

# UTKARSH SINGH

usingh7@jh.edu | +1 (410) 805-4462 | LinkedIn | ORCID | GitHub

## EDUCATION

<b>Johns Hopkins University</b> , MSE Applied Mathematics Honors: Dean's Master's Fellowship   GPA: 3.30/4.00 Relevant Coursework: Computing for Applied Mathematics, Data Science, Econometrics, Time Series Analysis	August 2025 – Expected December 2026
<b>Maharashtra Institute of Technology</b> , BTech in Electronics and Communication Engineering   GPA: 3.99/4.00 Relevant Coursework: Machine Learning, Natural Language Processing, Optimization Techniques, Artificial Neural Networks, Pattern Recognition	July 2019 – July 2023

## TECHNICAL SKILLS

- Programming and Data:** Python, SQL, C++, R, MATLAB, Pandas, NumPy, PostgreSQL, MySQL
- Machine Learning Libraries:** Scikit-learn, PyTorch, TensorFlow, Keras, LightGBM, XGBoost
- Deep Learning and LLM frameworks:** Transformers, Hugging Face, LangChain, LangGraph, PydanticAI, Ollama, Faiss, Inspect AI
- ML Systems and MLOps:** Git, Docker, Kubernetes, OpenShift AI, MLflow, Grafana, Prometheus, Jupyter
- Platforms and Big Data:** Google Cloud, IBM Cloud, IBM WatsonX (AI/Data), Apache Spark, Apache Hive
- Data Processing and Visualization:** PyPDF, Docing, Streamlit, Matplotlib, Seaborn, MS Excel, Stata, Power BI
- Certifications:** IBM Machine Learning Specialist, IBM Developer Profession, Redbooks Gold Author

## EXPERIENCE

<b>Research Assistant</b>   Johns Hopkins Bloomberg School of Public Health	August 2025 – Present
• Designed class-imbalance mitigation strategy (cost-sensitive learning) for <b>99:1</b> -skewed EHR data; enhanced rare-event recall while controlling false positives.	
• Fitted tree-based classifiers for longitudinal prediction tasks, optimizing predictive performance and calibration under data sparsity and label imbalance.	
• Iteratively refined model stability and generalization by incorporating additional data sources and targeted feature and hyperparameter refinements.	
<b>Machine Learning Engineer</b>   IBM Systems Development Lab	July 2023 – July 2025
• Deployed autoencoder-based anomaly detection models (TransformerAE, LSTM) on multivariate time-series data across <b>32</b> KPI groups and <b>150+</b> features.	
• Built end-to-end data pipelines for telemetry ingestion, feature engineering, and anomaly scoring, persisting scored windows for downstream analysis.	
• Calibrated detection thresholds using sequence-aware statistical tuning methods to balance recall and false positives under production constraints.	
• Architected a real-time semantic triage engine using <b>BERT</b> and <b>Faiss indexing</b> to correlate live system anomalies with time-aligned logs and past tickets.	
• Applied clustering and re-ranking techniques to compress high-dimensional incident windows and surface top resolution candidates within <b>sub-60s</b> latency.	
• Created an instruction-tuned <b>LLM assistant</b> for Storage Insights using Granite/Llama with <b>QLoRA</b> -based alignment and light <b>RLHF</b> -style preference shaping.	
• Implemented structured query translation from natural-language inputs to monitoring and observability APIs, improving reliability of automated diagnostics.	
• Engineered an <b>open-source MCP</b> service exposing observability signals to agent-based workflows, reducing end-to-end issue identification latency by <b>15%</b> .	
<b>Software Engineer Intern</b>   IBM Systems Development Lab	January 2023 – July 2023
• Optimized large-scale ingestion pipelines for multi-tenant storage telemetry, processing <b>2+ TB/day</b> of Protobuf-based configuration data on IBM Cloud.	
• Reduced small-file overhead by <b>40%</b> through <b>dynamic batching</b> and <b>size-aware flush logic</b> , improving overall pipeline efficiency and space utilization.	
• Launched scalable event-driven consumers on Kubernetes to stream, filter, and normalize millions of system metrics per hour for downstream ML workloads.	
• Accelerated analytical queries by <b>2–3×</b> by optimizing schema layouts and ingestion paths for distributed SQL and Spark-based engines.	
• Orchestrated Prometheus metric exporters across <b>10</b> ML pipelines and automated Grafana dashboard provisioning to enable monitoring of performance.	
• Established standardized observability and performance baselines to detect bottlenecks and reliability issues across production data pipelines.	
<b>Data Scientist Intern</b>   Sisai Technologies	November 2021 – February 2022
• Developed temporal convolution-based sequence models to process noisy, sparse IoT time-series data, improving feature extraction for forecasting tasks.	
• Devised simulation-based stress tests spanning <b>100+</b> thermal-failure scenarios to evaluate model behavior under edge conditions and distribution shifts.	

## PROJECTS

<b>Chain-of-Thought Tampering Detection in Frontier LLMs</b>   Johns Hopkins	November 2025 – February 2026
• Evaluated tampering detection capabilities of <b>120B–235B parameter</b> reasoning models by architecting an <b>Inspect AI</b> -based evaluation pipeline to systematically perturb chain-of-thought traces (step deletion, cross-model substitution, semantic injection) across multi-domain reasoning benchmarks.	
• Quantified post-completion and in-generation detection performance using controlled baselines and nonparametric bootstrap confidence intervals, measuring sensitivity to <b>structured reasoning interventions</b> and robustness differentials between reasoning traces and final outputs.	
<b>Expert Specialization in Multilingual MoE Transformers</b>   Johns Hopkins	October 2025 – December 2025
• Constructed a sparse <b>Mixture-of-Experts</b> (MoE) Transformer for multilingual language modeling across 4 languages, implementing top-k routing and load-balancing loss in PyTorch to scale model capacity with <b>3x</b> higher parameter count at constant inference FLOPs.	
• Analyzed token-level expert routing behavior over <b>1M+</b> multilingual samples, quantifying expert specialization using routing entropy, utilization imbalance, and KL divergence, and identified language and script-specific expert preferences that reduced routing entropy by <b>18–25%</b> compared to random assignment.	
<b>Workload Placement Advisor</b>   IBM	March 2025 – May 2025
• Formulated and trained a time-series forecasting engine for block-level I/O demand, benchmarking classical statistical models against <b>PatchTST</b> and <b>TTM</b> across <b>100+</b> <b>FlashSystem arrays</b> , increasing out-of-sample prediction accuracy by <b>25%</b> .	
• Operationalized the decision optimization layer that translated probabilistic forecasts into placement and migration actions, explicitly modeling capacity, latency, and risk tradeoffs through hybrid compatibility scores (model predictions + configuration rules) to prevent nearly <b>10,000</b> potential <b>SLA violations</b> .	
<b>Ransomware Threat Detection for FlashSystem</b>   IBM	November 2024 – March 2025
• Trained <b>SnapML</b> -based ensemble classifiers on <b>500k+</b> FlashCore Module traces, detecting ransomware with < <b>1%</b> <b>false-positive rate</b> on live workloads.	
• Automated pipelines for class re-balancing, labeling, and <b>cross-family validation</b> to ensure the system generalized well to unseen ransomware variants.	
• Integrated the inference layer into Storage Virtualize stack to trigger early alerts, immutable snapshots, and forensic trace dashboards for support teams.	
<b>Ticket Deflection</b>   IBM	April 2024 – September 2024
• Developed a Gen-AI deflection workflow that sanitized <b>GDPR</b> -sensitive fields and clustered <b>26,000+</b> historical tickets into representative knowledge groups.	
• Generated auto-validated FAQs from these groups, reducing Level 2 support workload by nearly <b>30%</b> and cutting MTTR by <b>20%</b> by deflecting repetitive issues.	