

# YASH KHARE

Boston, MA | (857) 506-4305 | [ykhare.y@northeastern.edu](mailto:ykhare.y@northeastern.edu) | [linkedin.com/in/yash-khare20](https://linkedin.com/in/yash-khare20) | [github.com/YashKhare20](https://github.com/YashKhare20)  
[yashkhare20.github.io/portfolio](http://yashkhare20.github.io/portfolio)

## Education

<b>Northeastern University</b> , GPA: 4.0 / 4.0	Jan. 2024 – Present
<i>Khoury College of Computer Sciences, Candidate for Master of Science in Data Science</i>	<i>Boston, MA</i>
<b>Relevant Coursework:</b> Supervised/Unsupervised Machine Learning, Natural Language Processing, Algorithms	
<b>MIT World Peace University</b> , GPA: 9.54 / 10.0	July 2017 – May 2021
<i>B.Tech in Computer Science &amp; Engineering</i>	<i>Pune, India</i>
<b>Relevant Coursework:</b> Data Structures & Algorithms, Database Management, Object Oriented Design	

## Experience

<b>L.L. Bean</b>	Jan. 2025 – Present
<i>Data Architect Intern</i>	<i>Freeport, Maine</i>
<ul style="list-style-type: none"><li>Architected a <b>GCS-to-BigQuery</b> ingestion pipeline using <b>Dataflow</b> and <b>Cloud Composer</b>, processing <b>10K+</b> daily files and enhancing raw data availability by <b>90%</b>.</li><li>Developed a <b>YAML-driven configuration framework</b> and implemented <b>partial record ingestion</b> with BigQuery audit tables, reducing pipeline setup time by <b>80%</b> and reruns by <b>45%</b>.</li><li>Designed a <b>reusable multi-agent framework</b> using Google's <b>Vertex AI Agent Development Kit (ADK)</b> for automated failure analysis across GCP services, achieving <b>99.9%</b> uptime with plug-and-play extensibility.</li><li>Implemented <b>source-controlled CI/CD pipelines</b> with <b>GitHub</b> and <b>Harness</b> for agent deployment, integrating <b>Splunk API</b> for real-time log extraction and automated <b>Microsoft Teams/email alerting</b> with comprehensive test suites.</li></ul>	
<b>Quantiphi Analytics Solution Pvt. Ltd.</b>	July 2021 - Nov. 2023
<i>Senior Data Engineer</i>	<i>Mumbai, India</i>
<ul style="list-style-type: none"><li>Integrated <b>SonarQube</b> &amp; <b>GitHub Actions</b> to automate quality checks in <b>Composer DAGs</b> and streamlined migration of <b>10K+</b> SKUs from <b>Snowflake</b> to BigQuery, boosting code security and data workflow productivity.</li><li>Directed the development of a <b>Kafka-PubSub</b> real-time <b>ETL</b> pipeline for a financial institute, integrating Cloud Composer, BigQuery, &amp; Dataflow to handle <b>2TB</b> daily, which resulted in a <b>40%</b> increase in operational efficiency.</li><li>Automated molecular enumeration of <b>1B</b> molecules using <b>Dataflow &amp; Workflows</b> for a healthcare innovator and subsequently implemented <b>CI/CD</b> pipelines for model training.</li><li>Orchestrated <b>MySQL &amp; MongoDB</b> to BigQuery migration for a drone racing league, utilizing <b>Datastream</b> and <b>Cloud Run</b>, streamlined with <b>Cloud Scheduler</b> and secured via <b>Secret Manager</b>, impacting <b>5K</b> drone pilots.</li></ul>	
<b>Maker's Lab, Tech Mahindra</b>	Nov. 2020 - May 2021
<i>Data Engineering Research Intern</i>	<i>Pune, India</i>
<ul style="list-style-type: none"><li>Developed an automatic image annotation system using <b>YOLOv4-Tiny</b>, processing data with <b>Apache Spark</b> on <b>HDFS</b>, achieving <b>40 mAP @.5</b> on the MS COCO dataset and reducing annotation time by <b>66%</b>.</li><li>Deployed the system using <b>Docker</b> and <b>Kubernetes</b>, optimizing CI/CD in <b>GitLab</b>, and cutting annotation costs by <b>89%</b>.</li></ul>	

## Projects

<b>HIMAS: Healthcare Intelligence Multi-Agent System for Federated Learning</b>	Sept. 2025 - Dec. 2025
<ul style="list-style-type: none"><li>Architected a <b>federated learning</b> pipeline using <b>Apache Airflow</b> to process <b>60K+</b> ICU records from <b>MIMIC-IV</b> across 3 hospitals, implementing <b>differential privacy</b> (<math>\epsilon=0.1</math>) and <b>k-anonymity</b> for HIPAA-compliant mortality prediction.</li><li>Trained a <b>TensorFlow</b> mortality prediction model using <b>Flower</b> framework with <b>FedAvg strategy</b> and <b>Bayesian optimization</b>, achieving <b>0.92 ROC-AUC</b>, and <b>74% recall</b> across <b>13K+</b> test samples with <b>11.6%</b> class imbalance.</li><li>Designed a <b>multi-agent system</b> using <b>Vertex AI ADK</b> with <b>A2A protocol</b> for cross-hospital coordination, and implemented <b>CI/CD</b> with <b>GitHub Actions</b>, <b>DVC</b> on <b>GCS</b>, achieving <b>zero data leakage</b> across data splits.</li></ul>	
<b>Tableau</b> , <b>Power BI</b> , <b>Matplotlib</b> , <b>Looker</b> , <b>Plotly</b> , <b>Seaborn</b>	
<b>BigQuery</b> , <b>Snowflake</b> , <b>MySQL</b> , <b>SQL Server</b> , <b>Hive</b> , <b>HBase</b> , <b>Informatica</b> , <b>Apache NiFi</b> , <b>Google Firebase</b>	
<b>Google Cloud Certified Associate Cloud Engineer</b> , <b>Professional Data Engineer</b>	

## Technical Skills

<b>Languages:</b> Python, SQL, Shell Script, R, Javascript, C#, HTML/CSS
<b>Frameworks &amp; Tools:</b> Apache Hadoop, Spark, Kafka, Kubernetes, Airflow, Docker, Git, Flask, DVC, Flower, Pytest
<b>ML/AI Platforms:</b> TensorFlow, Keras, Pytorch, Hugging Face, Google ADK, Scikit-learn, Pandas, Numpy
<b>Data Visualization:</b> Tableau, Power BI, Matplotlib, Looker, Plotly, Seaborn
<b>Databases/ETL:</b> BigQuery, Snowflake, MySQL, SQL Server, Hive, HBase, Informatica, Apache NiFi, Google Firebase
<b>Certifications:</b> <b>Google Cloud Certified Associate Cloud Engineer</b> , <b>Professional Data Engineer</b>