

Limin Yang

540-998-9158 | liminy2@illinois.edu | <https://liminyang.web.illinois.edu> | GitHub: [whyisyong](#) | [Google Scholar](#)

EDUCATION

University of Illinois (UIUC), Ph.D. in Computer Science, advisor: Gang Wang	Aug. 2019 – Aug. 2023
Virginia Tech, Ph.D. in Computer Science, advisor: Gang Wang	Aug. 2018 – Jul. 2019
East China Normal University, M.S. Study in Computer Science	Sep. 2015 – Jun. 2018
East China Normal University, B.Eng. in Computer Science	Sep. 2011 – Jun. 2015

INTERNSHIPS

IBM Research, Visiting Scholar (Research Intern), Hybrid Cloud Team	May 2022 – Aug. 2022
<ul style="list-style-type: none">Cleaned a noisy real-world network intrusion dataset (25 million traffic/day) from National Supercomputing CenterSemi-automatically labeled the dataset, defined 65 features for network logs, and summarized 279 security incidentsBuilt anomaly detection models (per host) with LSTM autoencoder and applied on 2 past real-world attacks	
TikTok, Security Engineering Intern, Threat XDR Team	May 2021 – Aug. 2021
<ul style="list-style-type: none">Lark Email Spam Rule System: detected 5,000 more spam/week by adding 25 patterns and 20 rulesProtected 50 client companies and 1 million emails/week by allowlisting 300 domains with semi-automationClustered similar emails based on user actions and increased the size of user labeled ground-truth by 100%Contributed to 2 urgent incidents response: a vendor offline and a spam campaign	

SELECTED PROJECTS

Selective ML Backdoor Attack Published in IEEE S&P'23	Sep. 2020 – Jan. 2022
<ul style="list-style-type: none">Proposed a new backdoor attack against Android malware classifiers: only a specific malware family is misclassifiedAchieved 80%–98% attack success rates by alternate optimization with a customized loss functionDefeated 4 recent backdoor defenses while traditional backdoor cannot bypass the detections	
Concept Drift Detect and Explain Published in USENIX Security'21	Jun. 2019 – Jun. 2020
<ul style="list-style-type: none">Leveraged contrastive learning and autoencoder to detect concept drift samples from previously unseen classesExplainable AI: proposed a novel distance-based explanation to find 45/1000 features making a sample outlierIncreased detection rate to $F_1 = 96\%$ versus state-of-the-art ($F_1 \leq 80\%$) on malware and network datasetsIdentified 161/165 unseen families on a company Blue Hexagon's Windows PE malware database	
VirusTotal Reliability Published in IMC'19 and USENIX Security'20	Feb. 2019 – Nov. 2019
<ul style="list-style-type: none">Controlled 66 phishing sites to measure the label inconsistencies and dynamics between vendors and VirusTotalMeasured the label dynamics of 14,000+ PE malware from 65 vendors via daily snapshots over one yearOffered insights and suggestions on a more proper use of VirusTotal (received 100+ citations in total)	

SELECTED PUBLICATIONS (CITATIONS: 400+, H-INDEX: 9)

- [Submitted to **IEEE S&P'24**] **1st author**. *Title anonymized for double-blind submission.*
- [**IEEE S&P'23**] **1st author**. “*Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers*”.
- [**IEEE S&P'23**] 3rd author. “*Practitioners' Perception of ML-Based Security Tools and Explanations*”.
- [**USENIX Security'21**] **1st author**. “*CADE: Detecting and Explaining Concept Drift Samples for Security Applications*”. Artifact Evaluated.
- [Deep Learning and Security'21] **1st author**. “*BODMAS: An Open Dataset for Learning based Temporal Analysis of PE Malware*”. Dataset requested by **85 institutions (about 120 research groups)**.
- [**USENIX Security'20**] 3rd author. “*Measuring and Modeling the Label Dynamics of Online Anti-Malware Engines*”.
- [**IMC'19**] 2nd author. “*Opening the Blackbox of VirusTotal: Analyzing Online Phishing Scan Engines*”.
- [**USENIX Security'18**] 3rd author. “*Understanding the Reproducibility of Crowd-reported Security Vulnerabilities.*”.

TECHNICAL SKILLS

Languages:	Python, C++, C, Shell, SQL (Postgres, Hive), Ruby
Security:	ML-based malware detection, Network IDS, Spam, Phishing, Bug reproduction
Deep Learning:	Keras, Tensorflow, PyTorch
Frameworks:	Hadoop, PySpark, MongoDB, Flask, Scrapy, Alexa Skills, Elasticsearch, Ruby on Rails
Developer Tools:	Linux, Git, VS Code, Jupyter Notebook, tmux, AWS, Vim, Docker, VirusTotal, Nmap
Libraries:	Scikit-learn, LightGBM, NumPy, pandas, Matplotlib