

# Filip Skogh

📍 Zürich, Switzerland

✉️ filipskogh99@hotmail.com

☎️ +41 76 269 4039



## PROFILE

Developer with 3 years of experience in C++. Industry experience in optimizing machine learning algorithms and working with LLMs in production.

## EDUCATION

<b>ETH Zürich</b> SEMP student, M.Sc. Computer Science	Zürich, Switzerland
• <b>Thesis:</b> Weakly Supervised Video Object Segmentation	Sep 2022 - Sep 2023
• <b>Activities:</b> ETH Analytics Club	
• <b>Courses:</b> Natural Language Processing, Advanced Machine Learning, Computer Vision Project	
<b>Chalmers University of Technology</b> M.Sc. Data Science and AI	Gothenburg, Sweden
• <b>GPA:</b> 4.8/5.0	Aug 2021 - Sep 2023
• <b>Courses:</b> Non-linear Optimization, Stochastic Processes and Bayesian Inference, Algorithms Computer Vision, Reinforcement Learning, Image Analysis, Large-scale Distributed Computation	
<b>Nanyang Technological University</b> Exchange Student, B.Sc. Computer Science	Singapore, Singapore
• <b>Courses:</b> Digital Signal Processing, Cryptography, Operating Systems, Computer Networks	Jan-Jun 2020
<b>Luleå University of Technology</b> B.Sc. Computer Science and Engineering	Luleå, Sweden
• <b>GPA:</b> 5.0/5.0	Aug 2018 - Jun 2021

## EXPERIENCE

<b>Machine Learning Intern</b> IBM Research	Zürich, Switzerland
• Got a deep understanding of different transformer variants in the Huggingface transformers library, along with optimizations like GQA and KV-Cache.	May 2024 - Oct 2024
• Open source contributions to drug discovery repos and huggingface spaces.	
• Integrated the Sinkhorn algorithm in to Transformer self attention to produce doubly-stochastic attention matrices.	
<b>Machine Learning Engineer</b> Logiblox AG	Zürich, Switzerland
• Increased token generation speed by over 300% on LLM inference server by implementing quantization, parallelization and deciding on the optimal hardware.	Oct 2023 - Apr 2024
• Set up from scratch a dockerized server with CI/CD pipeline running llama.cpp for LLM inference.	
• Optimized (caching and operation merging) compiler that converts visual no-code representation in to executable python code.	
• Set up a database and code-sandbox for an user facing LLM app that clean datasets.	
<b>Master Thesis Student</b> Computer Vision Lab ETH Zürich	Zürich, Switzerland
• <b>Supervisors:</b> Prof. Fisher Yu and Dr. Martin Danelljan	Feb 2023 - Sep 2023
• <b>Objective:</b> Reduce the annotation burden for video object segmentation.	
• <b>Solution:</b> Implemented a loss function that utilize spatial and temporal information between video frames to derive a consistency based loss.	
• Achieved a 90% relative performance to the fully-supervised model without having access to the video mask annotations.	
<b>Optimization Research</b> University of Massachusetts	Massachusetts, United States
• Our paper received the Best Paper Runner-up award at the 14 <sup>th</sup> IEEE IGSC 2023 conference in Toronto, Canada.	Jun 2022 - Sep 2022
• Mixed-integer optimization in a distributed systems setting.	
• Developed a server load scheduler that route requests to data-centers such that carbon is minimized while satisfying latency constraint.	
<b>Security Software Engineer</b> Orange Cyberdefense	Stockholm, Sweden
• Developed automated threat response scripts in Python that block ransomware, C&C servers and take snapshots for forensics. Deployed globally on 50.000+ end points in 70+ countries.	Summers 2019 - 2021
• Offensive security: metasploit, kerberos security, email security.	
• Network security, ssh and SMB server security, developed software to periodically scan network using nmap and masscan.	

## PROJECTS

---

**Packet intercept proxy:** C++ project developed continuously for three years. Taught myself java internals, java native interface and reflection. Reverse engineered encryption protocols and ciphers to intercept traffic at packet level.

**Teaching Transformers arithmetic:** Trained a GPT-style Transformer in PyTorch to learn addition with different tokenizers. Showed that if the digits are tokenized in reverse order the problem is easier.

**Neural network certification:** Developed custom network layers in PyTorch to propagate intervals through a network allowing us to deterministically prove properties about robustness, certain fairness guarantees and that adversarial attacks are not possible. This project was part of the course Reliable and Trustworthy AI at ETH.

**Transformer inference:** Implemented the decoder-only Transformer in PyTorch with a KV-cache for improved inference speed.

**Blockchain implementation:** Implemented parts of the Bitcoin protocol from scratch to create, (i) a wallet address derived from an elliptic curve public key, (ii) a signed transaction which can be broadcasted to the network, and (iii) a block verifier.

**Google Developer Student Club:** Built an interactive 3D learning game in Unity3D C# by building a hospital simulator. The project idea was conceived by medical professors prompted by the pandemic and was aimed to simulate medical students' practicum. During the project I worked in close contact with medical professionals and translated medical procedures into implementable scenarios in-game.

**RANSAC:** Iterative parameter estimation using RANSAC with optimal hypothesis testing that minimizes the number of tests performed. The Project was motivated by the scarcity of available implementation and was based on the original white paper.

## TECHNICAL SKILLS

---

**Languages:** Python, C++, Java, C, Matlab, English, Swedish

**Frameworks:** PyTorch, Hugging Face, Map-Reduce, PySpark, OpenCV, Flask, Java Native Interface, MySQL

**Miscellaneous:** Vim, Slurm, FastAPI, Docker, Git, Regex