

# Sim2Rec: A Simulator-based Decision-making Approach to Optimize Real-World Long-term User Engagement in Sequential Recommender Systems

Xiong-Hui Chen<sup>1,3</sup>, Yang Yu<sup>1,3,\*</sup>, Qingyang Li<sup>2</sup>, Bowei He<sup>4</sup>, Zhiwei Qin<sup>2</sup>, Wenjie Shang<sup>2</sup>, Jieping Ye<sup>2</sup>

<sup>1</sup> National Key Laboratory of Novel Software Technology, Nanjing University, Nanjing, China

<sup>2</sup> Didi Chuxing, <sup>3</sup> Polixir.ai, <sup>4</sup> City University of Hong Kong

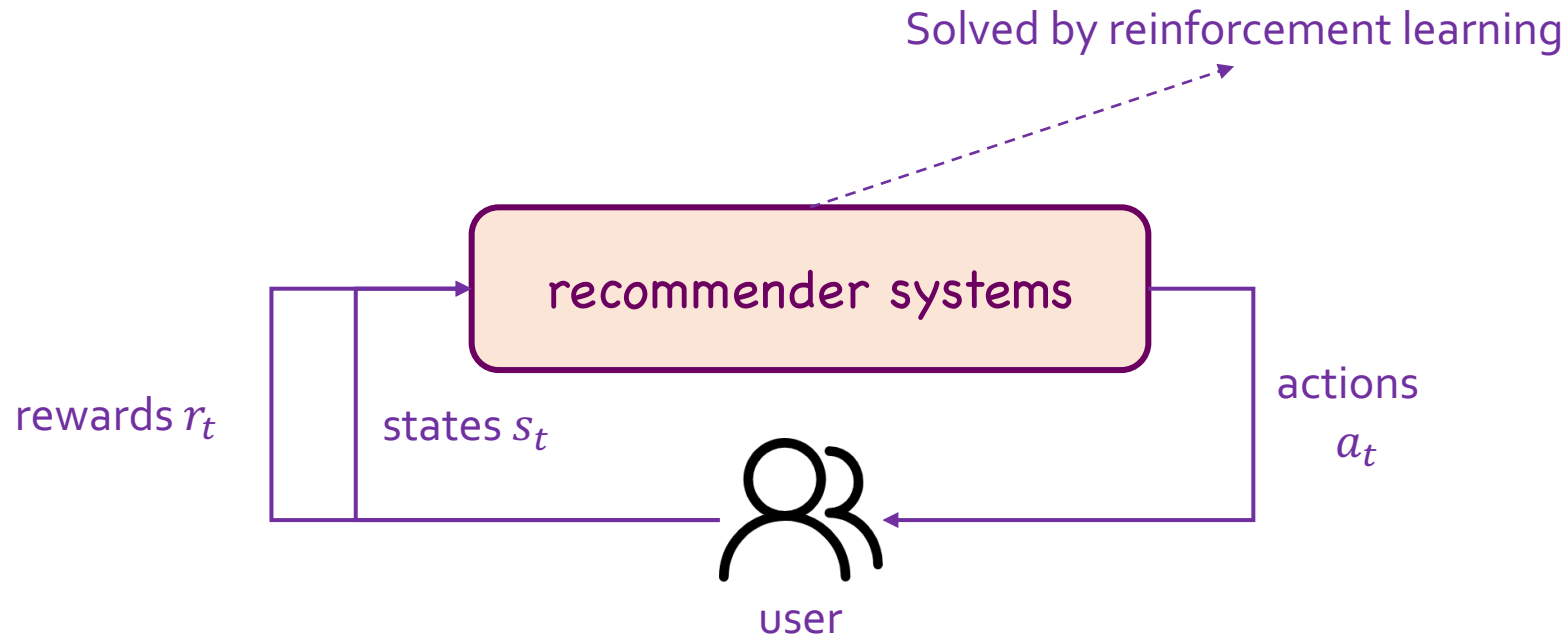
\* Corresponding author

chenxh@lamda.nju.edu.cn, yuy@nju.edu.cn, qingyangli@didiglobal.com, boweihe2-c@my.cityu.edu.hk

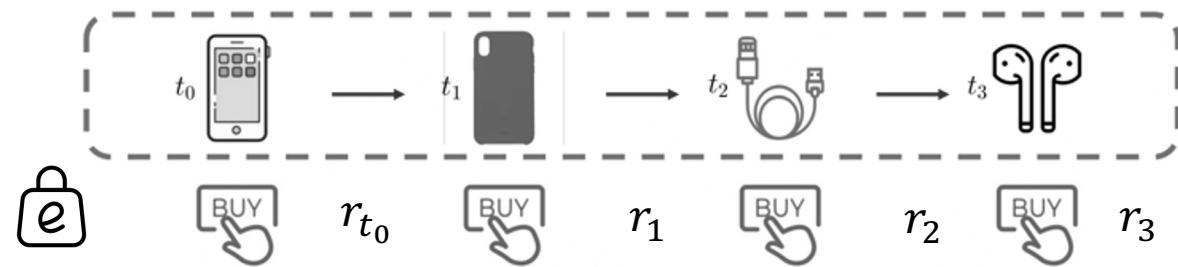
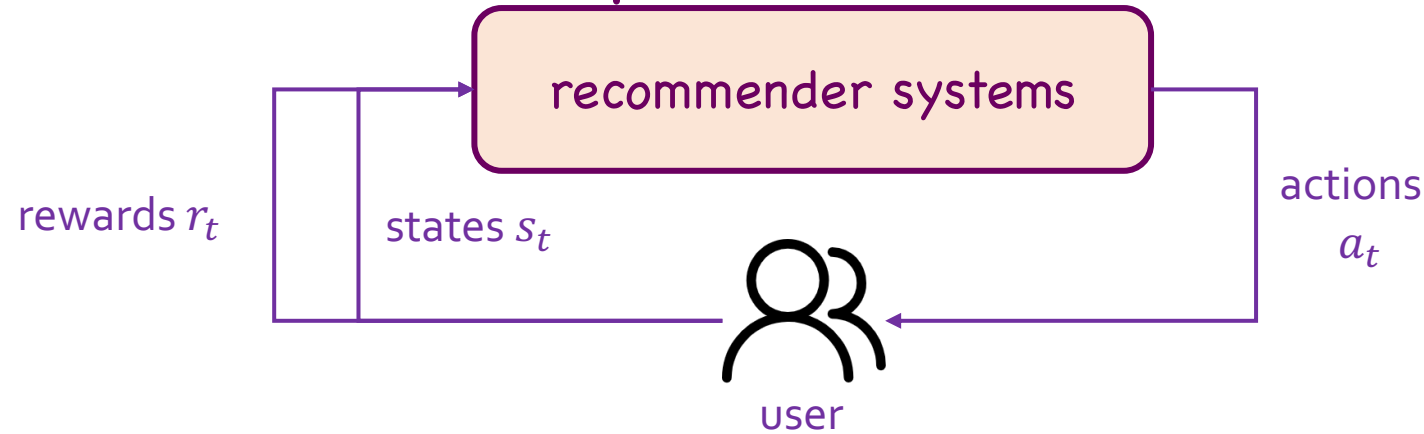
# Outline

1. Background and Motivation
2. Simulation to Recommender Systems (Sim2Rec)
3. Experiment
4. Take-home Messages

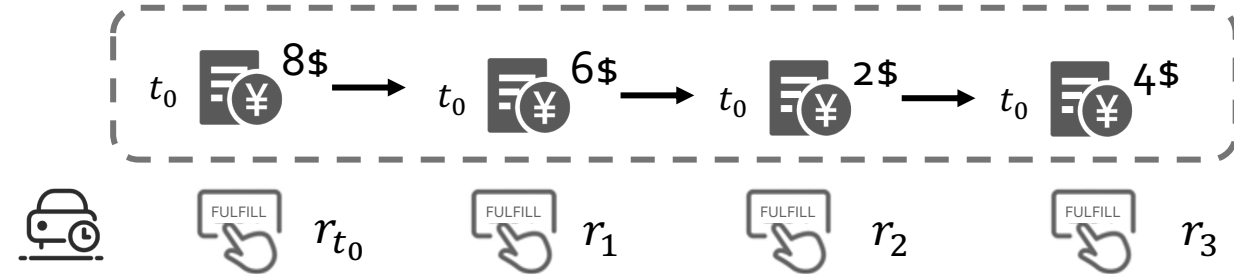
# Reinforcement Learning (RL) to Optimize Real-World Long-term User Engagement (LTE) in Sequential Recommender Systems (SRS)



# Reinforcement Learning (RL) to Optimize Real-World Long-term User Engagement (LTE) in Sequential Recommender Systems (SRS)



product recommendation in e-commerce platforms

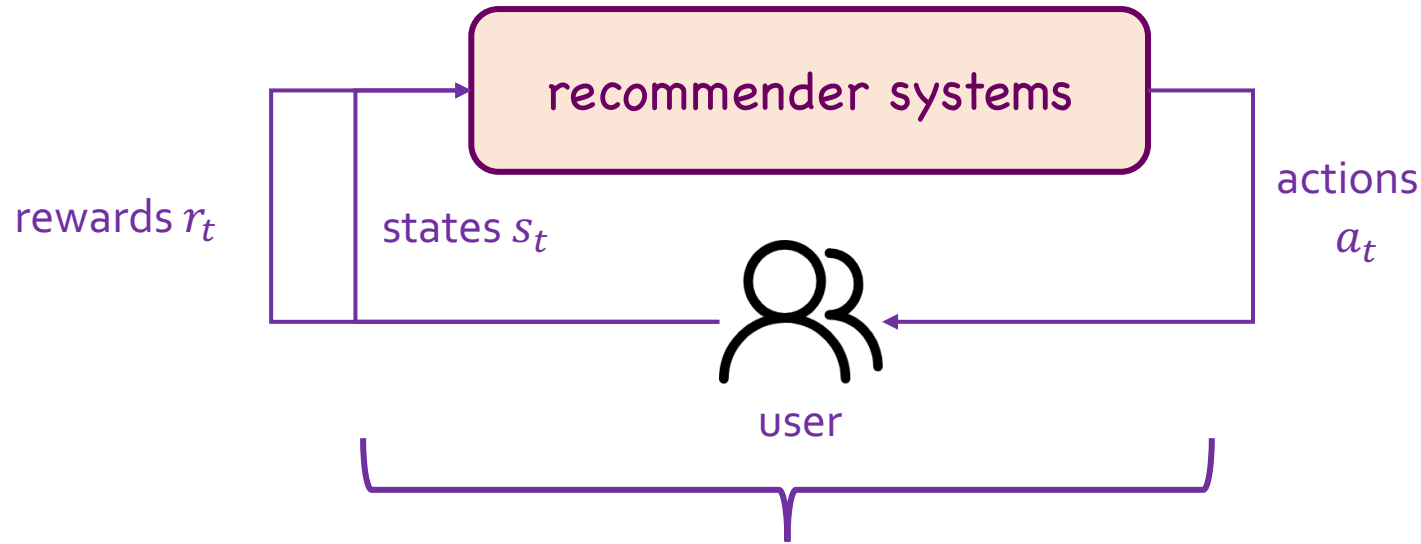


order dispatching in ride-hailing platforms

objective: maximize cumulative user engagement instead of immediate rewards

# Reinforcement Learning (RL) to Optimize Real-World Long-term User Engagement (LTE) in Sequential Recommender Systems (SRS)

objective of RL for SRS: maximize cumulative rewards instead of immediate rewards



need massive online interactions for exploration to find an optimal recommender policy.



costly: time-consuming and unsafe (might cause economic losses) in many real-world applications

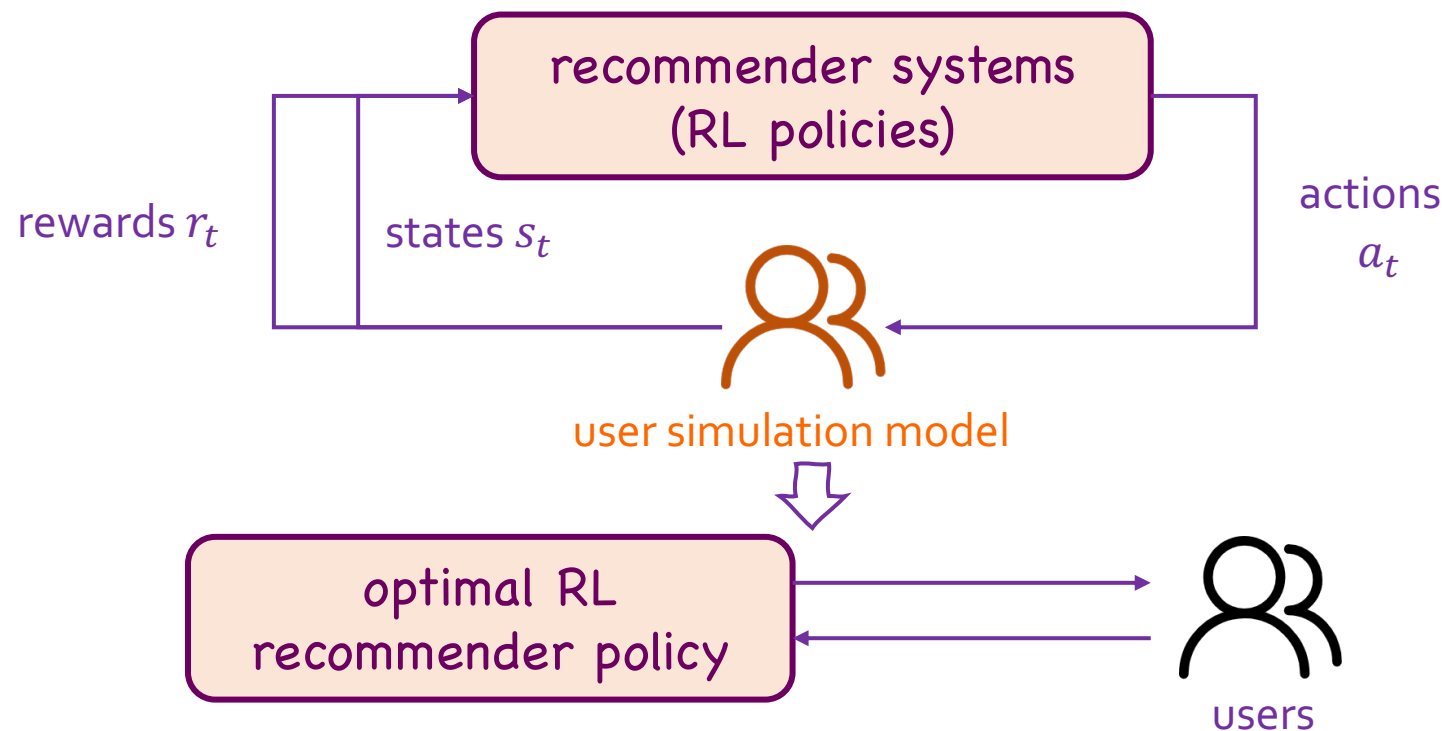
# Simulation-based RL to Optimize LTE in SRS

objective of RL for SRS: maximize cumulative rewards instead of immediate rewards

-> Need massive online interactions for exploration to find an optimal recommender policy, which is costly.

✓ ideal solution: simulation-based RL for SRS

1. training RL policies in customer simulation models learned by machine learning techniques.
2. deploy after the policy learned to obtain the optimal cumulative rewards in the simulation models.



# The Reality-gaps Challenge in Simulation-based RL Workflow

objective of simulation-based RL for SRS: maximize cumulative rewards in the simulation models



user simulation model

Learned by machine learning techniques  
(e.g., neural network or forest)



inevitably suffer prediction error  
to the real world

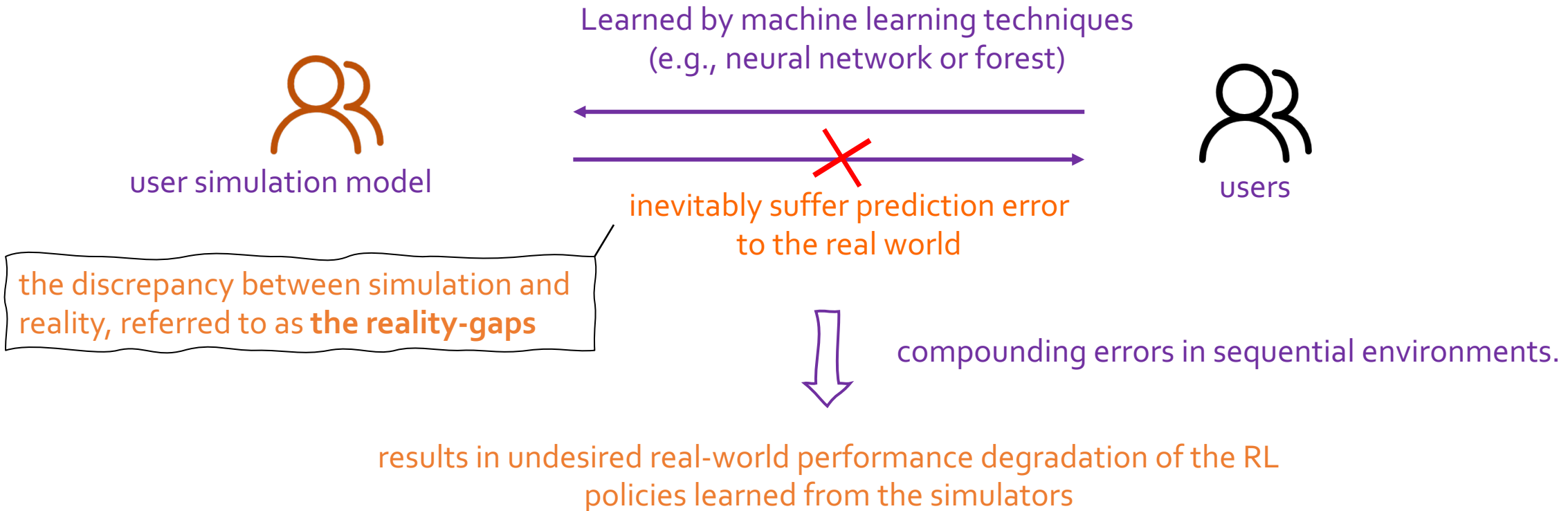


user

the discrepancy between simulation and  
reality, referred to as **the reality-gaps**

# The Reality-gaps Challenge in Simulation-based RL Workflow

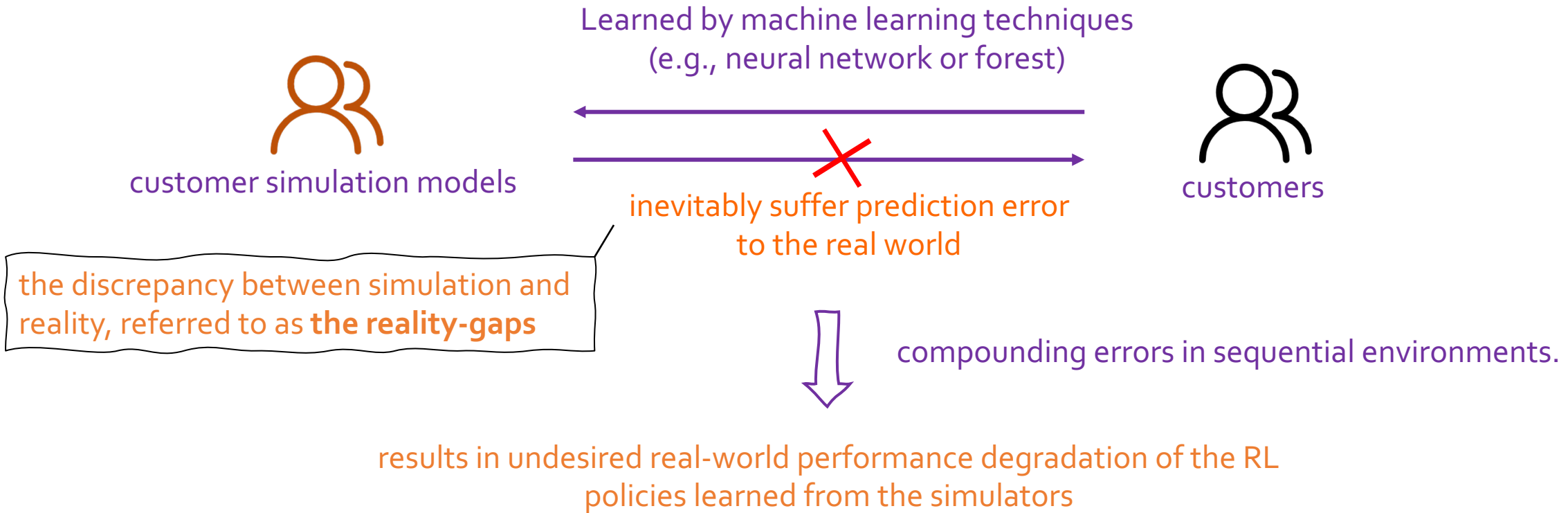
objective of simulation-based RL for SRS: maximize cumulative rewards in the simulation models





# The Reality-gaps Challenge in Simulation-based RL Workflow

objective of simulation-based RL for SRS: maximize cumulative rewards in the simulation models



*rarely been discussed explicitly!*

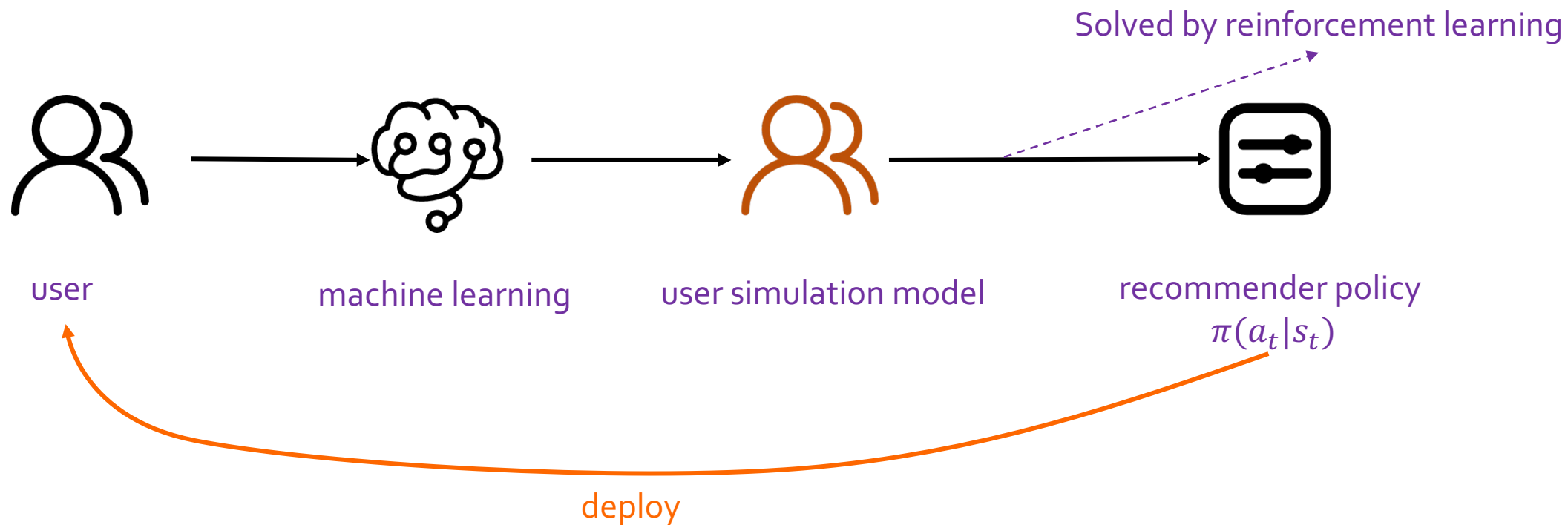
# Table of Contents

1. Background and Motivation
2. Simulation to Recommender Systems (Sim2Rec)
3. Experiment
4. Take-home Messages

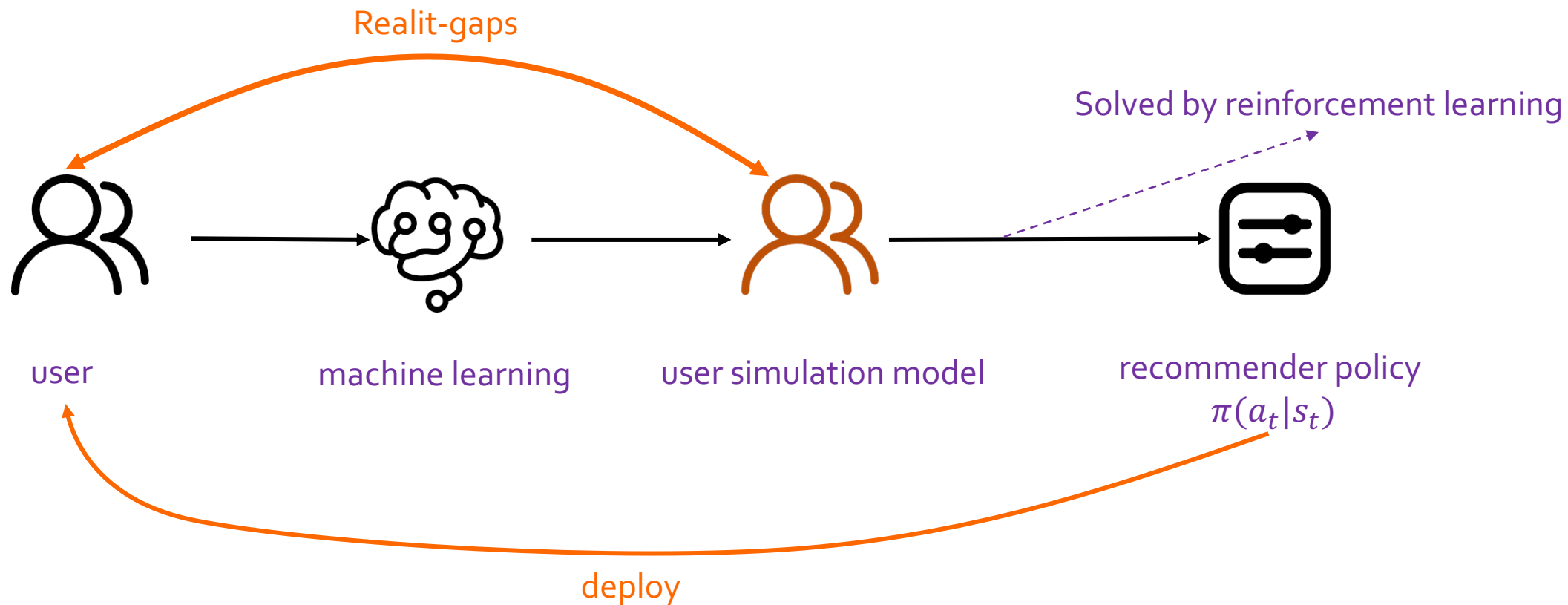
# Core contributions

1. Introduce **zero-shot policy transfer framework to SRS**, which is a popular solution in Robotics to solve the reality-gap problem between the physical simulator and the real world.
2. We identify and handle several extra challenges when adopting the framework to SRS.

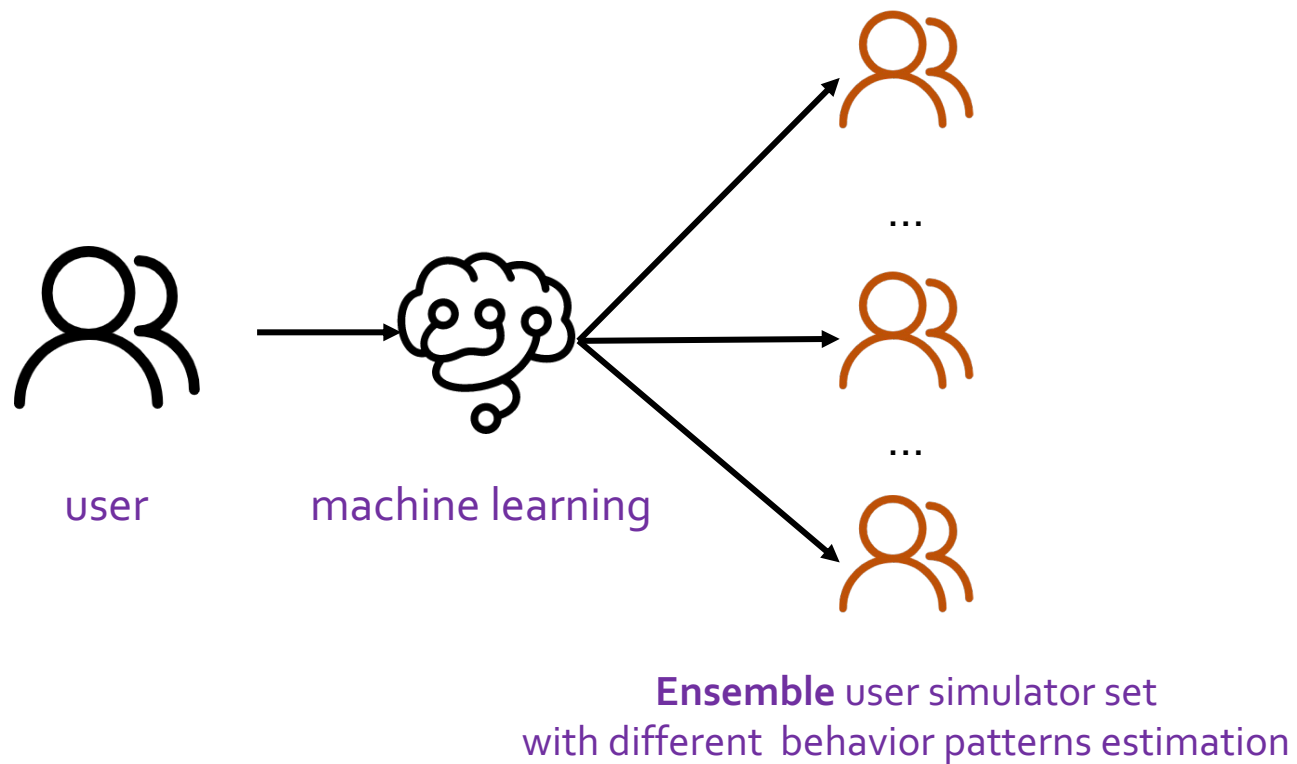
# Zero-shot Policy Transfer Framework for SRS: recap



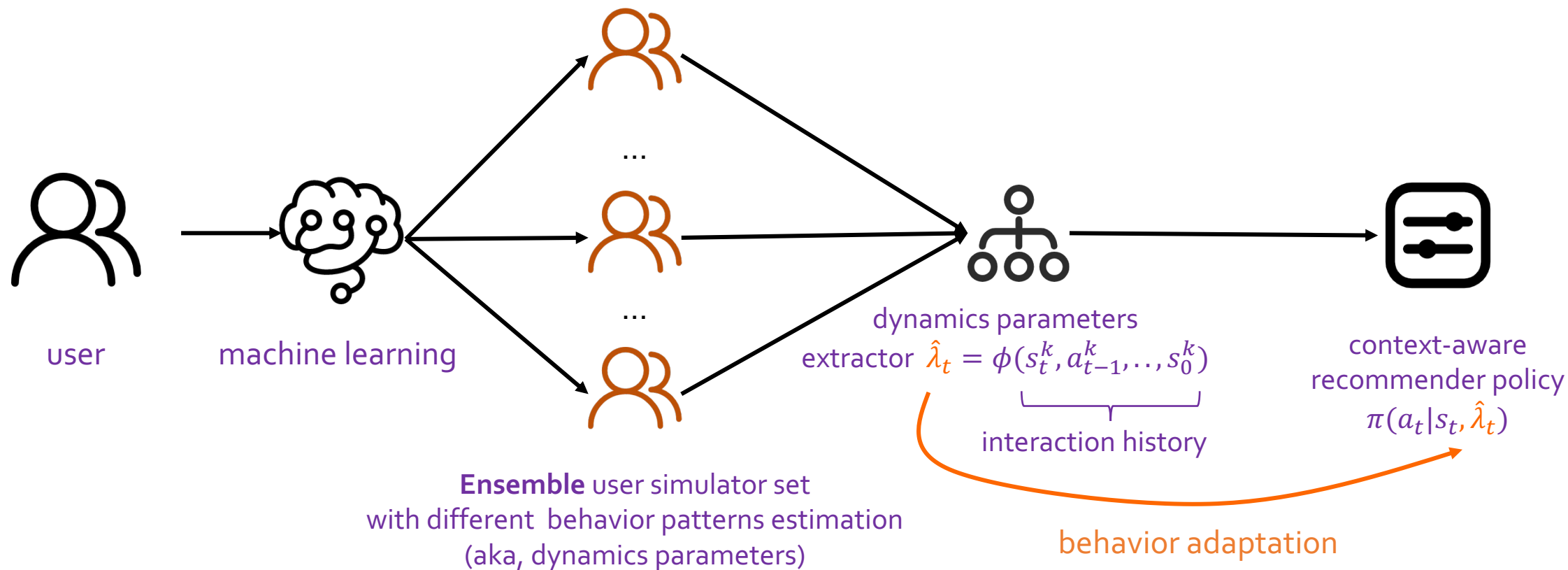
# Zero-shot Policy Transfer Framework for SRS: recap



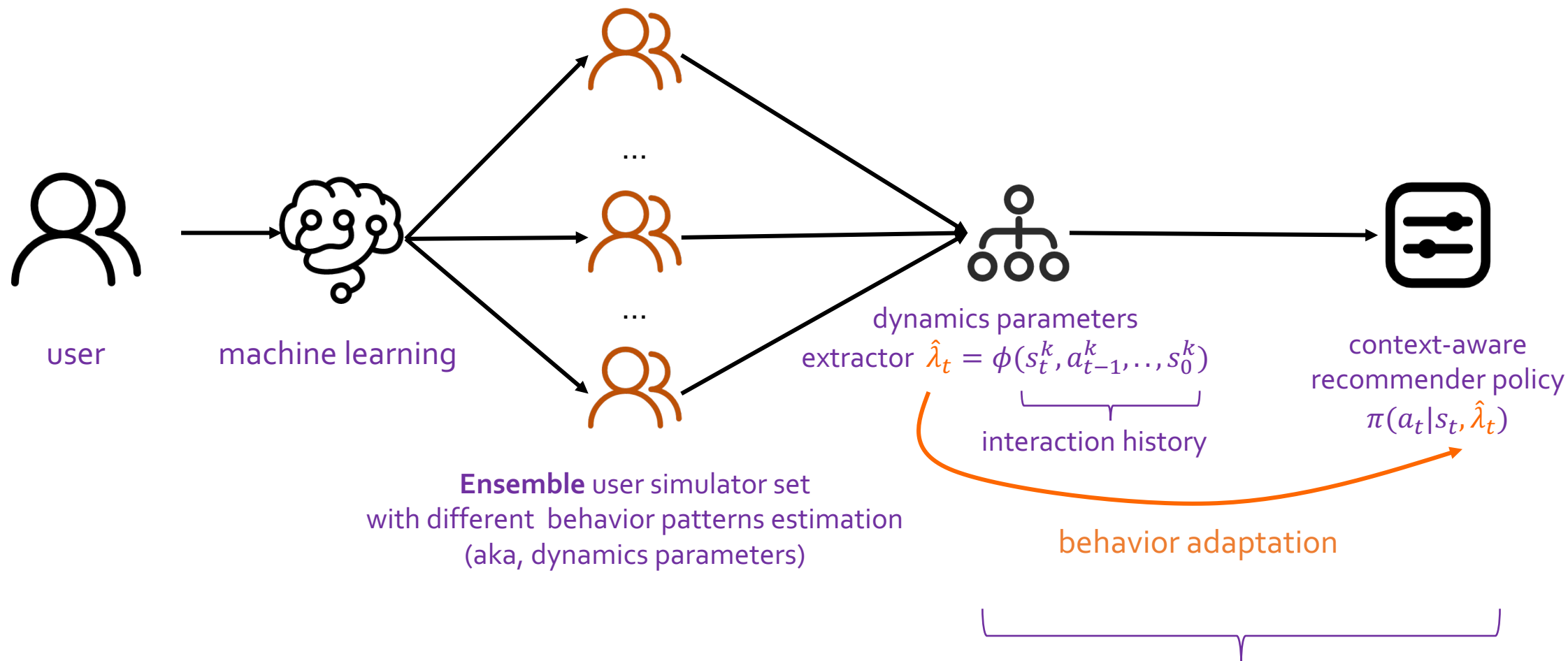
# Zero-shot Policy Transfer Framework for SRS



# Zero-shot Policy Transfer Framework for SRS

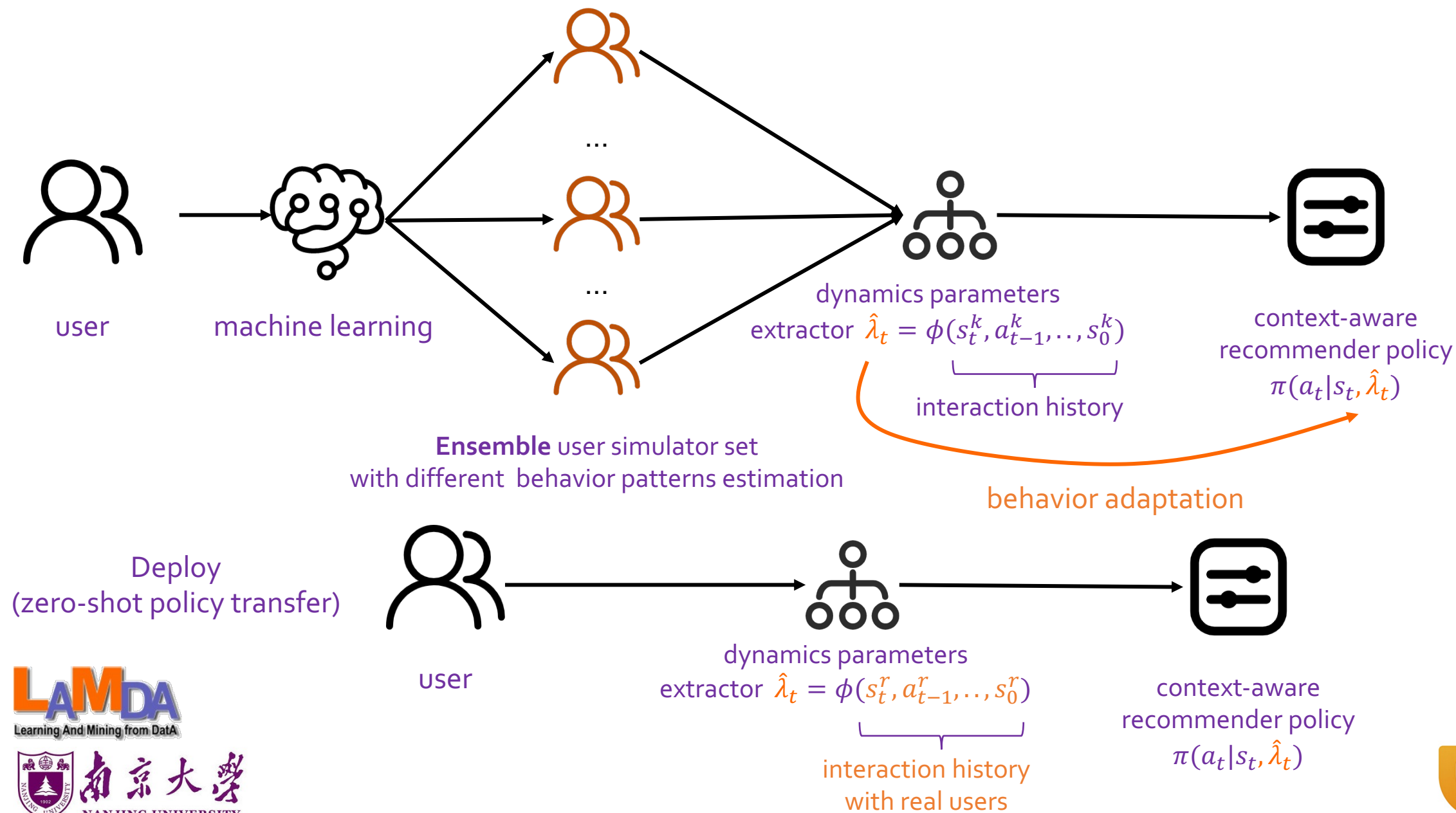


# Zero-shot Policy Transfer Framework for SRS

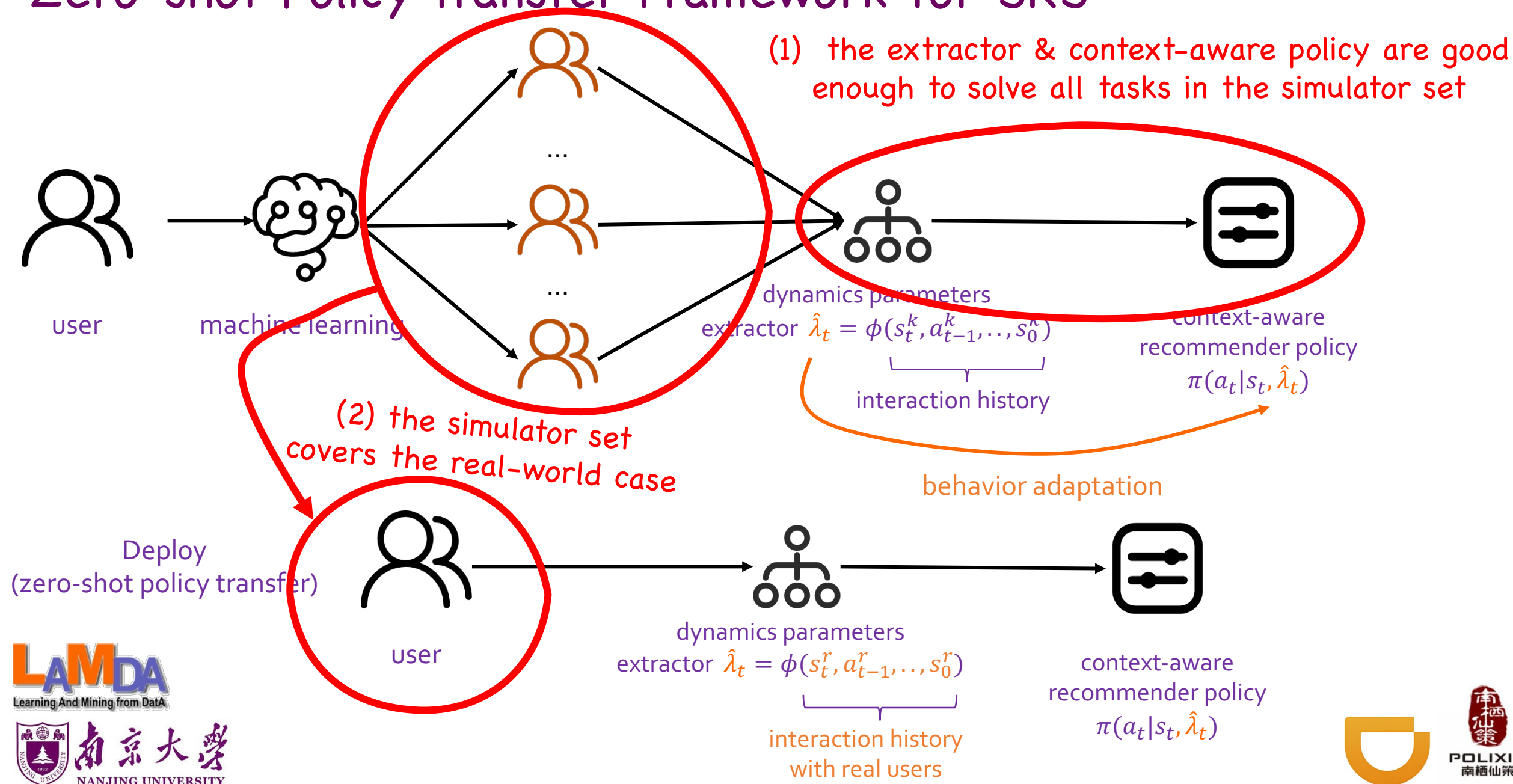




# Zero-shot Policy Transfer Framework for SRS



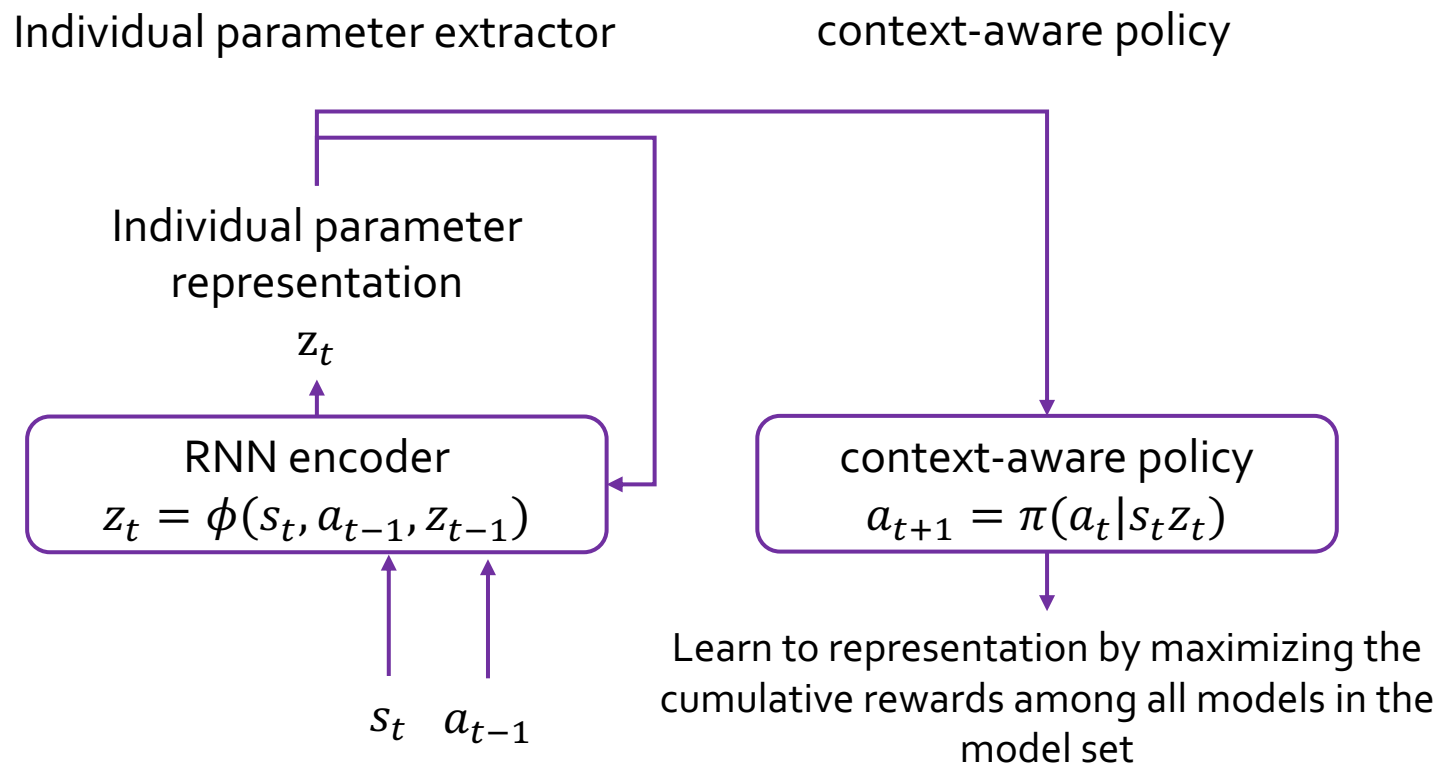
# Zero-shot Policy Transfer Framework for SRS



# The Challenge of zero-shot policy transfer framework to handle the Reality-gaps in SRS

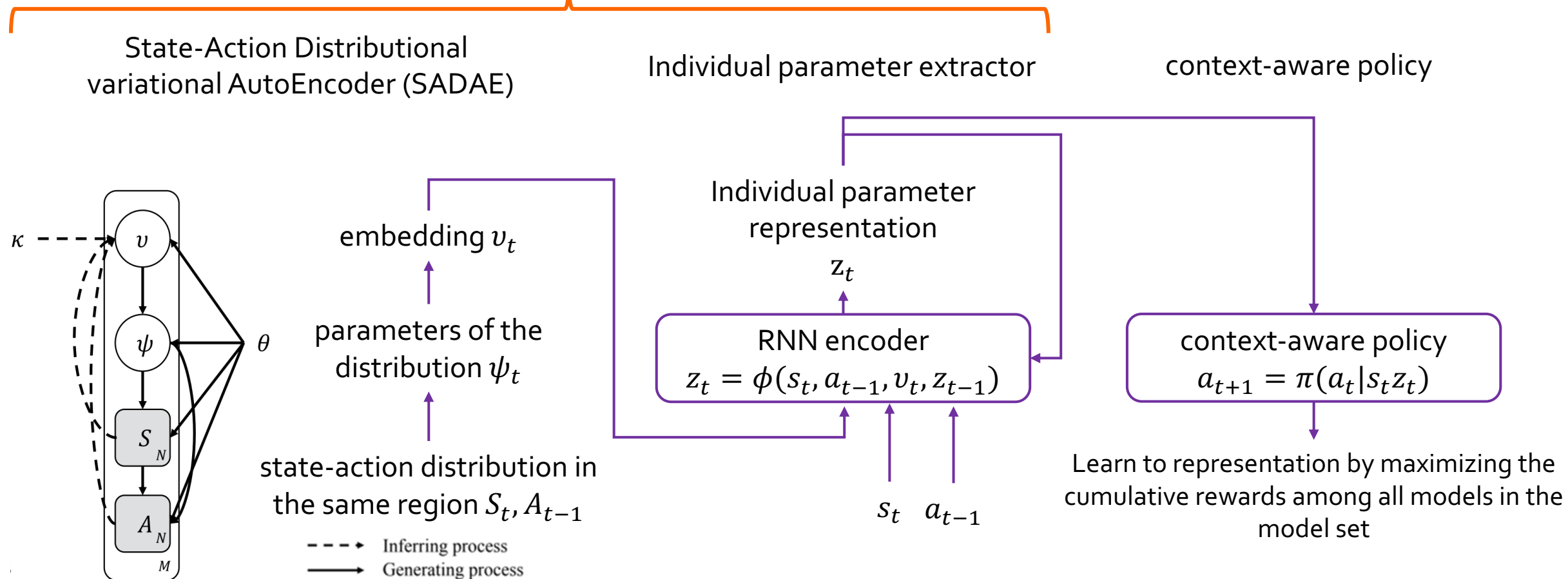
1. The complexity of extractor  $\phi$  to identify correct representations is much larger than in previous applications in robotics.
  - Need to identify **numerous user behavior patterns in different regions** from historical interactions.
    - In robotics, we often only need to infer **a single robot's** dynamics parameter.
  - » It is non-trivial to identify the representation of all users in all simulators
2. It is generally impractical to construct a simulator set that can cover all of the real-world user behaviors.
  - Limitation of the presentation capacity of the learning techniques.
  - Or rely on **extremely large size** of ensemble models, which is also impractical.
    - In robotics, simulators are built by physical laws, then we can generate a simulator set with coverability by adjusting the explicitly defined parameters (friction coefficients).
  - » Need to build a reliable policy learning paradigm for the simulator set which is uncoverable to reality.

# Sim2Rec Solution



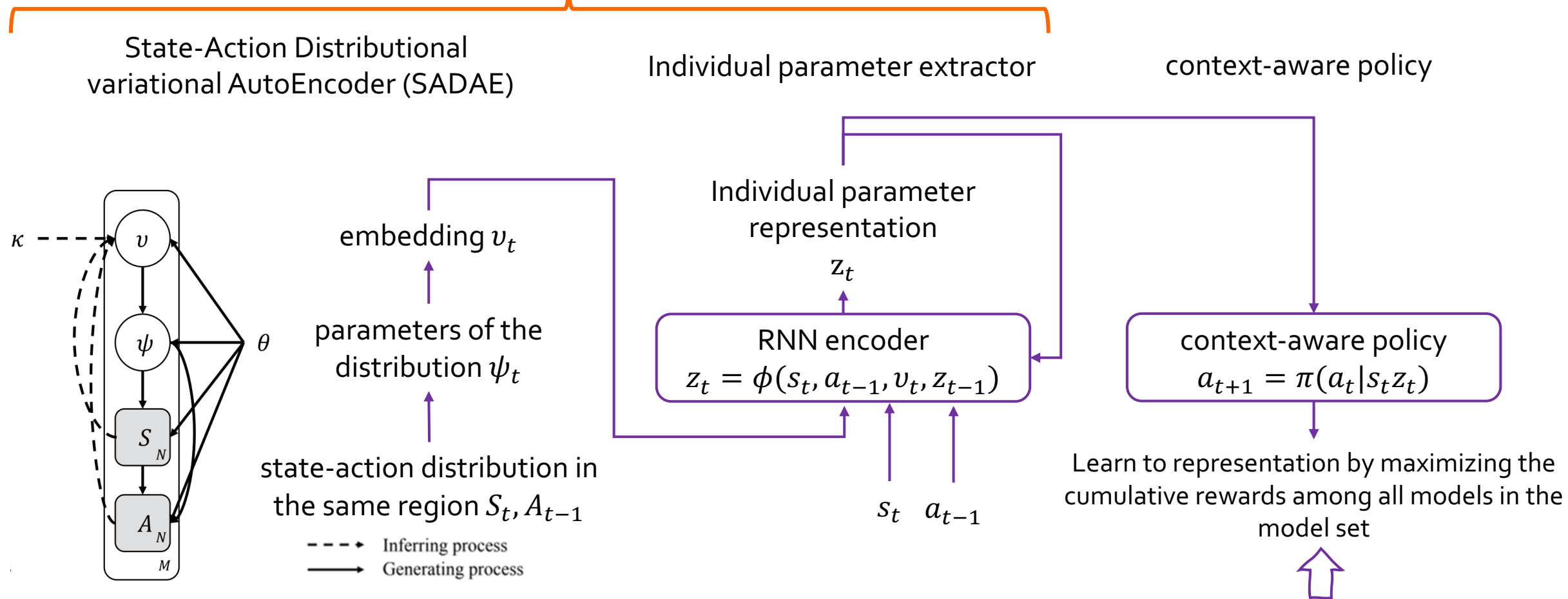
# Sim2Rec Solution

(1) efficient representation for extractor learning by utilizing the group information of different users in the same regions



# Sim2Rec Solution

(1) efficient representation for extractor learning by utilizing the group information of different users in the same regions



branching rollout steps to reduce compounding error

reward penalty via inconsistency of ensemble model predictions

pseudo-intervention test to filter out simulation data with severe extrapolation error

(2) reliable policy learning paradigm

# Table of Contents

1. Background and Motivation
2. Simulation to Recommender Systems
- 3. Experiment**
4. Take-home Messages

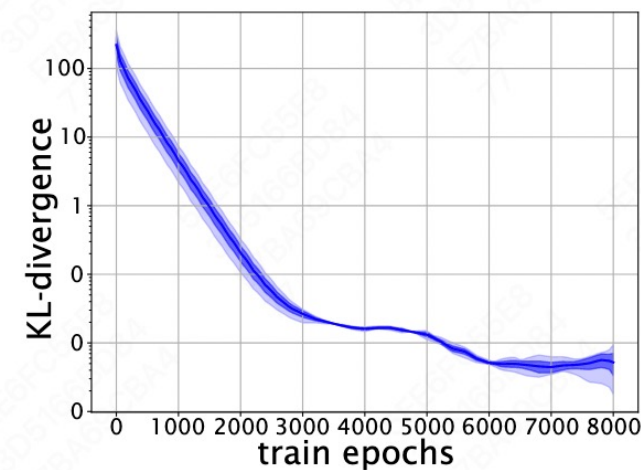
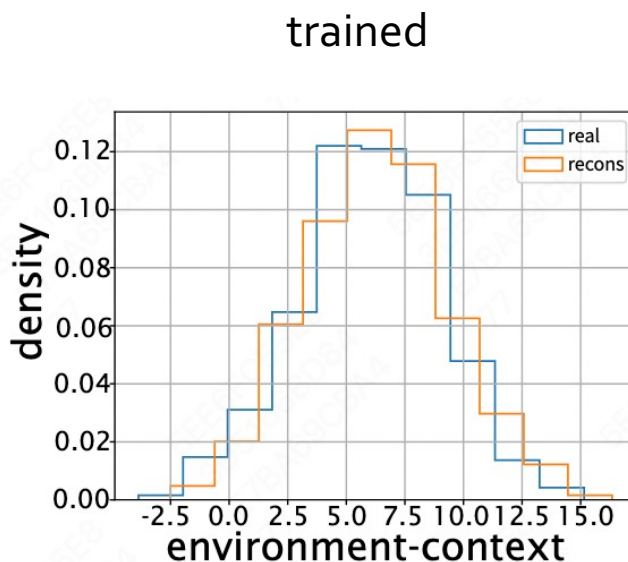
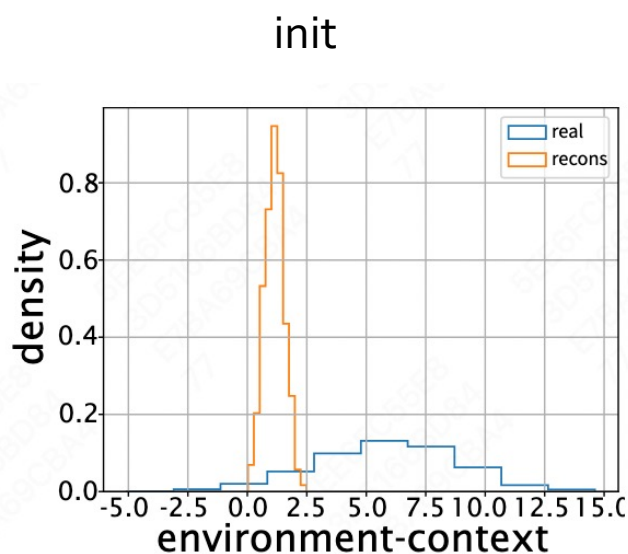
# Selected Experiment Results

RecSim

Real-world  
(Didichuxing)



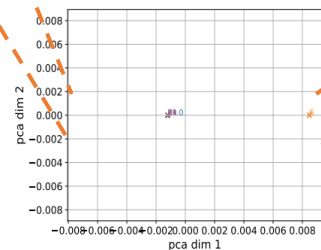
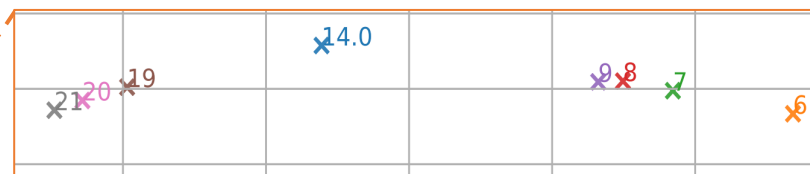
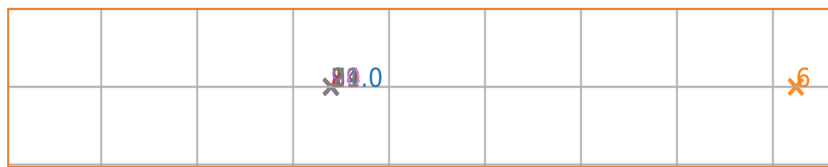
# Selected Experiment Results



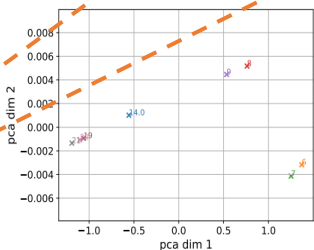
Real-world  
(Didichuxing)

# Selected Experiment Results

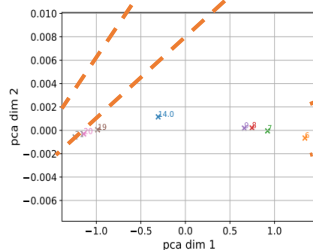
RecSim



(a) epoch 0

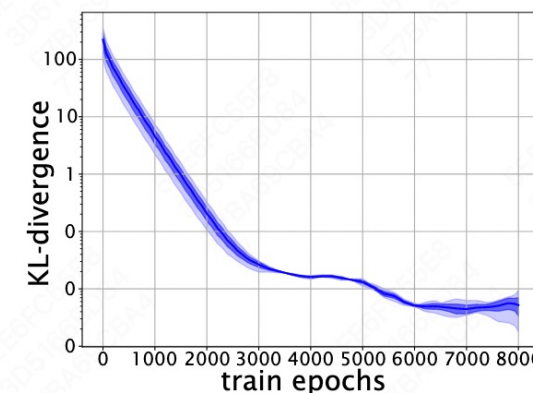
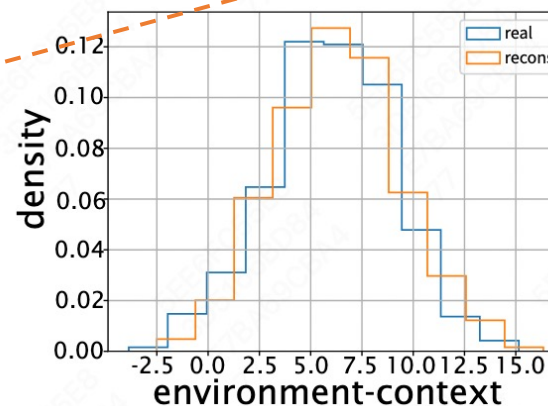


(b) epoch 4000



(c) epoch 8000

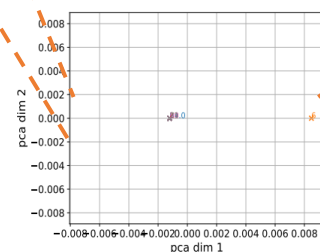
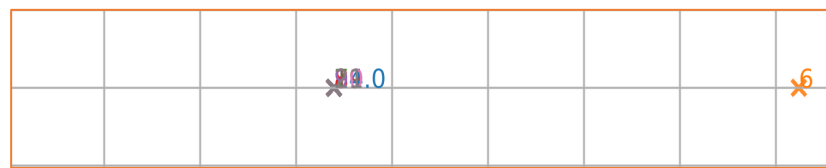
**Figure 1: Illustration of the visualization on  $v$ .** The X-axis denotes the first principal component, and the Y-axis denotes the second one. Each cross point denotes the projection of the latent code for the state distribution. The numbers with the same color to the point denote the ground-truth environment



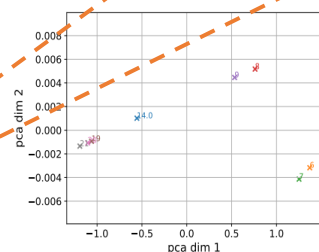
Real-world  
(Didichuxing)

# Selected Experiment Results

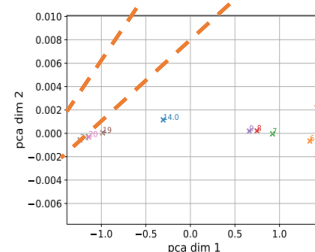
RecSim



(a) epoch 0

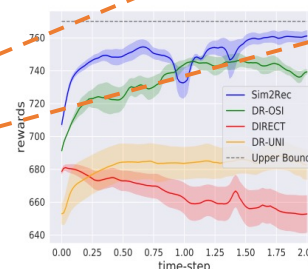
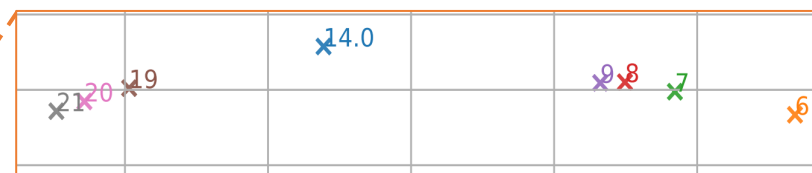


(b) epoch 4000

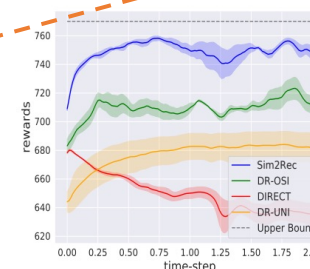


(c) epoch 8000

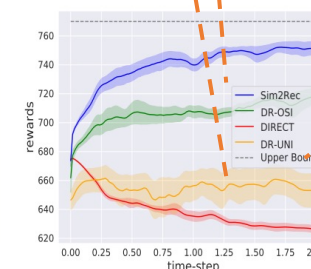
**Figure 1: Illustration of the visualization on  $v$ .** The X-axis denotes the first principal component, and the Y-axis denotes the second one. Each cross point denotes the projection of the latent code for the state distribution. The numbers with the same color to the point denote the ground-truth environment



(a) LTS1



(b) LTS2

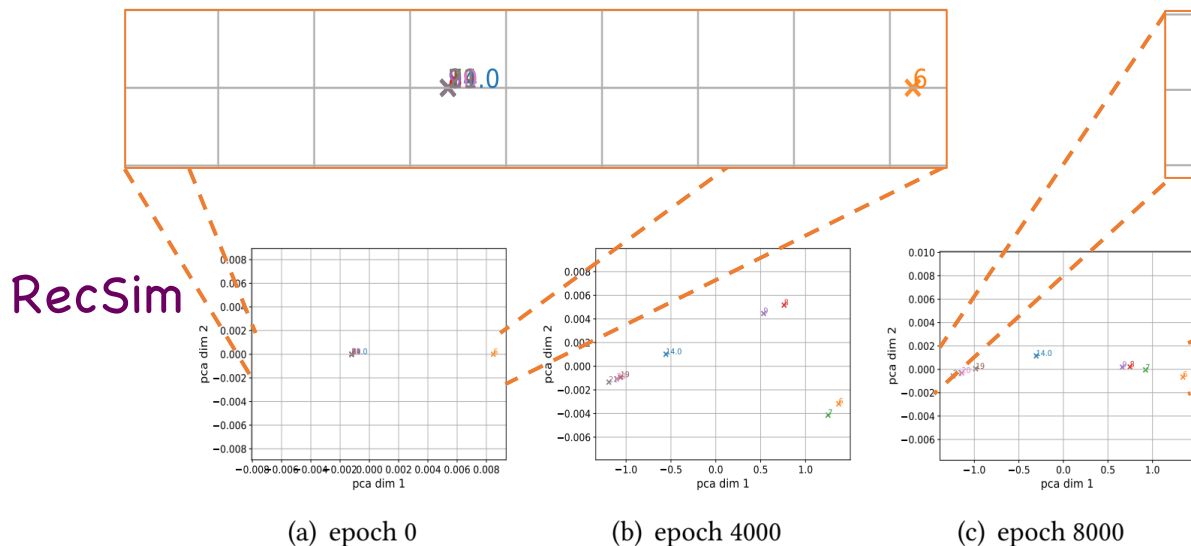


(c) LTS3

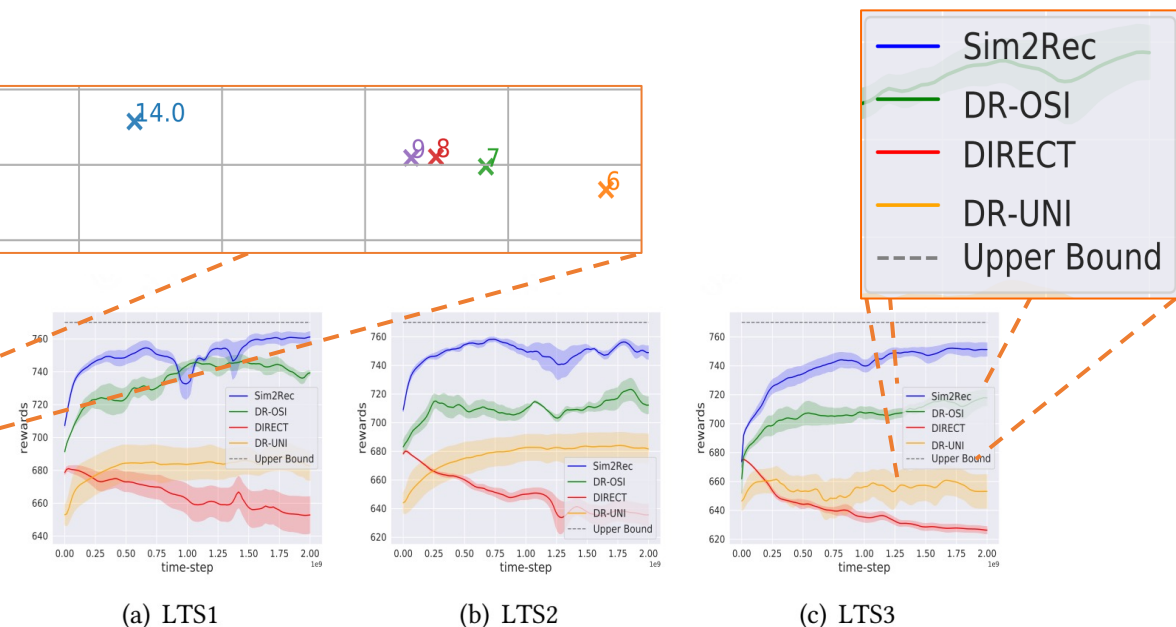
**Figure 2: Illustration of the performance in synthetic environments.** The solid curves are the mean reward and the shadow is the standard error of three seeds.

Real-world  
(Didichuxing)

# Selected Experiment Results



**Figure 1: Illustration of the visualization on  $v$ .** The X-axis denotes the first principal component, and the Y-axis denotes the second one. Each cross point denotes the projection of the latent code for the state distribution. The numbers with the same color to the point denote the ground-truth environment



**Figure 2: Illustration of the performance in synthetic environments.** The solid curves are the mean reward and the shadow is the standard error of three seeds.

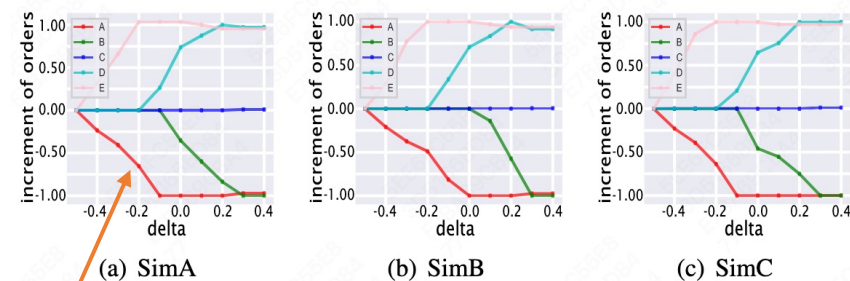
## Real-world (Didichuxing)

**TABLE III**  
THE PERFORMANCE OF POLICIES LEARNED WITH DIFFERENT POLICY LEARNING TECHNIQUES. THE PERFORMANCE IS TESTED IN SIMA.

	orders (test)	orders (train)	cost (test)	cost (train)
Sim2Rec	2.0%	1.6%	0.9%	4.5%
Sim2Rec-PE	1.3%	2.3%	-8.0%	-4.0%
Sim2Rec-EE	8.1%	8.2%	-10.0%	-11.1%

performance decline  
when deployed

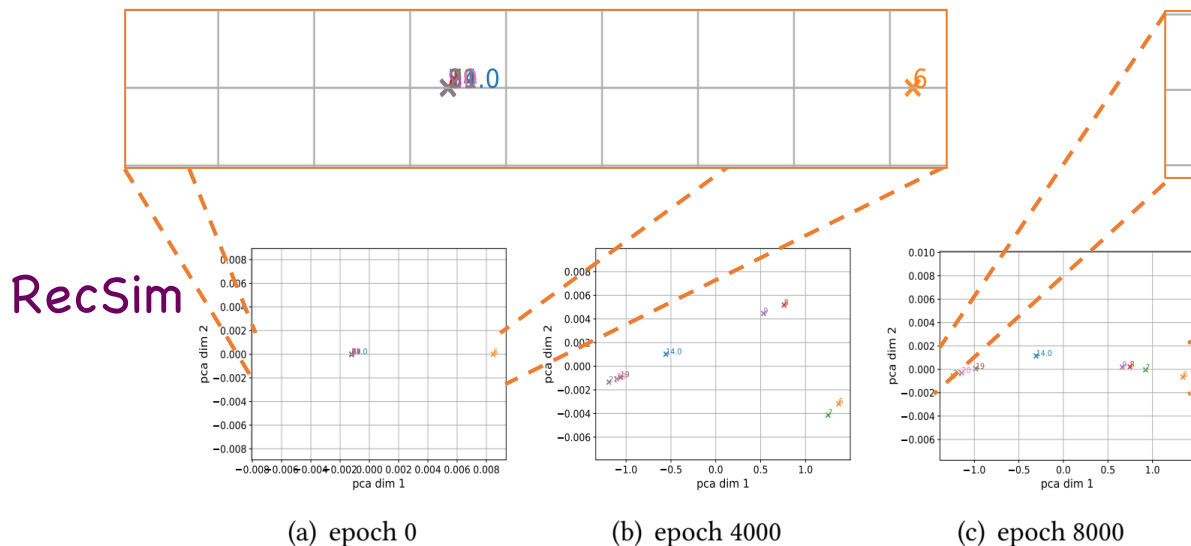
exploit the model to reach  
high performance through unreasonable actions



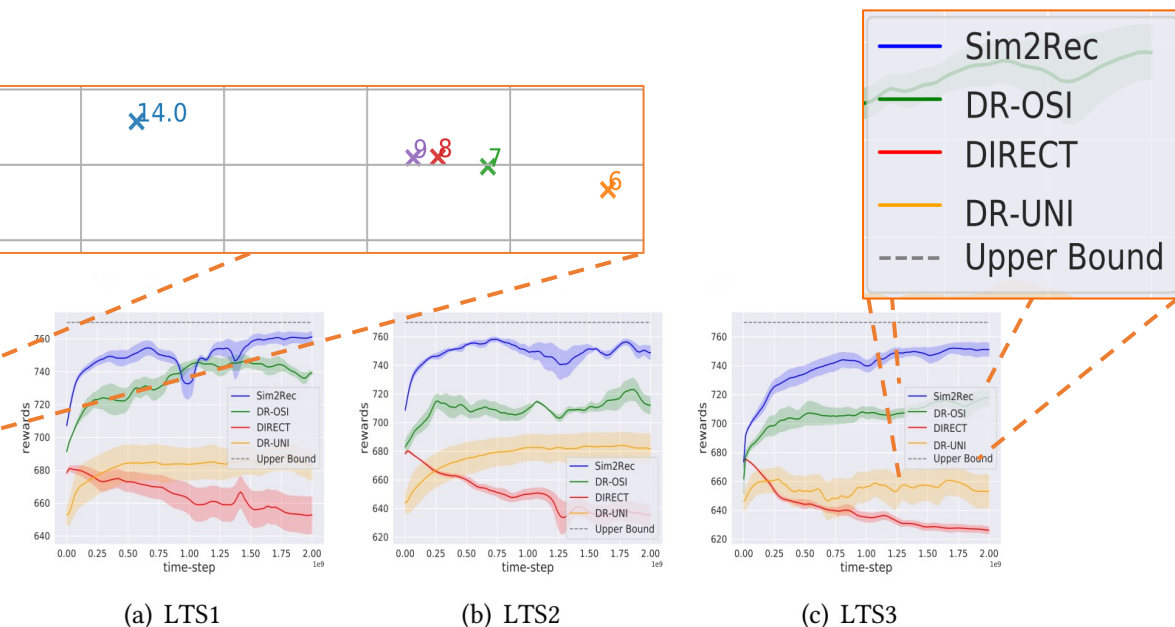
**Fig. 10.** Illustration of the increment of orders on intervention test. Each figure plots the clustering centers of the drivers' response vectors in a simulator. Each line denotes a cluster center. The X-axis is the value of  $\Delta B$ . The increment of orders of each point is subtracted to the value in  $\Delta B = -0.5$  of the corresponding cluster.



# Selected Experiment Results



**Figure 1: Illustration of the visualization on  $v$ .** The X-axis denotes the first principal component, and the Y-axis denotes the second one. Each cross point denotes the projection of the latent code for the state distribution. The numbers with the same color to the point denote the ground-truth environment



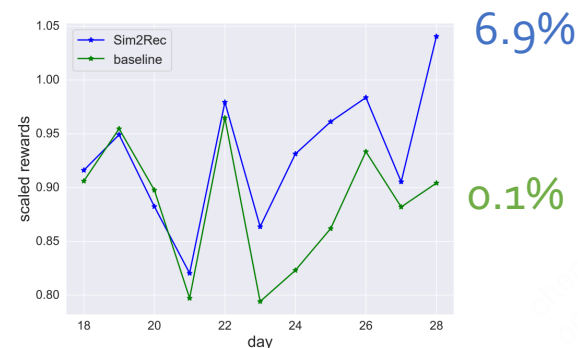
**Figure 2: Illustration of the performance in synthetic environments.** The solid curves are the mean reward and the shadow is the standard error of three seeds.

## Real-world (Didichuxing)

TABLE III

THE PERFORMANCE OF POLICIES LEARNED WITH DIFFERENT POLICY LEARNING TECHNIQUES. THE PERFORMANCE IS TESTED IN SIMA.

	orders (test)	orders (train)	cost (test)	cost (train)
Sim2Rec	2.0%	1.6%	0.9%	4.5%
Sim2Rec-PE	1.3%	2.3%	-8.0%	-4.0%
Sim2Rec-EE	8.1%	8.2%	-10.0%	-11.1%



**Figure 4: Illustration of the online test.** The X-axis is the date. The Y-axis is the average daily reward.

# Take-home Messages & Thanks

## KEY points:

1. The reality-gaps problem is important and should be taken into consideration for designing a power RL-based SRS
2. We identify and give a possible way to handle the several extra challenges in adopting standard zero-shot policy transfer solutions to SRS scenarios. We think more powerful and theoretical methods can be developed in future work

>> Thanks

