# Junhui Li

junhui@cs.cmu.edu • 734-882-9191 • ophelialjh.github.io • linkedin.com/in/junhui-li-ljh

## EDUCATION

**Carnegie Mellon University**, School of Computer Science                               Pittsburgh, PA
M.S. in Artificial Intelligence and Innovation | GPA: 3.96/4.00                               May 2023
- Core Coursework: Machine Learning with Large Datasets, Deep Learning, Multimodal Machine Learning
- Instructional Aid, Machine Learning, Jan.-May 2022

**University of Michigan** (UM)                               Ann Arbor, MI
B.S.E. in Computer Science | GPA: 3.85/4.00                               May 2021
- Honors: Multidisciplinary Design Program Summer Fellowship, Dean's List, University Honors
- Core Coursework: Machine Learning, Computer Vision, Natural Language Processing, Operating System, Java

## SKILLS

**Programming Languages:** Python, C++, C, SQL, Java, Scala, Kotlin, MATLAB, Arduino, Assembly code
**Tools:** AWS, Microsoft Azure, Databricks, Jupyter Notebook, MongoDB, MySQL, Maven, Linux, Dataswarm, LaTex
**Frameworks:** PyTorch, Apache Spark, TensorFlow, Sklearn, NLTK, Pandas

## Work Experience

**Meta |** Applied Research Scientist Intern                               Menlo Park, CA | May 2022 – Aug. 2022
- Improved and tuned a cross-lingual language model (XLM) that saved $480M+ in the last two years and serves 99 internal teams by classifying Meta's daily spending into 400+ invoice categories with ~$70B annual traffic
- Designed and deployed 2-layer hierarchical model architecture to mitigate a severe class imbalance problem within a dataset of 10M invoices improving production accuracy from 92% to 98%; delivered tech talk to team
- Experimented with 10 useful new XLM features by creating Dataswarm pipelines to extract features from raw data and enhanced model f1-score by >2%
- Suggested the cross-functional team with modified category taxonomy based on analyzing confusion matrix
- Resolved data leakage by re-splitting train/test datasets so that production accuracy reflects model generality

**Intel |** Deep Learning Intern                               Shanghai, China | Apr. 2021 – Aug. 2021
- Built a vertical federated linear regression model that enables companies to collaboratively train a predictive model of 78% accuracy with 5 GB local data of different feature spaces to preserve customer privacy
- Designed and coded a joint computation mechanism of gradient for loss function among several participants
- Wrote scripts to encapsulate LibOS and big data applications for training federated learning (FL) models
- Evaluated and debugged an LSTM-based horizontal FL model with Intel Software Guard Extensions

**ProQuest LLC |** Deep Learning Intern                               Ann Arbor, MI | Jan. 2020 – Dec. 2020
- Proposed a hierarchical deep learning (DL) model for patent code classification to capture class distribution
- Trained and assembled single-label DL models DistilBert of different classification levels on 1M Google public patents with optimal reweighting and resampling techniques; achieved 90% precision
- Analyzed coverage error of multi-label RoBERTa and located tough classes with high error; identified hyperparameters to relieve the extreme multi-label text classification problem; achieved 80% precision
- Obtained the MDP Summer Fellowship; the optimized system adopted after product delivery

**NetEase Inc.** | Android Developer Intern                               Hangzhou, China | Jun. 2020 – Aug. 2020
- Added new Lottie animations to UI interface and beautified user profile display in the music app with Kotlin
- Developed 10+ new features such as "follow anchorman"; partnered with artists, QA to maintain best practices
- Analyzed and optimized edge cases, usability, and general reliability; tested and debugged code for robustness

## PROJECT EXPERIENCE

**Database Management System**                               Jan. 2021 – Feb. 2021
- Designed an ER model and translated it to relational schemas; created data tables for 80k+ Facebook records
- Programmed an efficient Java application with JDBC and Oracle to execute each SQL database query in 3ms
- Created a mongo data collection and executed queries with MongoDB; implemented Grace Hash Join in C++

**Network File Server with Customized Thread Library**                               Feb. 2020 – Mar. 2020
- Coded a thread library from scratch in C++ to ensure concurrency and atomicity for multiprocessor systems
- Built a multi-threaded file server from scratch enhancing its concurrency with upgradable reader/writer locks
- Realized encrypted TCP network communication with sockets for the client-server system

## PUBLICATION

**S2-CAN: Sufficiently Secure Controller Area Network**
Mert D. Pesé, Jay W. Schauer, Junhui Li, Kang G. Shin.
2021 Annual Computer Security Applications Conference, ACSAC 2021 [Webpage]