



HARVARD

**School of Engineering
and Applied Sciences**

CDC 2020

Online Residential Demand Response via Contextual Multi-Armed Bandit

Xin Chen, Yutong Nie, Na Li

**School of Engineering and Applied Sciences, Harvard University
chen_xin@g.harvard.edu**

Power Systems

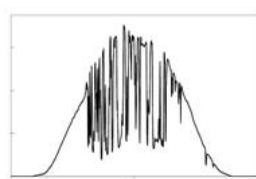
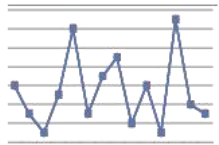
Generation

uncertainty



intermittency

Power Plants



Large-scale renewable generation

Transmission System

=

Load

growing load peak

State of emergency declared as California faces historic heat, possible power outages

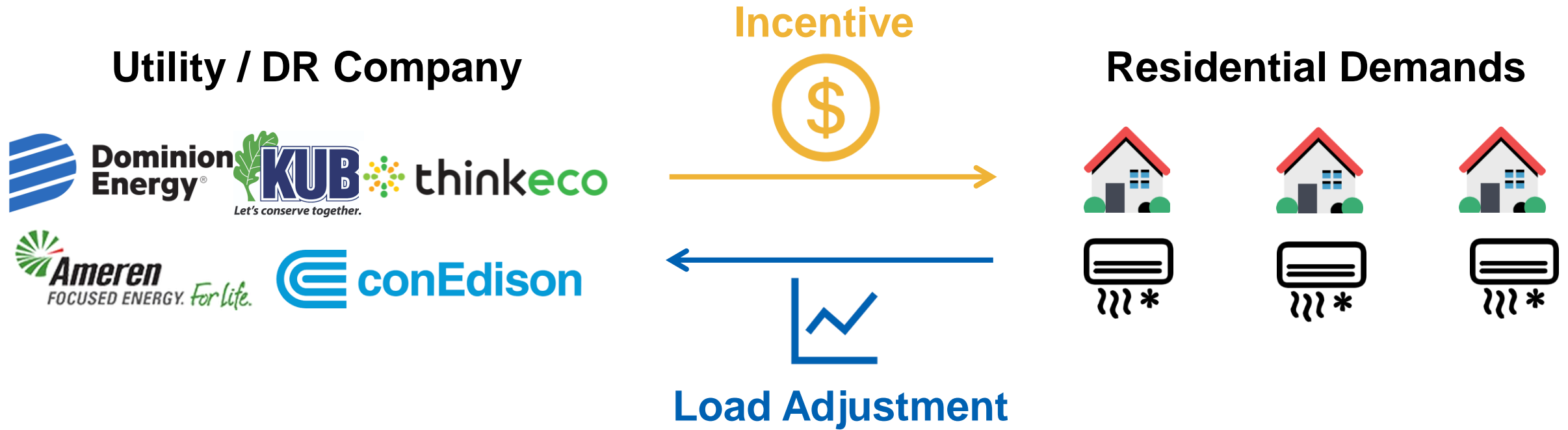
Los Angeles Times

Heat Wave Roasts Southern California With Record of 121 Degrees

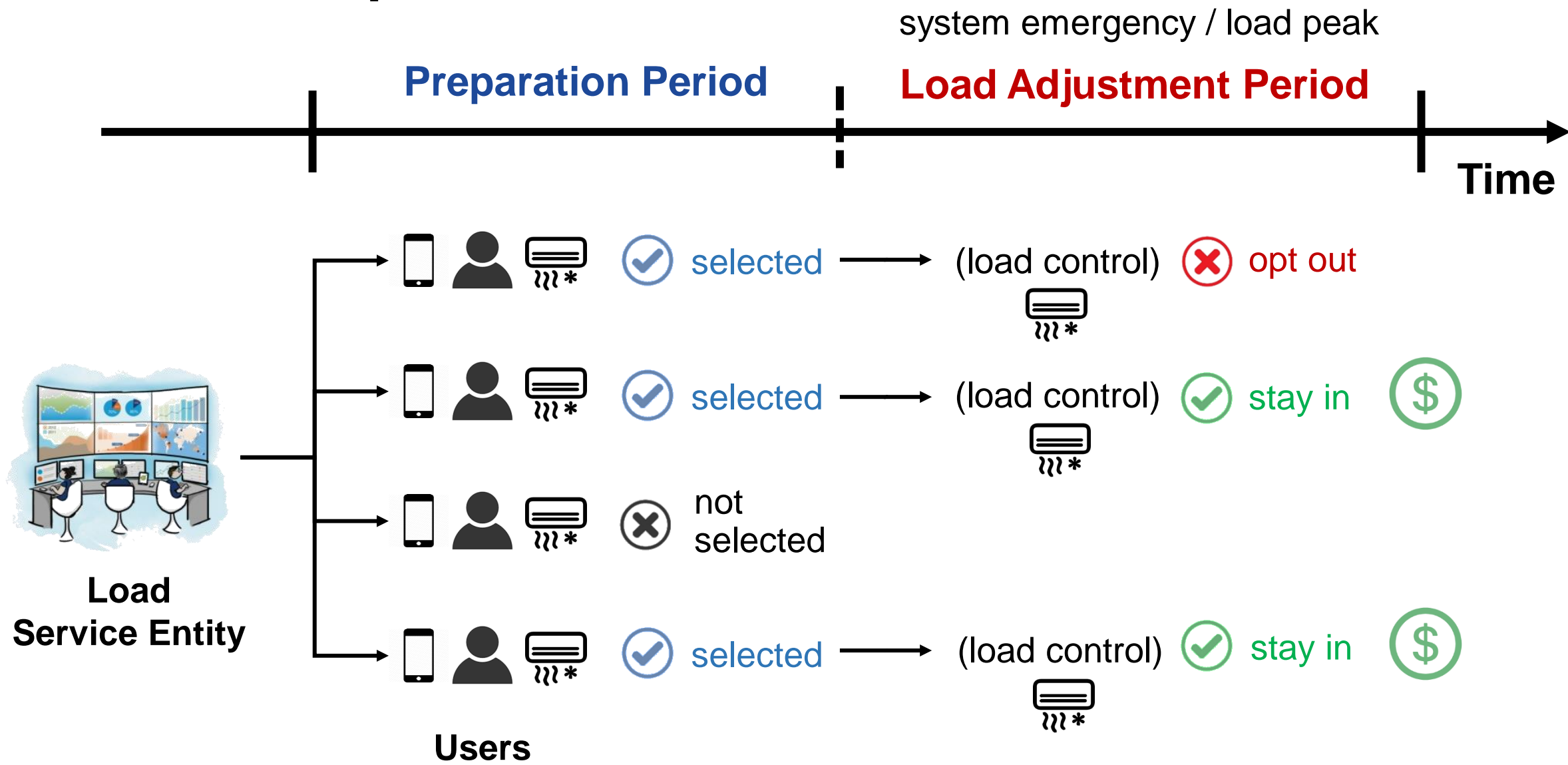
BRIEF

ERCOT calls 2 energy emergencies in one week, 3rd in 5 years

Residential Demand Response



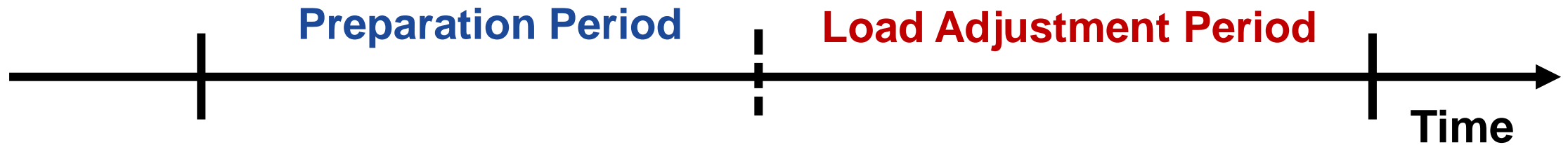
Demand Response Event



[1] X. Chen, Y. Li, J. Shimada, and N. Li, “Online Learning and Distributed Control for Residential Demand Response”, arXiv:2010.05153, 2020.

Q2- How to optimally control load in real-time?

A2: A follow-up paper [1].



Q1- How to select right users for DR?

A1: This talk.

Challenge: Uncertain and unknown user opt-out behaviors.

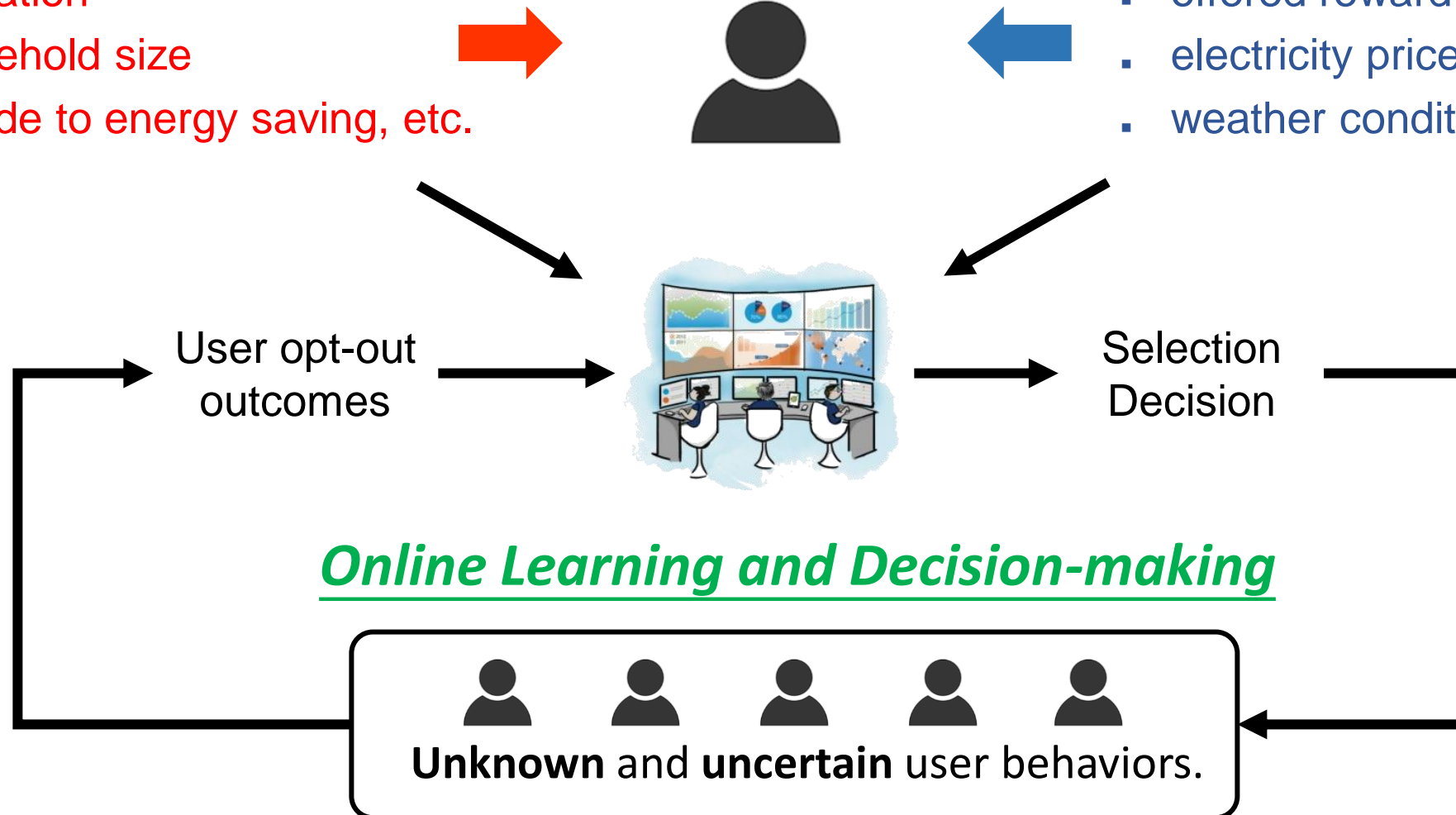
Individual Preference

- age
- education
- household size
- attitude to energy saving, etc.

Opt Out?

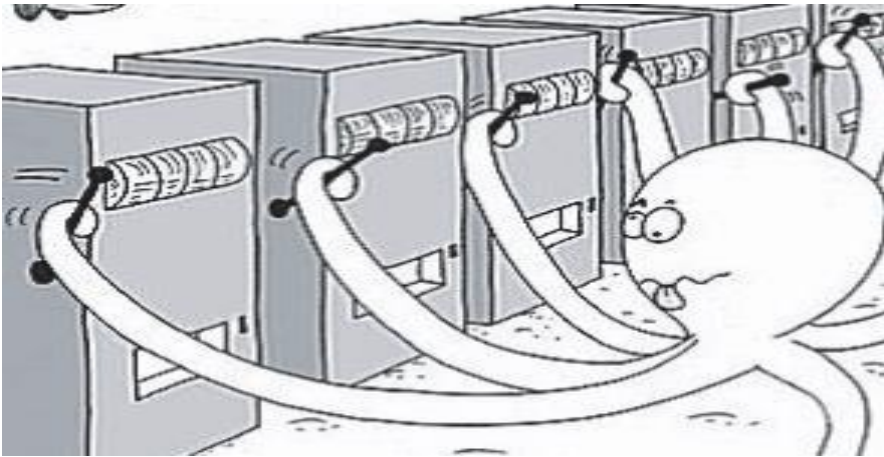
Environmental Factors

- indoor temperature
- offered reward
- electricity price
- weather conditions, etc.



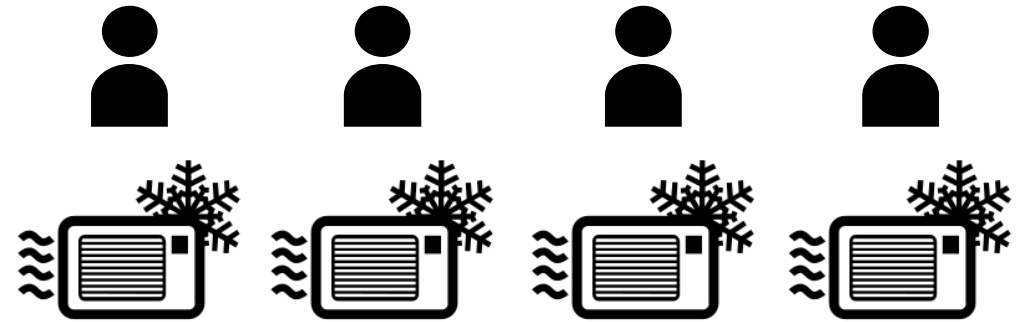
Contextual Multi-armed bandit (CMAB)

Slot Machine



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

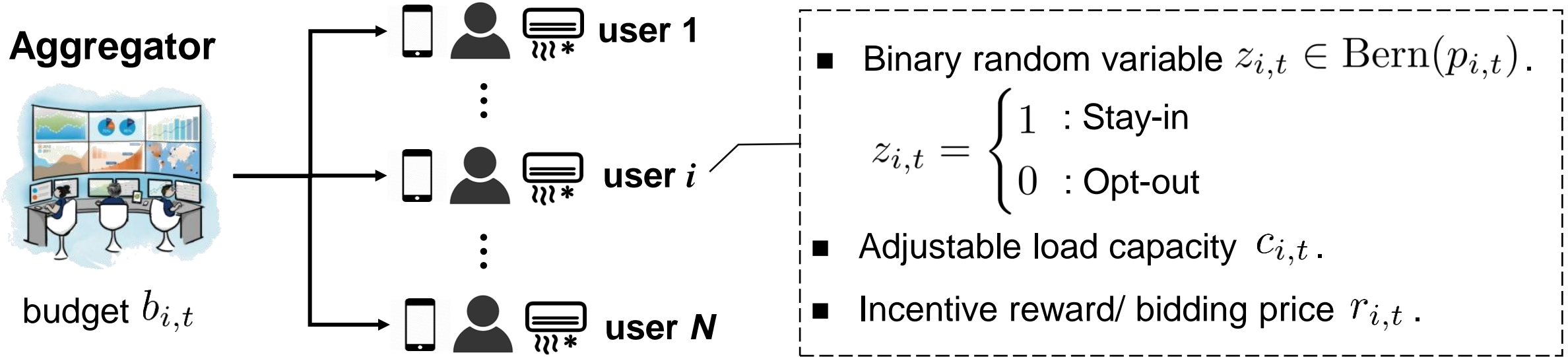
Demand Response



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

Problem Formulation

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



**Optimal User
Selection Model**

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

Unknown

User Behavior Learning with Contexts

- **Logistic model** to predict $p_{i,t} := p_i(t)$ for user i :

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

Personal Weight

$$\boldsymbol{\theta}_i = (\theta_{i,1}, \theta_{i,2}, \dots, \theta_{i,d})$$

individual preference response to context

Context at time t

$$\mathbf{x}_i(t) = (1, x_{i,2}, \dots, x_{i,d})(t)$$

environmental factors



Goal : to learn parameter $\boldsymbol{\theta}_i$ for each user.

- **Thompson sampling** is used to balance *exploration* and *exploitation*.

Thompson Sampling (Bayesian learning)

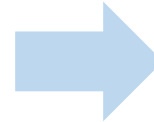
Assume unknown θ_i be a random variable with Gaussian prior $\mathbb{P}_{\theta_i} = \mathcal{N}(\mu_i, \Sigma_i)$.

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

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

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 [2]

Regret Analysis

- (Expected) Reward Function: $f_{\boldsymbol{\theta}}(\mathcal{S}_t, t) = \mathbb{E}(\sum_{i \in \mathcal{S}_t} c_{i,t} z_{i,t}) = \sum_{i \in \mathcal{S}_t} \frac{c_{i,t}}{1 + \exp(-\mathbf{x}_{i,t}^\top \boldsymbol{\theta}_i)}$
- T-time Regret: $\text{Regret}(T, \boldsymbol{\theta}) = \sum_{t=1}^T \mathbb{E}[f_{\boldsymbol{\theta}}(\mathcal{S}_t^*, t) - f_{\boldsymbol{\theta}}(\mathcal{S}_t, t) \mid \boldsymbol{\theta}]$
- T-time Bayesian Regret: $\text{BayesRegret}(T) = \mathbb{E}_{\boldsymbol{\theta} \sim P_0} [\text{Regret}(T, \boldsymbol{\theta})]$

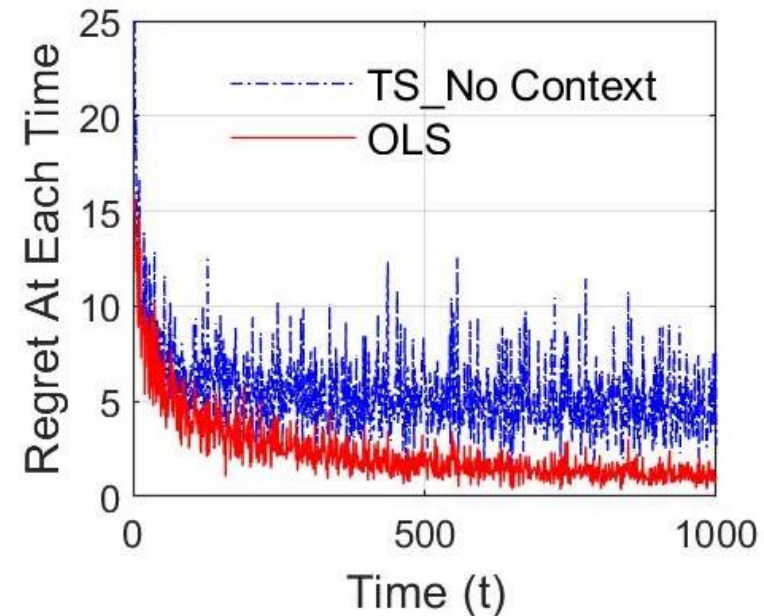
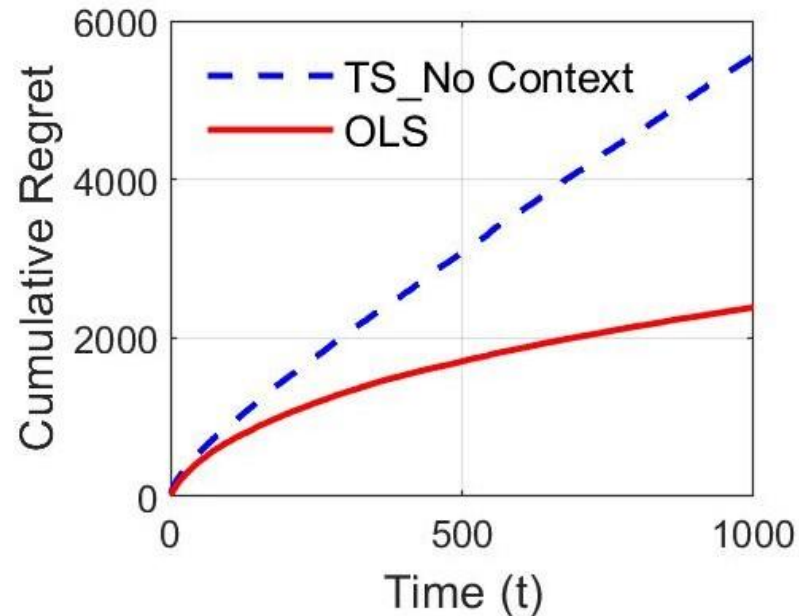
Theorem (informal): When T is sufficiently large, the Bayesian regret is

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

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

Simulation:

- $N = 1000$ users; $m = 9$ environmental factors.
- A Gaussian prior distribution $\mathcal{N}(\boldsymbol{\theta}_i^* + 0.3\mathbf{u}_i, 0.09\mathbf{I})$ for each user.
- Regret comparison between the proposed algorithm (OLS) and Thompson sampling without context information.



Thank you!