

# Saransh Gupta

[Email](#)[LinkedIn](#)[Github](#)[Website](#)[+91-9530277421](#)

## ACADEMIC PROFILE

Year	Institution	Degree	CGPA
2022	Indian Institute of Technology Kharagpur	B. Tech. + M. Tech. (Engineering Product Design)	8.09 / 10.00

## PUBLICATIONS

- S. Gupta *et al.*, "Integrative Network Modeling Highlights the Crucial Roles of Rho-GDI Signaling Pathway in the Progression of Non-Small Cell Lung Cancer," in *IEEE - JBHI*, 2022, doi: **10.1109/JBHI.2022.3190038**
- Entity-aware Question-Answer Extraction for Shopping Guidance, **Amazon Machine Learning Conference - 2022**
- (Gupta et al.) Development of a virtual reality-based fire training simulator and machine learning-based path guidance system (working paper), **IHIET-AI, 2020**, Centre Hospitalier Universitaire Vaudois, **Lausanne, Switzerland**

## INTERNSHIPS AND PROJECTS

Amazon Development Center (India)   Applied Scientist - Intern		Jan'22 – June '22
<b>Project - 1:</b> Build a demo tool to help in the navigation and exploration of the Pre-curated Question Bank (PCQB)		
▪ Created a dashboard using <b>streamlit</b> enabling a user to input their query and get relevant questions accordingly		
▪ Integrated the frontend with the backend and a <b>BERT</b> based model to fetch relevant questions based on queries input		
▪ Demonstrated the coverage of PCQB with respect to user queries using the query-question relevance feature		
<b>Project - 2:</b> Generate Pre-curated Question Bank (PCQB) Question and Answer extraction from articles		
▪ Developed a <b>Transformers</b> based two-step model for the Question Generation followed by the answer extraction		
▪ Scrapped Texts, <b>People Also Ask (PAA)</b> questions and answers using certain queries related to E-Commerce domain		
▪ Increased the size of training dataset by <b>20</b> times by paraphrasing the dataset using T5 Text to Text Generator model		
▪ Achieved a <b>Perplexity score</b> of <b>82.3</b> on Question Generation by fine-tuning pre-trained <b>T5</b> model on the PAA dataset		
▪ Attained an <b>F-1 score</b> of <b>0.79</b> on the answer extraction task by fine-tuning encoders of <b>T5-large</b> model on PAA dataset		
▪ Deployed the two step model pipeline on the <b>streamlit</b> based demo web-application that accept user input as text		
<b>Tools and Software:</b> streamlit, Python, PyTorch, Transformers, BeautifulSoup, BERT, T5 (text to text generator)		
ZS Associates Inc.   Data Science Associate - Intern		Jan'21 – June '21
<b>Project - 1:</b> Extract biomedical text dataset, identify entities, and classify if there exists a relation between entities		
▪ Created a pipeline to extract texts from PubMed database, identifying the entities using <b>Selenium</b> and <b>PubTator</b>		
▪ Implemented <b>Binary Classification rules</b> , devised <b>four</b> labeling functions using bio-verbs, co-occurrence of entities		
▪ Generated a training dataset utilizing the four labeling functions in <b>Snorkel</b> by applying the <b>Weak Supervision</b>		
▪ Achieved F1 score of 0.88 on the gold-standard dataset in relation-classification by training <b>RoBERTa base</b> model		
<b>Project - 2:</b> Identify the type of relationship between two entities if it exists from the results of the Project-1		
▪ Created a new set of <b>three</b> labelling functions for <b>relation-type identification</b> by using the results of the project-1		
▪ Attained <b>F1 score</b> of <b>0.83</b> on the gold-standard dataset using <b>XGBoost</b> Model followed by feature engineering		
<b>Tools and Software:</b> Python, TensorFlow, Transfer Learning, Medline-Plus API, PubTator, Selenium, Snorkel		
Osaka University, Japan   Remote Research Assistant		Jan '20 – Dec '20
<b>Guide:</b> <a href="#">Dr. Kenji Mizuguchi</a> , <b>Mizuguchi Lab, Osaka University, Osaka, Japan</b>		
<b>Project:</b> Predict the Non-Small Cell Lung Cancer (NSCLC) using Machine Learning, identify its potential drug targets		
▪ Extracted <b>412</b> essential genes out of <b>10,077</b> by applying <b>Boruta</b> Feature selection on their gene expression dataset		
▪ Obtained <b>F-1 score</b> of <b>1.0</b> on validation and <b>0.98</b> on test dataset by using the <b>XGBoost</b> model to predict NSCLC		
▪ Predicted drug targets for the NSCLC by simulating a <b>Bayesian Network Model</b> on the Rho-GDI signaling pathway		
▪ Discovered methodology leads to an accurate treatment of the disease impacting <b>85%</b> of the lung cancer patients		
<b>Tools and Software:</b> Python, TargetMine, scikit-learn, smote, NetworkX, NumPy, pandas, Plotly, joblib		
ACHIEVEMENTS		
▪ Featured as one of the <b>Top 30 Undergraduate Achievers</b> of IIT Kharagpur in the UG Achievers Directory 2020		
▪ Conferred merit-based scholarship of <b>2200 €</b> by The A*Midex Foundation of <b>Aix-Marseille University, France</b>		
▪ Selected among <b>Top 5%</b> out of all for the summer fellowship at <b>The Institute of Science &amp; Technology Austria</b>		
▪ Got featured in the ISE Newsletter Autumn-2020 under the Department Spotlight of <b>ISE fights COVID-19, 2020</b>		
▪ Awarded as <b>Intern of The Month</b> for my contribution as a Data Analyst at Sapio Analytics by the CEO in July 2020		
COMPETITIONS / CONFERENCES		
▪ Annual Amazon Machine Learning Conference (AMLC) – Bengaluru, Karnataka		<b>[Aug 2022]</b>
▪ 23rd World Business Dialogue, Creation Lab at Evonik - Cologne, Germany		<b>[Jun 2022]</b>
▪ International Conference on Human Interaction & Emerging Technologies: Future Applications		<b>[Aug 2020]</b>
▪ Young Data Scientists annual meetup at Kaggle - days, Dubai World Trade Centre		<b>[Mar 2020]</b>
▪ Winner   Databuzz(Data Analytics Competition) DoMS IIT Madras		<b>[Jan 2020]</b>
Problem Statement: Prediction of the defaulters on lending credit cards to minimize loss incurred to the banks		