



蚂蚁集团
ANT GROUP

Robust Fraud Transaction Detection: A Two-Player Game Approach

**Qi Tan, Yi Zhao, Laizhong Cui, Qi Li, Min Zhu, Xing Fu, Weiqiang Wang,
Xiaotong Lin, and Ke Xu**

February, 2026



Outline



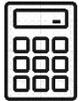
Background



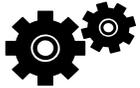
Existing Techniques



Causal Models of Feature Falsification



Two-Player Games in Fraud Detection



Design of GAMER



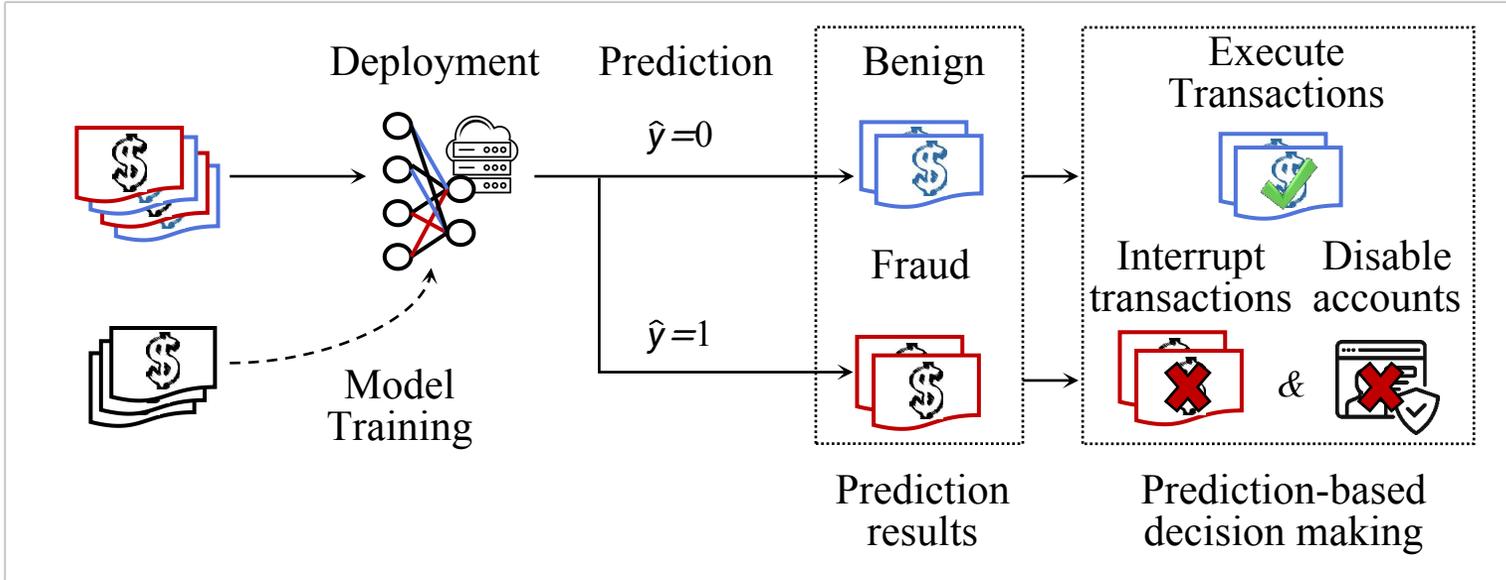
Evaluation



Conclusion



Background



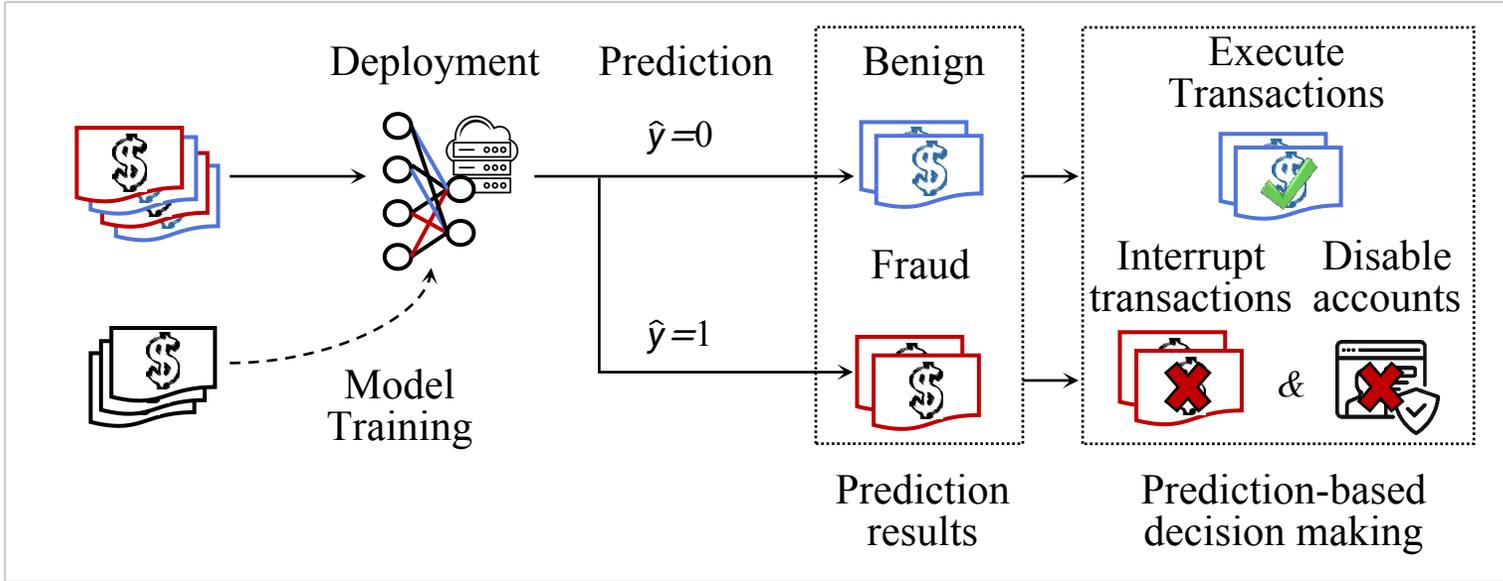
AI-based fraud detection

Detection: Training AI models to detect fraudulent activities

Decision-making: Processing the transactions and accounts based on predictions



Background

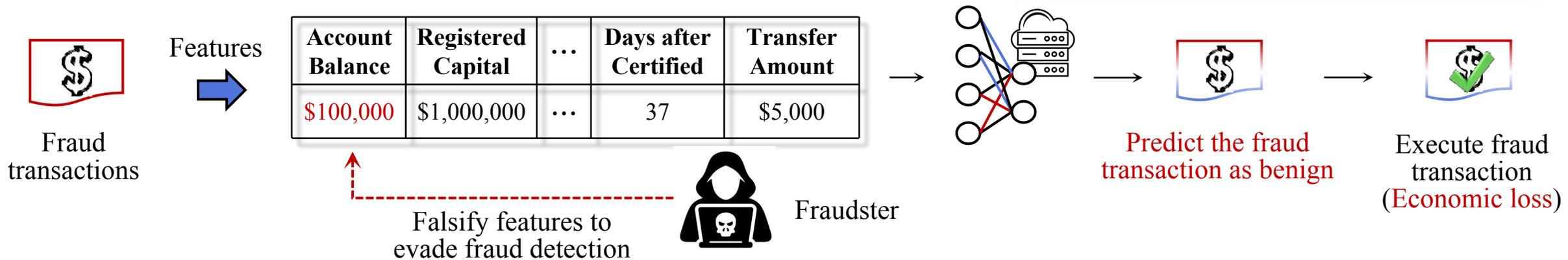


AI-based fraud detection

Detection: Training AI models to detect fraudulent activities

Decision-making: Processing the transactions and accounts based on predictions

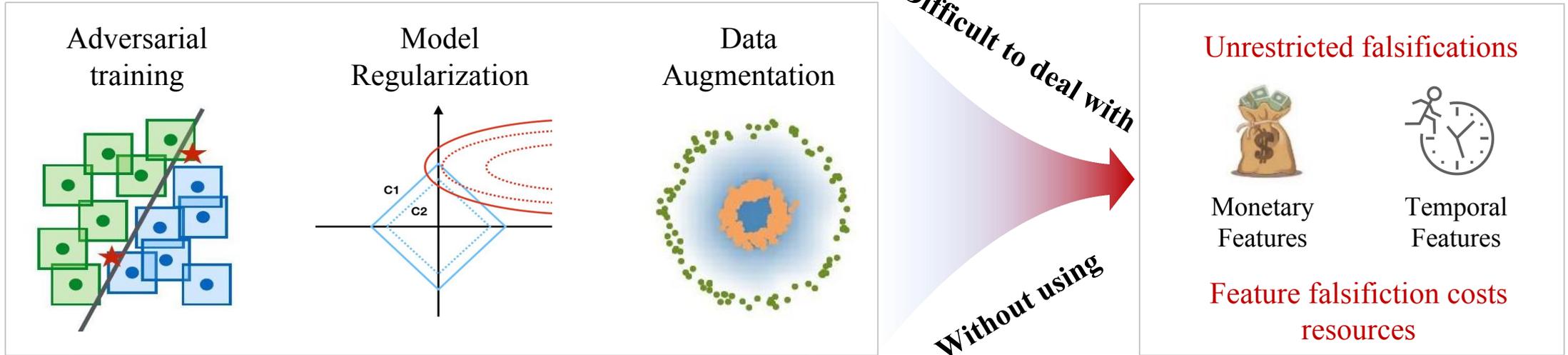
Fraudsters attack the deployed model, evading fraud detection system via feature falsification:





Existing Techniques

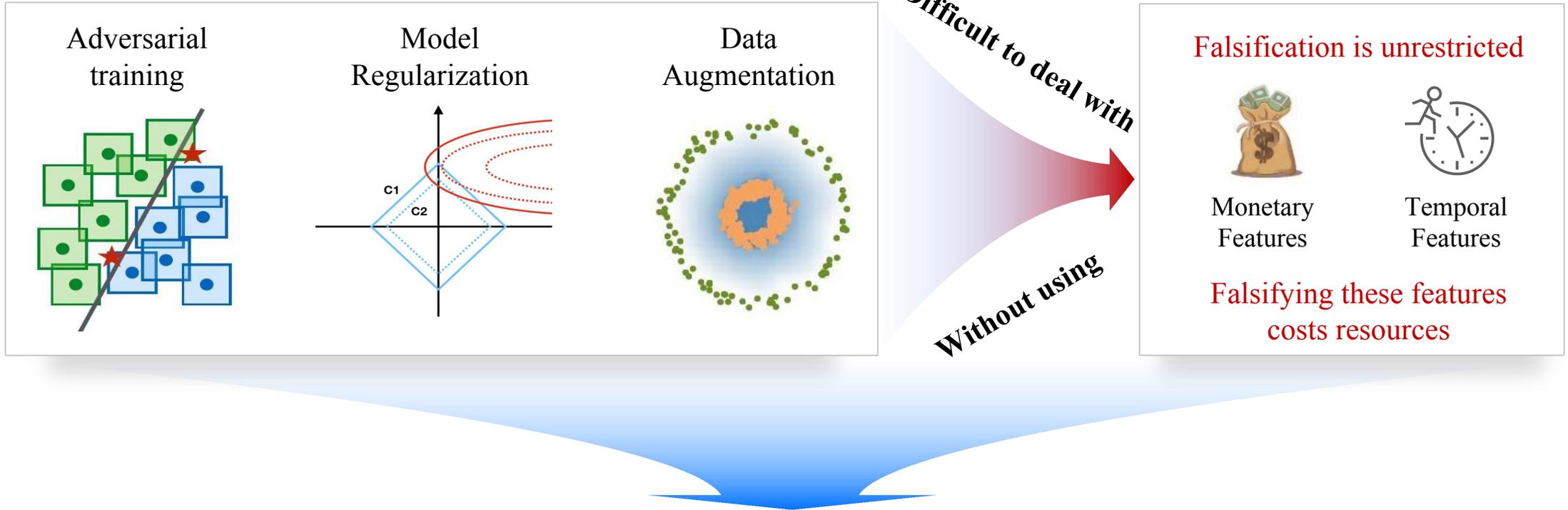
Existing Techniques for enhancing the robustness





Existing Techniques

Existing Techniques for enhancing the robustness



Rethinking the **principle of falsifying features** to evade fraud detection, and design new detection method to **combat intelligent fraudsters**

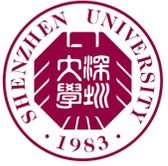


Key issues

Key issues in dealing with feature falsification:

- ◆ Why fraudsters can evade fraud detection via feature falsification?
- ◆ How to defend against potential feature falsification when collected data is limited?
- ◆ How to overcome the deficiencies in defense methods when transactions are benign?



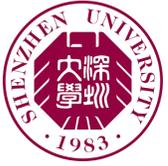


Key issues

Key issues in dealing with feature falsification:

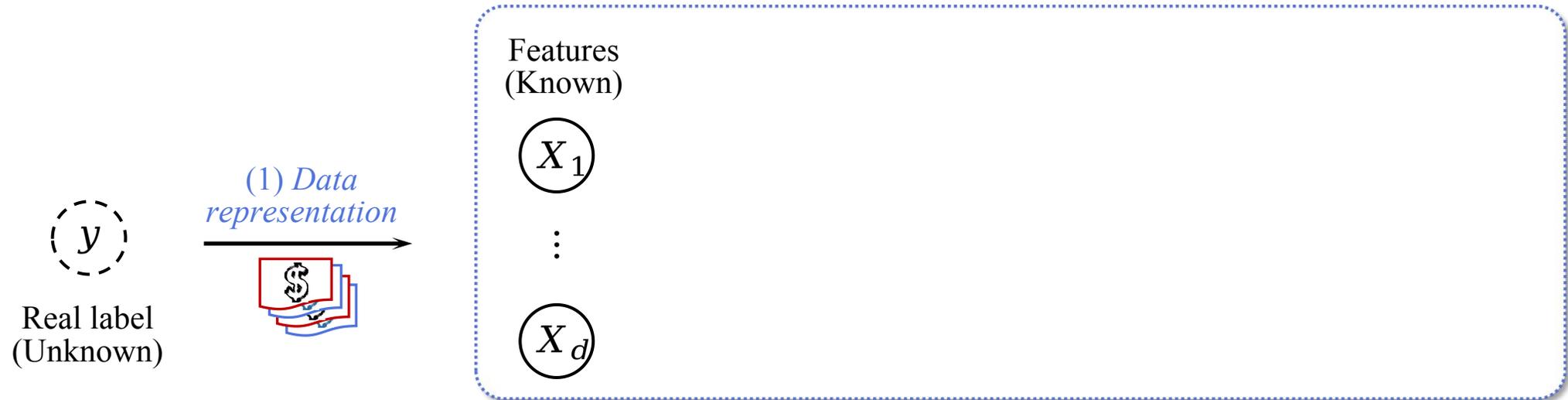
- ◆ **Why fraudsters can evade fraud detection via feature falsification?**
- ◆ How to defend against potential feature falsification when collected data is limited?
- ◆ How to overcome the deficiencies in defense methods when transactions are benign?





Causal Models of Feature Falsification

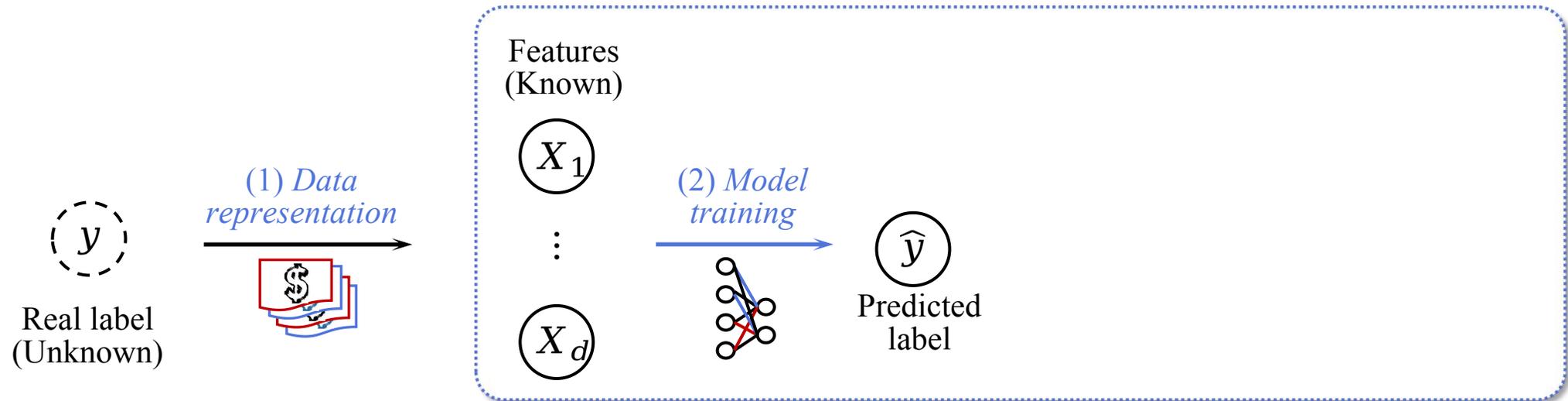
◆ The entire lifecycle of fraud detection





Causal Models of Feature Falsification

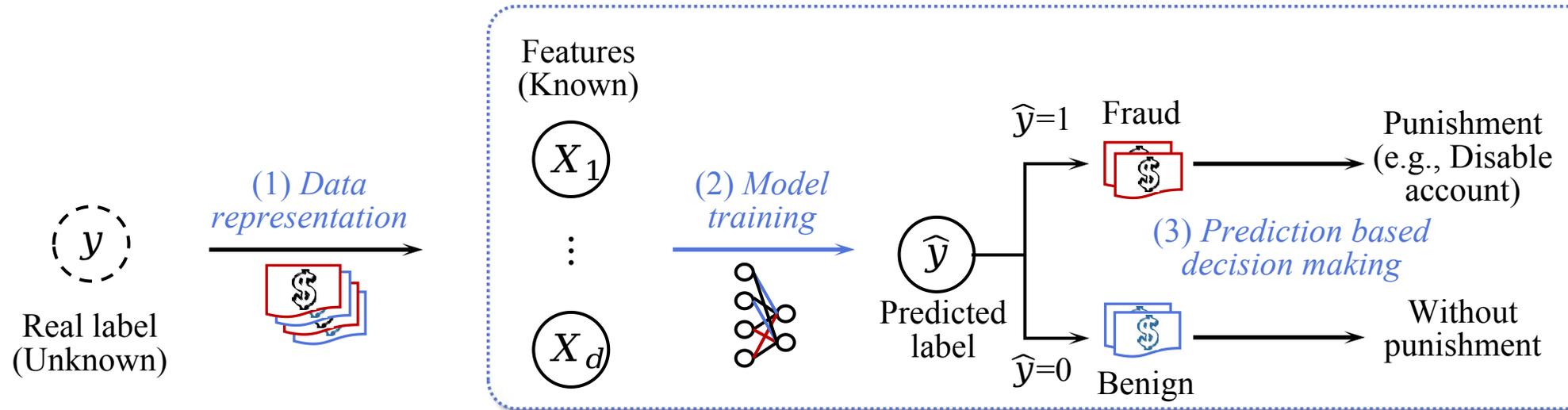
◆ The entire lifecycle of fraud detection

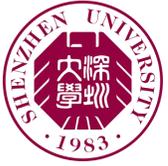




Causal Models of Feature Falsification

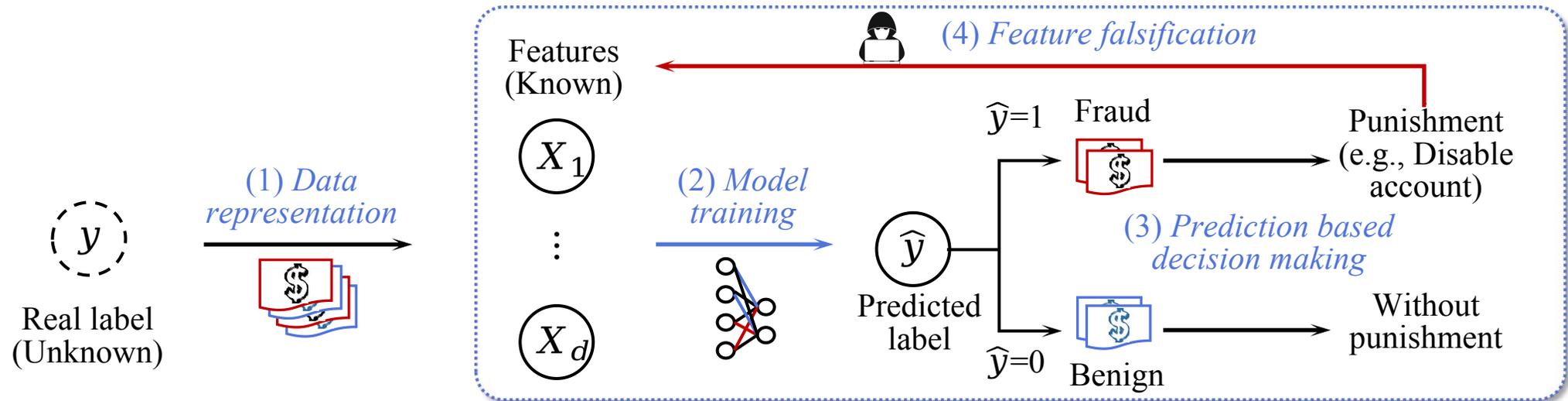
◆ The entire lifecycle of fraud detection





Causal Models of Feature Falsification

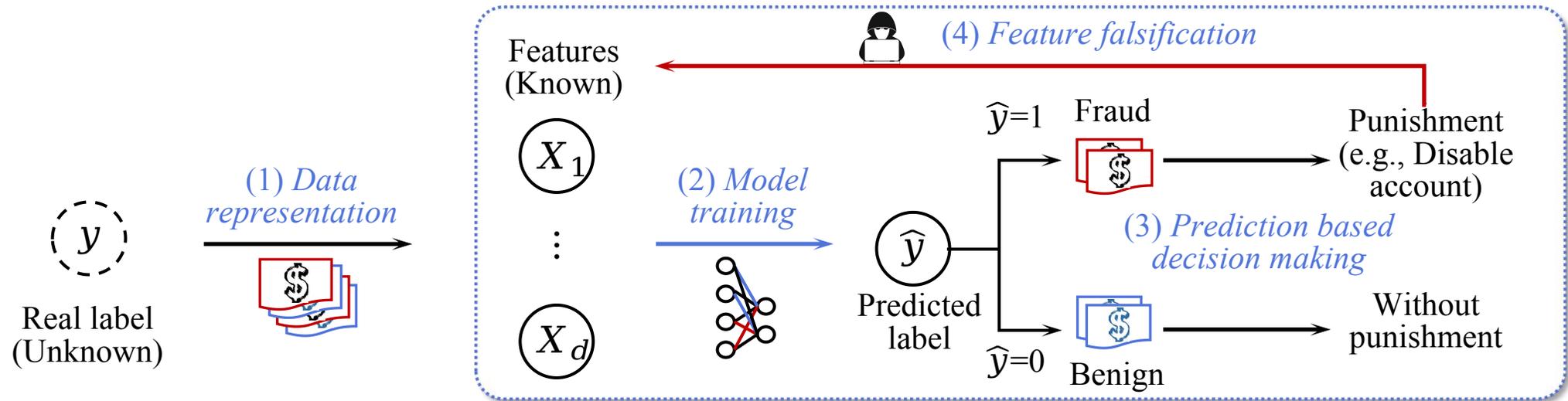
◆ The entire lifecycle of fraud detection





Causal Models of Feature Falsification

◆ The entire lifecycle of fraud detection

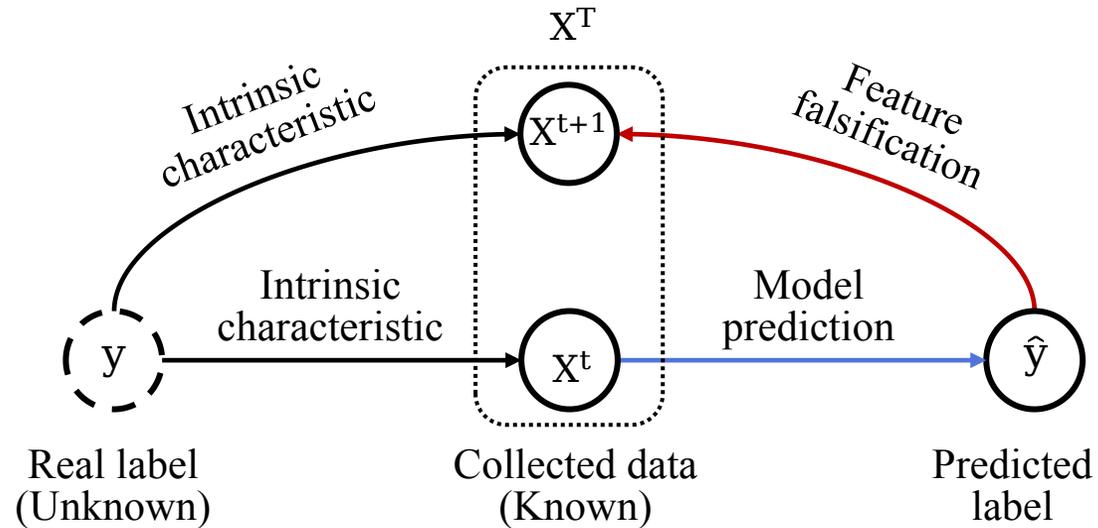


Formalize these processes in fraud detection, leveraging **causal analysis** to investigate the feature falsification



Causal Models of Feature Falsification

◆ Causal diagram of fraud detection

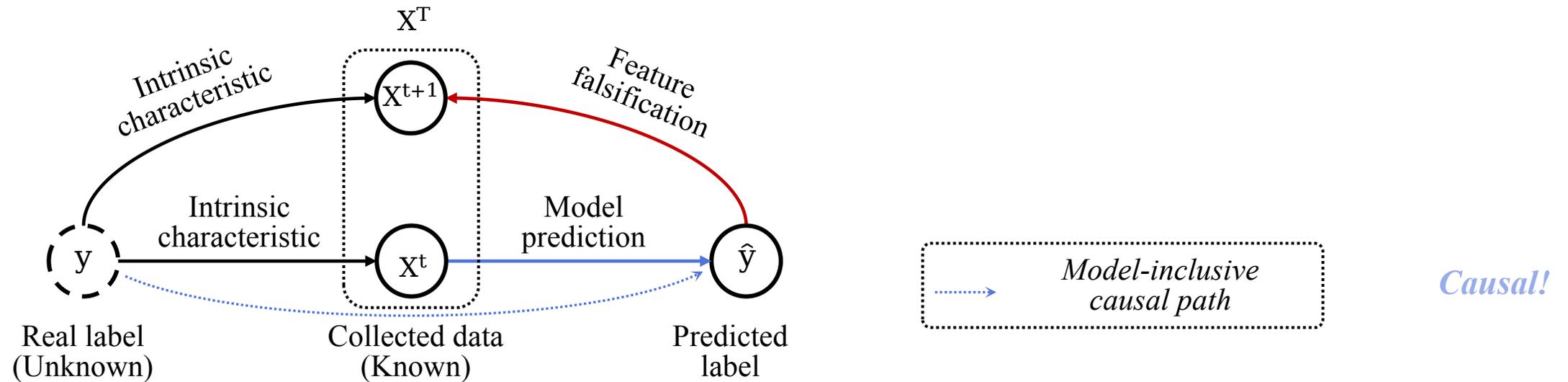


The **four processes** in the detection lifecycle can be formalized as a **causal diagram**, our target is to investigate the **correlation between real label and predicted label**



Causal Models of Feature Falsification

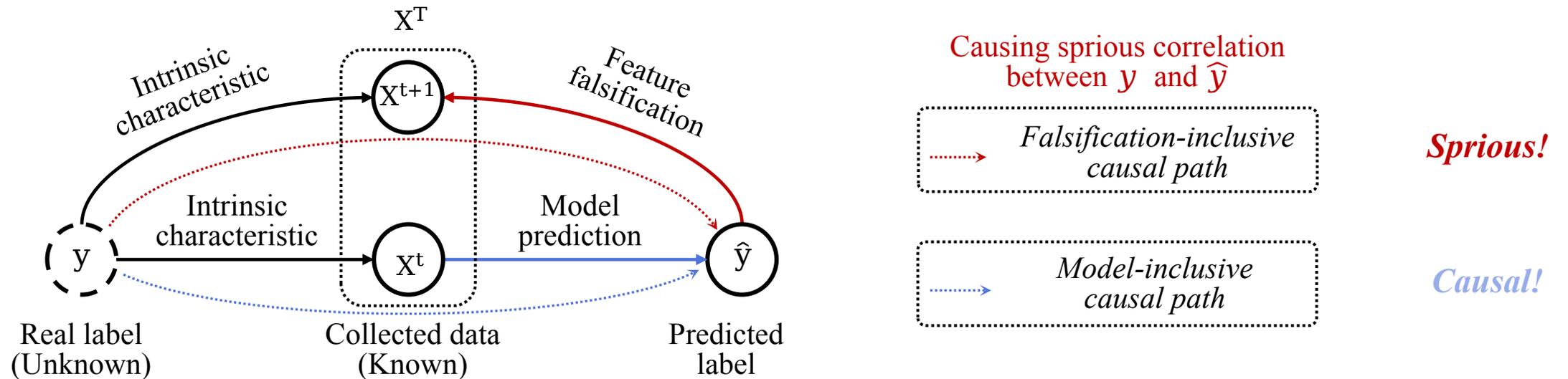
◆ Causal diagram of fraud detection





Causal Models of Feature Falsification

◆ Causal diagram of fraud detection



$y = 1 \rightarrow \hat{y} = 0$ (i.e., successfully evading the detection system) is a spurious correlation

Selection bias (i.e., the insufficient training data) results in spurious correlations



Key issues

Key issues in dealing with feature falsification:

- ◆ Why fraudsters can evade fraud detection via feature falsification?
- ◆ **How to defend against potential feature falsification when collected data is limited?**
- ◆ How to overcome the deficiencies in defense methods when transactions are benign?





Two-player game in the detection



Motivation: Feature selection is a classical method for addressing selection bias



Two-player game in Fraud detection



Motivation: Feature selection is a classical method for addressing selection bias

Key issue: Which feature should be selected?

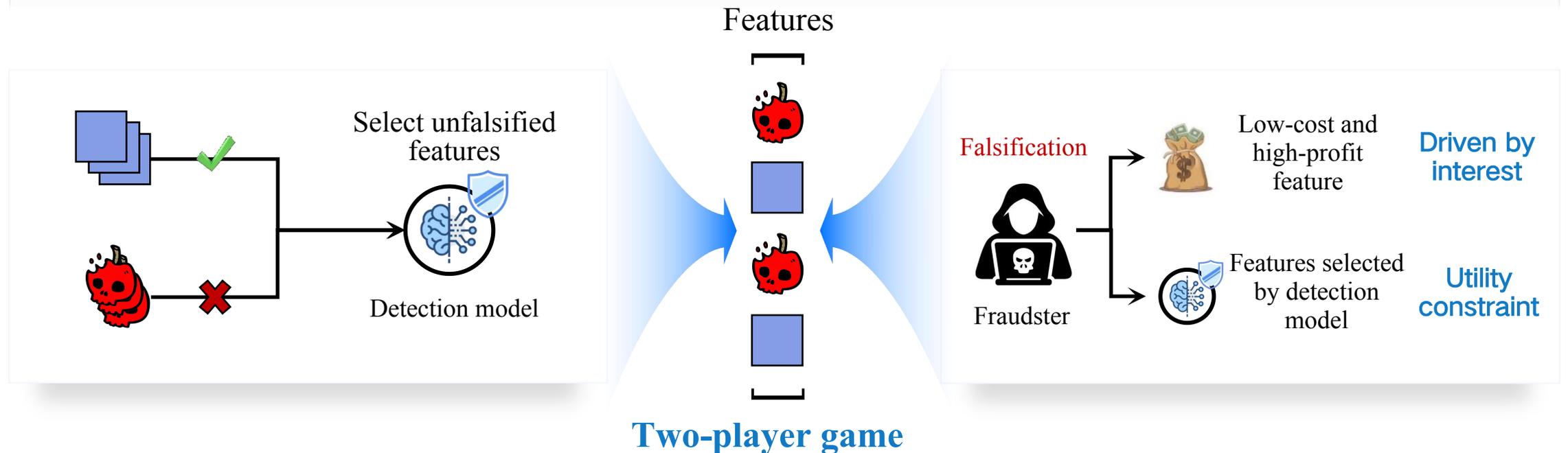


Two-player game in Fraud detection



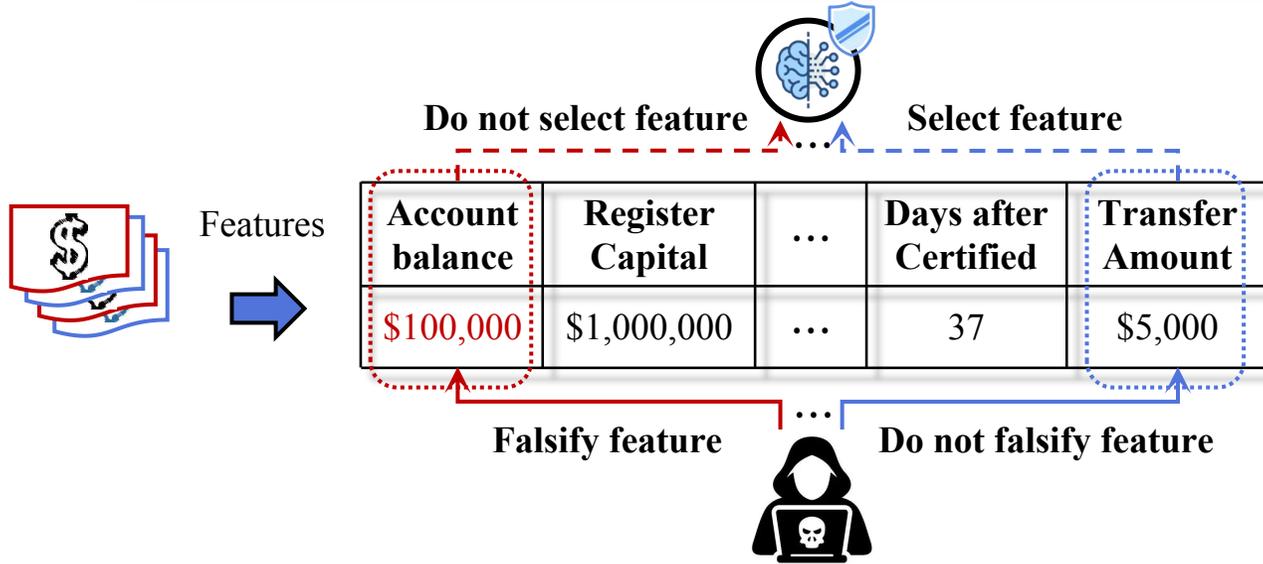
Motivation: Feature selection is a classical method for addressing selection bias

Key issue: Which feature should be selected?



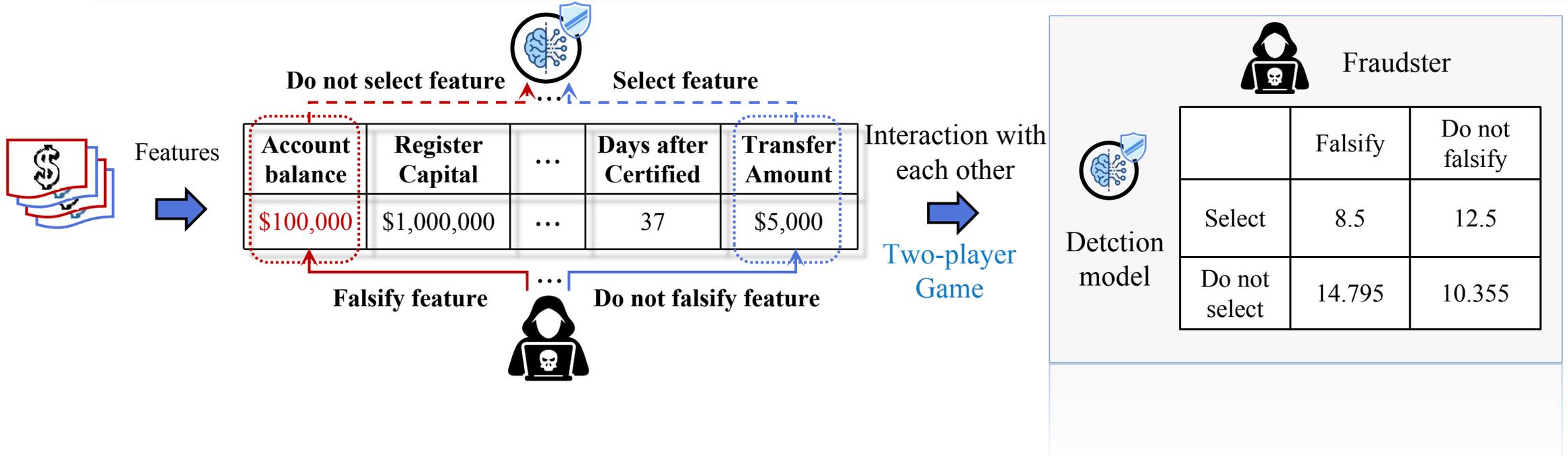


Two-player game in Fraud detection



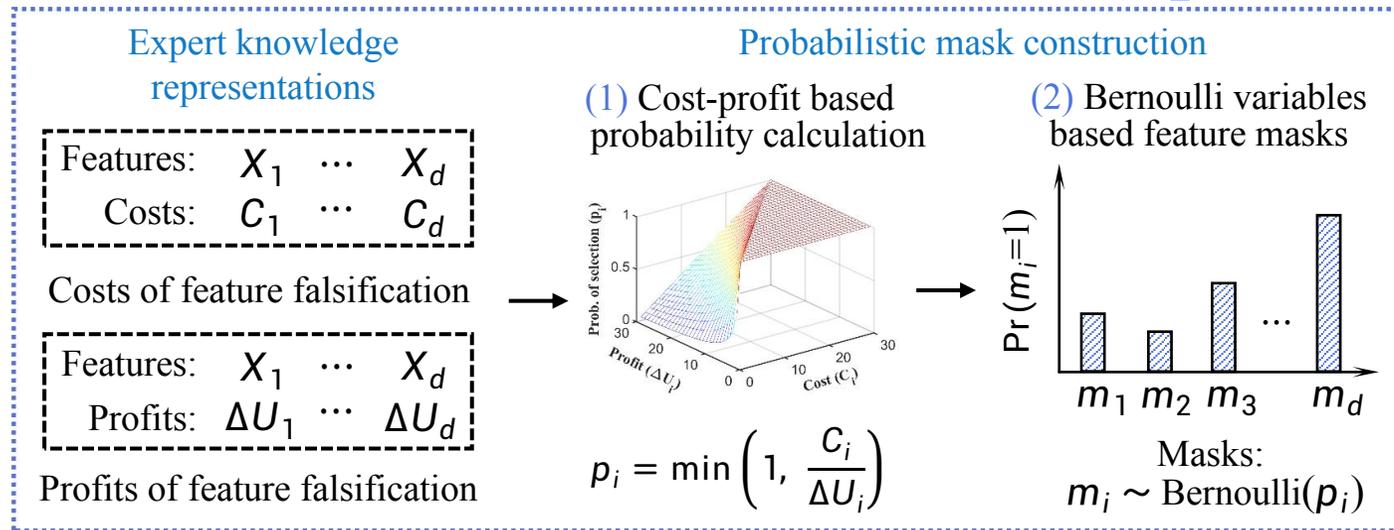
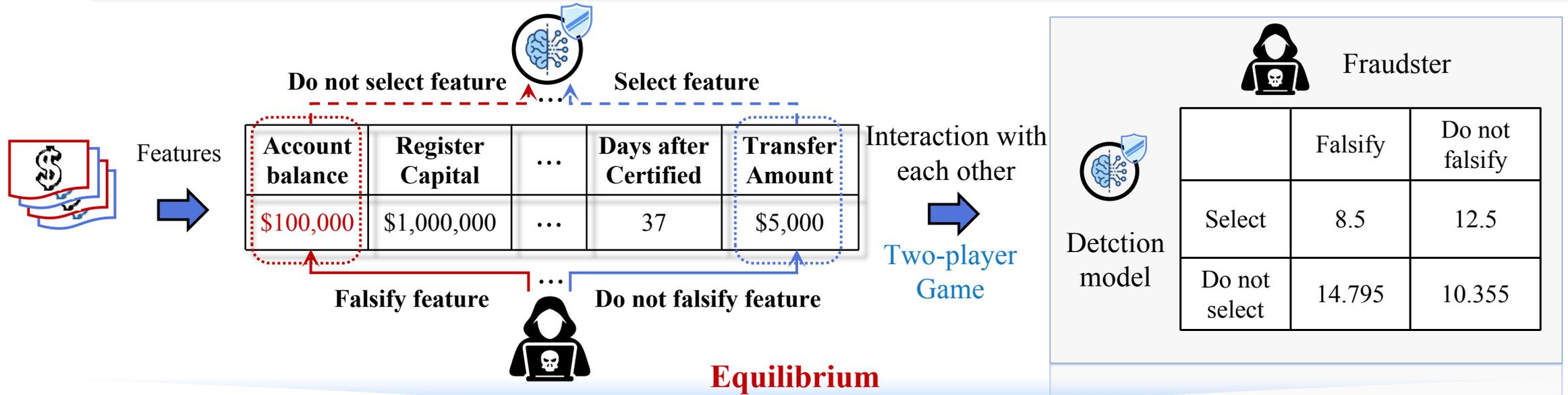


Two-player game in Fraud detection



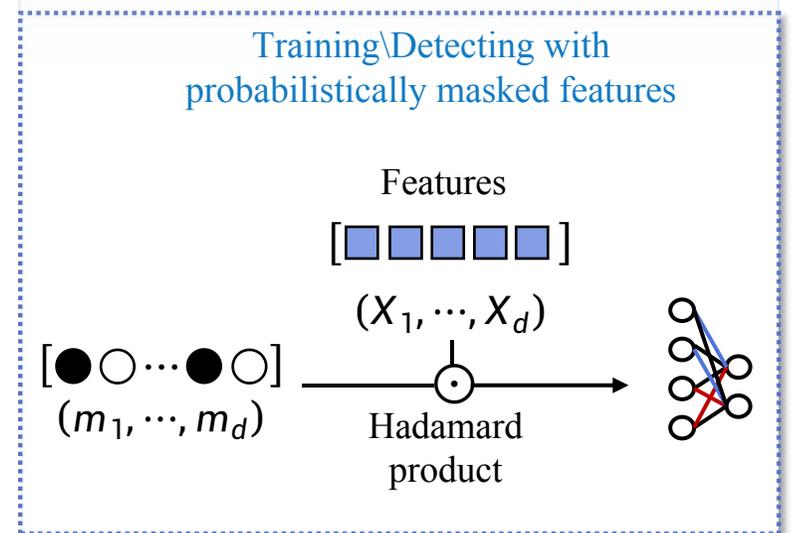
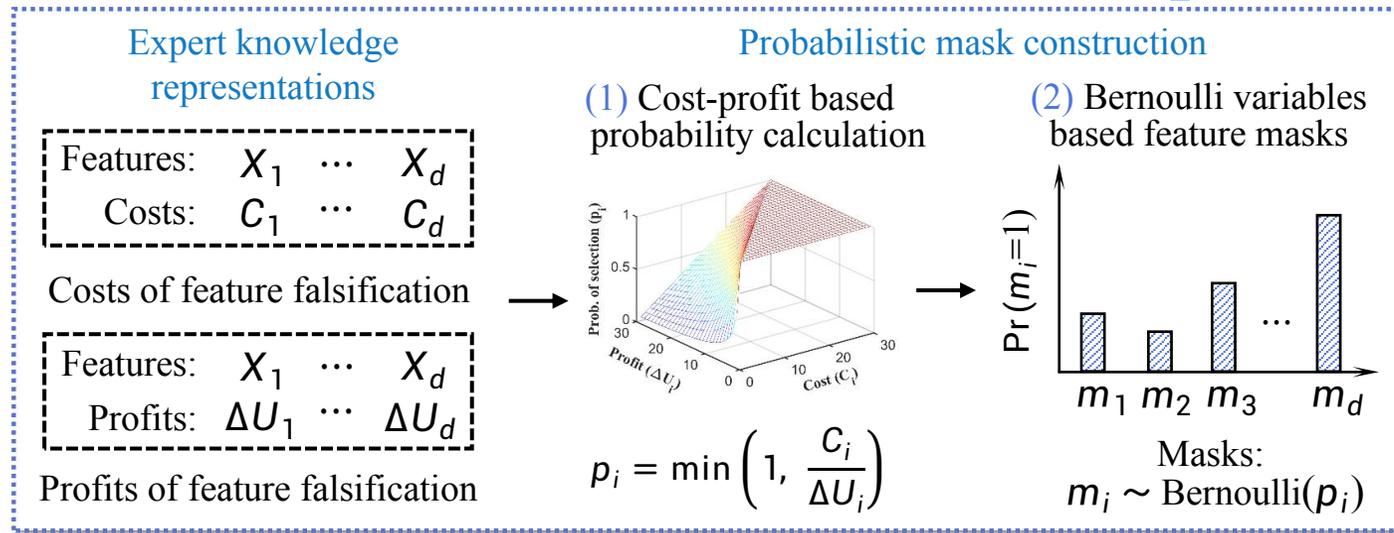
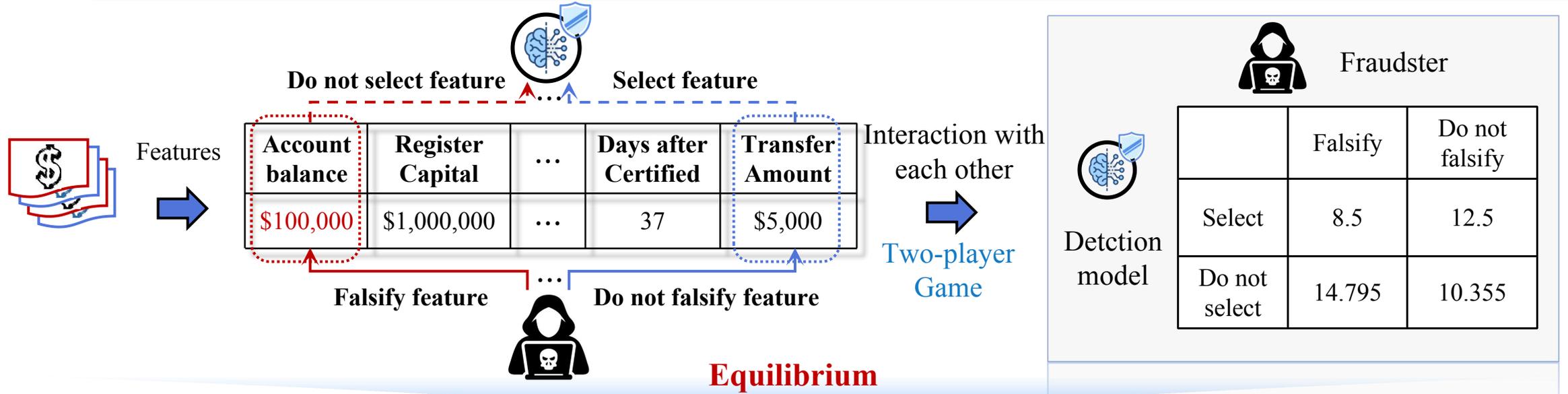


Two-player game in Fraud detection





Two-player game in Fraud detection

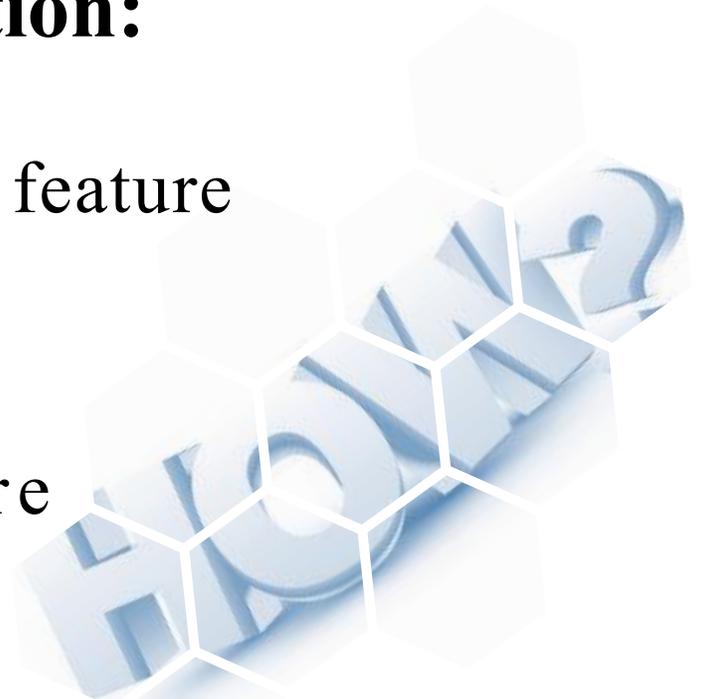




Key issues

Key issues in dealing with feature falsification:

- ◆ Why fraudsters can evade fraud detection via feature falsification?
- ◆ How to defend against potential feature falsification when collected data is limited?
- ◆ **How to overcome the deficiencies in defense methods when transactions are benign?**





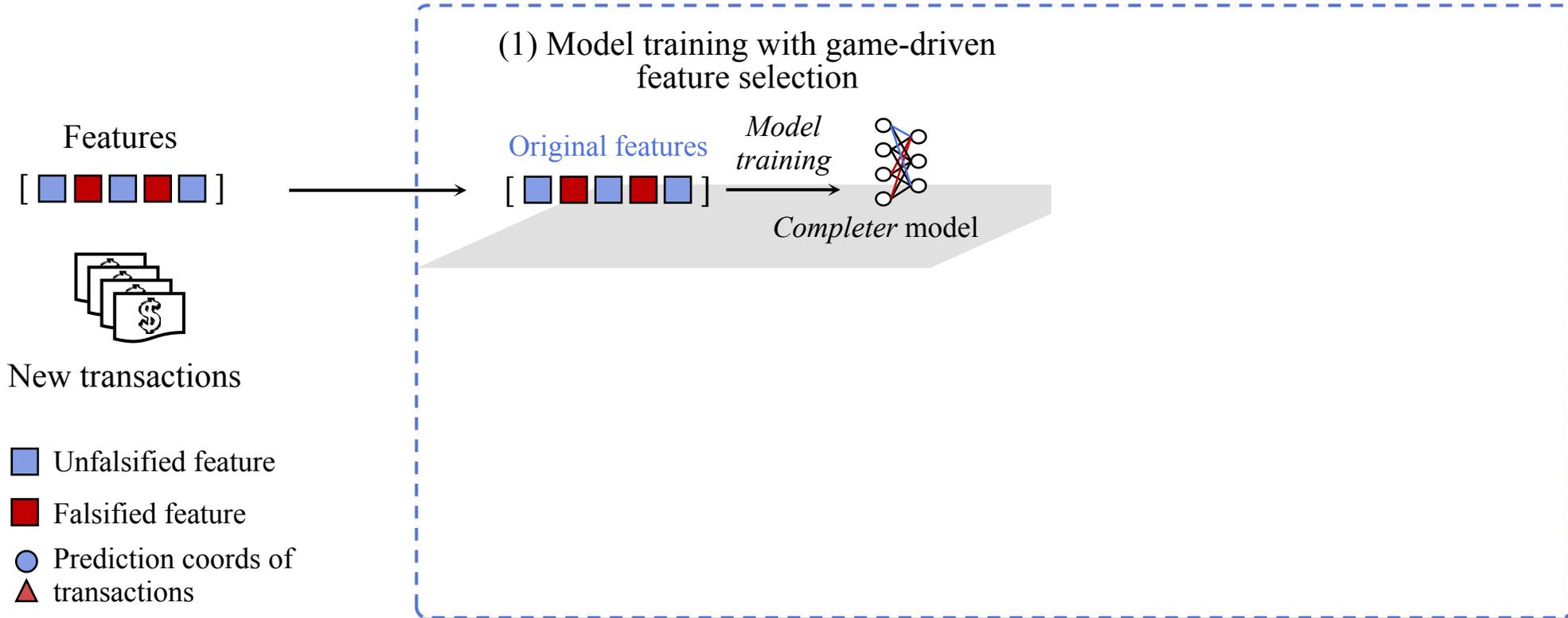
Design of GAMER

Solution: Use different models to detect transactions that they specialized in



Design of GAMER

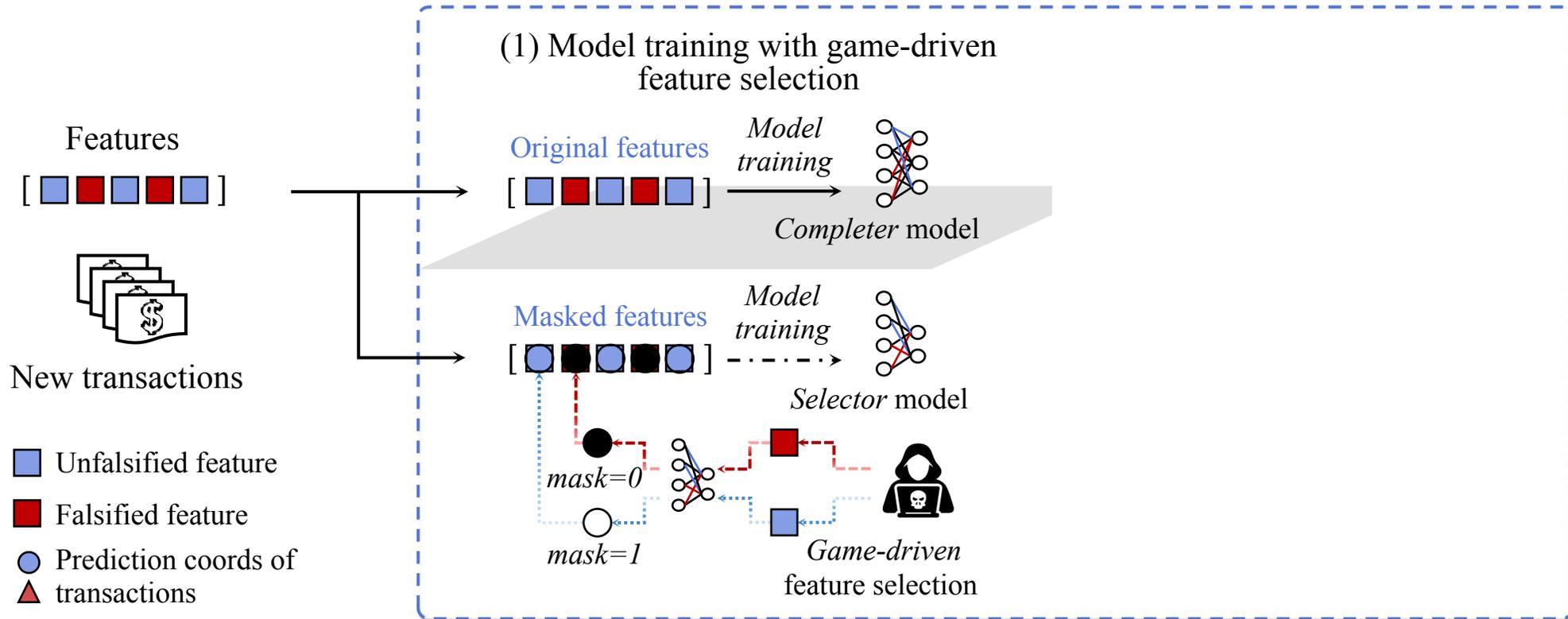
Solution: Use different models to detect transactions that they specialized in





Design of GAMER

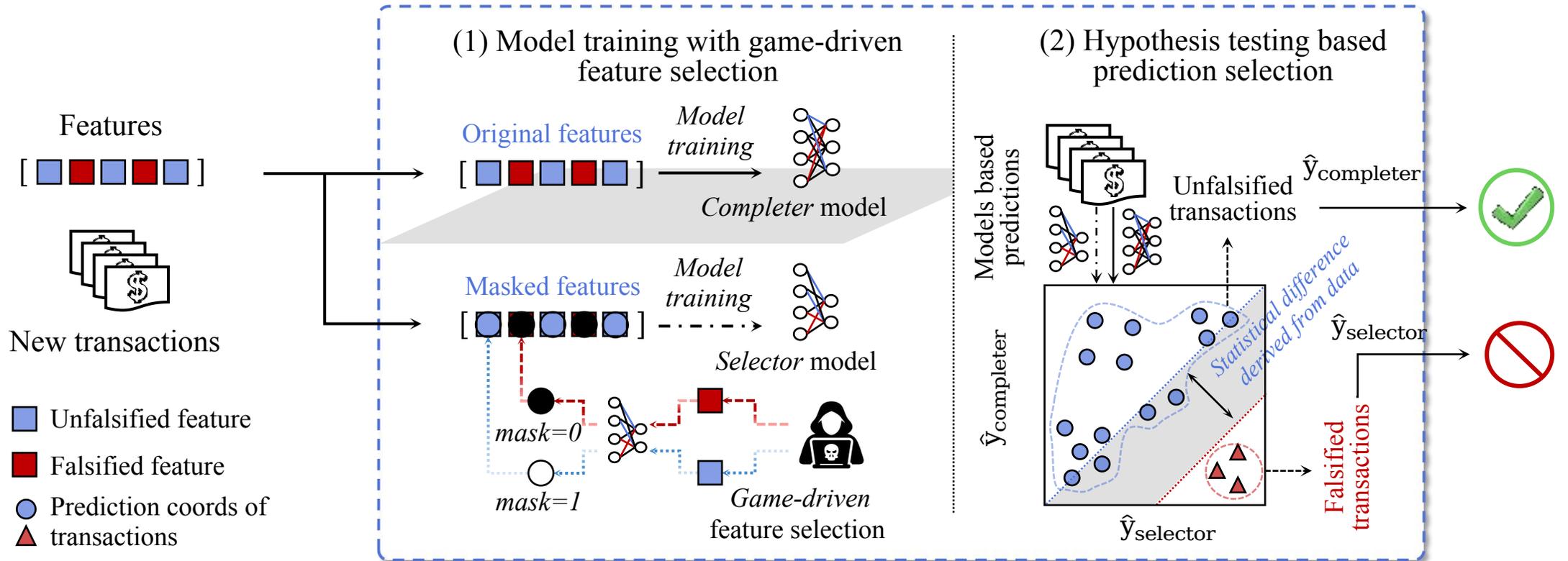
Solution: Use different models to detect transactions that they specialized in

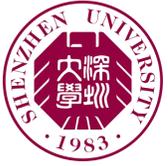




Design of GAMER

Solution: Use different models to detect transactions that they specialized in

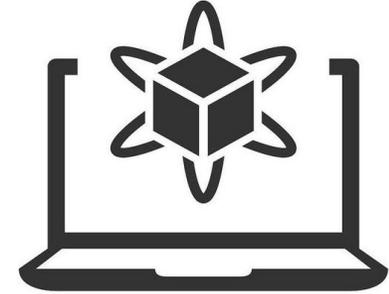




Evaluation

◆ Simulation:

- We **randomly** set costs to each feature and repeat each experiment 30 times to **avoid the bias introduced by cost assignment**



◆ Real-world data:

- Real-world data are collected from **inter-enterprise transactions**
- The labels are collected from **whether transactions are complained about by users**
- Each data has **302 features**, which capture various factors such as transaction amount, transaction frequency, etc



**Real-world
Transactions**



Simulation

Rational Fraudsters (i.e., taking the cost-profit into account)

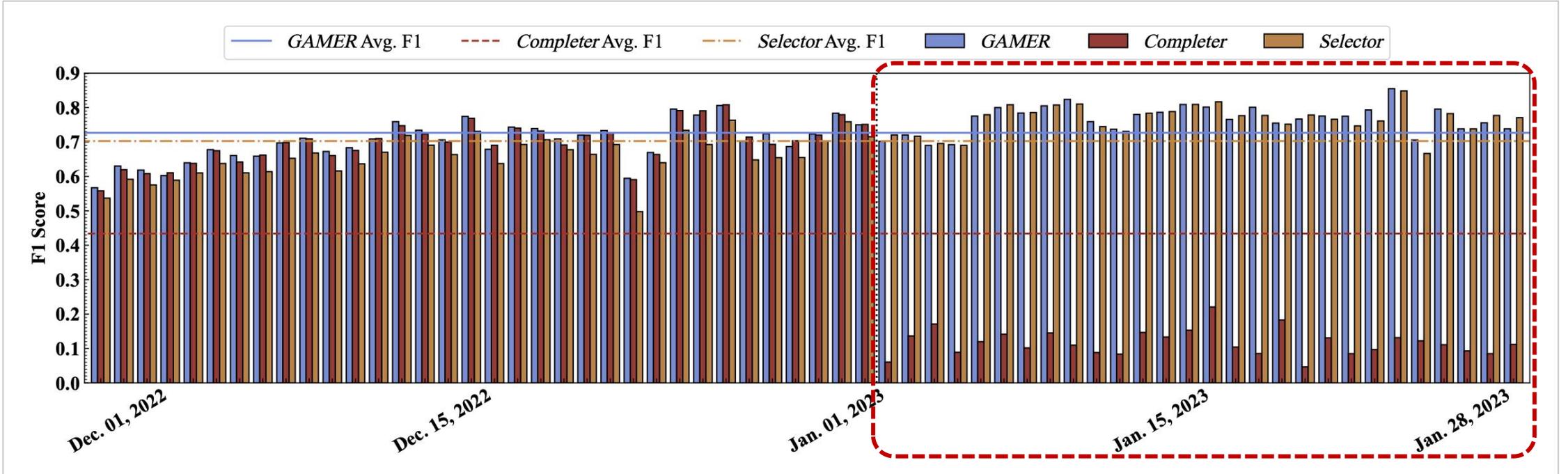
Dataset		AUC					F1 Score				
		At. Cost									
Method		100	200	300	400	500	100	200	300	400	500
CreditCard	Non-robust	0.9629 ₍₋₎	0.8547 ₍₋₎	0.7883 ₍₋₎	0.7199 ₍₋₎	0.6609 ₍₋₎	0.8729 ₍₋₎	0.7504 ₍₋₎	0.5483 ₍₋₎	0.3712 ₍₋₎	0.2827 ₍₋₎
	CW [12]	0.9970 _(↑)	0.9566 _(↑)	0.8914 _(↑)	0.8269 _(↑)	0.7651 _(↑)	0.9107 _(↑)	0.8802 _(↑)	0.8476 _(↑)	0.7937 _(↑)	0.7342 _(↑)
	TARDES [85]	0.9042 _(↓)	0.8748 _(↑)	0.8436 _(↑)	0.8178 _(↑)	0.8025 _(↑)	0.7565 _(↓)	0.7323 _(↓)	0.7163 _(↑)	0.7077 _(↑)	0.7041 _(↑)
	PGD [45]	0.9265 _(↓)	0.8995 _(↑)	0.8746 _(↑)	0.8493 _(↑)	0.8306 _(↑)	0.7467 _(↓)	0.7286 _(↓)	0.7150 _(↓)	0.7032 _(↓)	0.6955 _(↓)
	FGSM [25]	0.9297 _(↓)	0.8940 _(↑)	0.8625 _(↑)	0.8482 _(↑)	0.8405 _(↑)	0.8590 _(↓)	0.8452 _(↑)	0.8365 _(↑)	0.8335 _(↑)	0.8317 _(↑)
	GAMER	0.9828 _(↑)	0.9584 _(↑)	0.9402 _(↑)	0.9247 _(↑)	0.9164 _(↑)	0.9403 _(↑)	0.9115 _(↑)	0.8954 _(↑)	0.8785 _(↑)	0.8687 _(↑)

Irrational Fraudsters

Attack method	Without attack	AutoAttack [19]			Square attack [2] (Black-box attack)		
Metric	F1	F1	ASR	$\frac{At. Cost}{ASR}$	F1	ASR	$\frac{At. Cost}{ASR}$
Def. method							
Non-robust	0.8298 ₍₋₎	0.0323 ₍₋₎	0.9812 ₍₋₎	11.55 ₍₋₎	0.0446 ₍₋₎	0.9738 ₍₋₎	11.58 ₍₋₎
CW [12]	0.6090 _(↓26.6%)	0.3743 _(↑)	0.7210 _(↓26.5%)	4.75 _(↓58.9%)	0.4142 _(↑)	0.6964 _(↓28.5%)	3.79 _(↓67.3%)
TARDES [85]	0.5922 _(↓28.6%)	0.4749 _(↑)	0.4891 _(↓50.1%)	4.34 _(↓62.4%)	0.4796 _(↑)	0.4751 _(↓51.2%)	4.41 _(↓61.9%)
PGD [45]	0.6006 _(↓27.6%)	0.4678 _(↑)	0.5442 _(↓44.5%)	4.50 _(↓61.0%)	0.4677 _(↑)	0.5381 _(↓44.8%)	4.76 _(↓58.8%)
FGSM [25]	0.5234 _(↓36.9%)	0.0323 ₍₋₎	0.9670 _(↓1.44%)	9.77 _(↓15.4%)	0.0367 _(↓)	0.9609 _(↓1.32%)	10.20 _(↓11.9%)
GAMER	0.8012 _(↓3.45%)	0.6356 _(↑)	0.4145 _(↓57.8%)	23.26 _(↑101.4%)	0.7308 _(↑)	0.2826 _(↓70.9%)	34.67 _(↑199.4%)



Evaluation on real-world data



Using the **equilibrium** in detection, GAMER increases the **F1 score** by **67.5%** on average for two-month fraud detection



Conclusion

- ◆ Insufficient training data (i.e., **selection bias**) causes the spurious correlations between real label and predicted label ($y = 1 \rightarrow \hat{y} = 0$)
 - Using **feature selection** to overcome selection bias
- ◆ Using **Two-player game** to decide the **selection strategies** of detection model
- ◆ Using **different model** to detect the transactions **they specilized in**

Email: tanqi@szu.edu.cn

Thanks! Questions?