



THE **WEB**
CONFERENCE

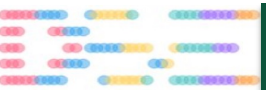
UserSim: User Simulation via Supervised Generative Adversarial Network

Xiangyu Zhao¹, Long Xia², Lixin Zou³

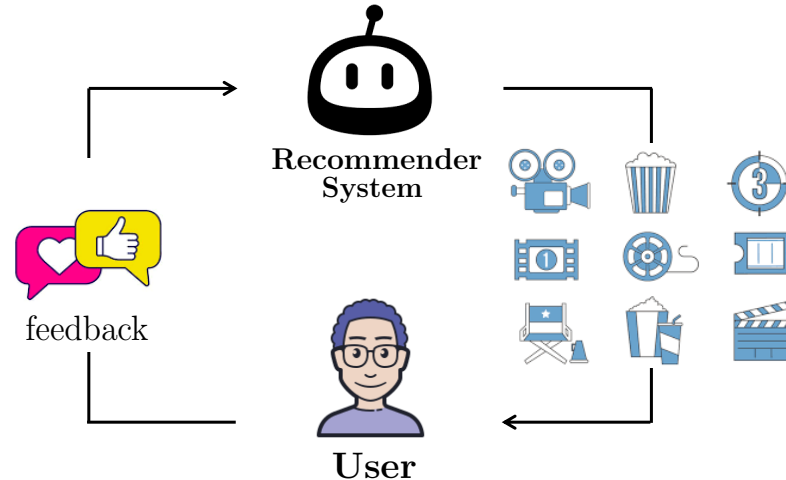
Hui Liu¹, Dawei Yin³, Jiliang Tang¹

1: Data Science and Engineering Lab, Michigan State University

2: York University 3: Baidu

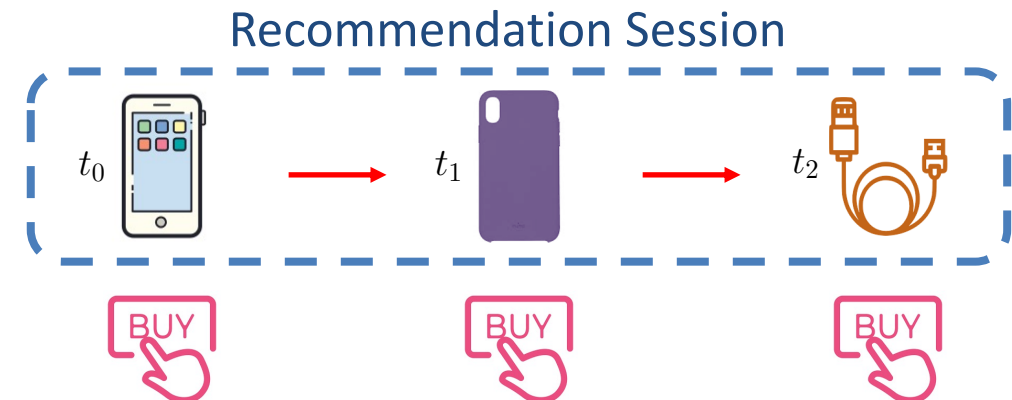
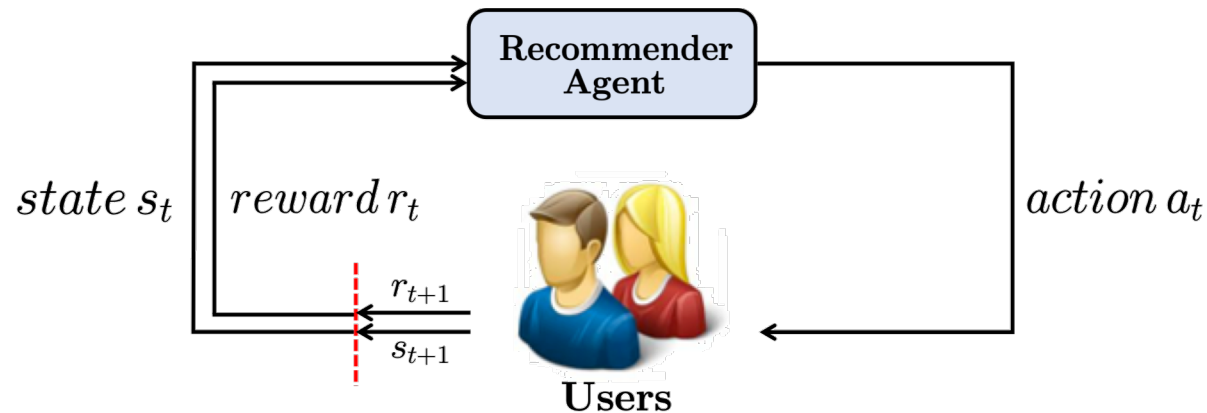


- Increasing interests in applying Reinforcement Learning for recommendations

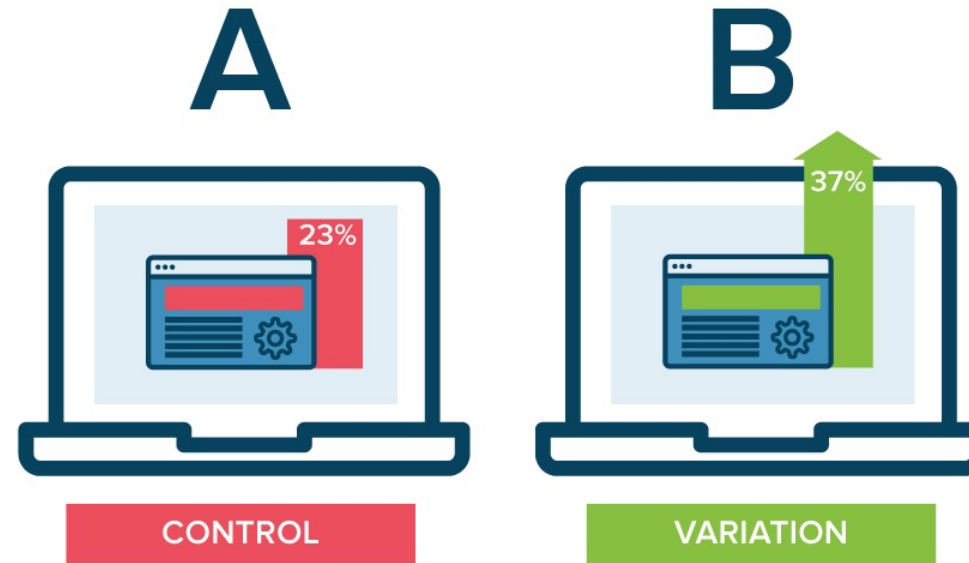


- Advantages

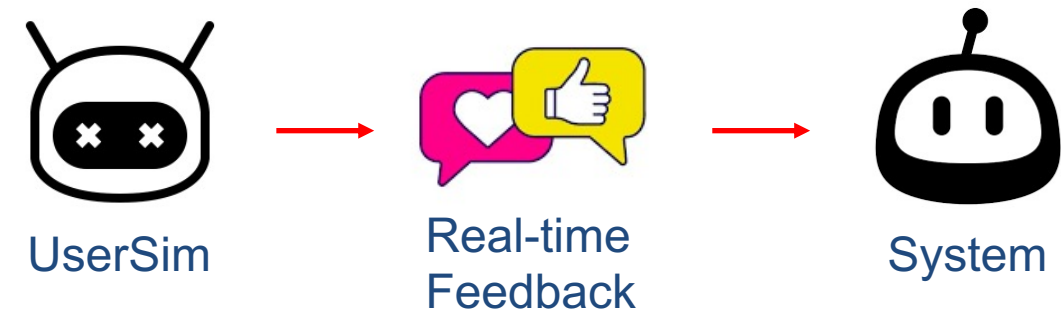
- Continuously updating the recommendation strategies during the interactions
- Maximizing the long-term reward from users



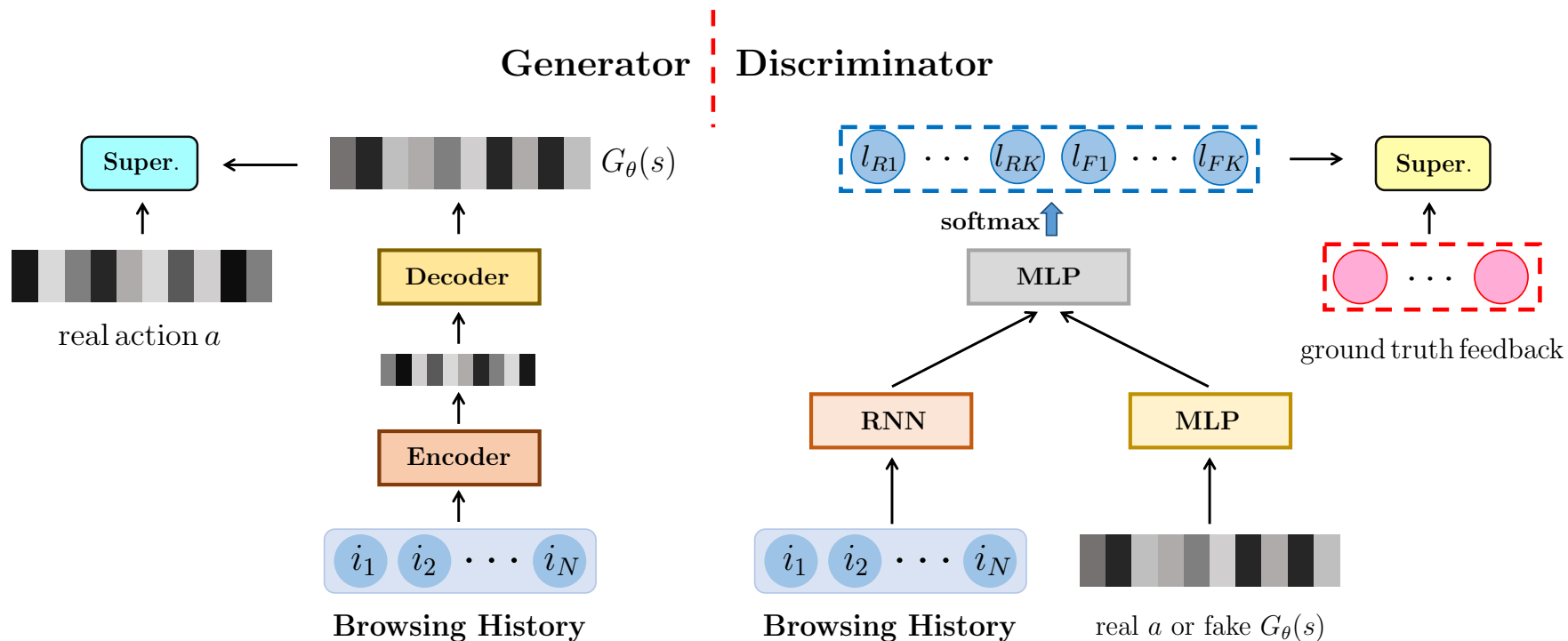
- The most practical and precise way is online A/B test



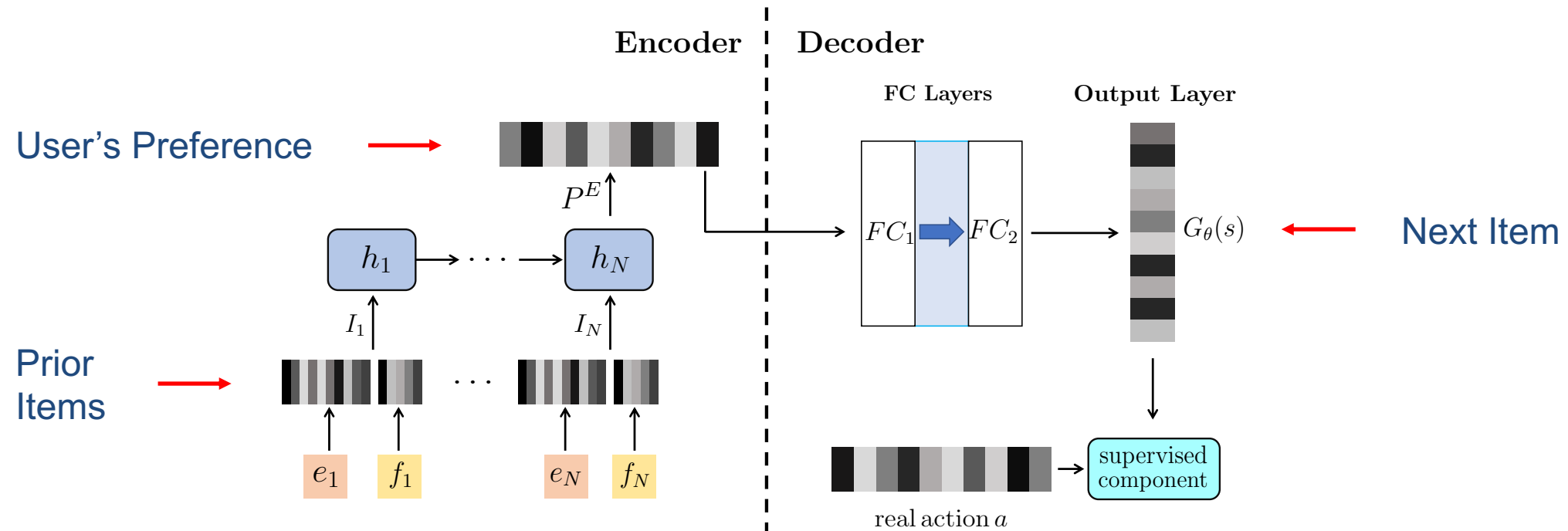
- Online A/B test is inefficient and expensive
 - Taking several weeks to collect sufficient data
 - Numerous engineering efforts
 - Bad user experience



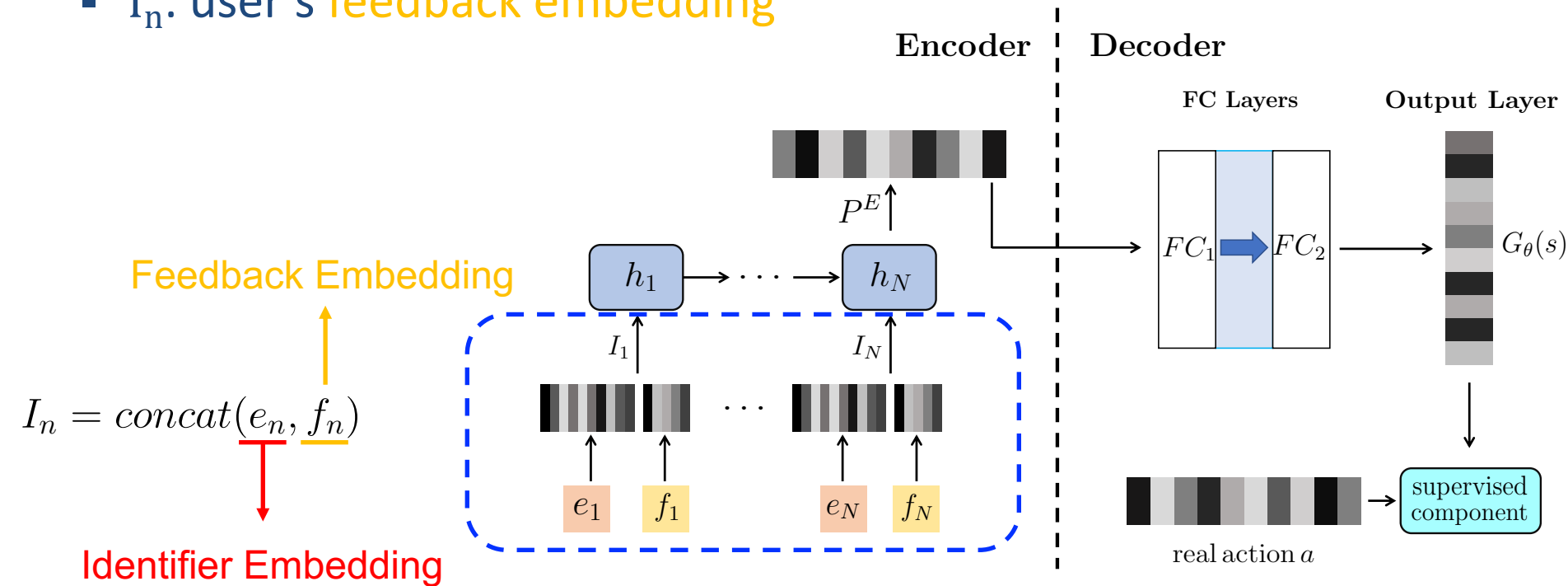
- Simulating users' real-time feedback is challenging
 - Underlying distribution of item sequences is extremely complex
 - Data available to each user is rather limited



- Learning the data distribution
- Generating indistinguishable logs based on users' browsing history

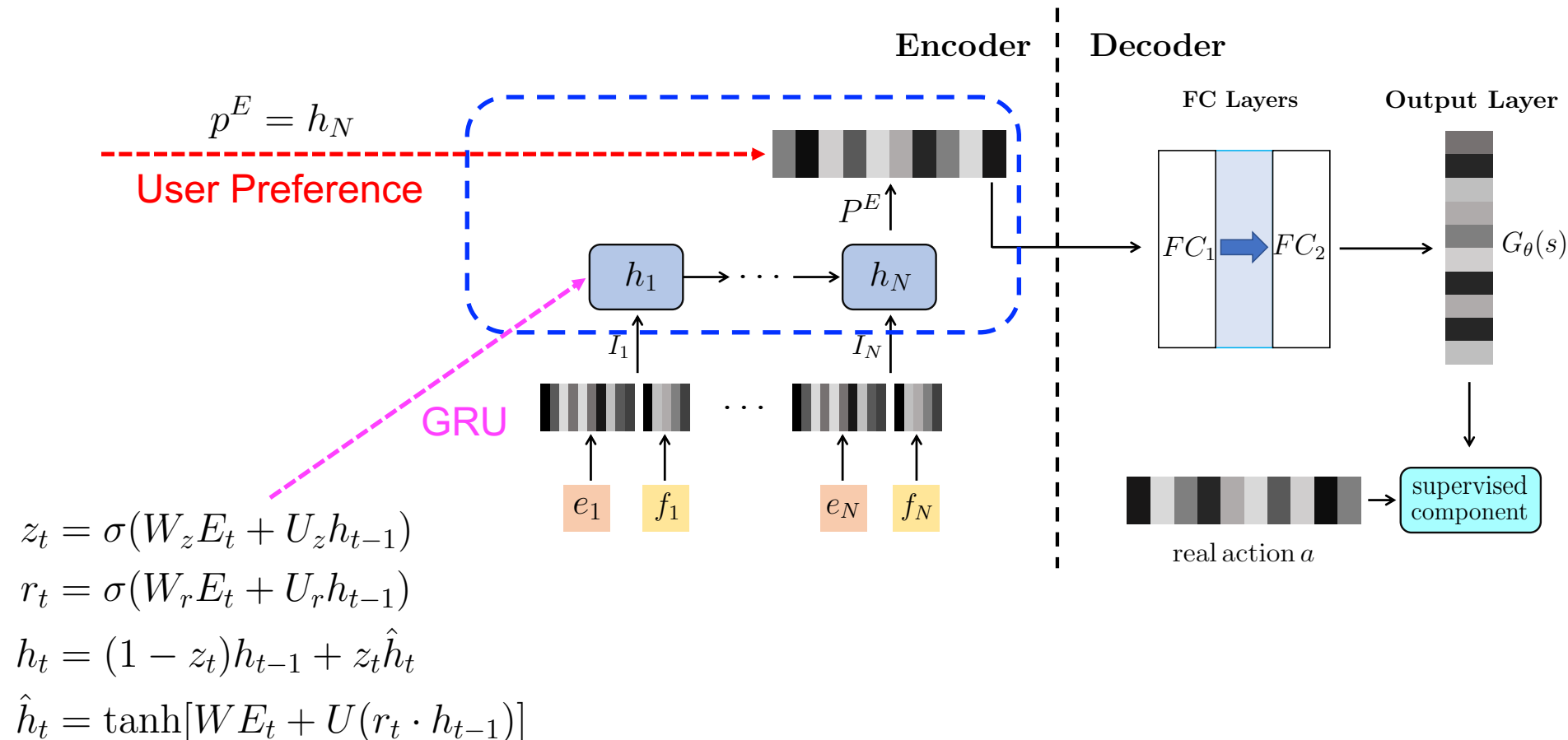


- Input layer
 - e_n : item's identifier embedding
 - f_n : user's feedback embedding

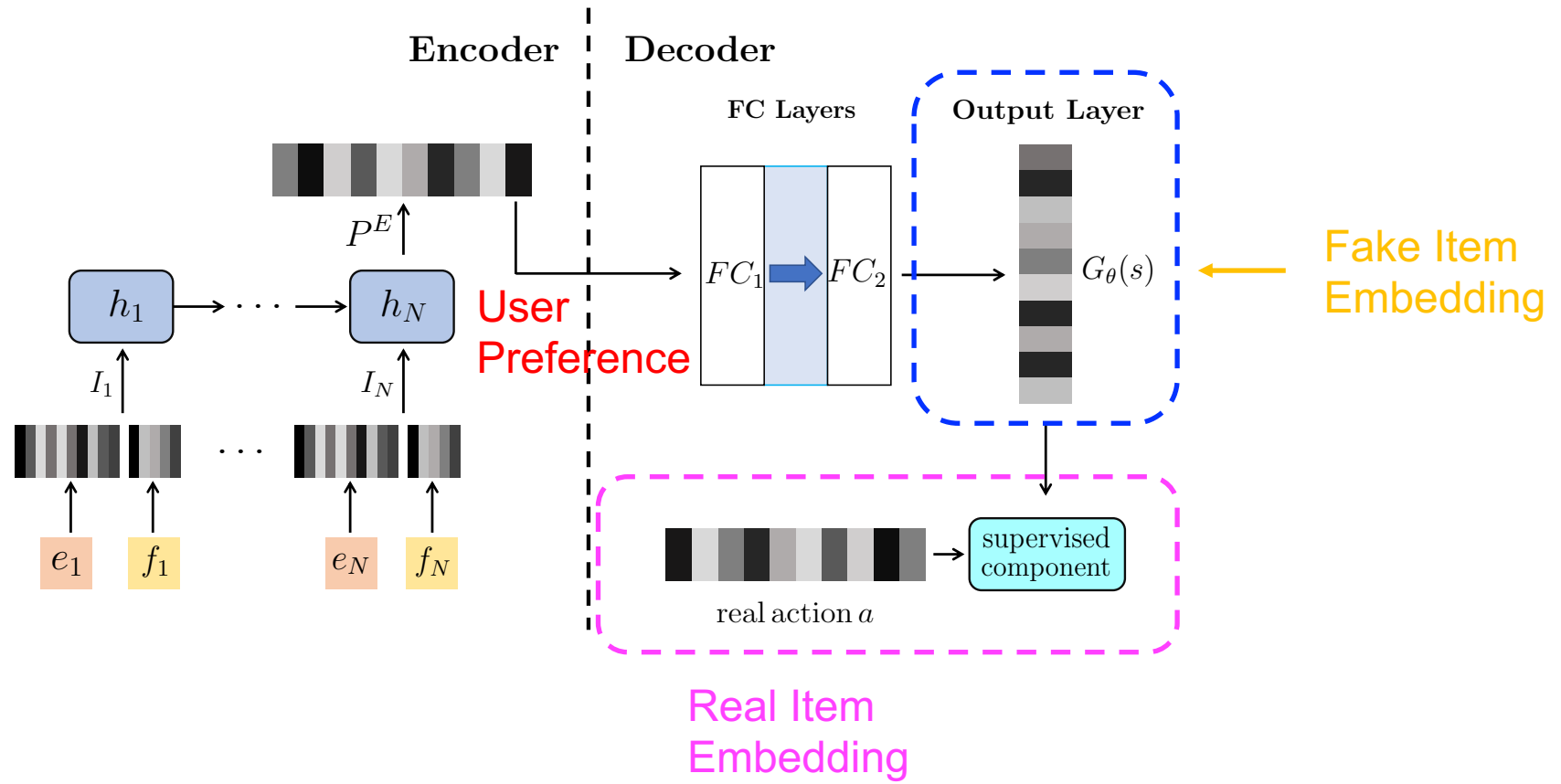


Encoder

- GRU layer:
 - Capturing user's preference from the sequence of items

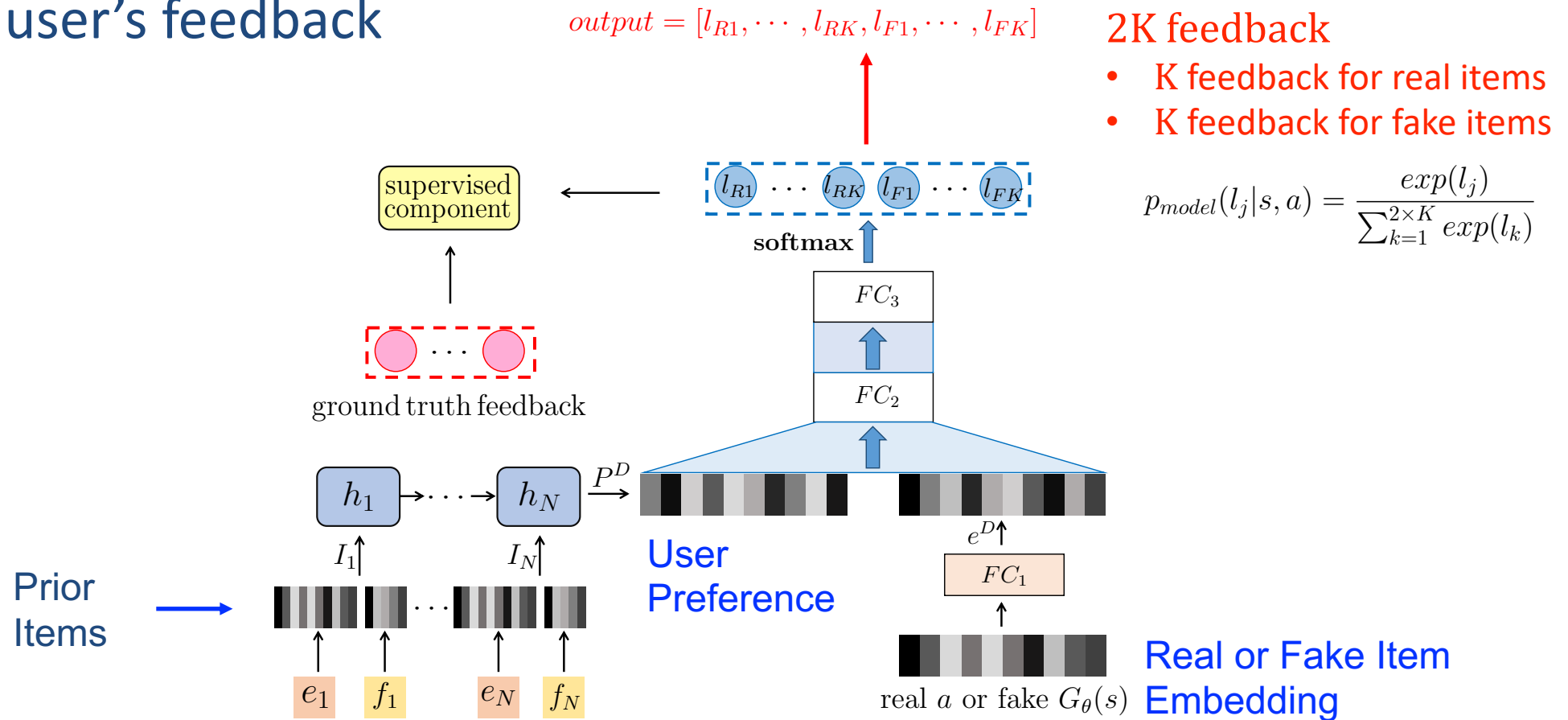


- Goal:
 - Predicting the item to be recommended



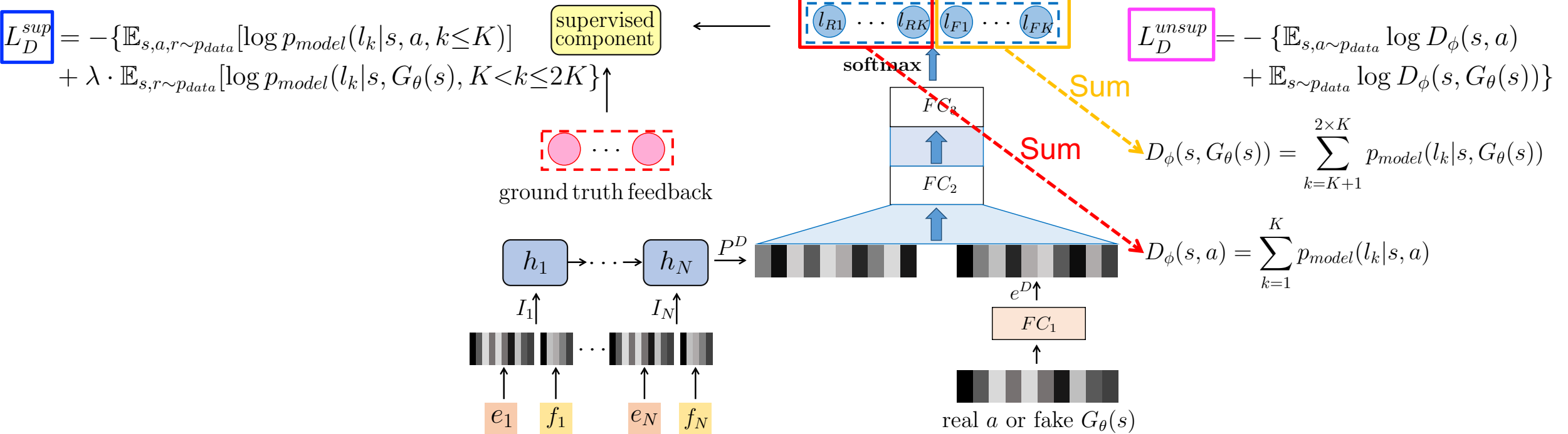
Discriminator

- Distinguishing real/fake items
- Predicting user's feedback

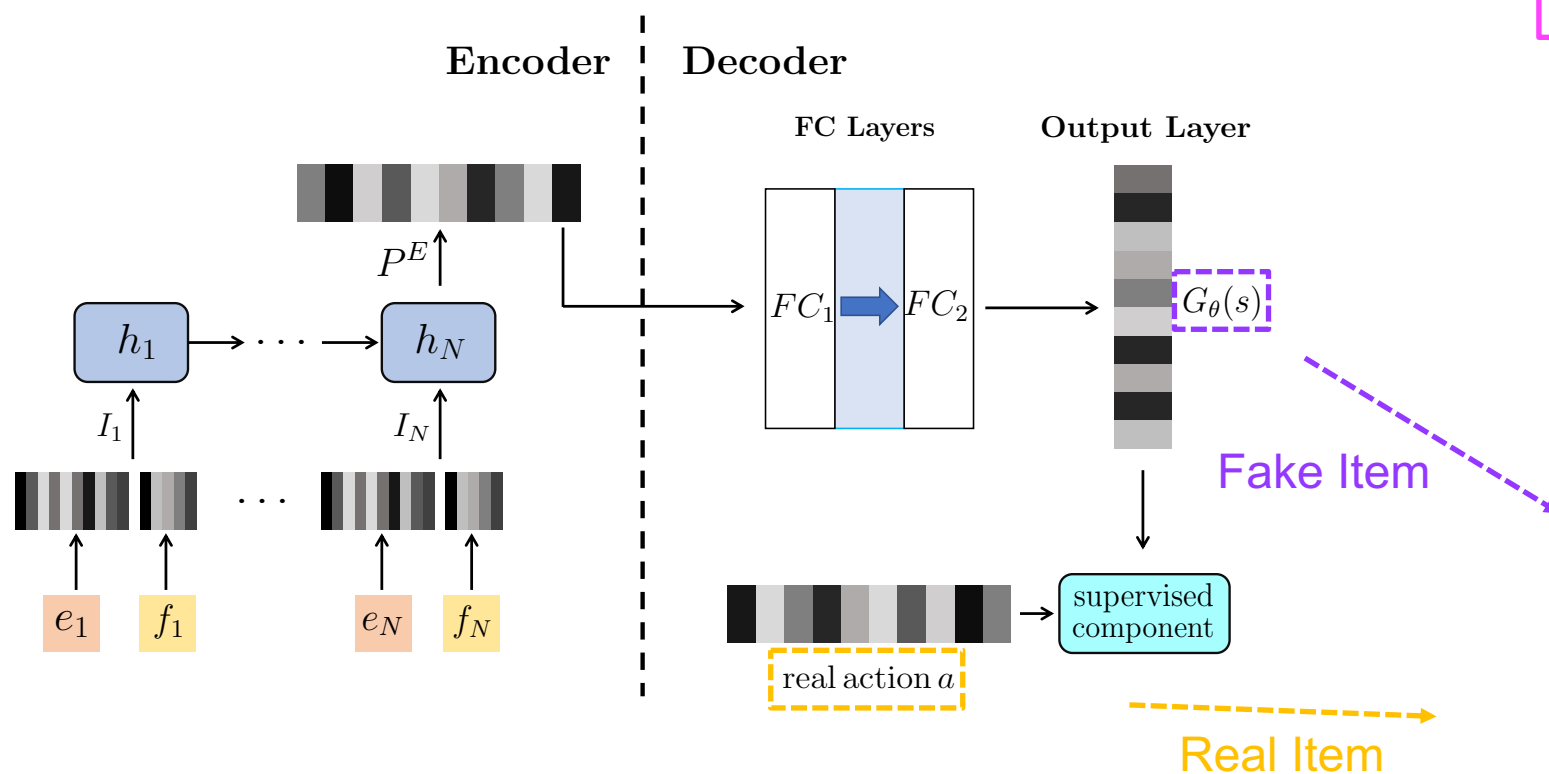


Discriminator

$$L_D = L_D^{unsup} + \alpha \cdot L_D^{sup}$$



Generator



$$L_G^{unsup} = \mathbb{E}_{s \sim p_{data}} [\log D_\phi(s, G_\theta(s))]$$

$$L_G = L_G^{unsup} + \beta \cdot L_G^{sup}$$

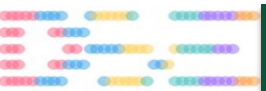
$$L_G^{sup} = \mathbb{E}_{s, a \sim p_{data}} \|a - G_\theta(s)\|_2^2$$

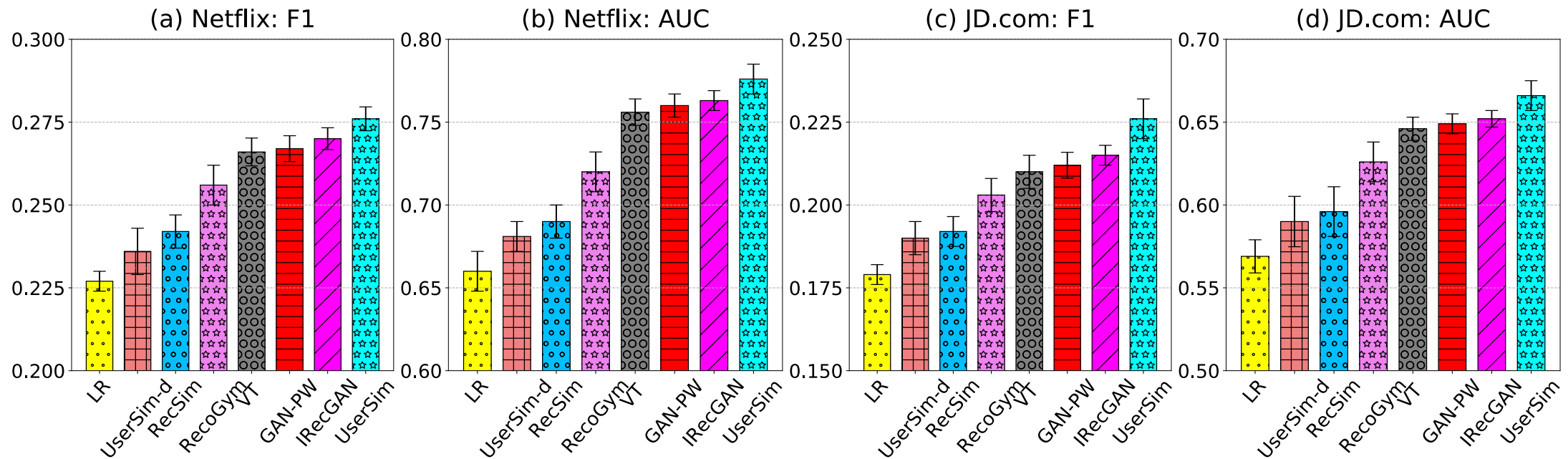
- Public benchmark datasets
 - Netflix and JD.com
 - 70%: training/validation set
 - 30%: test set
- 4 types of feedback
 - Real-positive
 - Real-negative
 - Fake-positive
 - Fake-negative
 - Real: real item from data
 - Fake: fake item from generator

Object	Netflix Prize	JD.com
# user (session)	480,189	283,228
# item	17,770	1,355,255
# interaction	100,480,507	97,713,660
# ave. length	209	345
# feedback	rating 1~5	skip, click

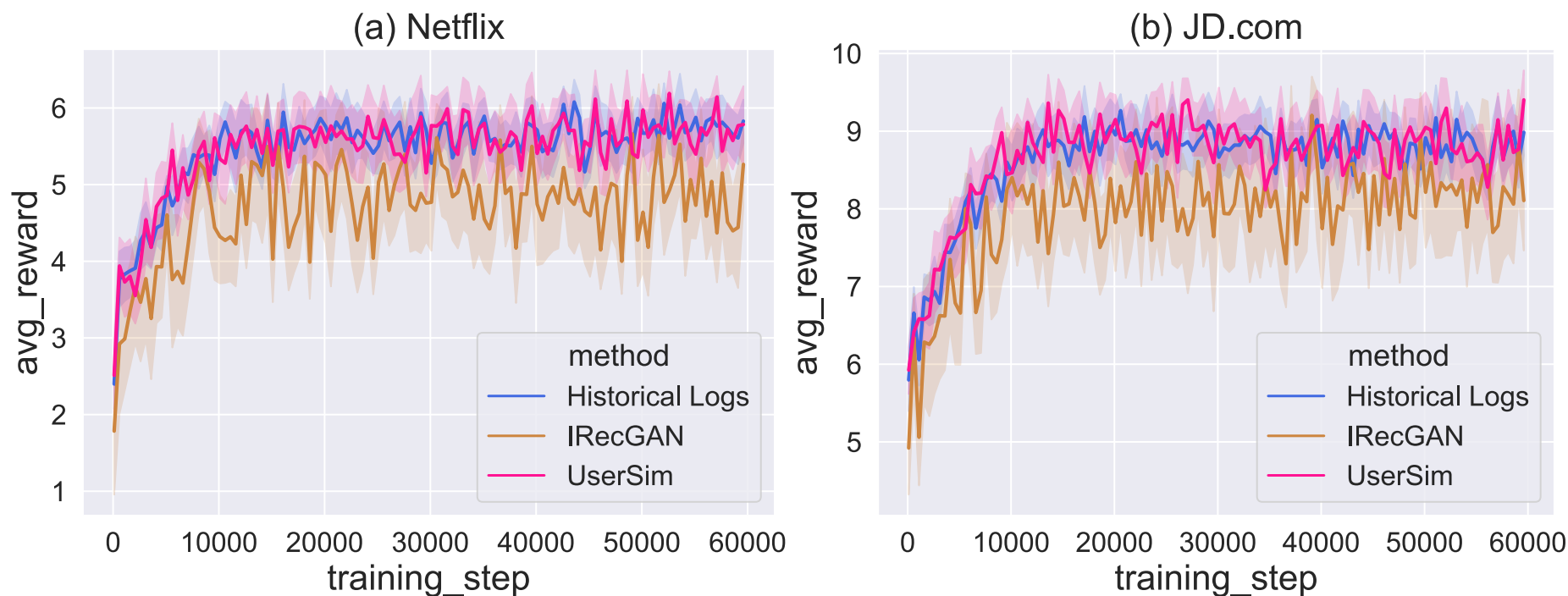
4~5: positive
1~3: negative

click: positive
skip: negative





- Metric: F1-score
- Baselines: LR, UserSim-d, RecSim, RecoGym, Virtual-Taobao, GAN-PW, IRecGAN
- Generator can learn the item distribution, and generate fake items
- Discriminator can distinguish real and fake items, and predict user's feedback



- Metric: average reward of a session
- Baselines: Historical Logs, IRecGAN
- UserSim converges to the similar avg_reward with the one upon historical data
- UserSim performs much more stably than the one trained based upon IRecGAN

On-policy RL algorithms such as SARSA cannot be directly trained on historical data

- We propose a novel user simulator based on Generative Adversarial Network
 - Generating real-time feedback like real users
 - Pre-training and evaluating new recommendation algorithms before launching them online

zhaoxi35@msu.edu

