

# Smart-LLaMA-DPO: 强化大模型进行可解释智能合约漏洞检测

于磊, 黄智榕, 袁航, 成仕骐, 杨立\*, 张凤军\*, 沈宸杰,  
马佳佳, 张竞元, 陆俊逸, 左春

The 34th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2025), 2025

联系人: 于磊、杨立、张凤军

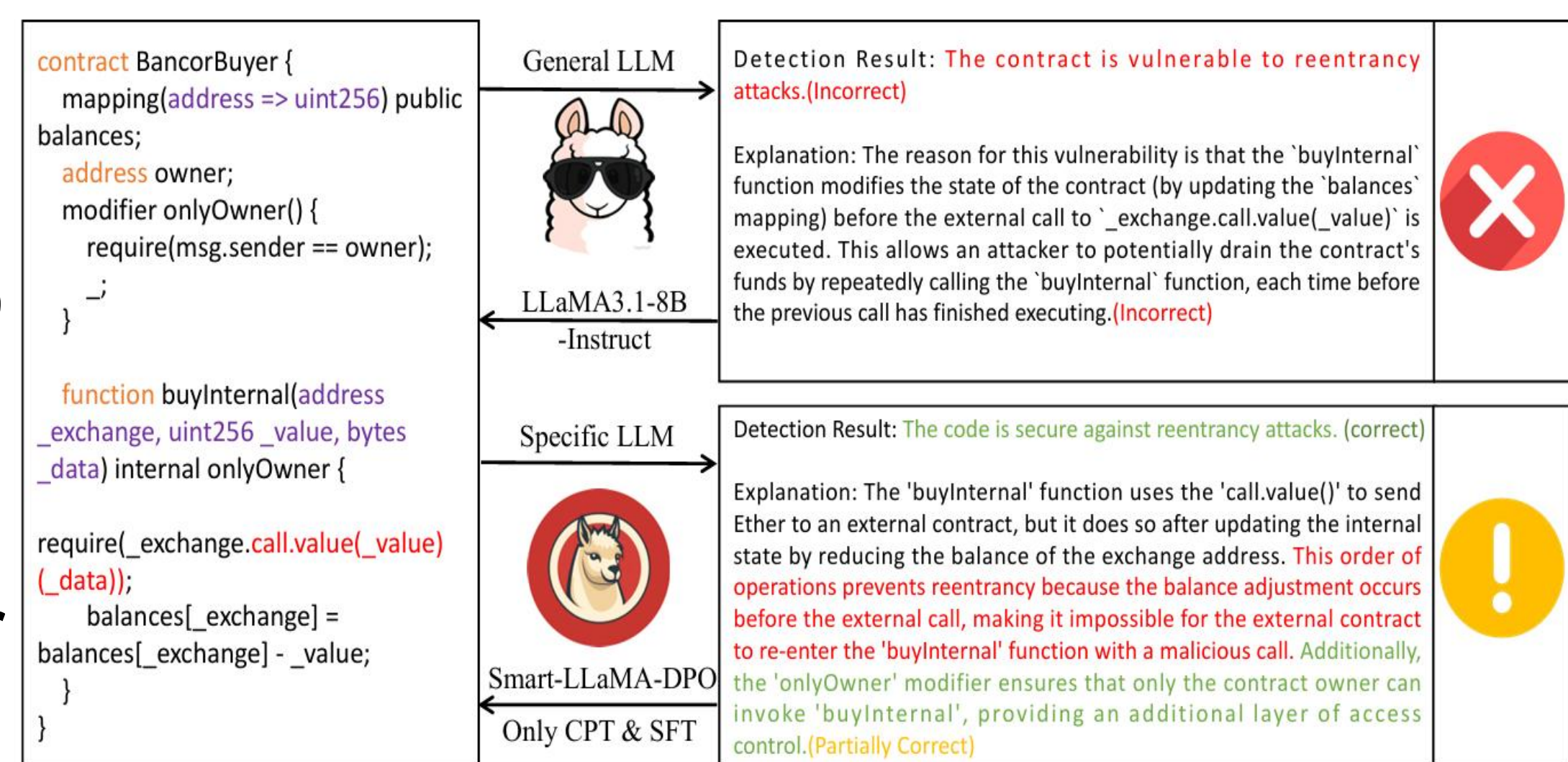
联系方式: {yulei2022,yangli2017,fengjun}@iscas.ac.cn

## Background

- Security vulnerabilities in smart contracts (e.g., Reentrancy, Timestamp, Delegatecall, Price Oracle Manipulation) have caused huge financial losses.

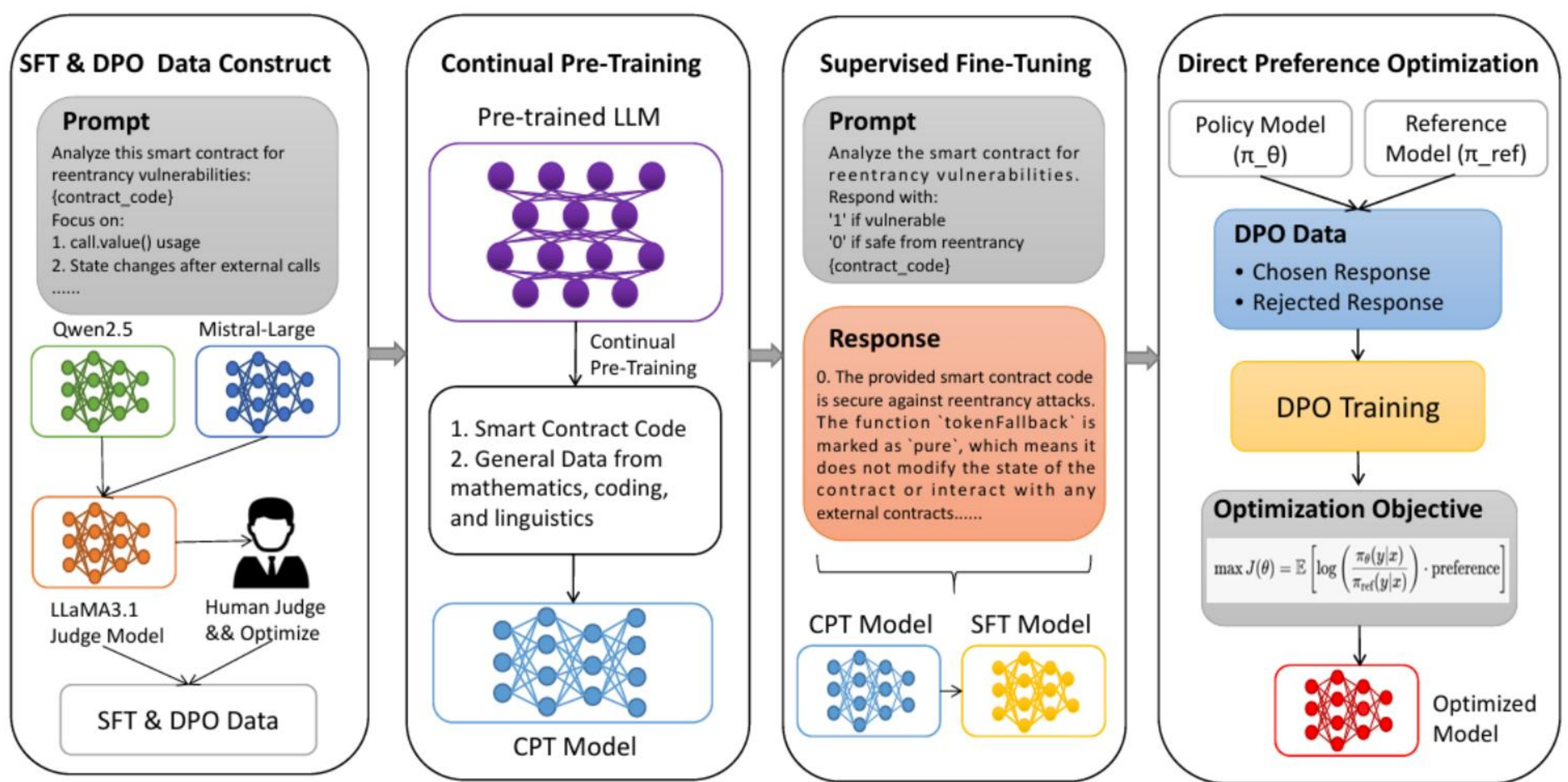
## Motivation

- Existing works lack comprehensive, high-quality datasets and struggle to generate **accurate** explanations.
- LLMs often misinterpret **execution order** and provide **inconsistent** or superficial analysis.



## Approach: Smart-LLaMA-DPO

- CPT**: Large-scale Solidity contracts, deepens model's domain knowledge.
- SFT**: Expert-verified data, labels, detailed explanations, vulnerability locations for 4 main types + 7 machine-unauditable types.
- DPO**: Trains on expert-preferred vs. suboptimal explanation pairs, maximizing probability of human-like, high-quality explanations.



## Evaluation

- RQ1&RQ2: Detection Performance**
  - Smart-LLaMA-DPO **outperforms all baselines on 4 major vulnerability types and 7 machine-unauditable types.**
- RQ3: Ablation Study**
  - Both **Continual Pre-training** and **Direct Preference Optimization** are essential.
- RQ4: Human Evaluation**
  - Ours generates more **accurate, thorough, and clear** explanations than baselines.
- RQ5: Real-World Applicability**
  - Case studies show Smart-LLaMA-DPO avoids false positives/negatives and delivers context-aware security advice.

Types	Metric	Base	Base+CoT	w/o dpo	w/o cpt	w/o dpo & cpt
RE	Acc(%)	94.47	94.47	83.40	86.38	90.21
	F1(%)	88.50	88.50	73.47	77.78	79.65
TD	Acc(%)	95.54	93.75	80.80	82.14	69.64
	F1(%)	96.43	95.30	85.32	86.39	69.64
IO	Acc(%)	94.65	93.42	89.71	83.95	85.19
	F1(%)	88.29	86.21	81.75	53.01	56.10
DE	Acc(%)	94.12	94.12	93.53	93.53	91.76
	F1(%)	84.85	84.85	83.08	83.08	78.12
MU	Acc(%)	90.74	91.53	78.84	86.24	72.22
	F1(%)	83.41	85.19	71.22	80.60	66.88

	Correctness				Thoroughness				Clarity			
	1	2	3	4	1	2	3	4	1	2	3	4
LLM Evaluation												
LLaMA3.1-8B	116	201	165	579	42	229	332	458	30	204	578	249
FTSmartAudit	101	234	161	565	53	167	351	490	41	165	454	401
iAudit	135	129	84	713	48	216	211	586	27	47	240	747
Ours	56	86	85	834	9	96	225	731	2	25	469	565
Human Evaluation												
LLaMA3.1-8B	76	303	326	356	70	266	457	268	29	212	623	197
FTSmartAudit	127	234	352	348	126	184	523	228	28	165	482	386
iAudit	61	255	357	388	48	241	584	188	27	143	418	473
Ours	19	181	215	646	18	153	346	544	9	48	435	569