

# HOANG ANH JUST

Google Scholar  $\diamond$  [justhoanganh.com](https://justhoanganh.com)  $\diamond$  [just@vt.edu](mailto:just@vt.edu)

## EDUCATION

---

Virginia Polytechnic Institute and State University, Blacksburg, VA  
Doctor of Philosophy in Computer Engineering

Fall 2021 - Present

❖ Advisor: Prof. Ruoxi Jia

Gettysburg College, Gettysburg, PA

Bachelor of Science in Computer Science and Mathematics

Fall 2017 - Spring 2021

❖ GPA: 4.1/4.3 CS GPA: 4.1/4.3 MATH GPA: 4.1/4.3

## FOCUS AREAS

---

Data {Valuation, Selection, Quality, Efficiency}, Predicting Model/Prediction, AI Privacy, Additive Combinatorics

## RESEARCH EXPERIENCE

---

Virginia Tech Graduate Research Assistant

Fall 2021 - Present

Mentor: Prof. Ruoxi Jia

Paper Title: Get more for less: Principled Data Selection for Warming Up Fine-Tuning in LLMs

- ➔ Developed a scalable data selection method to *pre-fine-tune* a pretrained large language model (LLM) by selecting (unlabeled) data that can shift the source distribution to better align with the target distribution.
- ➔ Published at the International Conference on Learning Representations (ICLR), 2024

Paper Title: Performance Scaling via Optimal Transport: Enabling Data Selection from Partially Revealed Sources

- ➔ Proposed a performance estimator for a model trained on any data composition given only sample information and a scaling law to predict performance on larger scales, which effectively finds the optimal composition of data sources for any target data size.
- ➔ Published at the Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS), 2023
  - ASR Data Selection from Multiple Sources: A Practical Approach on Performance Scaling  
Extension of performance scaling method to automatic speech recognition (ASR) data.
  - Published at the 3<sup>rd</sup> Efficient Natural Language and Speech Processing (ENLSP-III) Workshop @ NeurIPS, 2023

Paper Title: NARCISSUS: A Practical Clean-Label Backdoor Attack with Limited Information

- ➔ Launched an efficient (poisoning 0.5% of the target class and 0.05% of the entire training dataset) and stealthy (hard to detect) backdoor attack, which requires only knowledge of the target class to successfully deploy the attack.
- ➔ Published at the ACM Conference on Computer and Communications Security (CCS), 2023

Paper Title: PRIVMON: Real-Time Platform-Agnostic Privacy Leakage Detection for Machine Learning Models

- ➔ Established an efficient real-time detection system to membership inference attacks which prevents attackers from inferring sensitive data used for model training.
- ➔ Published at the International Symposium on Research in Attacks, Intrusions and Defenses (RAID), 2023

Paper Title: 2D-Shapley: A Framework for Fragmented Data Valuation Algorithms

- ➔ Proposed a novel, efficient approach to fine-grained data analysis, which values the quality of each feature of each data point with theoretical grounding.
- ➔ Published at the International Conference on Machine Learning (ICML), 2023

Paper Title: LAVA: Data Valuation without Pre-Specified Learning Algorithms

- ➔ Introduced an efficient data quality valuation method through adopting a modified class-wise Wasserstein distance, which is robust to noisy, mislabeled, and poisoned data without requiring any model training.
- ➔ Published at the International Conference on Learning Representations (ICLR), (Spotlight), 2023

Paper Title: ModelPred: A Framework for Predicting Trained Model from Training Data

- ➔ Developed a set-function based neural network which can predict model weights from the training dataset of any size. This method enables efficient applications for data valuation, data selection, or data memorization, which requires multiple model re-trainings.

- ✦ Published at the IEEE Conference on Secure and Trustworthy Machine Learning (SatML), 2023

Paper Title: Label-Only Model Inversion Attacks via Boundary Repulsion

- ✦ Designed a novel practical model inversion attack which recovers sensitive data by accessing only labels of the model output without additional information.
- ✦ Published at the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022

Intermediate Research in Mathematics, Gettysburg College

January 2020 - May 2021

Mentor: Prof. Béla Bajnok

Paper Title: On perfect bases in finite Abelian groups

- ✦ Proved that for sets of size greater than 3, there are no perfect restricted 2-basis in  $\mathbb{Z}_n$ . Showed that for only sets of size smaller equal to 3 there exists a perfect restricted 2-basis in  $\mathbb{Z}_n$ , proving by contradiction knowing that  $\mathbb{Z}_n$  is closed under both addition and subtraction.
- ✦ Published at Involve, a Journal of Mathematics: 12/2022

## CURRENT PROJECTS

---

- 🔗 Enhancing LLM Reasoning through a Debiasing Framework  
In Progress
- 🔗 Efficient Instruction Tuning Data Creation and Selection for Generalized LLM Performance
- 🔗 Efficient Instruction Tuning Data Creation and Selection for Generalized LLM Performance

## TEACHING EXPERIENCE

---

The Bradley Department of Electrical and Computer Engineering, Virginia Tech  
Graduate Teaching Assistant

Fall 2021 - Fall 2022

Mathematics Department, Gettysburg College  
Peer Learning Associate

Fall 2018 - Spring 2021

Computer Science Department, Gettysburg College  
Teaching Assistant and Grader

Fall 2018 - Spring 2021

## WORK EXPERIENCE

---

Musselman Library, Gettysburg College

May 2019 - July 2019

Digital Scholarship Summer Fellow

Supervisor: R. C. Miessler

- | Created the Augmented Reality College tour for mobile application that recognizes the College buildings, shows historical facts about the campus and provides archival photos of the buildings.
- | Developed in Unity with Wikitude SDK to support object recognition and image recognition.
- | Presented at the Digital Scholarship Student Symposium at Lafayette College: 06/2019.

IT Department, Gettysburg College

May 2018 - July 2019

Digital Technology Student Fellow

Supervisor: Eric Remy

- | Created the Virtual Tour of the Lincoln Cemetery in Gettysburg in Virtual Reality in Unity.
- | Designed the 3D model of the Lincoln Cemetery using Blender.
- | Created the 3D model of the Lincoln Cemetery using drone photogrammetry and PIX4D.
- | Presented the project at the Annual Meeting of the Pennsylvania Geographical Society at Villanova University: 11/2018.

## PUBLICATIONS

---

- 📄 Get more for less: Principled Data Selection for Warming Up Fine-Tuning in LLMs  
Feiyang Kang, Hoang Anh Just, Yifan Sun, Himanshu Jahagirdar, Yuanzhi Zhang, Rongxing Du, Anit Sahu, Ruoxi Jia  
International Conference on Learning Representations (ICLR), 2024

- ¶ Performance Scaling via Optimal Transport: Enabling Data Selection from Partially Revealed Source  
Feiyang Kang\*, Hoang Anh Just\*, Anit Sahu, Ruoxi Jia  
Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS), 2023
- ¶ NARCISSUS: A Practical Clean-Label Backdoor Attack with Limited Information  
Yi Zeng\*, Minzhou Pan\*, Hoang Anh Just, Lingjuan Lyu, Meikang Qiu and Ruoxi Jia  
ACM Conference on Computer and Communications Security (CCS), 2023
- ¶ PRIVMON: Real-Time Platform-Agnostic Privacy Leakage Detection for Machine Learning Models  
Myeongseob Ko, Xinyu Yang, Zhengjie Ji, Hoang Anh Just, Peng Gao, Ruoxi Jia  
International Symposium on Research in Attacks, Intrusions and Defenses (RAID), 2023
- ¶ 2D-Shapley: A Framework for Fragmented Data Valuation  
Liu Zhihong\*, Hoang Anh Just\*, Xiangyu Chang, Xi Chen, Ruoxi Jia  
International Conference on Machine Learning (ICML), 2023
- ¶ LAVA: Data Valuation without Pre-Specified Learning Algorithms  
Hoang Anh Just\*, Feiyang Kang\*, Tianhao Wang, Yi Zeng, Myeongseob Ko, Ming Jin and Ruoxi Jia  
International Conference on Learning Representations (ICLR), 2023 (Spotlight)
- ¶ ModelPred: A Framework for Predicting Trained Model from Training Data  
Yingyan Zeng, Tianhao Wang, Si Chen, Hoang Anh Just, Ran Jin, Ruoxi Jia  
IEEE Conference on Secure and Trustworthy Machine Learning (SatML), 2023
- ¶ Label-Only Model Inversion Attacks via Boundary Repulsion  
Mostafa Kahla, Si Chen, Hoang Anh Just, Ruoxi Jia  
IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022
- ¶ On perfect bases in finite Abelian groups  
Béla Bajnok, Connor Berson and Hoang Anh Just  
Involve, a Journal of Mathematics, 12/2022
- ¶ Opponent Hand Estimation in the Game of Gin Rummy  
Peter Francis\*, Hoang Anh Just\*, Todd Neller  
AAAI Conference on Artificial Intelligence (AAAI), 2021

SERVICE

---

Reviewer

CVPR 2023, NeurIPS 2023, ENLSP 2023, ICLR 2023, ICML 2024

ACADEMIC HONOR AND AWARD

---

Virginia Tech College of Engineering Fellowship

Spring 2023

Pi Mu Epsilon Mathematics Society Membership

Spring 2021-Present

Phi Beta Kappa Honor Society Membership

Fall 2020-Present

J. Roger Musselman Award, Provost, Gettysburg College

Fall 2020

Paul Mugabi Problem Solving Award, Department of Mathematics, Gettysburg College

Fall 2019, Fall 2020