

# Vedant Palit

(+91) 9163389534 | [vedantpalit@kgpian.iitkgp.ac.in](mailto:vedantpalit@kgpian.iitkgp.ac.in) | [vedantpalit.github.io](https://vedantpalit.github.io) | [github.com/vedantpalit](https://github.com/vedantpalit) | [google scholar](https://scholar.google.com/citations?user=vedantpalit)

## EDUCATION

<b>Indian Institute of Technology(IIT) Kharagpur, India</b> <i>Btech in Industrial &amp; Systems Engineering and Mtech in Financial Engineering</i>	2021 – 2026
<b>Birla High School</b> <i>Subjects: English, Physics, Chemistry, Mathematics, Computer Science</i>	2008 – 2021

## PUBLICATIONS

<b>Towards Vision-Language Mechanistic Interpretability: A Causal Tracing Tool for BLIP</b> [Link] <i>Published at the ICCV 2023 Workshop on Closing The Loop Between Vision and Language - <b>First Author</b></i>	
<b>Knowledge Graph Guided Semantic Evaluation of Language Models For User Trust</b> [Link] <i>Published and Presented at the IEEE Conference on Artificial Intelligence 2023</i>	
<b>Benchmark of Wellness Dimensions for Robustness and Explainability Evaluation in LMs</b> <i>Accepted at the ACL ARR 2024 and Under Review at EMNLP 2024</i>	
<b>A Mechanistic Interpretability Pipeline for Noise-free Text-Image Corruption and Evaluation</b> <i>Under Review at EMNLP 2024</i>	
<b>Leveraging QA-based Chunk Formatting for Improving Retriever Quality during RAG</b> <i>Under Review at Machine Learning and Knowledge Extraction Journal 2024</i>	

## RESEARCH EXPERIENCE

<b>Mechanistic Interpretability of Vision-Language Models</b> <i>As a part of the Eickhoff AI Lab, Brown University</i>	Providence, RI (Remote) Feb 2024 – June 2024
<ul style="list-style-type: none"><li>Developed and created an extensive pipeline for an in-depth mechanistic interpretability study of the BLIP VL model through path patching and knockouts.</li><li>Introduced the novel semantic minimal pair and symmetric text replacement corruption scheme demonstrating more reliable results from causal mediation analysis over the pre-existing gaussian noise corruption scheme.</li></ul>	
<b>QA-based Chunk Formatting for Retrieval Improvement</b> <i>In Collaboration with Kaushik Roy Phd, University of South Carolina</i>	Columbia, SC (Remote) March 2024 – June 2024
<ul style="list-style-type: none"><li>Implemented and benchmarked Vanilla and Sentence-Window RAG on the MultiHopRAG dataset, using multiple open-source models such as Llama2, OrcaMini-3B and evaluation metrics such as BLUE, ROUGE-L and NUBIA.</li><li>Devised a novel paradigm of generating independent and closed context question-answer pairs to improve retrieval capability of vanilla RAG.</li></ul>	
<b>Causal Intervention on the BLIP Architecture</b> <i>In Collaboration with Rohan Pandey, Carnegie Mellon University</i>	Pittsburgh, PA (Remote) April 2023 – Sep 2023
<ul style="list-style-type: none"><li>Created a pipeline adapting causal mediation analysis to interpret blackbox architectures of VL transformers.</li><li>Implemented the method on the BLIP transformer and used the COCO-QA dataset to study the effect of various layers on the final outputs.</li></ul>	
<b>Wellness Dimensions Benchmark for Explainability of LMs</b> <i>Under the guidance of Prof Manas Gaur, University of Maryland, Baltimore</i>	Baltimore, MD (Remote) Nov 2022 – Dec 2023
<ul style="list-style-type: none"><li>Trained various general and domain-specific models for suicide risk assessment, using the gamblers and cross entropy loss functions on annotated datasets containing social media posts classified into 6 different wellness dimensions.</li><li>Utilised singular value decomposition to analyse the impact of the loss function on the attention scores of the models.</li></ul>	
<b>Knowledge Graph Guided Semantic Evaluation of LMs For User Trust</b> <i>In Collaboration with Kaushik Roy Phd, University of South Carolina</i>	Columbia, SC (Remote) Feb 2023 – March 2023
<ul style="list-style-type: none"><li>Developed a novel evaluation method to measure error in reconstruction of masked knowledge graph structures from outputs by LLMs.</li><li>Analysed and benchmarked the performance of GPT-3.5, GPT-J and GPT-NeoX in reconstructing KG paths using the evaluation metric.</li></ul>	

## RELEVANT COURSEWORK

---

- **University:** Regression and Time Series Models(MA60280), Machine Learning(AI41002) Transform Calculus(MA20202), Operations Research-I(IM21201), Programming and Data Structures(CS10003), Linear Algebra-Numerical and Complex Analysis(MA11004)
- **MOOCs:** Natural Language Processing with Deep Learning(CS224N), Introduction to Algorithms(MIT 6.006), Neural Networks and Deep Learning, Machine Learning

## TECHNICAL SKILLS

---

**Programming Languages:** C/C++, Python, MATLAB      **ML-DL:** TensorFlow, Pytorch, Torchvision, Sklearn, Caffe  
**CV-NLP:** Transformers, OpenCV, PIL, Llama-Index      **Miscellaneous:** Mysql, LaTeX, HTML, Markdown, Git

## AWARDS AND ACHIEVEMENTS

---

**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 out of 1000+ candidates to represent India at the Asia Pacific Forum for Science Talented 2019.

**Case Study:** Stood 1st amongst 5000+ participants in the BITS APOGEE, CaseQuesta challenge 2022.

## EXTRACURRICULARS

---

**Technical Writing:** Writer of a series of blogs reviewing papers on ML, DL and AI. [Medium]

**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.