# Vedant Nanda

✉ nvedant07@gmail.com | ⌂ nvedant07.github.io | ⬛ nvedant07 | 🐦 @_nvedant_

## Work Experience

**Aleph Alpha** *Heidelberg, DE*
AI Researcher *June 2024 - Present*
**Fast and controllable LLM inference:** We build our custom LLM serving platform that achieves high overall throughput using efficient KV Cache management (eg: PagedAttention) and provides novel inference-time capabilities such as steering and explainability.

**Amazon AWS** *WA, USA*
Applied Science Intern *June 2023 - August 2023*
**AWS Bedrock:** I built an inference time de-biasing algorithm for text-to-image diffusion models.

**Amazon AWS** *Cambridge, UK*
Applied Science Intern *November 2022 - January 2023*
**AWS Clarify:** I built an actionable causal explainability method using manifold constraints that outperformed existing methods for the feasibility of acting upon the explanations and could be applied to a broad family of ML models on tabular data.

## Education

**University of Maryland, College Park** *MD, USA & SB, DE*
Ph.D. in Computer Science *August 2019 - May 2024*
- Research Interests: Trustworthy Deep Learning: Fairness, Robustness, and Efficiency
- Part of Maryland-Max Planck joint program through which I spent time at MPI-SWS in Saarbrücken, DE.
- Advisors: Krishna P. Gummadi (MPI-SWS) and John P. Dickerson (University of Maryland)

**University of Maryland, College Park** *College Park, MD, USA*
M.S. in Computer Science *August 2019 - May 2022*
- Advisor: John P. Dickerson
- TA for CMSC 320 (Intro to Data Science), Fall 2019

**Indraprastha Institute of Information Technology (IIIT) Delhi** *New Delhi, India*
B.Tech. in Computer Science and Engineering *August 2015 - May 2019*
- **GPA: 9.47/10**, in top 5% of institute
- Part of dean's list for academic excellence for all years
- Selected Coursework: Numerical Methods, Calculus-I, Calculus-II, Machine Learning, Collaborative Filtering, Information Retrieval, Designing Human-Centered Systems, System Administration, Network Administration

## Publications

**The Impact of Inference Acceleration Strategies on Bias of Large Language Models** *NAACL*
Elisabeth Kirsten, Ivan Habernal, **Vedant Nanda**, Muhammad Bilal Zafar *2025*
An earlier version appeared at SafeGenAI workshop at NeurIPS 2024

**Lawma: The Power of Specialization for Legal Tasks** *ICLR*
Ricardo Dominguez-Olmedo, **Vedant Nanda**, Rediet Abebe, Stefan Bechtold, Christoph Engel, Jens
Frankenreiter, Krishna Gummadi, Moritz Hardt, Michael Livermore *2025*
An earlier version appeared at Symposium on CS&Law

**Towards Reliable Latent Knowledge Estimation in LLMs: Zero-Prompt Many-Shot Based Factual Knowledge Extraction** *WSDM*
Qinyuan Wu, Mohammad Aflah Khan, Soumi Das, **Vedant Nanda**, Bishwamittra Ghosh, Camila Kolling, Till
Speicher, Laurent Bindschaedler, Krishna P. Gummadi, Evimaria Terzi *2025*
Code: https://github.com/QinyuanWu0710/ZeroPrompt_LKE, HF Dataset: https://huggingface.co/datasets/QinyuanWu/T-Rex-MC

**Understanding the Role of Invariance in Transfer Learning** *TMLR*
Till Speicher, **Vedant Nanda**, Krishna P. Gummadi *2024*
Code: github.com/tillspeicher/representation-invariance-transfer

**Diffused Redundancy in Pre-trained Representations**

*NeurIPS*

VEDANT NANDA, TILL SPEICHER, JOHN P. DICKERSON, KRISHNA P. GUMMADI, SOHEIL FEIZI, ADRIAN WELLER

*2023*

CODE: GITHUB.COM/NVEDANT07/DIFFUSED-REDUNDANCY

**What Happens During Finetuning of Vision Transformers: An Invariance Based Investigation**

*Conference on Lifelong Learning Agents (CoLLAs)*

GABRIELE MERLIN, VEDANT NANDA, RUCHIT RAWAL, MARIYA TONEVA

*2023*

**Do Invariances in Deep Neural Networks Align with Human Perception?**

*AAAI (Oral)*

VEDANT NANDA, AYAN MAJUMDAR, CAMILA KOLLING, JOHN P. DICKERSON, KRISHNA P. GUMMADI, BRADLEY C. LOVE, ADRIAN WELLER

*2023*

CODE: GITHUB.COM/NVEDANT07/HUMAN-NN-ALIGNMENT

**Rawlsian Fairness in Online Bipartite Matching: Two-sided, Group, and Individual**

*AAAI*

SEYED A. ESMAEILI, SHARMILA DUPPALA, DAVIDSON CHENG, VEDANT NANDA, ARAVIND SRINIVASAN, JOHN P. DICKERSON

*2023*

EARLIER VERSION APPEARED AS EXTENDED ABSTRACT AT AAMAS 2022

**Measuring Representational Robustness of Neural Networks Through Shared Invariances**

*ICML (Long Oral)*

VEDANT NANDA, TILL SPEICHER, CAMILA KOLLING, JOHN P. DICKERSON, KRISHNA P. GUMMADI, ADRIAN WELLER

*2022*

CODE: GITHUB.COM/NVEDANT07/STIR

**Fairness Through Robustness: Investigating Robustness Disparity in Deep Learning**

*FAccT*

VEDANT NANDA*, SAMUEL DOOLEY*, SAHIL SINGLA, SOHEIL FEIZI, JOHN P. DICKERSON

*2021*

* EQUAL CONTRIBUTION; CODE: GITHUB.COM/NVEDANT07/FAIRNESS-THROUGH-ROBUSTNESS

**Balancing the Tradeoff between Profit and Fairness in Rideshare Platforms during High-Demand Hours**

*AAAI*

VEDANT NANDA, PAN XU, KARTHIK A. SANKARARAMAN, JOHN P. DICKERSON, ARAVIND SRINIVASAN

*2020*

ALSO PRESENTED AT AIES 2020 (ORAL); CODE: GITHUB.COM/NVEDANT07/RIDESHARE-FAIRNESS-PEAK

**On the Long-term Impact of Algorithmic Decision Policies: Effort Unfairness and Feature Segregation through Social Learning**

*ICML*

HODA HEIDARI *, VEDANT NANDA *, KRISHNA P. GUMMADI

*2019*

* EQUAL CONTRIBUTION; CODE: GITHUB.COM/NVEDANT07/EFFORT_REWARD_FAIRNESS

**Leveraging Facebook's Free Basics Engine for Web Service Deployment in Developing Regions**

*ICTD*

SIDDHARTH SINGH*, VEDANT NANDA*, RIJUREKHA SEN, SATADAL SENGUPTA, PONNURANGAM KUMARAGURU, KRISHNA P. GUMMADI

*2017*

* EQUAL CONTRIBUTION

# Workshops and Posters

**Learning to Explain Machine Learning**

*CHI workshop on Human-Centered Explainable AI*

VEDANT NANDA*, DUNCAN MCELFRESH*, JOHN P. DICKERSON

*2021*

* EQUAL CONTRIBUTION

**Technical Challenges for Training Fair Neural Networks**

*ICLR workshop on Responsible AI*

VALERIIA CHEREPANOVA*, VEDANT NANDA*, MICAH GOLDBLUM, JOHN P. DICKERSON, TOM GOLDSTEIN

*2021*

* EQUAL CONTRIBUTION

**Unifying Model Explainability and Robustness via Reasoning Labels**

*NeurIPS workshop on Safety and Robustness in Decision Making*

VEDANT NANDA, JUNAID ALI, KRISHNA P. GUMMADI, MUHAMMAD BILAL ZAFAR

*2019*

**Stop the KillFies! Using Deep Learning Models to Identify Dangerous Selfies**

*WWW workshop on Modelling Social Media*

VEDANT NANDA, H.LAMBA, D.AGARWAL, M.ARORA, N.SACHDEVA, P.KUMARAGURU

*2018*

**Empirical Analysis of Facebook's Free Basics**                                    *SIGMETRICS (poster)*

S.Singh*, **Vedant Nanda***, R.Sen, S.Ahmad, S.Sengupta, A.Phokeer, Z.A.Farooq, T.A.Khan, P.Kumaraguru, I.A.Qazi,    *2017*
D.Choffnes, K.P.Gummadi

* Equal Contribution

## Other Experience

**University of Maryland, College Park**                                                      *MD, USA*

Teaching Assistant, CMSC320: Intro to Data Science                                    *Aug 2019 - Dec 2019*

Advisor: John P. Dickerson

**Max Planck Institute for Software Systems**                                       *Saarbrücken, Germany*

Research Intern                                                                       *May 2018 - Aug 2018*

Advisor: Krishna P. Gummadi

**Precog, IIITD**                                                                       *New Delhi, India*

Research Intern                                                                       *May 2017 - Aug 2017*

Advisor: Ponnurangam Kumaraguru

## Presentations/Talks

| | |
|---|---|
| 2024 | Talk @ Sysnets MPI-SWS on Efficient and Controllable LLM Inference. Hosted by Laurent Bindschaedler. |
| 2024 | Talk @ AI Fest on Efficient and Controllable LLM Inference. Hosted by Arthur AI. |
| 2024 | Talk @ Ruhr Universität Bochum on Efficient and Controllable LLM Inference. Hosted by Muhammad Bilal Zafar. |
| 2024 | Talk @ ETH Center for Law and Economics on Specializing LLMs for Legal Tasks. Hosted by Stefan Bechtold. |
| 2024 | Thesis Defense at University of Maryland. |
| 2024 | Talk @ Huawei Research "Towards Foundations of Trustworthy Deep Learning: Fairness, Robustness and Efficiency". |
| 2024 | Talk @ Bosch Center for AI "Towards Foundations of Trustworthy Deep Learning: Fairness, Robustness and Efficiency". |
| 2023 | Thesis Proposal at University of Maryland. |
| 2022 | Talk at University of Cambridge Machine Learning Group. Hosted by Adrian Weller. |
| 2022 | Oral Talk at International Conference on Machine Learning (ICML), Baltimore, Maryland. |
| 2022 | Talk at Computer Vision and Machine Learning seminar @ MPI-INF, virtual. |
| 2022 | Talk at ML Tea @ MPI-SWS, virtual. |
| 2021 | Talk at UMD Fairness in AI Seminar, *joint with Valeriia Cherepanova*, virtual. Link. |
| 2021 | Paper QnA at Conference on Fairness Accountability and Transparency (FAccT), virtual. Link. |
| 2020 | Oral talk at Conference on AI, Ethics and Society (AIES), NYC, USA |

## Service

| | |
|---|---|
| Reviewer | ICML 2021, 2023, 2024, 2025 |
| | ICLR 2023, 2025 |
| | NeurIPS 2021 |
| | AAAI 2021 |
| | CVPR 2021 |
| | ICCV 2021 |
| | WWW 2020, 2021 |
| | ASONAM 2019 |
| Other | UMD Graduate Admission Reviewer 2020 |
| | ELLIS PhD Admission Reviewer 2023 |

## Skills

| | |
|---|---|
| ML | PyTorch, Lightning/LitGPT, Transformers, Accelerate, Numpy, Triton/CUDA, DeepSpeed |
| Other | Matplotlib, Pandas, Git, C/C++ |

## References

1. **Prof. Krishna P. Gummadi**
   SCIENTIFIC DIRECTOR
   MAX PLANCK INSTITUTE FOR SOFTWARE SYSTEMS

2. **Prof. John P. Dickerson**
   ASSOCIATE PROFESSOR, COMPUTER SCIENCE
   UNIVERSITY OF MARYLAND, COLLEGE PARK