

Experience

Senior Machine Learning Engineer @ Delivery Hero SE

Berlin - 05/2022 to Present

I am a senior machine learning engineer at Delivery Hero that creates tools and infrastructure to improve the productivity of the data scientists and data analysts in our vertical (over 50). My main responsibilities include:

- Deploying and maintaining a modern workflow management service, [Metaflow](#), and the production workflow orchestrator, Argo Workflows + Argo Events ([Outerbounds Talk on Metaflow at Delivery Hero](#))
- Creating and maintaining a project template for easily interacting with CI/CD systems (Drone and Spinnaker) and implementing best coding practices (dependency management with poetry, containerization, linting and autoformatting with Ruff)
- Designing and organizing a custom python package to unify experimentation (AB testing and pretest power analysis) for our vertical

In addition to my primary responsibilities I have also lead our vertical to transition to a modern infra stack and have made many contributions to the platform teams tooling including:

- Enabling a custom cluster type for ml workloads (terraform)
- Enabling features in our application deployment pipeline (golang)
- Pioneering best practices for enabling externally accessible applications with the traffic engineering team
- Pioneering best practices to utilize our new CI/CD system (drone + ArgoCD)

The project template was born out of necessity at the start of my time at DH. I created a data science project template that could be continuously updated once with a simple command that allowed our data science projects to follow best practices and fully utilize our kubernetes infrastructure for experimentation and

production deployments. This included writing a suite of small libraries that allow data scientists customize their projects with minimal dependencies, which replaced a monolith library that would often install unneeded dependencies and cause unnecessary dependency conflicts. As part of this effort, I standardized the dependency management, linters and autoformatters (previously a hodgepodge of outdated programs) that we used to a modern stack (poetry + ruff), which improved the maintainability, reliability and usability of all our projects.

Another issue, we had was time from ideation to production with data science models. Specifically, the Airflow setup made it difficult to iterate on developing and training models (see my [Outerbounds metaflow talk](#) mentioned above). Our team decided to explore alternatives to Airflow and settled on Metaflow, which is a modern orchestrator that makes running a DAG/flow locally and in the cloud with Kubernetes easy for our data scientists. I was in charge of developing the POC to deploy the infrastructure with IaC. We went from zero usage at the end of 2022 to now almost all projects use Metaflow for model development and some for production deployments. Our tribe's model training setup was so successful that I am not tasked with developing this as a global ML platform product. I also created custom dashboards and monitoring of these flows with Prometheus and Grafana, which required an [upstream contribution to add kubernetes labels](#) to metaflow. In addition, I developed a kubernetes mutating webhook in Go to assist in adding cost tagging to all of our BigQuery queries made from metaflow. A concrete example of the benefits of this change is when I converted one of our model training Airflow DAGs to a Metaflow flow and reduced the training time from 28 hours on GPU machines to 5 hours on CPU machines.

I really enjoy solving problems even if they are not my primary responsibility. In early 2023, we had lots of teams that wanted to do AB testing for new models / features and all of them would reinvent the wheel to do the statistical analysis. So our team made it a priority to unify these efforts with a single python package that was easy to use and extend, while being efficient and producing statistically correct results. The challenge here was less the coding and more managing the various parties involved. This required a lot of project management to get the various parties to encourage them to contribute to the project, so we wouldn't have to write the entire thing. All the while managing the slightly different needs of each team. The final results is a well documented library that does statistical analysis including standard testing and pretesting that is easy enough to use for even non-technical

users such as product managers that is maintained by a council of different data scientists and data analysts.

Generally, I find myself drawn to more infrastructure-heavy projects. I am proficient with kubernetes (helm / helmfile), terraform (and terragrunt) and CI/CD (mostly Drone + Spinnaker, but also Github Actions). With my contributions to the platform team, I am often writing terraform, Golang, bash scripting, jq, and all kinds of yaml since these are the primary languages in use in our infra. Lastly, I have also worked on a few Kotlin projects to integrate our products into JVM-based consumers.

Senior Machine Learning Engineer @ Solvemate GmbH

Berlin - 02/2020 to 04/2022

I led the project to implement free-text input into an existing decision tree-based click-bot. One project I am particular proud of is a word vectorization microservice that allows for configurable pipelines of word vectorization models, preprocessing and postprocessing, and sentiment analysis. This included a solution to benchmark the different pipelines (model + preprocessor combinations) used by this microservice. I also trained custom models in multiple languages for our word vectorization. Currently, I am tasked with troubleshooting and maintaining of our kubernetes-based architecture and develop a custom python-based cli-tool to automate tasks used by the entire engineering team. My latest project is adding a voice channel to our bots so our product can be accessed via a telephone or another voice-based channel.

From a software engineering / dev-ops perspective, I introduced a modern Python stack based on [fastAPI](#) and [Pydantic](#), full pytest-based testing, and automatic linting/formatting via [isort](#), [black](#), flake8 to projects throughout the company. In addition to maintaining kubernetes, I also maintain our Google Cloud Platform, Jenkins instances, and various other legacy services as part of my dev-ops duties.

I made a [large OSS-contribution](#) to the [Microsoft's presidio](#) anonymization library to completely rewrite the python analyzer engine. The primary feature that I added was the ability to use multiple language models at once, which enabled multilingual capabilities to the library. But I also completely modernized their unit tests, optimized the code, and removed a lot of hard-coded variables.

Machine Learning Engineer @ i2x GmbH

Berlin - 04/2019 to 02/2020

I primarily worked on NLP projects at i2x. This included coordinating with a team of over 15 labelers to create datasets from our raw data, creating tooling to manage the data pipeline from raw data to production datasets, and even performing labeling myself (all data was German). For these completed datasets, I wrote tooling to train models and gather metrics with libraries like huggingface transformers (BERT), fasttext, sklearn and others. I integrated pytorch models into our GRPC production, which was previously tensorflow only. Another project that I worked on was a sentence level similarity search using word embeddings and approximate nearest neighbors, tool capable of searching tens of millions of utterances in milliseconds. This project was also dockerized and GRPC-based for easy of deployment.

Although my primarily responsibilities were NLP, I experimented with several newer ASR systems (fairseq, wav2letter, and NVIDIA's NeMo) and even contributed several PRs to NeMo. Regarding general software development, code optimizations was a task that I excelled at specifically, I often built optimized dataloaders that could process data online during training, which sped up training and greatly reduced memory requirements for large datasets. I also believe in contributing back to upstream OSS projects and had PRs accepted at huggingface transformers, facebook's FBGEMM, pytorch, pytorch tutorials, and PyThaiNLP.

Lastly, due to my proficiency with linux, I was the ML team's de facto systems administrator for our in-house multi-GPU dev machines, so I'm proficient with the *nix cli ecosystem.

Lead Developer @ PyTorch Audio

Berlin - 06/2017 to 01/2019

I was the lead developer for PyTorch's official audio loading library, torchaudio, and continue to be involved in the project. I work directly with the PyTorch team at Facebook on this project. I have implemented input error checking testing, audio datasets, variable length input collate functions, audio IO functionality via SoX (Sound eXchange), and added audio transformations (both implemented directly in PyTorch and using PyTorch's c++ hooks for SoX's effects chain). Additionally, I transferred the main PyTorch project's code flaking and documentation standards to

this project. In October 2018, I was invited to the PyTorch developer conference for my contributions to the project. I continue to follow the project and contribute when I can. I am an author on the paper, "torchaudio: Building Blocks for Audio and Speech Processing" (submitted to ICASSP Oct. 2021).

Machine Learning Fellow @ Fellowship.ai

Berlin - 01/2019 to 04/2019

Once again, I developed an end-to-end solution to do image classification on a fashion dataset. We achieved state of the art results in the task classifying a fashion style from the Fashion14 dataset. In preparation for a weekly reading group, I optimized an implementation of memory networks on the babi question answer dataset. Worked on implementing semi-supervised semantic segmentation utilizing Deep Extreme Cut and DeepLab-v3+.

Machine Learning Engineer (Student Job) @ YEAY GmbH

Berlin - 01/2018 to 06/2018

My primary role was to use object recognition on videos to identify different classes of clothing. I utilized a Mask R-CNN-based network with custom additions for blur detection to select the clearest frames. Data prep included finding an appropriate dataset, writing scripts to translate between various annotation formats (PASCAL, COCO, YOLO) and add missing fields such as segmentation naively or algorithmically (i.e. GrabCut). Additionally, I dockerized my solution to make it more portable.

Co-Founder and CTO @ Cygnus Association Management

Atlanta (Remote) - 08/2014 to 01/2019

Founded an real estate management company to manage residential and commercial properties in Atlanta, GA. Responsible for the organization's technological solutions, including but not limited to Google Apps infrastructure, company web presence, and payment services. Collaborate with other cofounders regarding all strategic aspects of the company. Work 100% remotely.

Senior Consultant @ Navigant

Atlanta - 09/2006 to 04/2013

Completed re-underwrites of over 750 residential home loans for several of the TBTF (Too Big To Fail) US banks in disputes related to the US housing crisis. Managed and reviewed team members' work product. Developed software to automate repetitive tasks that cut the per-loan prep time by an order of a magnitude. Promoted in 2007 and 2010.

Education

- Master of Science in Economics and Management Science from Humboldt Universität in 2018
- Bachelors of Arts in Economics from Emory University in 2005

Languages

- English (native speaker)
- German (professional working proficiency - C1)

| Straßburger Straße 38, 10405 Berlin | +49 152 0476 3223 | david@sologourmand.com |