# Vedant Palit

(+91) 9163389534 | vedantpalit@kgpian.iitkgp.ac.in | vedantpalit.github.io | github.com/vedantpalit | google scholar

## EDUCATION

**Indian Institute of Technology(IIT) Kharagpur, India**　　　　　　　　　　　　　**2021 − 2026**

*Btech in Industrial & Systems Engineering and Mtech in Financial Engineering*

**Relevant Coursework:** Big Data Analysis, Scalable Data Mining, Regression and Time Series Models, Statistical Learning, Transform Calculus, Operations Research, Programming and Data Structures, Linear Algebra - Numerical and Complex Analysis

## PUBLICATIONS

- Towards Vision-Language Mechanistic Interpretability: A Causal Tracing Tool for BLIP - **Vedant Palit\***, Rohan Pandey\*, Aryaman Arora, Paul Pu Liang **(*ICCV 2023 CLVL Workshop, First Author*)**.

- Forgotten Polygons: Multimodal Large Language Models are Shape-Blind - William Rudman\*, Michal Golovanevsky\*, Amir Bar, **Vedant Palit**, Yann LeCun, Carsten Eickhoff, Ritambhara Singh **(*ACL 2025 Findings Track*)**.

- What Do VLMs NOTICE? A Mechanistic Interpretability Pipeline for Gaussian-Noise-free Text-Image Corruption and Evaluation - Michal Golovanevsky\*, William Rudman\*, **Vedant Palit**, Ritambhara Singh, Carsten Eickhoff **(*NAACL 2025*)**.

- WellDunn: On the Robustness and Explainability of Language Models and Large Language Models in Identifying Wellness Dimensions - Seyedali Mohammadi, Edward Raff, Jinendra Malekar, **Vedant Palit**, Francis Ferraro, Manas Gaur **(*EMNLP 2024 Blackbox NLP Workshop*)**.

- Knowledge Graph Guided Semantic Evaluation of Language Models For User Trust - Kaushik Roy, Tarun Garg, **Vedant Palit** **(*IEEE CAI 2023*)**.

- Exploring The Potential of Large Language Models for Assisting with Mental Health Diagnostic Assessments - Kaushik Roy, Harshul Surana, Darssan Eswaramoorthi, Yuxin Zi, **Vedant Palit**, Ritvik Garimella, Amit Sheth **(*ACM Transactions on Computing for Healthcare*)**.

## PREPRINTS

- Adaptive Federated Learning Defences via Trust-Aware Deep Q-Networks - **Vedant Palit** (*Under Review at **ICLR 2026**, Bachelor Thesis Project*).

## RESEARCH EXPERIENCE

**Research Collaborator | Brown University**　　　　　　　　　　　　　**Providence, RI (Remote)**

*Advised by Michal Golovanevsky, William Rudman and Dr Carsten Eickhoff*　　　　　*Feb 2024 – Present*

- Analysed the effectiveness of a novel semantic minimal pair and symmetric text replacement corruption scheme, against Gaussian noise-based causal mediation analysis through activation patching and knockouts.
- Evaluated the multi-step math reasoning capability for the LLaVA-1.5, LLaVA-Next, LLaVA-OneVision, Qwen2-VL, Molmo, InternVL, LLaMA-3.2 and Math-PUMA vision-language models.

**Research Intern | Inria, University of Lille**　　　　　　　　　　　　　**Lille, France (Remote)**

*Advised by Dr Debobrata Basu*　　　　　　　　　　　　　　　　　　*Dec 2023 – Present*

- Proposed and implemented a standardised gymnasium environment for the evaluation of algorithmic fairness in the approval of loans through credit worthiness.
- Benchmarked standard RL algorithms such as Linear-Q, Discrete-Q, Proximal Policy Gradient and performative RL algorithms across social welfare, rawlsian, fairness aware and utilitarian reward functions.

**Research Collaborator | Carnegie Mellon University**　　　　　　　　　　**Pittsburgh, PA (Remote)**

*Advised by Rohan Pandey, Aryaman Arora and Dr Paul Pu Liang*　　　　　　　　*April 2023 – Sep 2023*

- Proposed and developed a pipeline adapting causal mediation analysis utilising Gaussian Noise injection and activation patching to interpret blackbox architectures of VL transformers.

- Demonstrated the methodology on BLIP and used the COCO-QA dataset to study the internal layer effect of the multimodal text-encoder on model outputs as well as analysed the variation of patching effects with the strength of the noise.

### Research Intern | University of Maryland, Baltimore
*Advised by Dr Manas Gaur*

Baltimore, MD (Remote)

*Nov 2022 – Sep 2024*

- Trained various general and domain-specific models for suicide risk assessment, using the cross entropy and gamblers loss functions on annotated datasets containing social media posts classified into 6 different wellness dimensions.
- Utilised singular value decomposition to analyse the impact of the loss function on the attention scores of the models.

### Research Collaborator | University of South Carolina
*Advised by Dr Kaushik Roy and Dr Amit Sheth*

Columbia, SC (Remote)

*Feb 2023 – May 2023*

- Developed a novel evaluation method to measure error in reconstruction of masked knowledge graph structures from outputs by LLMs.
- Devised a novel paradigm of generating independent and closed context question-answer pairs to improve retrieval capability of vanilla RAG.

## INDUSTRY EXPERIENCE

### Machine Learning Intern | JP Morgan Chase
*Model Risk Governance and Review Division*

Bangalore, India

*May 2025 – July 2025*

- Investigated the internal mechanisms of a LLaMA-3.1 model fine-tuned on confidential JP Morgan data, utilising causal mediation analysis and activation patching.
- Implemented a rank-one editing pipeline to benchmark the performance improvements of layer-wise edits on model decisions.
- Reviewed quantitative models through backtesting and classical interpretability techniques like SHAP and Partial Dependence Plots.

## RESEARCH ALIGNED ACTIVITIES

**Conference Reviewing**: Reviewed for ACM SIGKDD Workshop 2024, ICLR 2025, NeurIPS Mechanistic Interpretability Workshop 2025, ICLR 2026
**Mentoring**: Mentored over 25 students regarding research, career and academics.
**Technical Writing**: Writer of a series of blogs reviewing papers on ML, DL and AI. [Medium]
**Research Communities**: Member of ACM (Association for Computing Machinery)

## TECHNICAL SKILLS

**Programming Languages**: C/C++, Python, MATLAB
**CV-NLP**: Transformers, OpenCV, PIL, Llama-Index

**ML-DL**: TensorFlow, Pytorch, Torchvision, Sklearn, Caffe
**Miscellaneous**: Mysql, LaTeX, HTML, Markdown, Git

## AWARDS AND ACHIEVEMENTS

**Adobe AI Challenge**: Captained the team of IIT KGP to Gold among 23 competing IITs at Inter IIT Tech Meet 2024.
**JEE Advanced**: Placed in the top 0.5% nationally among candidates appearing in JEE Advanced, 2021.
**JEE Mains**: Placed in the top 0.8% nationally among candidates appearing in JEE MAIN 2021.
**WBJEE**: Placed in the top 0.1% in the state among candidates appearing in WBJEE 2021
**Scientific Forum**: Selected as a delegate to represent India at the Asia Pacific Forum for Science Talented 2019.

## EXTRACURRICULARS

**NSS Volunteer**: Recipient of the gold medal for exceptional service work as an active participant in cleanliness drives, clothes distribution drives and education camps conducted by the NSS in villages near Kharagpur.
**Blogs on history**: Writer of a series of blogs discussing events of history and prehistory and their implications on the world and ecosystem. [Blog] **Quizzing**: An Active member of the IIT Kharagpur Quiz Club.