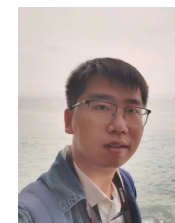
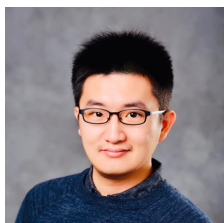


Trustworthy Recommender Systems



Wenqi Fan¹, Xiangyu Zhao², Lin Wang¹, Xiao Chen¹, Jingtong Gao², Qidong Liu², Shijie Wang¹

¹The Hong Kong Polytechnic University

²City University of Hong Kong

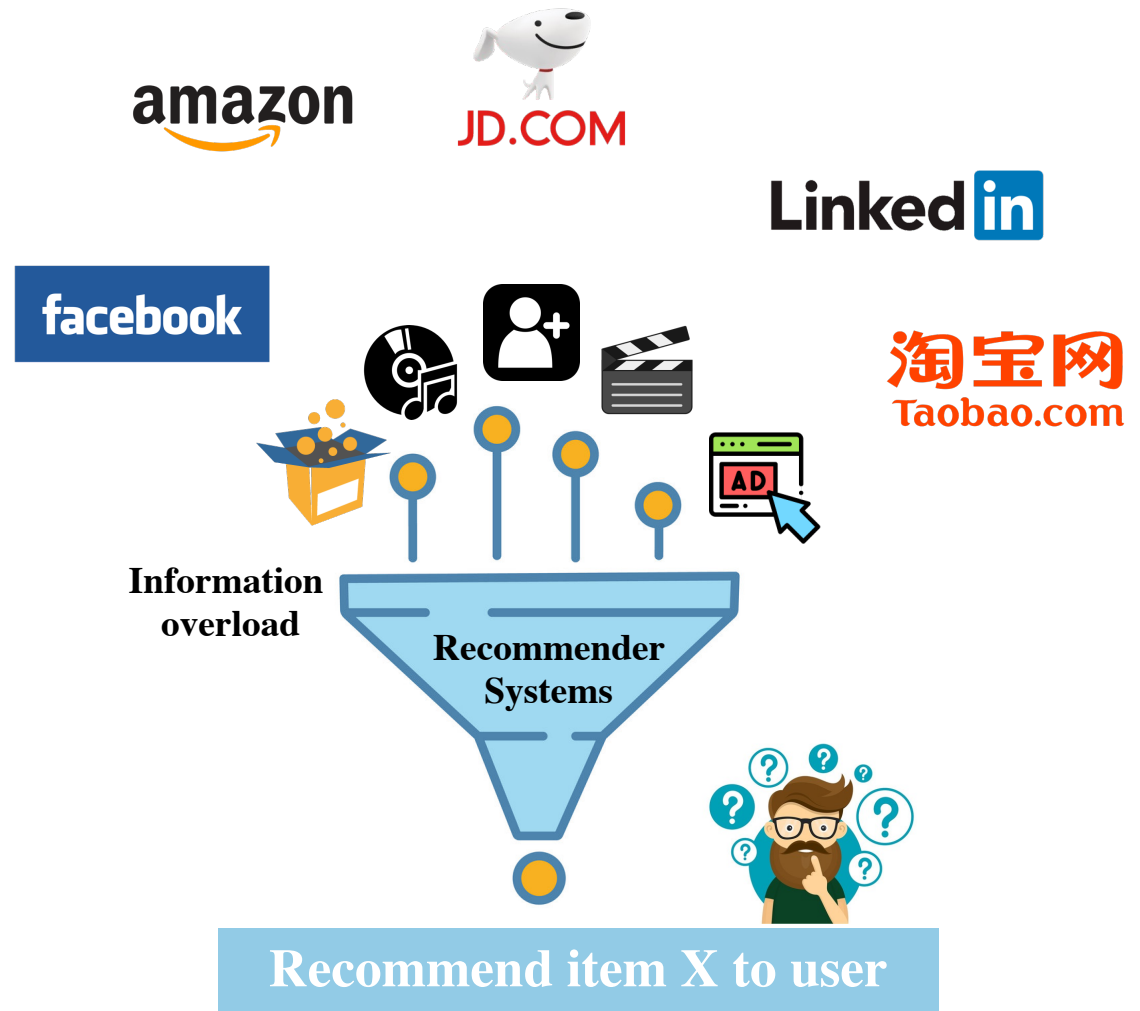


Website (Slides): <https://advanced-recommender-systems.github.io/trustworthiness-tutorial/>

Survey: A Comprehensive Survey on Trustworthy Recommender Systems, arXiv:2209.10117, 2022.

Recommender Systems

Age of Information Explosion



Items can be: Products, Friends, News, Movies, Videos, etc.

Recommender Systems

Recommendation has been widely applied in online services:

- E-commerce, Content Sharing, Social Networking ...



Product Recommendation

Frequently bought together



Total price: \$208.9

Add all three to Cart

Add all three to List

Recommender Systems

Recommendation has been widely applied in online services:

- E-commerce, **Content Sharing**, Social Networking ...

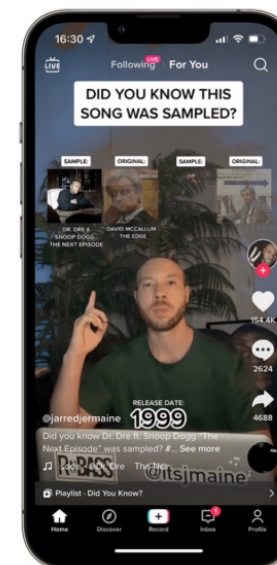
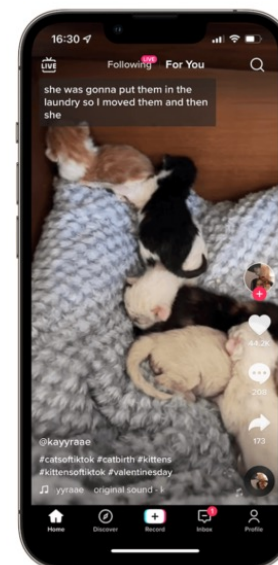
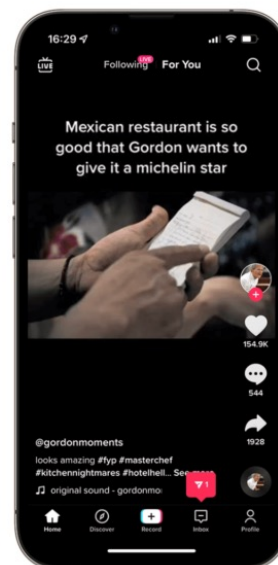


News/Video/Image Recommendation

MIT
Technology
Review

Top 10 Global Breakthrough
Technologies in 2021

TikTok's recommendation algorithm



Recommender Systems

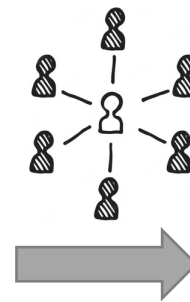
Recommendation has been widely applied in online services:

- E-commerce, Content Sharing, **Social Networking** ...

facebook

新浪微博
weibo.com

LinkedIn



Social Recommendations

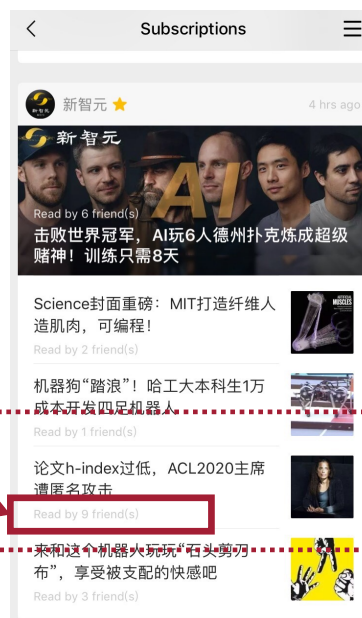


Like what you read?
Share with your friends!



Subscriptions
(訂閲號信息)

Read by 9
friends



Top Stories (看一看)
Wow (朋友在看)



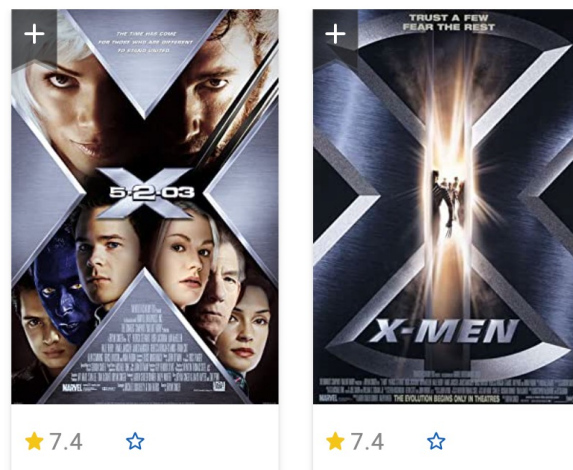
Recommender System is Everywhere



Business



Healthcare



Entertainment



Education

The Good and The Bad

The Good



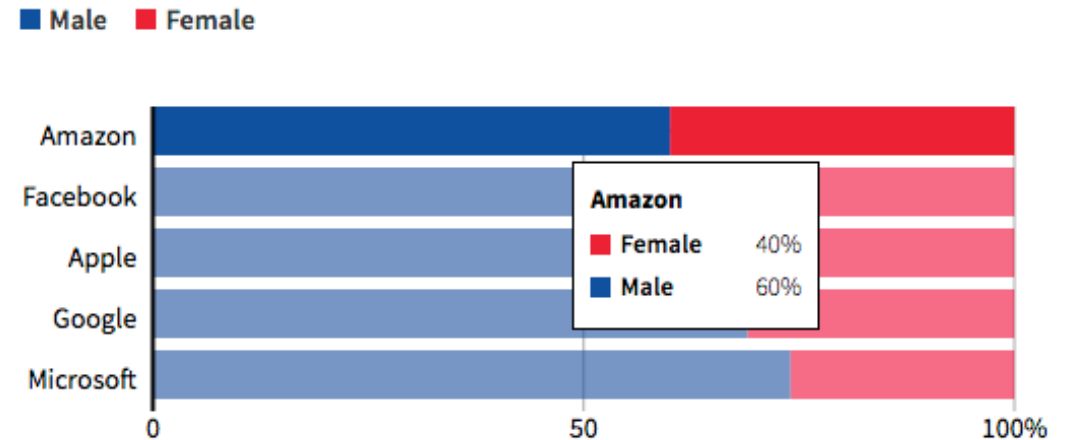
The Bad

Discrimination & Fairness Issue



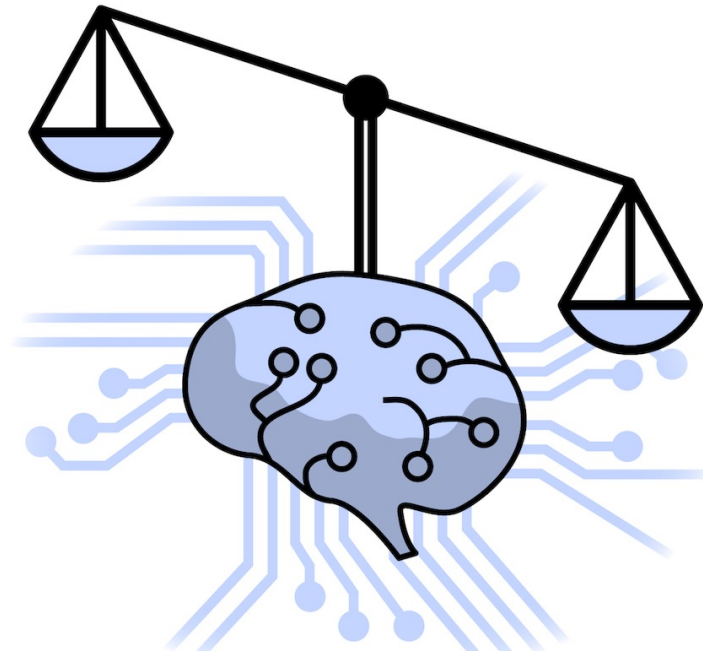
Job recommendation
(Lambrecht et al., 2019)

GLOBAL HEADCOUNT

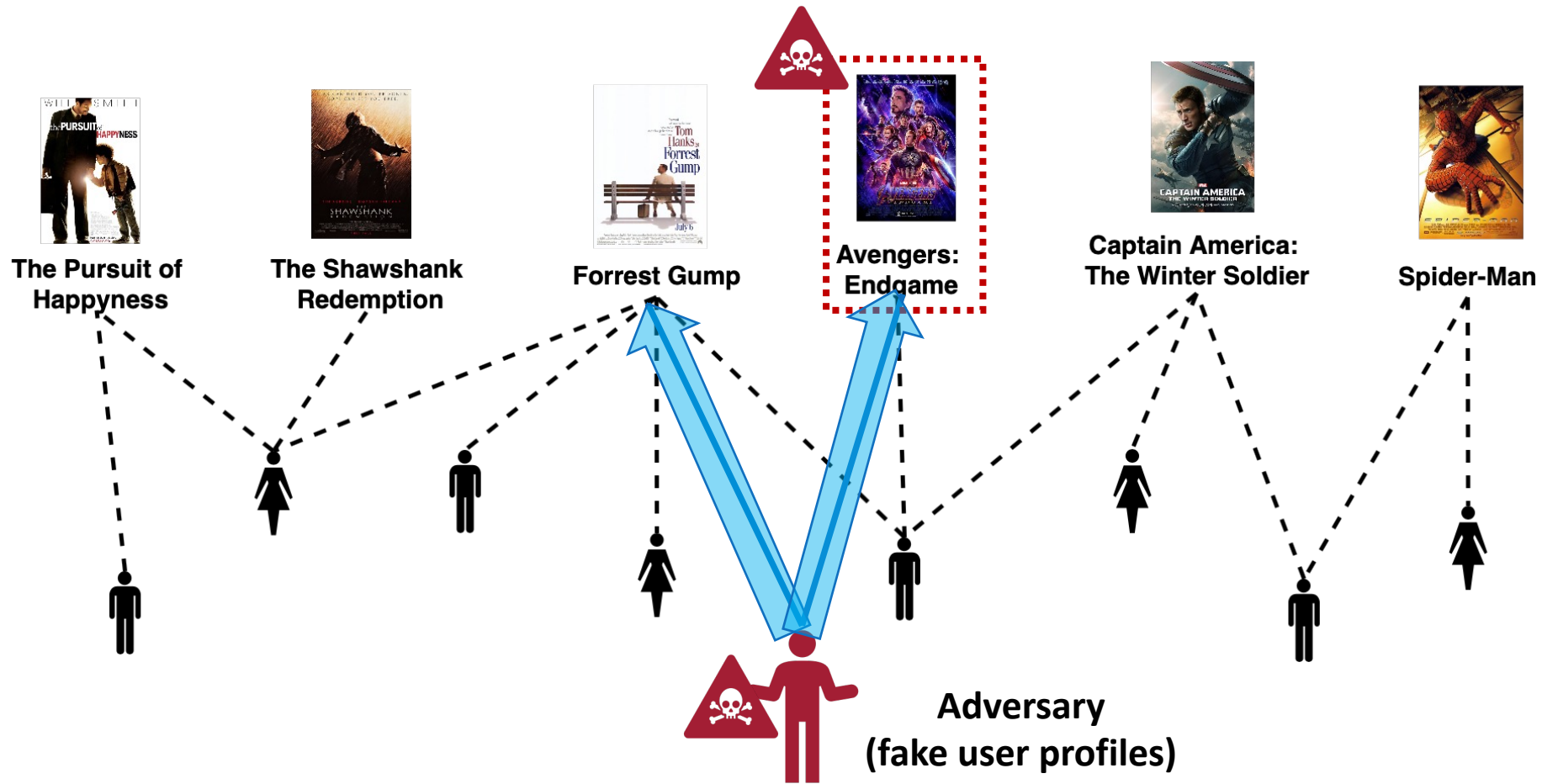


Non-discrimination & Fairness

- A recommender system should avoid discriminatory behaviors in human-machine interaction.
- A recommender system should ensure fairness in decision-making.



Safety & Robustness Issue



Attacks can happen in Recommender Systems



Business | Market Data | New Economy | New Tech Economy | Companies | Entrepreneurship | Technology of Business | Business of Sport | Global Education | Economy | Global Car Industry

Amazon 'flooded by fake five-star reviews' - Which? report

16 April 2019



Home > Competition

Press release

Facebook and eBay pledge to combat trading in fake reviews

Following action from the CMA, Facebook and eBay have committed to combatting the trade of fake and misleading reviews on their sites.

From: [Competition and Markets Authority](#)
Published 8 January 2020



“More than three-quarters of people are influenced by reviews when they shop online.”



Understand system's vulnerability and how attacks can be performed



Defend against potential adversarial attacks



“The Impact of Fake Reviews on Online Visibility: A Vulnerability Assessment of the Hotel Industry”, Information Systems Research, 2016

<https://www.bbc.com/news/business-47941181>

<https://www.gov.uk/government/news/facebook-and-ebay-pledge-to-combat-trading-in-fake-reviews>

Black-box Issue

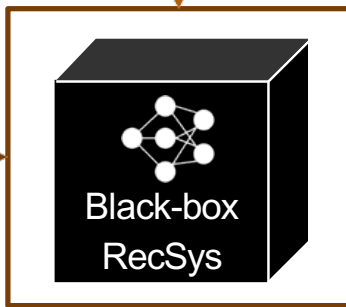
How recommender systems work?

Today

				...		
	5	4	?	...	?	?
	?	?	5	...	?	?
...	...					
	5	?	\hat{r}_{ij}	...	5	1
	?	?	?	...	2	5

Training Data

Learning Process



Learned Function



Output

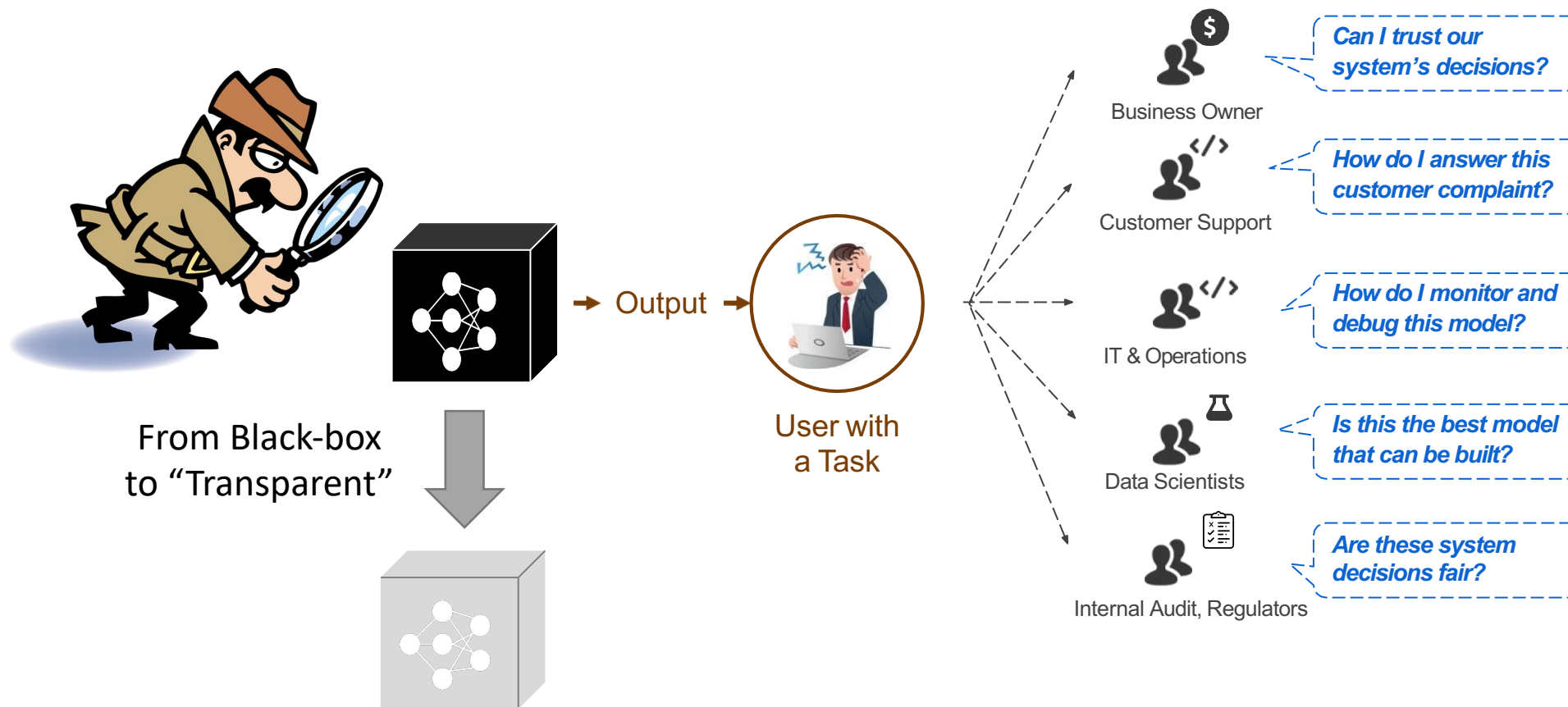


User with a Task

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

Explainability

Black-box system creates confusion and doubt



The Need for Explainable Recommendation

Privacy Issue



- ❑ The success of recommender systems heavily relies on data that might contain private and sensitive information.
- ❑ Can we still take the advantages of data while effectively protecting the privacy?

Environmental Issue

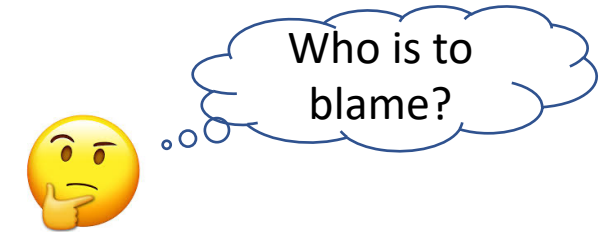


GPU Power Consumption Comparison

Dataset	XDL	DLRM	FAE
Criteo Kaggle	61.83W	58.91W	55.81W
Alibaba	56.39W	60.21W	56.62W
Criteo Terabyte	59.71W	62.47W	57.03W
Avazu	60.2W	58.03W	56.4W

Estimated carbon emissions from training common recommendation models

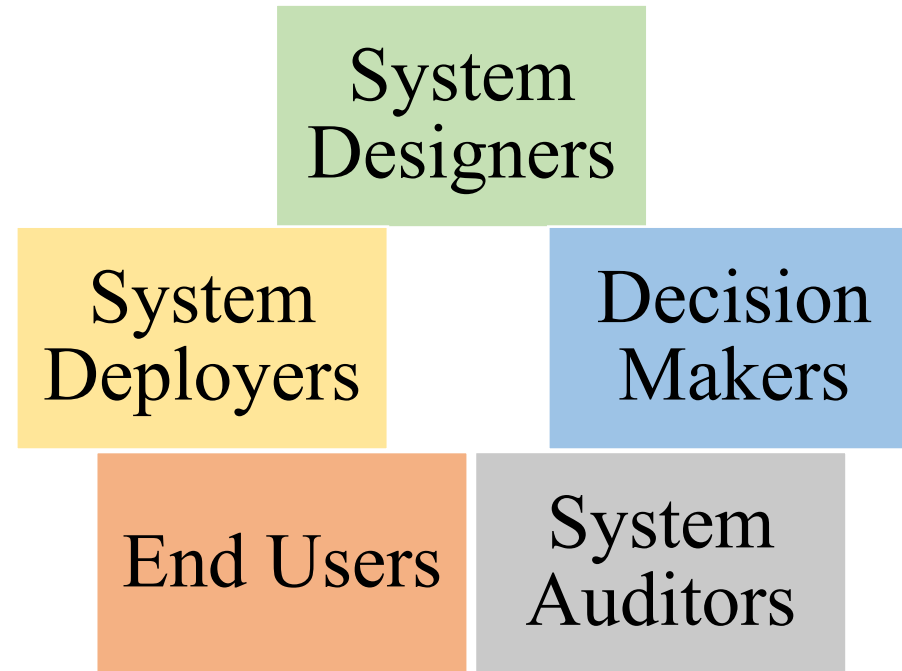
Auditability & Accountability



A clear responsibility distribution, which focuses on who should take the responsibility for what impact of recommender systems.

Auditability & Accountability

- Five roles in Recommender Systems



It is necessary to determine the roles and the corresponding responsibility of different parties in the function of a recommender system.

Interactions Among Different Dimensions



Privacy



Safety
& Robustness



Explainability



Non-discrimination
& Fairness



Environmental
Well-being



Accountability
& Auditability



How do these **SIX** dimensions influence each other?

There exist both **accordance** and the **conflicts** among the six dimensions.

Trustworthy Recommender Systems



A Survey on The Computational Perspective

A Comprehensive Survey on Trustworthy Recommender Systems

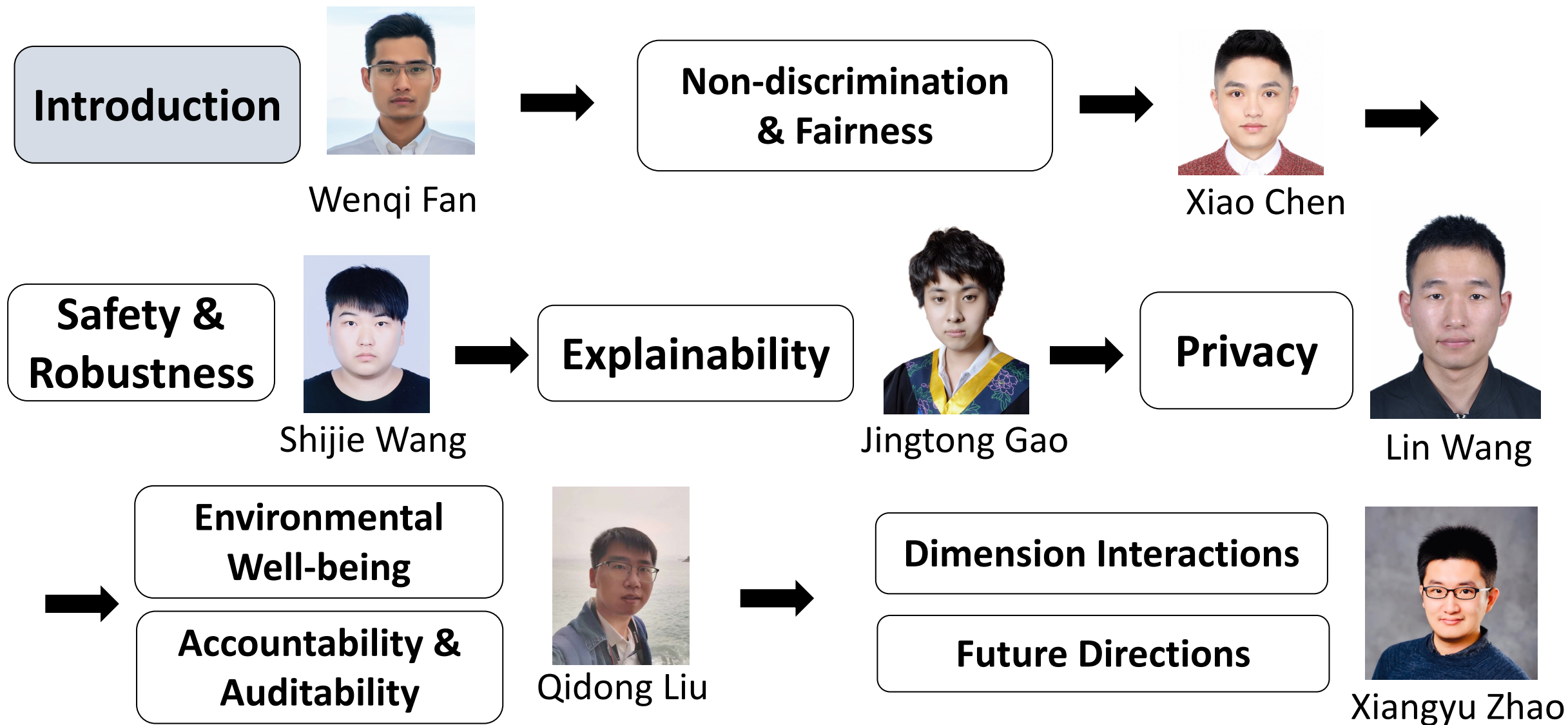
WENQI FAN, The Hong Kong Polytechnic University, Hong Kong
XIANGYU ZHAO*, City University of Hong Kong, Hong Kong
XIAO CHEN, The Hong Kong Polytechnic University, Hong Kong
JINGRAN SU, The Hong Kong Polytechnic University, Hong Kong
JINGTONG GAO, City University of Hong Kong, Hong Kong
LIN WANG, The Hong Kong Polytechnic University, Hong Kong
QIDONG LIU, City University of Hong Kong, Hong Kong
YIQI WANG, Michigan State University, USA
HAN XU, Michigan State University, USA
LEI CHEN, The Hong Kong University of Science and Technology, Hong Kong
QING LI, The Hong Kong Polytechnic University, Hong Kong

<https://arxiv.org/abs/2209.10117>

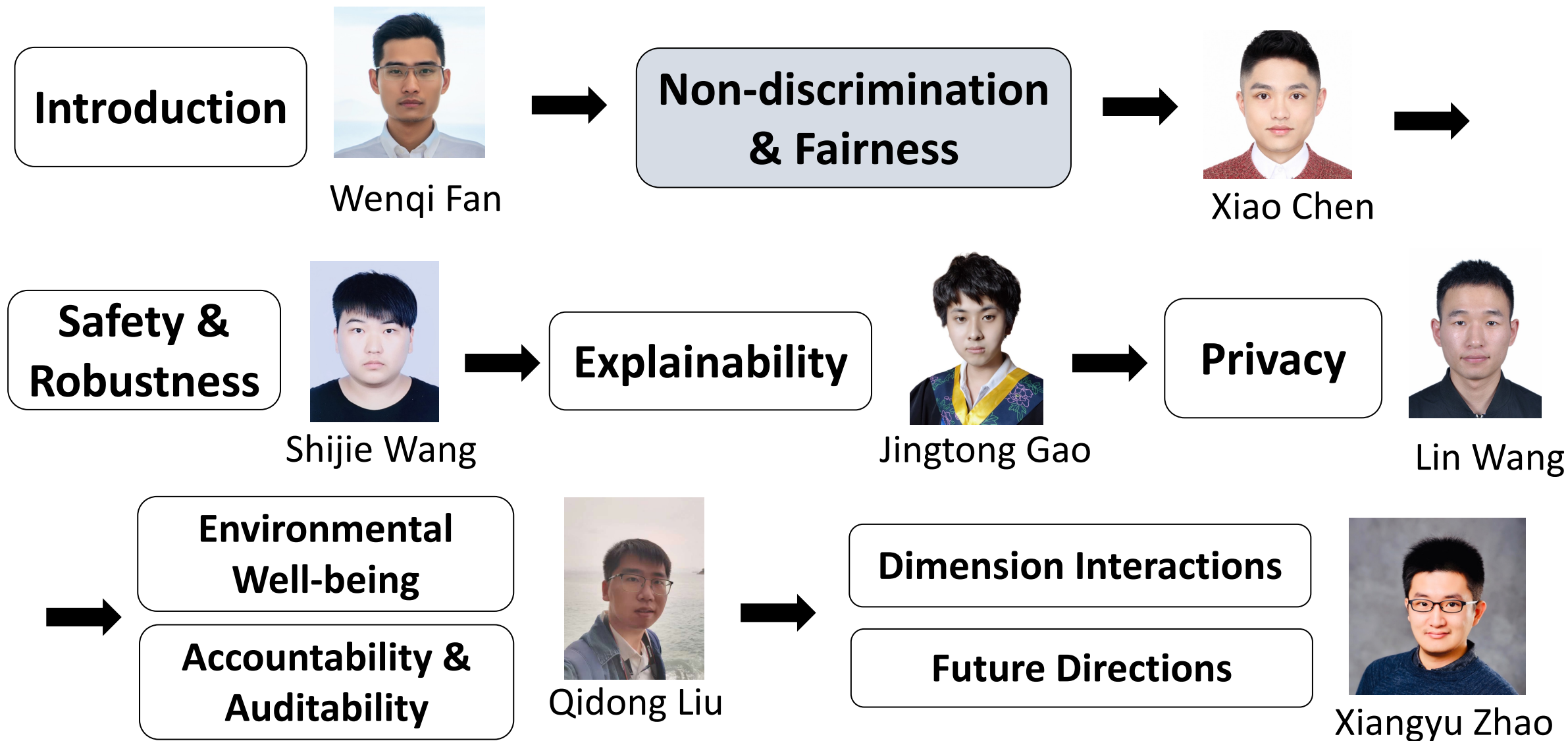
WWW'2023
Tutorial
Website (Slides)



Trustworthy Recommender Systems



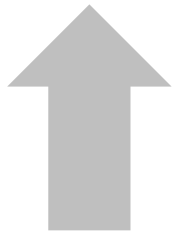
Trustworthy Recommender Systems



Contents



**CONCEPTS AND
TAXONOMY**



METHODOLOGY



APPLICATIONS



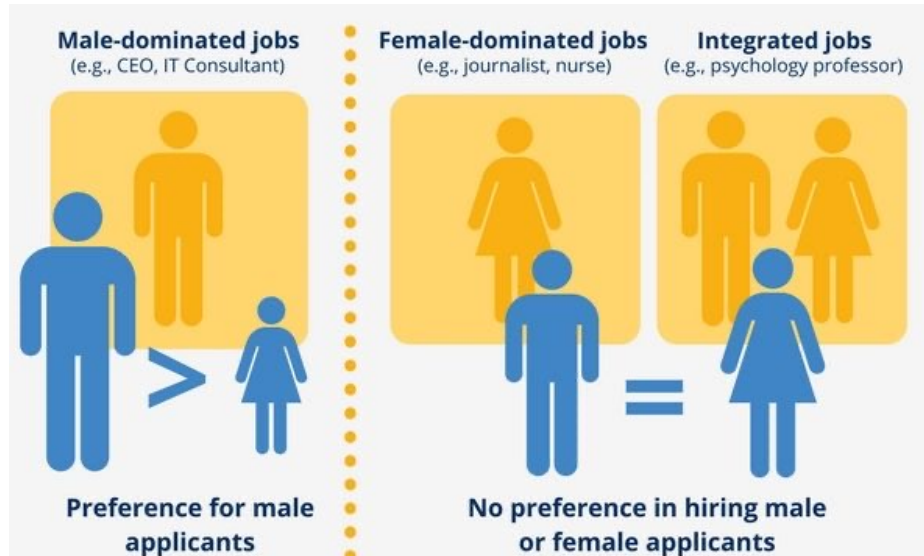
SURVEYS AND
TOOLS



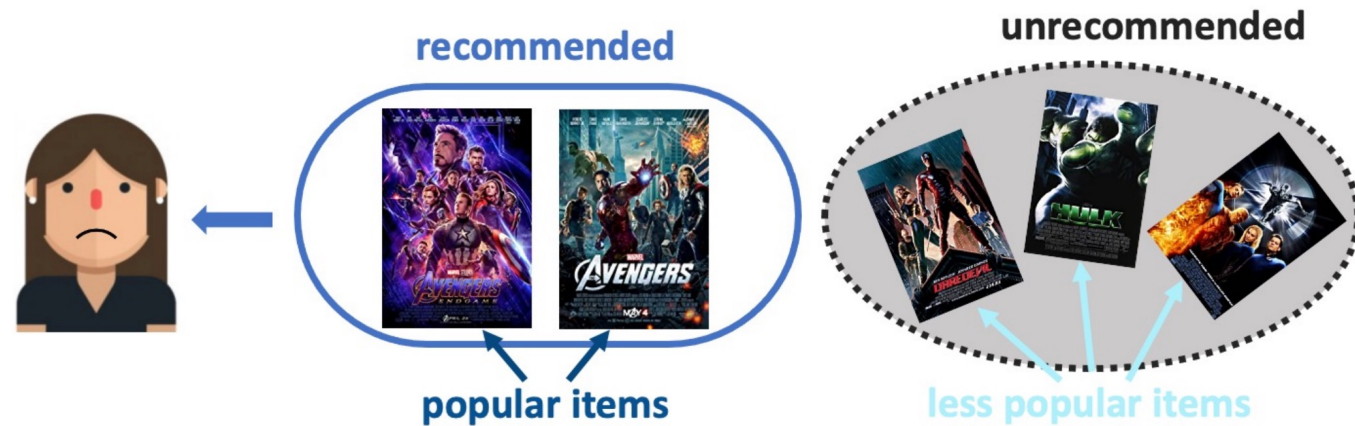
FUTURE
DIRECTIONS

Potential discrimination and bias in RecSys

- Recommender Systems make unfair decisions for specific user/item groups



Gender Discriminatory Bias [1]

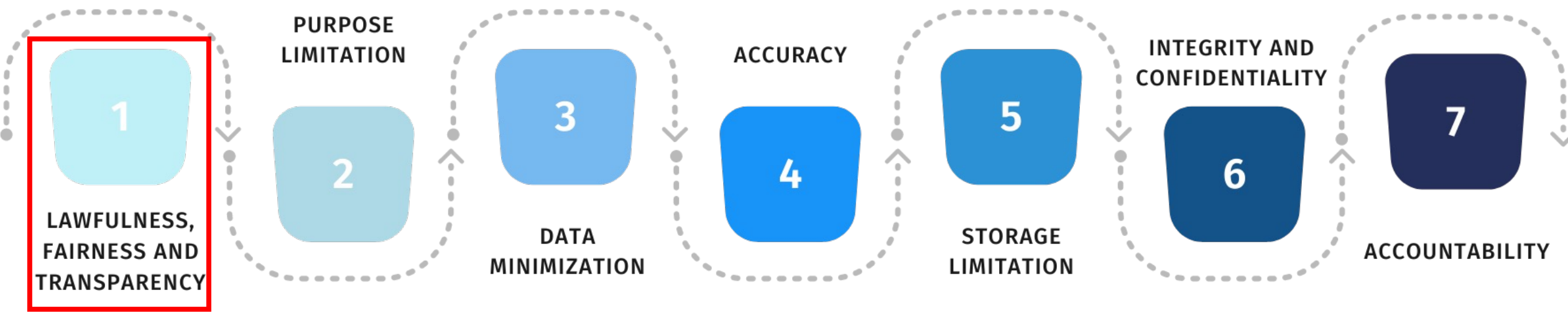


Popularity Bias [2]

[1] Lambrecht, et al. "Algorithmic bias? An empirical study of apparent gender-based discrimination in the display of STEM career ads." 2019.
 [2] Abdollahpouri, et al. "Popularity bias in ranking and recommendation." 2019.

Why Need Fairness in RecSys: From the Ethics Perspective

- 7 principles of EU GDPR regulation



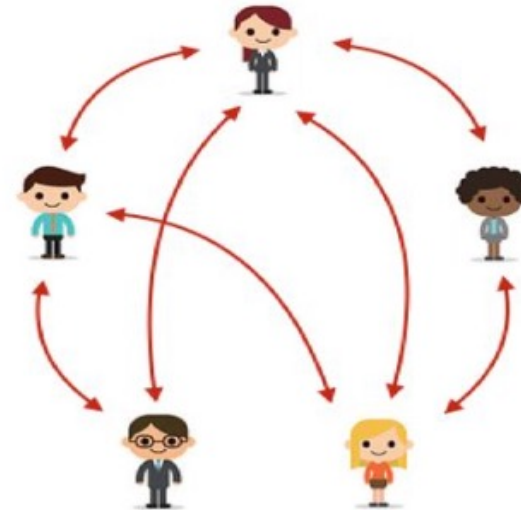
Fairness often couples with other responsible AI perspectives (e.g., explainability).

Why Need Fairness in RecSys: From the Utility Perspective

- Fair exposure opportunity guarantees the sustainable development of the RecSys platform



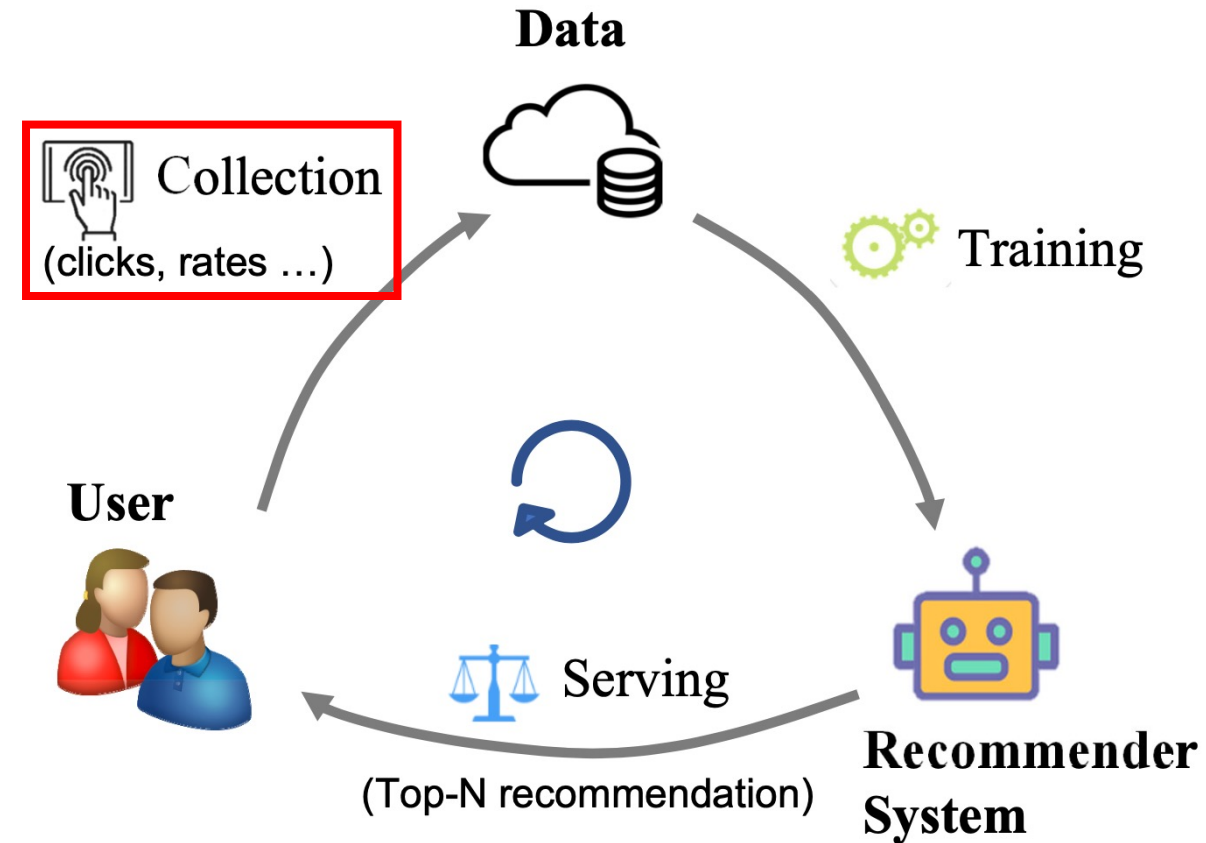
Big retailers vs. Small retailers
in the e-commerce system



Star accounts vs. Grassroot accounts
in the social recommendation system

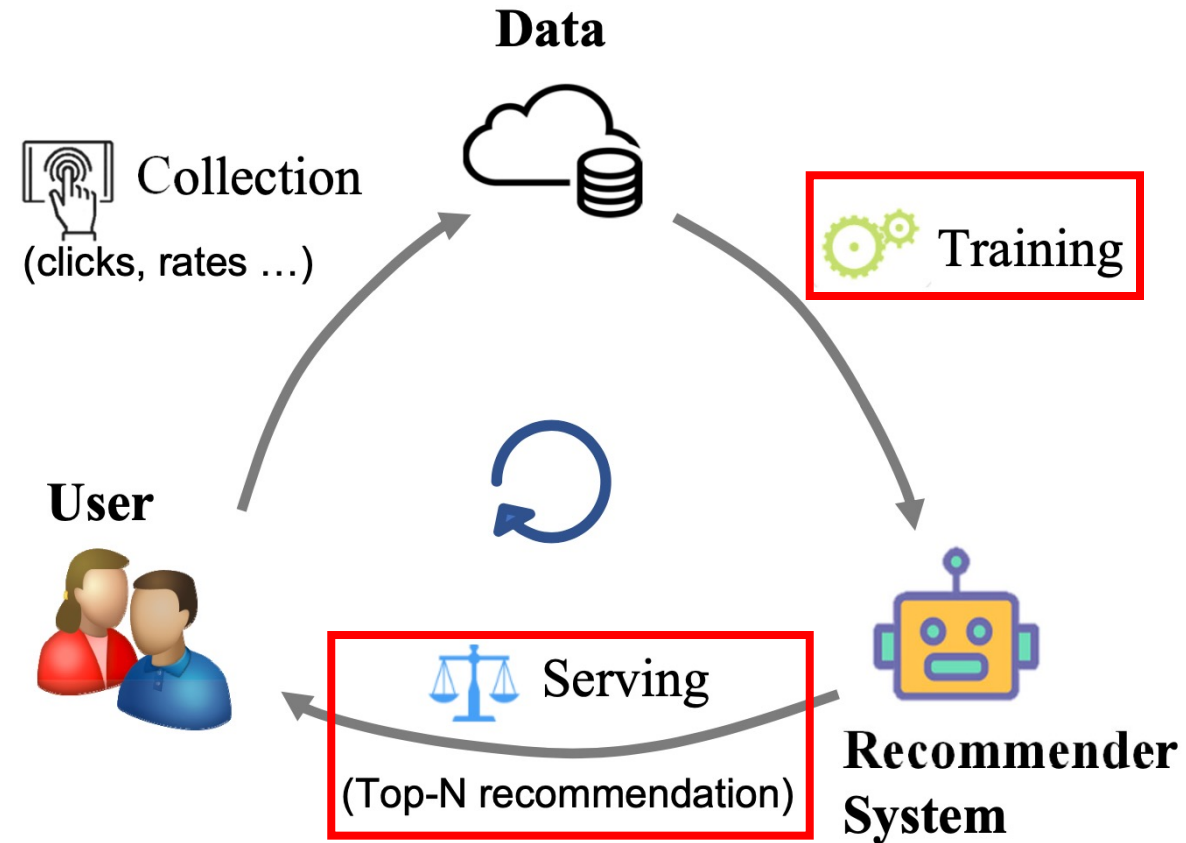
Sources of Bias

- **Data bias**
 - **Selection Bias:**
selecting rating behavior of users
 - **Exposure Bias:**
unobserved interactions may not fully represent the disliked items of users
 - **Conformity Bias:**
users behave similarly to other group members
 - **Position Bias:**
the higher positions on a recommendation list tends to receive more interaction



Sources of Bias

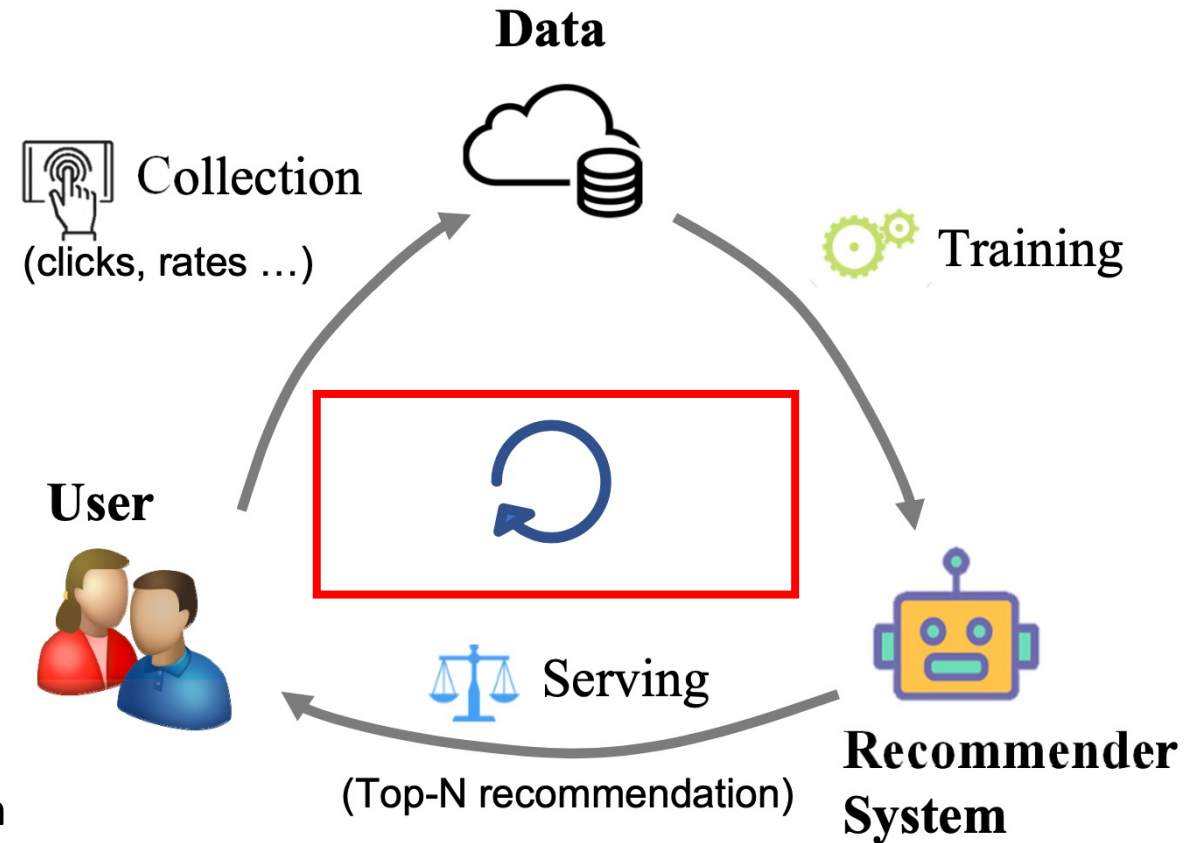
- Data bias
 - Selection Bias
 - Exposure Bias
 - Conformity Bias
 - Position Bias
- Model and result bias
 - **Popularity Bias:**
popular items are over-recommended compared to what their popularity warrant



Sources of Bias

- **Data bias**
 - Selection Bias
 - Exposure Bias
 - Conformity Bias
 - Position Bias
- **Model and result bias**
 - Popularity Bias
- **Feedback loop bias**
 - **Reinforced RS Feedback Loop Bias:**
 Unfair recommendations would influence users' behaviors in the online serving process

 Biased user behavior data enlarges model discrimination



Fairness Definition

- **Procedural Fairness:** procedural justice in decision-making processes
- **Outcome Fairness:** fair outcome performance

User Fairness vs. Item Fairness

Group Fairness vs. Individual Fairness

Causal Fairness vs. Associative Fairness

Static Fairness vs. Dynamic Fairness

Fairness Evaluation Metrics

- **Absolute Difference (AD):** group-wise utility difference

$$AD = |u(G_0) - u(G_1)|$$

- **Variance:** performance dispersion at the group/individual-level

$$\text{Variance} = \frac{1}{|\mathcal{V}|^2} \sum_{v_i \neq v_j} (u(v_i) - u(v_j))^2$$

- **Min-Max Difference:** the difference between the maximum and the minimum score value of all allocated utilities
- **Entropy**
- **KL-Divergence ...**

Contents



CONCEPTS AND
TAXONOMY



METHODOLOGY



APPLICATIONS



SURVEYS AND
TOOLS



FUTURE
DIRECTIONS

Method category

Pre-processing

Transform the data to remove the data bias before training

In-processing

Modify the learning algorithms to remove discrimination during the model training process

Post-processing

Perform post-processing by evaluating a holdout set that was not involved during model training

Pre-processing methods

- **Resampling**

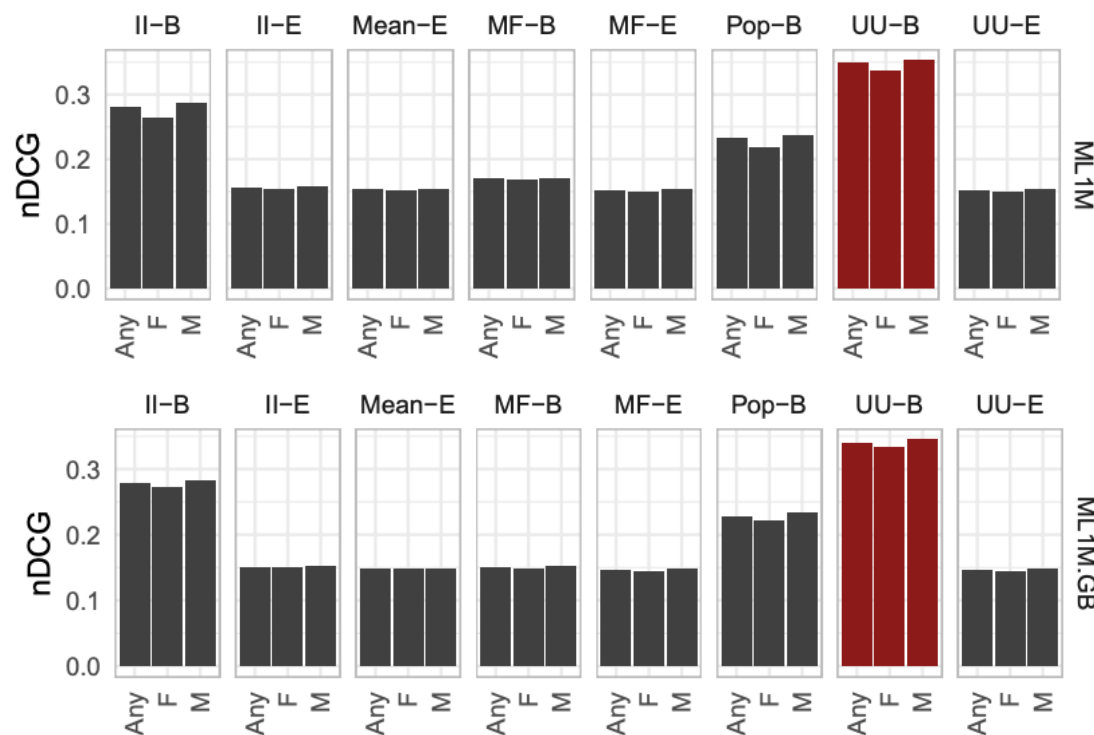
Rebalance the dataset distribution w.r.t the sensitive attribute

- **Data Augmentation**

Generating additional data for promoting the fairness of recommender systems

Pre-processing method (Resampling)

Idea: Different demographic groups obtain different utilities due to imbalanced data distribution. Balance the ratio of various user groups via a re-sampling strategy.



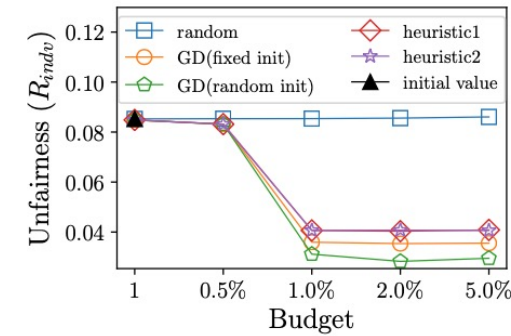
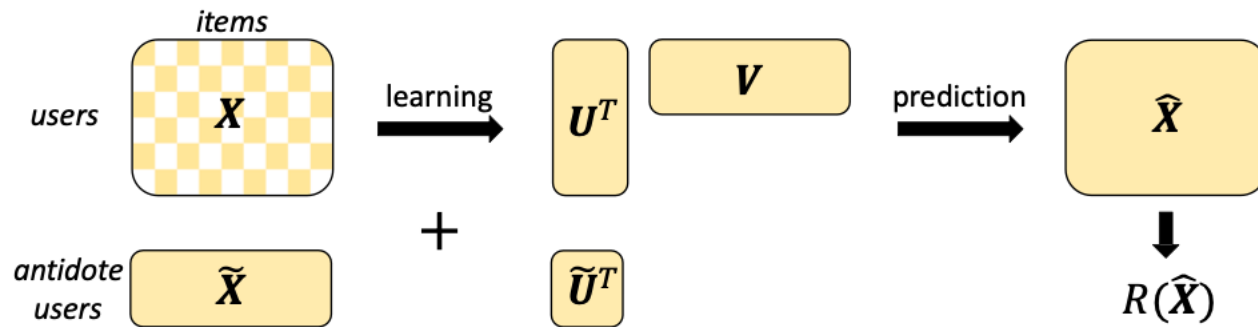
statistically-significant differences between gender groups

results on gender-balanced dataset

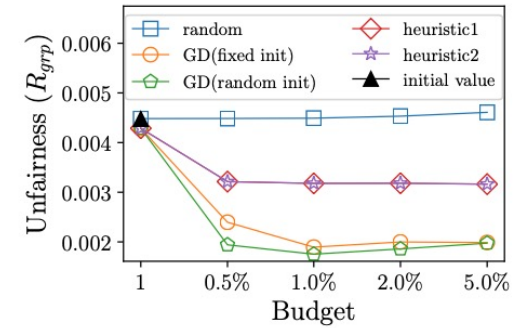
All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. ICFAT 2018.

Pre-processing method (Adding Antidote Data)

Idea: Improving the social desirability of recommender system outputs by adding more “antidote” data to the input.



(a) Individual fairness



(b) Group fairness

Matrix Factorization:
$$\arg \min_{U, V} ||P_{\Omega}(X - U^T V)||_F^2 + \lambda(||U||_F^2 + ||V||_F^2)$$

Objectives:
$$\arg \min_{\tilde{X} \in \mathcal{M}} R(\hat{X}(\Theta(X; \tilde{X})))$$

\swarrow fairness objective
 \searrow antidote data

Summary of Pre-processing methods



Flexibility, decoupled with the recommender systems



Performance gains might be degraded by the following steps

In-processing method

- **Regularization and constrained optimization**
- **Adversary Learning**
- **Causal graph**
- **Reinforcement Learning**
- **Others**

In-processing method (Regularization)

Idea: propose four new metrics that address different forms of unfairness. These metrics can be optimized by adding fairness terms to the learning objective [1].

$$U_{abs} = \frac{1}{n} \sum_{i=1}^n \left| |E_{adv}[y]_i - E_{adv}[r]_i| - |E_{\neg adv}[y]_i - E_{\neg adv}[r]_i| \right|,$$

$$\min_{P, Q, u, v} J(P, Q, u, v) + U.$$

Idea: a novel pairwise regularizer for pairwise ranking fairness [2].

$$\min_{\theta} \left(\sum_{(q, j, y, z) \in \mathcal{D}} \mathcal{L}_{rec}(f_{\theta}(q, v_j), (y, z)) \right) + |\text{Corr}_{\mathcal{P}}(A, B)|,$$

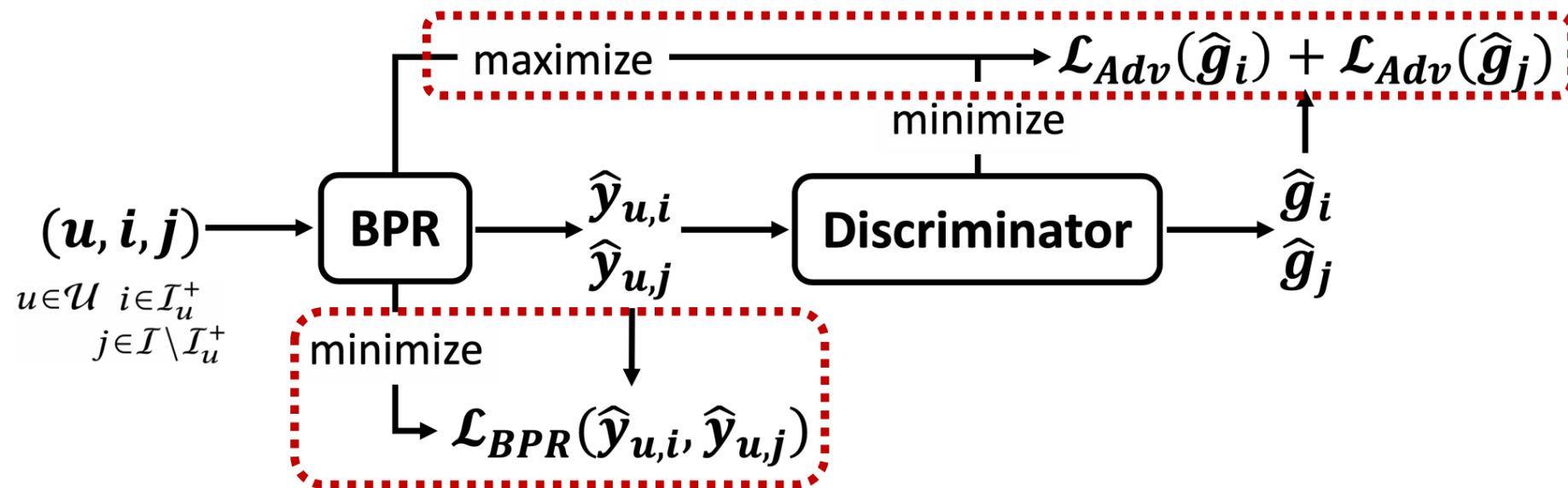
[1] Beyond Parity: Fairness Objectives for Collaborative Filtering. NeurIPS17

[2] Fairness in recommendation ranking through pairwise comparisons. KDD19

In-processing method (Adversary Learning)

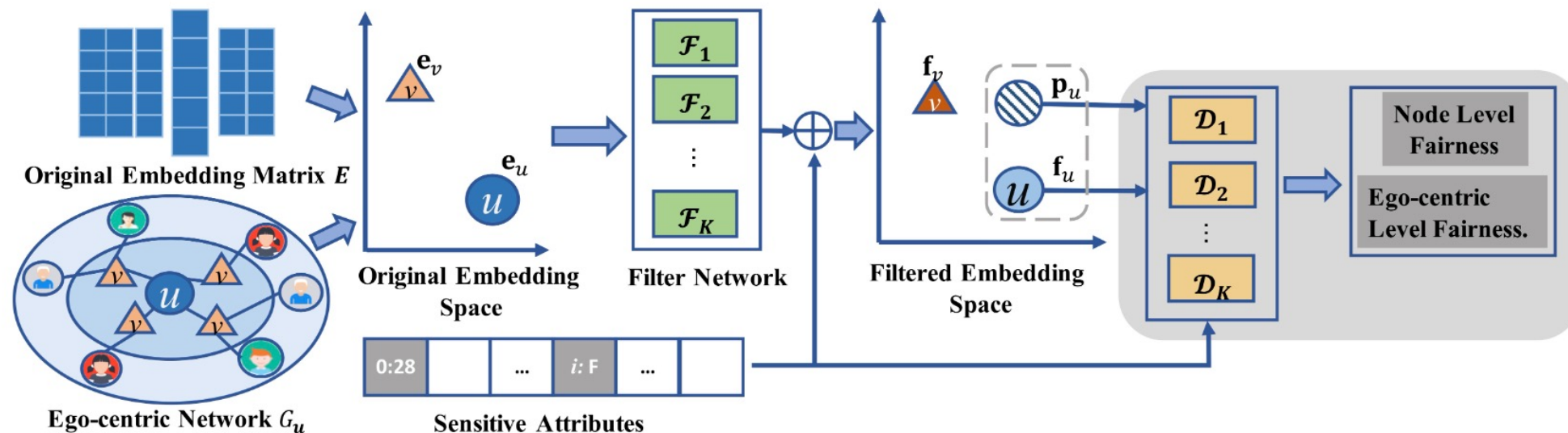
Idea: decouple the predicted score with the group attribute.

normalize the score distribution for each user to align predicted score with ranking position.



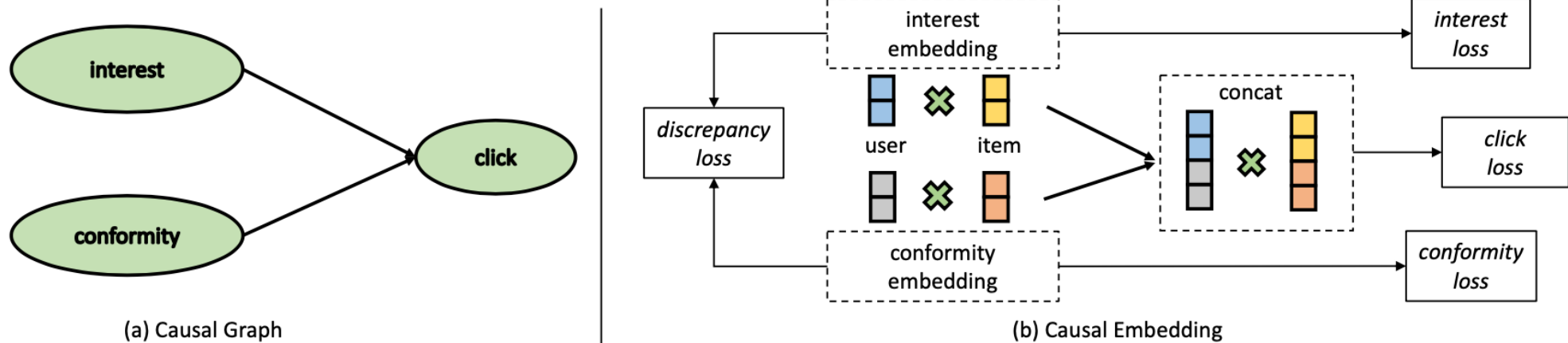
In-processing method (Adversary Learning)

Idea: propose a graph-based perspective for fairness-aware representation learning of any recommendation models. Adversarial learning of a user-centric graph.



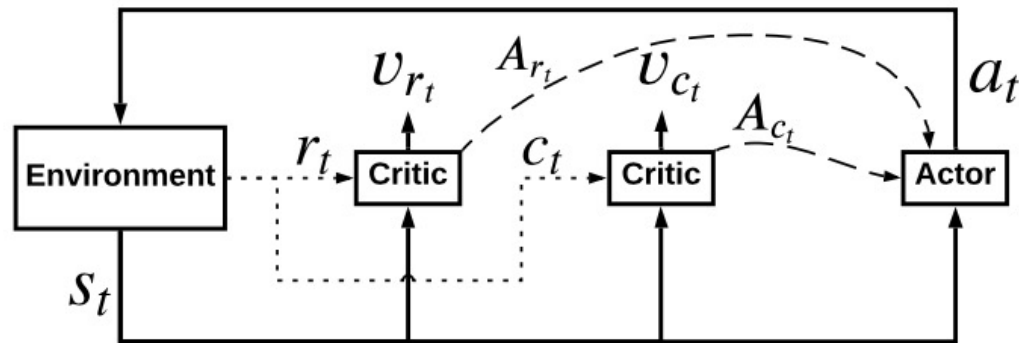
In-processing method (Causal Graph)

Idea: Disentangling Interest and Conformity with Causal Embedding (DICE).
 Separate embeddings are adopted to capture the two causes, and are trained with cause-specific data.



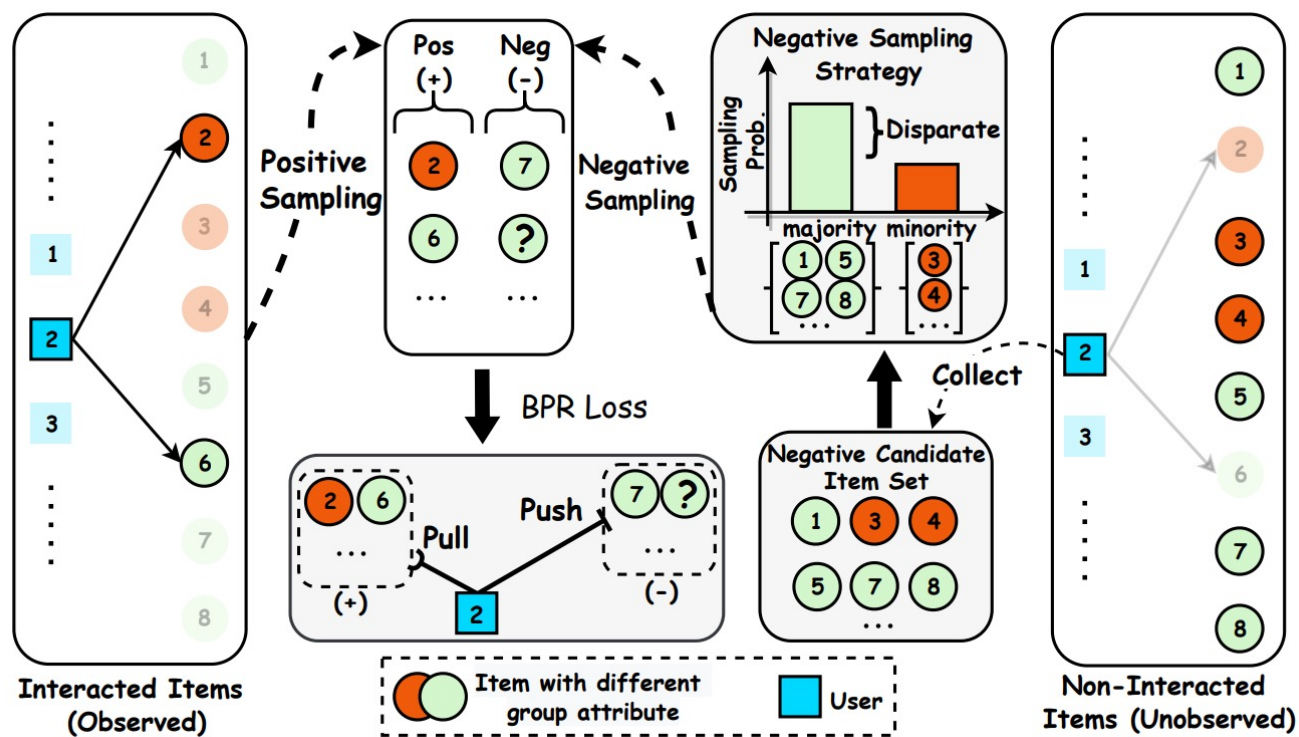
In-processing method (Reinforcement Learning)

Idea: propose a fairness-constrained reinforcement learning algorithm, which models the recommendation problem as a Constrained Markov Decision Process (CMDP). Dynamically adjust the recommendation policy for the fairness requirement.



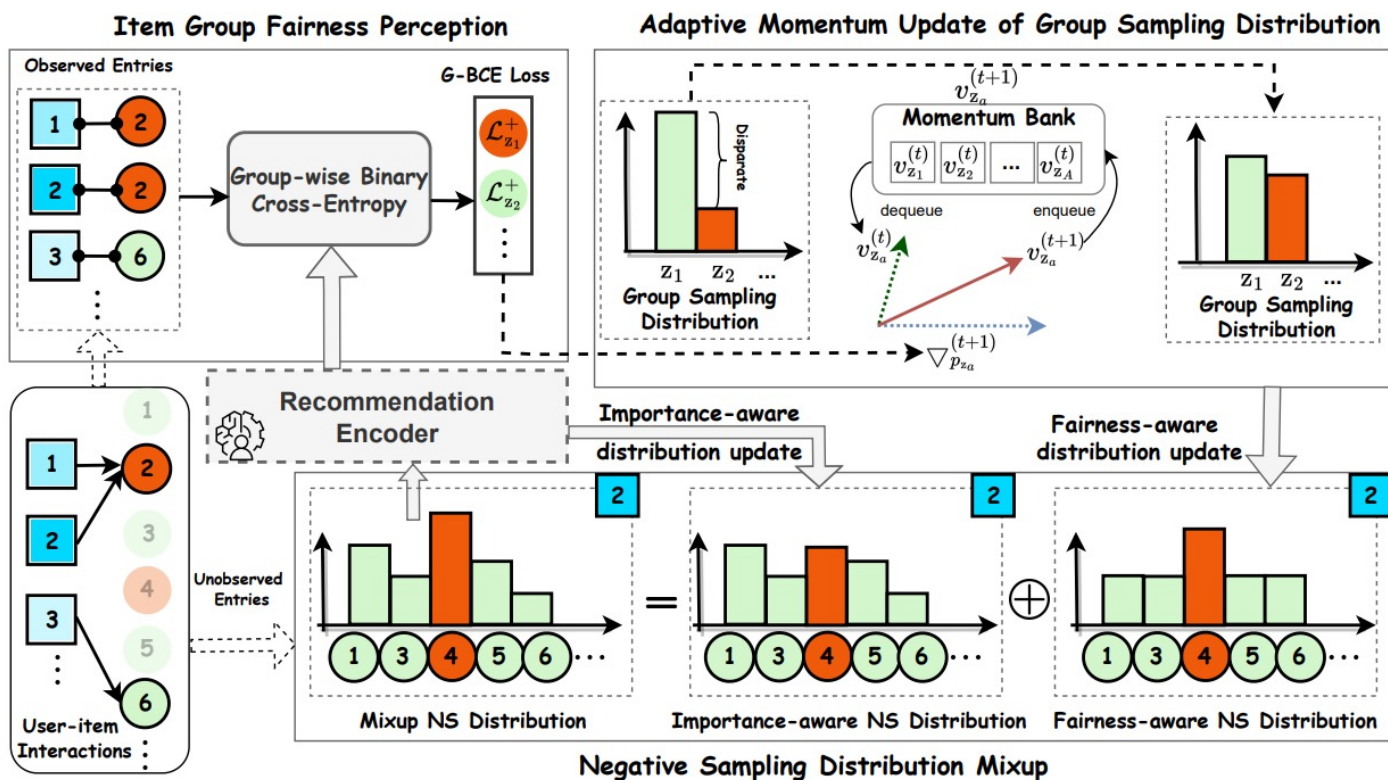
In-processing method (Negative Sampling)

- **Observation:** the majority item group obtains low (biased) prediction scores via the BPR loss (group-wise performance disparity)



In-processing method (Negative Sampling)

- **Idea:** adjust the negative sampling distribution (group-wise) adaptively in the training process for meeting the item group fairness objective



In-processing method (Negative Sampling)

- Bi-level Optimization of FairNeg

The optimization of the group-wise negative sampling distribution is nested within the recommendation model parameters optimization

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} \mathcal{L}_{\text{Recall-Disp}}(\Theta, \mathbf{p}) := \sum_{z_a \in Z} \left| \mathcal{L}_{z_a}^+ - \frac{1}{|A|} \sum_{z \in Z} \mathcal{L}_z^+ \right|,$$
$$\Theta_{\mathbf{p}}^* = \arg \min_{\Theta} \mathcal{L}_{\text{utility}}(\Theta, \mathbf{p}) := - \sum_{u \in \mathcal{U}} \sum_{i \in \mathcal{V}_u^+, j \in \mathcal{V}_u^-} \mathcal{L}_{\text{BPR}}(u, i, j; \Theta, \mathbf{p}),$$

- Updating Group Sampling Distribution

(1) Group-wise gradient calculation

$$\nabla_{\mathbf{p}_{z_a}}^{(t)} := \mathcal{L}_{z_a}^{+(t)} - \frac{1}{|A|} \sum_{z \in Z} \mathcal{L}_z^{+(t)},$$

(2) Adaptive momentum update

$$\mathbf{v}_{z_a}^{(t+1)} = \gamma \mathbf{v}_{z_a}^{(t)} + \alpha \cdot \nabla_{\mathbf{p}_{z_a}}^{(t+1)},$$
$$\mathbf{p}_{z_a}^{(t+1)} = \mathbf{p}_{z_a}^{(t)} - \mathbf{v}_{z_a}^{(t+1)},$$

Summary of In-processing methods



Substantial fairness improvements



Fairness and utility trade-off

Resource-intensive

Post-processing method

- **Slot-wise reranking**
- **Global-wise reranking**
- **User-wise reranking**

Slot-wise Re-ranking

Idea: propose a personalized re-ranking algorithm to achieve a fair microlending RS.

A combination of personalization score and a fairness term.

$$\max_{v \in R(u)} \underbrace{(1 - \lambda)P(v | u)}_{\text{personalization}} + \lambda \underbrace{\sum_c P(\mathcal{V}_c) \mathbb{1}_{\{v \in \mathcal{V}_c\}} \prod_{i \in S(u)} \mathbb{1}_{\{i \notin \mathcal{V}_c\}}}_{\text{fairness}},$$

User-wise Re-ranking

Idea: formulate fairness constraints on rankings in terms of exposure allocation. Find rankings that maximize the utility for the user while provably satisfying a specific notion of fairness.

$$\begin{aligned}
 \mathbf{P} &= \operatorname{argmax}_{\mathbf{P}} \mathbf{u}^T \mathbf{P} \mathbf{v} && \text{(expected utility)} \\
 \text{s.t. } \mathbb{1}^T \mathbf{P} &= \mathbb{1}^T && \text{(sum of probabilities for each position)} \\
 \mathbf{P} \mathbb{1} &= \mathbb{1} && \text{(sum of probabilities for each document)} \\
 0 \leq \mathbf{P}_{i,j} &\leq 1 && \text{(valid probability)} \\
 \mathbf{P} &\text{ is fair} && \text{(fairness constraints)}
 \end{aligned}$$

$$\text{Exposure}(G_0 | \mathbf{P}) = \text{Exposure}(G_1 | \mathbf{P}) \quad (4)$$

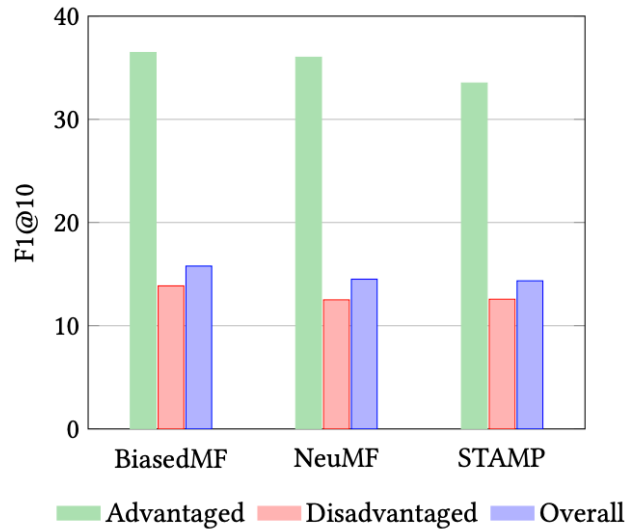
$$\Leftrightarrow \frac{1}{|G_0|} \sum_{d_i \in G_0} \sum_{j=1}^N \mathbf{P}_{i,j} \mathbf{v}_j = \frac{1}{|G_1|} \sum_{d_i \in G_1} \sum_{j=1}^N \mathbf{P}_{i,j} \mathbf{v}_j \quad (5)$$

$$\Leftrightarrow \sum_{d_i \in \mathcal{D}} \sum_{j=1}^N \left(\frac{\mathbb{1}_{d_i \in G_0}}{|G_0|} - \frac{\mathbb{1}_{d_i \in G_1}}{|G_1|} \right) \mathbf{P}_{i,j} \mathbf{v}_j = 0 \quad (6)$$

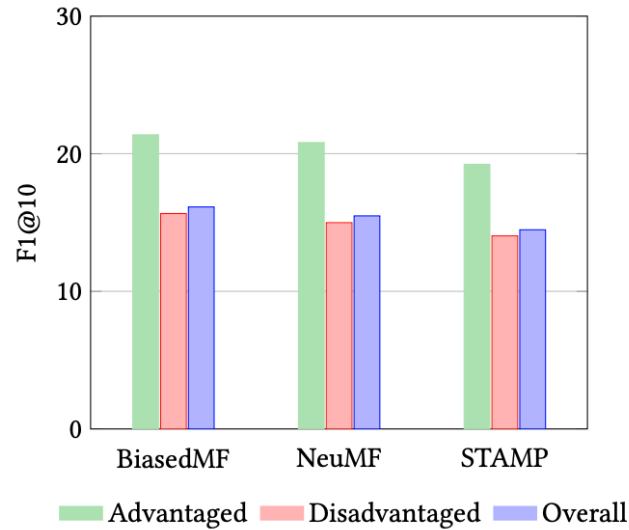
$$\Leftrightarrow \mathbf{f}^T \mathbf{P} \mathbf{v} = 0 \quad \left(\text{with } \mathbf{f}_i = \frac{\mathbb{1}_{d_i \in G_0}}{|G_0|} - \frac{\mathbb{1}_{d_i \in G_1}}{|G_1|} \right)$$

Global-wise Re-ranking

Idea: a re-ranking approach to mitigate this unfairness problem by adding constraints over evaluation metrics.



(a) Original



(b) Fair Method

$$\begin{aligned}
 & \max_{\mathbf{W}_{ij}} \sum_{i=1}^n \sum_{j=1}^N \mathbf{W}_{ij} S_{i,j} \\
 & \text{s.t. } \text{UGF} (Z_1, Z_2, \mathbf{W}) < \varepsilon \\
 & \sum_{j=1}^N \mathbf{W}_{ij} = K, \mathbf{W}_{ij} \in \{0, 1\}
 \end{aligned}$$

Summary of Post-processing methods



Can be applied to any recommendation systems



Constrained to unfair recommendation model outputs

• Summary of existing methods

Taxonomy	Method type	Related research
Pre-processing	Data Re-sampling	[95]
	Adding Antidote Data	[289]
In-processing	Regularization & Constrained Optimization	[26, 351, 393, 409, 461]
	Adversarial Learning	[33, 207, 215, 221, 285, 379, 380]
	Reinforcement Learning	[120, 122, 244]
	Causal Graph	[121, 162, 387, 452]
	Others	[31, 110, 167, 224]
Post-processing	Slot-wise Re-ranking	[124, 185, 189, 243, 262, 300, 305] [306, 323, 328, 405, 419]
	User-wise Re-ranking	[28, 253, 304, 318]
	Global-wise Re-ranking	[87, 114, 219, 250, 279, 335, 384, 462]

Contents



CONCEPTS AND
TAXONOMY



METHODOLOGY



APPLICATIONS



SURVEYS AND
TOOLS



FUTURE
DIRECTIONS

Applications

- **Ecommerce (Amazon, Etsy)**
- **Social Media (Twitter, LinkedIn)**
- **Content Streaming (Spotify, Youtube)**
- **Ride-hailing (Uber, Lyft)**



Contents



CONCEPTS AND
TAXONOMY



METHODOLOGY



APPLICATIONS



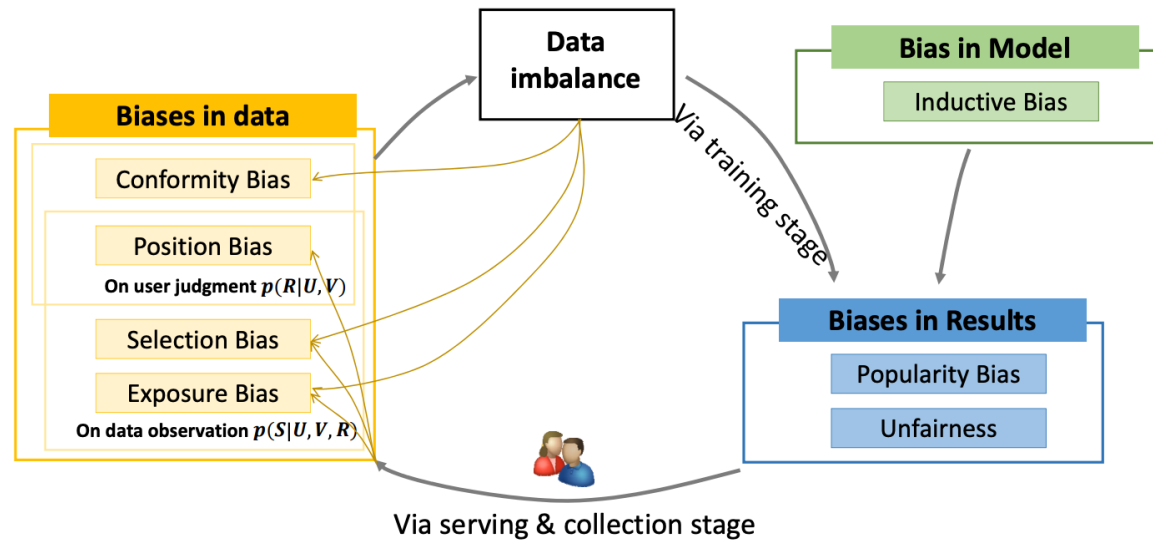
**SURVEYS AND
TOOLS**



FUTURE
DIRECTIONS

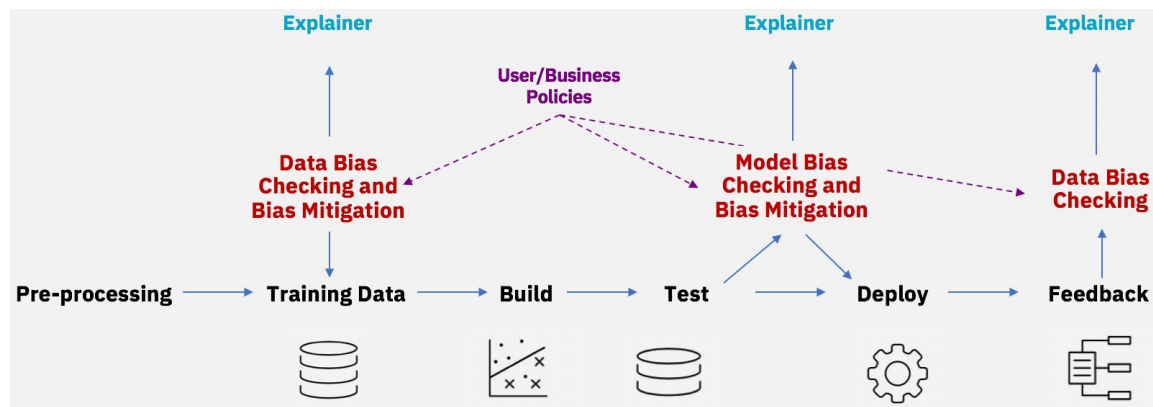
Surveys

- TOIS 23' Bias and Debias in Recommender System: A Survey and Future Directions
- Arxiv 22' A Comprehensive Survey on Trustworthy Recommender Systems

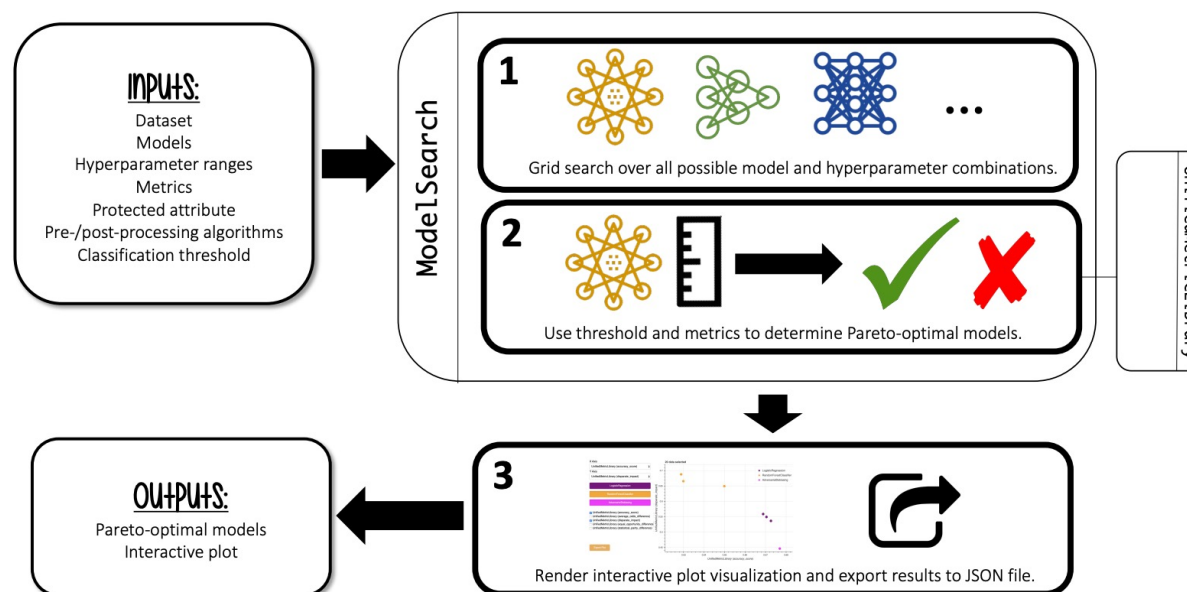


Tools

- IBM Fairness 360



- Fairkit-learn



Contents



CONCEPTS AND
TAXONOMY



METHODOLOGY



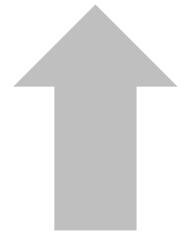
APPLICATIONS



SURVEYS AND
TOOLS



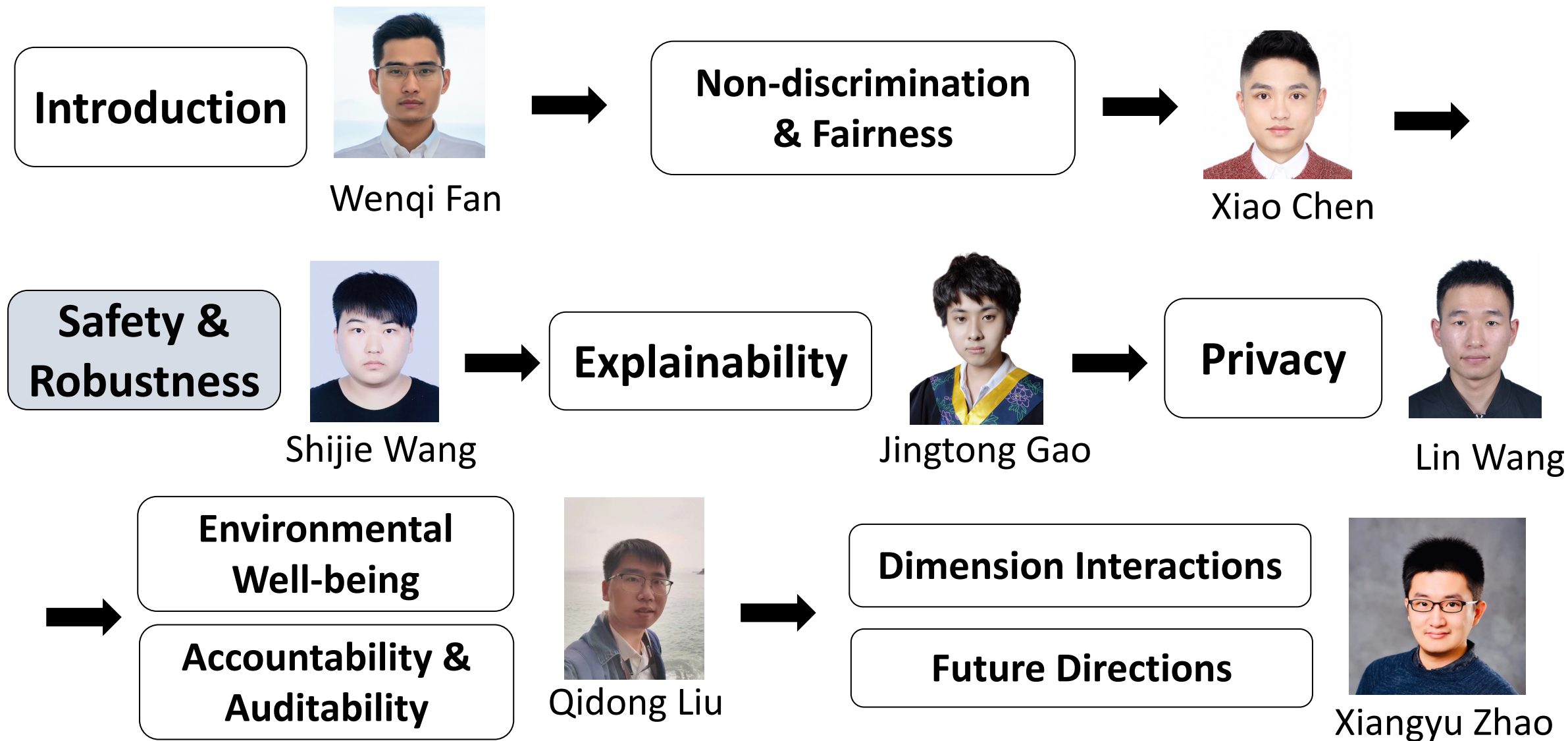
**FUTURE
DIRECTIONS**



Future Directions

- **Consensus on Fairness Definition**
- **Fairness-Utility tradeoff**
- **Fairness-aware algorithm design**
- **Better evaluation**

Trustworthy Recommender Systems



Real World Attacks in Recommender Systems

DIGITAL LIVING | JULY 26, 2022

Amazon's War on Fake Reviews

By Matt Stieb, *Intelligencer* staff writer



Photo-Illustration: Intelligencer; Photos: Getty Images/Amazon

BUSINESS

How merchants use Facebook to flood Amazon with fake reviews

By Elizabeth Dvoskin and Craig Timberg

April 23, 2018 at 1:26 p.m. EDT



An Amazon distribution center in Madrid, shown in November. (Emilion Naranjo/EPA-EFE/Shutterstock)

Safety and Robustness

“A decision aid, no matter how sophisticated or ‘intelligent’ it may be, may be rejected by a decision maker who does not trust it, and so its potential benefits to system performance will be lost.”

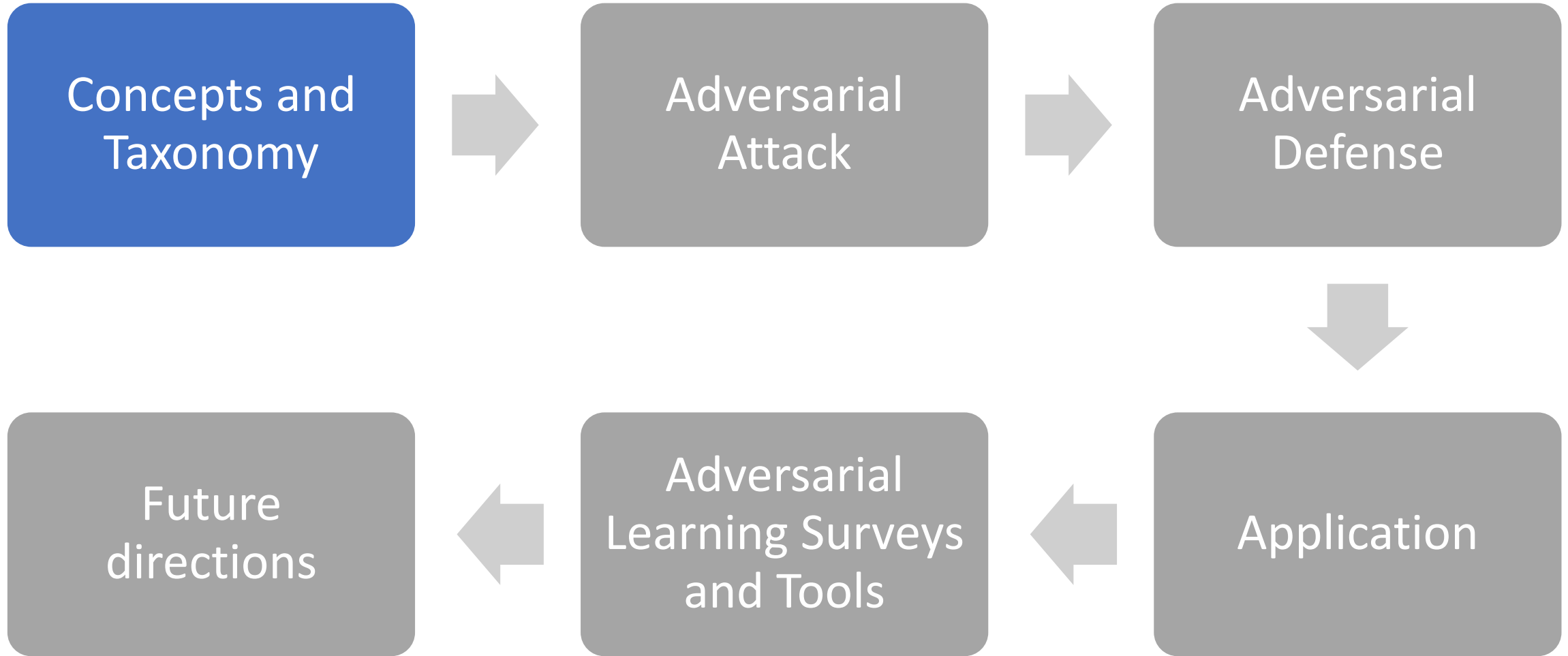
—Bonnie M. Muir, psychologist at University of Toronto

Safety and Robustness

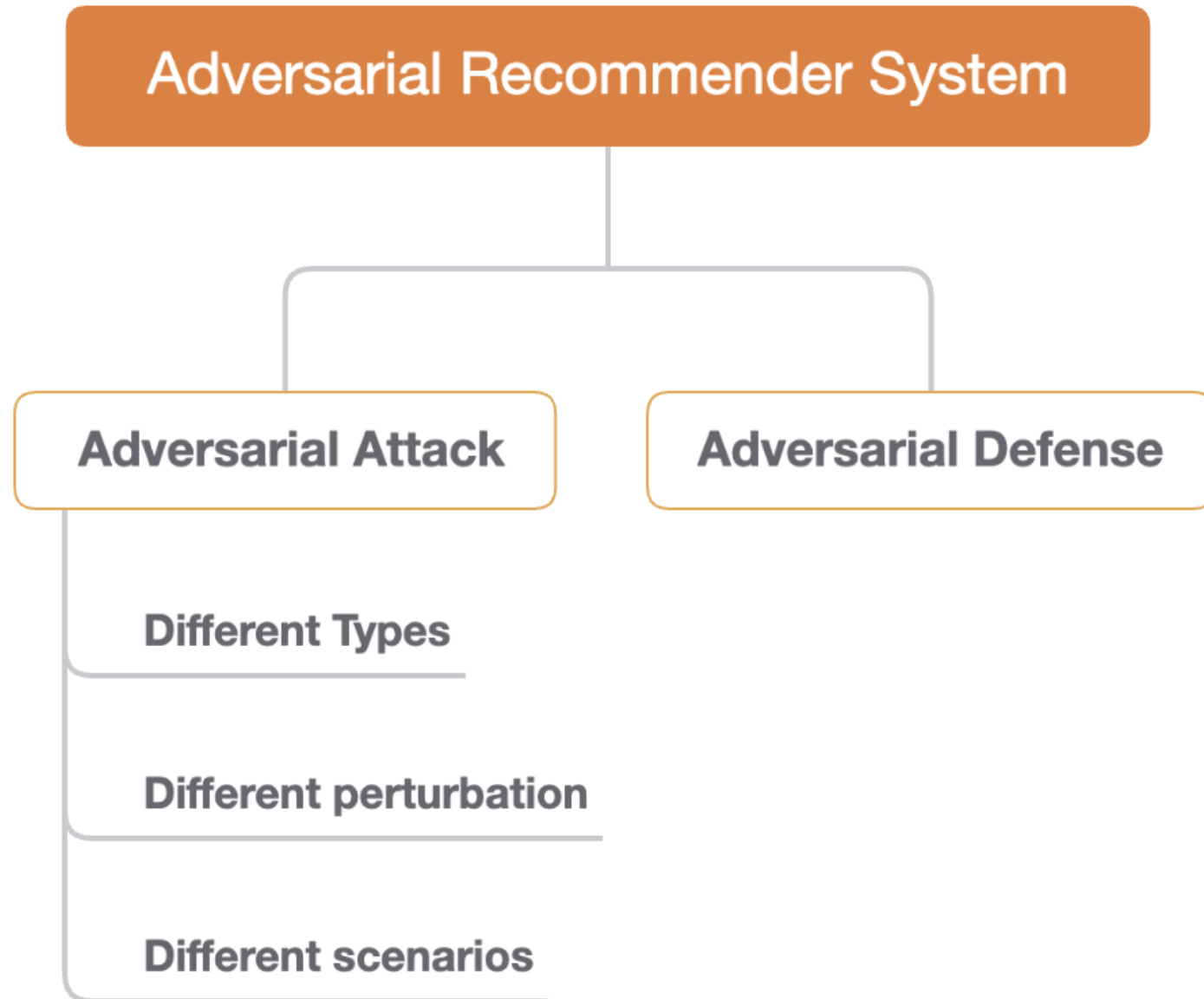
By examining Adversarial Robustness,
we expect the recommender system to:

- Be reliable, secure and stable

Outline



Taxonomy

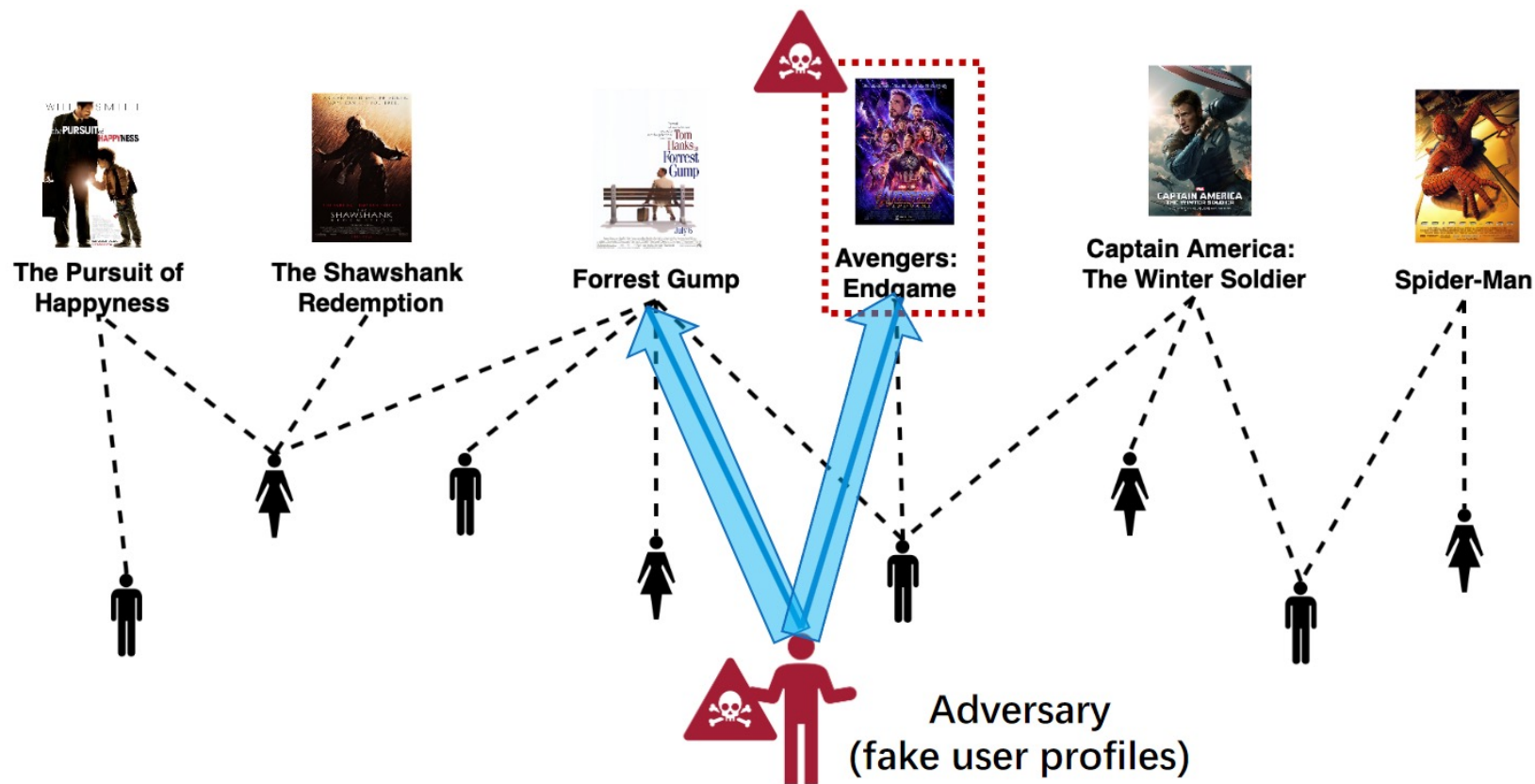


Adversarial Attack

- Poisoning Attacks vs. Evasion Attacks
 - They happen in **training phase**/ happen in **test/inference phase**
- White-box attacks vs. Grey-box attacks vs. Black-box attacks
 - They have **all knowledge** of the recommender system / have **partial knowledge**/ have **no knowledge** or limit knowledge
- Targeted Attacks vs. Untargeted Attacks
 - They aim to **promote/demote** a set of **target items**/ aim to **degrade** a recommendation system's **overall performance**

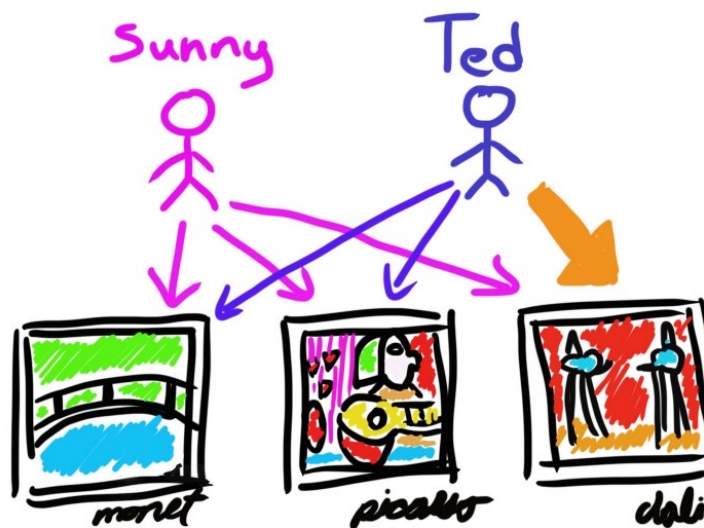
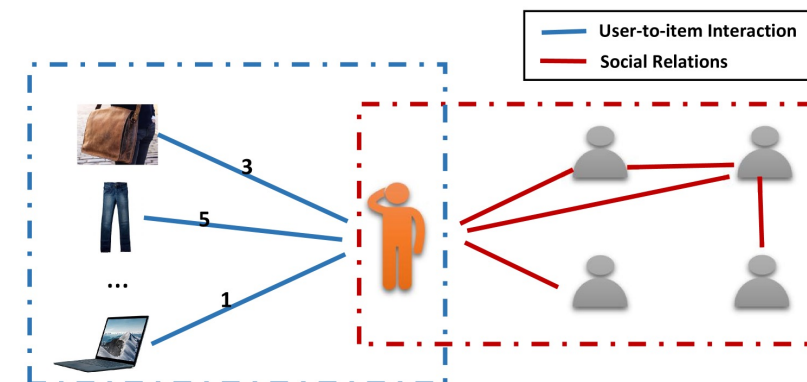
Adversarial in Different Perturbation

- Adding fake user profiles into user-item interactions, modifying user attributes information, adding social relations, etc



Adversarial in Different Scenarios

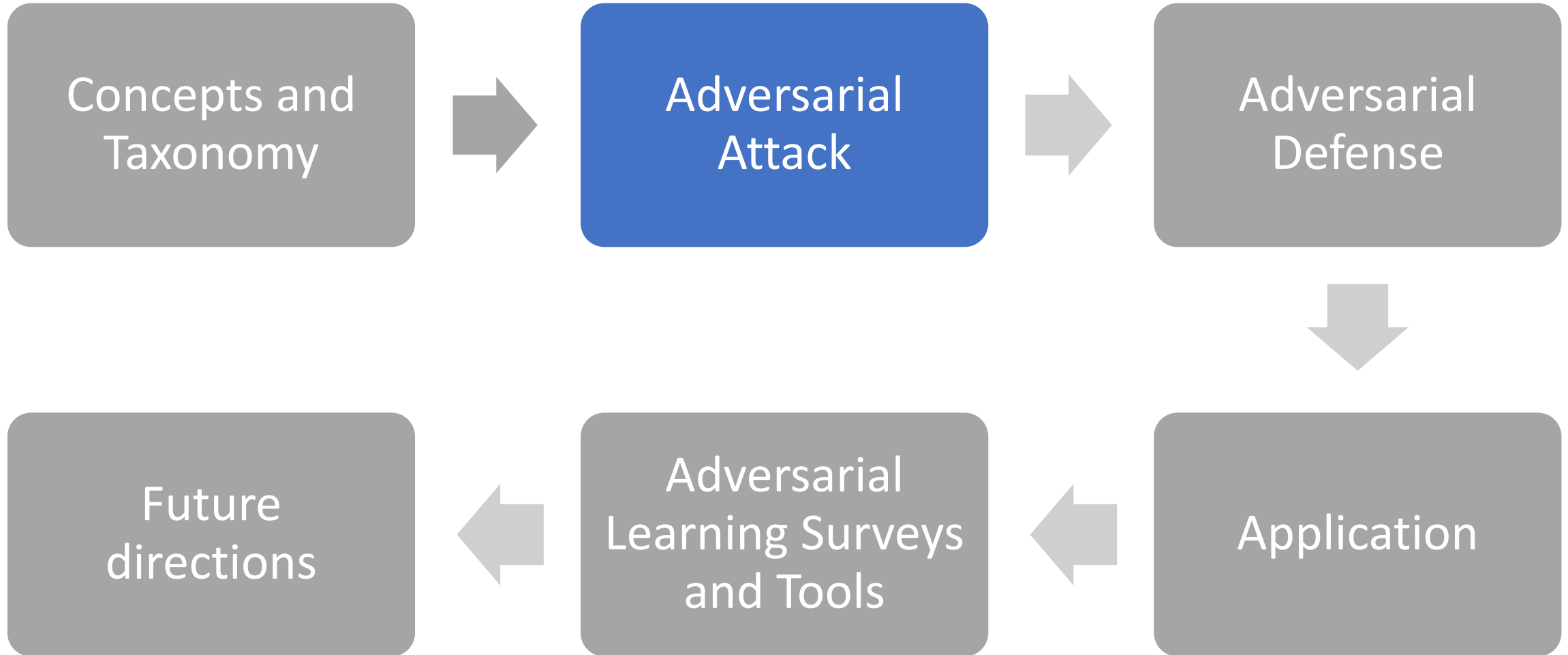
- Collaborative Filtering Recommender System
- Social Recommender System
- Content-based Recommender System
- ...



Adversarial Defenses

- Perturbations Detection vs. Adversarial Training
 - It is to **identify perturbations** data and remove them/ **enhances the robustness** of recommender systems

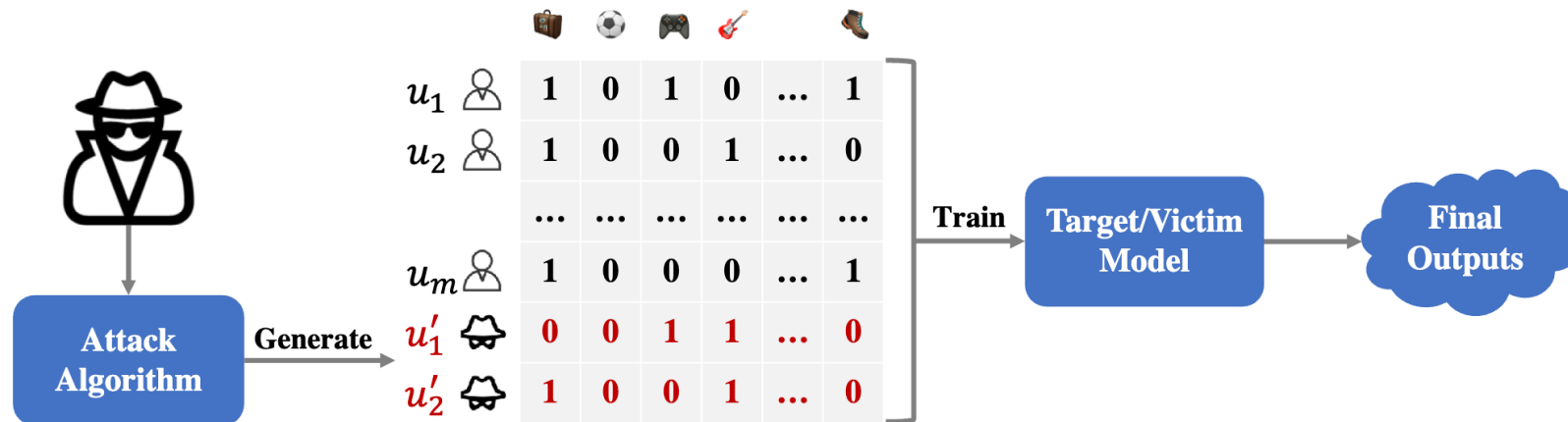
Outline



Adversarial Attack for Recommender System

- A Unified Formulation of Poisoning Attack

$$\min_{\hat{U}} \mathcal{L}_{adv}(\theta^*), \quad \text{s.t.} \quad \theta^* = \arg \min_{\theta} (\mathcal{L}_{rec}(\mathbf{R}, \mathbf{O}_{\theta}) + \mathcal{L}_{rec}(\hat{\mathbf{R}}, \mathbf{O}_{\theta}))$$



Heuristic Attack

- Heuristic Attack Method
 - It assigns high scores to target items
 - Give a low score to random others
 - It interacts with some popular items
 - Include random attack, average attack, bandwagon attack, and segment attack
 - ...

Heuristic Attack



Heuristic Attack

- Random Attack

- Attacker's Goal: promote certain items availability of being recommended

	Item1	Item2	Item3	Item4	Item5	Item6
User1	4	3	4	-	3	4
User2	5	5	1	4	1	3
User3	1	5	2	5	4	2
User4	5	1	5	3	-	5
User5	3	5	4	4	1	0
User6	-	5	5	4	-	2
Attacker1	1	-	1	1	5	-
Attacker2	-	1	1	1	5	-

high scores to target item

low score to random others

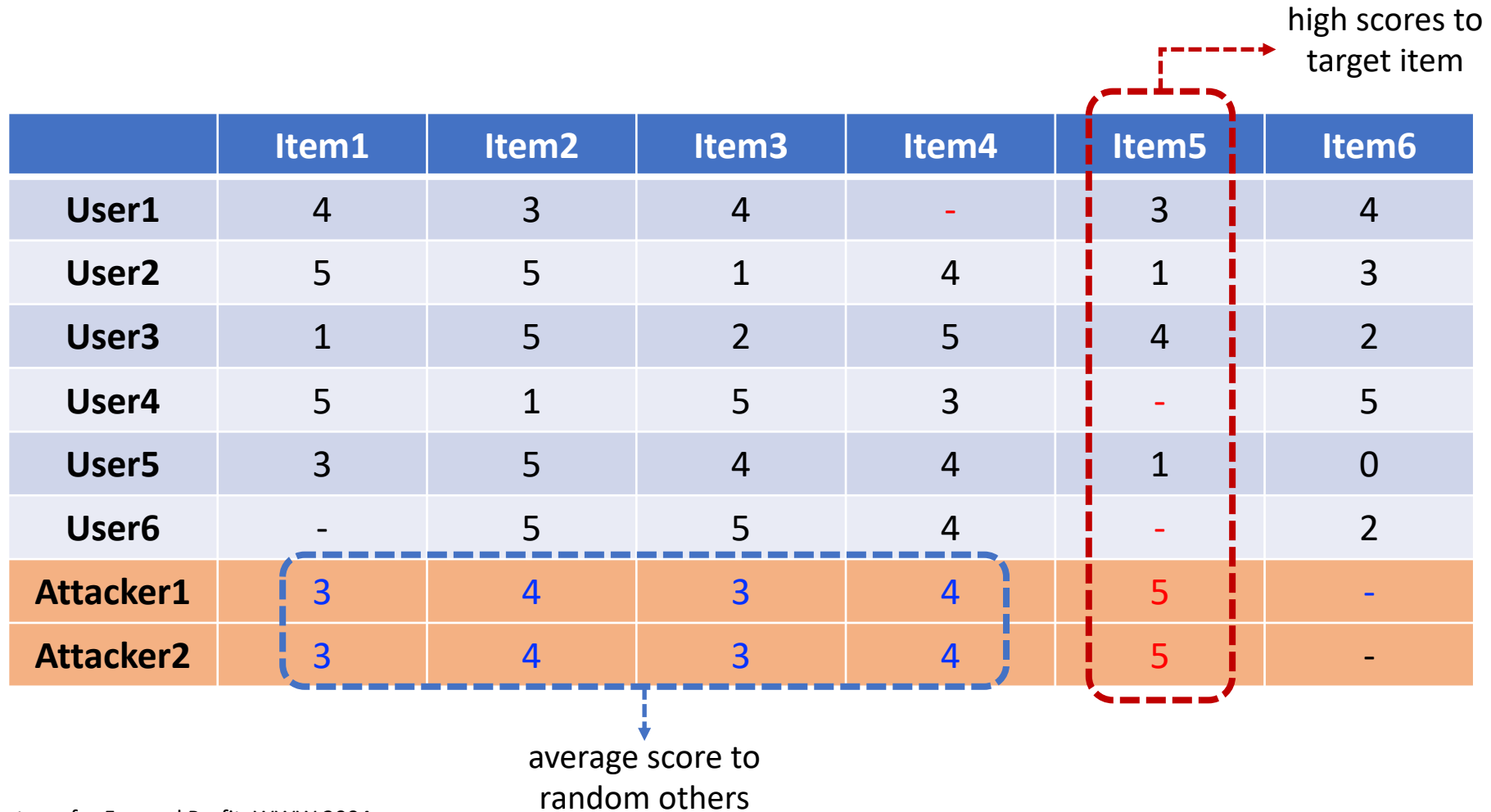
Heuristic Attack

- Average Attack

	Item1	Item2	Item3	Item4	Item5	Item6
User1	4	3	4	-	3	4
User2	5	5	1	4	1	3
User3	1	5	2	5	4	2
User4	5	1	5	3	-	5
User5	3	5	4	4	1	0
User6	-	5	5	4	-	2
Attacker1	3	4	3	4	5	-
Attacker2	3	4	3	4	5	-

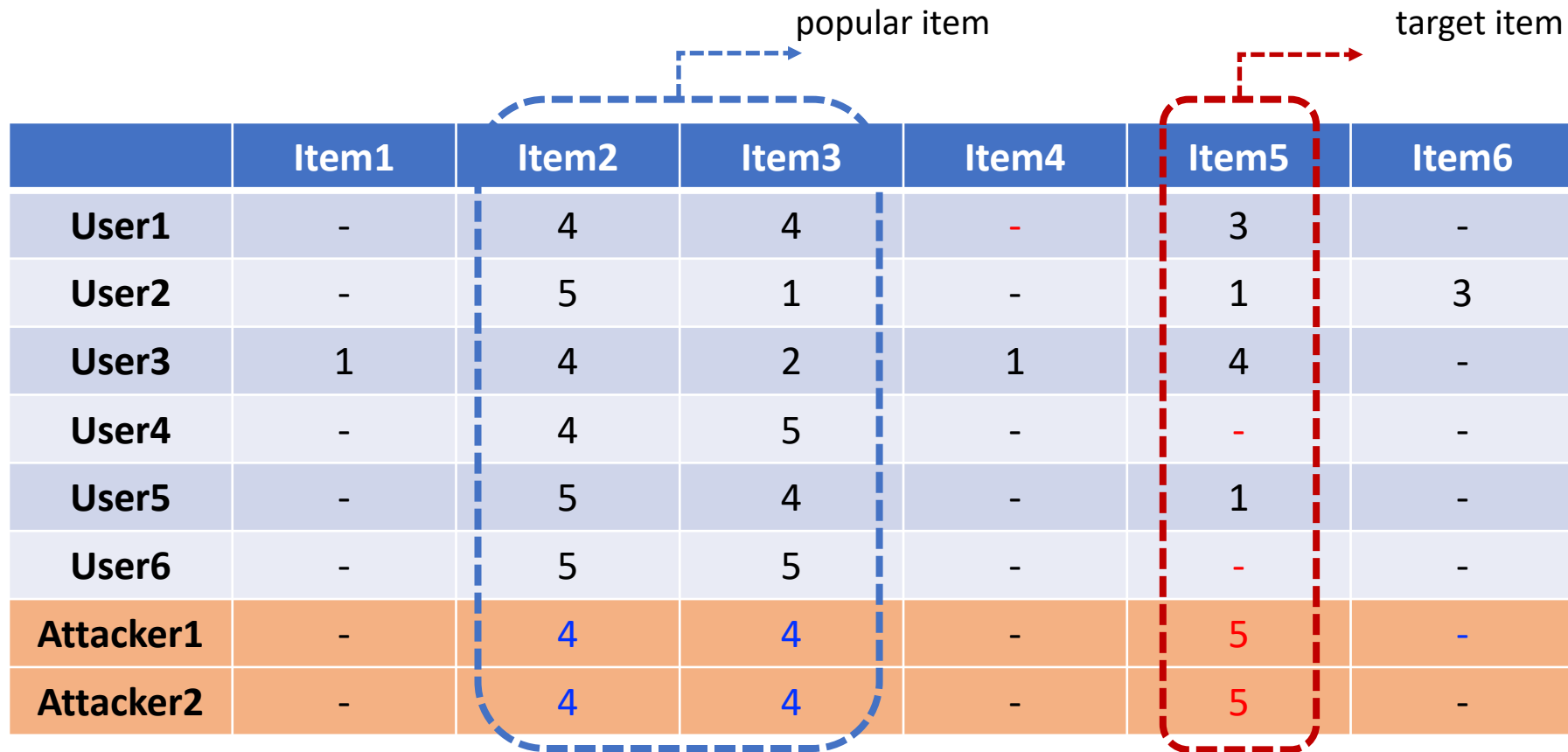
high scores to target item

average score to random others



Heuristic Attack

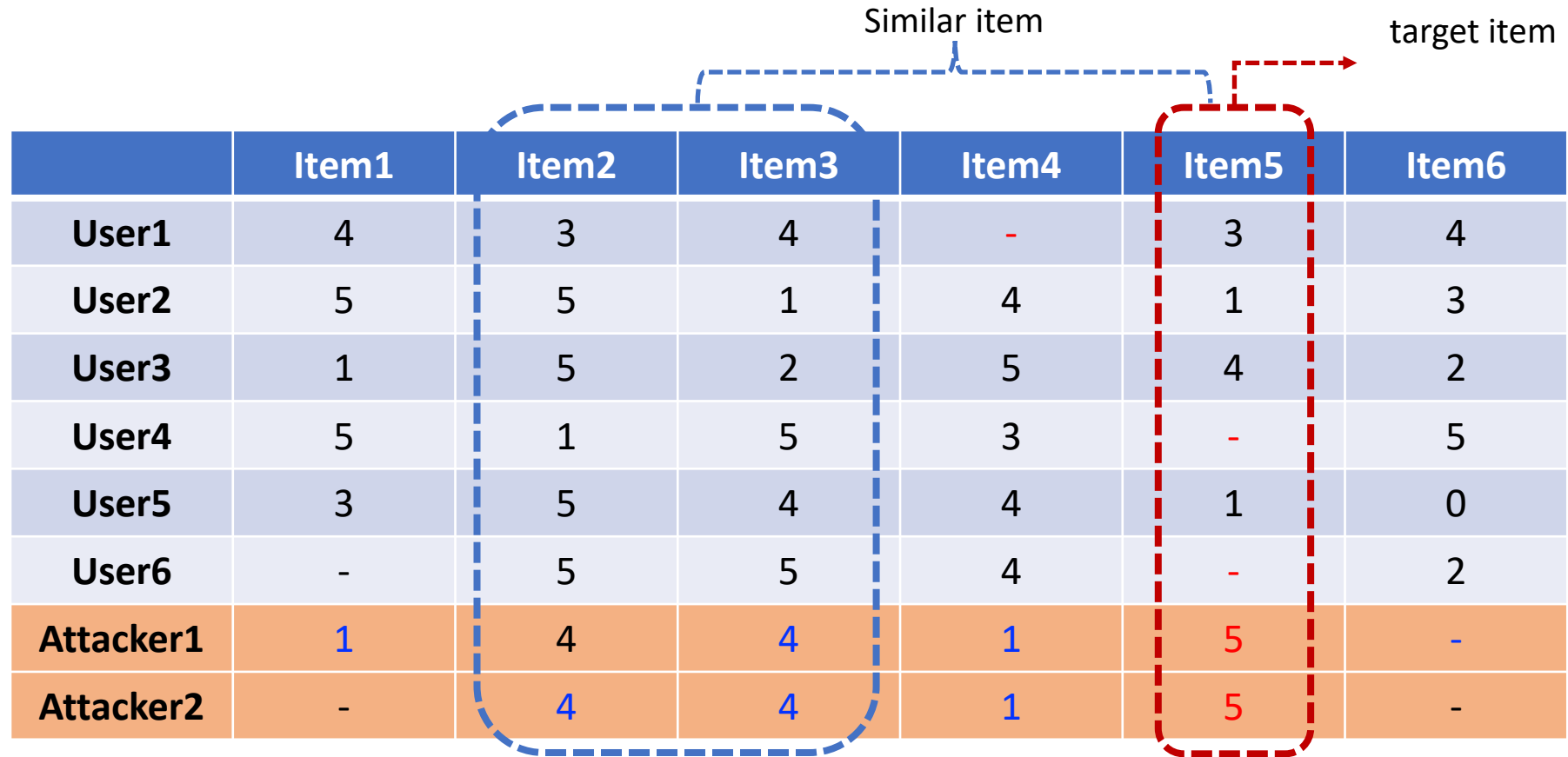
- Bandwagon attack



	Item1	Item2	Item3	Item4	Item5	Item6
User1	-	4	4	-	3	-
User2	-	5	1	-	1	3
User3	1	4	2	1	4	-
User4	-	4	5	-	-	-
User5	-	5	4	-	1	-
User6	-	5	5	-	-	-
Attacker1	-	4	4	-	5	-
Attacker2	-	4	4	-	5	-

Heuristic Attack

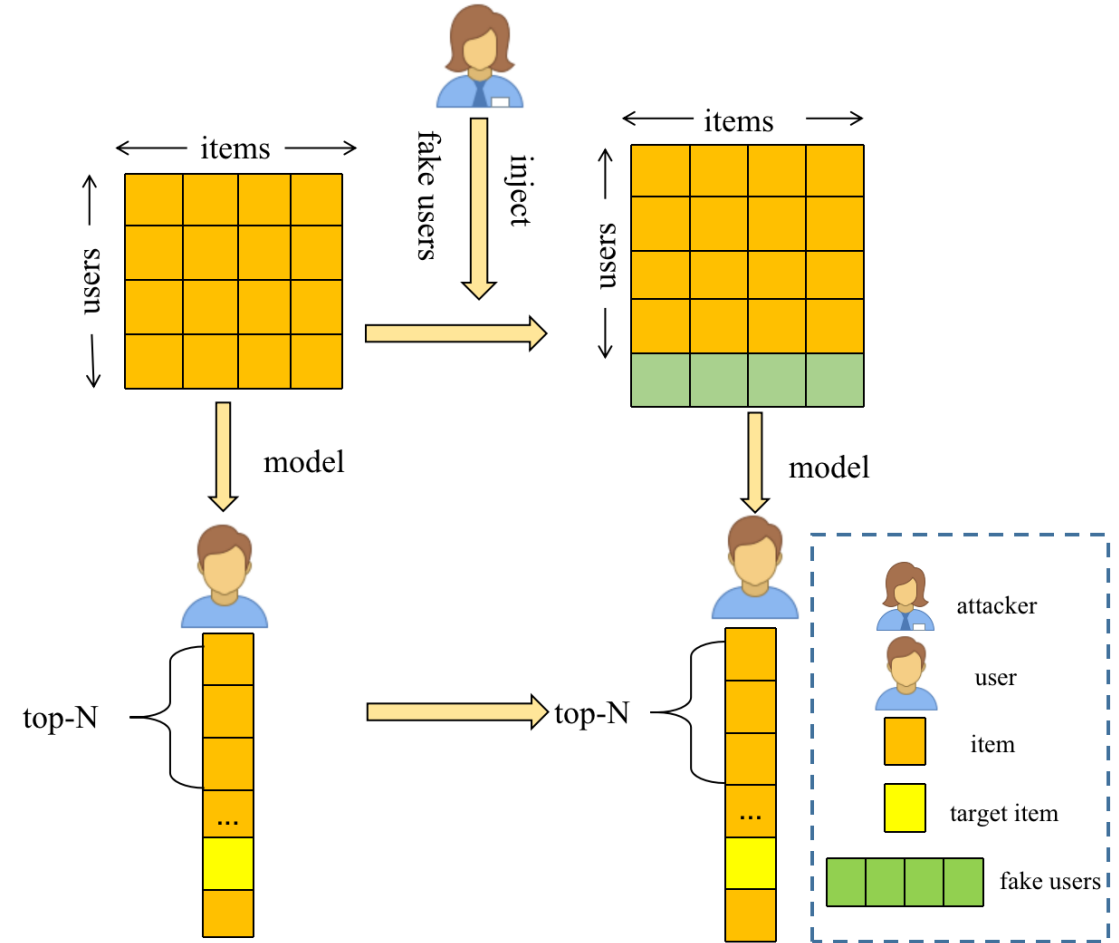
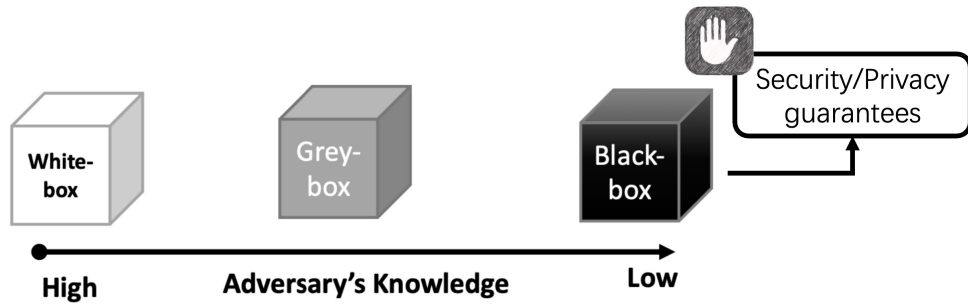
- Segment attack



	Item1	Item2	Item3	Item4	Item5	Item6
User1	4	3	4	-	3	4
User2	5	5	1	4	1	3
User3	1	5	2	5	4	2
User4	5	1	5	3	-	5
User5	3	5	4	4	1	0
User6	-	5	5	4	-	2
Attacker1	1	4	4	1	5	-
Attacker2	-	4	4	1	5	-

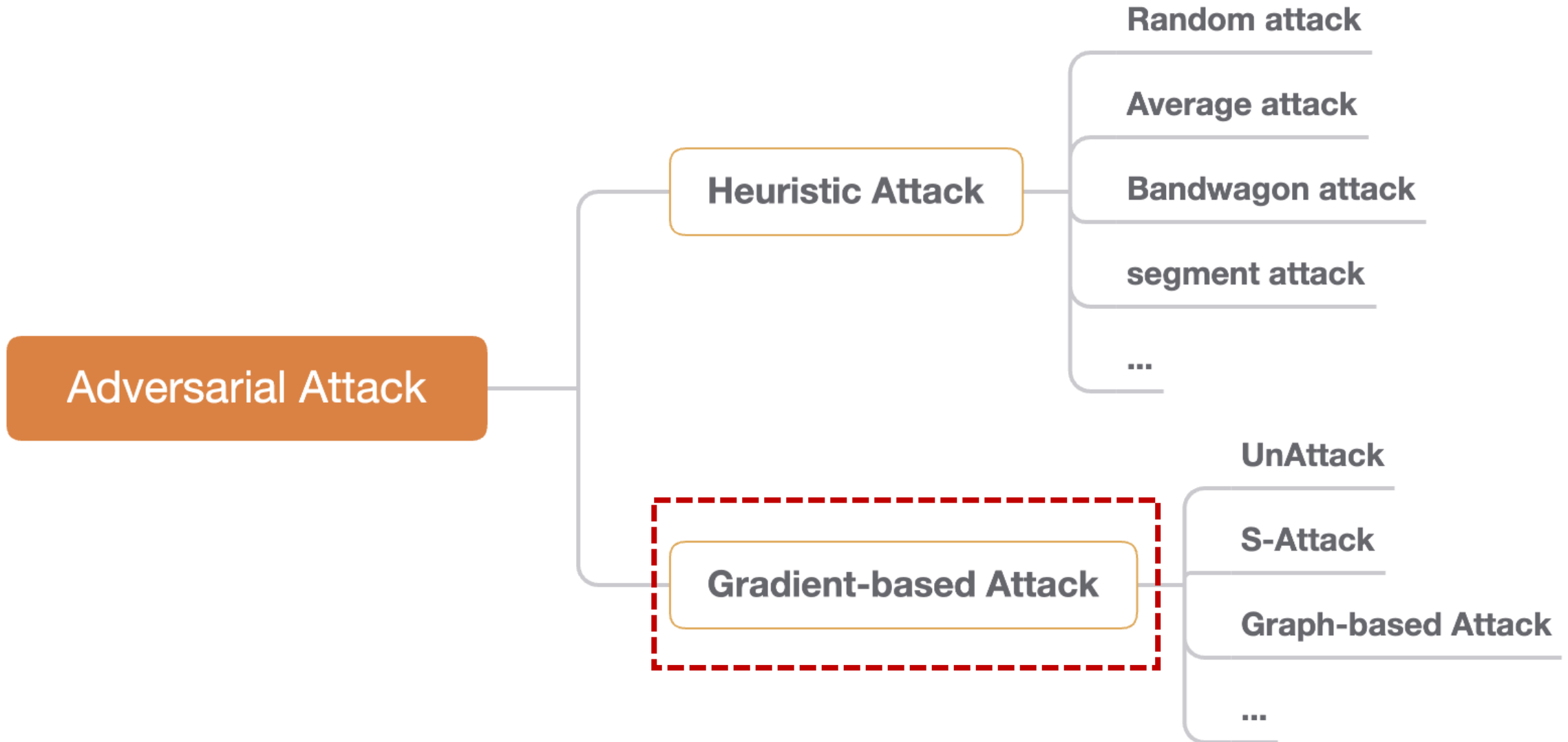
Gradient-based Attack

- Gradient-based Methods
 - White-Box Attack: Optimization



$$\min_{\hat{U}} \mathcal{L}_{adv}(\theta^*), \quad \text{s.t.} \quad \theta^* = \arg \min_{\theta} (\mathcal{L}_{rec}(\mathbf{R}, \mathbf{O}_{\theta}) + \mathcal{L}_{rec}(\hat{\mathbf{R}}, \mathbf{O}_{\theta}))$$

Gradient-based Attack



UNAttack

- UNAttack
 - Optimize the ratings of fake users one by one rather than for all m fake users at the same time
 - Borrow the strategy from the ranking problem to construct pairwise loss function

$$\begin{aligned}
 loss_1 &= \sum_{v \in S(u, K)} \sigma(s_{uv} - s_{uf}) \\
 loss_2 &= \sum_{i \in L_u} \sigma(p_{ui} - p_{ut}) \\
 loss_u &= (1 - \lambda)loss_1 + \lambda loss_2
 \end{aligned}$$

$$p_{ui} = \sum_{v \in S(u, K) \cap U_i^+} s_{uv} X_{vi}$$

Minimize $(F(X_f) = loss)$

s. t. $|X_f| \leq z,$

$X_{fi} \in \{0, 1, \dots, r_{max}\}$

$$loss = \sum_{u \in U_t^-} loss_u \quad loss = \sum_{u \in U_t^-} loss_u$$

Make the fake user be in the top-K nearest neighbours of user,
which can be expressed as $s_{uf} > s_{uv}$.

UNAttack

- UNAttack
 - Choosing the optimal filler-items for fake users

$$X_f^{(t)} = \text{Project}(X_f^{(t-1)} - \eta \frac{\partial F(X_f)}{\partial X_f})$$

where $\text{Project}(x)$ is the project function that cuts each X_{fi} into the range $[0, 1, \dots, r_{max}]$.

$$\frac{\partial F(X_f)}{\partial X_f} = \sum_{u \in U_f} (1 - \lambda) \frac{\partial \text{loss}_1}{\partial X_f} + \lambda \frac{\partial \text{loss}_2}{\partial X_f}$$

Gradient

$$\frac{\partial (\text{loss}_1)}{\partial X_f} = \sum_{v \in S(u, k)} \frac{\partial \sigma(Q)}{\partial Q} \left(\frac{\partial s_{uv}}{\partial X_f} - \frac{\partial s_{uf}}{\partial X_f} \right)$$

$$\frac{\partial (\text{loss}_2)}{\partial X_f} = \sum_{i \in L_u} \sum_{v \in W} \frac{\partial \sigma(P)}{\partial P} \left(\frac{\partial s_{uv} X_{vi}}{\partial X_f} - \frac{\partial s_{uf} X_{ft}}{\partial X_f} \right)$$

similarity

$$\frac{\partial s_{uf}}{\partial X_f} = \frac{X_u}{\|X_u\| \|X_f\|} - \frac{X_u X_f}{\|X_u\| \|X_f\| \|X_f\|^2}$$

UNAttack

- UNAttack

Algorithm 1. UNAttack

Input: Matrix $R_{m \times n}$

Parameter: λ, K, N, z, j

Output: j fake users

- 1: **for** each fake user f **do**
- 2: Solve the problem in Equation 6 with current rating matrix R to get X_f
- 3: Let $X_{ft} = r_{\max}$
- 4: Select z items with highest value in X_{fi} as filler items.
- 5: For each filler-items j , $X_{fj} \sim \mathcal{N}(\mu_j, \sigma_j^2)$
- 6: $R_{m \times n} = R_{m \times n} \cup X_f$
- 7: **end for**

Give the target items the maximum ratings.

Inspired by the ranking problem, all items will be ranked according to X_{fi} , and top- z items with the highest values will be chosen as the filler-items.

The rating score assigned to each filler-item is drawn from a normal distribution of the normal users' rating data of this item.

S-Attack

- Attack matrix factorization based recommender systems
 - Attacker's Goal: promote certain items availability of being recommended
 - Attacker's knowledge: fully (partial) observable dataset
 - Challenge:
 - User ratings are discrete
 - Excessive number of users

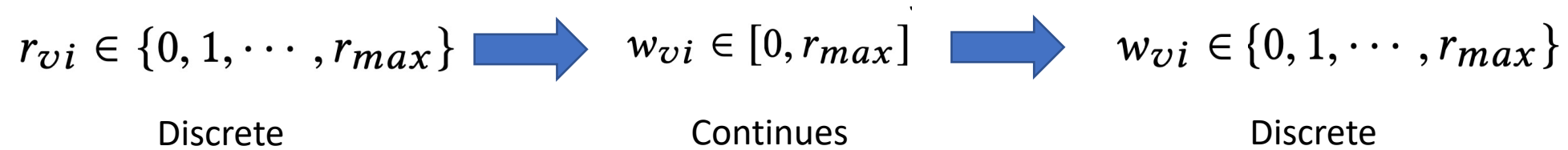
$$\arg \min_{X, Y} \sum_{(u, i) \in \mathcal{E}} (r_{ui} - \mathbf{x}_u^\top \mathbf{y}_i)^2 + \lambda \left(\sum_u \|\mathbf{x}_u\|_2^2 + \sum_i \|\mathbf{y}_i\|_2^2 \right)$$

$$\begin{aligned} & \max h(t) \\ & \text{s.t. } |\Omega_v| \leq n + 1, \quad \forall v \in \mathcal{M}, \\ & \quad r_{vi} \in \{0, 1, \dots, r_{max}\}, \quad \forall v \in \mathcal{M}, \forall i \in \Omega_v. \end{aligned}$$

S-Attack

- Step 1: Optimize one by one
- Step 2: Relax the discrete ratings to continuous

$$\mathbf{w}_v = [w_{vi}, i \in \Omega_v]^\top$$



S-Attack

- Step 3: Approximating the Hit Ratio
- Step 4: Determining the Set of Influential Users

$$\min_{\mathbf{w}_v} \mathcal{L}_{\mathcal{U}}(\mathbf{w}_v) = \sum_{u \in \mathcal{U}} \sum_{i \in \Gamma_u} g(\hat{r}_{ui} - \hat{r}_{ut}) + \eta \|\mathbf{w}_v\|_1$$

s.t. $w_{vi} \in [0, r_{max}]$, Γ_u Top-k list

Influential Users

$$\min_{\mathbf{w}_v} \mathcal{L}_{\mathcal{S}}(\mathbf{w}_v) = \sum_{u \in \mathcal{S}} \sum_{i \in \Gamma_u} g(\hat{r}_{ui} - \hat{r}_{ut}) + \eta \|\mathbf{w}_v\|_1$$

s.t. $w_{vi} \in [0, r_{max}]$.

Graph-Based Attack

- Attack graph-based recommender systems
 - Attack using random walk algorithm

Random walk:

$$p_u = (1 - \alpha) \cdot Q \cdot p_u + \alpha \cdot e_u$$

$$Q_{xy} = \begin{cases} \frac{r_{xy}}{\sum_{z \in \Gamma_x} r_{xz}} & \text{if } (x, y) \in E \\ 0 & \text{otherwise} \end{cases}$$

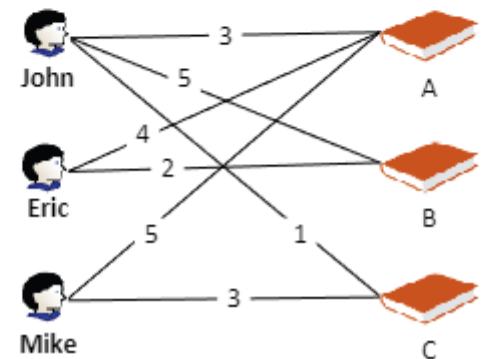
Loss function:

$$l_u = \sum_{i \in L_u} g(p_{ui} - p_{ut})$$

$$g(x) = \frac{1}{1 + \exp(-x/b)}$$

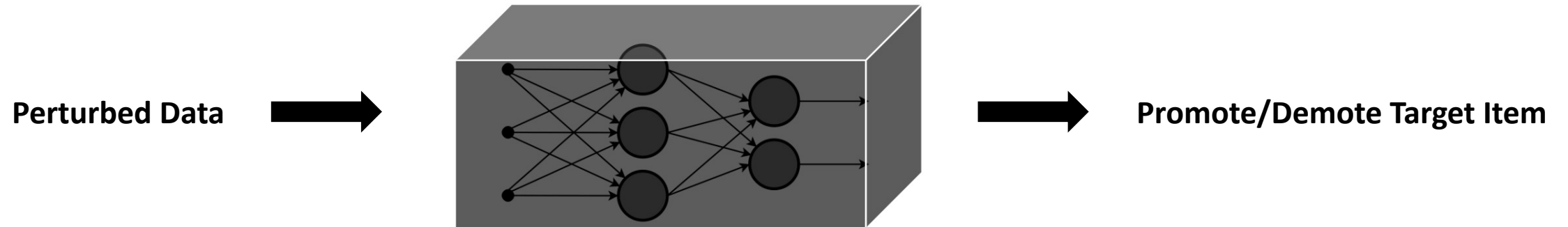


User preference graph



Black-Box Attack

- Black-Box Attack



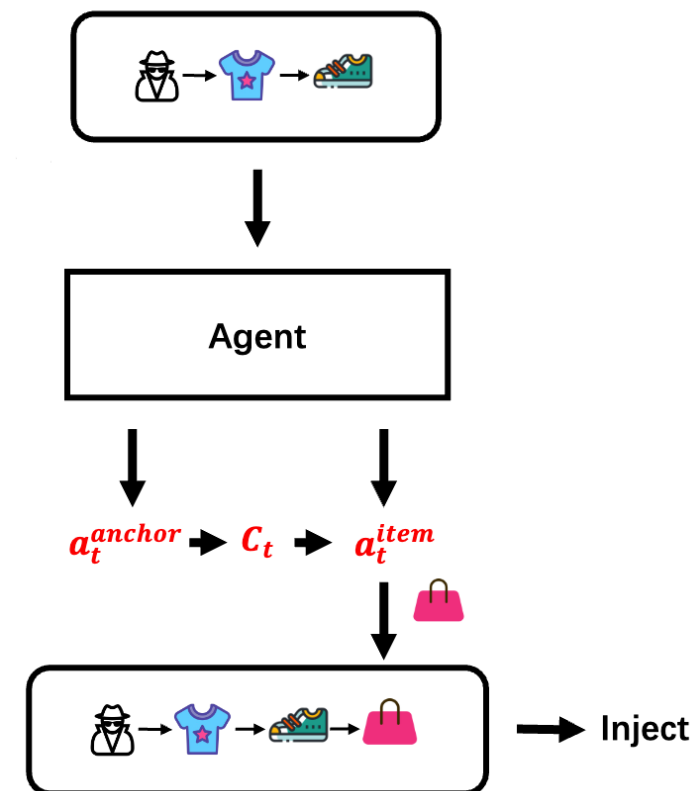
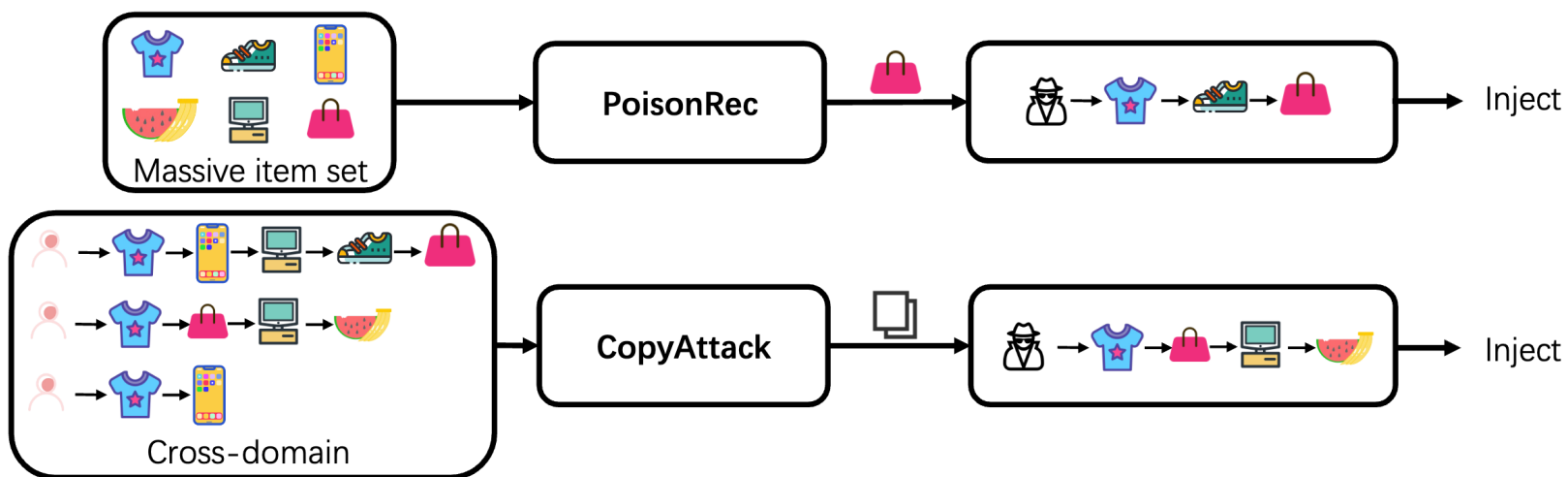
Reinforcement Learning-based Attack

- Challenges in existing attacking methods:
 - Model structure, parameters and training data are unknown
 - Unable to get user-item interactions
 - **Black-box setting**
 - Reinforcement Learning (RL) -- Query Feedback (Reward)

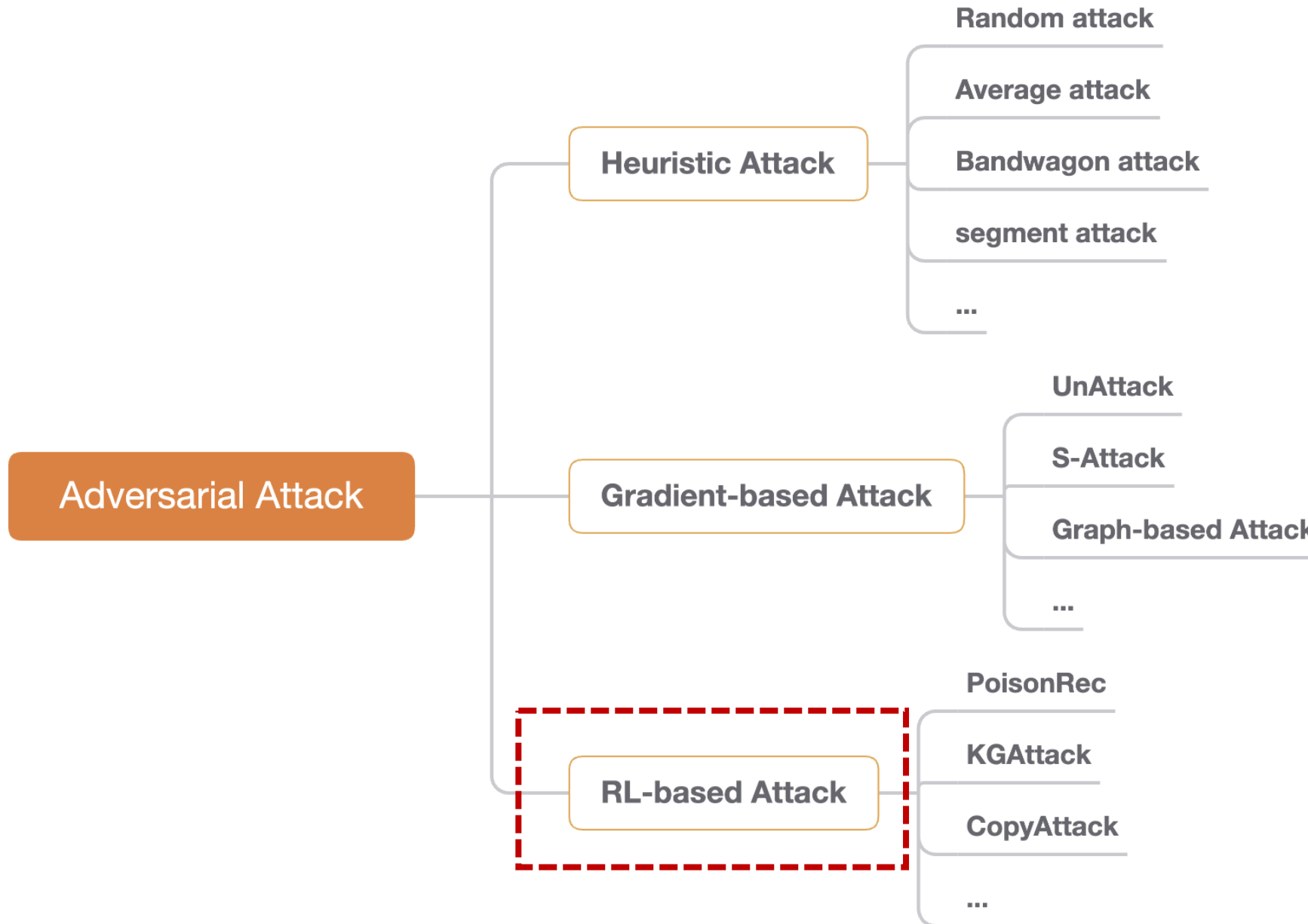
Reinforcement Learning-based Attack

- Reinforcement Learning-based Methods

- PoisonRec
- KGAttack
- CopyAttack



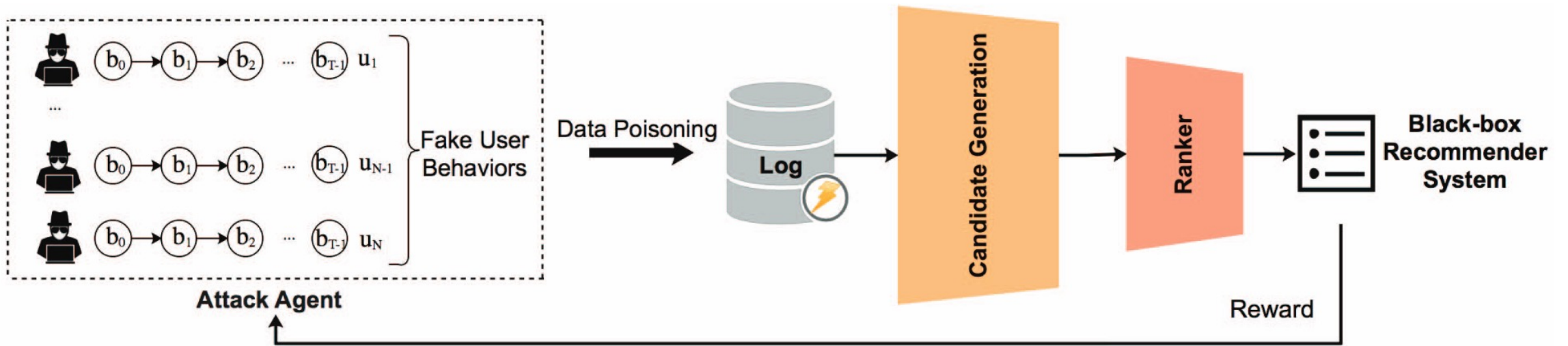
Reinforcement Learning-based Attack



PoisonRec

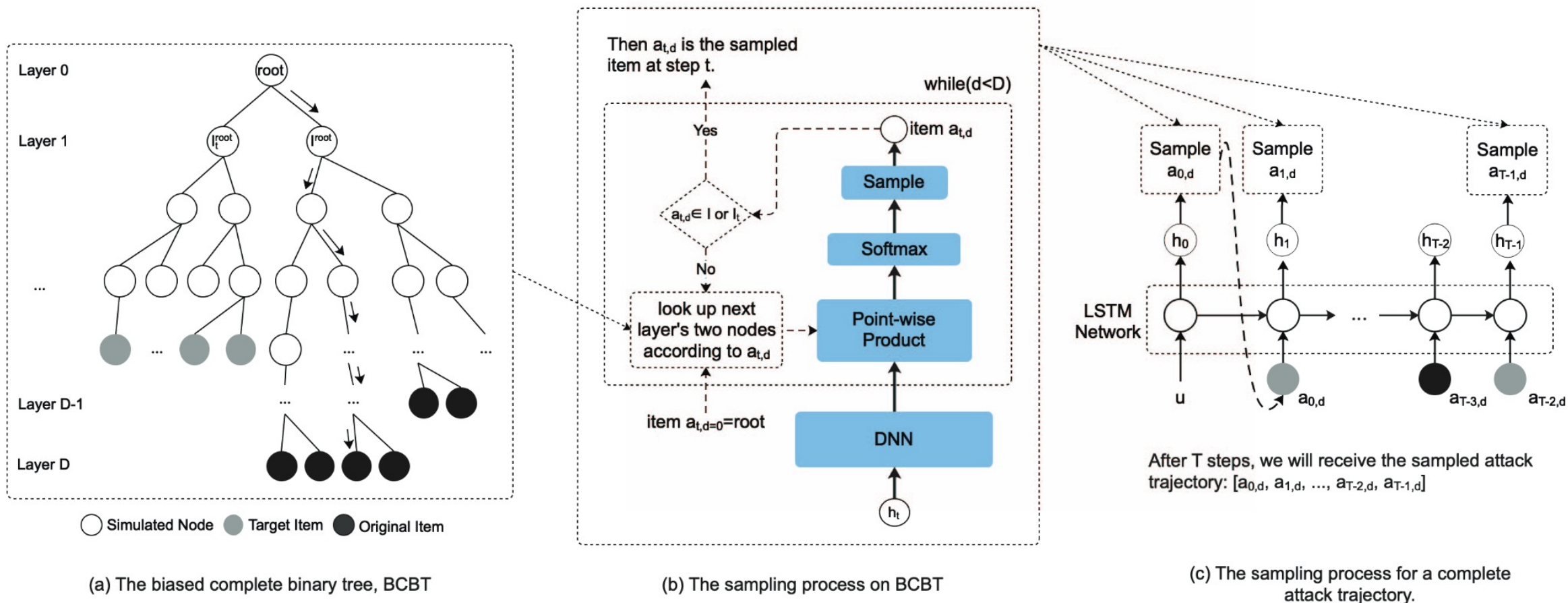
- Target: $RecNum = \sum_u |L_u \cap I_t|$

- DNN + PPO



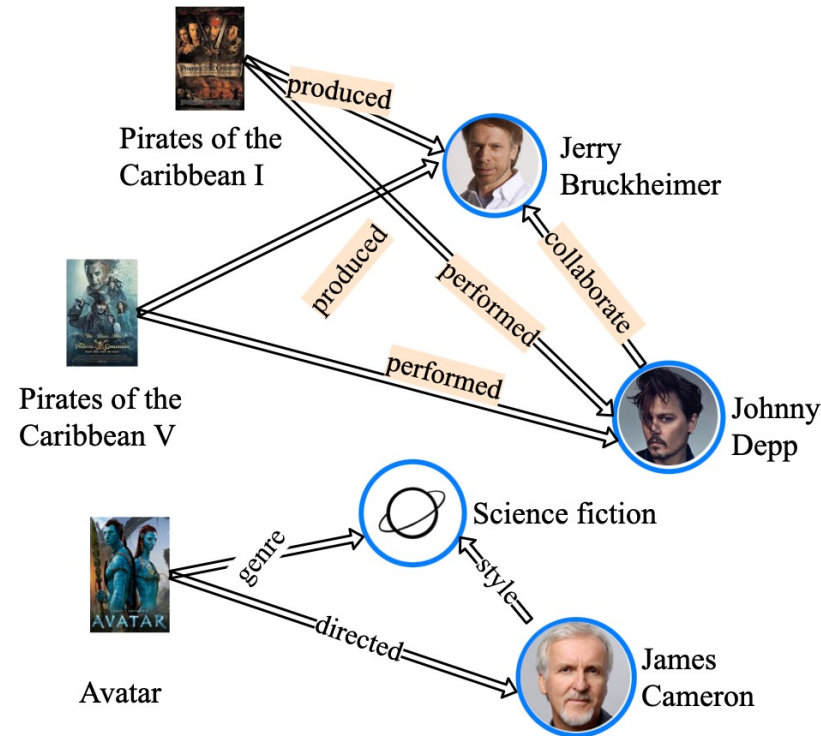
PoisonRec

- Introduce (Biased Complete Binary Tree) BCBT to reduce action space



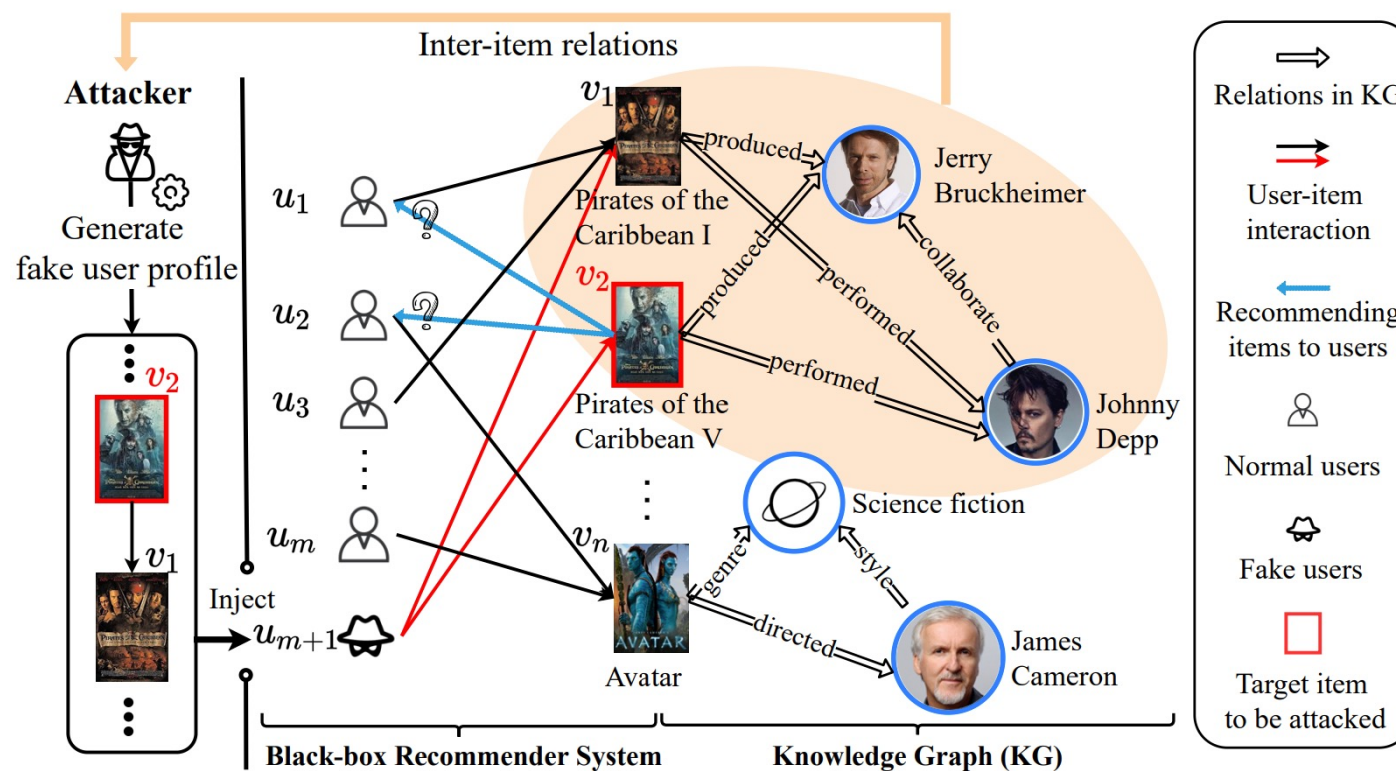
KGAttack

- Side-information: Knowledge Graph (KG)
 - Rich auxiliary knowledge: relations among items and real-world entities
 - The underlying relationships between **Target items** and other items



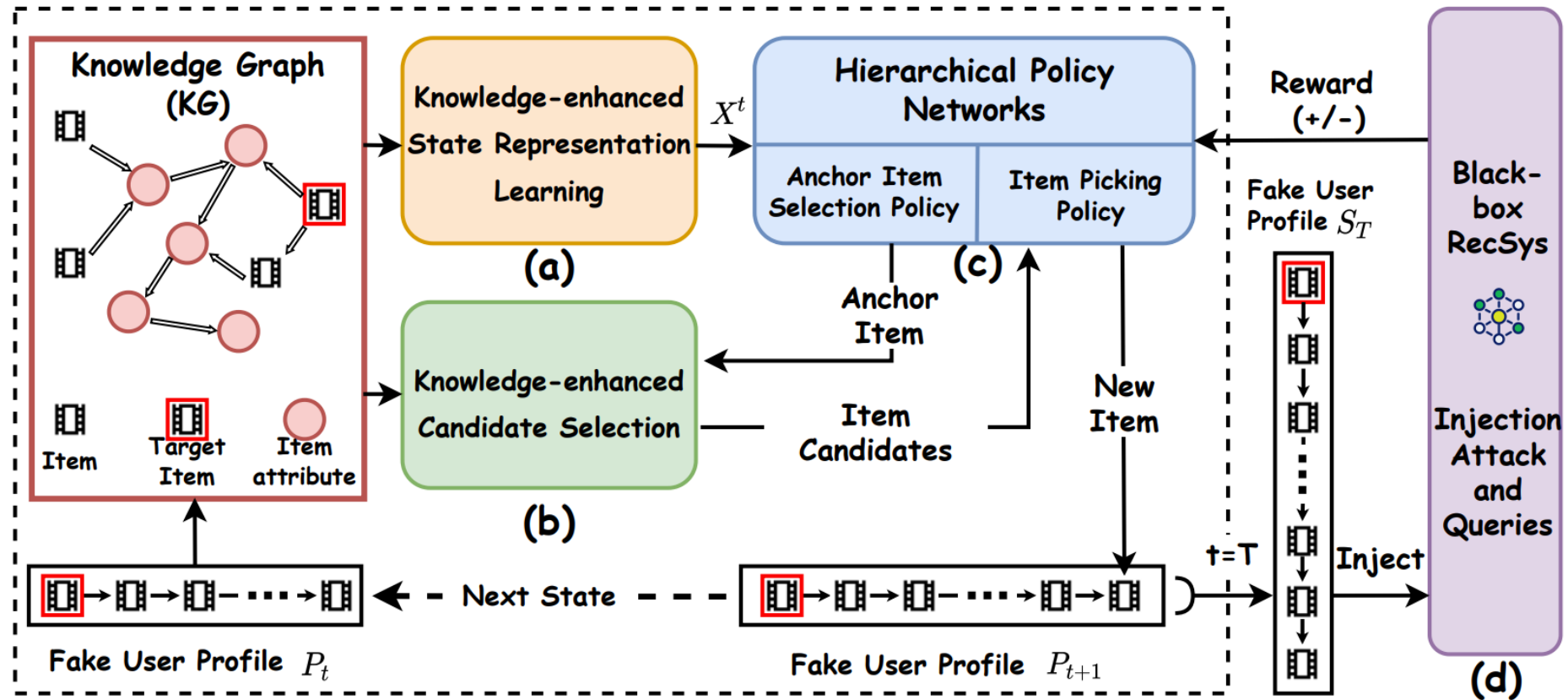
KGAttack

- Employs the KG to enhance the generation of fake user profiles from the massive item sets



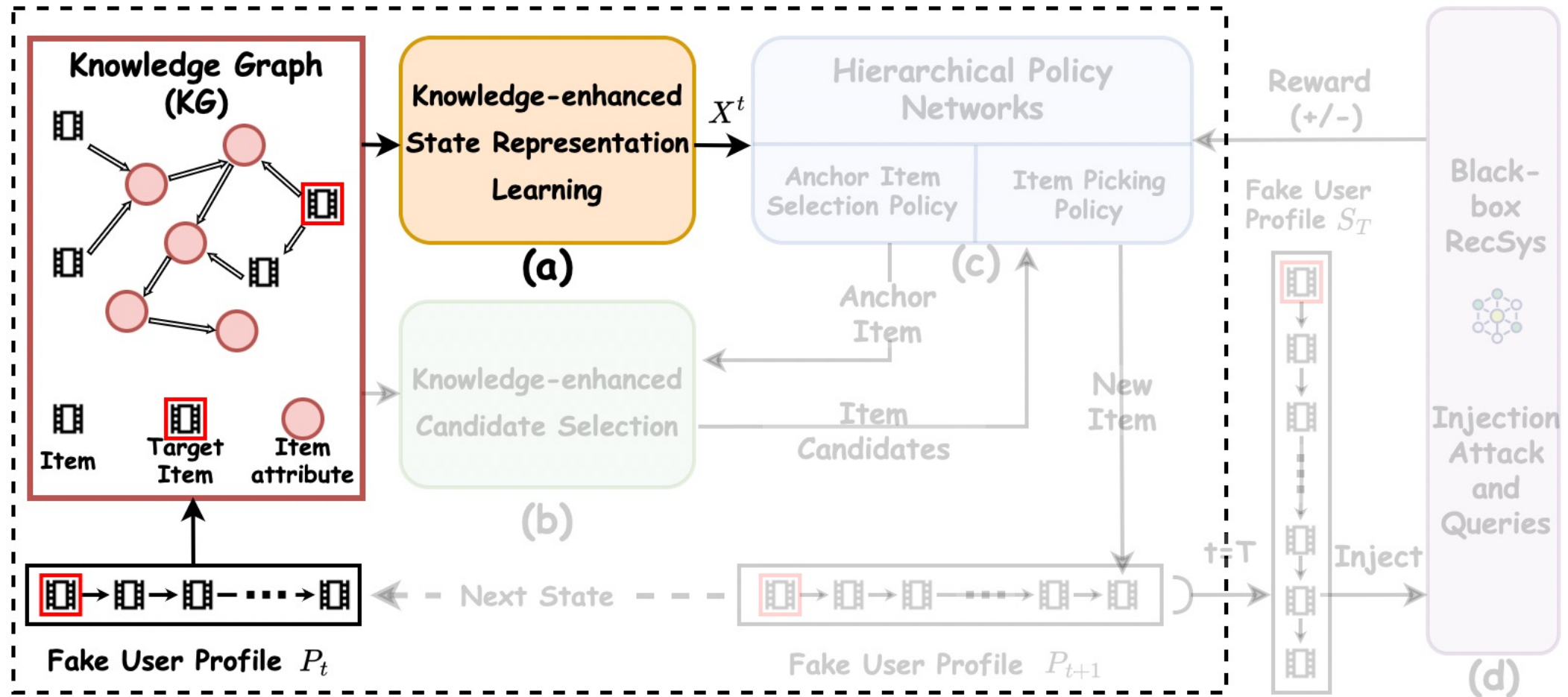
KGAttack

- Using KG to enhance the representation of state
- RL agent, generate user profiles



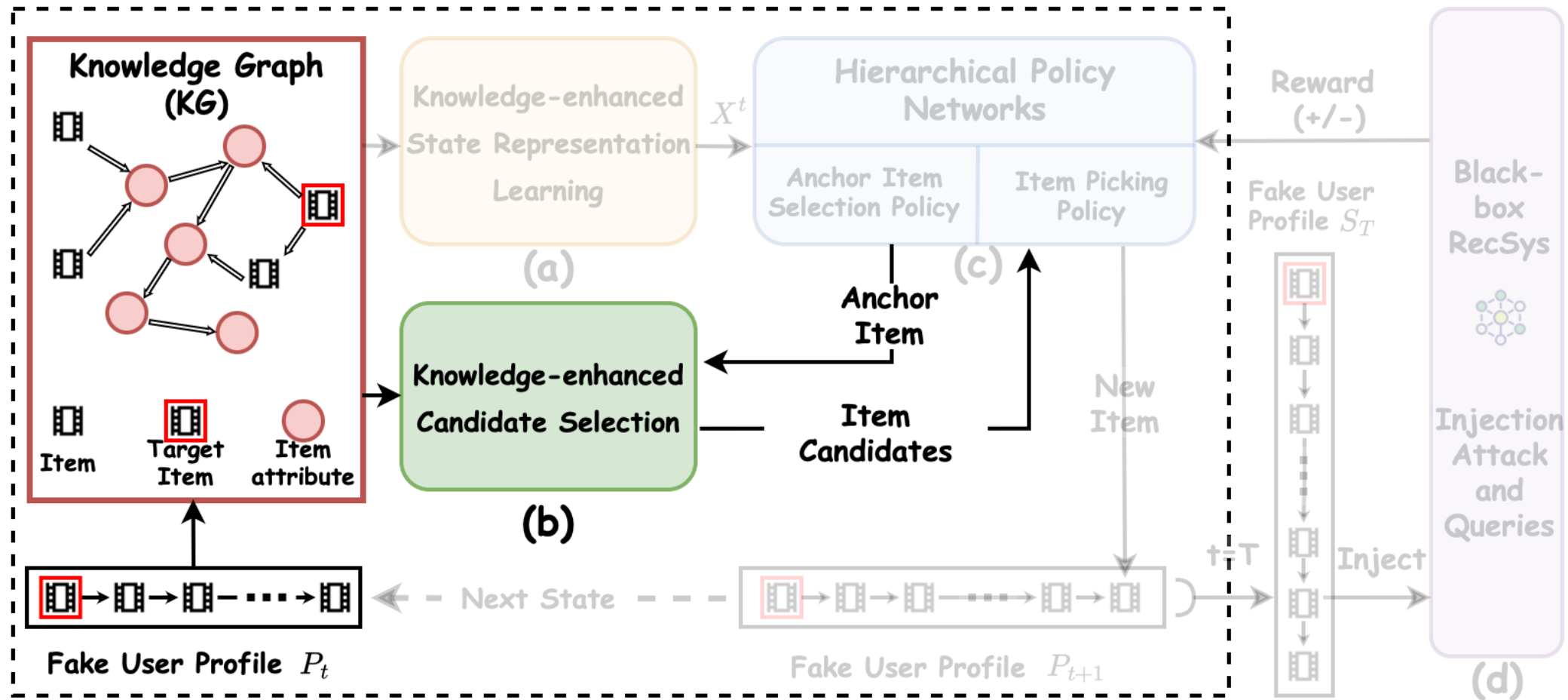
KGAttack

- (a): Using **KG** to enhance the representation of state



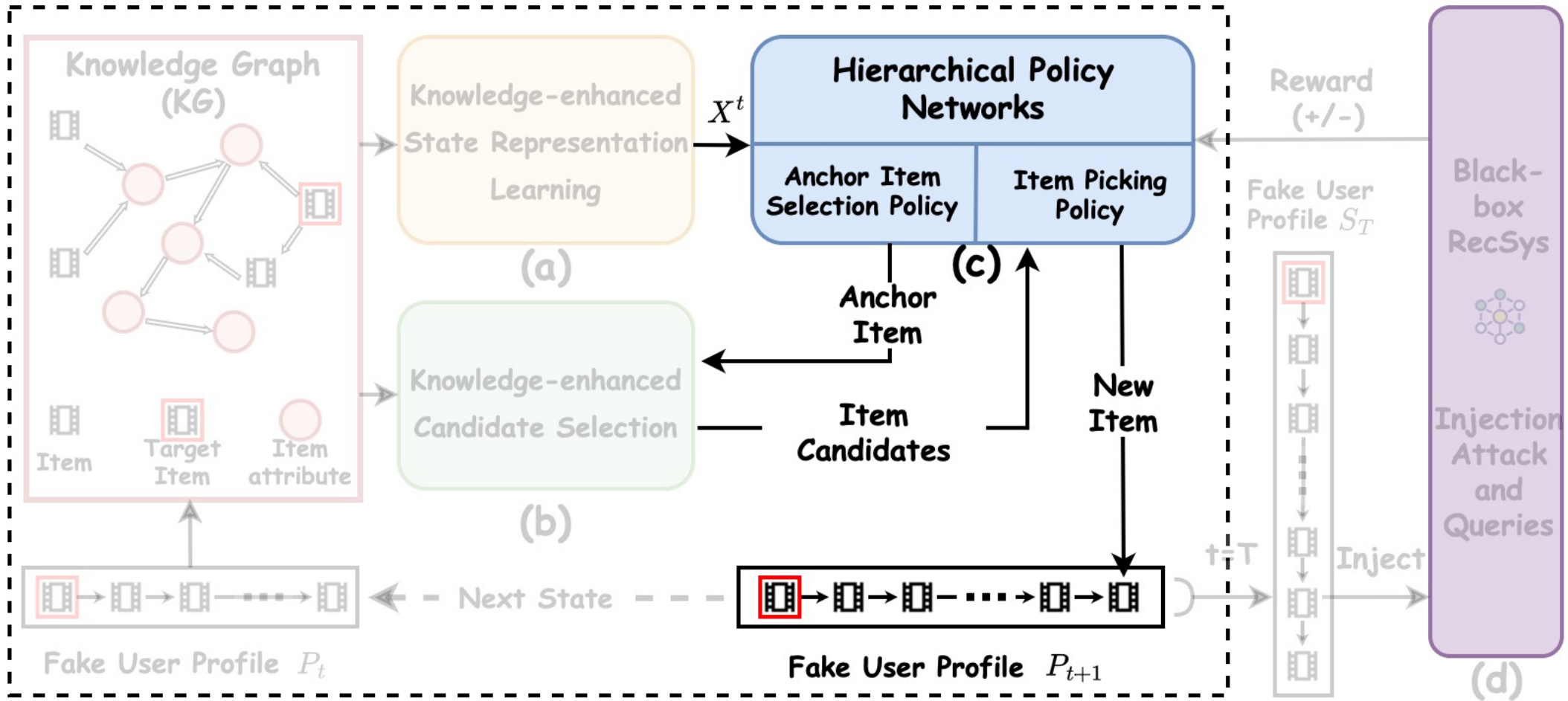
KGAttack

- (b): Using **KG** to localize relevant item candidates



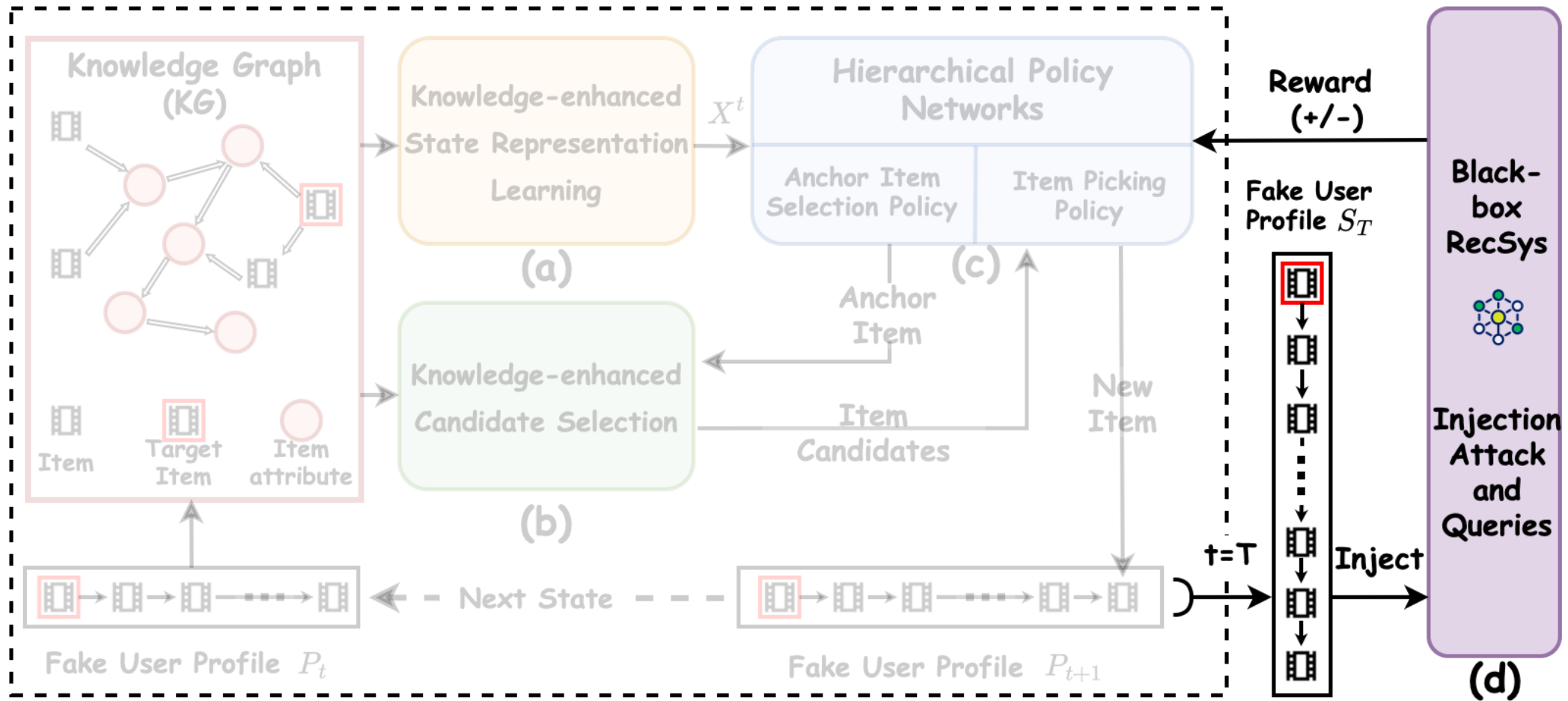
KGAttack

- (c): Using **KG** to localize relevant item candidates



KGAttack

- (d): Injection attacks and query

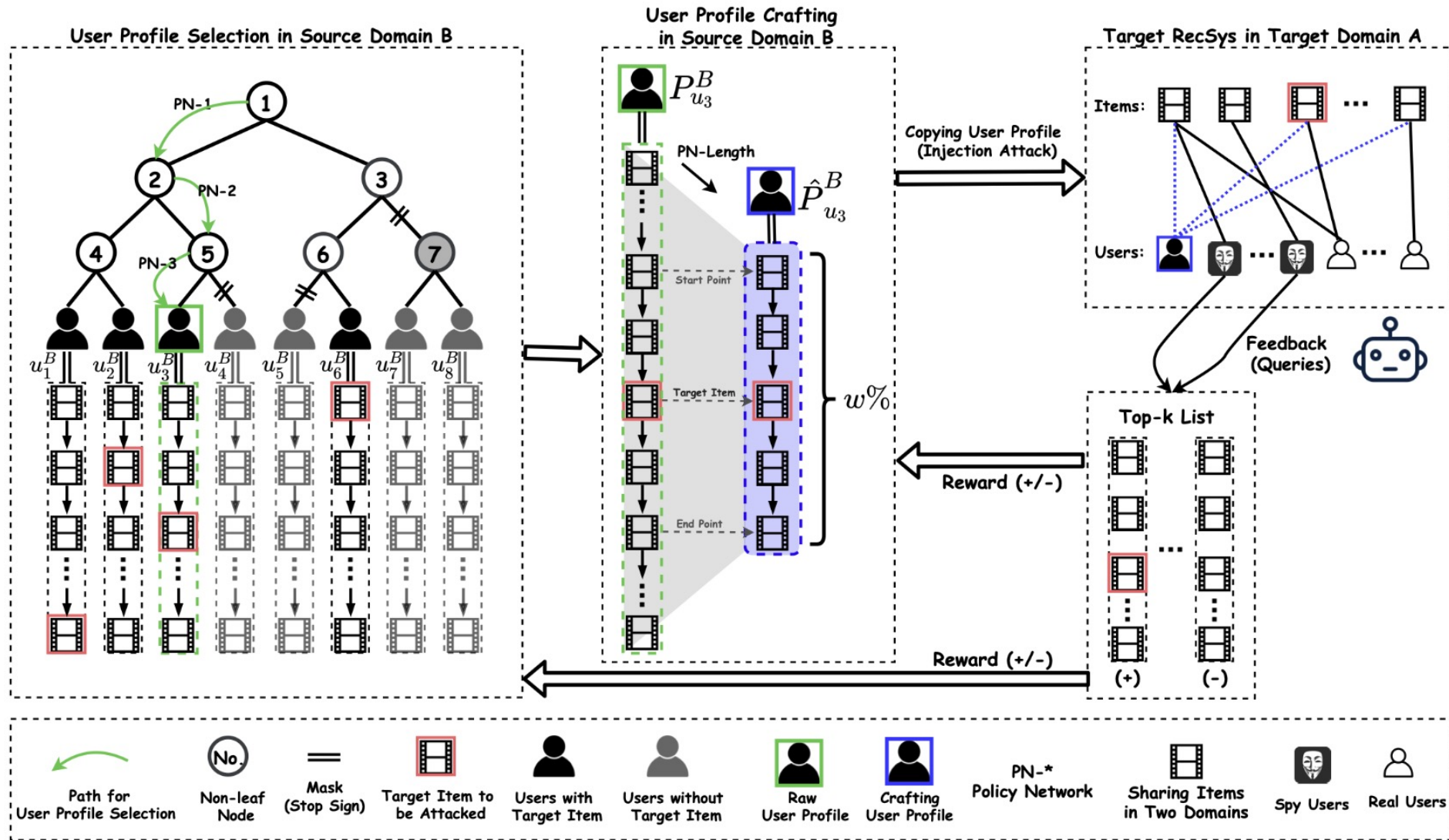


CopyAttack

- Cross-domain Information
 - Share a lot of items
 - Users from these platforms with similar functionalities also share similar behavior patterns/preferences



CopyAttack



CopyAttack

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

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

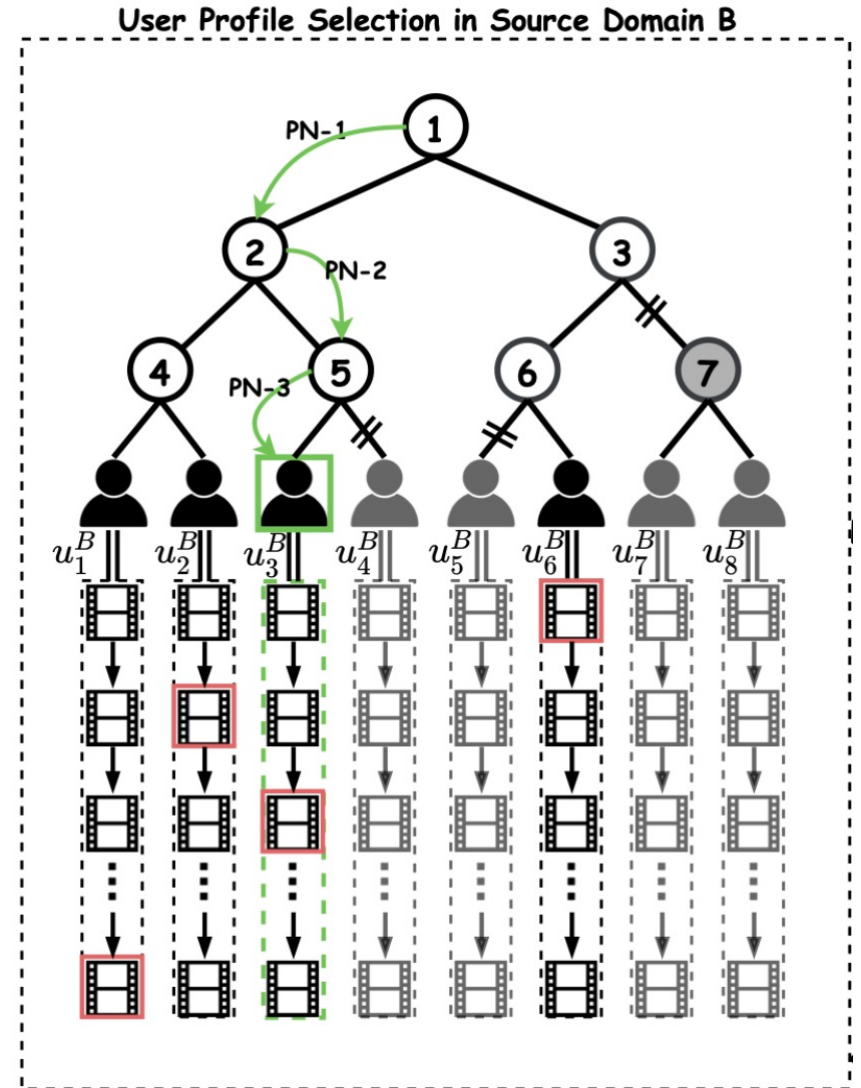
$$p^u(a_t^u | s_t^u) = \prod_d p_d^u(a_t^u | \cdot, s_t^u)$$

$$= p_d^u(a_{[t,d]}^u | s_t^u) \cdot p_{d-1}^u(a_{[t,d-1]}^u | s_t^u) \cdots p_1^u(a_{[t,1]}^u | s_t^u)$$

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

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

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



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)

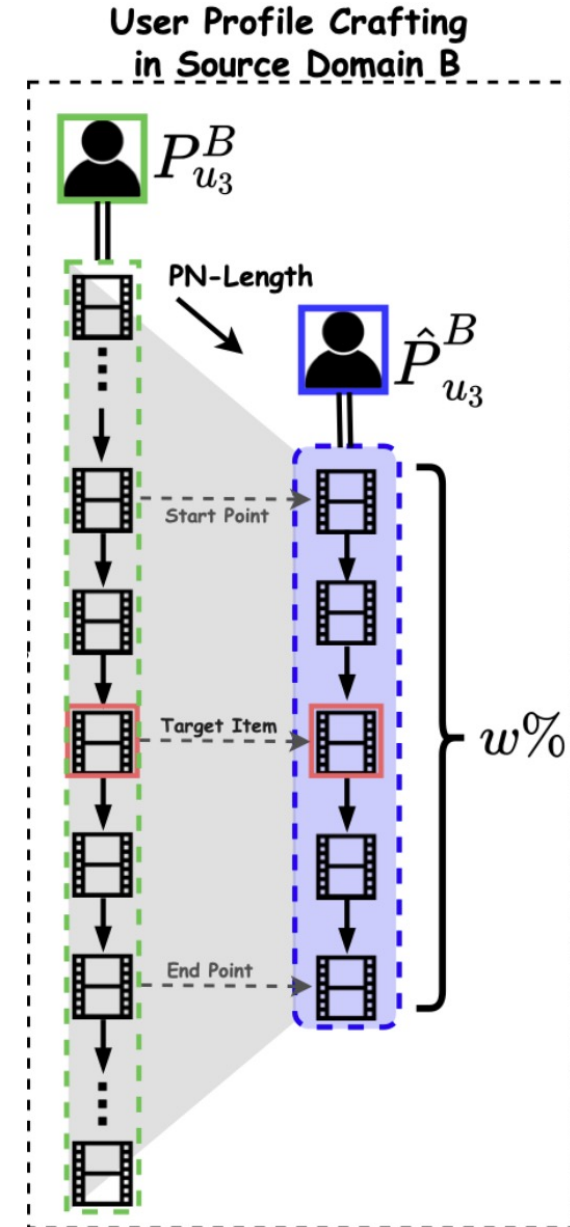
Example:

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

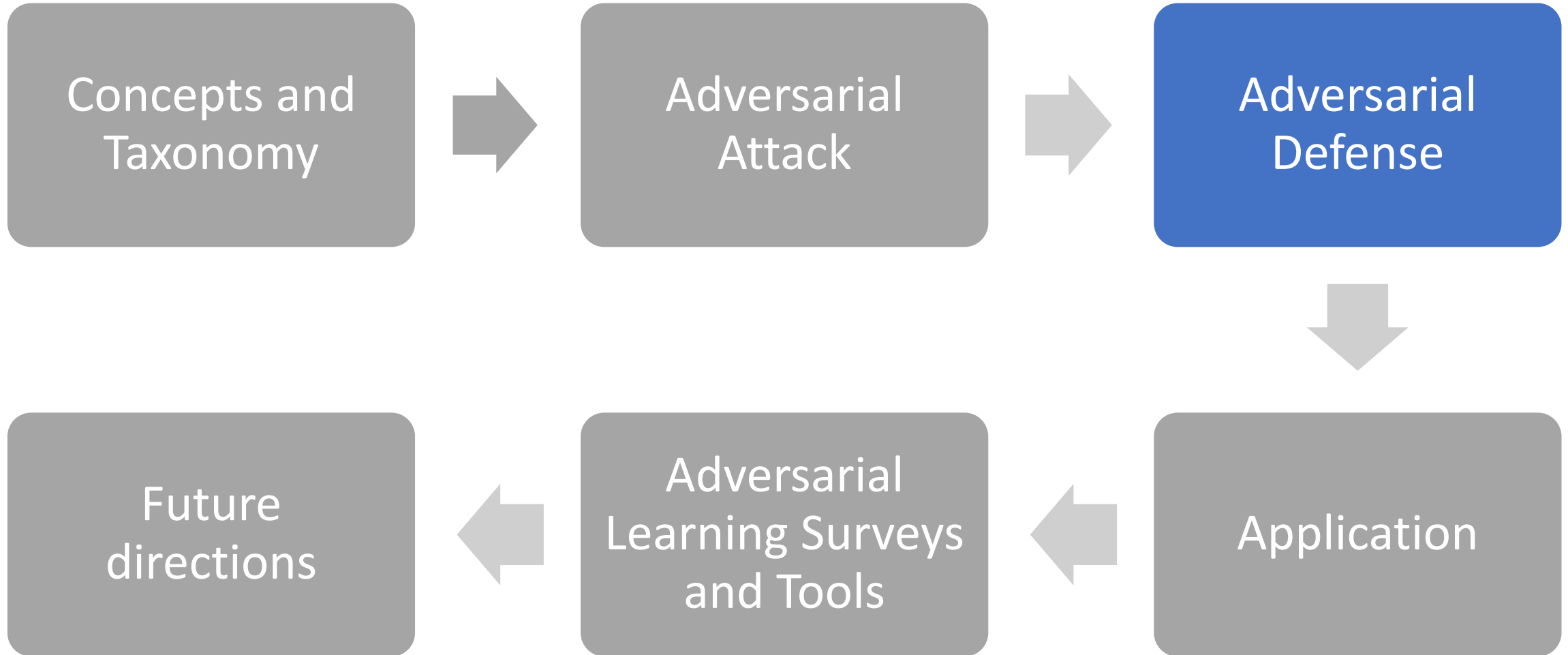
$w = 50\%$

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

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



Outline



Detection

- Exceptions and outliers in the recommendation system
 - Discrepancies between user's ratings and item's average ratings
 - Spectrum-based features of series rate values of each user
 - Cluster instances
 - User behaviors
 - The process of learning users and items representations
 - The distribution of normal users' behaviors over a partial dataset
 - ...

Detection

Adversarial Defense

Detection

DegSim and RDMA

PPu and Du

TSGR, RSF, and TBR

...

Detection

- Detection of shilling attacks in online recommender systems
 - Detecting Process:
 - Extract the supposed characteristics, DegSim and RDMA

Degree of similarity with Top Neighbors:

$$\text{Degsim}_u = \frac{\sum_{v=1}^k W_{u,v}}{k}$$

Rating Deviation from Mean Agreement:

$$\text{RDMA}_j = \frac{\sum_{i=0}^{N_j} \frac{|r_{i,j} - \text{Avg}_i|}{NR_i}}{N_j}$$

Detection

- Detection of shilling attacks via selecting patterns analysis
 - Detecting Process:
 - Extract the supposed characteristics, popularity profile and popularity distribution

A set of item popularity values of rated items:

$$P P_u = (d_{u,1}, d_{u,2}, \dots, d_{u,N_u})$$

Popularity distribution:

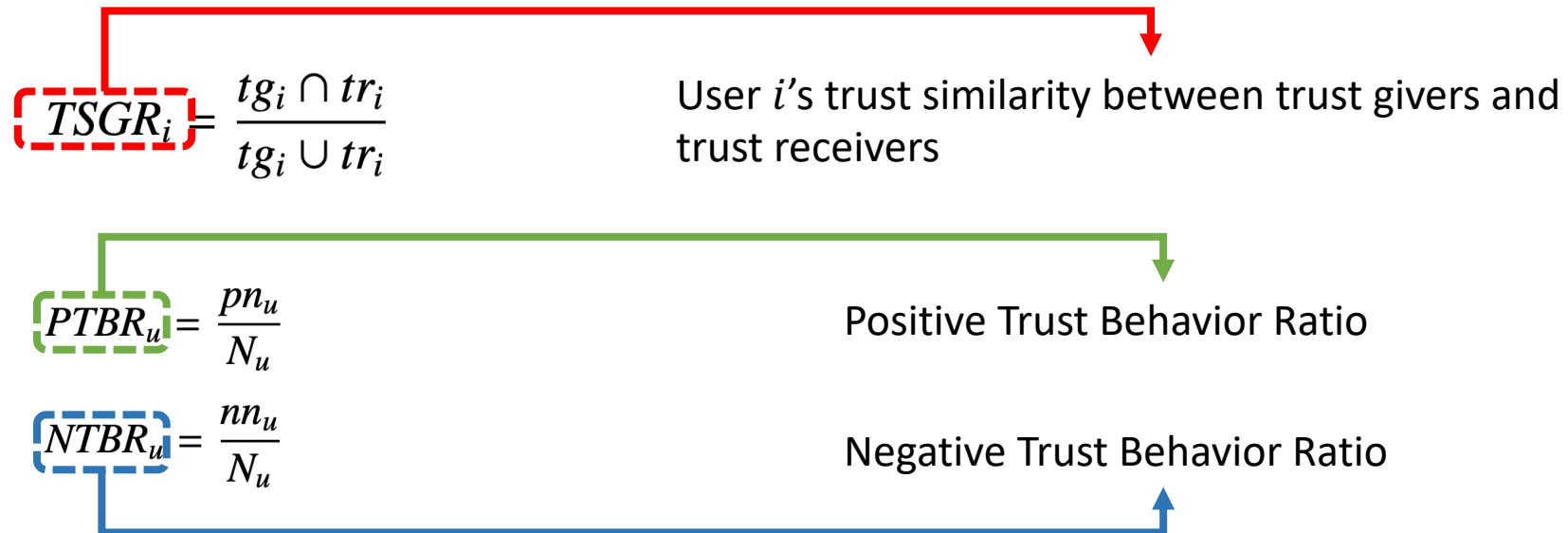
$$D_u = (p_{u,1}, p_{u,2}, \dots, p_{u,d_{\max}})$$

Detection

- Detection of trust shilling attacks in recommender systems

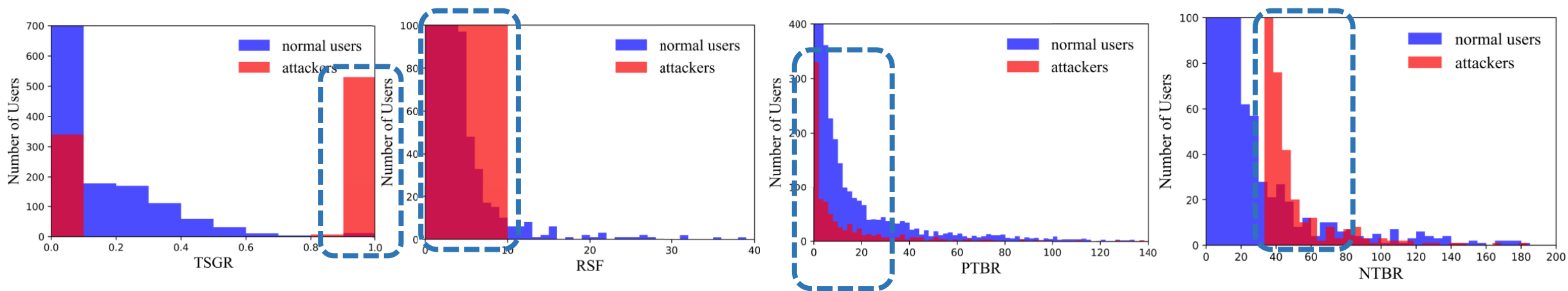
- Detecting Process:

- Extract the supposed characteristics, TSGR, RSF, and TBR



Detection

- Normal vs. attackers distributions for each feature:

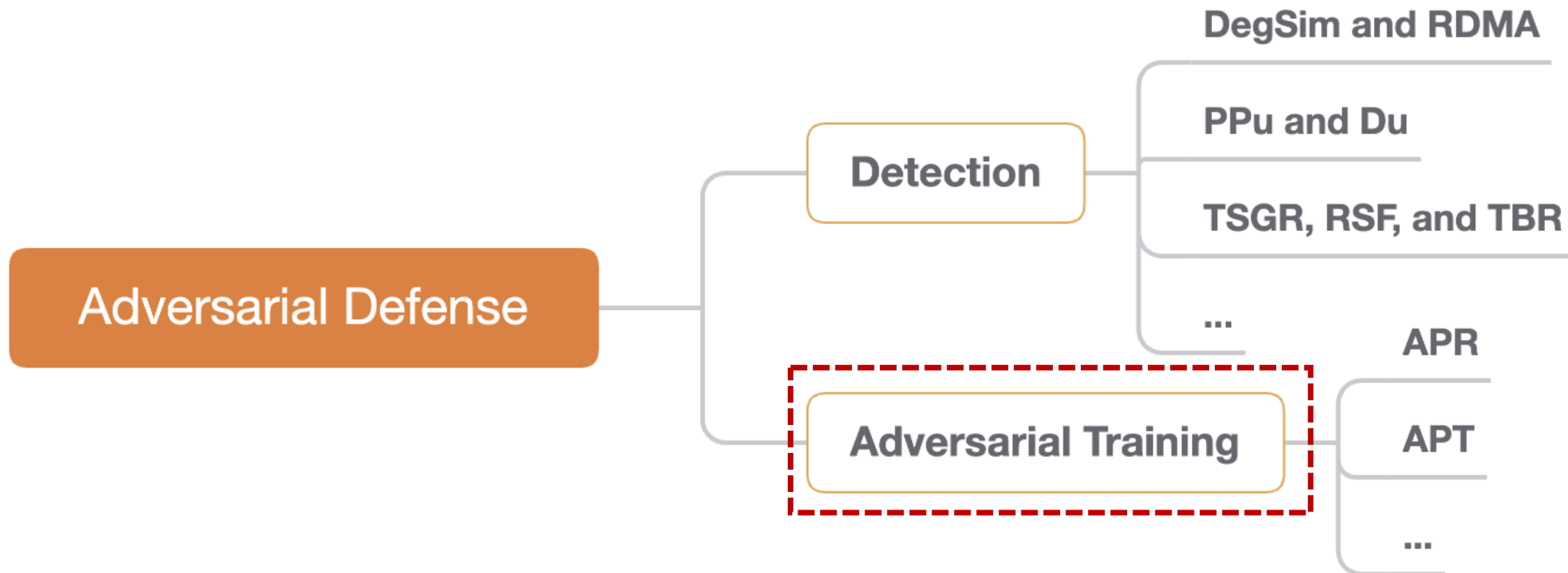


Adversarial Training

- Adversarial training contains two alternating processes:
 - Generating perturbations that can confuse a recommendation model
 - Training the recommendation model along with generated perturbations

$$\min_{\theta} \max_{\eta} \mathcal{L}(\mathcal{X} + \eta, \theta)$$

Adversarial Training



Adversarial Training

- Adversarial Personalized Ranking (APR)

Optimization objectives against noise:

$$\Delta_{adv} = \arg \max_{\Delta, \|\Delta\| \leq \epsilon} L_{BPR}(\mathcal{D} | \hat{\Theta} + \Delta)$$

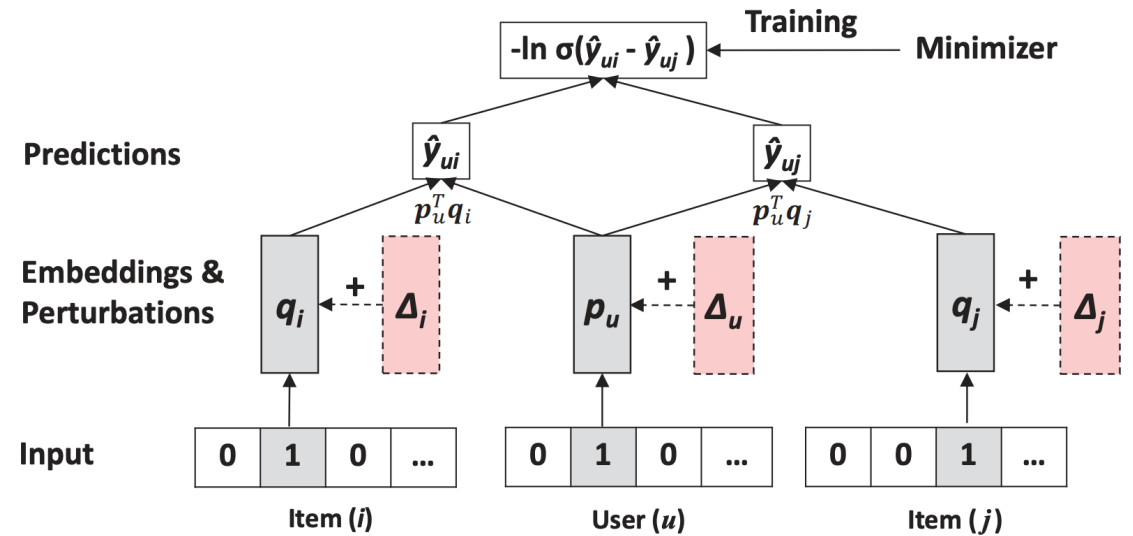
Adversarial Personalized Ranking (APR):

$$L_{APR}(\mathcal{D} | \Theta) = L_{BPR}(\mathcal{D} | \Theta) + \lambda L_{BPR}(\mathcal{D} | \Theta + \Delta_{adv})$$

where $\Delta_{adv} = \arg \max_{\Delta, \|\Delta\| \leq \epsilon} L_{BPR}(\mathcal{D} | \hat{\Theta} + \Delta)$

The training process of APR:

$$\Theta^*, \Delta^* = \arg \min_{\Theta} \max_{\Delta, \|\Delta\| \leq \epsilon} L_{BPR}(\mathcal{D} | \Theta) + \lambda L_{BPR}(\mathcal{D} | \Theta + \Delta)$$



Adversarial Training

- Adversarial poisoning training (APT)

$$\min_{\theta_R} \min_{\mathcal{D}^* | n^*} \mathcal{L}(\mathcal{D} \cup \mathcal{D}^*, \theta_R)$$

$\mathcal{D}^* = \{r_1^*, \dots, r_{n^*}^*\}$ is a set of n^* fake users dedicated to minimizing the empirical risk.

Algorithm 1: Adversarial Poisoning Training

Input: The epochs of training T , pre-training T_{pre} , and poisoning interval T_{inter} .

- 1 Randomly initialize the user set \mathcal{D}^* defined in Definition 3.1. ①
- 2 **for** T_{pre} epochs **do** ②
- 3 Do standard training on the dataset \mathcal{D} ;
- 3 **end**
- 4 $\mathcal{D}' = \mathcal{D}$;
- 5 **for** $T - T_{pre}$ epochs **do**
- 6 **for** per T_{inter} epochs **do** ③
- 7 Calculate the influence vector \mathcal{I} according to Eq. 5;
- 8 **for** each ERM user in \mathcal{D}^* **do** ④
- 9 Select m^* items in Φ with probability $\frac{\exp(-tI_i)}{\sum_{j \in \Phi} \exp(-tI_j)}$ and rate the selected items with normal distribution $(\mu_i + r^+, \sigma_i)$ at random;
- 10 **end**
- 11 $\mathcal{D}' = \mathcal{D} \cup \mathcal{D}^*$; ⑤
- 12 **end**
- 13 Do standard training on the dataset \mathcal{D}' ;
- 14 **end**

Summary

Adversarial Recommender System

Adversarial Attack

Adversarial Defense

Heuristic Attack

Random attack

Average attack

Bandwagon attack

segment attack

...

Gradient-based Attack

UnAttack

S-Attack

Graph-based Attack

...

RL-based Attack

PoisonRec

KGAttack

CopyAttack

...

Detection

DegSim and RDMA

PPu and Du

TSGR, RSF, and TBR

...

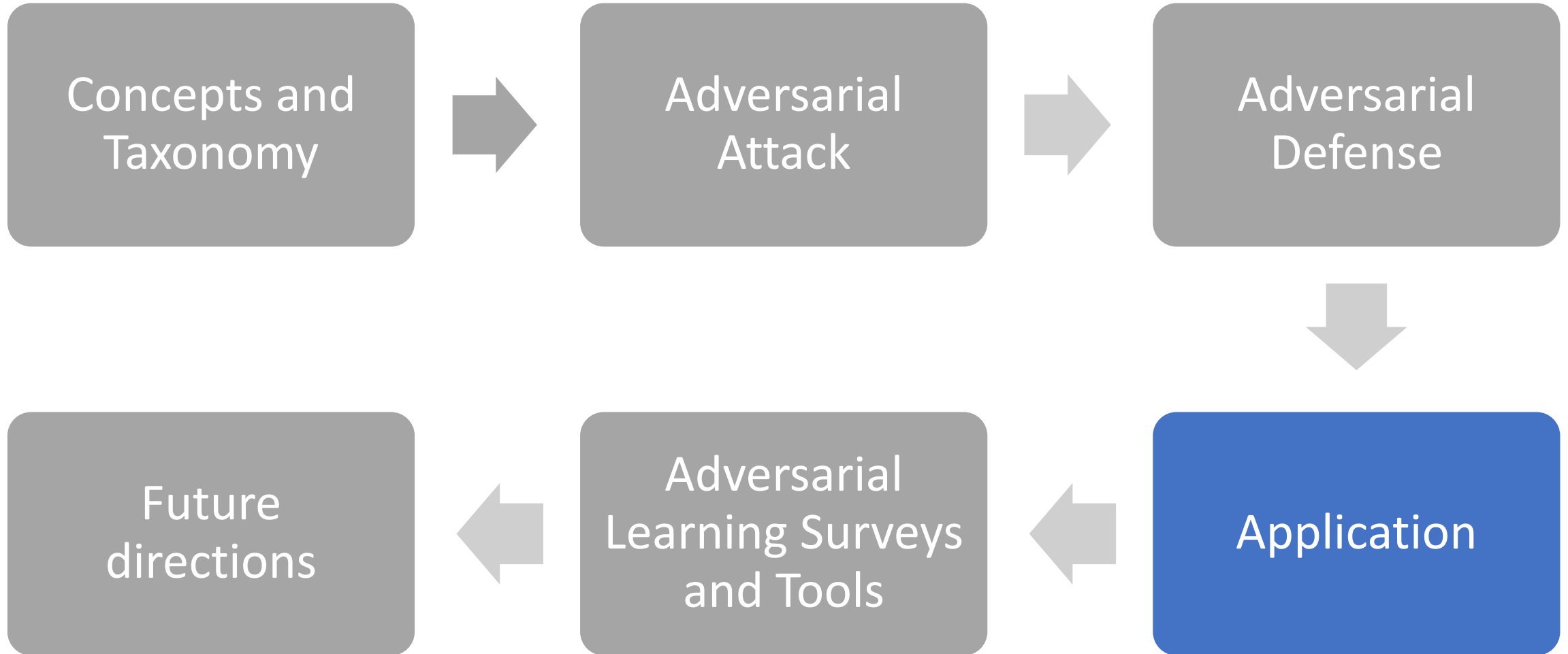
Adversarial Training

APR

APT

...

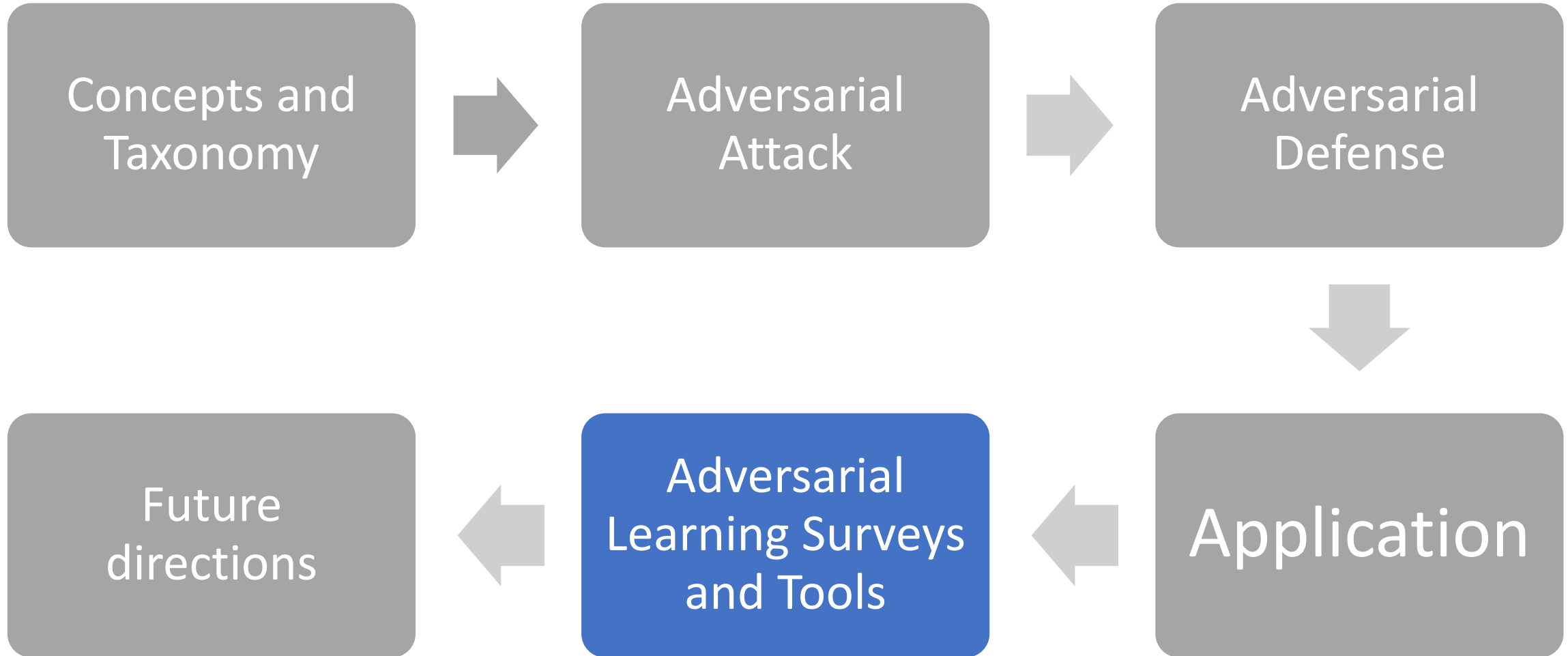
Outline



Application

- The application of adversarial training can help improve the trustworthiness and reliability of recommendation systems in various domains, including:
 - E-health recommendation
 - E-commercial recommendation
 - ...

Outline



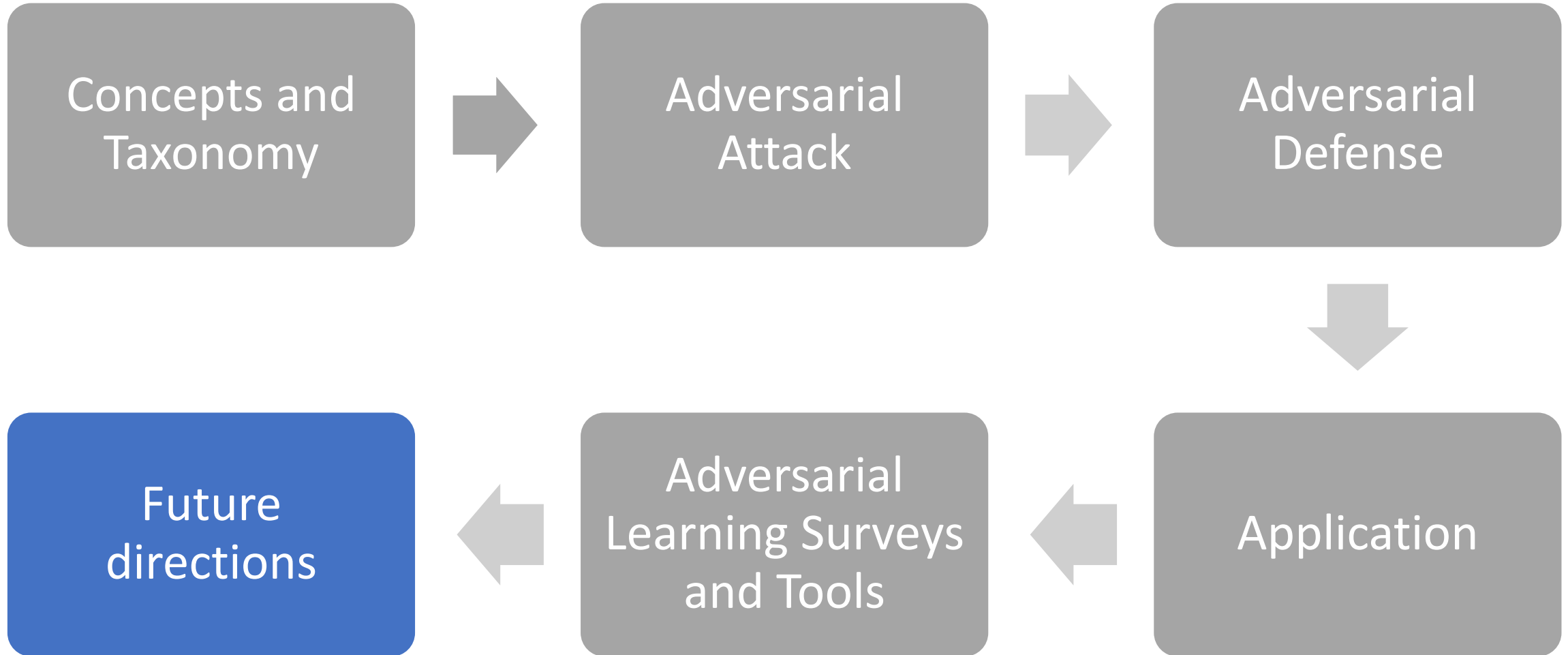
Adversarial Learning Surveys

- Attack:
 - Zhang, Fuguo. "A survey of shilling attacks in collaborative filtering recommender systems." 2009
 - Gunes, Ihsan, et al. "Shilling attacks against recommender systems: A comprehensive survey." 2014
 - Si, Mingdan, and Qingshan Li. "Shilling attacks against collaborative recommender systems: a review." 2020
- Adversarial recommender systems:
 - Truong, Anh, Negar Kiyavash, and Seyed Rasoul Etesami. "Adversarial machine learning: The case of recommendation systems." 2018
 - Deldjoo, Yashar, Tommaso Di Noia, and Felice Antonio Merra. "A survey on adversarial recommender systems: from attack/defense strategies to generative adversarial networks." 2021

Adversarial Learning Tools

- RGRRecSys (Ovaisi et al., 2022)

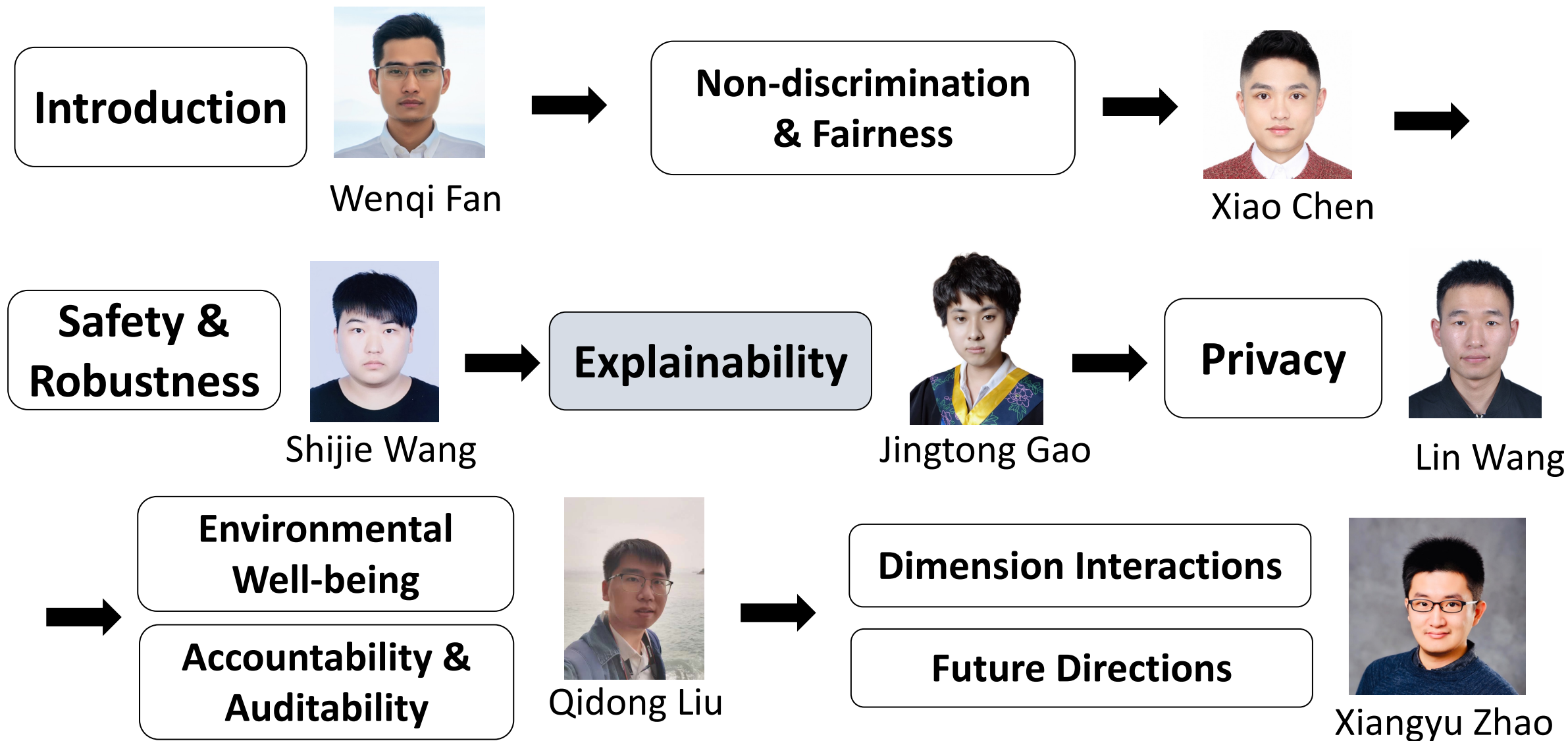
Outline



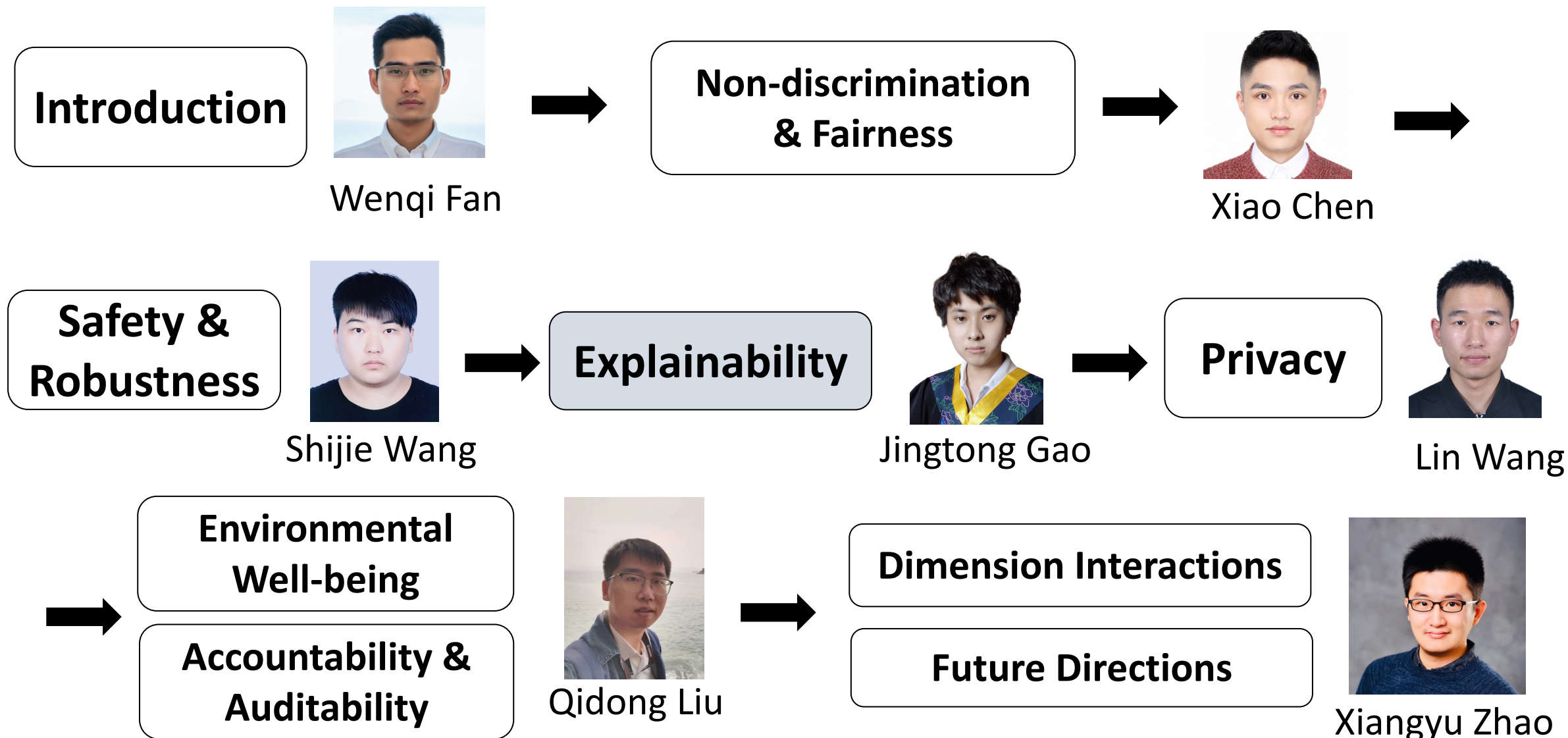
Future Directions

- Investigate vulnerability of different recommender systems
- Generate adversarial perturbations on user-item interactions for adversarial robust training
- Address open problems and challenges in robustness in recommendation

Trustworthy Recommender Systems



Trustworthy Recommender Systems



Explainability

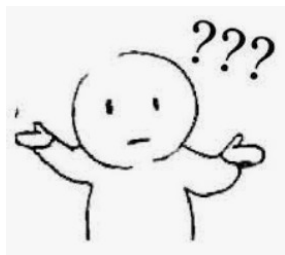
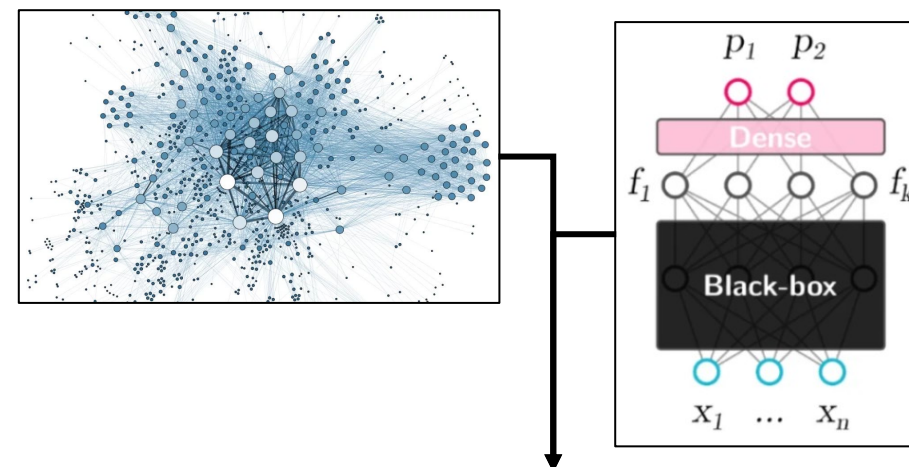
- **What's explainability in Rec, or to say explainable recommendations?**
 - It refers to the recommendation algorithms focusing on **providing explanation for recommendation results**



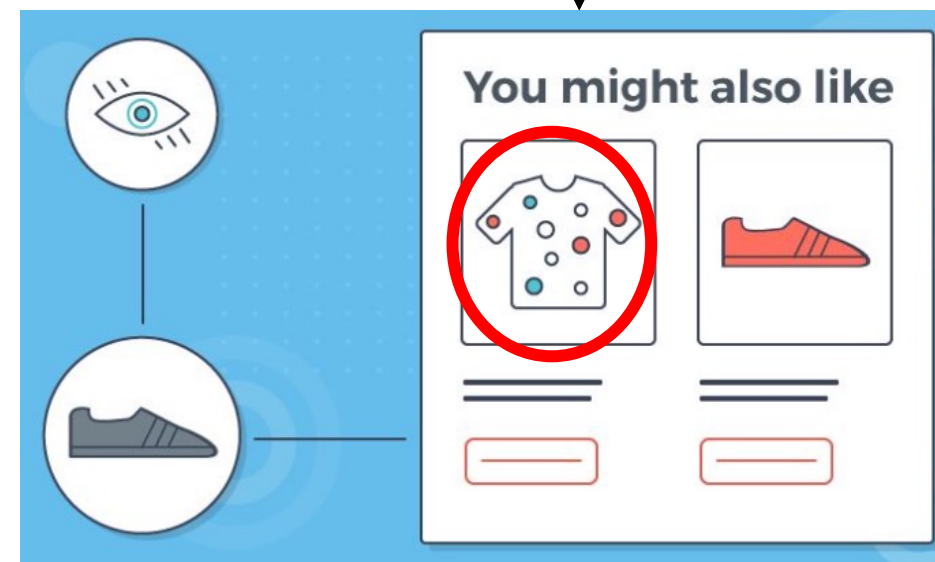
Explainability

- **Why do we need explainability in a trustworthy Rec system?**

- Complicated modeling & Black-box module:

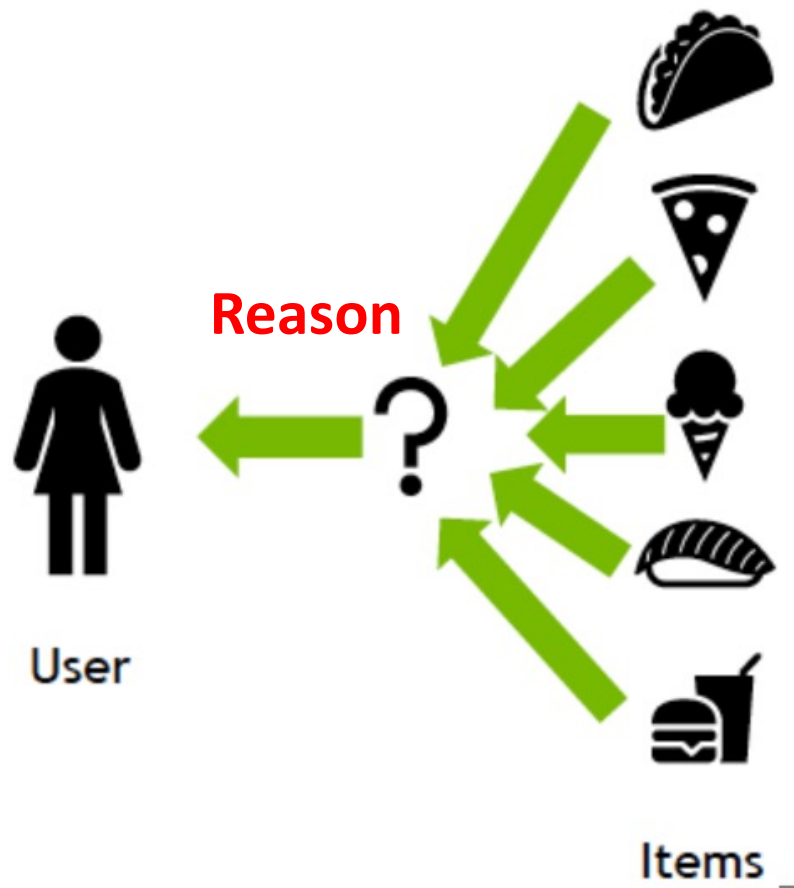


- Why would you recommend this to me?
- Similar style, same brand, or just a mis-recommendation?



Concepts

- The ability to explain or to present in understandable terms to a human



Explainability



METHODS



EVALUATIONS



APPLICATIONS



**FUTURE
DIRECTIONS**

Taxonomy

- How to **produce explanations**: model-intrinsic based (mostly used) or post-hoc
- How the **explanations are presented**: structured or unstructured

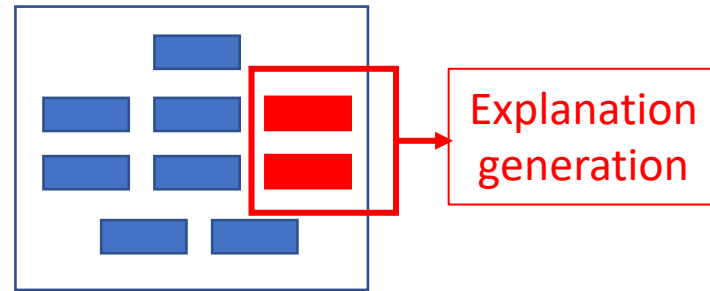
	Model-intrinsic based	Post-Hoc	<i>Characteristics</i>
Structured	[48, 114, 364, 389, 390, 396]	[280, 319]	Logical, Visible
Unstructured	[63, 64, 291]	[211, 315, 338]	Diversified, Fragmented
<i>Focus</i>	Model's reasoning process	Instances' relationship	-

Note: Since some studies construct models from multiple perspectives at the same time, these different classifications are not completely antithetical

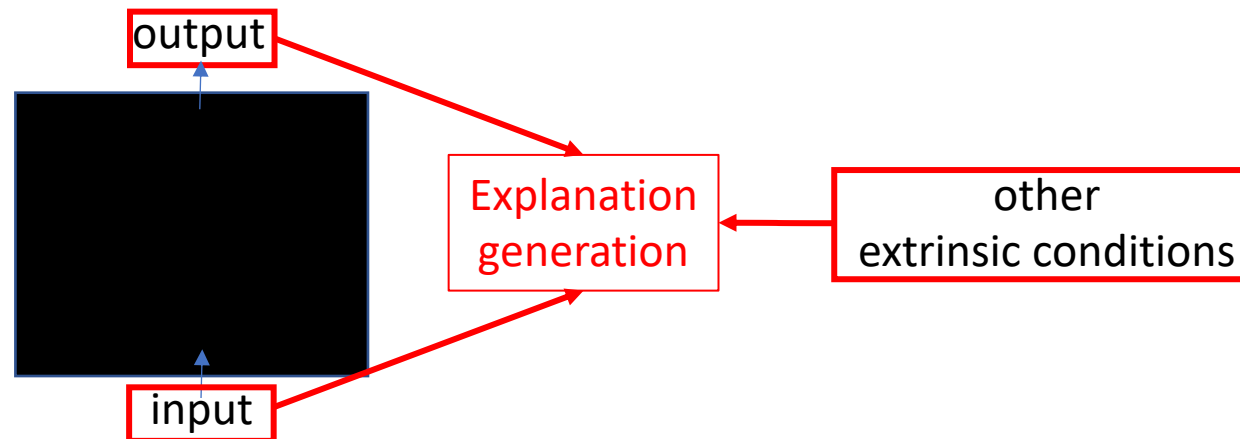
Taxonomy

- **The first criteria: How to produce explanations**

- Model-intrinsic based methods: seek to derive explanations from the **intrinsic structure** of the model



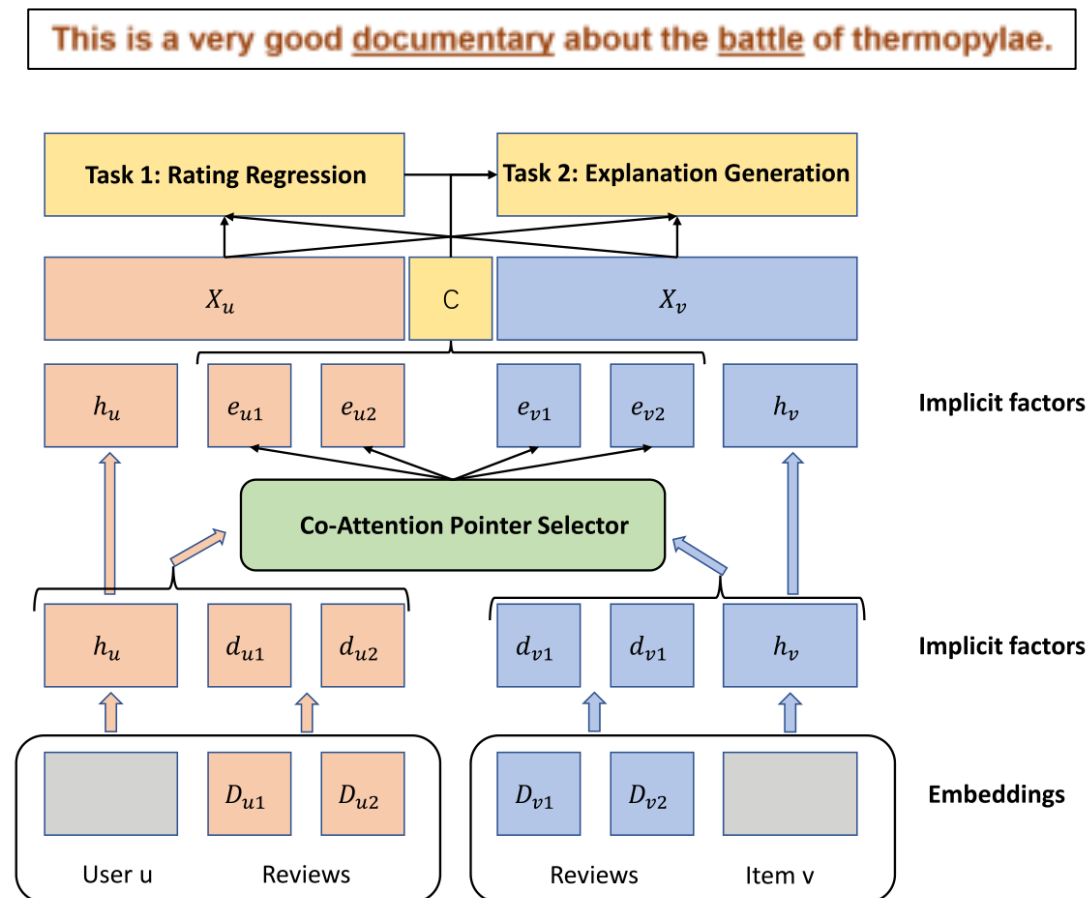
- Post-hoc methods: provide explanations based only on the inputs, outputs and extrinsic conditions of the model



Model-intrinsic based methods

• CAML

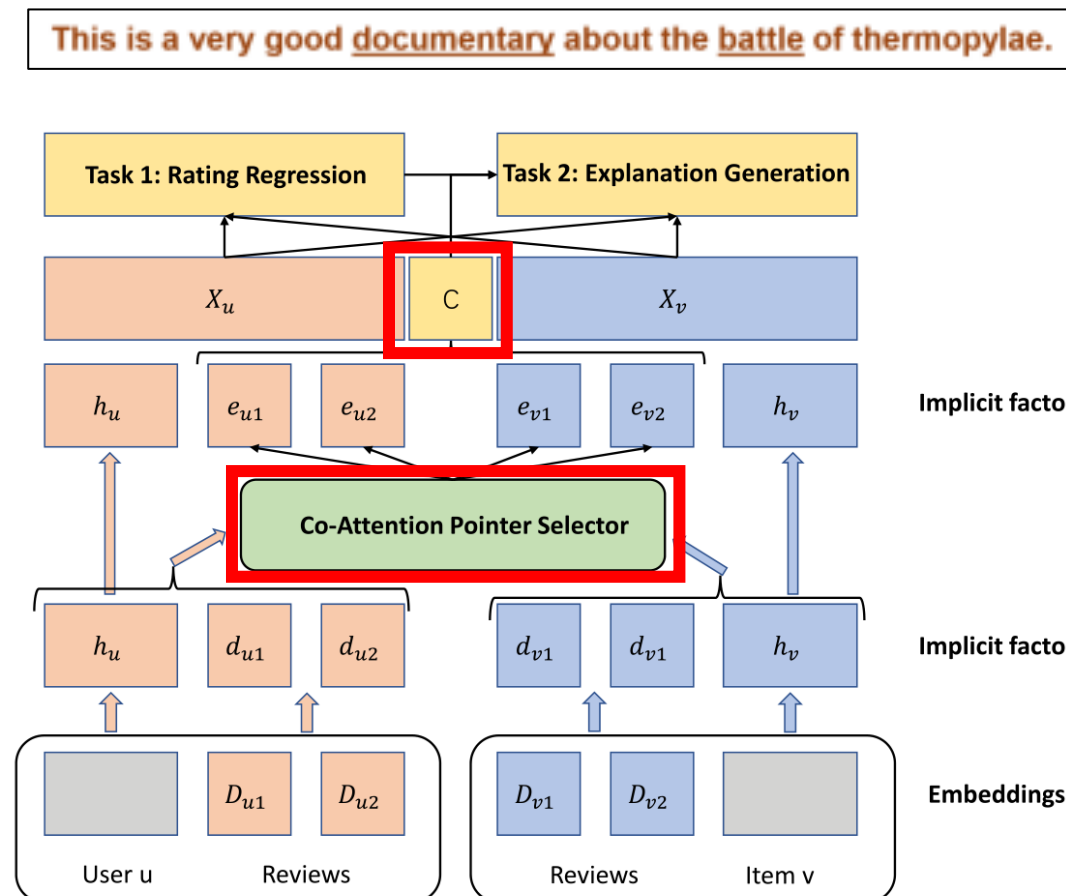
- The explanation is one of the major tasks and modeling goals
- Only effective for the embedded models and cannot simply be reused in other models



Model-intrinsic based methods

• CAML

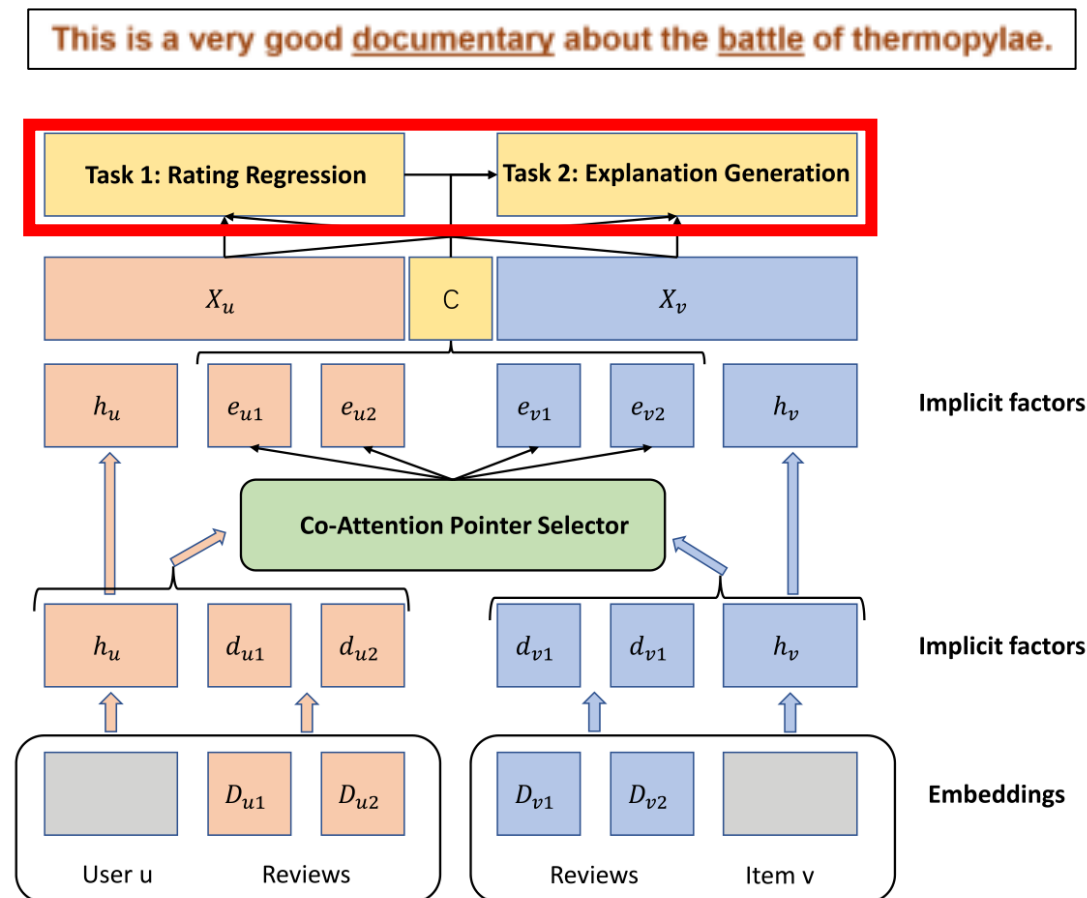
- The explanation is one of the major tasks and modeling goals
- Only effective for the embedded models and cannot simply be reused in other models



Model-intrinsic based methods

• CAML

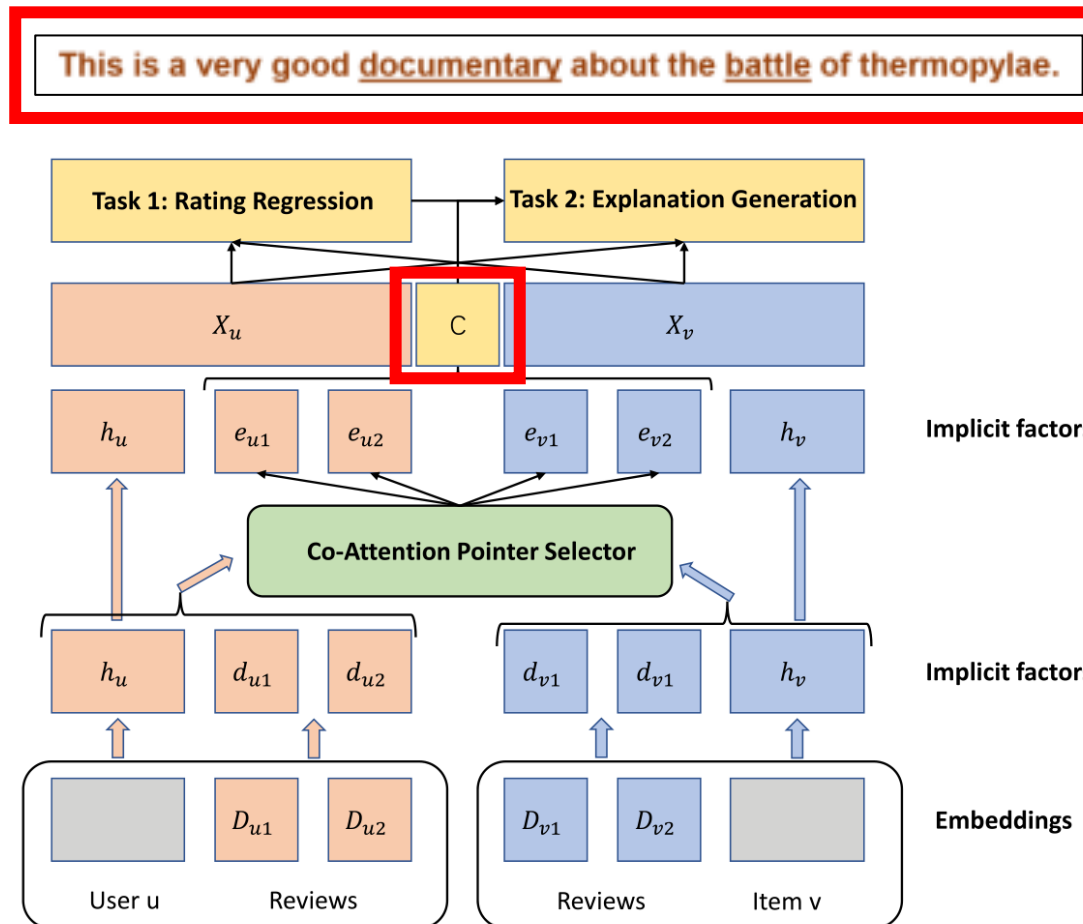
- The explanation is one of the major tasks and modeling goals
- Only effective for the embedded models and cannot simply be reused in other models



Model-intrinsic based methods

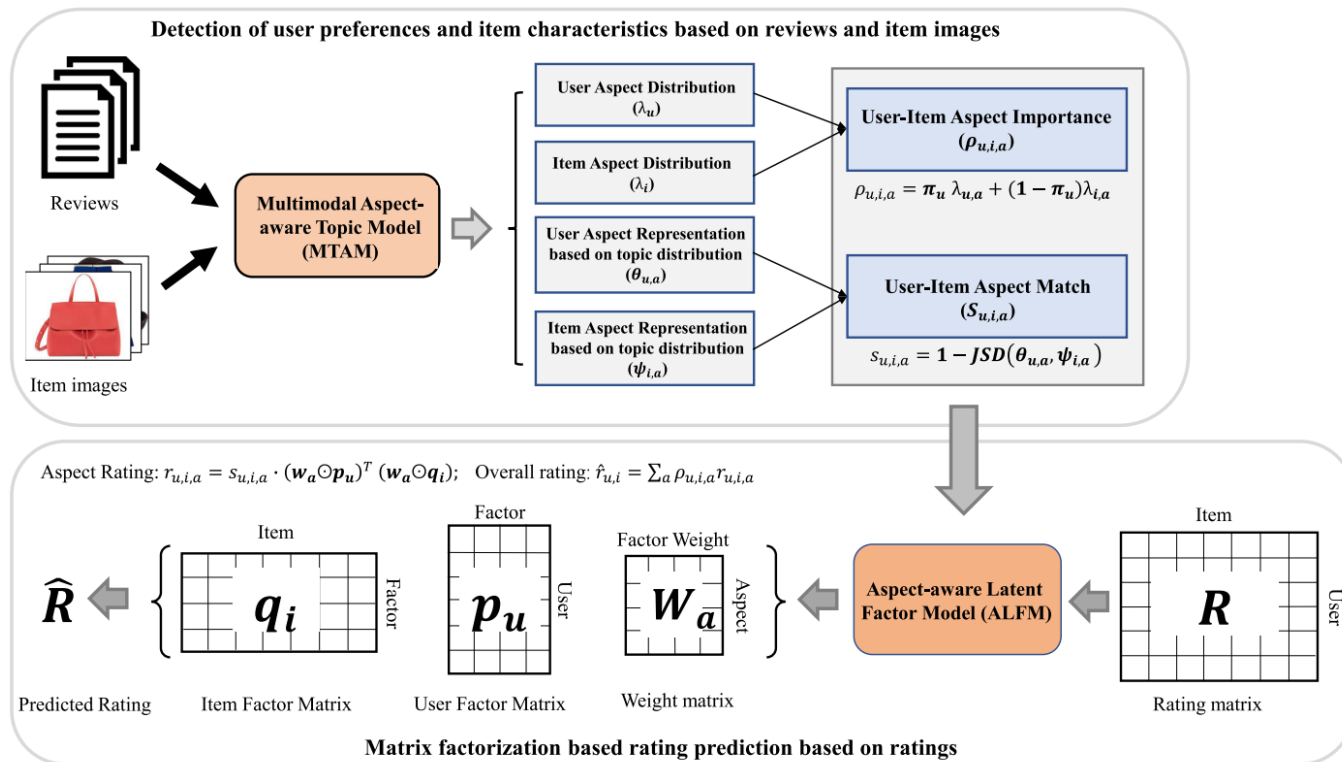
• CAML

- The explanation is one of the major tasks and modeling goals
- Only effective for the embedded models and cannot simply be reused in other models



Model-intrinsic based methods

• MMALFM



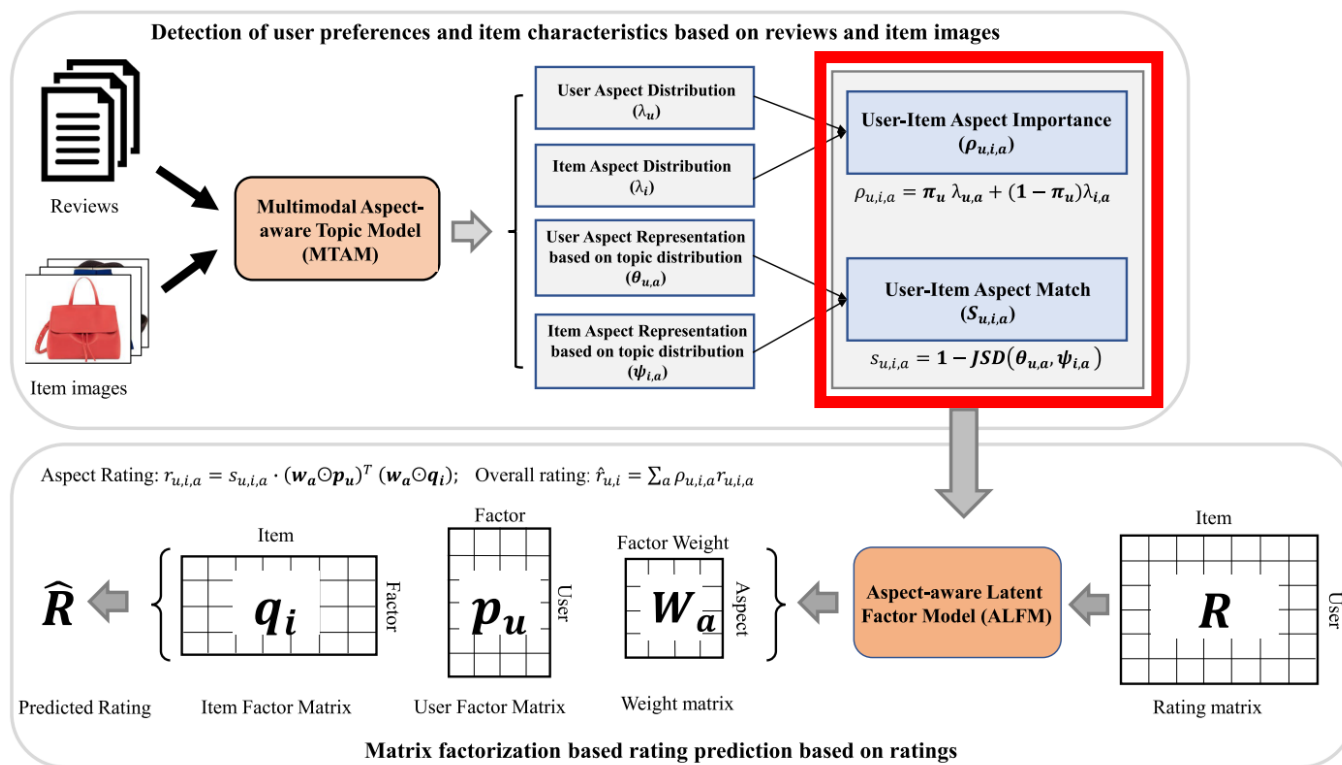
User_2397	Food Ambience Price Service Misc.	sauce, fried, bread, fresh, huge, flavor, shrimp, dessert, dish nice, bar, atmosphere, location, friendly, inside, decor, staff, music expensive, high, cheap, pricey, decent, pay, reasonable, priced, deal table, server, friendly, minutes, nice, staff, asked, make, seated never, give, restaurant, times, stars, friends, night, places, dinner
Item_137	Food Ambience Price Service Misc.	sauce, salad, fries, dish, cheese, dishes, burger, fresh, crab bar, atmosphere, patio, area, inside, wine, small, cool, decor price, worth, prices, better, bit, meal, sauce, dishes, quality table, bar, friendly, wait, server, staff, minutes, beer, atmosphere eat, dinner, Vegas, experience, wait, friends, times, never, give
Item_673	Food Ambience Price Service Misc.	nigiri, sake, tempura, shrimp, sauce, items, poke, crab, chef atmosphere, friendly, bar, staff, inside, area, spot, monta, feel price, worth, prices, nigiri, sake, tempura, items, lunch, special service, table, server, friendly, minutes, staff, nice, asked, seated restaurant, times, give, favorite, night, places, stars, friends, Vegas

Table 6. Interpretation for Why the “User 2397” Rated “Item 137” and “Item 673” with 5 and 2, Respectively

Item	Aspect	Food	Ambience	Price	Service	Misc.
Item_137	Importance	0.3815	0.1034	0.0723	0.2038	0.2390
	Matching	0.5672	0.4523	0.5329	0.6021	0.7138
	Polarity	+	+	-	+	+
Item_673	Importance	0.3726	0.0794	0.0853	0.2076	0.2551
	Matching	0.1813	0.6535	0.4512	0.6018	0.7093
	Polarity	-	-	+	+	-

Model-intrinsic based methods

- MMALFM



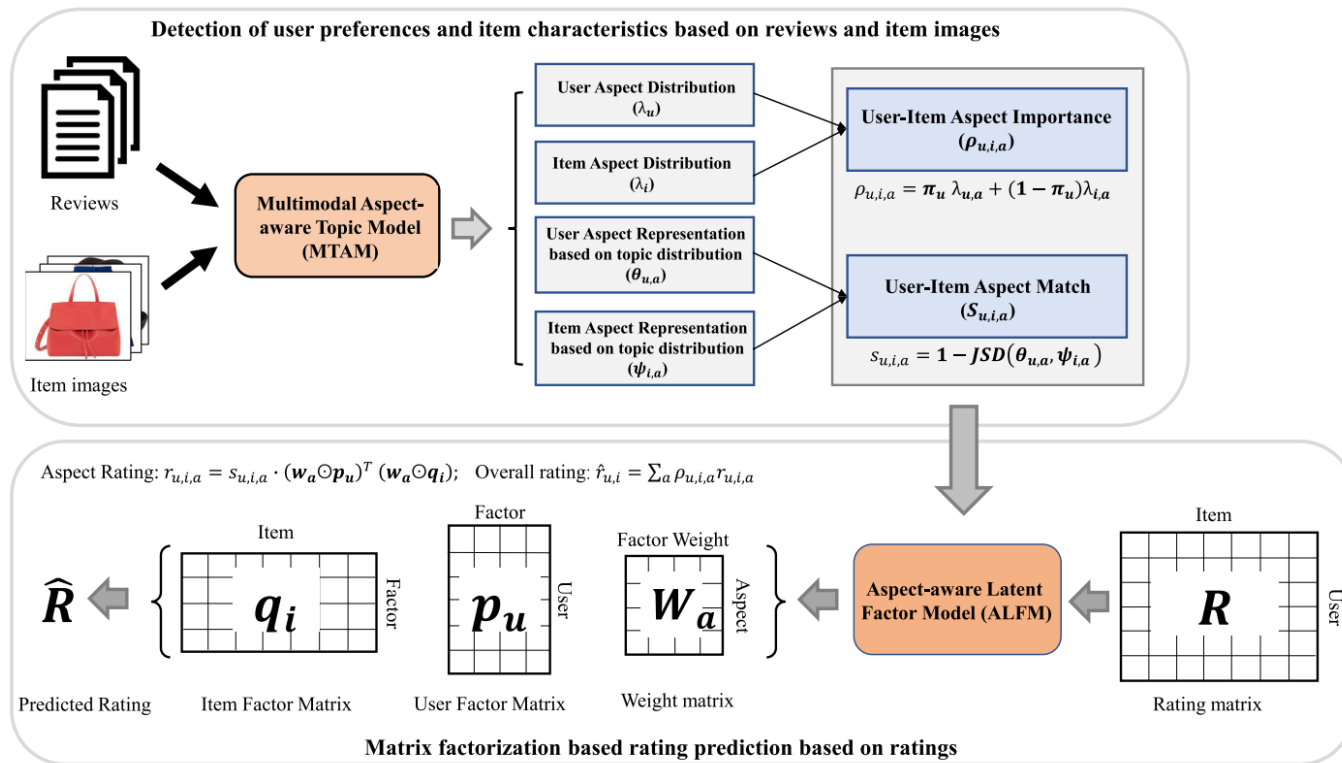
User_2397	Food	sauce, fried, bread, fresh, huge, flavor, shrimp, dessert, dish
	Ambience	nice, bar, atmosphere, location, friendly, inside, decor, staff, music
	Price	expensive, high, cheap, pricey, decent, pay, reasonable, priced, deal
	Service	table, server, friendly, minutes, nice, staff, asked, make, seated
	Misc.	never, give, restaurant, times, stars, friends, night, places, dinner
Item_137	Food	sauce, salad, fries, dish, cheese, dishes, burger, fresh, crab
	Ambience	bar, atmosphere, patio, area, inside, wine, small, cool, decor
	Price	price, worth, prices, better, bit, meal, sauce, dishes, quality
	Service	table, bar, friendly, wait, server, staff, minutes, beer, atmosphere
	Misc.	eat, dinner, Vegas, experience, wait, friends, times, never, give
Item_673	Food	nigiri, sake, tempura, shrimp, sauce, items, poke, crab, chef
	Ambience	atmosphere, friendly, bar, staff, inside, area, spot, monta, feel
	Price	price, worth, prices, nigiri, sake, tempura, items, lunch, special
	Service	service, table, server, friendly, minutes, staff, nice, asked, seated
	Misc.	restaurant, times, give, favorite, night, places, stars, friends, Vegas

Table 6. Interpretation for Why the “User 2397” Rated “Item 137” and “Item 673” with 5 and 2, Respectively

Item	Aspect	Food	Ambience	Price	Service	Misc.
Item_137	Importance	0.3815	0.1034	0.0723	0.2038	0.2390
	Matching	0.5672	0.4523	0.5329	0.6021	0.7138
	Polarity	+	+	-	+	+
Item_673	Importance	0.3726	0.0794	0.0853	0.2076	0.2551
	Matching	0.1813	0.6535	0.4512	0.6018	0.7093
	Polarity	-	-	+	+	-

Model-intrinsic based methods

- MMALFM



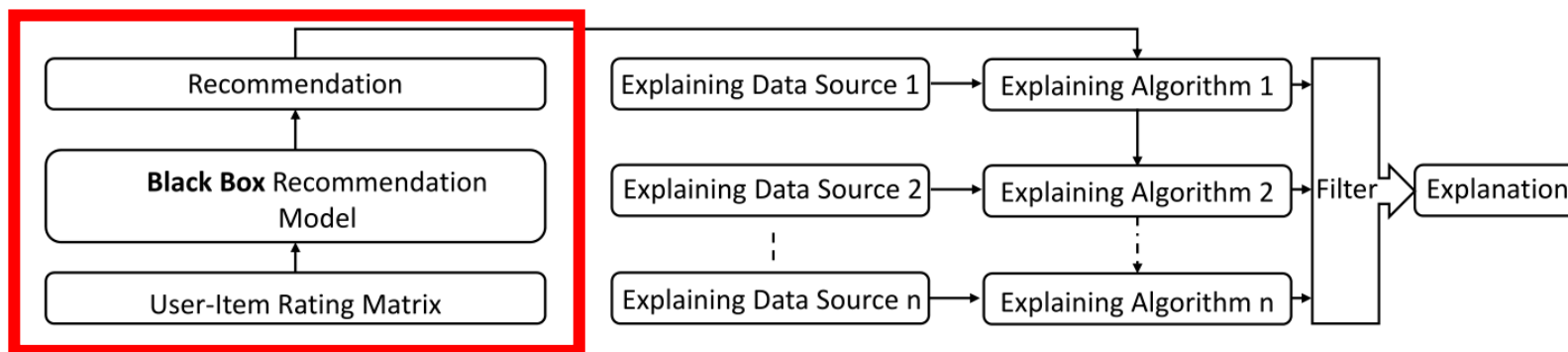
User_2397	Food	sauce, fried, bread, fresh, huge, flavor, shrimp, dessert, dish
	Ambience	nice, bar, atmosphere, location, friendly, inside, decor, staff, music
	Price	expensive, high, cheap, pricey, decent, pay, reasonable, priced, deal
	Service	table, server, friendly, minutes, nice, staff, asked, make, seated
	Misc.	never, give, restaurant, times, stars, friends, night, places, dinner
Item_137	Food	sauce, salad, fries, dish, cheese, dishes, burger, fresh, crab
	Ambience	bar, atmosphere, patio, area, inside, wine, small, cool, decor
	Price	price, worth, prices, better, bit, meal, sauce, dishes, quality
	Service	table, bar, friendly, wait, server, staff, minutes, beer, atmosphere
	Misc.	eat, dinner, Vegas, experience, wait, friends, times, never, give
Item_673	Food	nigiri, sake, tempura, shrimp, sauce, items, poke, crab, chef
	Ambience	atmosphere, friendly, bar, staff, inside, area, spot, monta, feel
	Price	price, worth, prices, nigiri, sake, tempura, items, lunch, special
	Service	service, table, server, friendly, minutes, staff, nice, asked, seated
	Misc.	restaurant, times, give, favorite, night, places, stars, friends, Vegas

Table 6. Interpretation for Why the “User 2397” Rated “Item 137” and “Item 673” with 5 and 2, Respectively

Item	Aspect	Food	Ambience	Price	Service	Misc.
Item_137	Importance	0.3815	0.1034	0.0723	0.2038	0.2390
	Matching	0.5672	0.4523	0.5329	0.6021	0.7138
	Polarity	+	+	-	+	+
Item_673	Importance	0.3726	0.0794	0.0853	0.2076	0.2551
	Matching	0.1813	0.6535	0.4512	0.6018	0.7093
	Polarity	-	-	+	+	-

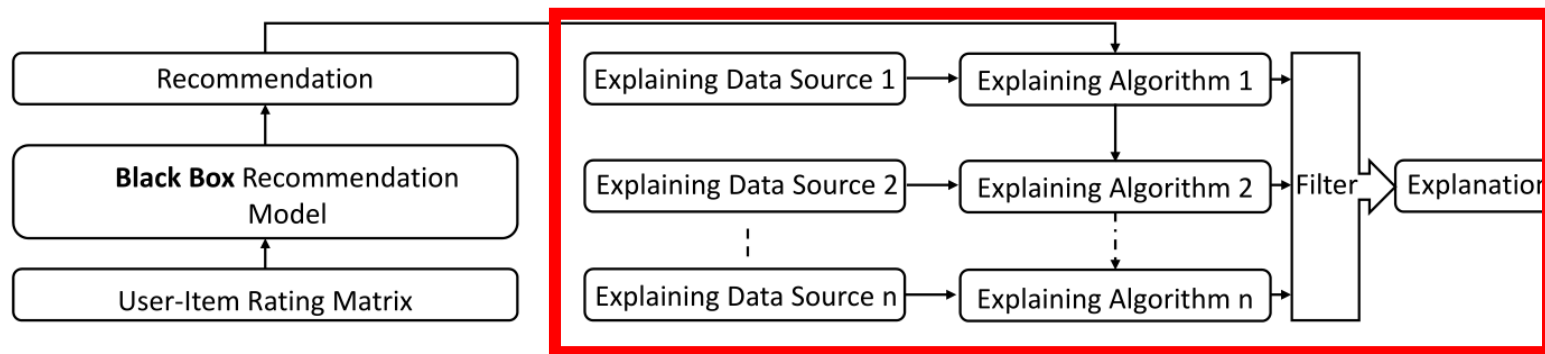
Post-hoc methods

- **An example from Shmaryahu et al.**
 - It generates explanations directly from the recommendation and explaining data source



Post-hoc methods

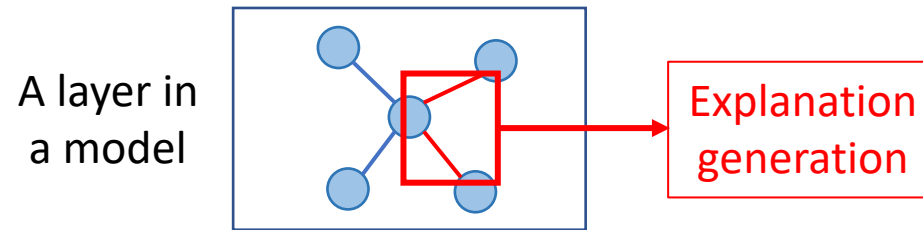
- An example from Shmaryahu et al.
 - It generates explanations directly from the recommendation and explaining data source



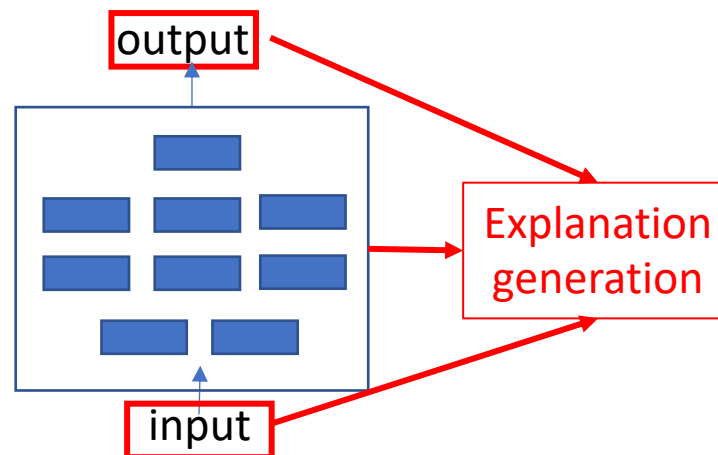
Taxonomy

- **The second criteria: How the explanations are presented**

- Structured methods: present explanations in the form of **logical reasoning** based on some particular structures, such as a graph, or a knowledge graph



- Unstructured methods: provide explanations based on the inputs, outputs and models, do not rely on, or explicitly rely on logical reasoning

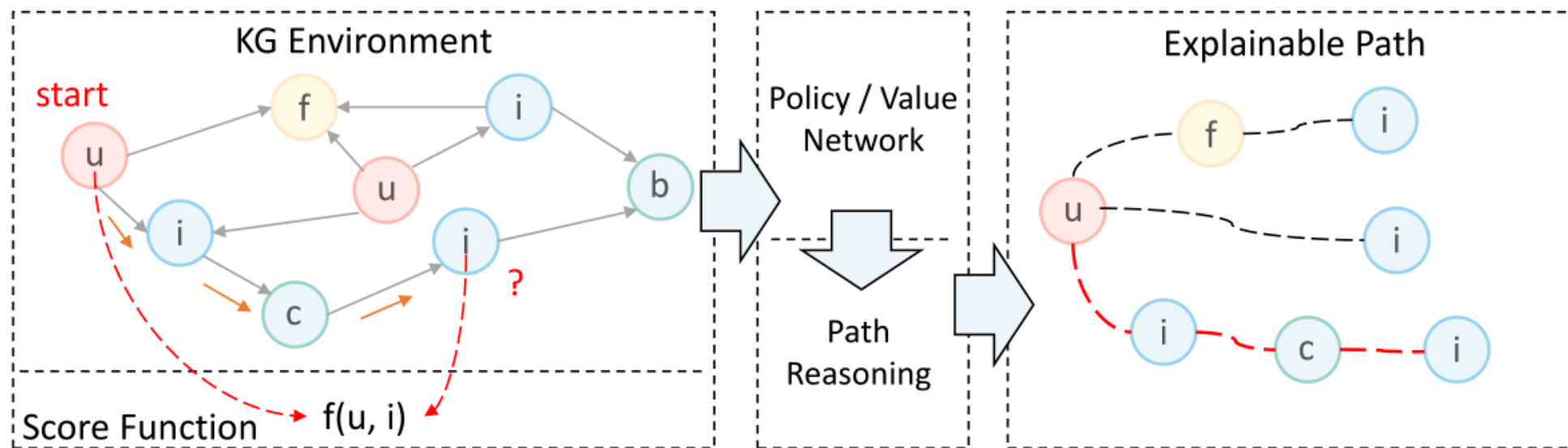


Structured methods

- **PGPR**

- An explanation path graph generated with knowledge graph

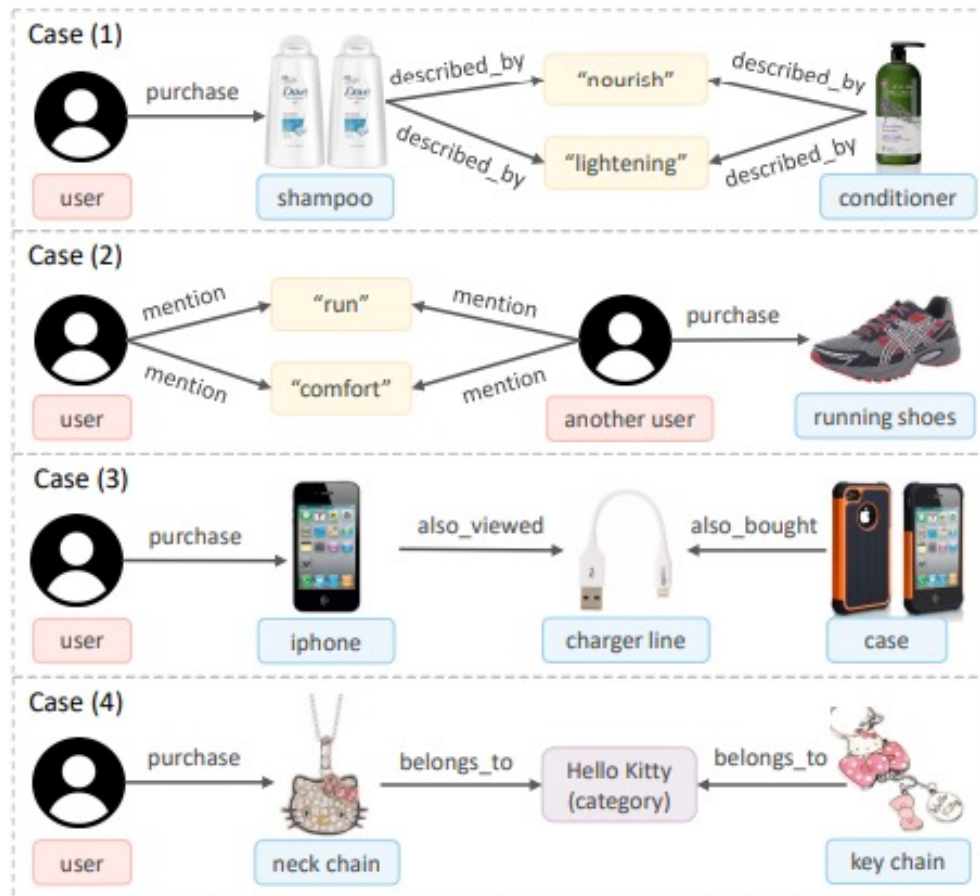
- Path definition: $p_k(e_0, e_k) = \{e_0 \xleftrightarrow{r_1} e_1 \xleftrightarrow{r_2} \dots \xleftrightarrow{r_k} e_k\}$



Structured methods

- PGPR

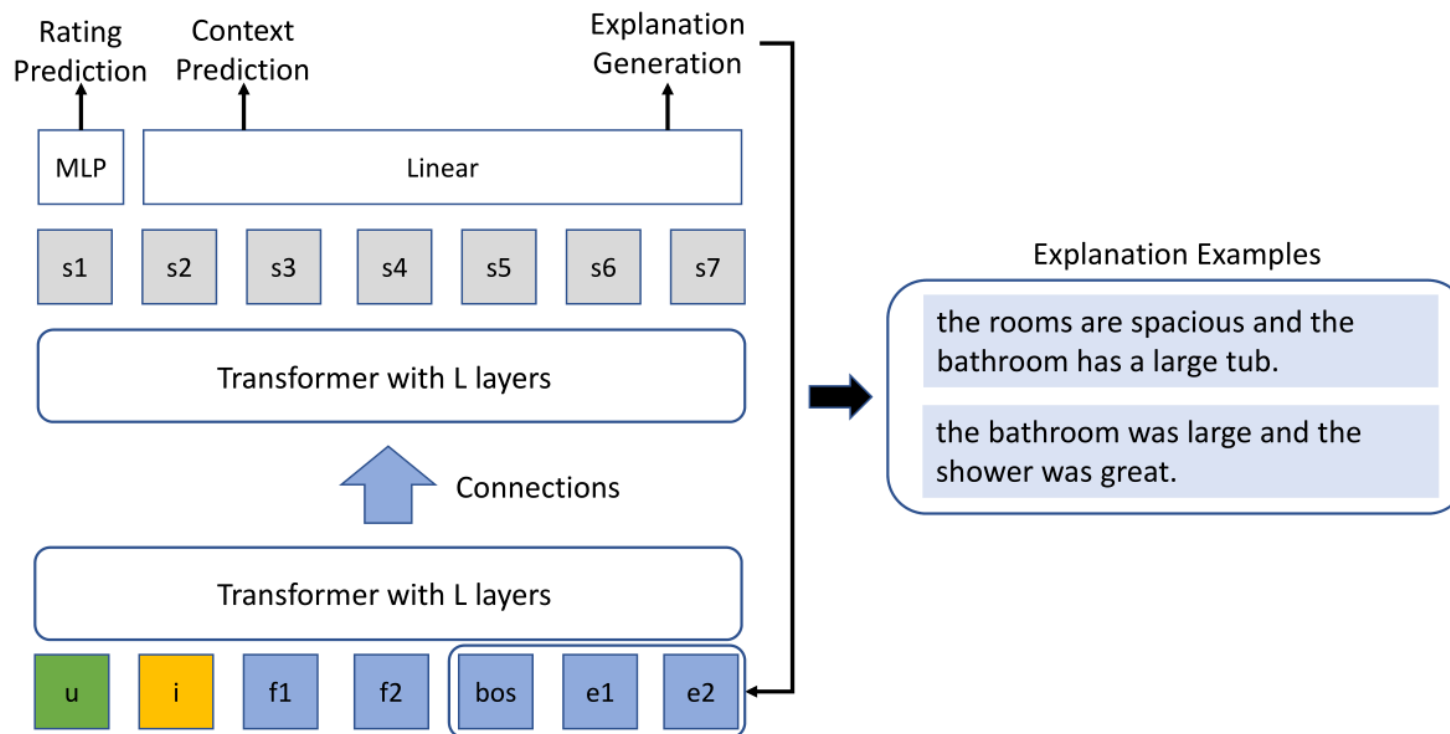
- Explanation path



Unstructured methods

- **PETER**

- Generate explanation sentence word by word
- The final explanation is a sentence based on probability, not the sole reason deduced according to deterministic rules or structures



Unstructured methods

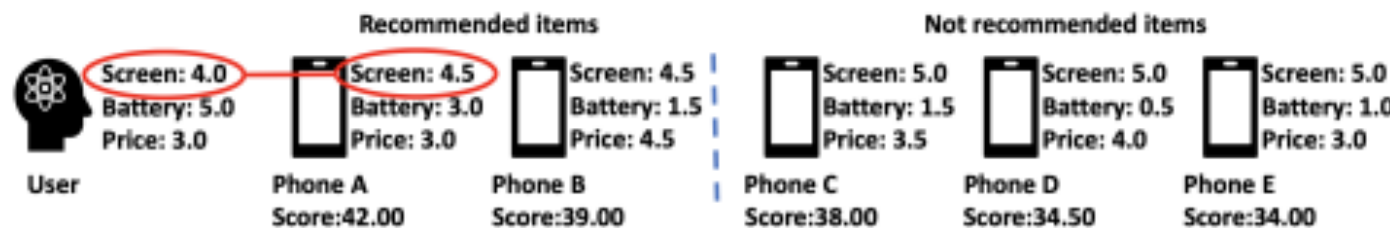
- **CountER**

- It tries to use small changes in item aspects to reverse the decision

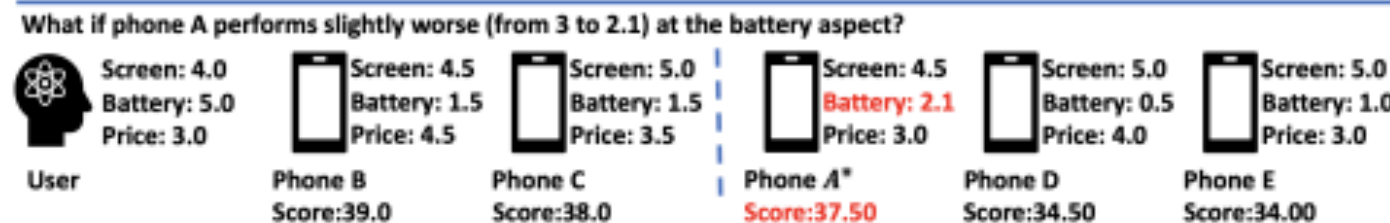
If the item had been slightly worse on [aspect(s)], then it will not be recommended.

minimize Explanation Complexity
s.t., Explanation is Strong Enough

Matching-based:



Counterfactual reasoning:



Unstructured methods

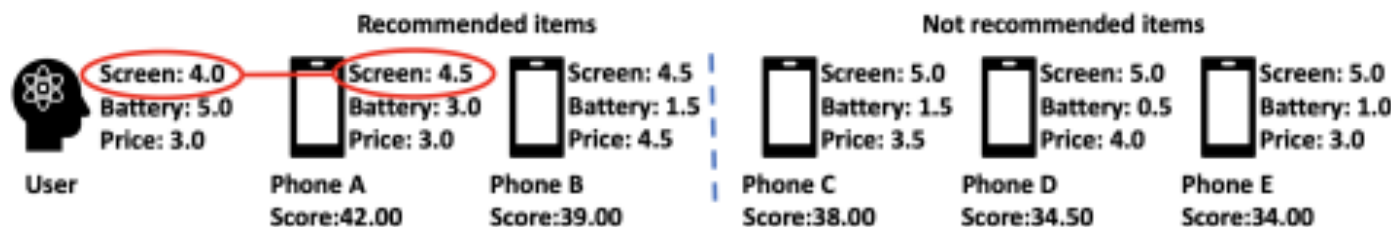
- **CountER**

- It tries to use small changes in item aspects to reverse the decision

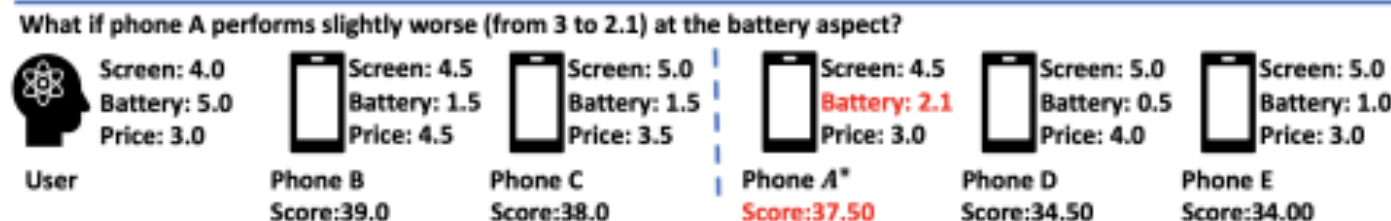
If the item had been slightly worse on [aspect(s)], then it will not be recommended.

minimize Explanation Complexity
s.t., Explanation is Strong Enough

Matching-based:



Counterfactual reasoning:



Explainability



METHODS



EVALUATIONS



APPLICATIONS



FUTURE
DIRECTIONS

Taxonomy of research on evaluations

- **Evaluation perspectives**

- Effectiveness
- Transparency
- Scrutability

- **Evaluation form**

- Quantitative metrics
- Case study
- Real-world performance
- Ablation Study

Taxonomy of Evaluation

- **Evaluation perspectives**
 - Effectiveness
 - Transparency
 - Scrutability

Evaluation perspective	Evaluation criteria	Related research
Effectiveness	Whether the explanations are useful to users? (e.g. Decision making, Recommendation results)	[8, 58, 337]
Transparency	Whether the explanations can reveal the working principles of the model?	[18, 144, 225]
Scrutability	Whether the explanations contribute to the prediction of the model?	[327, 347, 362]

Taxonomy of Evaluation

- **Evaluation form**
 - **Quantitative:** ROUGE score, BLEU, USR, FMR...
 - **Case study:** Whether the explanation conforms to human logic
 - **Real-world performance:** The practical effects of the explanation
 - **Ablation study:** How algorithmic modules provide explanations and how these modules enhance the recommendation model

Evaluation form	Corresponding perspectives	Related research
Quantitative metrics	Effectiveness; Scrutability	[337, 338]
Case study	Effectiveness; Transparency	[225, 362, 396]
Real-world performance	Effectiveness; Scrutability; Transparency	[58, 347, 392]
Ablation Study	Effectiveness; Transparency	[64, 211, 327]

Explainability



METHODS



EVALUATIONS



APPLICATIONS



FUTURE
DIRECTIONS

E-commercial Recommendation



Social Media



Explainability



METHODS



EVALUATIONS



APPLICATIONS



**FUTURE
DIRECTIONS**

Natural Language Generation

- **Templated based (now)**

I recommend Iron Man to you because you've seen The Avengers

- **Full paragraph interpretation generation (currently exist but their effectiveness has yet to improve)**

Since you've seen movies like The Avengers, and your recent interest is in the TV series, we recommend something similar for you: Agents of S.H.I.E.L.D.

Summary

- **Concept of explainability in Rec**
 - The ability to explain or to present in understandable terms to a human
- **Taxonomy of methods**
 - How to produce explanations: model-intrinsic based (mostly used) or post-hoc
 - How the explanations are presented: structured or unstructured
- **Taxonomy of evaluations**
 - Evaluation perspectives: Effectiveness, Transparency, Scrutability
 - Evaluation forms: Quantitative, Case study, Real-world performance, Ablation study
- **Application**
 - E-commercial Recommendation
 - Social Media
- **Future directions**
 - Natural Language Generation for Explanation
 - Explainable recommendations in more fields