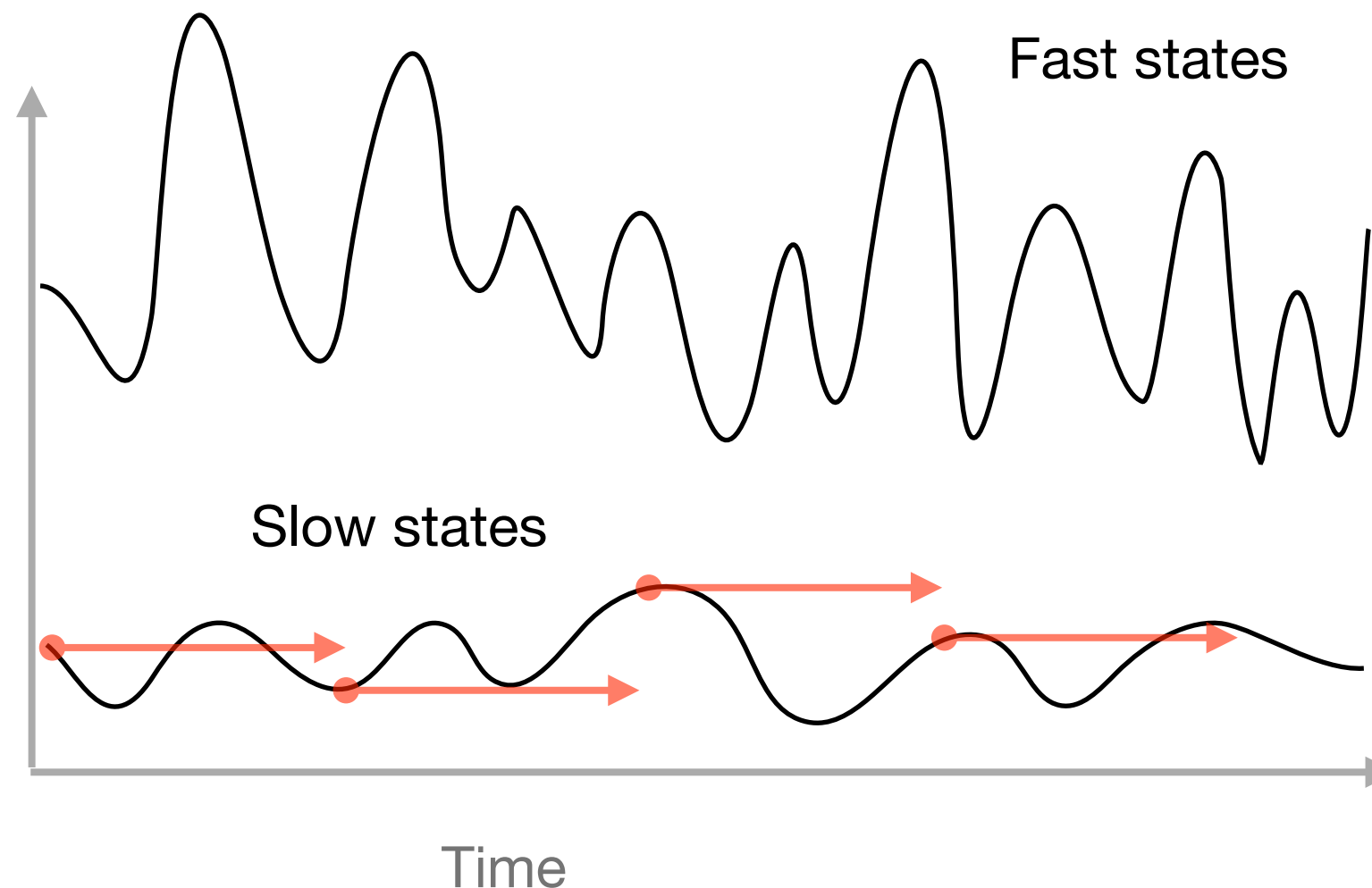
Proposition. The optimal values are equal: $U^*(s) = \bar{U}^*(s)$. Therefore, we can use the hierarchical reformulation as a basis for our approximation.

4. Frozen-state approximation and its regret

Frozen-state approximation



Main idea: Since slow states don't change much, let's **freeze** them for some number (T) of periods. Easier sub-problem with a smaller state space.



Frozen-state approximation



Main idea: Since slow states don't change much, let's **freeze** them for some number (T) of periods. Easier sub-problem with a smaller state space.

Implementation

1. At $t = 0$, take a “upper-level” action (using $\tilde{\mu}$), i.e., an action that considers the γ^T timescale
2. At $t = 1$, observe slow state and pretend it is frozen until $t = T$ and that $t = T$ is the end of the horizon
3. Solve this *easier* lower-level finite horizon problem.
4. Execute this T -period lower-level policy ($\tilde{\pi}_1, \tilde{\pi}_2, \dots, \tilde{\pi}_{T-1}$) in the real system
5. Repeat

Computation

- **Pre-compute** finite-horizon lower-level policy with frozen slow states
- Re-use pre-computed lower-level policy to solve infinite-horizon upper-level problem, which **takes advantage of γ^T**

Frozen-state, lower-level problem



The true problem

(Or, we can use any VFA we would like as the terminal value.)

Frozen-state lower-level MDP

$$J_1^\star(x, y) = \max_{\tilde{\pi}} \mathbb{E} \left[\sum_{t=1}^{T-1} \gamma^{t-1} r(\mathbf{x}_1, y_t, \tilde{\pi}_t) \mid (x_1, y_1) = (x, y) \right]$$

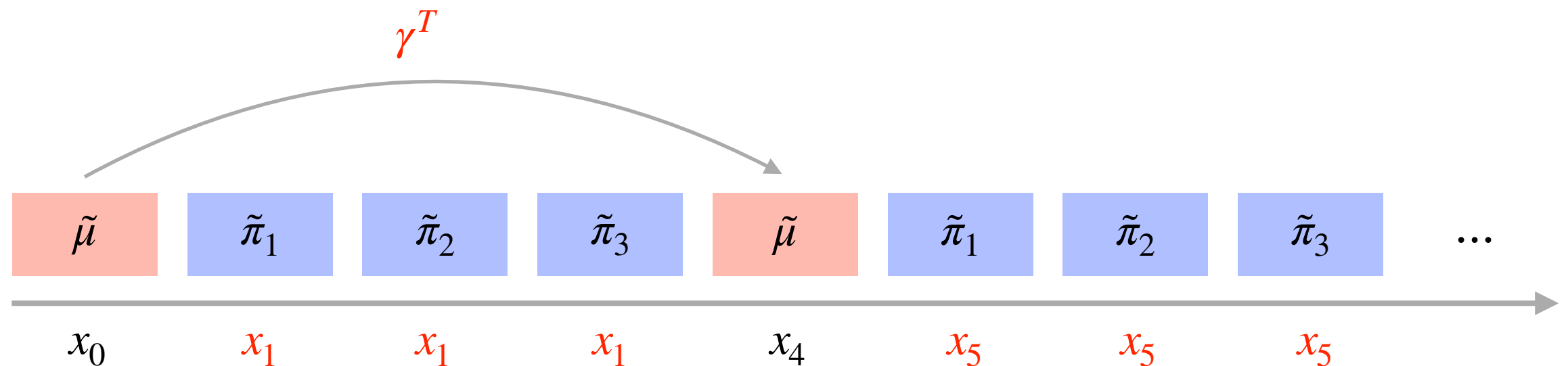
$$J_t^\star(x, y) = \max_a r(x, y, a) + \gamma \mathbb{E}[J_{t+1}^\star(x, y')], \quad J_T^\star \equiv 0$$

$$\tilde{\pi}_t^\star(x, y) = \operatorname{argmax}_a r(x, y, a) + \gamma \mathbb{E}[J_{t+1}^\star(x, y')].$$

Computational benefits

- Small number of successor states (since slow state is frozen)
 - $\mathcal{O}(S^2A) \rightarrow \mathcal{O}(XY^2A)$
- Independent across x
- Independent from upper-level problem (replaced U^\star by 0)

Frozen-state, upper-level problem



Frozen-state upper-level MDP

Let $(\tilde{\pi}^*, J_1^*)$ be the optimal policy/value of the lower-level problem.

$$\tilde{R}(s_0, a, J_1^*) = r(s_0, a) + \gamma J_1^*(f(s_0, a, w))$$

$$V^*(s_0, J_1^*, \tilde{\pi}^*) = \max_a \mathbb{E} [\tilde{R}(s_0, a, J_1^*) + \gamma^T V^*(s_T, J_1^*, \tilde{\pi}^*)] \text{ (transitions according to } \tilde{\pi}^*)$$

After solving both levels, let $(\tilde{\mu}^*, \tilde{\pi}^*)$ be the solution of the frozen-state approximation.

In the exact reformulation, we were maximizing over policies, now it is just a single action.

Per-cycle reward approximation error

Proposition. The difference between true and approximate T -horizon rewards:

$$\begin{aligned}
 & \left| \underbrace{\mathbb{E}[R(s_0, a, \pi^\star)]}_{\text{True}} - \underbrace{\mathbb{E}[\tilde{R}(s_0, a, J_1^\star)]}_{\text{Frozen}} \right| \\
 & \leq \underbrace{\alpha d_{\mathcal{Y}} \left(L_r \sum_{i=1}^{T-2} \gamma^i \sum_{j=0}^{i-1} L_f^j \right)}_{\text{error from freezing}} + \underbrace{\gamma^{T-1} L_U \left[\alpha d_{\mathcal{Y}} \sum_{j=0}^{T-2} L_f^j + \gamma d_{\mathcal{Y}} (\alpha + 2)(T - 1) \right]}_{\text{end of horizon error}}
 \end{aligned}$$

Main ideas.

$$1. \quad \mathbb{E}[R(x_0, y_0, a, \pi^\star)] = \mathbb{E} \left[r(x_0, y_0, a) + \gamma (H^{T-1} U^\star)(x_1, y_1) - \gamma^T U^\star(x_T, y_T) \right]$$

where $(HU)(x, y) = \max_a r(x, y, a) + \gamma \mathbb{E}[f(s, a, w)]$ (true Bellman operator)

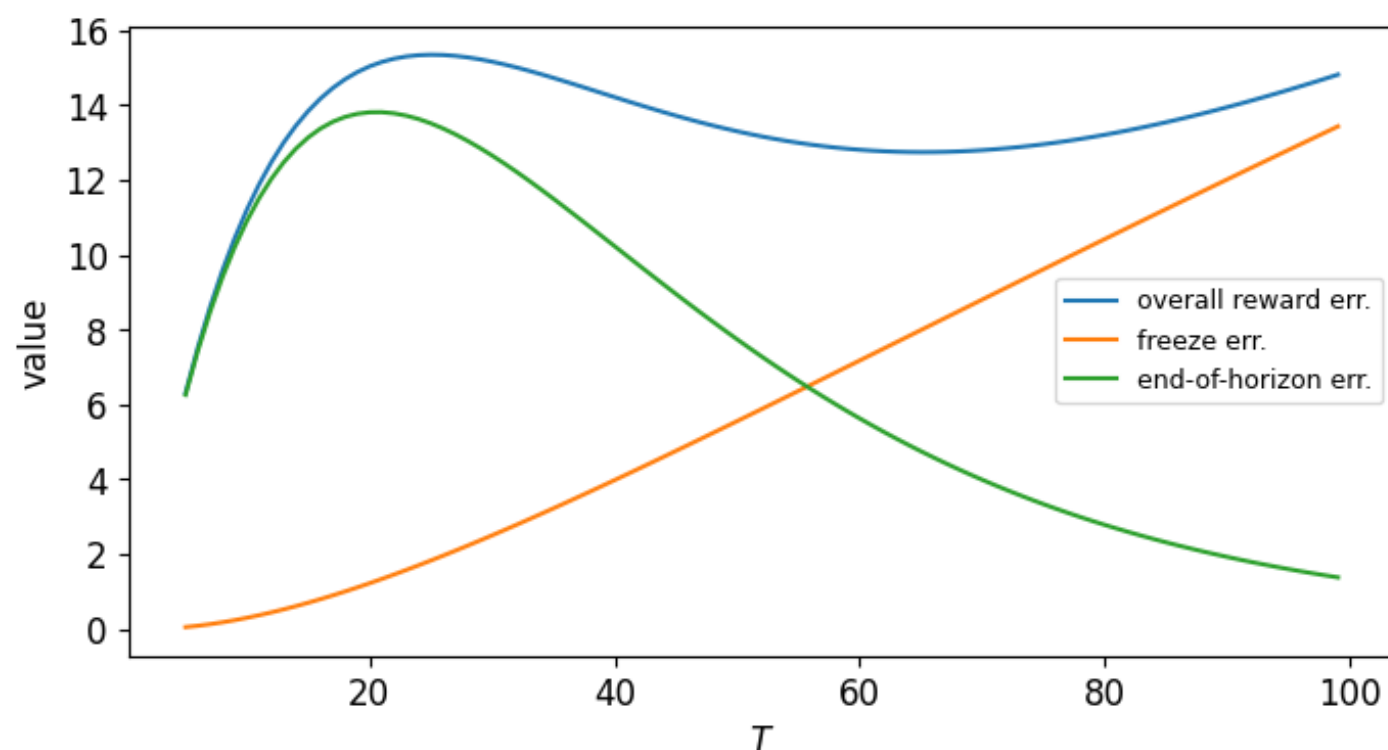
$$2. \quad \mathbb{E}[\tilde{R}(x_0, y_0, a, J_1^\star)] = r(x_0, y_0, a) + \gamma (\tilde{H}^{T-1} \mathbf{0})(x_1, y_1)$$

where $(\tilde{H}J_{t+1})(x, y) = \max_a r(x, y, a) + \gamma \mathbb{E}[J_{t+1}(x, f_{\mathcal{Y}}(x, y, a, w))]$ (frozen Bellman operator)

Per-cycle reward approximation error

Proposition. The difference between true and approximate T -horizon rewards:

$$\begin{aligned}
 & \left| \underbrace{\mathbb{E}[R(s_0, a, \pi^\star)]}_{\text{True}} - \underbrace{\mathbb{E}[\tilde{R}(s_0, a, J_1^\star)]}_{\text{Frozen}} \right| \\
 & \leq \underbrace{\alpha d_{\mathcal{Y}} \left(L_r \sum_{i=1}^{T-2} \gamma^i \sum_{j=0}^{i-1} L_f^j \right)}_{\text{error from freezing}} + \underbrace{\gamma^{T-1} L_U \left[\alpha d_{\mathcal{Y}} \sum_{j=0}^{T-2} L_f^j + \gamma d_{\mathcal{Y}} (\alpha + 2)(T - 1) \right]}_{\text{end of horizon error}}
 \end{aligned}$$



5. A new algorithm: *Frozen-state value iteration*

Standard value iteration on the base model

Recall: Given an MDP and Bellman operator H , where $(HU)(s) = \max_a r(s, a) + \gamma \mathbb{E} U(f(s, a, w))$, the *value iteration* algorithm is $U^k = H^k U^0$

• Convergence to optimal value function: $\lim_{t \rightarrow \infty} H^t U = U^*$ for any initial estimate U

• $\|U^{\nu^k} - U^*\|_\infty \leq \frac{2r_{\max}\gamma^{k+1}}{(1-\gamma)^2}$, where $\nu^k(s) = \operatorname{argmax}_a r(s, a) + \gamma \mathbb{E}[U^k(f(s, a, w))]$

Depends on

- $\|U^k - U^*\|_\infty \leq \gamma^k \|U^0 - U^*\|_\infty$
- $\|U^0 - U^*\|_\infty \leq \frac{r_{\max}}{1-\gamma}$
- $\|U^{\nu^k} - U^*\|_\infty \leq \frac{2\|U^k - U^*\|_\infty}{1-\gamma}$

Algorithm 1: Exact VI for the Base Model

Input: Initial values U_0 , number of iterations k .

Output: Approximation to the optimal policy ν^k .

```
1 for  $i = 1, 2, \dots, k$  do
2   for  $s$  in the state space  $\mathcal{S}$  do
3      $U^i(s) = \max_a r(s, a) + \gamma \mathbb{E}[U^{i-1}(f(s, a, w))]$ .
4   end
5 end
6 for  $s$  in the state space  $\mathcal{S}$  do
7    $\nu^k(s) = \operatorname{argmax}_a r(s, a) + \gamma \mathbb{E}[U^k(f(s, a, w))]$ .
8 end
```

Frozen-state value iteration (FSVI)

Algorithm 2: Frozen-State Value Iteration (FSVI)

Input: Initial values $J_T^* \equiv 0$ and V^0 , number of iterations k .

Output: Approximation of the T -periodic frozen-state policy $(\tilde{\mu}^k, \tilde{\pi}^*)$ and J_1^* .

```

1 for  $t = T - 1, T - 2, \dots, 1$  do
2   for each slow state  $x \in \mathcal{X}$  do
3     for each fast state  $y \in \mathcal{Y}$  do
4        $J_t^*(x, y) = \max_a r(x, y, a) + \gamma \mathbb{E}[J_{t+1}^*(x, f_{\mathcal{Y}}(x, y, a, w))]$ .
5        $\tilde{\pi}_t^*(x, y) = \arg \max_a r(x, y, a) + \gamma \mathbb{E}[J_{t+1}^*(x, f_{\mathcal{Y}}(x, y, a, w))]$ .
6     end
7   end
8 end

9 for  $i = 1, 2, \dots, k$  do
10  for  $s_0 = (x_0, y_0)$  in the state space  $\mathcal{X} \times \mathcal{Y}$  do
11     $V^i(x_0, y_0, J_1^*, \tilde{\pi}^*) = \max_a \mathbb{E}[\tilde{R}(s_0, a, J_1^*) + \gamma^T V^{i-1}(x_T, y_T, J_1^*, \tilde{\pi}^*)]$ .
12  end
13 end

14 for  $s_0 = (x_0, y_0)$  in the state space  $\mathcal{X} \times \mathcal{Y}$  do
15    $\tilde{\mu}^k(x_0, y_0) = \arg \max_a \mathbb{E}[\tilde{R}(s_0, a, J_1^*) + \gamma^T V^k(x_T, y_T, J_1^*, \tilde{\pi}^*)]$ .
16 end

```

Note: Freezing the state only happens “within” the algorithm to more efficiently compute J_1^*

Solving the lower level
incurs a **one time fixed cost**

Pre-compute lower-level problem, a finite-horizon DP:

- To solve lower-level DP: $\mathcal{O}(XY^2AT)$
- To compute multi-step transition: $\mathcal{O}(S^2T)$

Upper-level problem (infinite-horizon VI on slow-timescale MDP with γ^T discounting):

- Per upper-level VI iteration: $\mathcal{O}(S^2A)$

$\mathcal{O}(S^2A)$ per iteration is the same as Base VI, but now the discount factor is γ^T instead of γ !

Regret of a periodic policy (μ, π)

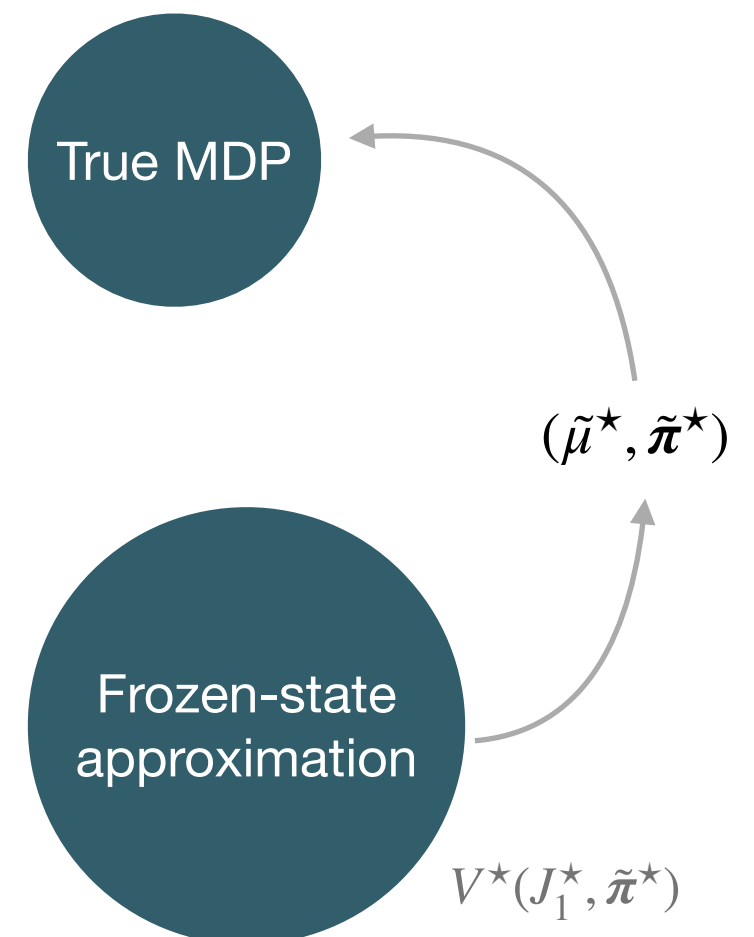
Definition. Suppose the optimal policy is ν^\star . The regret is

$$\mathcal{R}(s, \mu, \pi) = U^{\nu^\star}(s) - \bar{U}^\mu(s, \pi) = \bar{U}^\star(s) - \bar{U}^\mu(s, \pi) \quad \text{and} \quad \mathcal{R}(\mu, \pi) = \max_s \mathcal{R}(s, \mu, \pi),$$

where we have used the equivalence between the base and hierarchical formulations.

Remarks:

- We always measure regret with respect to the *true* MDP.
 - Although (μ, π) is computed **with the help of frozen states**, it is evaluated in the original MDP with true dynamics.
- Consider $\mathcal{R}(\tilde{\mu}^\star, \tilde{\pi}^\star)$, notice that $V^\star(J_1^\star, \tilde{\pi}^\star)$ does not directly enter the regret definition.
 - It is the optimal value of the approximation, but doesn't reflect the performance of $(\tilde{\mu}^\star, \tilde{\pi}^\star)$ in the true model.



Main idea behind regret analysis

Lemma (Approximation to FSVI).

- Suppose we *approximately* solve the lower-level problem and obtain π, J_1 , instead of the optimal solutions π^\star, U^\star .
- Suppose we approximately solve the upper-level problem and obtain V instead of $V^\star(J_1, \pi)$.
- Let μ be greedy with respect to both J_1 and V :
 - $\mu(s_0) = \operatorname{argmax}_{a \in \mathcal{A}} \mathbb{E} [\tilde{R}(s_0, a, J_1) + \gamma^T V(s_T(a, \pi))]$.
- Then,

$$\mathcal{R}(\mu, \pi) \leq \left(\frac{2\gamma^T}{(1 - \gamma^T)^2} + \frac{2}{1 - \gamma^T} \right) \epsilon_r(\pi^\star, J_1) + \left(\frac{2\gamma^{2T}}{(1 - \gamma^T)^2} + \frac{2\gamma^T}{1 - \gamma^T} \right) L_U d(\alpha, d_{\mathcal{Y}}, T) + \frac{2\gamma^T}{1 - \gamma^T} \|V^\star(J_1, \pi) - V\|_\infty.$$

End of horizon error

V-approximation error

Reward error

Regret of FSVI

Theorem. The regret of FSVI after k upper-level iterations is:

$$\begin{aligned} \mathcal{R}(\mu, \pi) \leq & \left(\frac{2\gamma^T}{(1 - \gamma^T)^2} + \frac{2}{1 - \gamma^T} \right) \epsilon_r(\pi^\star, J_1) \\ & + \left(\frac{2\gamma^{2T}}{(1 - \gamma^T)^2} + \frac{2\gamma^T}{1 - \gamma^T} \right) L_U d(\alpha, d_{\mathcal{Y}}, T) + \frac{2r_{\max}\gamma^{(k+1)T}}{(1 - \gamma)(1 - \gamma^T)}, \end{aligned}$$

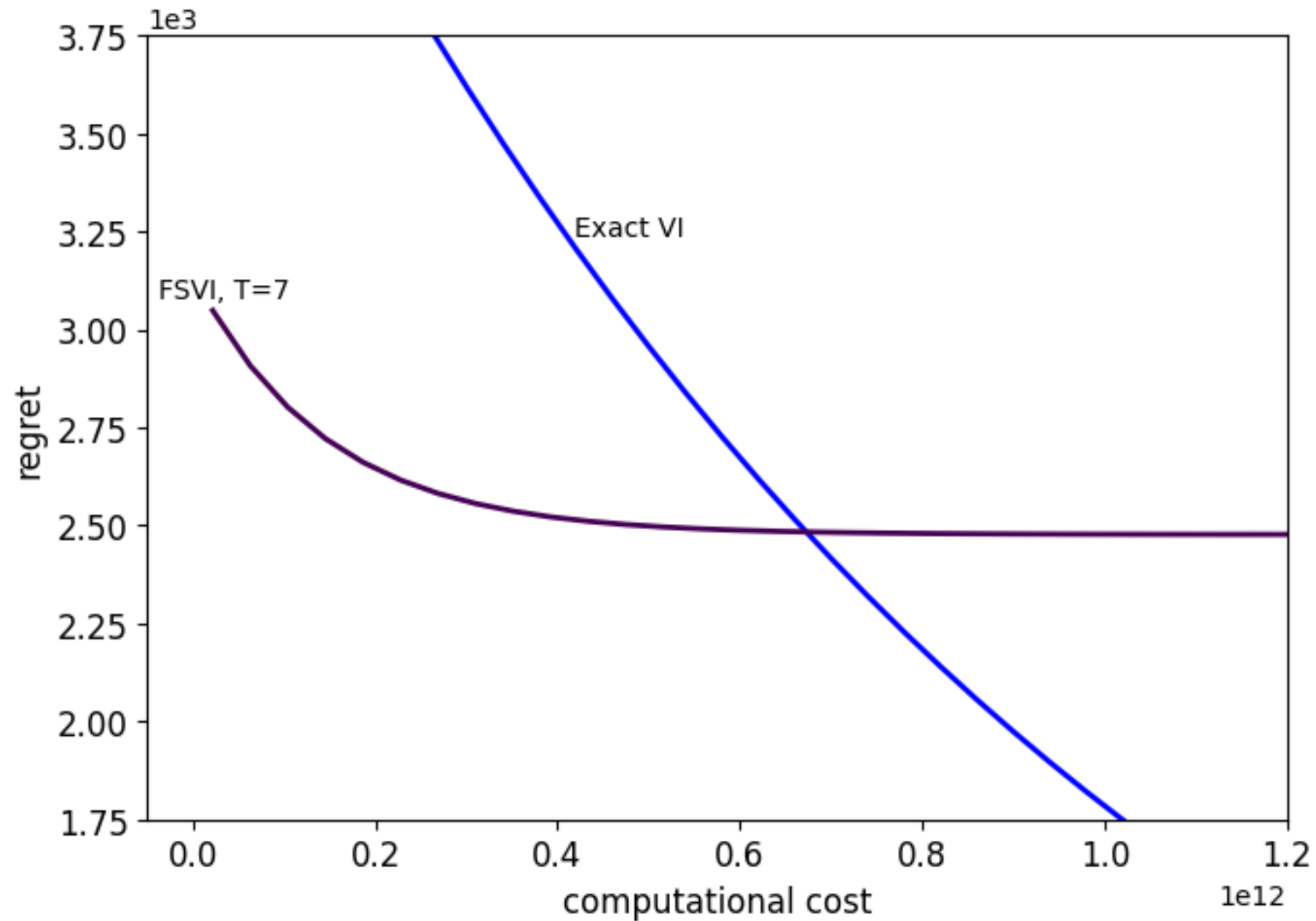
which replaces the V-approximation error term with the VI error.

Main question to answer using the regret analysis

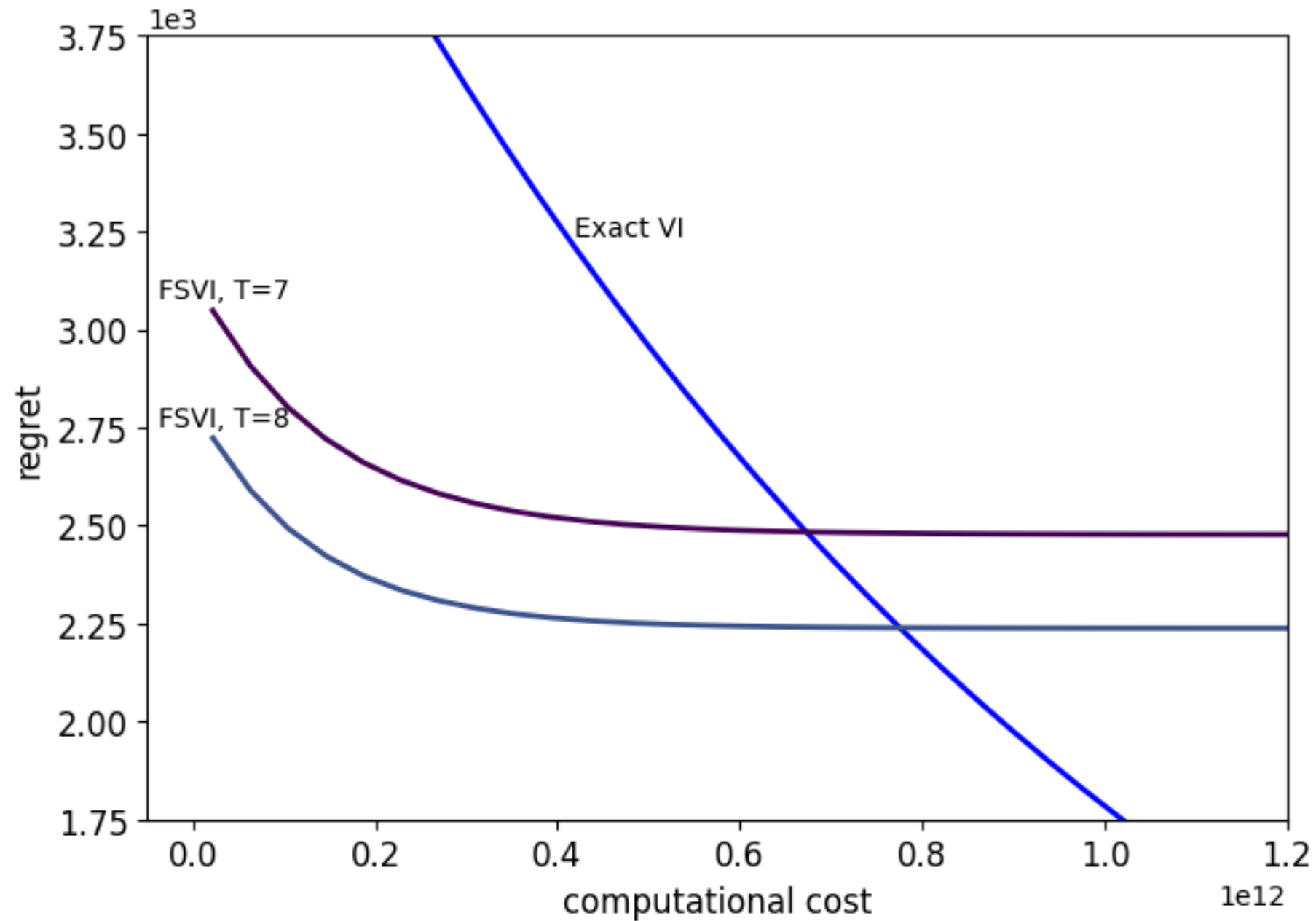
Given some computational budget, should we use the FSVI algorithm at all?

If so, how should we choose T (the number of periods to freeze) and k (number of upper-level steps iterations)?

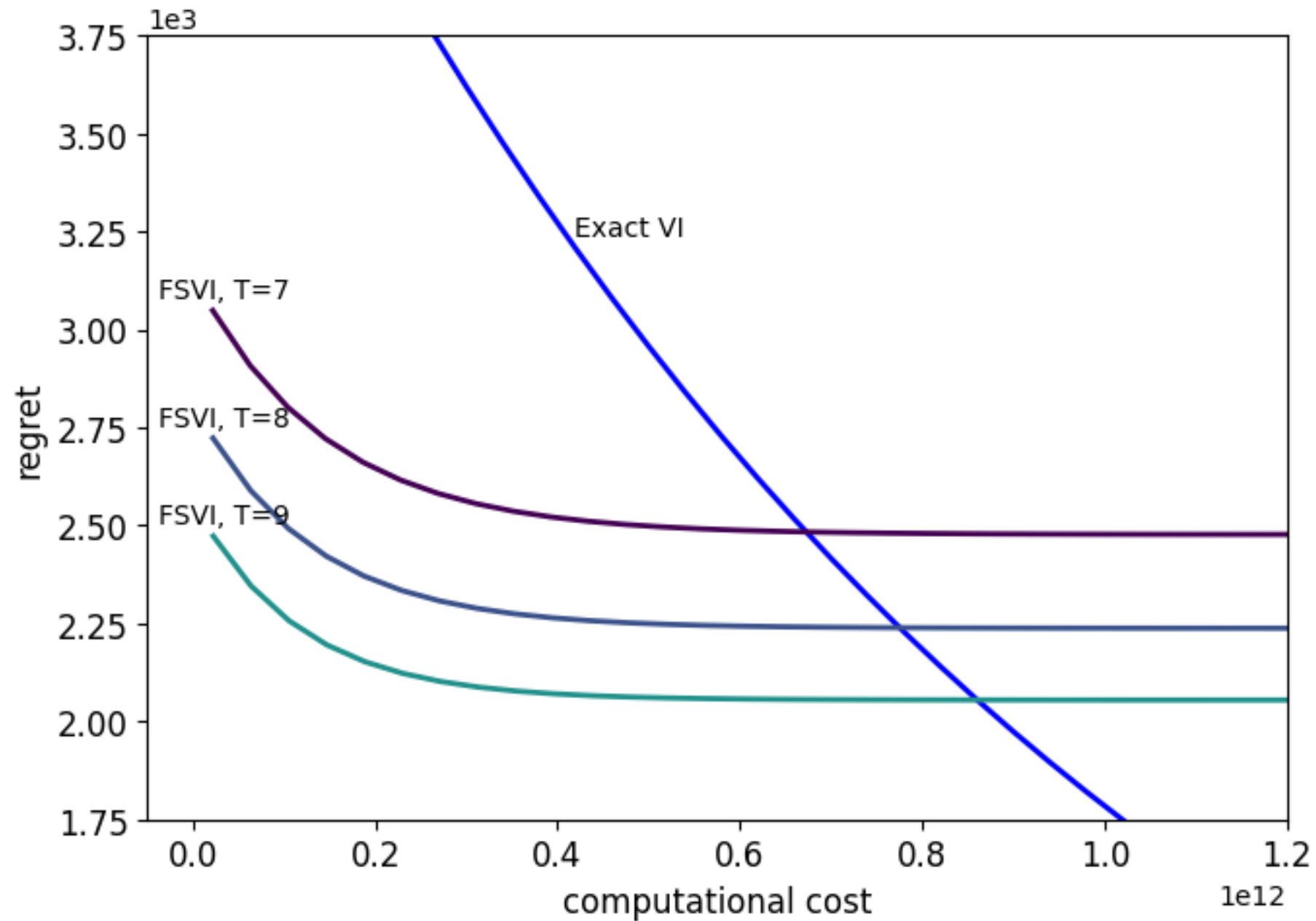
How can we use the bound to choose T?



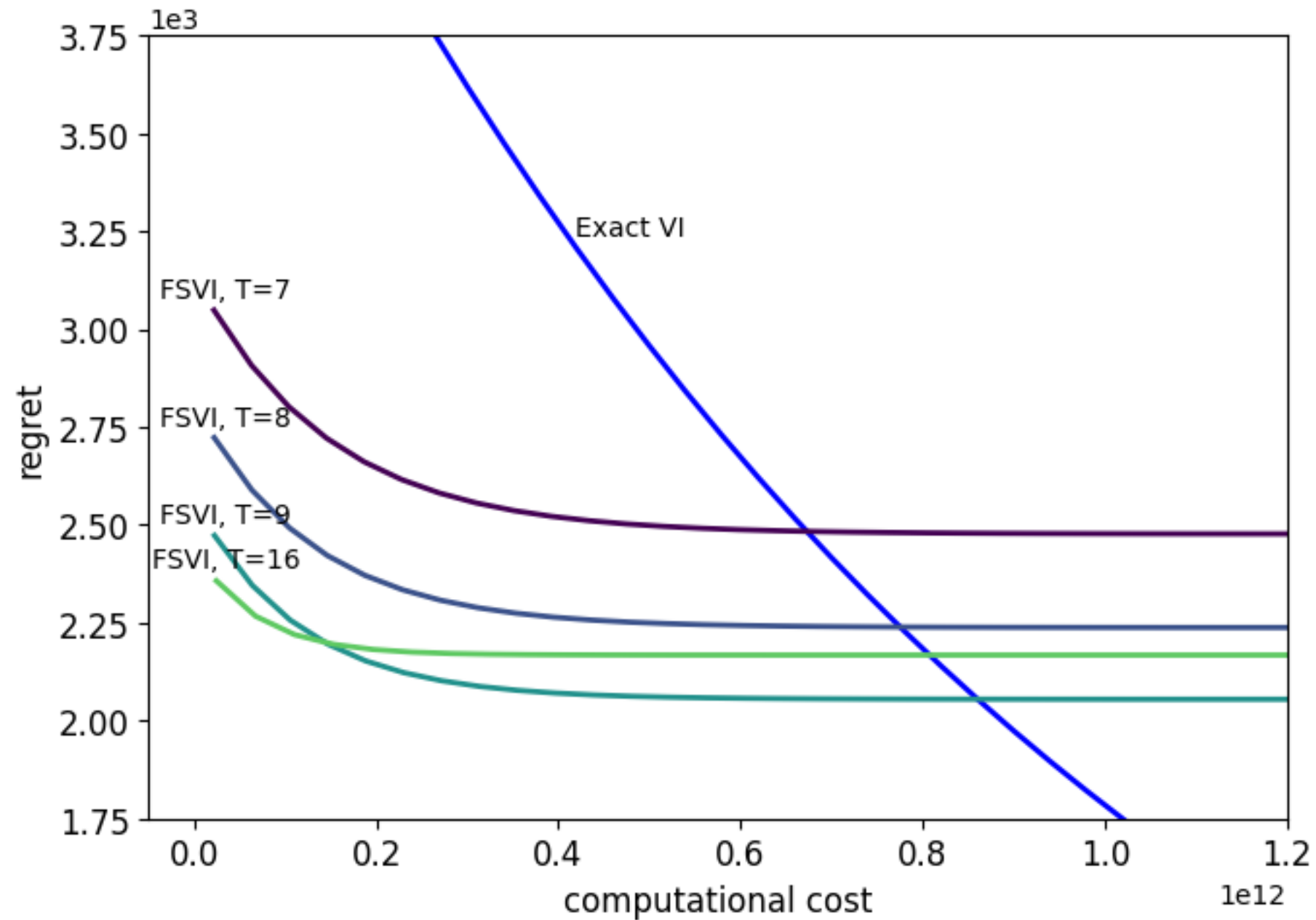
How can we use the bound to choose T?



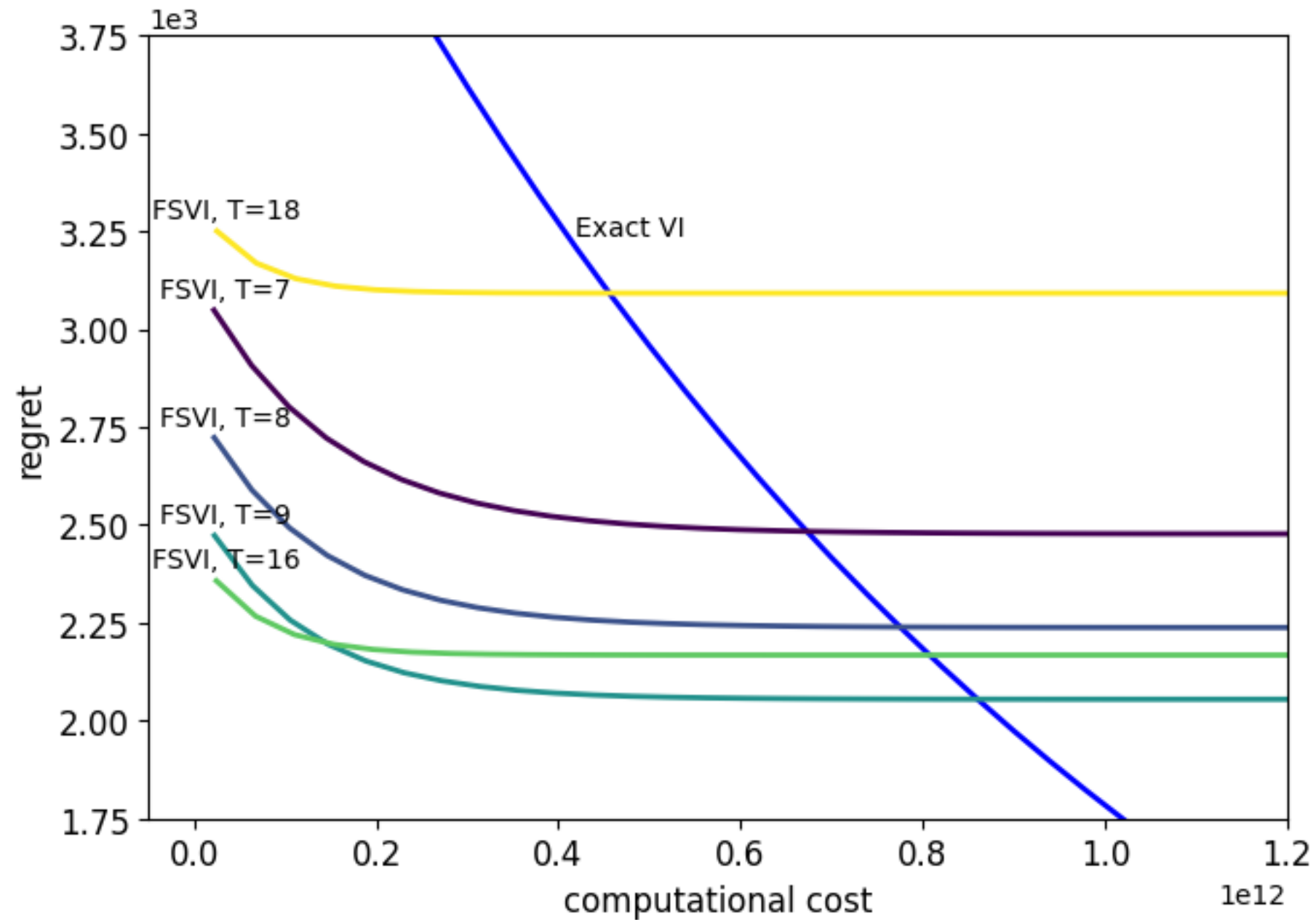
How can we use the bound to choose T?



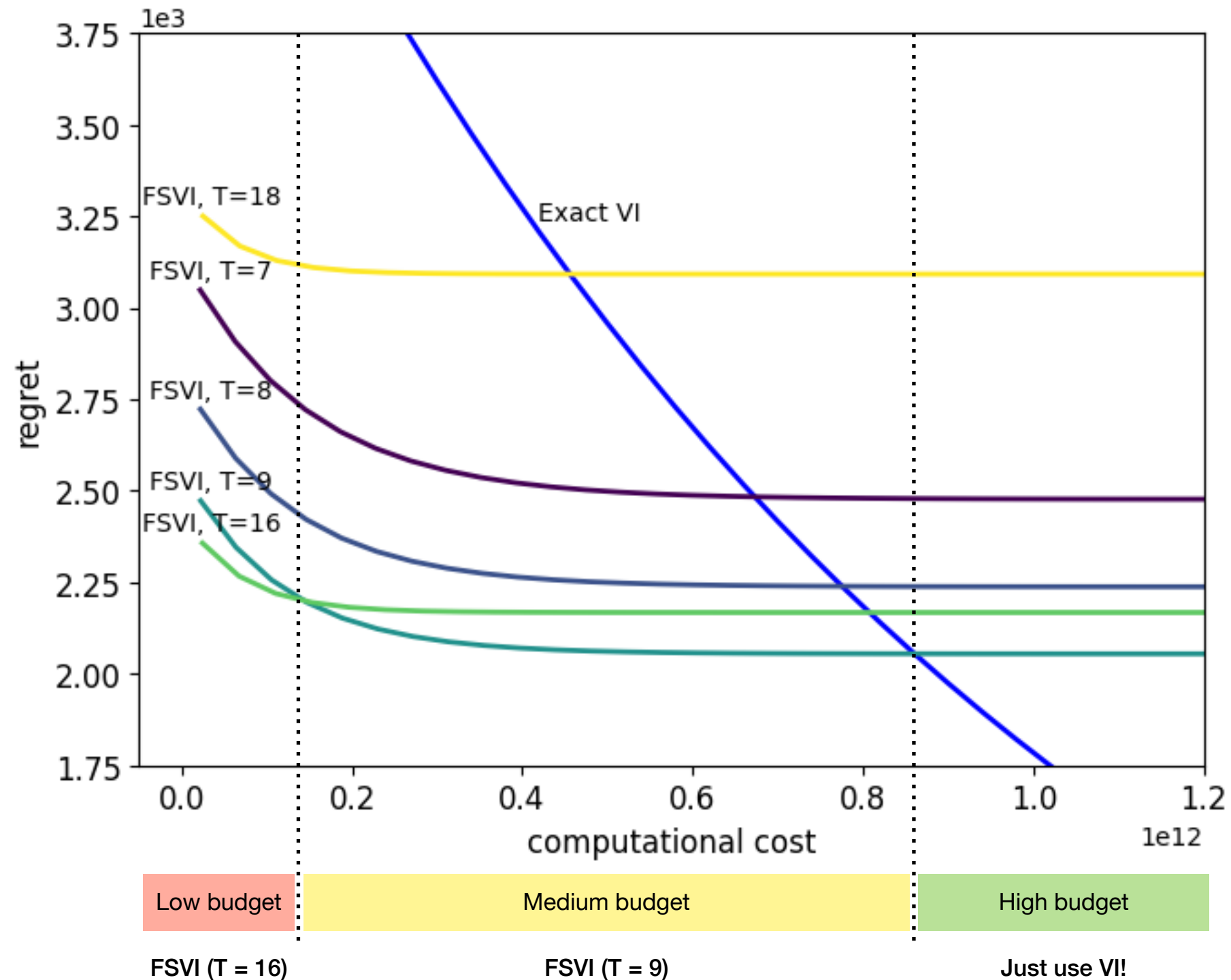
How can we use the bound to choose T?



How can we use the bound to choose T?



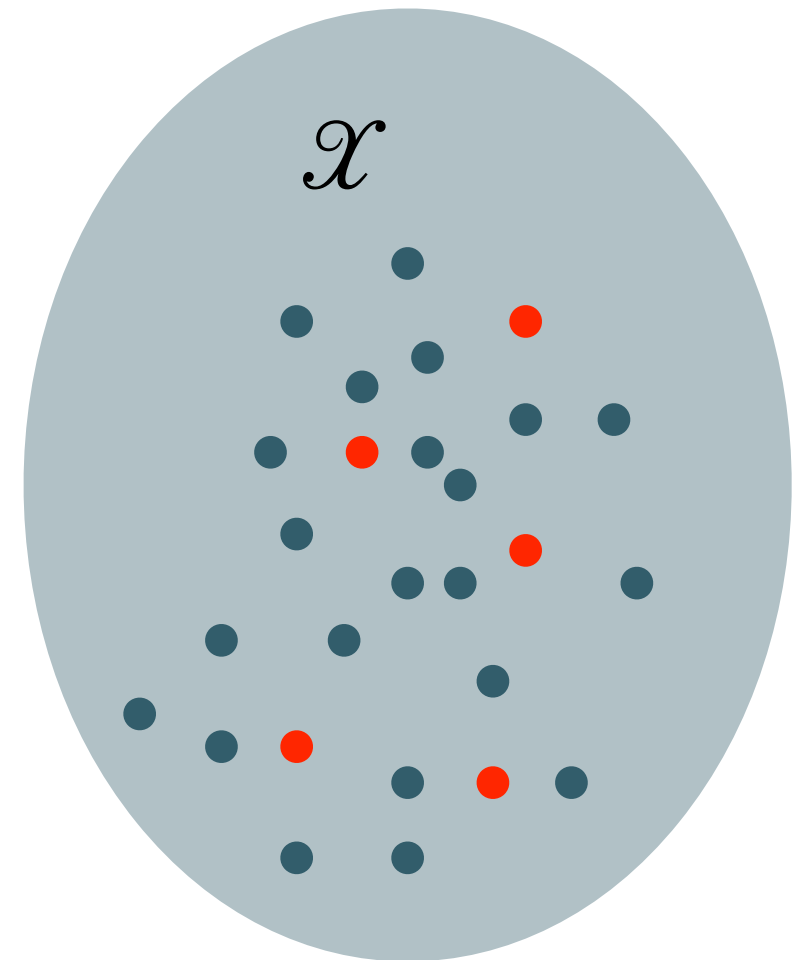
How can we use the bound to choose T?



6. Extensions

Extension: Nominal state version of FSVI

- In FSVI, we have to solve the lower-level MDP for each x
- We can further do approximations by solving the lower-level MDP for a few nominal states only
 - $\mathcal{O}(S^2A) \rightarrow \mathcal{O}(XY^2A) \rightarrow \mathcal{O}(X_{\text{nom}}Y^2A)$
- Later, extrapolate to nearby nominal states
- Theoretical results can be adapted given an additional assumption on the MDP rewards



Extension: Scaling to larger state spaces using feature-based approximate value iteration

Architecture:

- Consider M pre-selected states $\tilde{\mathcal{S}} = \{s_1, s_2, \dots, s_M\}$.
- Consider an M -dimensional feature vector $\phi(s)$, where $\phi(s_m)$ are linearly independent.
- Assume there exists $\gamma' \in [\gamma, 1)$ s.t. for any s , there exists $\theta_m(s)$, where

$$\sum_m |\theta_m(s)| \leq 1 \text{ and } \phi(s) = \frac{\gamma'}{\gamma} \sum_{m=1}^M \theta_m(s) \phi(s_m).$$

- Lower level: $\hat{J}(s, \omega_t) = \phi^\top(s) \omega_t$.
- Upper level: $\hat{V}(s, \beta^k) = \phi^\top(s) \beta^k$.
- Update procedure:
 1. Compute Bellman update at pre-selected states only: $y(s_m)$.
 2. Compute next parameter vector (ω_{t-1} or β^{k+1}) such that the updated value function evaluated at the pre-selected states is equal: e.g., $\hat{J}(s_m, \omega_{t-1}) = y(s_m)$.

Limited “expansion” after going to parameter space and back:

$$\|(\Phi\Phi^\dagger)(J) - (\Phi\Phi^\dagger)(J')\|_\infty \leq \kappa \|J - J'\|_\infty \quad (\kappa = \gamma'/\gamma)$$

Extension: Scaling to larger state spaces using feature-based approximate value iteration

Algorithm 4: Frozen-State Approximate Value Iteration (FSAVI)

Input: $\tilde{\mathcal{S}} = \{s_1, s_2, \dots, s_M\}$, ϕ , initial weights $\omega_T = \beta_0 = \mathbf{0}$, number of iterations k .

Output: Approximation of the T -periodic frozen-state policy $(\hat{\mu}_{(\beta^k, \omega^*)}, \hat{\pi}_{\omega^*})$ and $\hat{J}_1(\omega^*)$

```
1 for  $t = T - 1, T - 2, \dots, 1$  do
2   for each pre-selected state  $s = (x, y) \in \tilde{\mathcal{S}}$  do
3      $J_t(x, y) = \max_a r(x, y, a) + \gamma \mathbb{E}[\hat{J}_{t+1}(x, f_Y(x, y, a, w), \omega_{t+1})]$ .
4   end
5   Set remaining entries of  $J_t$  to zero. Update parameter vector:  $\omega_t^* = \Phi^\dagger J_t$ .
6 end
7 Let  $\hat{\pi}_{\omega^*}$  be greedy with respect to  $\hat{J}_t(\omega_t^*) = \Phi \omega_t^*$ , similar to (23).
8 for  $i = 1, 2, \dots, k$  do
9   for each pre-selected state  $s_0 \in \tilde{\mathcal{S}}$  do
10     $V^i(s_0) = \max_a \mathbb{E}[\tilde{R}(s, a, \hat{J}_1(\omega_1^*)) + \gamma^T \hat{V}(s_T(a, \tilde{\pi}_{\text{avi}}), \beta_{i-1})]$ .
11    Set remaining entries of  $V^i$  to zero. Update parameter vector:  $\beta_i = \Phi^\dagger V^i$ .
12  end
13 end
14 for  $s_0$  in the state space  $\mathcal{S}$  do
15    $\hat{\mu}_{(\beta^k, \omega^*)}(s_0) = \arg \max_a \mathbb{E}[\tilde{R}(s_0, a, \hat{J}_1(\omega_1^*)) + \gamma^T \hat{V}(s_T(a, \tilde{\pi}_{\omega^*}), \beta_k)]$ .
16 end
```

Regret of FSAVI

Theorem. The regret of FSAVI after k upper-level iterations is:

$$\begin{aligned} \mathcal{R}(\mu, \pi) \leq & \left(\frac{2\gamma^T}{(1-\gamma^T)^2} + \frac{2}{1-\gamma^T} \right) \epsilon_r(\pi^*, \hat{J}_1(\omega_1^*)) \\ & + \left(\frac{2\gamma^{2T}}{(1-\gamma^T)^2} + \frac{2\gamma^T}{1-\gamma^T} \right) L_U d(\alpha, d_{\mathcal{Y}}, T) + \underbrace{\left(\frac{1+\kappa}{1-\kappa\gamma^T} \right) \epsilon_{\text{up}}}_{\|V_{\omega^*}^* - \hat{V}(\beta^*)\|_{\infty}} + \underbrace{(\kappa\gamma^T)^k \left(\frac{\kappa^2 - \kappa^2(\kappa\gamma)^{T+1}}{(1-\kappa\gamma^T)(1-\kappa\gamma)} \right) r_{\max}}_{\|\hat{V}(\beta^*) - \hat{V}(\beta_k)\|_{\infty}}, \end{aligned}$$

which $\epsilon_r(\pi^*, \hat{J}_1(\omega_1^*)) = \epsilon_r(\pi^*, J_1^*) + \underbrace{\left(\frac{1+\kappa}{1-\kappa\gamma} - \frac{(\kappa\gamma)^T(1+\gamma)}{\gamma - \kappa\gamma^2} \right) \epsilon_{\text{low}}}_{\|J_1^* - \hat{J}_1(\omega_1^*)\|_{\infty}}.$

Upper-level feature approximation error

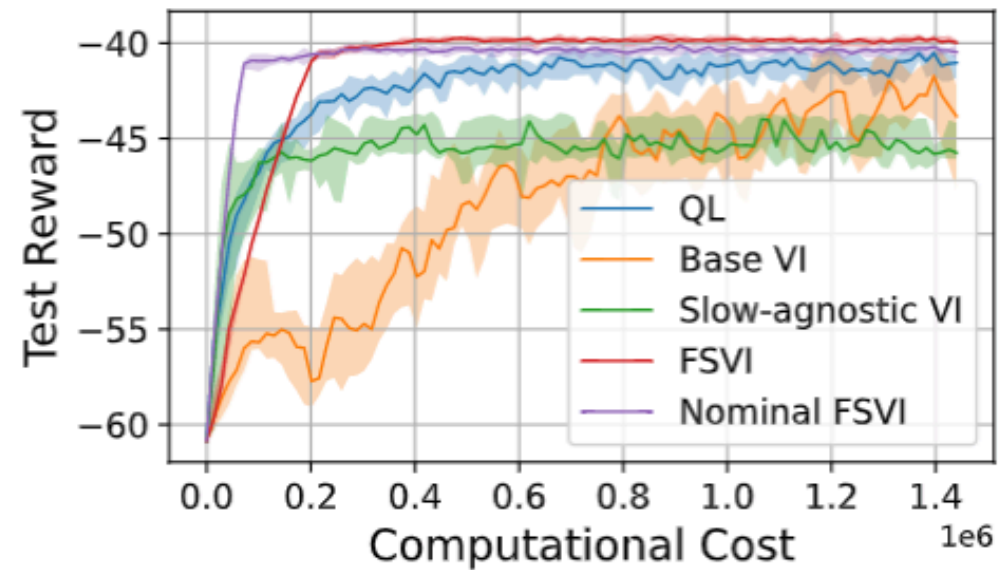
Lower-level feature approximation error

7. Numerical results

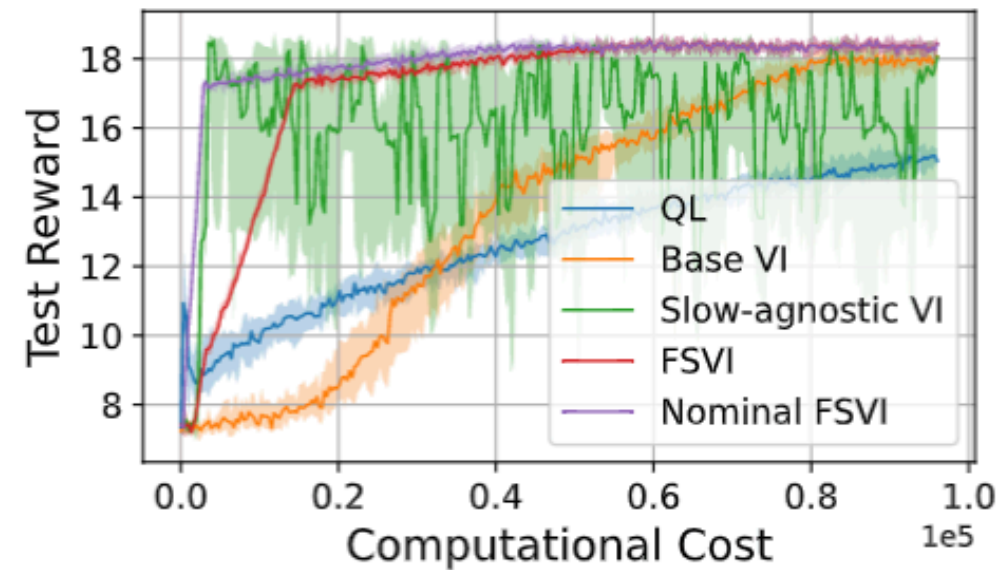
Baseline algorithms

- VI / AVI
- Slow-agnostic VI / AVI
 - Average over slow states during learning
 - Upon implementation, ignore slow state
- Q-learning (QL)
- Deep Q-networks (DQN)
- **Ours: FSVI / Nominal FSVI**
- **Ours: FSAVI / Nominal FSAVI**

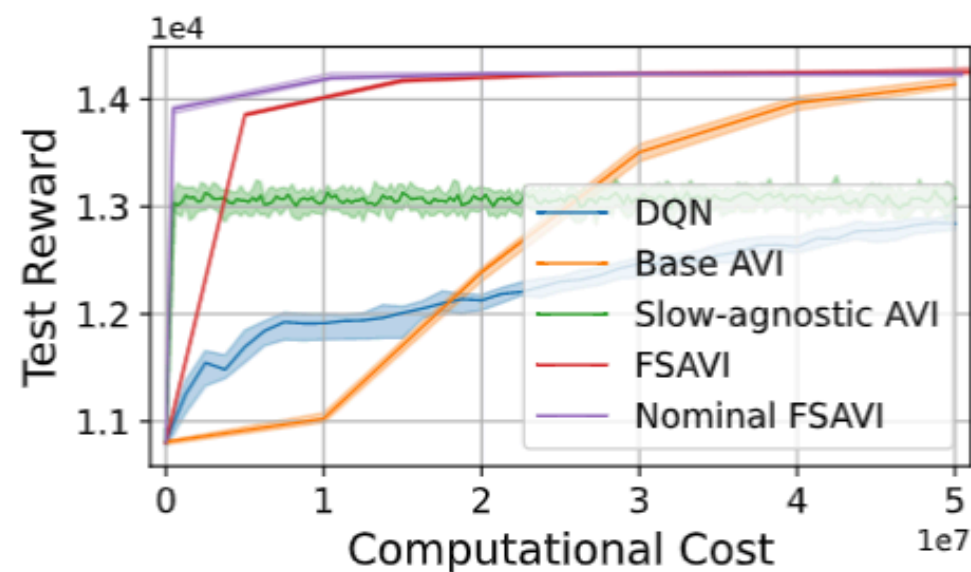
Overall performance comparison



(a) Multi-class service allocation



(b) Restless two-armed bandit



(c) Energy demand response

Questions

Please feel free to email me at danielrjiang@gmail.com for additional comments and discussion.

