



Online Learning for Residential Demand Response via Advanced Multi-Armed Bandits

Dr. Xin Chen

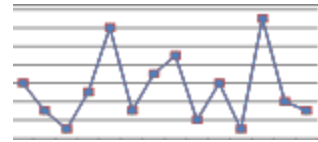
Assistant Professor
Electrical & Computer Engineering
Texas A&M University
Email: xin_chen@tamu.edu

09/03/2024

Residential Demand Response



Generation



=

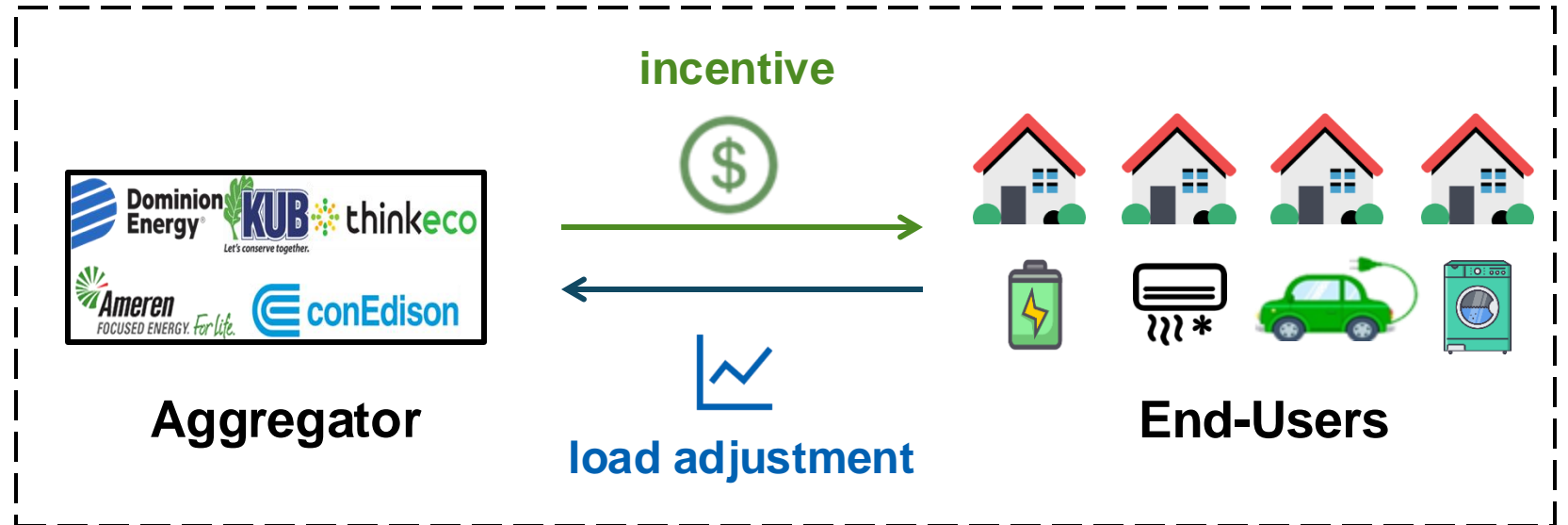
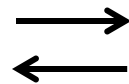
Demand

+ Loss

(adjustment)



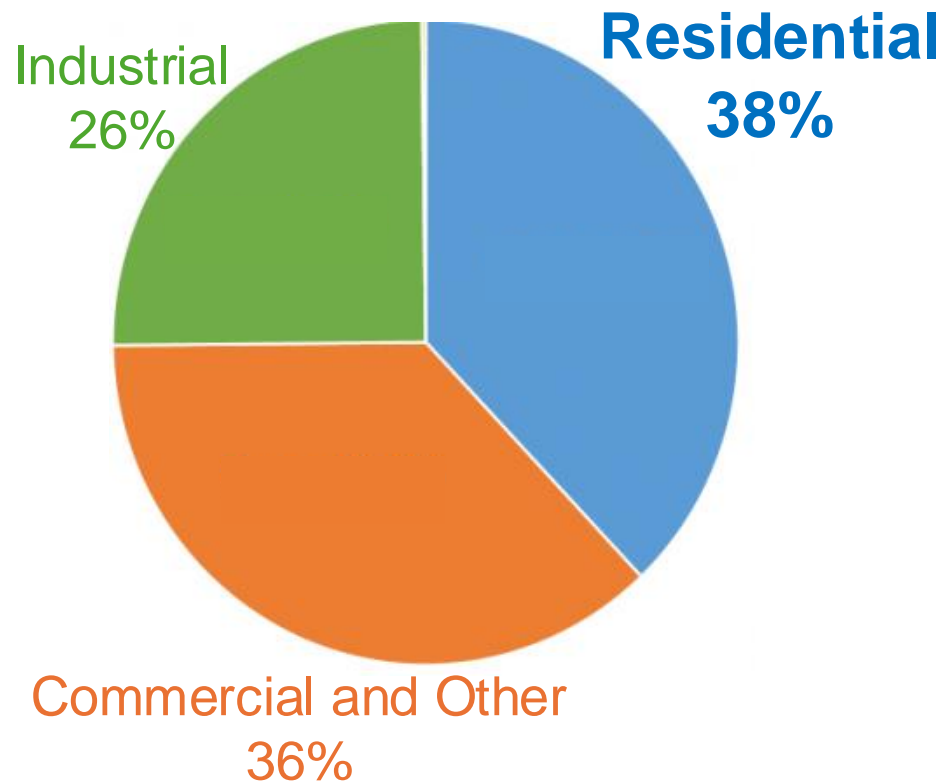
Power System



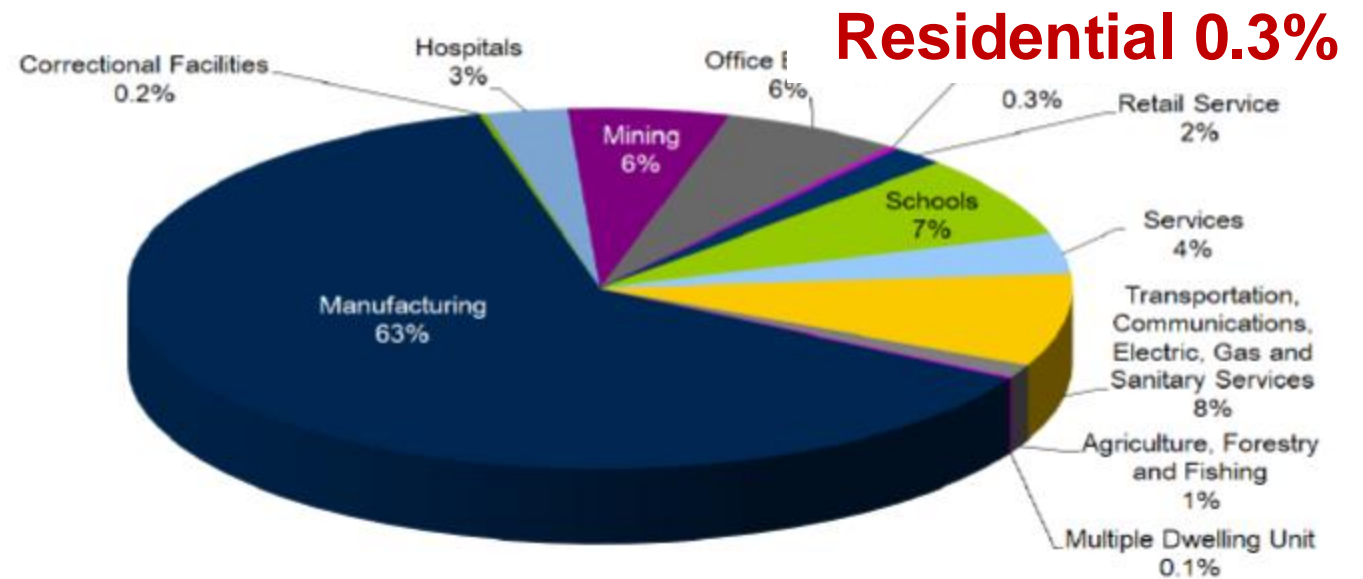
Why Study Residential Demand Response?

Because **residential** demand takes the **largest share**;

Huge potential but **underutilized**.



20/21 Confirmed DR registrations business segments



Source: 2022 U.S. Energy Information Administration (EIA), "Annual Electric Power Industry Report".

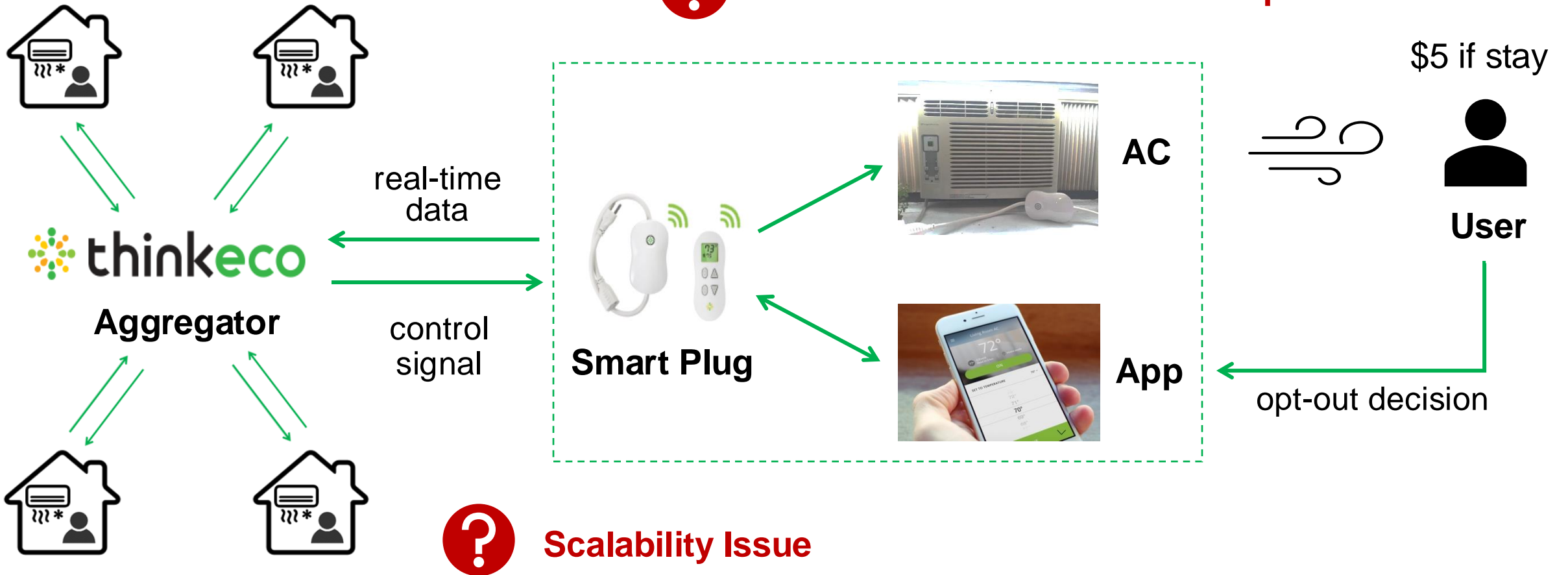
Source: PJM, "Demand response operations market activity report," 2022. Percent of load capacity (MWs)



Goal: Select Right Users to Signal

Due to budget constraints, need to select a subset of users (e.g., 1k) from the user pool (e.g., 10k)

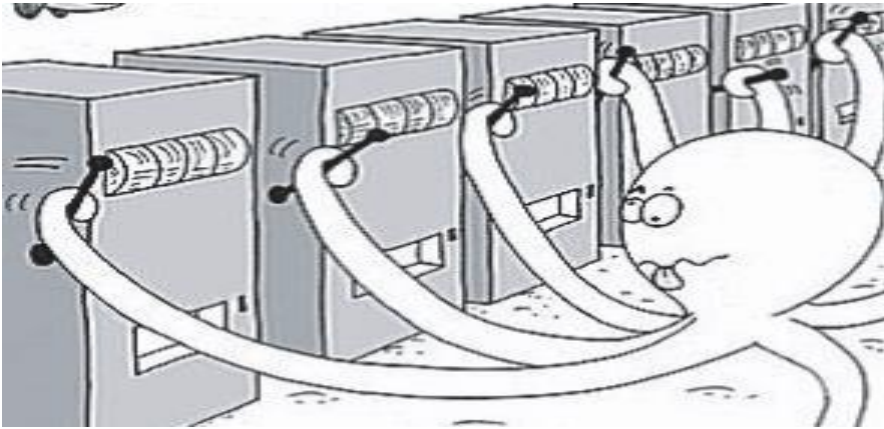
? Uncertain and Unknown User Opt-out Behaviors



Multi-Armed Bandits (MAB) Framework

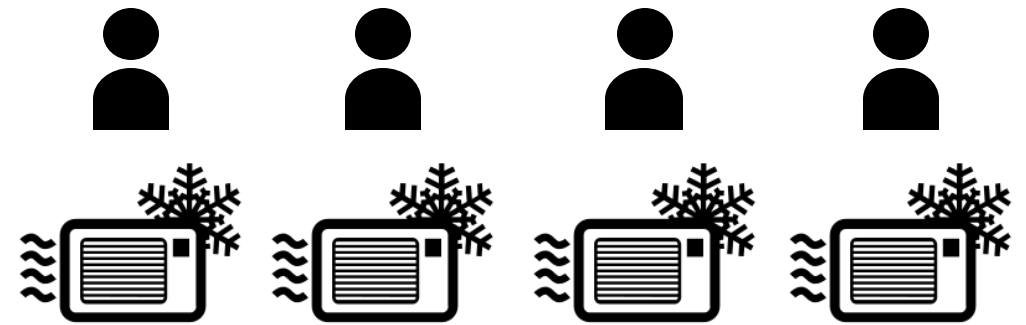
- ❖ Key Features: 1) Arms are different but independent; 2) Uncertain and unknown behaviors

Bandit Slot Machine



- Select one arm to maximize the profits;
- Observe the reward of the selected arm;
- Improve play strategies from feedback.

Demand Response



- Select a subset of users for DR;
- Observe responses from selected users;
- Learn users' behaviors from responses.

Application Examples



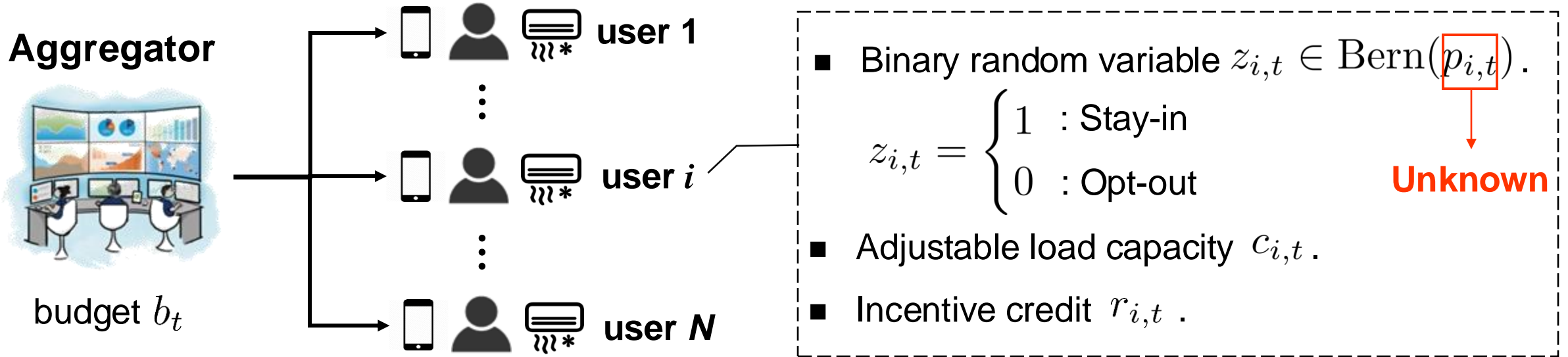
EV charging management



Residential load control

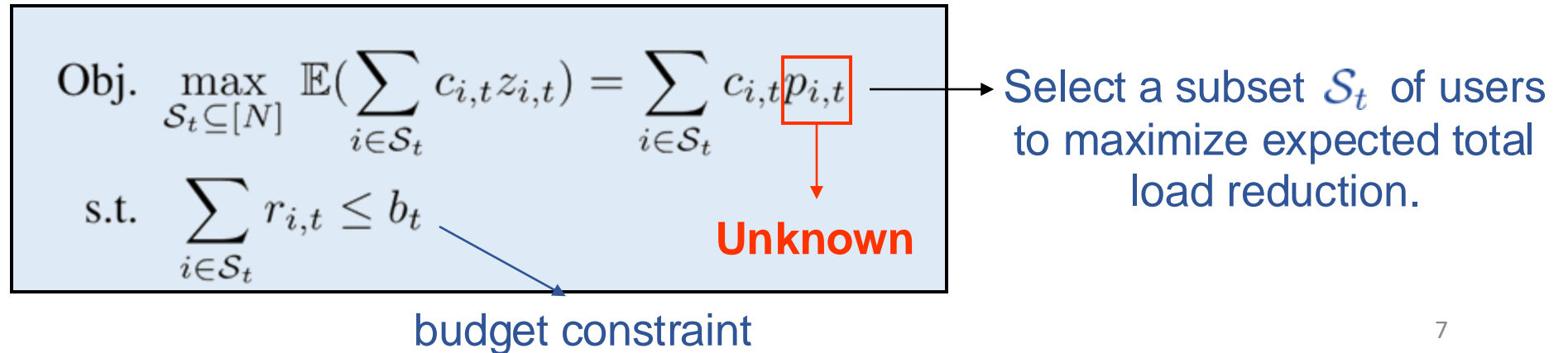
Problem Formulation

Consider a time horizon $[T] = \{1, 2, \dots, T\}$
 Each time $t \in [T]$ denotes a DR event.



↓

Optimal User Selection Model



Solution: Contextual Multi-Armed Bandits

- ❖ **Logistic regression** to predict $p_{i,t}$ for each user i at time t under *contextual influence*:

$$p_{i,t} = g(\boldsymbol{\theta}_i^\top \mathbf{x}_{i,t}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}_i^\top \mathbf{x}_{i,t})}$$

(Unknown) Individual Preference

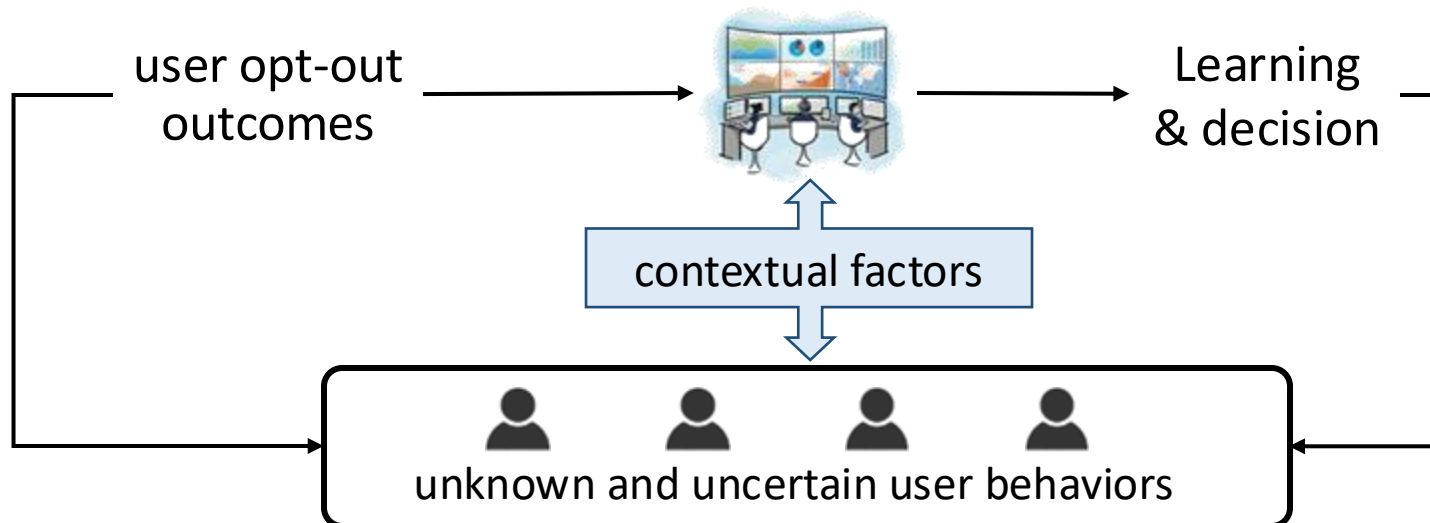
Contextual Factors

$$\boldsymbol{\theta}_i = (\theta_i^{(1)}, \theta_i^{(2)}, \dots, \theta_i^{(m)})$$

$$\mathbf{x}_{i,t} = (1, x_{i,t}^{(1)}, \dots, x_{i,t}^{(m)})$$

(electricity price, credit, weather, temperature, ...)

- ❖ **Online Learning and Human-In-the-Loop Decision:**



Thompson Sampling
to learn unknown $\boldsymbol{\theta}_i$ with
balance of **exploitation**
and **exploration**

Online Algorithm Based on Thompson Sampling.

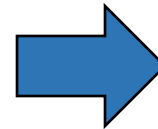
- Assume unknown θ_i be a random variable with Gaussian prior $\mathbb{P}_{\theta_i} = \mathcal{N}(\mu_i, \Sigma_i)$.
- In each demand response event t ,

Step 1: Sample $\hat{\theta}_i$ from its distribution \mathbb{P}_{θ_i} .

$$p_{i,t} = \frac{1}{1 + \exp\left(-\hat{\theta}_i^\top \mathbf{x}_{i,t}\right)}$$

Step 2: Select users by solving

$$\begin{aligned} \text{Obj. } & \max_{S_t \subseteq [N]} \mathbb{E}\left(\sum_{i \in S_t} c_{i,t} z_{i,t}\right) = \sum_{i \in S_t} c_{i,t} p_{i,t} \\ \text{s.t. } & \sum_{i \in S_t} r_{i,t} \leq b_t \end{aligned}$$



$$\begin{aligned} \text{Obj. } & \max_{\alpha_{i,t} \in \{0,1\}} \sum_{i=1}^N c_{i,t} p_{i,t} \alpha_{i,t} \\ \text{s.t. } & \sum_{i=1}^N r_{i,t} \alpha_{i,t} \leq b_t \end{aligned} \quad \text{Binary Optimization}$$

Step 3: Update posterior $\mathbb{P}_{\theta_i} \leftarrow \mathbb{P}_{\theta_i}(\cdot | \mathbf{x}_{i,t}, z_{i,t})$ with the observation $\mathbf{x}_{i,t}, z_{i,t}$.

variational Bayesian inference approach [3]

Regret Analysis

■ T-time Regret:
$$\text{Regret}(T, \boldsymbol{\theta}) = \sum_{t=1}^T \mathbb{E} [f_{\boldsymbol{\theta}}(\mathcal{S}_t^*, t) - f_{\boldsymbol{\theta}}(\mathcal{S}_t, t) \mid \boldsymbol{\theta}]$$

Optimal objective with true $\boldsymbol{\theta}$.

Objective using the proposed algorithm.

■ T-time Bayesian Regret:
$$\text{BayesRegret}(T) = \mathbb{E}_{\boldsymbol{\theta} \sim P_0} [\text{Regret}(T, \boldsymbol{\theta})]$$

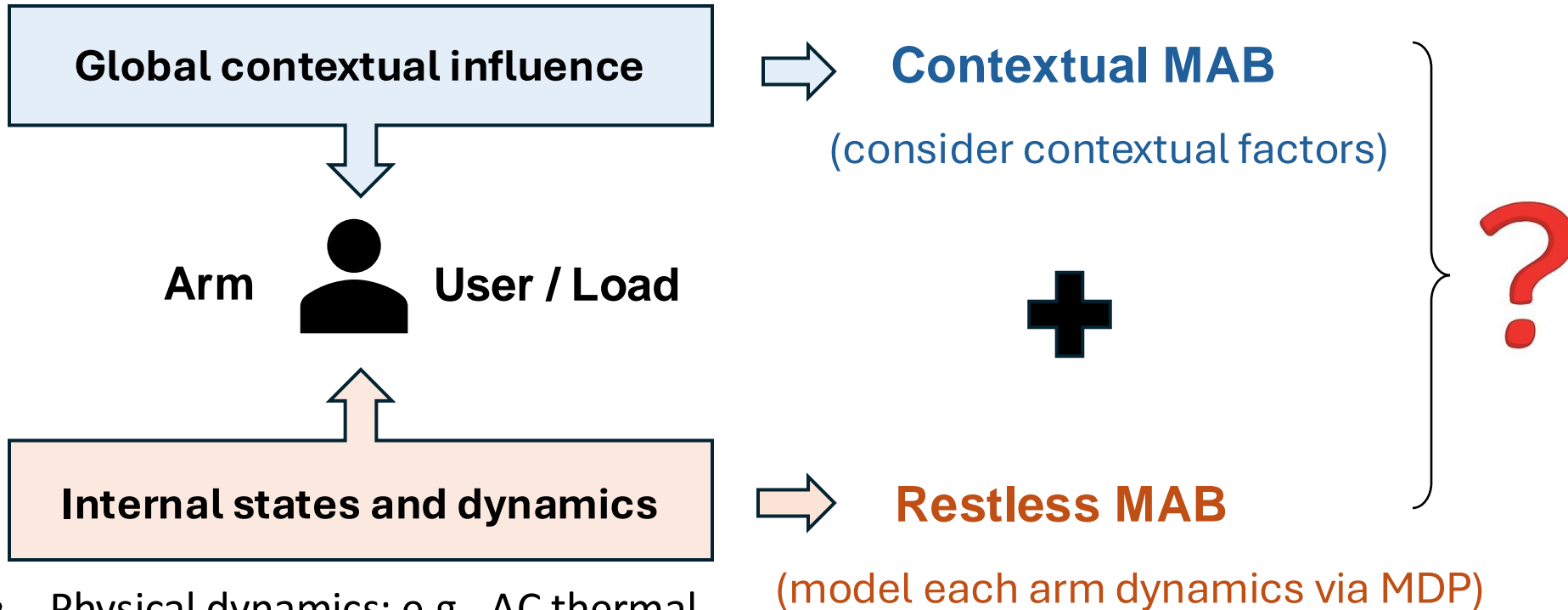
Theorem (informal): When T is sufficiently large, the Bayesian regret is

$$\text{BayesRegret}(T) \leq O\left(N^2 \gamma^d \sqrt{T \log T (d + \log T)}\right) \sim O(\log(T) \sqrt{T})$$

where $\gamma = \exp(2 \sup_{i \in [N]} \|\boldsymbol{\theta}_i\|_{\infty})$ and d is the dimension of $\boldsymbol{\theta}_i$. **sublinear**

Recent Extension: Contextual Restless Bandits

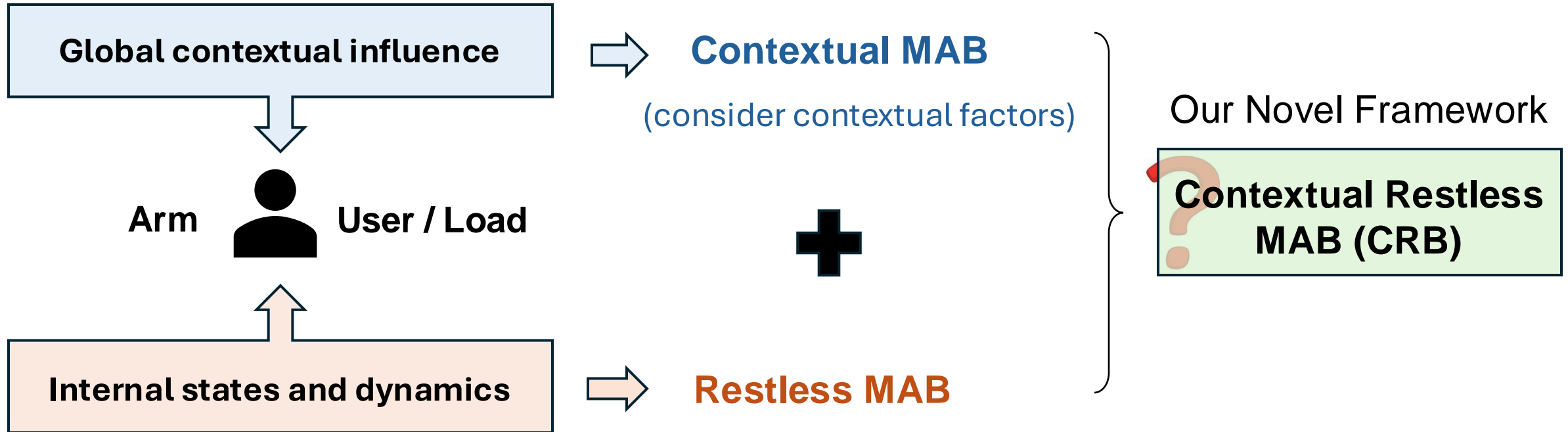
- temperature, weather, price, time ...



- Physical dynamics: e.g., AC thermal dynamics, EV charging SOC dynamics
- User “fatigue effect” ...

Recent Extension: Contextual Restless Bandits

- temperature, weather, price, time ...



- Physical dynamics: e.g., AC thermal dynamics, EV charging SOC dynamics
- User “fatigue effect” ...

- **X. Chen**, I. Hou, “Contextual Restless Multi-Armed Bandits with Application to Demand Response Decision-Making”, IEEE CDC, 2024.

CRB Problem Formulation

Consider N arms and an infinite time horizon $t = 0, 1, 2, \dots$

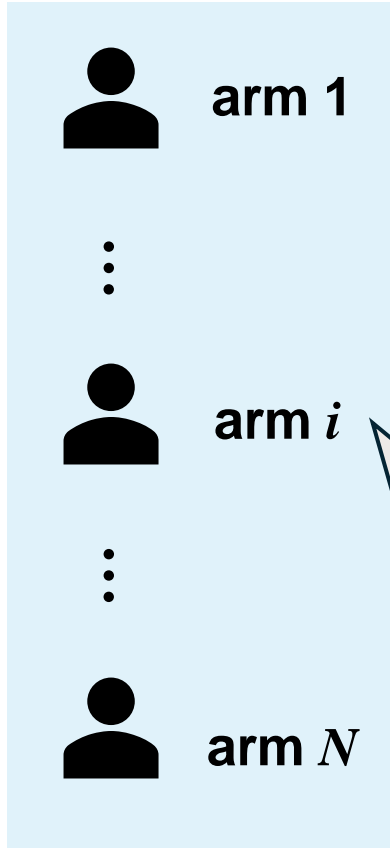
At each time t , observe the states $s_{i,t}$ of all arms and global context g_t , selects at most C_{g_t} arms

Decision-maker



find optimal selection policy $\pi(\mathbf{a} | \mathbf{s}, g)$

$$\begin{aligned}
 \text{(Primal)} : \quad & \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \beta^t r_{i,t} \right] \\
 \text{s.t.} \quad & \sum_{i=1}^N a_{i,t} \leq C_{g_t}, \quad \forall t \in \mathcal{T}
 \end{aligned}$$



Each arm i : **Context-augmented MDP**

- **Action:** $a_{i,t} \in \mathcal{A} := \{0, 1\}$
- **State:** $s_{i,t} \in \mathcal{S}$
- **Transition** (context-dependent): $P_i(s' | g, s, a) := \mathbb{P}(s_{i,t+1} = s' | g_t = g, s_{i,t} = s, a_{i,t} = a)$
- **Reward** (context-dependent): $r_{i,t} = R_i(g_t, s_{i,t}, a_{i,t})$

- **Global context:** $g_t \in \mathcal{G}$
follows positive-recurrent Markov Process $G(g' | g) := \mathbb{P}(g_{t+1} = g' | g_t = g)$

Solution Algorithm: Dual Decomposition

Very high-dimension MDP

$$\begin{aligned}
 \text{(Primal)} : \quad & \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \beta^t r_{i,t} \right] \\
 \text{s.t.} \quad & \sum_{i=1}^N a_{i,t} \leq C_{g_t}, \quad \forall t \in \mathcal{T}
 \end{aligned}$$

Relax budget constraint

expected per-global-context budget constraint

$$\begin{aligned}
 \text{(Relaxed)} : \quad & \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \beta^t r_{i,t} \right] \\
 \text{s.t.} \quad & \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \beta^t \mathbb{I}(g_t = g) \left(\sum_{i=1}^N a_{i,t} \right) \right] \\
 & \leq \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \beta^t \mathbb{I}(g_t = g) \right] C_g, \quad \forall g \in \mathcal{G}
 \end{aligned}$$

Challenge ?

Lagrange multiplier $\boldsymbol{\lambda} := (\lambda_g)_{g \in \mathcal{G}}$

Dual Decomposition Solution

- Gradient descent update of $\boldsymbol{\lambda}$
- Each arm i solves a **local** problem

$$\text{(Arm}_i(\boldsymbol{\lambda})) : \max_{\rho_i} \mathbb{E}_{\rho_i} \left[\sum_{t=0}^{\infty} \beta^t (r_{i,t} - \lambda_{g_t} a_{i,t}) \right]$$

$$\text{(Dual)} : \min_{\boldsymbol{\lambda} \geq \mathbf{0}} \left[\max_{\pi} L(\pi, \boldsymbol{\lambda}) \right].$$

$$L(\pi, \boldsymbol{\lambda}) = \sum_{i=1}^N \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \beta^t (r_{i,t} - \lambda_{g_t} a_{i,t}) \right] + \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \beta^t \lambda_{g_t} C_{g_t} \right].$$

Solution Algorithm: Dual Decomposition

- Gradient descent update of λ

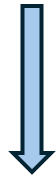
$$\lambda_g^{(k+1)} = \left[\lambda_g^{(k)} + \delta_k \mathbb{E}_{\pi^*(\lambda^{(k)})} \left[\sum_{t=0}^{\infty} \beta^t \mathbb{I}(g_t = g) \left(\sum_{i=1}^N a_{i,t} - C_g \right) \right] \right]^+$$

$\lambda^{(k)} := (\lambda_g^{(k)})_{g \in \mathcal{G}}$

Each arm i solves a **local unconstrained** problem

$$(\text{Arm}_i(\lambda)) : \max_{\rho_i} \mathbb{E}_{\rho_i} \left[\sum_{t=0}^{\infty} \beta^t (r_{i,t} - \lambda_{g_t} a_{i,t}) \right] \quad \text{the policy for arm } i: \rho_i(a|g, s, \lambda)$$

Solve Bellman
Optimality Equation



Interpretation of λ : global “price” paid to select an arm

✓ Optimal Q-function: $Q_{i,\lambda}^*(g, s, a) = R_i(g, s, a) - \lambda_g a + \beta \sum_{g' \in \mathcal{G}, s' \in \mathcal{S}} G(g'|g) P_i(s'|g, s, a) V_{i,\lambda}^*(g', s')$.

Optimal local policy:

$$\rho_i^*(\lambda) : a_{i,t} = \begin{cases} 1, & \text{if } Q_{i,\lambda}^*(g_t, s_{i,t}, 1) > Q_{i,\lambda}^*(g_t, s_{i,t}, 0) \\ 0, & \text{otherwise.} \end{cases}$$

Index Policy with **Known Models**

Algorithm 1 Index Policy Algorithm for Solving the CRB Problem with Known Arm Models.

- 1: **Initialization:** $\lambda^{(0)} := (\lambda_g^{(0)})_{g \in \mathcal{G}} \leftarrow \mathbf{0}$; $k \leftarrow 0$; $\epsilon > 0$.
 - 2: **repeat**
 - 3: Perform the following three steps for each arm $i \in [N]$ in parallel:
 - 1) Solve $\mathbf{LP}_i(\lambda^{(k)})$ (9) to obtain $V_{i,\lambda^{(k)}}^*(g, s)$;
 - 2) Compute the Q -function $Q_{i,\lambda^{(k)}}^*(g, s, a)$ by (10);
 - 3) Construct the optimal policy $\rho_i^*(\lambda^{(k)})$ by (11).
 - 4: Let $\pi^*(\lambda^{(k)}) := (\rho_i^*(\lambda^{(k)}))_{i \in [N]}$, and perform the update (7) to obtain $\lambda^{(k+1)}$, where the expectation is computed using (12)-(14). Let $k \leftarrow k + 1$.
 - 5: **until** the convergence $\|\lambda^{(k)} - \lambda^{(k-1)}\| \leq \epsilon$ is met.
 - 6: Let $\lambda^* \leftarrow \lambda^{(k)}$.
 - 7: **for** time $t = 0, 1, 2, \dots$ **do**
 - 8: Compute the index of each arm $i \in [N]$:

$$I_{i,t} \leftarrow Q_{i,\lambda^*}^*(g_t, s_{i,t}, 1) - Q_{i,\lambda^*}^*(g_t, s_{i,t}, 0). \quad (15)$$
 - 9: Sort all arms such that $I_{(1),t} \geq I_{(2),t} \geq \dots \geq I_{(N),t}$.
 - 10: Activate the top C_{g_t} arms $(1), (2), \dots, (C_{g_t})$.
 - 11: **end for**
-

Real-time application:

1. Compute index of each arm
2. Rank arms by index values
3. Select top C_{g_t} arms

Index of each arm ←

Solve Relaxed Problem to obtain λ^* via dual decomposition

Implement Index Policy (ensure budget constraint)

Online Learning with **Unknown** Models

Algorithm 2 Online Learning Algorithm for Solving the CRB Problem with Unknown Arm Models.

1: **Initialization:** Time window T ; $\epsilon_n > 0$; initial transition kernel $P_i^0(\cdot)$ of each arm $i \in [N]$.

2: **for** epoch $n = 0, 1, 2, \dots$ **do**

3: Based on the up-to-date transition kernel model $P_i^n(\cdot)$ of each arm $i \in [N]$, follow Steps **2-6** in Algorithm **1** to compute the optimal λ_n^* .

4: **for** time $t = nT, nT+1, \dots, nT+T-1$ **do**

5: With probability of $1 - \epsilon_n$,

- Compute the index $I_{i,t}$ (**16**) of each arm $i \in [N]$;
- Sort all arms such that $I_{(1),t} \geq I_{(2),t} \geq \dots \geq I_{(N),t}$ and activate the top C_{g_t} arms.

With probability of ϵ_n , randomly activate C_{g_t} arms.

6: **end for**

7: For each arm $i \in [N]$, based on the observed state transitions, update its transition kernel model by:

$$P_i^{n+1}(s_{i,t+1} = s' | g_t = g, s_{i,t} = s, a_{i,t} = a) = \frac{M_{s',g,s,a}^i}{M_{g,s,a}^i},$$

where $M_{g,s,a}^i$ and $M_{s',g,s,a}^i$ are arm i 's cumulative historical counts of the context-state-action tuple (g, s, a) and the state transition $(g, s, a) \rightarrow s'$.

8: **end for**

Solve Relaxed Problem

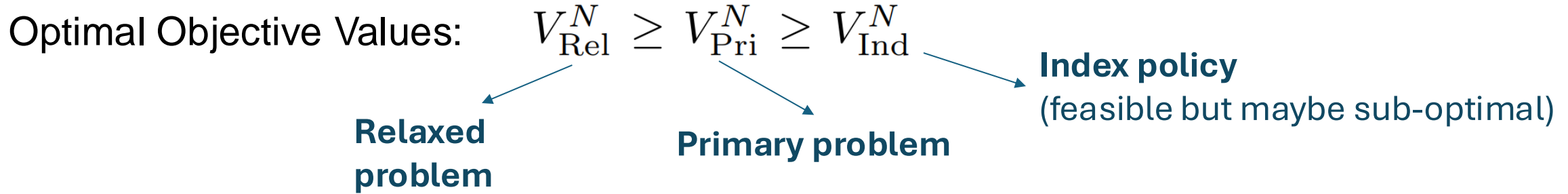
Implement Index Policy
with ϵ -exploration

Better online learning algorithms:
Thompson Sampling, UCB, Q-learning

**Update estimation of
transition kernel models**

Function approximation for Scalability

Theoretical Analysis: Asymptotical Optimality



Theorem 1. *Suppose that the initial global context g_0 is chosen uniformly at random from \mathcal{G} and the initial state $s_{i,0}$ of each arm $i \in [N]$ is chosen independently with the distribution $\mathbb{P}(s_{i,0} = s | g_0 = g) = m_g^*(s)$, then, under Assumption 1, we have*

$$\underline{V_{\text{Rel}}^N \geq V_{\text{Pri}}^N \geq V_{\text{Rel}}^N - \mathcal{O}(\sqrt{N})}. \quad (20)$$

as $N \rightarrow \infty$

$$V_{\text{Rel}}^N / N \longrightarrow V_{\text{Pri}}^N / N$$

Numerical Simulations

Residential Demand Response:

- Number of users (arms) $N = 500$

- Discrete global context

$$g \in \mathcal{G} = \{1, 2, \dots, 6\}$$

- Selection budget $C = 100$

- State: $s_{i,t} := (z_{i,t}, x_{i,t})$

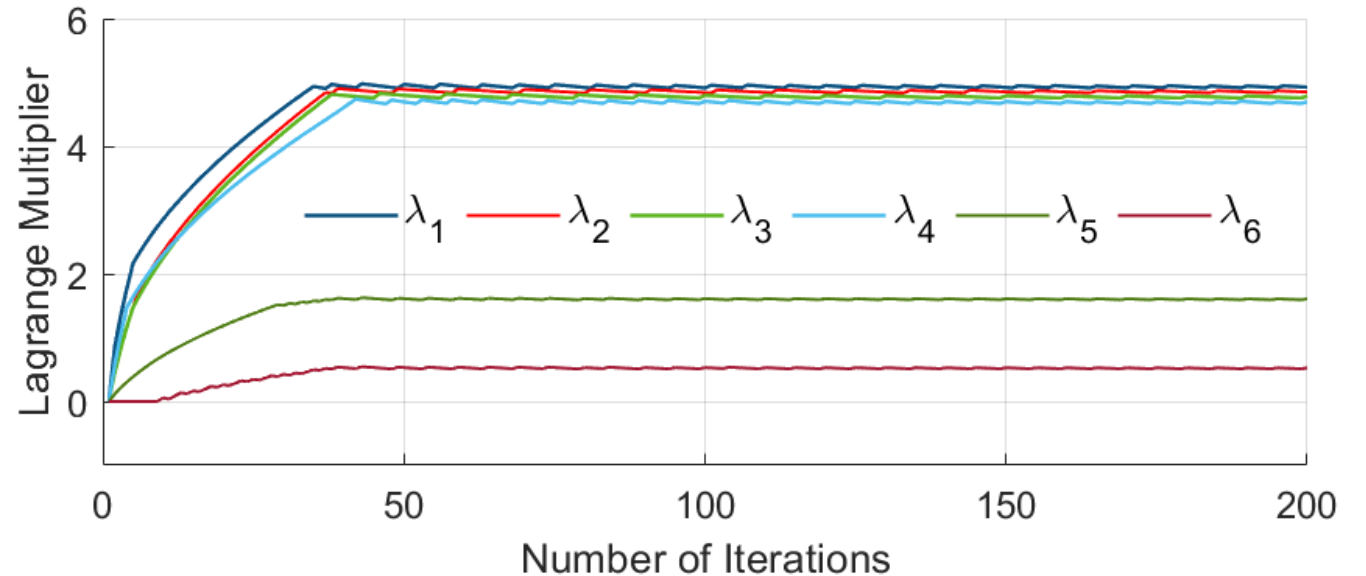
$$x_{i,t} \in \{1, 2, \dots, 4\}$$

- Reward: $R_i(\cdot) = \frac{a_{i,t} z_{i,t} l_i}{(g_t - x_{i,t})^2 + 1}$

- Discount factor $\beta = 0.97$

finite time horizon $T = 300$

Fig 1. Convergence of Dual Decomposition Alg.



Numerical Simulations

Fig 2. Asymptotical optimality of index policy.

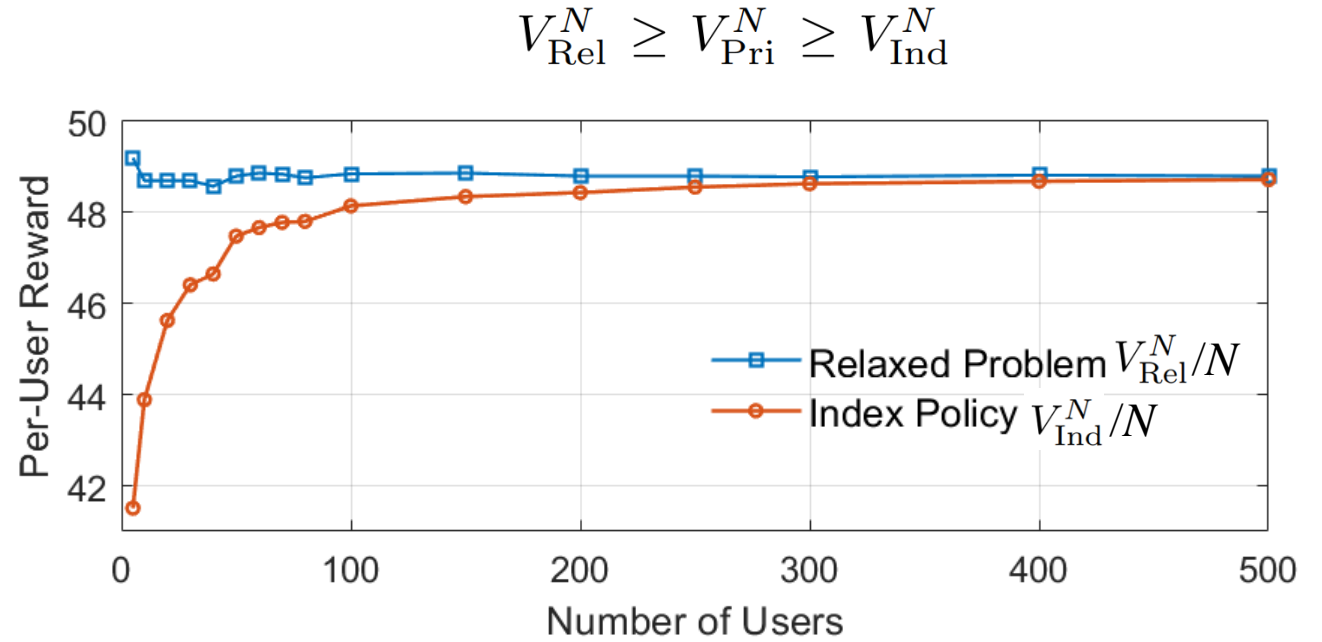
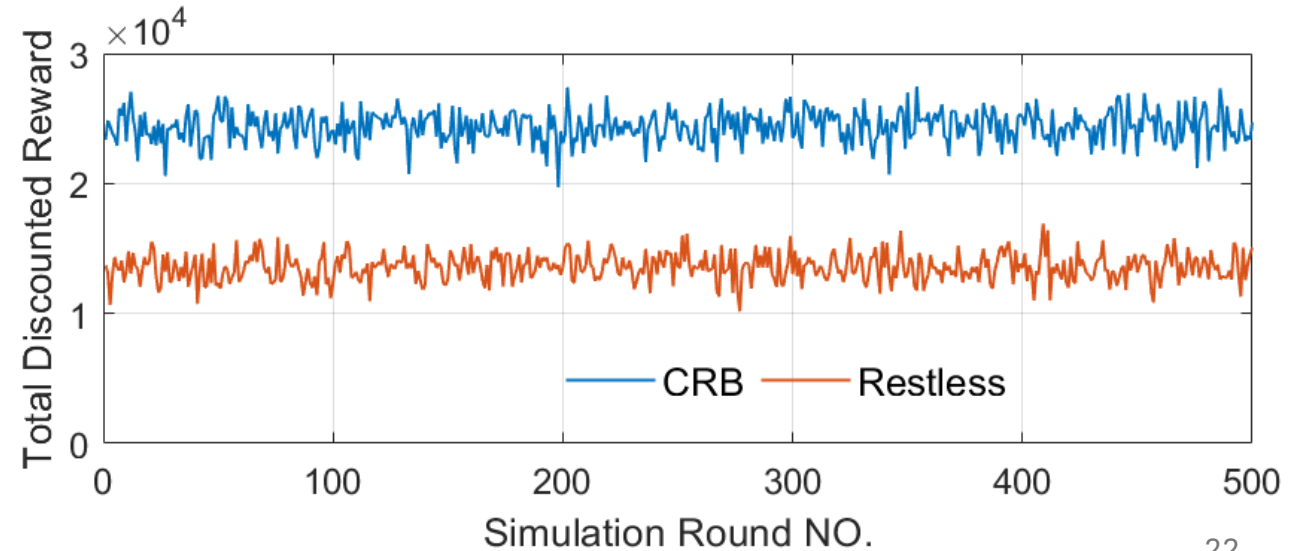
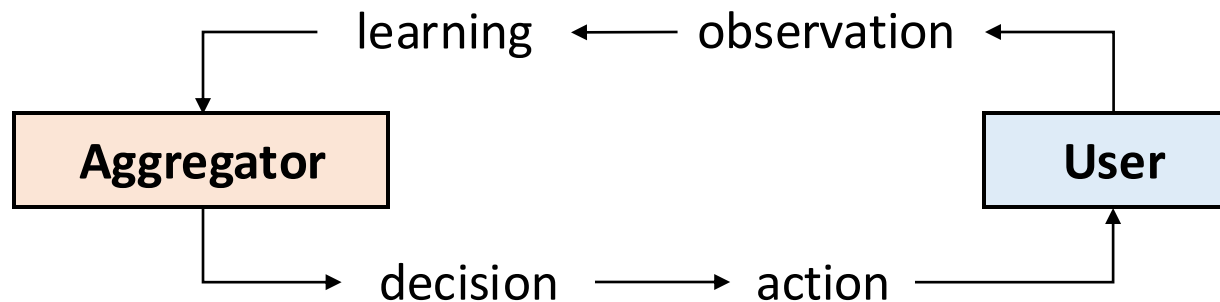


Fig 3. Comparison with Restless bandits.



Online learning for human-in-the-loop DR decision

Key Takeaway



- ❖ The **multi-armed bandits (MAB)** method is useful for large-scale DR online decision-making.
- ❖ A novel MAB framework, **Contextual Restless Bandits (CRB)**, models both the **dynamic state transitions of each arm** and the influence of **external global environmental context**.
- ❖ A scalable **index policy** algorithm based on dual decomposition is proposed to solve CRB.
- ❖ Simulation results demonstrate the **asymptotic optimality** and enhanced **modeling capability**.

References

- [1] A. Slivkins et al., “Introduction to multi-armed bandits,” *Foundations and Trends in Machine Learning*, vol. 12, no. 1-2, pp. 1–286, 2019
- [2] L. Zhou, “A survey on contextual multi-armed bandits,” arXiv preprint, arXiv:1508.03326, 2015.
- [3] P. Whittle, “Restless bandits: Activity allocation in a changing world,” *Journal of Applied Probability*, vol. 25, no. A, pp. 287–298, 1988.
- [4] R. R. Weber and G. Weiss, “On an index policy for restless bandits,” *Journal of Applied Probability*, vol. 27, no. 3, pp. 637–648, 1990.
- [5] X. Chen, I. Hou, “Contextual Restless Multi-Armed Bandits with Application to Demand Response Decision-Making“, arXiv:2403.15640, 2024.
- [6] X. Chen, Y. Li, J. Shimada, and N. Li, “Online learning and distributed control for residential demand response,” *IEEE Transactions on Smart Grid*, vol. 12, no. 6, pp. 4843–4853, 2021.
- [7] X. Chen, Y. Nie, and N. Li, “Online residential demand response via contextual multi-armed bandits,” *IEEE Control Systems Letters*, vol. 5, no. 2, pp. 433–438, 2020.
- [8] B. Liang, L. Xu, A. Taneja, M. Tambe, and L. Janson, “A Bayesian approach to online learning for contextual restless bandits with applications to public health,” arXiv preprint arXiv:2402.04933, 2024
- [9] J. A. Taylor and J. L. Mathieu, “Index policies for demand response,” *IEEE Transactions on Power Systems*, vol. 29, no. 3, pp. 1287–1295, 2013.
- [10] N. Gast, B. Gaujal, and C. Yan, “Exponential convergence rate for the asymptotic optimality of whittle index policy,” arXiv preprint arXiv:2012.09064, 2020.