

Jefferson Xu

me@jeffersonxu.com • (201) 788-0863 • <http://jeffersonxu.com>

WORK EXPERIENCE

Amazon, NY

Software Developer Engineer II (Personalization, Product Intent)

Apr 2025 - Present

- Scaled LLM inferencing system for backend ML infra by 5x handling > 1000 RPM with AWS Lambda + Step Functions
- Production-ized personalized recommendation system that supported 14M+ products in < 2 months for VP demo
- Increased LLM inferencing latency invocations by > 180% through addition of prompt caching
- Wired up AWS SageMaker endpoint to handle > 90 TPS through load testing with fallback DynamoDB workflow

Software Developer Engineer I (Personalization, Product Quality)

Sep 2023 - Mar 2025

- Produced ETL Spark data pipeline that ingested 10 shopping quality metrics (~1 TB daily) used to filter out low quality products leading to +\$111M CSales/yr and +\$29MM attributed OPS/yr across worldwide marketplaces
- Built internal auditing tool using React, AWS API Gateway, and Lambda to evaluate and manage 40k+ LLM annotations
- Generated image diversity score to deduplicate homepage product recommendations with > 90% image similarity
- Mentored intern with zero industry experience to return offer and served as main contact when onboarding 3 teammates
- Established organization wide metric to track product quality through third party collaboration across 3 teams
- Created an evaluation comparison table using Figma and React to seamlessly compare LLM performances

Software Developer Engineer Intern (Personalization, Product Graph)

Jun 2022 - Sep 2022

- Researched and clearly defined signals for model drift by running K-S statistic tests using Pandas, NumPy, and SQL
- Designed model drift detection workflow using AWS Step Functions, EMR Clusters, and Spark jobs
- Wrote PySpark jobs that processed petabytes of product catalog data that fetched model and statistical metrics
- Persisted drift signal metrics as JSON in S3 to be read from dashboard alerting when model retraining is needed

Amazon, WA

Software Developer Engineer Intern (Amazon Web Services, Kumo)

Jun 2021 - Sep 2021

- Developed REST API service connecting AWS support feedback data using API Gateway, DynamoDB, and Lambda
- Improved support ticket search rankings by adding metrics to existing API workflow using React.js and Java
- Wrote extensive integration and unit tests using Enzyme and JUnit test libraries reaching > 95% code coverage

Bizi, IL

Software Developer

Nov 2020 - May 2021

- Worked with team of 5 to develop web app displaying 50+ local businesses for socially conscious consumers
- Led front-end team converting Figma mockups to components using React.js, AWS Amplify, DynamoDB, and S3

System Software & Security Lab at Rutgers University, NJ

Undergraduate Research Assistant

Apr 2020 - Aug 2020

- Researched the usage of the Vapnik–Chervonenkis (VC) dimension as an evaluation metric for Neural Networks
- Ran polynomial regression tests using PyTorch and NumPy finding correlation between model accuracy and VC bound

BlackLapel, NY

Web Development Intern

Jun 2018 - May 2019

- Streamlined unit and integration tests using Selenium WebDriver and Node.js to allow for testing automation
- Refactored existing API endpoints using PHP, Neo4j, and Docker containers to be RESTful microservices

EDUCATION

Northwestern University

GPA: 3.7/4.0

Bachelor of Arts in Computer Science and Minor in Economics

Sep 2020 - Jun 2023

- **Relevant Courses:** Data Structures & Algorithms, Computer Systems, Intro to Networking, Software Construction, Human Computer Interface, Intro to AI, Intro to Data Science, Machine Learning, Scalable Software Architecture, Intro to Computational Linguistics
- **Awards:** Capital One Software Engineering Summit Participant (1/60 chosen from 500 applicants), Semifinalists for VentureCat (Northwestern's annual startup competition)

Bergen County Academies

GPA: 3.9/4.0

Academy for Technology and Computer Science

Sep 2015 - Jun 2019

SKILLS

Languages - Java | Python | Javascript | HTML | CSS

Technologies - Git | React.js | AWS | mySQL | Spark