

# Saransh Gupta

[Contact](#)[Email](#)[LinkedIn](#)[Github](#)[Website](#)

## ACADEMIC PROFILE

Year	Institution	Degree/Examination	Percentage/CGPA
2022	Indian Institute of Technology Kharagpur	M.Tech. (Engineering Product Design)	9.26 / 10.00
2022	Indian Institute of Technology Kharagpur	B. Tech. (Engineering Product Design)	8.07 / 10

## PUBLICATIONS

- (Gupta et al.) An integrative machine learning and Bayesian modeling approach highlights the crucial roles of the Rho-GDI signaling pathway in the progression of non-small cell lung cancer (NSCLC); implications for drug target discovery, (under review), **IEEE Journal of Biomedical and Health Informatics (JBHI-01644-2021), 2021**
- (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

<b>Amazon Inc.   Applied Scientist - Intern</b>	<b>Jan'22 – June '22</b>
<b>Project - 1:</b> <i>Build a demo tool to help in the navigation and exploration of the Pre-curated Question Bank (PCQB)</i> <ul style="list-style-type: none"><li>▪ Created a dashboard using <b>streamlit</b> enabling a user to input their query and get relevant questions accordingly</li><li>▪ Integrated the frontend with the backend and a <b>BERT</b> based model to fetch relevant questions based on queries input</li><li>▪ Demonstrated the coverage of PCQB with respect to user queries using the query-question relevance feature</li><li>▪ Developed a navigation feature to allow searching and navigating the PCQB based on pre-determined tags</li></ul>	
<b>Project - 2:</b> <i>Pre-curated Question Bank (PCQB) Question and Answer extraction from articles (ongoing)</i>	
<b>ZS Associates Inc.   Data Science Associate - Intern</b>	<b>Jan'21 – June '21</b>
<b>Project - 1:</b> <i>Extract biomedical text dataset, identify entities, and classify if there exists a relation between entities</i> <ul style="list-style-type: none"><li>▪ Created a pipeline to extract texts from PubMed database, identifying the entities using <b>Selenium</b> and <b>PubTator</b></li><li>▪ Implemented <b>Binary Classification rules</b>, devised <b>four</b> labeling functions using bio-verbs, co-occurrence of entities</li><li>▪ Generated a training dataset utilizing the four labeling functions in <b>Snorkel</b> by applying the <b>Weak Supervision</b></li><li>▪ Achieved F1 score of 0.88 on the gold-standard dataset in relation-classification by training <b>RoBERTa base</b> model</li></ul>	
<b>Project - 2:</b> <i>Identify the type of relationship between two entities if it exists from the results of the Project-1</i> <ul style="list-style-type: none"><li>▪ Created a new set of <b>three</b> labelling functions for <b>relation-type identification</b> by using the results of the project-1</li><li>▪ Attained <b>F1 score</b> of <b>0.83</b> on the gold-standard dataset using <b>XGBoost</b> Model followed by feature engineering</li></ul>	
<b>Tools and Software:</b> Python, TensorFlow, Transfer Learning, Medline-Plus API, PubTator, Selenium, Snorkel	
<b>Osaka University, Japan   Remote Research Assistant</b>	<b>Jan '20 – Dec '20</b>
<b>Guide:</b> <a href="#">Dr. Kenji Mizuguchi</a> , <b>Mizuguchi Lab, Osaka University, Osaka, Japan</b>	
<b>Project:</b> <i>Predict the Non-Small Cell Lung Cancer (NSCLC) using Machine Learning, identify its potential drug targets</i> <ul style="list-style-type: none"><li>▪ Extracted <b>412</b> essential genes out of <b>10,077</b> by applying <b>Boruta</b> Feature selection on their gene expression dataset</li><li>▪ 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</li><li>▪ Predicted drug targets for the NSCLC by simulating a <b>Bayesian Network Model</b> on the Rho-GDI signaling pathway</li><li>▪ Discovered methodology leads to an accurate treatment of the disease impacting <b>85%</b> of the lung cancer patients</li></ul>	
<b>Tools and Software:</b> Python, TargetMine, scikit-learn, smote, NetworkX, NumPy, pandas, Plotly, joblib	
<b>ACHIEVEMENTS</b>	
<ul style="list-style-type: none"><li>▪ Featured as one of the <b>Top 30 Undergraduate Achievers</b> of IIT Kharagpur in the UG Achievers Directory 2020</li><li>▪ Conferred merit-based scholarship of <b>2200 €</b> by The A*Midex Foundation of <b>Aix-Marseille University, France</b></li><li>▪ Selected among <b>Top 5%</b> out of all for the summer fellowship at <b>The Institute of Science &amp; Technology Austria</b></li><li>▪ Got featured in the ISE Newsletter Autumn-2020 under the Department Spotlight of <b>ISE fights COVID-19, 2020</b></li><li>▪ Awarded as <b>Intern of The Month</b> for my contribution as a Data Analyst at Sapio Analytics by the CEO in July 2020</li></ul>	
<b>COMPETITIONS / CONFERENCES</b>	
<b>Young Data Scientists annual meetup at Kaggle - days, Dubai World Trade Centre</b>	<b>[Mar 2020]</b>
<b>International Conference on Human Interaction &amp; Emerging Technologies: Future Applications</b>	<b>[Aug 2020]</b>
<ul style="list-style-type: none"><li>- Virtually Presented the research paper on "Development of a virtual reality based fire training simulator and machine learning based path guidance system: A case of hospital fire breakouts" at Lausanne, Switzerland</li></ul>	
<b>Runner Up   Databuzz (Data Analytics Competition) DoMS IIT Madras</b>	<b>[Jan 2020]</b>
<b>Objective :</b> <i>Prediction of the defaulters on lending credit cards to minimize loss incurred to the banks (credit risk)</i> <ul style="list-style-type: none"><li>- Visualised key features governing target variable, removed outliers, imputed missing values</li><li>- Implemented undersampling for handling imbalanced classes, applied decision trees as baseline model</li></ul>	
Improved the validation AUC-ROC from 0.56 to 0.79 by applying the Bayesian optimization model	
<b>Gold Winner   Inter-Hall Data-Analytics competition   IIT Kharagpur</b>	<b>[Feb 2019]</b>
<ul style="list-style-type: none"><li>- Worked in the Visualization of data-set, making interpretation of trends followed by different variables</li><li>- Applied <b>XGBoost</b> followed by feature engineering, obtained <b>82.55%</b> on training, <b>80.66%</b> on test data</li><li>- Awarded best in the category of best F-1 Score, Overall stood second in terms of Business Case Study</li></ul>	