# Furkan Ozyurt

( +1 (618) 802-0464

∑ furknozyurt@outlook.com

New York, United States

# EXPERIENCE

### Data Scientist (Contract)

#### Amgen

#### **a** 03/2022 06/2023 **Q** Cambridge, United States

- Fine-tuned large language models (e.g., BERT, ROBERTA, GPT 3) on company-specific documents to perform text summarization, generation, paraphrasing, classification, and question answering.
- Developed and deployed a machine learning pipeline to predict whether changes in documents required reporting.
- Integrated the model into an application and implemented performance monitoring. The pipeline reduced decision-making time by 95%.
- Built a search engine that retrieved documents most similar to user inputs. This cut the time spent searching for information by 99% after being integrated into the application.
- Developed robust, scalable, and automated ETL/data pipelines to provide the team with reliable, high-quality, and up-to-date data.
- Optimized a previously developed machine learning pipeline and reduced runtime from 18 hours to 4 hours (a 75% reduction) using Spark.

# Associate Engineer (Contract)

#### Amgen

**= 09/2020 03/2022** Cambridge, United States

- · Developed a deep learning pipeline that classified documents into categories with 80% accuracy and integrated it into an application which saved the department approximately \$500,000.
- Designed and implemented an end-to-end machine learning pipeline in AWS SageMaker to forecast product consumption for key company products across multiple locations.

# PROJECTS

#### Self-Learning

11/2023 02/2024

Attps://github.com/ozvurtf/self-learning

- Worked on a project with the goal of developing a reinforcement-learning free system that enables a simulated truck to back up to a target position from any initial location autonomously without collecting any data manually.
- Developed a custom loss function tailored to the challenges of the task because standard loss functions were not useful.
- · Built two separate models: one to create internal representation of the environment in which the truck operates and another to determine the optimal steering angle for the truck's next move based on the internal representation of the environment.
- · Successfully trained these models to enable the truck to consistently reach the target position smoothly no matter where it is initialized.

#### Attention in CUDA

🛱 03/2025 - present

Attps://github.com/ozyurtf/attention-cuda

- · Implementing multi-head attention mechanism in CUDA by utilizing shared memory, coalesced memory, warp shuffle, and tiling.
- · Profiling and optimizing CUDA kernels using Nsight Systems and Nsight Compute to reduce inference latency.

# EDUCATION

#### Master of Science - Computer Science New York University - Courant

New York, United States

#### Bachelor of Science - Industrial Engineering (Mathematics Minor) **Istanbul Technical University**

**m** 08/2016 05/2020 **Q** Istanbul, Turkey

- (Ranked in the top 0.9% of students in the national university exam)
- (Jointly completed the program with Southern Illinois University Edwardsville)

# **SUMMARY**

Master of Science student in Computer Science at New York University, with experience in building data pipelines and training/optimizing/deploying/monitoring machine learning and deep learning models. Possesses strong theoretical knowledge of various deep learning architectures, including CNNs, RNNs, LSTMs. Transformers, Autoencoders, VAES, GANS, and Diffusion Models. Has a strong background in GPU architecture and CUDA.

# KEY ACHIEVEMENTS

#### **Decision-Making Time Reduction**

Reduced the decision-making time by 95% with a deep learning pipeline.

#### Document Classification Savings

Built and deployed a deep learning document classifier (80% accuracy), saving the department \$500,000.

#### Pipeline Runtime Optimization

Reduced the runtime of a macine learning pipeline by 75% using Spark's parallel processing capabilities.

#### Information Search Time Reduction

Built and deployed a document search engine that retrieved results based on input similarity. This reduced information retrieval time by 99%.

# SKILLS

Programming Languages & Query Languages Python, SQL, C/C++, CUDA

#### Machine Learning & AI

Machine Learning, Deep Learning, Natural Language Processing (NLP), Large Language Models (LLMs), MLOPs

ML Frameworks & Libraries PyTorch, Tensorflow, MLFlow, Huggingface

Data Engineering & Big Data Databases, Databricks, Delta Lakes, Spark

#### Cloud Computing

AWS (EC2, S3, IAM, Athena, EMR, Glue, Redshift, Sagemaker), Microsoft Azure

**Development Tools** 

Git, Docker, Kubernetes

# FIND ME ONLINE

**Personal Website** https://ozyurtf.github.io/

Github https://github.com/ozyurtf

Linkedin http://www.linkedin.com/in/ozyurtf/