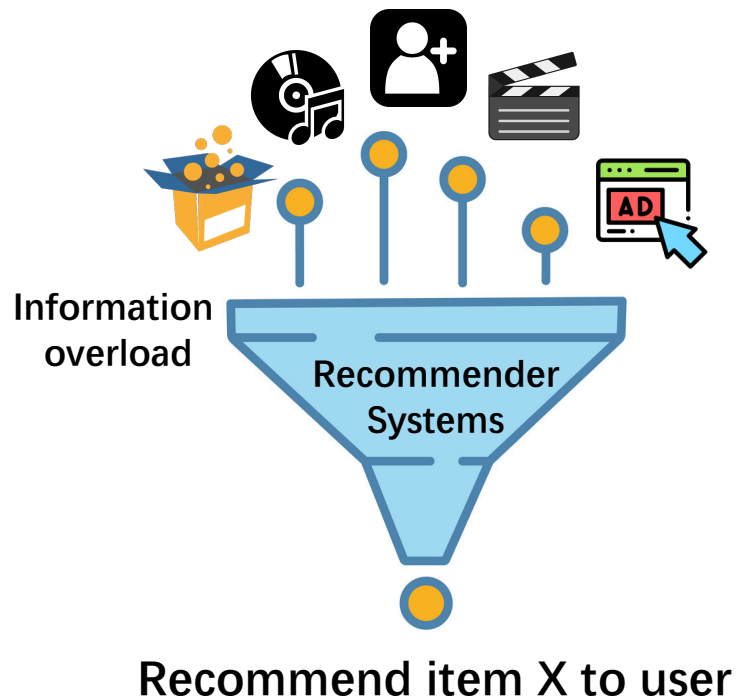


Attacking Black-box Recommendations via Copying Cross-domain User Profiles

Wenqi Fan, Tyler Derr, Xiangyu Zhao, Yao Ma,
Hui Liu, Jianping Wang, Jiliang Tang, and Qing Li

Recommender systems

- Goal: suggest items that best fit users' preferences



amazon

淘宝网
Taobao.com

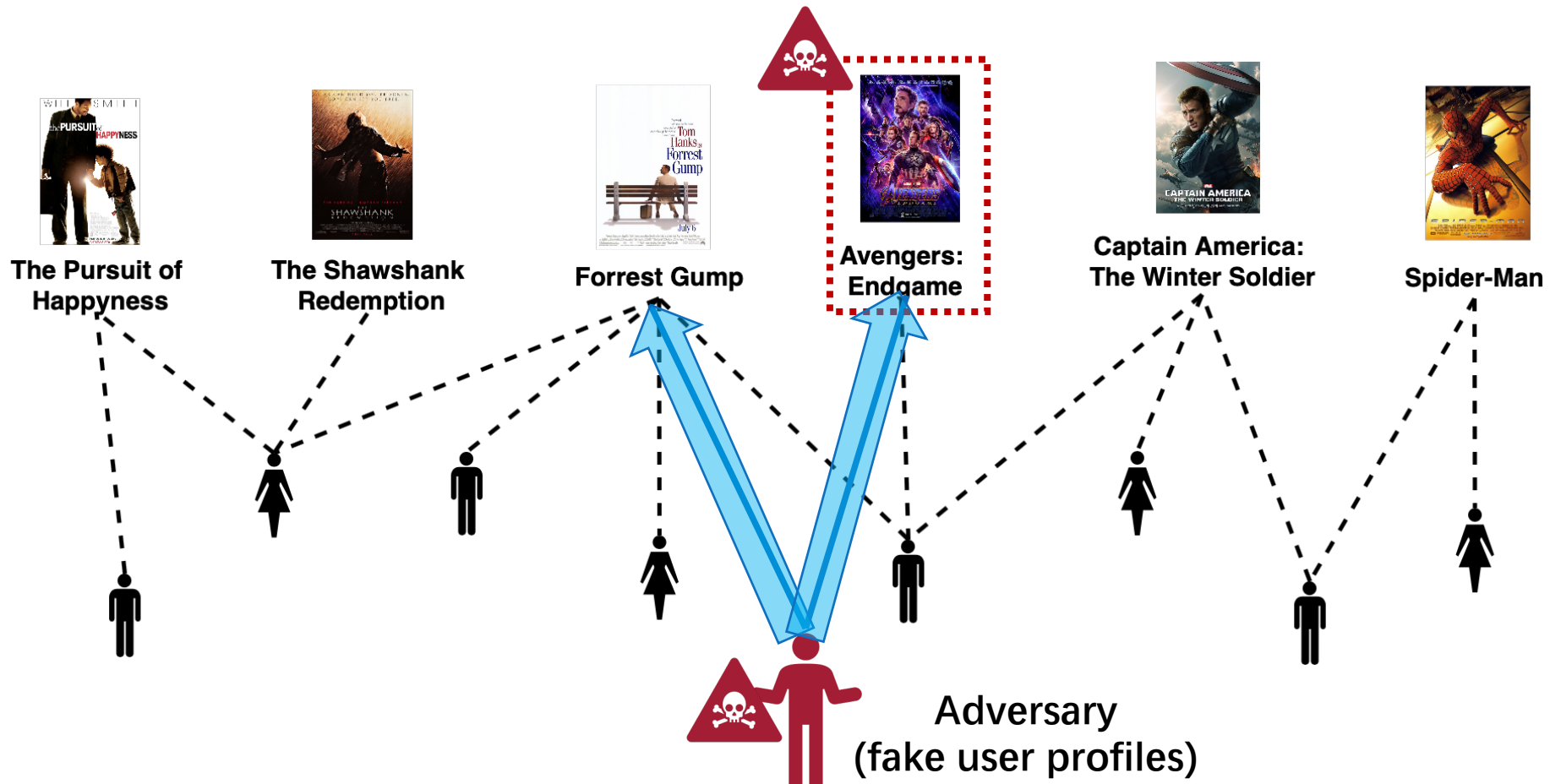
JD.COM

LinkedIn

facebook

Recommender systems

- Security (Attacking) in Recommender Systems
 - **Data poisoning attacks:** promote/demote a set of items



Attacking in recommender systems

- **Challenges in existing attacking methods:**
 - Less "realistic" user profiles (easily detected)

Attacking in recommender systems

- **Cross-domain Information**

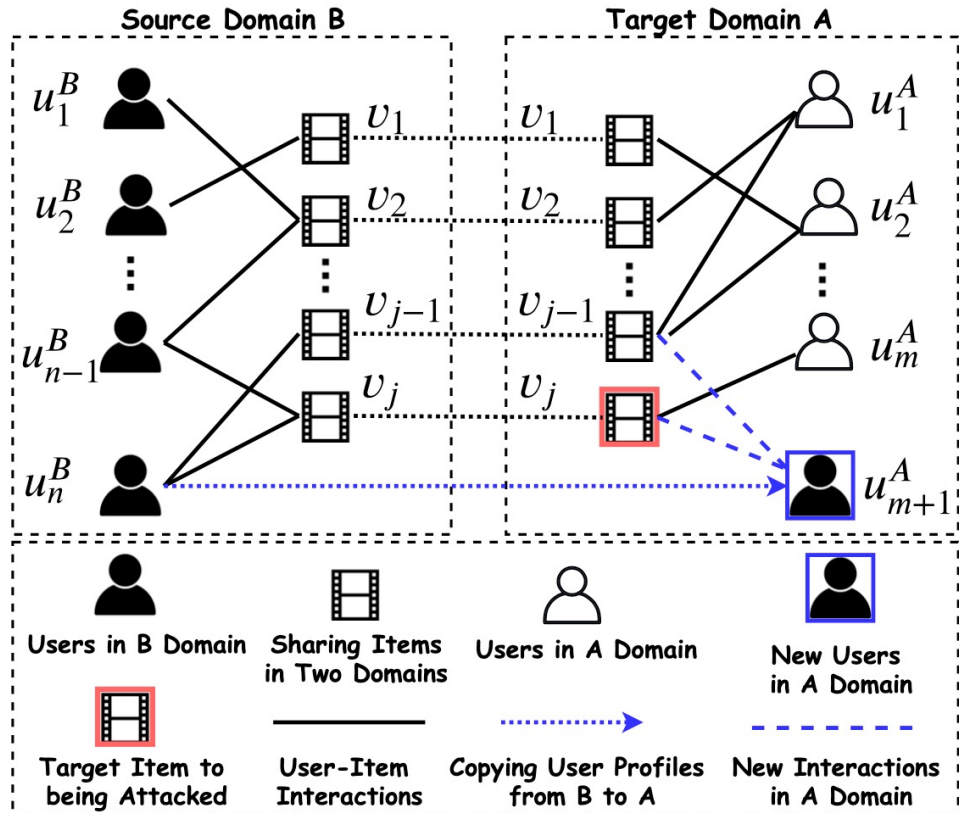
- Share a lot of items
- Users from these platforms with **similar functionalities** also share similar behavior patterns/preferences.

The image displays four e-commerce platforms: Taobao.com, JD.com, Amazon, and Best Buy. Each platform shows listings for the iPhone 12 series. Taobao and JD.com are Chinese sites, while Amazon and Best Buy are US-based. A large 'VS' graphic with lightning bolts is placed between the Chinese and US sites, indicating a comparison or competition. The Amazon and Best Buy listings are also connected by a 'VS' graphic.


Platform	Product Name	Price	Availability
Taobao.com	Apple iPhone 12 (128GB)	¥6799	Available
Taobao.com	Apple iPhone 12 Pro Max (256GB)	¥10099.00	Available
JD.com	Apple iPhone 12 (128GB)	¥6799.00	Available
JD.com	Apple iPhone 12 Pro Max (256GB)	¥10099.00	Available
Amazon	iPhone 12 (128GB)	\$699	Available
Amazon	iPhone 12 Pro (256GB)	\$1099	Available
Amazon	iPhone 12 Pro Max (256GB)	\$1299	Available
Best Buy	iPhone 12 mini (128GB)	\$499	Now available
Best Buy	iPhone 12 (128GB)	\$699	Now available
Best Buy	iPhone 12 Pro (256GB)	\$1099	Now available

Attacking in recommender systems

- Challenges in existing attacking methods:
 - Less "realistic" user profiles (easily detected)
 - **Copy cross-domain users with real profiles from other domains**



Attacking in recommender systems

- **Challenges in existing attacking methods:**
 - **Less "realistic" user profiles (easily detected)**
 -  Cross-domain Information
 - **White-box setting** (i.e., model architecture and parameters, and datasets)
 - impossible and unrealistic (**privacy and security**)
 - **Black-box setting**
 - Reinforcement Learning (RL) -- Query Feedback (Reward)

CopyAttack

• Problem Statement

- Target RecSys A Users: $\mathcal{U}^A = \{u_1^A, u_2^A, \dots, u_{n^A}^A\}$
 - User profile: $P_{u_i^A}^A = \{v_1 \rightarrow \dots \rightarrow v_j \rightarrow \dots \rightarrow v_l\}$
 - Item profile: $P_{v_j^A}^A = \{u_1^A \rightarrow \dots \rightarrow u_i^A \rightarrow \dots \rightarrow u_o^A\}$

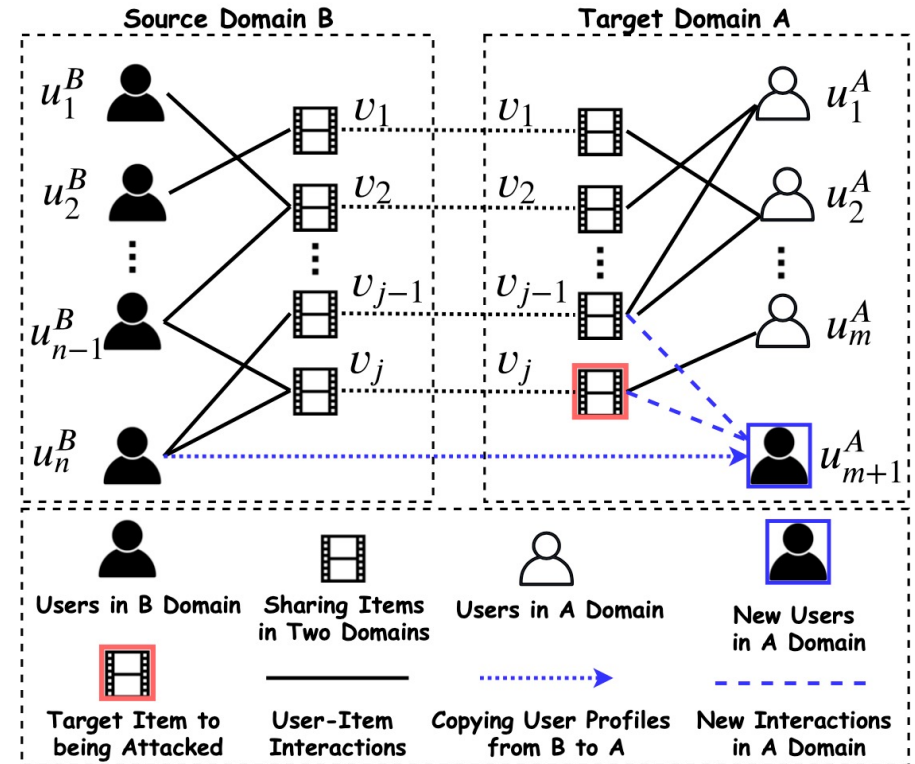
Items: $\mathcal{V}^A = \{v_1, v_2, \dots, v_{m^A}\}$

- Source RecSys B: Users: \mathcal{U}^B Items: \mathcal{V}^B
 - User Profile: $\mathcal{P}_{u_i^B}^B = \{v_1 \rightarrow \dots \rightarrow v_j \rightarrow \dots \rightarrow v_l\}$

- Overlapping items: $\mathcal{V} = \mathcal{V}^A \cap \mathcal{V}^B$

- Goal: $\mathcal{U}^{A'} = \mathcal{U}^A \cup \mathcal{U}^{B \rightarrow A}$

$$y_{i, >k}^A = \{v_{[1]}, v_{[2]}, \dots, v_{[k]}\} = \text{Rec}(P_{u_i^A}^A, \mathcal{P}_{\mathcal{V}^A}^A)$$



CopyAttack

- **Attacking RL Environment**

- Action A: user profiles in source domain B

- Reward R (Hitting Ratio, HR):

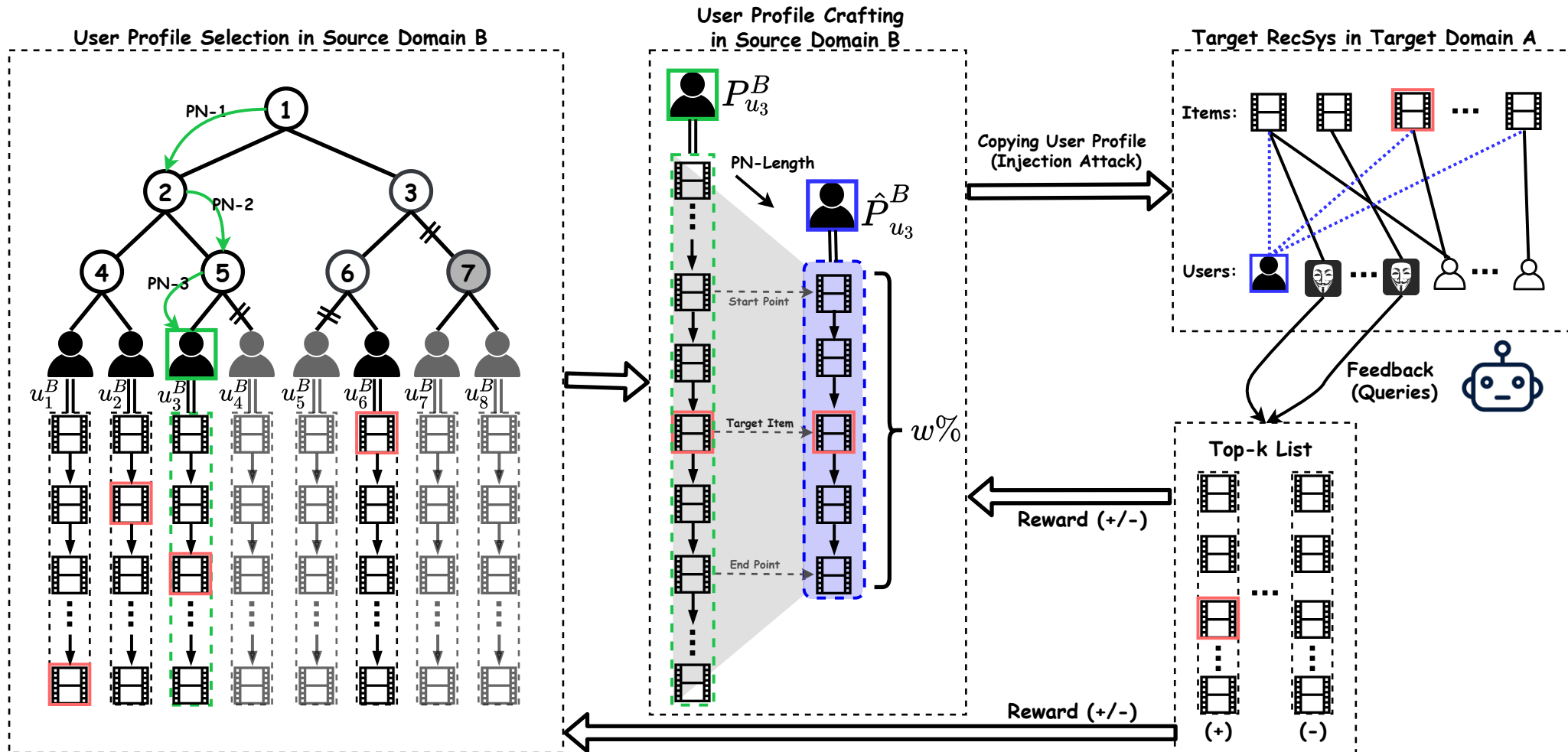
- **Spy users**

$$r(s_t, a_t) = \frac{1}{|\mathcal{U}_*^A|} \sum_{i=1}^{|\mathcal{U}_*^A|} HR(u_{i*}^A, v_*, k)$$

$$HR(u_{i*}^A, v_*, k) = \begin{cases} 1, & v_* \in y_{u^*, > k}, \\ 0, & v_* \notin y_{u^*, > k} \end{cases}$$

- Terminal: reach the budget or successfully satisfy the promotion task

CopyAttack



CopyAttack

- User Profile Selection

- Construct hierarchical clustering tree
- **Masking Mechanism** - specific target items
- Hierarchical-structure Policy Gradient

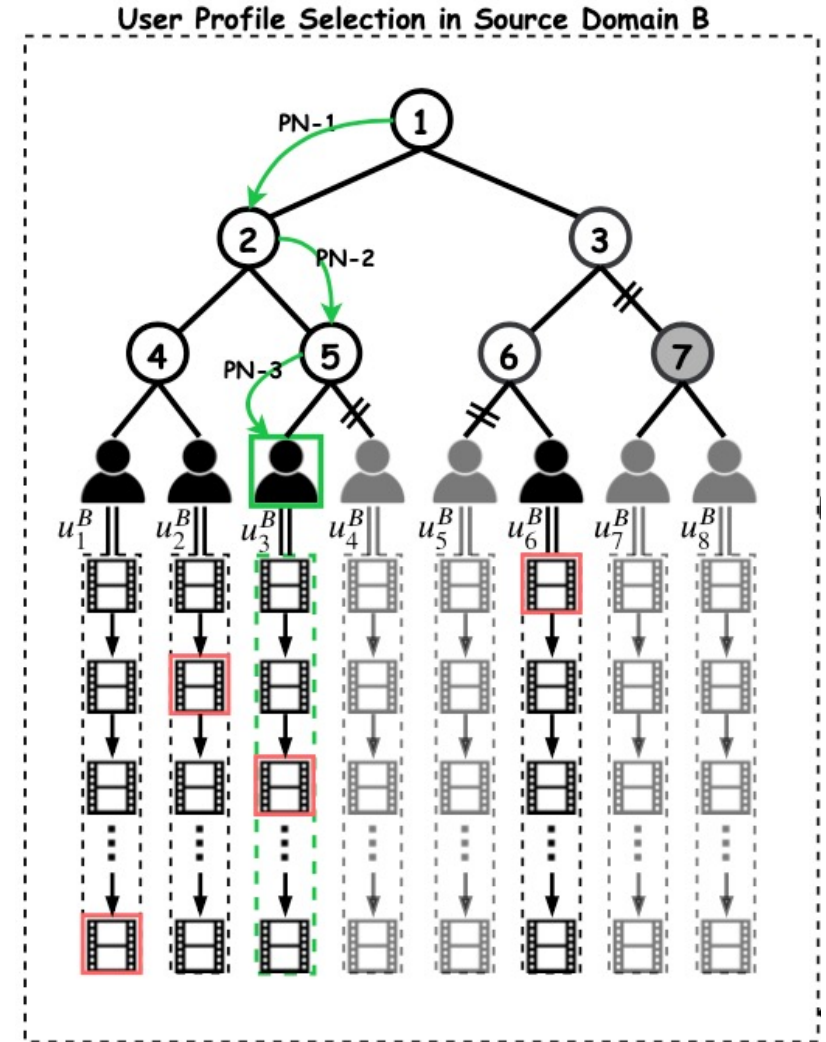
$$a_t^u = \{a_{[t,1]}^u, a_{[t,2]}^u, \dots, a_{[t,d]}^u\}$$

$$\begin{aligned} p^u(a_t^u | s_t^u) &= \prod_{d=1}^d p_d^u(a_{[t,d]}^u | \cdot, s_t^u) \\ &= p_d^u(a_{[t,d]}^u | s_t^u) \cdot p_{d-1}^u(a_{[t,d-1]}^u | s_t^u) \cdots p_1^u(a_{[t,1]}^u | s_t^u) \end{aligned}$$

$$\mathbf{x}_{v_*} = RNN(\mathcal{U}_t^{B \rightarrow A}),$$

$$p_i^u(\cdot | s_t^u) = \text{softmax}(MLP([\mathbf{q}_{v_*}^B \oplus \mathbf{x}_{v_*}] | \theta_i^u))$$

Time Complexity : $\mathcal{O}(|\mathcal{U}^B|) \longrightarrow \mathcal{O}(d \times |\mathcal{U}^B|^{1/d})$



CopyAttack

- User Profile Crafting

- Clipping operation to craft the raw user profiles

$$W = \{10\%, 20\%, 30\%, 40\%, 50\%, 60\%, 70\%, 80\%, 90\%, 100\%\}$$

- Sequential patterns (forward/backward)

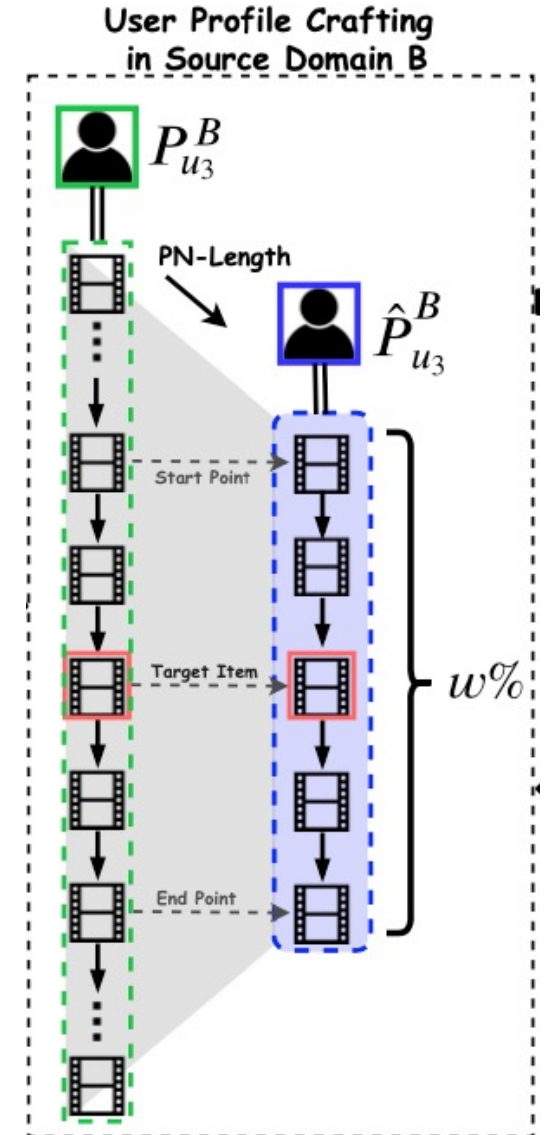
Example:

$$P_{u_i}^B = \{v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow v_4 \rightarrow v_{5*} \rightarrow v_6 \rightarrow v_7 \rightarrow v_8 \rightarrow v_9 \rightarrow v_{10}\}$$

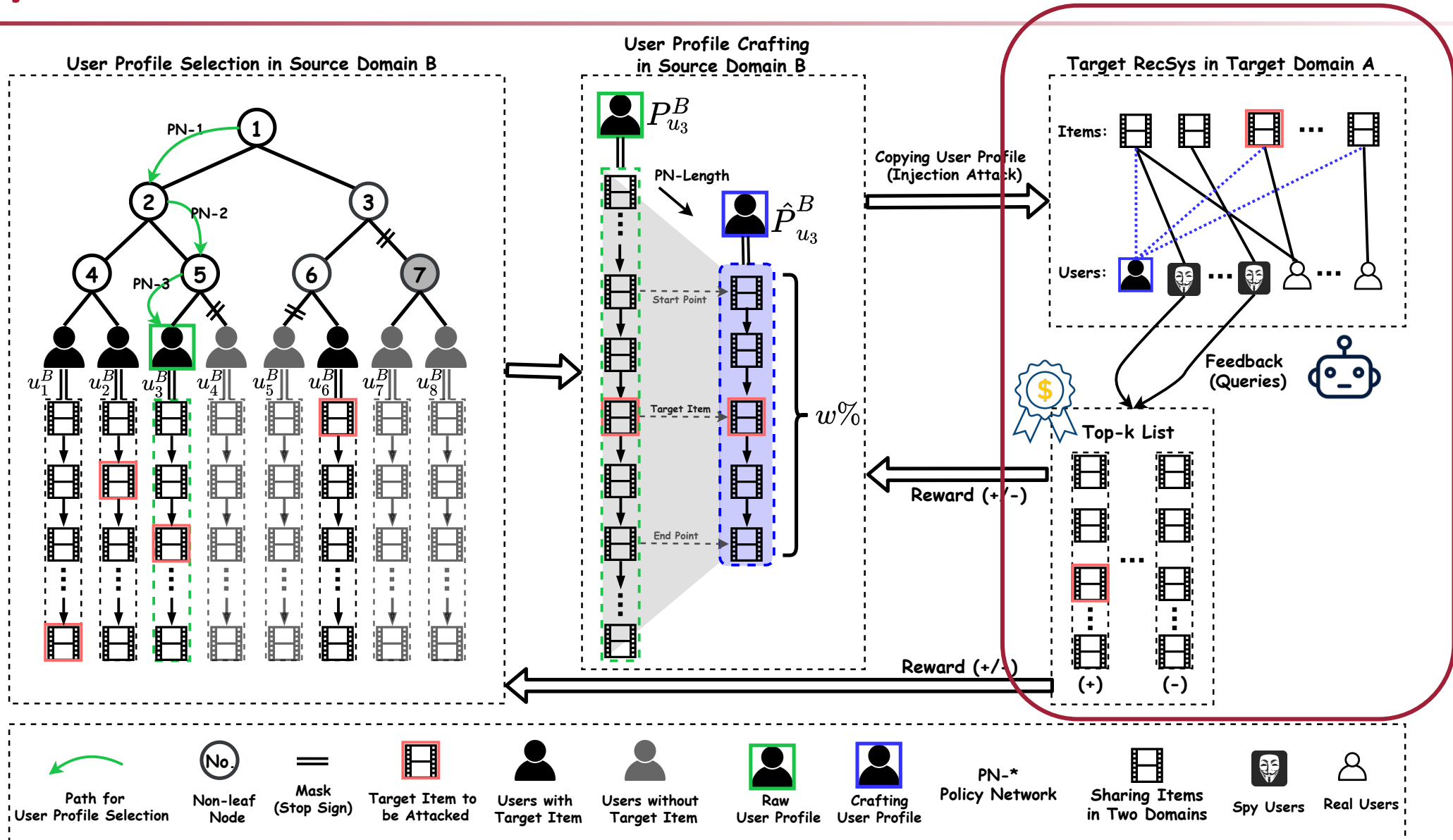
w = 50%

$$\hat{P}_{u_i}^B = \{v_3 \rightarrow v_4 \rightarrow v_{5*} \rightarrow v_6 \rightarrow v_7\}$$

$$p^l(\cdot | s_t^l) = \text{softmax}(\text{MLP}([\mathbf{p}_i^B \oplus \mathbf{q}_{v_*}^B] | \theta^l))$$



CopyAttack



Thank You

Wenqi Fan, wenqifan03@gmail.com

Please see my homepage for more details:
<https://wenqifan03.github.io>

Attacking Black-box Recommendations via Copying Cross-domain User Profiles

Wenqi Fan, Tyler Derr, Xiangyu Zhao, Yao Ma,
Hui Liu, Jianping Wang, Jiliang Tang, and Qing Li