

# Rachit Sabharwal

📍 10211 Camden Garden Lane, Katy, Texas, 77494    ✉ rachit-sabharwal@outlook.com    ☎ (585) 281-1928  
🌐 www.rachitsabharwal.me    🌐 rachitest

## Technical Skills

**Machine Learning:** Scikit-learn, Pytorch, Tensorflow, Transfer Learning, LLMs (Gemini, Haiku, o1), Graph NNs (PyG), Reinforcement Learning (OpenAI Gymnasium)

**Programming Languages:** Python, R, Javascript

**Databases:** RDBMS (PostgreSQL, SQLite, MySQL), NoSQL DBMS (MongoDB, Elasticsearch, Neo4J), BigQuery

**Cloud and Distributed Computing:** AWS, GCP, Azure, On-Prem HPC

**DevOps:** Git, GitHub, GitLab, Docker, GitHub/GitLab CI/CD, Jenkins

**Workflow Orchestration:** Airflow, Prefect, Cron

## Relevant Experience

Sept 2022 - Oct 2024

**Biostatistics and Data Science - Graduate Research Assistant (Doctoral)**, The University of Texas Health Science Center at Houston, School of Public Health – Houston, TX

- Established DevOps practices, architected repository structures, and implemented CI/CD pipelines. Designed and implemented ETL pipelines for datasets of various sizes, achieving a **24x speed increase over prior solution leading to faster model iteration and deliverable generation**.
- Applied statistical models and machine learning techniques for data analysis on vaccine non-response and Covid-19 pediatric comorbidities. **Contributed to the publishing of three peer-reviewed articles, with two more under review.**
- Developed a weekly report generation pipeline with CI/CD, ensuring seamless integration of new data. Planned and enacted a comprehensive test suite, **increasing team lead confidence and reducing publishing turnaround times.**

May 2021 - Jan 2022

**Research and Early Development, Development Sciences & Informatics - Informatics Intern**, Genentech – San Francisco, CA

- Contributed to the development of frameworks for an in-house data annotation tool, utilizing advanced deep learning NLP models and custom processing pipelines. Created a efficient document tokenization framework **allowing for fast and high-quality data annotation for NLP model building** and downstream informatics systems.
- Engineered a neural network using deep transfer learning, to predict adverse drug events. Focusing on Drug-Induced Liver Injury created a **low parameter model that matched the prediction accuracy of the SOTA model.**
- Created a Knowledge Graph with Neo4j and a Graph Neural Network using NetworkX and PyTorch to generate gene expression signature-likes for drugs. Began the creation of a knowledge repository — an accessible data source which **reduced the need to manually sift through dense primary sources.**

Jan 2021 - May 2021

**Consumer & Market Knowledge - Advanced Analytics Co-Op**, Procter & Gamble – Cincinnati, OH

- Developed models to uncover consumer behaviors, market trends, and retailer patterns using click-stream data, **increasing corporate attention on DTC online retail.**
- Created, maintained, and optimized big data ETL pipelines to ensure seamless data integration and processing leading to a **1.5x increase in analyst operational efficiency.**
- Implemented advanced DevOps practices by integrating unit-testing, containerization, and agile methodologies to improve team efficiency and code quality. **Increased code coverage from 0% to 70%.**

June 2020 - Aug 2020

**Data Engineering Intern**, Bristol Myers Squibb – San Francisco, CA

- Developed a multifeatured patent recommendation web application. The application implemented a mix of classical NLP algorithms (TF-IDF, BM 25) and deep learning algorithms (BERT). **Improved bench scientists' productivity by streamlining patent search and analysis workflow.**

- Engineered and maintained performant ETL pipelines with Python and Apache Airflow. Undertook extensive data cleaning and data wrangling on datasets of varying sizes (small, medium, and large) **ensuring data quality and readiness for analysis at a high velocity.**
- Designed and maintained robust relational and graph databases in PostgreSQL and Neo4j, **optimizing data storage and retrieval to support various internal projects and analyses.**

## Education

June 2023 - present	<b>The University of California, Berkeley</b> , Certificate in Software Development and Programming – Berkeley, CA
Aug 2022 - present	<b>The University of Texas Health Science Center at Houston</b> , Doctor of Philosophy in Biostatistics – Houston, TX <ul style="list-style-type: none"> <li>• Advanced Certificate in Data Science</li> </ul>
Jan 2020 - May 2022	<b>The University of Texas Health Science Center at Houston</b> , Master of Science in Biostatistics – Houston, TX <ul style="list-style-type: none"> <li>• Thesis — BioRec: A Biomedical Recommendation System for Academic Conferences and Journals</li> <li>• Certificate in Data Science</li> </ul>
Sept 2014 - May 2018	<b>University of Rochester</b> , Bachelor of Science in Environmental Health – Rochester, NY <ul style="list-style-type: none"> <li>• Minor in Psychology</li> </ul>

## Publications

Aug 2024	<b>Factors associated with elevated SARS-CoV-2 immune response in children and adolescents</b> Messiah SE, Abbas R, Bergqvist E, Swartz MD, Talebi Y, <b>Sabharwal R</b> , Han H, Valerio-Shewmaker MA, DeSantis SM, Yaseen A, Gandhi HA, Amavisca XF, Ross JA, Padilla LN, Gonzalez MO, Wu L, Silberman MA, Lakey D, Shuford JA, Pont SJ, Boerwinkle E <a href="https://doi.org/10.3389/fped.2024.1393321">10.3389/fped.2024.1393321</a> (Frontiers in Pediatrics)
May 2024	<b>Baseline characteristics of SARS-CoV-2 vaccine non-responders in a large population-based sample</b> Yaseen A, DeSantis SM, <b>Sabharwal R</b> , Talebi Y, Swartz MD, Zhang S, Leon Novelo L, Pinzon-Gomez CL, Messiah SE, Valerio-Shewmaker M, Kohl HW 3rd, Ross J, Lakey D, Shuford JA, Pont SJ, Boerwinkle E <a href="https://doi.org/10.1371/journal.pone.0303420">10.1371/journal.pone.0303420</a> (PLoS One)
Apr 2024	<b>An Interactive Online Dashboard with COVID-19 Trends and Data Analysis in Northeast and South Texas</b> Zhang Z, <b>Sabharwal R</b> , Lee M, Zhang K, McGaha P, Crum M, Bauer C, Fisher-Hoch SP, McCormick JB, Reininger BM, Thomas S, Guajardo E, Pinon D, Yaseen A <a href="https://research.ebsco.com/linkprocessor/plink?id=894625e1-7146-30bf-aa2c-9f5637dac41e">research.ebsco.com/linkprocessor/plink?id=894625e1-7146-30bf-aa2c-9f5637dac41e</a> (Texas Public Health Journal)
Oct 2023	<b>Long-term immune response to SARS-CoV-2 infection and vaccination in children and adolescents</b> Messiah SE, Talebi Y, Swartz MD, <b>Sabharwal R</b> , Han H, Bergqvist E, Kohl HW 3rd, Valerio-Shewmaker M, DeSantis SM, Yaseen A, Kelder SH, Ross J, Padilla LN, Gonzalez MO, Wu L, Lakey D, Shuford JA, Pont SJ, Boerwinkle E <a href="https://doi.org/10.1038/s41390-023-02857-y">10.1038/s41390-023-02857-y</a> (Pediatric Research)
June 2023	<b>Scholarly recommendation systems: a literature survey</b> Zhang Z, Patra BG, Yaseen A, Zhu J, <b>Sabharwal R</b> , Roberts K, Cao T, Wu H <a href="https://doi.org/10.1007/s10115-023-01901-x">10.1007/s10115-023-01901-x</a> (Knowledge and Information Systems)
May 2021	<b>Data Cleaning for eCommerce: Standardizing Data Handling Practices for eCommerce Datasets</b> <b>Sabharwal R</b> Procter & Gamble, Internal White Paper