

# Huzaifa Arif

(518) 961-8482 | huzaifaarif20@gmail.com | Troy, NY | huzaifa-arif.github.io

## PROFILE SUMMARY

A fifth-year Ph.D. candidate and researcher with multiple first-author publications in trustworthy AI. My work focuses on exposing and mitigating privacy and safety vulnerabilities in Large Language Models, demonstrated through research at IBM and LLNL. **I am seeking a full-time Research Scientist position starting December 2025** where I can apply my expertise in LLM alignment and privacy to build verifiably safe AI systems.

## EDUCATION

<b>Ph.D. in Electrical and Computer Systems Engineering</b> Rensselaer Polytechnic Institute   Advisor: Alex Gittens	<b>Expected Dec 2025</b>
<b>B.S. in Electrical Engineering</b> Lahore University of Management Sciences (LUMS)   Dean's Honor List	<b>Graduated with High Merit</b>

## RESEARCH EXPERIENCE

<b>IBM T.J. Watson Research Center</b>   <i>AI Research Extern - Trustworthy AI</i>	<b>May 2025 – Aug 2025</b>
<ul style="list-style-type: none"><li>Developed parameter-efficient steering method for LLM safety using lightweight prefix tuning (<math>&lt;0.01\%</math> parameters)</li><li>Combined SFT and DPO training to achieve 80% reduction in toxic responses across Llama and Aya model families</li><li>Extending framework to mitigate demographic bias and PII leakage (manuscript in preparation for ICLR 2026)</li></ul>	
<b>Rensselaer Polytechnic Institute</b>   <i>Research Assistant</i>	<b>Aug 2022 – Present</b>
<ul style="list-style-type: none"><li>Discovered novel "association leakage" vulnerability in LLMs where jail-broken models reveal sensitive data</li><li>Created SPARK attack method achieving 3x increase in private data recall using only 0.01% trainable parameters (manuscript in preparation for ICLR 2026)</li><li>Published work on federated fair learning via kernel regularization in TMLR 2025</li></ul>	
<b>Lawrence Livermore National Laboratory</b>   <i>Data Science Intern</i>	<b>May 2024 – Aug 2024</b>
<ul style="list-style-type: none"><li>Exposed critical vulnerabilities in FourCastNet weather prediction model through novel adversarial attacks</li><li>Demonstrated susceptibility to localized evasion attacks leading to faulty predictions (Accepted at AAAI 2025)</li></ul>	
<b>IBM T.J. Watson Research Center</b>   <i>AI Research Extern</i>	<b>Jun 2023 – Aug 2023</b>
<ul style="list-style-type: none"><li>Developed PEEL method for backward feature inversion in residual networks, achieving 10x MSE improvement</li><li>Published findings at IEEE SATML 2025; filed patent currently under review</li></ul>	
<b>IBM T.J. Watson Research Center</b>   <i>AI Research Extern</i>	<b>Jun 2022 – Aug 2022</b>
<ul style="list-style-type: none"><li>Created Reprogrammable-FL framework adapting model reprogramming to differentially private federated learning</li><li>Achieved 60% accuracy improvement under same privacy budget (Published at IEEE SATML 2023)</li><li>Patent filed (US20240256894A1); authored book chapter in "Federated Learning for Medical Imaging"</li></ul>	

## SELECTED PUBLICATIONS

- H. Arif**, A. Gittens, & P.-Y. Chen. "Reprogrammable-FL: Improving Utility-Privacy Tradeoff in Federated Learning via Model Reprogramming." In *IEEE Conference on Secure and Trustworthy Machine Learning (SATML)*, 2023.
- H. Arif**, A. Gittens, K. Murugesan, P. Das, & P.-Y. Chen. "Peel the Layers and Find Yourself: Revisiting Inference-time Data Leakage for Residual Neural Networks." In *IEEE Conference on Secure and Trustworthy Machine Learning (SATML)*, 2025.
- H. Arif**, A. Gittens, J. Diffenderfer, P.-Y. Chen, & B. Kailkhura. "Forecasting Fails: Unveiling Evasion Attacks in Weather Prediction Models." In *AAAI Workshop on Artificial Intelligence for Science and Engineering (AI2SE)*, 2025.
- H. Arif**, P.-Y. Chen, K. Murugesan, & A. Gittens. "Fair Federated Learning via Stochastic Kernel Regularization." *Transactions on Machine Learning Research (TMLR)*, 2025. (Journal)
- H. Arif**, A. Gittens, K. Murugesan, I. Ko, P. Das, & P.-Y. Chen. "The Safety Patch: Lightweight Prefix Tuning with DPO for Controllable LLMs." (In Preparation: ICLR 2026)
- H. Arif**, A. Gittens, K. Murugesan, P. Das, & P.-Y. Chen. "SPARK: Amplifying Association Leakage in LLMs via Soft Prompting and Attention Head Steering." (In Preparation: ICLR 2026)

## PATENTS and INTELLECTUAL PROPERTY

---

- US20240256894A1: "Differentially Private Federated Learning Using Model Reprogramming" (Under Review)
- Patent Application: "Method for Quantifying Private Leakage in Pretrained Neural Networks" (Filed)
- **Book Chapter:** Utility Privacy Tradeoff in the book "*Federated Learning for Medical Imaging*"

## HONORS & AWARDS

---

- Belsky Award for Computational Science and Engineering 2025 (Graduate Research Excellence)
- Top-5 Candidate in ECSE Research Qualifier Exam at RPI
- Travel Support Awards: IEEE SATML 2023 & 2025

## Technical Skills

---

**Programming & Languages:** Python, C++, SQL, MATLAB

**AI/ML Frameworks:** PyTorch, TensorFlow, Keras, Scikit-learn

**Data Science Libraries:** NumPy, Pandas, SciPy, Matplotlib, Seaborn

**Research Areas:** Trustworthy AI, Large Language Models, Federated Learning, Differential Privacy, Reinforcement Learning, Constrained Optimization, Stochastic Optimization, Transfer Learning, Spatio-temporal Data Analysis, NLP, Deep Learning, Alignment, RLHF, Prompt Tuning, Soft Prompt Engineering, Model Reprogramming

**Tools & Other:** Git/GitHub, Jupyter Notebooks, LaTeX, Tableau

## SERVICE

---

**Reviewer:** ICLR 2025, ICASSP 2023/2025, AISTATS 2023, IEEE MLSP 2023

**Teaching Assistant (4.5/5 avg rating):** Machine Learning, Embedded Control, Signals & Systems (2021-2024)