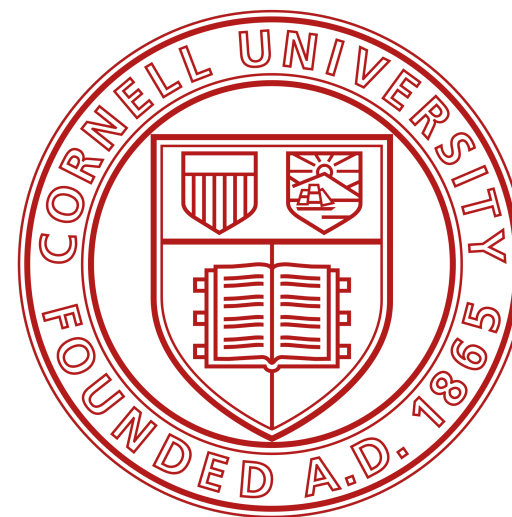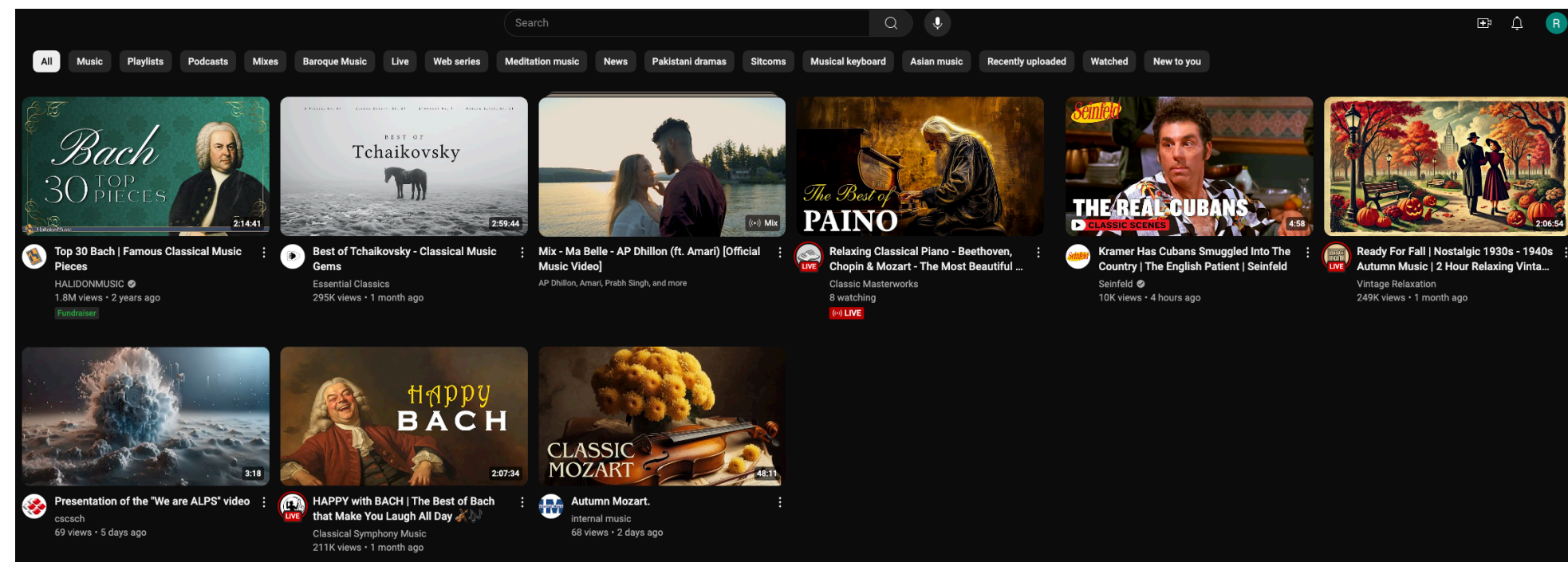# MultiScale Contextual Bandits for Long Term Objectives

**Richa Rastogi, Yuta Saito, Thorsten Joachims**
**Cornell University**

# Motivation

- In many interactive AI systems, (recommender, conversational systems), there is abundant short term feedback (e.g., clicks, generated response quality)
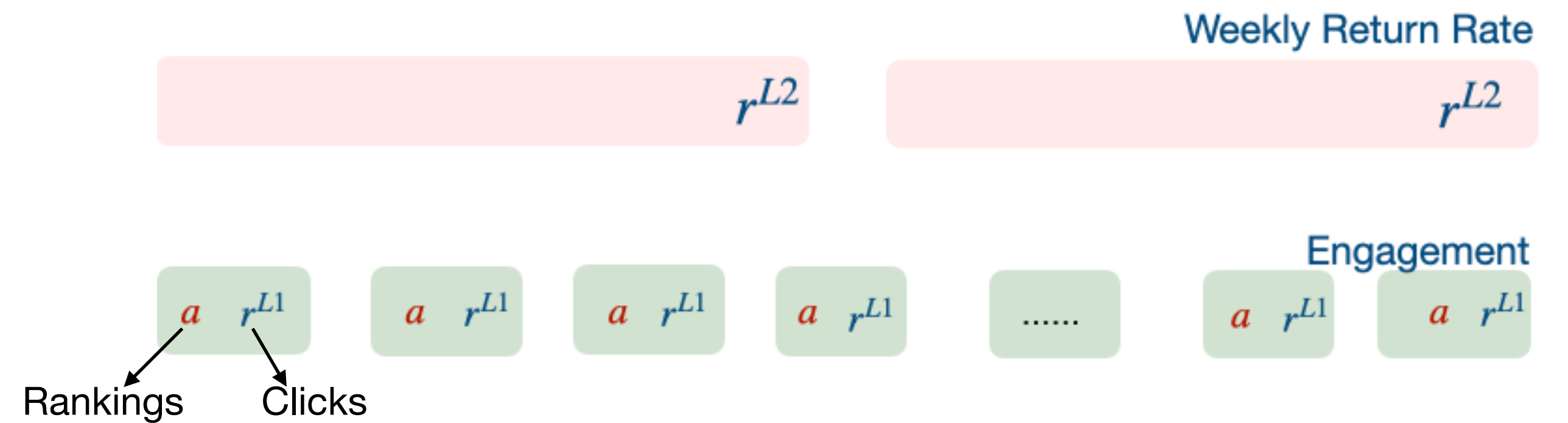


User Retention



Beneficial dialogue outcomes

- Prior work shows that optimizing for short term feedback does not necessarily achieve the desired long term objective (e.g., clickbait feeds do not lead to user retention)

# Motivation

A key problem — long-term feedback is at a different timescale than the short-term interventions
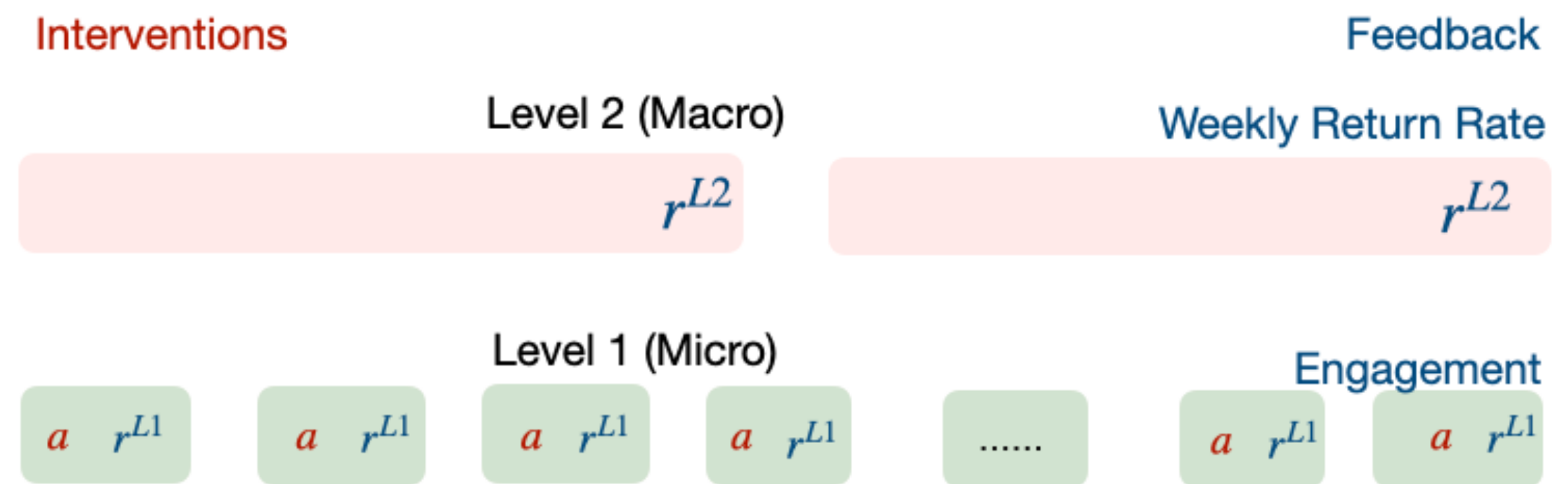
Weekly Return Rate

$r^{L2}$ $r^{L2}$

Engagement

$a$ $r^{L1}$ $a$ $r^{L1}$ $a$ $r^{L1}$ $a$ $r^{L1}$ ...... $a$ $r^{L1}$ $a$ $r^{L1}$

Rankings    Clicks

We address it by contextually reconciling this disconnect in timescales

# MultiScale Policy Framework
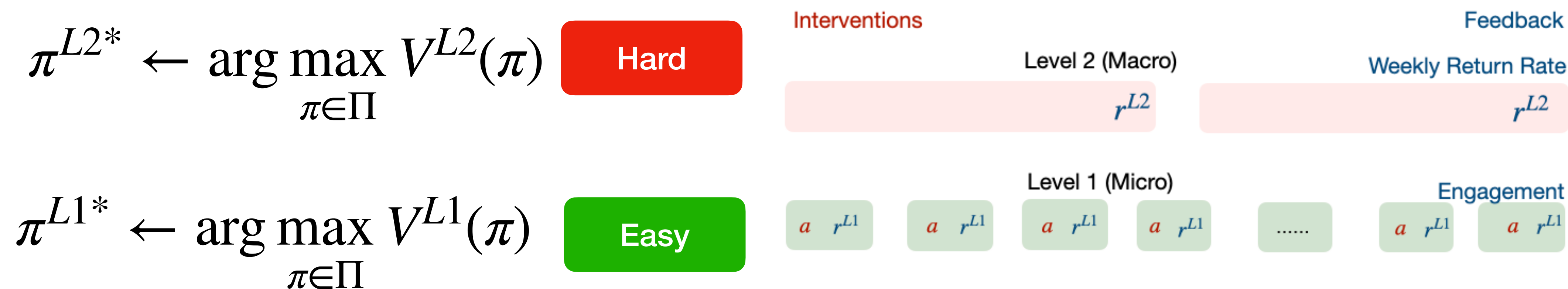
Consider two levels

- A micro level that operates at faster timescale, e.g., clicks, response quality

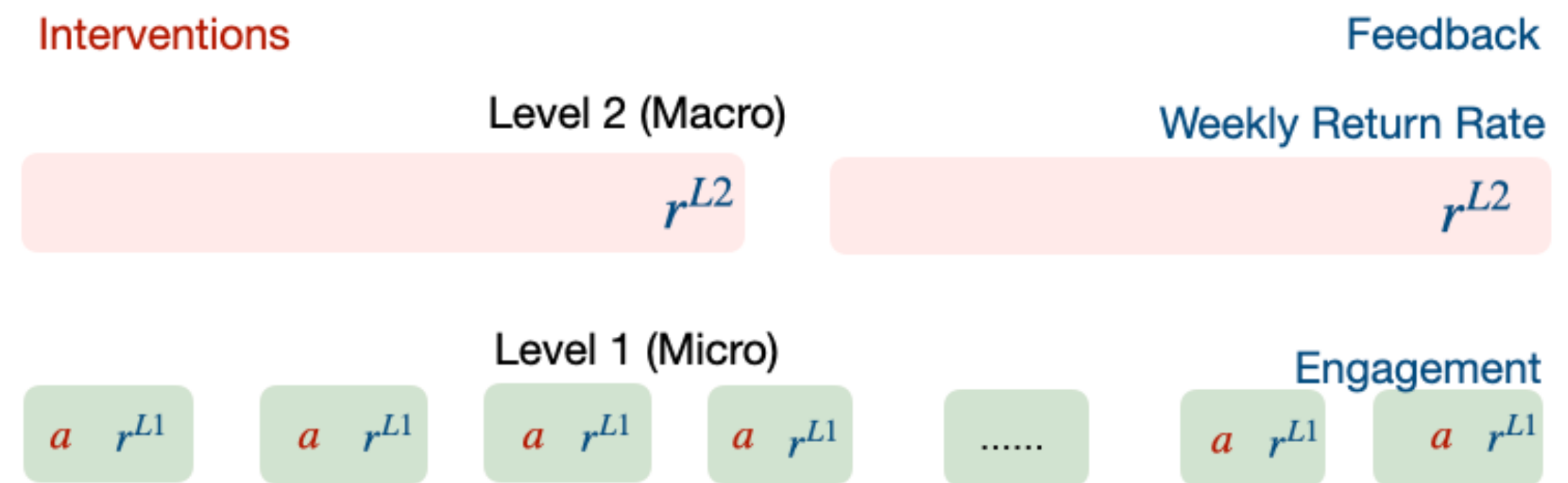- A macro level that operates at slower timescale, e.g., user retention

Interventions

Feedback

Level 2 (Macro)

Weekly Return Rate

$r^{L2}$   $r^{L2}$

Level 1 (Micro)

Engagement

$a$ $r^{L1}$   $a$ $r^{L1}$   $a$ $r^{L1}$   $a$ $r^{L1}$   ......   $a$ $r^{L1}$   $a$ $r^{L1}$

# MultiScale Policy Framework

$$\pi^{L2*} \leftarrow \arg\max_{\pi \in \Pi} V^{L2}(\pi)$$

Hard

$$\pi^{L1*} \leftarrow \arg\max_{\pi \in \Pi} V^{L1}(\pi)$$

Easy

Interventions

Level 2 (Macro)

Feedback

Weekly Return Rate

$r^{L2}$

$r^{L2}$

Level 1 (Micro)

Engagement

| $a$ | $r^{L1}$ | $a$ | $r^{L1}$ | $a$ | $r^{L1}$ | $a$ | $r^{L1}$ | ...... | $a$ | $r^{L1}$ | $a$ | $r^{L1}$ |

Even though $V^{L2}(\pi^{L1*}) < V^{L2}(\pi^{L2*})$, $\pi^{L1*}$ is typically much better than a random policy from $\Pi$

# MultiScale Policy Framework

Interventions

Level 2 (Macro)

Feedback

Weekly Return Rate

$r^{L2}$

$r^{L2}$

Level 1 (Micro)

Engagement

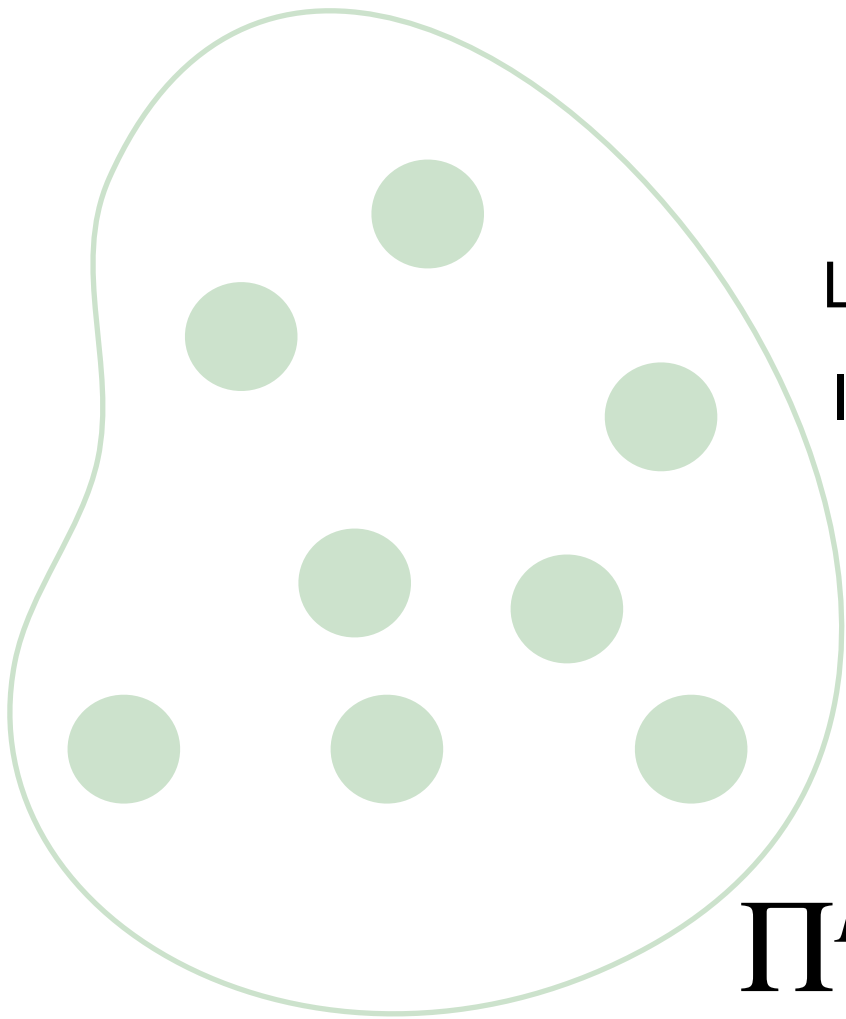| $a$ $r^{L1}$ | $a$ $r^{L1}$ | $a$ $r^{L1}$ | $a$ $r^{L1}$ | ...... | $a$ $r^{L1}$ | $a$ $r^{L1}$ |

Can we exploit feedback at the micro level to learn the long term optimal policy?
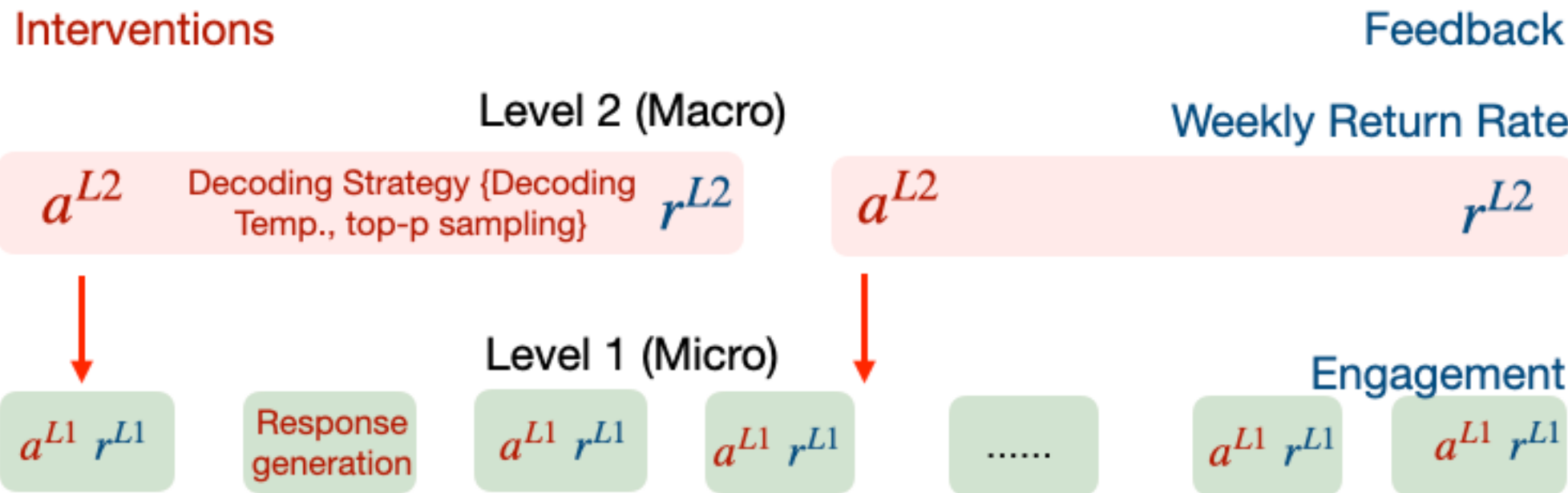
# MultiScale Policies

Factorization of policies

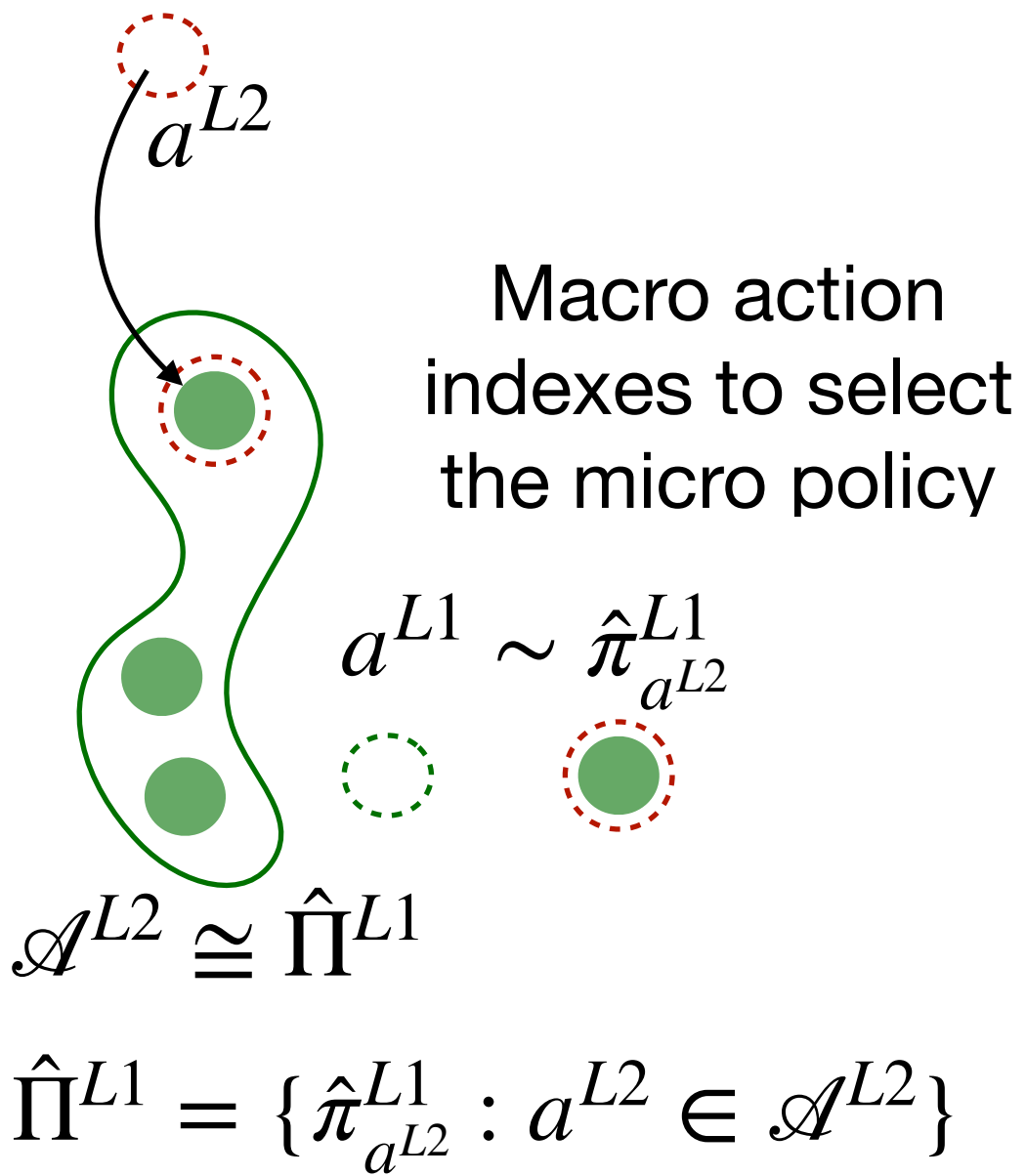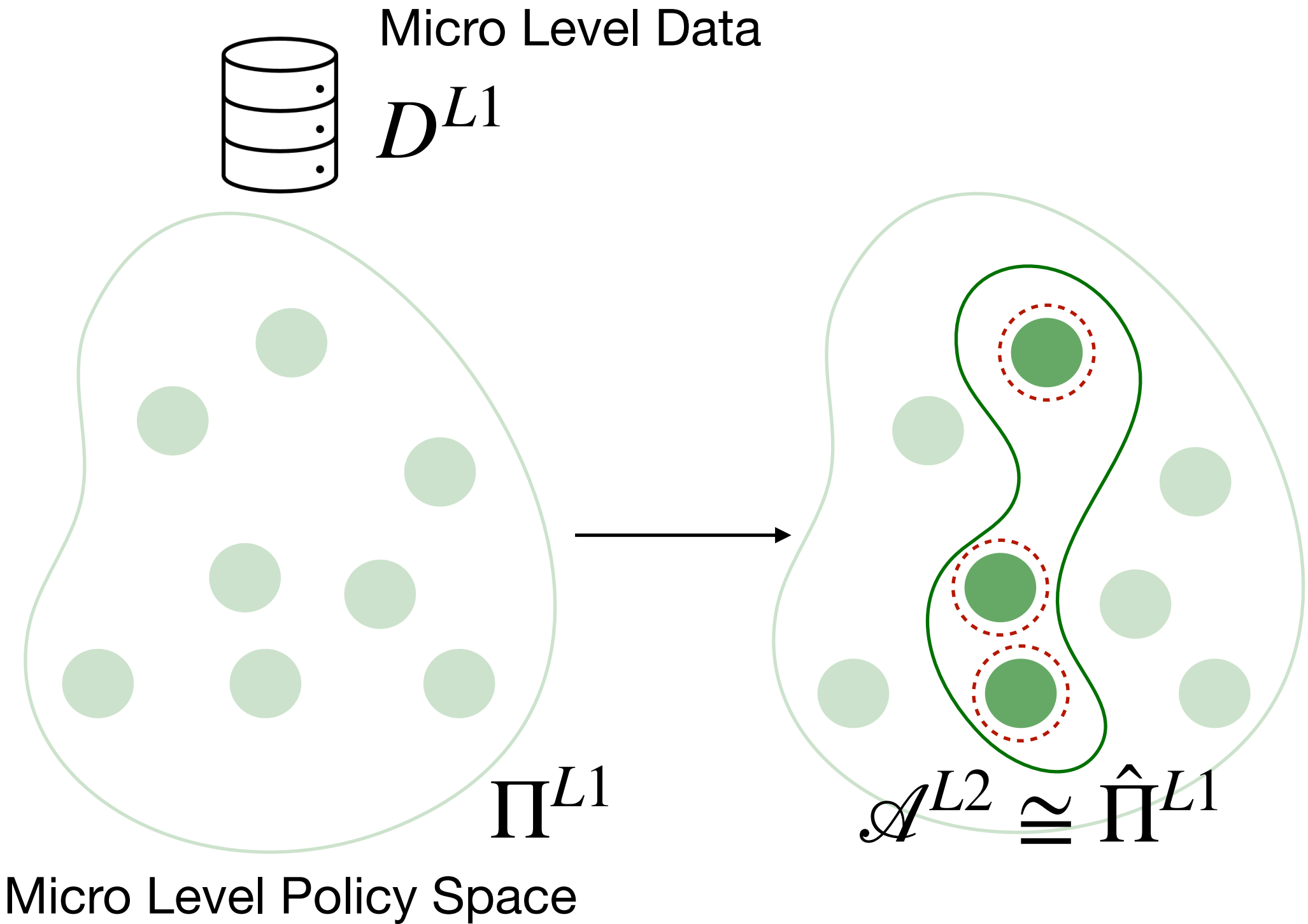$$\Pi \triangleq \Pi^{L1} \cdot \Pi^{L2}$$

Interventions

Feedback

Level 2 (Macro)

Weekly Return Rate

$a^{L2}$ | Decoding Strategy {Decoding Temp., top-p sampling} | $r^{L2}$        $a^{L2}$        $r^{L2}$

Level 1 (Micro)

Engagement

$a^{L1}$ $r^{L1}$    Response generation    $a^{L1}$ $r^{L1}$    $a^{L1}$ $r^{L1}$    ......    $a^{L1}$ $r^{L1}$    $a^{L1}$ $r^{L1}$

Learns a large part of the parameter space

Inductive bias for long term optimal policy

$\Pi^{L1}$

Simplified learning at macro level

Small policy space

$\Pi^{L2}$

# Policy Learning at micro level



Micro Level Data
$D^{L1}$

$\Pi^{L1}$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

$\Pi^{L1}$

Micro Level Policy Space

$a^{L2}$

Macro action indexes to select the micro policy

$a^{L1} \sim \hat{\pi}^{L1}_{a^{L2}}$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

$\hat{\Pi}^{L1} = \{\hat{\pi}^{L1}_{a^{L2}} : a^{L2} \in \mathscr{A}^{L2}\}$

$\pi^{L2}$

# Policy Learning at micro level

Micro Level Data
$D^{L1}$

Policy or
Feedback
Modification

Micro Level Policy Space

$\Pi^{L1}$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

$\Pi^{L1}$

Example:

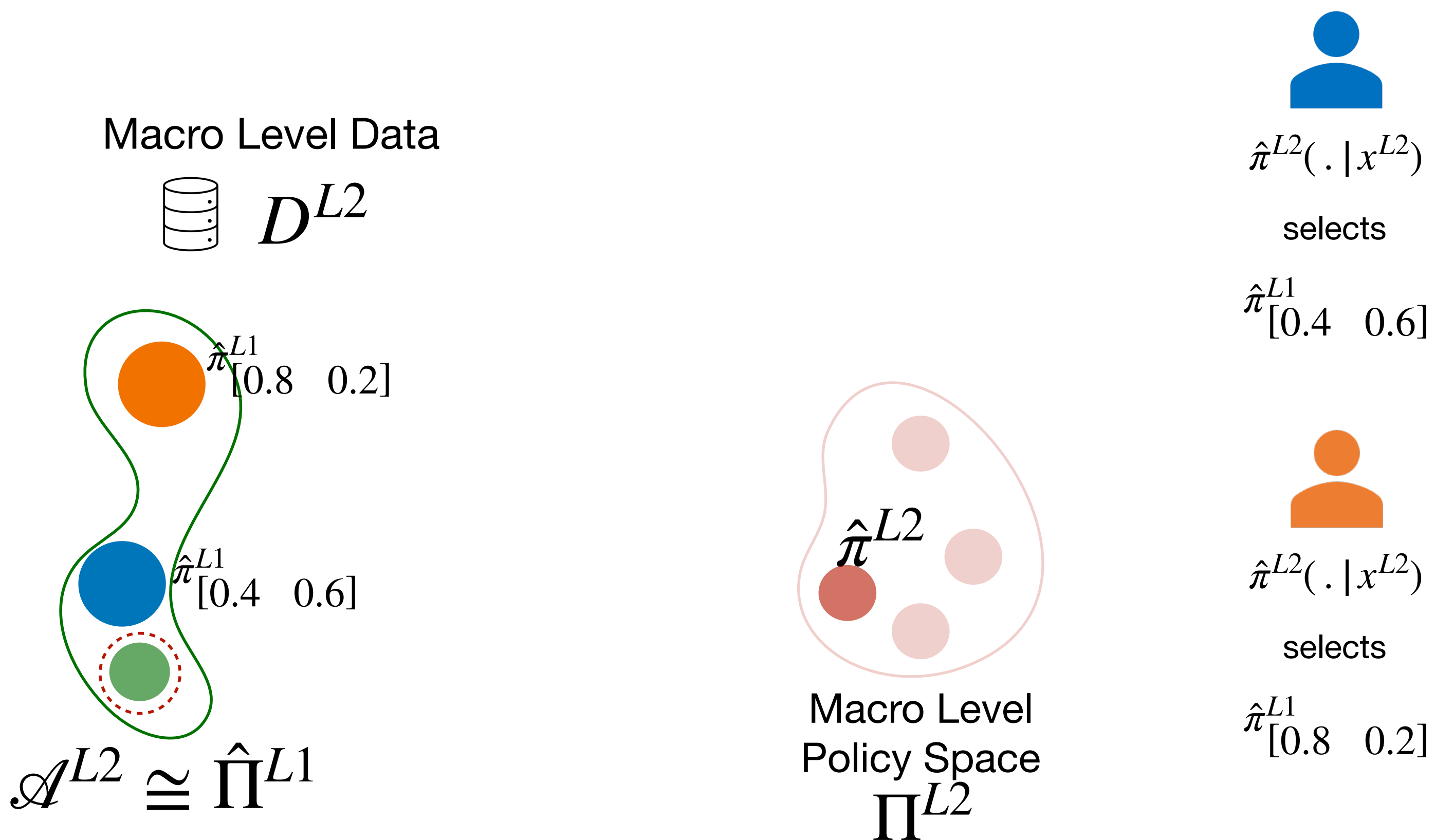$$r^{L1} \triangleq \begin{bmatrix} \text{Clicks} \\ \text{Likes} \end{bmatrix}$$

$$\begin{bmatrix} 0.8 & 0.2 \end{bmatrix} \begin{bmatrix} \text{Clicks} \\ \text{Likes} \end{bmatrix}$$

$$\begin{bmatrix} 0.4 & 0.6 \end{bmatrix} \begin{bmatrix} \text{Clicks} \\ \text{Likes} \end{bmatrix}$$

$$\cdots$$

$$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1} = \{\hat{\pi}^{L1}_{[0.8 \quad 0.2]}, \hat{\pi}^{L1}_{[0.4 \quad 0.6]}, \cdots\}$$

$$\pi^{L2}$$

# Policy Learning at macro level

Macro Level Data

$\boxminus$ $D^{L2}$

$\hat{\pi}^{L1}$
$[0.8 \quad 0.2]$

$\hat{\pi}^{L1}$
$[0.4 \quad 0.6]$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

$\hat{\pi}^{L2}$

Macro Level
Policy Space
$\Pi^{L2}$

$\hat{\pi}^{L2}(\,.\,|x^{L2})$

selects

$\hat{\pi}^{L1}$
$[0.4 \quad 0.6]$

$\hat{\pi}^{L2}(\,.\,|x^{L2})$

selects

$\hat{\pi}^{L1}$
$[0.8 \quad 0.2]$

$\Pi^{L1}$

$\pi^{L2}$

# MultiScale Contextual Bandits Algorithm

---

**Algorithm 1** MultiScale Training: Off-Policy Contextual Bandits

---

**Procedure** $PolicyLearning(\pi_0^{L2}, \pi_0^{L1})$

Collect Micro Logged dataset $D^{L1} := \{(x_i^{L1}, a_i^{L1}, r_i^{L1}, p_i^{L1})\}_{i=1}^{n^{L1}} \sim \pi_0^{L1}$

Learn Micro policies $\hat{\Pi}^{L1}$(Eq. (5) or (6) using $D^{L1}$)

---

# MultiScale Contextual Bandits Algorithm

---
**Algorithm 1** MultiScale Training: Off-Policy Contextual Bandits

---
**Procedure** $PolicyLearning(\pi_0^{L2}, \pi_0^{L1})$

Collect Micro Logged dataset $D^{L1} := \{(x_i^{L1}, a_i^{L1}, r_i^{L1}, p_i^{L1})\}_{i=1}^{n^{L1}} \sim \pi_0^{L1}$

Learn Micro policies $\hat{\Pi}^{L1}$(Eq. (5) or (6) using $D^{L1}$)

Collect Macro Logged dataset $D^{L2} := \{(x_j^{L2}, a_j^{L2}, r_j^{L2}, p_j^{L2})\}_{j=1}^{n^{L2}} \sim \pi_0^{L2}$
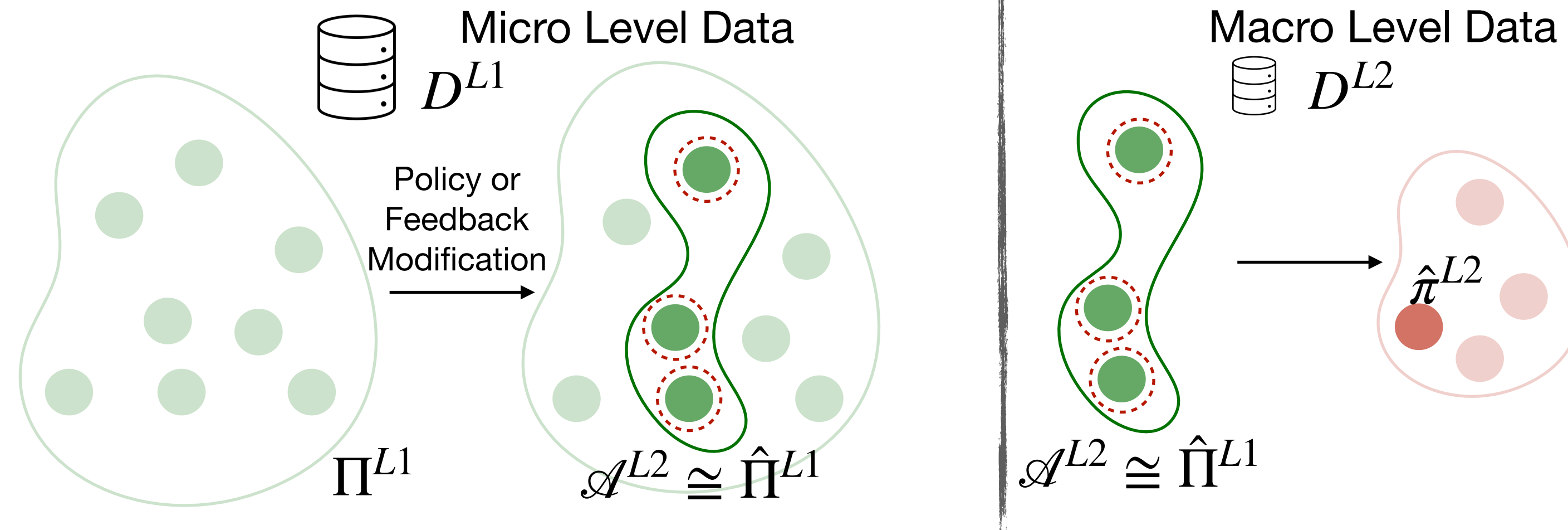
Learn Macro Policy $\hat{\pi}^{L2} \leftarrow \arg\max_{\pi^{L2}} \hat{V}^{L2}(\pi^{L2}; D^{L2})$ (Eq. (7))

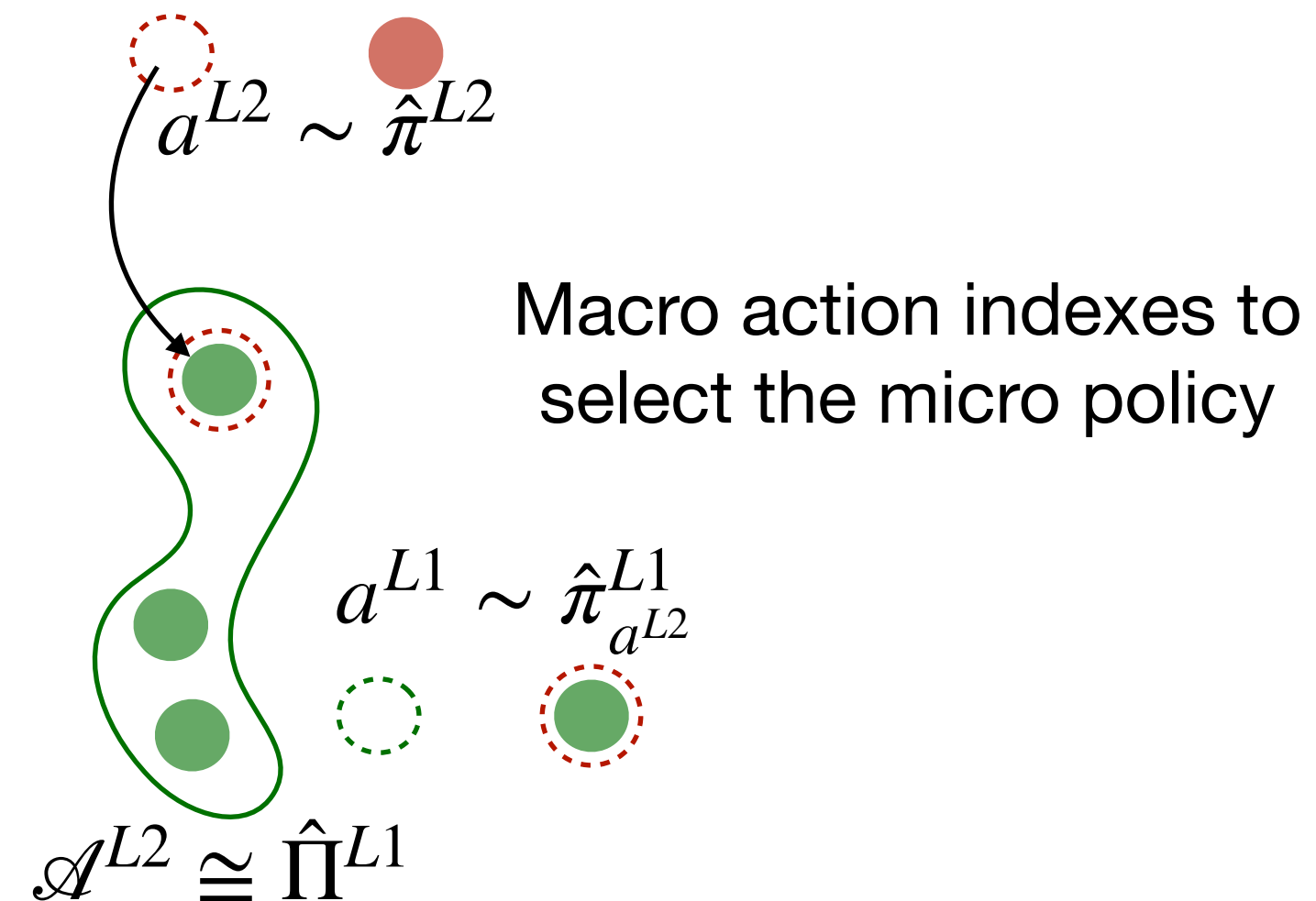**return** learned policies $\hat{\pi}^{L2}, \hat{\Pi}^{L1}$

---

This procedure can be recursively called for extending to arbitrary number of levels

# MultiScale Contextual Bandits

- Training

  - Bottom up



Micro Level Data
$D^{L1}$

Policy or
Feedback
Modification

$\Pi^{L1}$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

Macro Level Data
$D^{L2}$

$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

$\hat{\pi}^{L2}$

- Inference

  - Top down

$a^{L2} \sim \hat{\hat{\pi}}^{L2}$

Macro action indexes to
select the micro policy

$a^{L1} \sim \hat{\pi}^{L1}_{a^{L2}}$
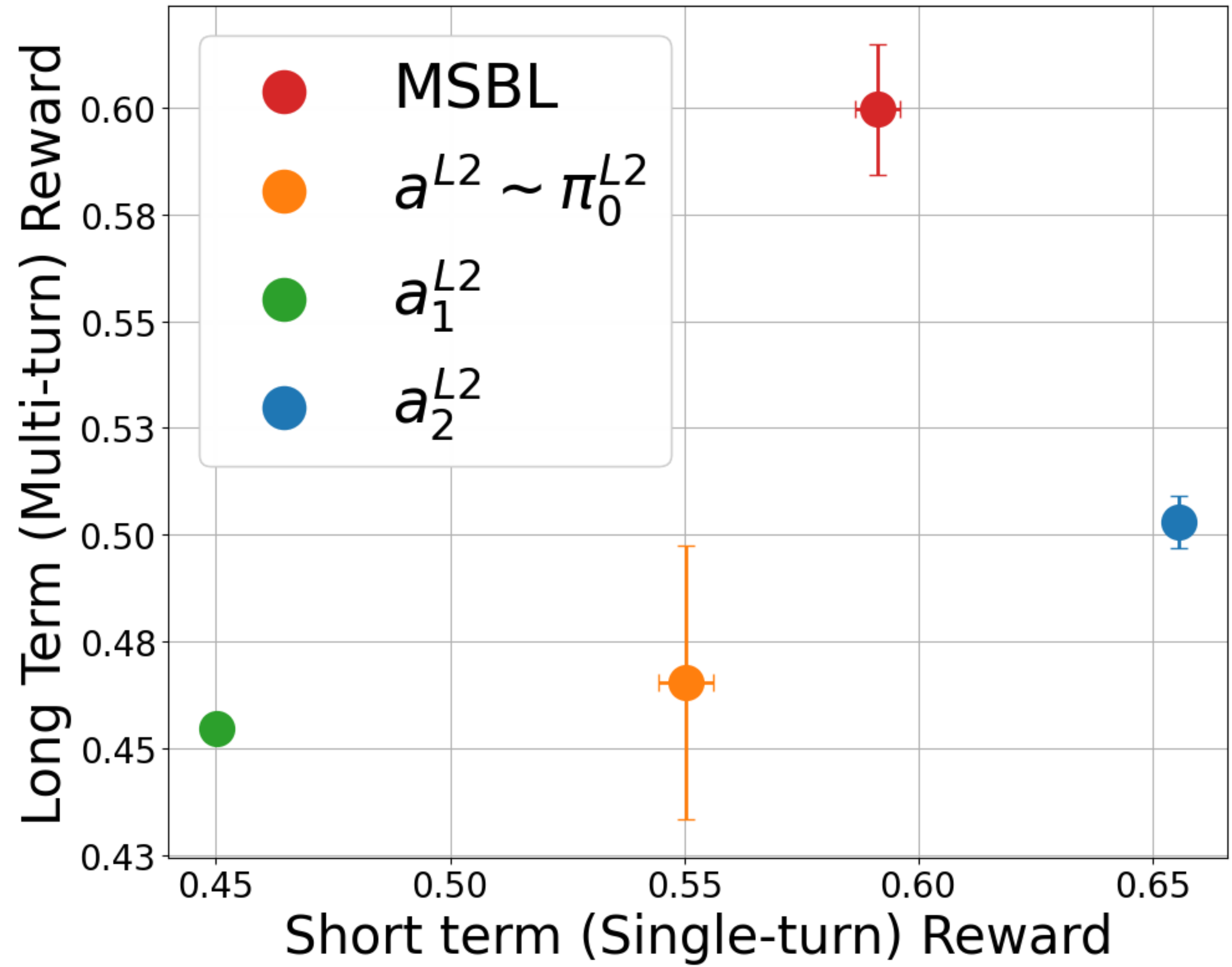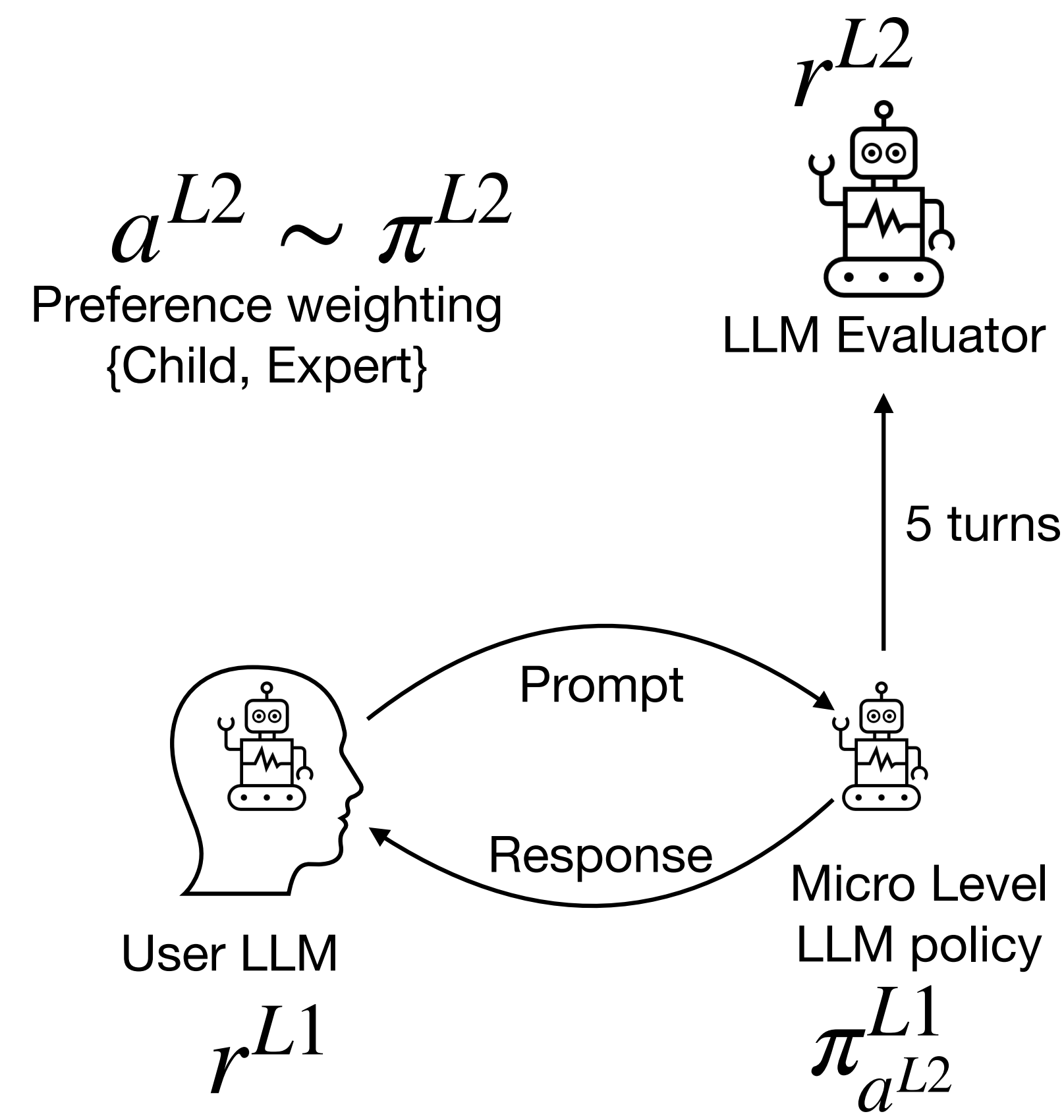
$\mathscr{A}^{L2} \cong \hat{\Pi}^{L1}$

# Experiments

- Multi turn Conversation

- Conversational recommender system

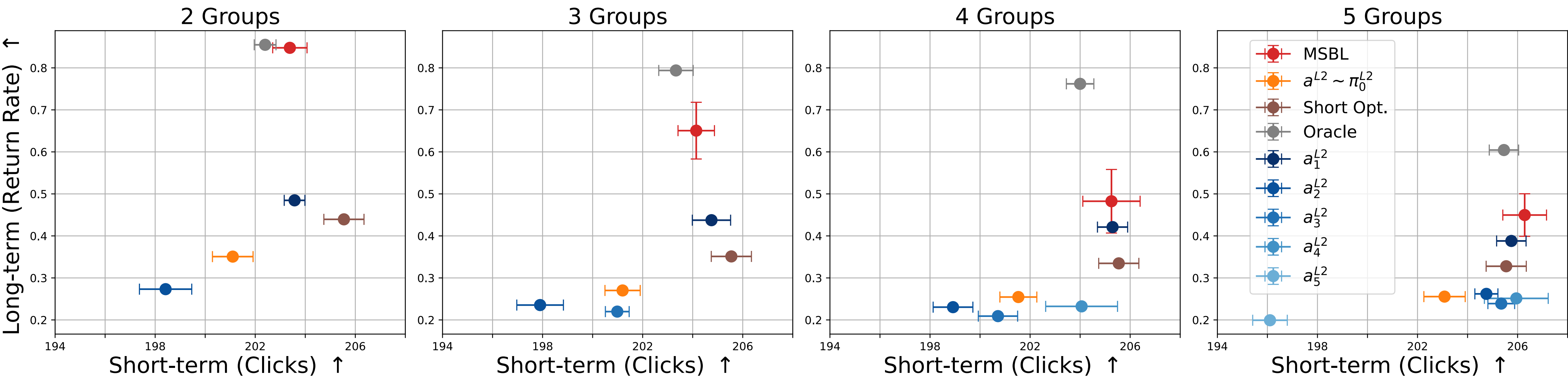- Large Scale Recommender System

# Experiments

## Multi turn Conversation

$$a^{L2} \sim \pi^{L2}$$

Preference weighting
{Child, Expert}

$$r^{L2}$$

LLM Evaluator

5 turns

Prompt

Response

User LLM

$$r^{L1}$$

Micro Level
LLM policy

$$\pi^{L1}_{a^{L2}}$$

# Experiments

## Conversational Recommender System

# Experiments

## Large Scale Recommender System

# Summary

- Introduce a principled framework to optimize for long term objectives.

- Motivated by using plentiful short-term data for faster learning with scarcer long term feedback

- We discuss two ways - policy and feedback modification to learn a family of policies at micro level.

- Propose a practical bandit algorithm for recursively learning policies at multiple interdependent levels.

- Checkout the paper for more results and analysis - PAC Bayesian motivation, updating micro and macro policies after deployment when new data is available, scaling action space and more !