

Aleksander Lempinen

**MLOps approach for application specific performance  
tuning for machine learning systems**

Master's Thesis in Information Technology

January 25, 2023

University of Jyväskylä

Faculty of Information Technology

**Author:** Aleksander Lempinen

**Contact information:** aleksander.lempinen@gmail.com

**Supervisor:** TODO supervisor

**Title:** MLOps approach for application specific performance tuning for machine learning systems

**Työn nimi:** TODO samma på finska

**Project:** Master's Thesis

**Study line:** Educational Technology

**Page count:** 15+0

**Abstract:** TODO abstract

**Keywords:** L<sup>A</sup>T<sub>E</sub>X, gradu3, Master's Theses, Bachelor's Theses, user's guide

**Suomenkielinen tiivistelmä:** TODO tiivistelmä suomeksi

**Avainsanat:** MLOPS, TODO

## **Glossary**

ML

Machine Learning

MLOps

Machine Learning Operations

TODO

TODO

# Contents

1	INTRODUCTION .....	1
2	MACHINE LEARNING OPERATIONS .....	2
2.1	Introduction Machine Learning .....	2
2.1.1	Overview .....	2
2.1.2	Machine learning algorithms .....	2
2.1.3	Hyperparameter optimization .....	2
2.2	DevOps .....	3
2.2.1	Overview .....	3
2.2.2	Performance metrics .....	3
2.2.3	Performance tuning .....	4
2.3	MLOps .....	4
2.3.1	Overview .....	4
2.3.2	AutoML .....	5
2.3.3	Performance prediction and early stopping .....	5
3	METHODS .....	6
4	RESULTS .....	7
5	DISCUSSION .....	8
5.1	Research Questions revisited .....	8
5.2	Limitations and Validity .....	8
5.3	Related Work .....	8
5.4	Future Work .....	8
6	CONCLUSIONS .....	9
	BIBLIOGRAPHY .....	10

# 1 Introduction

Machine learning (ML) systems are widely adopted and many organizations successfully have ML models running in production.

- Problem with MLOps/ML tools
  - traditional ML performance metrics such as accuracy
  - fancy features such as neural architecture search, AutoML, performance tuning
  - fancy techniques such as early stopping, grid search, bayesian optimization search etc.
  - little support for non-ML metrics: CPU util, memory used, latency, throughput (images/s etc.), hardware required (CPU, GPU, TPU etc.)
  - ML in production has many objectives besides accuracy for example: satellite image processing model A took 4-5h and model B took 5min. Model A is infeasible in production despite being more accurate.
  - "Better" depends on the specific application
- Hypothesis: Early stopping will speed up computing non-ML metrics (CPU util, Memory use, GPU/TPU requirement, latency, throughput)

TODO smaller and smaller devices, limited resources

Real-world ML systems in addition to ML performance metrics will have similar performance metrics as traditional software systems. The aim of the study is to tune performance metrics particularly relevant to real-world ML systems.

The main contribution of this master's thesis is using ML performance tuning techniques such as early stopping for tuning a wider range of real-world ML system performance metrics.

## **2 Machine Learning Operations**

This chapter introduces the basic concepts of ML in section 2.1 and basic concepts of DevOps in section 2.2. Finally in section 2.3 DevOps and ML are combined to form a new concept of MLOps.

### **2.1 Introduction Machine Learning**

#### **2.1.1 Overview**

TODO what is ML

TODO supervised vs unsupervised vs reinforcement

TODO training vs serving

#### **2.1.2 Machine learning algorithms**

TODO Model evaluation

#### **2.1.3 Hyperparameter optimization**

Model parameters given as part of a configuration to the machine learning model are called hyperparameters (Yang and Shami November 2020).

Hyperparameter optimization techniques include grid search, random search, gradient based optimization and Bayesian optimization and they have different benefits and limitations (Yang and Shami November 2020)

Neural Architecture optimization and Meta modeling are similar to hyperparameter optimization where model structure or modeling algorithm is treated as a tunable parameter (Baker et al. November 8, 2017)

## **2.2 DevOps**

### **2.2.1 Overview**

There is little consensus on the exact definition of DevOps, but especially collaboration between development and operation is emphasized (Mishra and Otaiwi November 1, 2020; Waller, Ehmke, and Hasselbring April 3, 2015). DevOps can be studied from different points of view such as culture, collaboration, automation, measurements and monitoring (Mishra and Otaiwi November 1, 2020; Waller, Ehmke, and Hasselbring April 3, 2015). This thesis is mostly focused on the automation, measurements and monitoring parts of DevOps.

TODO: picture about devops

Continuous integration, continuous deployment and continuous monitoring are well known practices in DevOps (Waller, Ehmke, and Hasselbring April 3, 2015) describing the automatic nature of integrating, deploying and monitoring code changes. Performance profiling and monitoring are similar activities and the main difference is whether it's done during the development process or during operations respectively (Waller, Ehmke, and Hasselbring April 3, 2015). DevOps bridges the gap between evaluating performance during the development process and during operations (Brunnert et al. August 18, 2015).

TODO resource allocation/resource consumption, small memory software, benchmarking

### **2.2.2 Performance metrics**

Performance metrics are fundamental to all activities involving performance evaluation such as profiling or monitoring (Brunnert et al. August 18, 2015). Common metrics involve measuring the CPU, but other metrics such as memory usage, network traffic or I/O usage are not as well defined as a CPU metric (Brunnert et al. August 18, 2015).

- Task Completion time
- Throughput
- Latency
- CPU usage

- GPU usage
- RAM usage
- VRAM usage
- I/O usage
- Network traffic

### **2.2.3 Performance tuning**

Preprocessing

Training

Serving Latency

Resource demands might change depending on the inputs (Brunnert et al. August 18, 2015) making it important to systematically measure performance not only based on code changes but also on configuration changes or even data changes.

## **2.3 MLOps**

### **2.3.1 Overview**

Performance measuring software is not new, but ML brings additional challenges in the form of models and data which requires a modified approach (Breck et al. 2017). It is also important to note, that not every data scientist or machine learning engineer working on machine learning systems has a software engineering background (Finzer 2013) and might lack the necessary knowledge to apply software engineering best practices to machine learning systems.

TODO what DevOps brings to ML

TODO Continuous Training



### **2.3.2 AutoML**

Machine learning systems in addition to machine learning performance metrics and system performance metrics will have their performance metrics tied to product or organization metrics such as user churn rate or click-through rate (Shankar et al. September 16, 2022). Choosing the right metrics to evaluate a machine learning system is important and the metrics will be different for different machine learning systems (Shankar et al. September 16, 2022).

Automated Machine Learning (AutoML) aims to minimize human intervention in completing data analytics tasks using machine learning algorithms (Yang and Shami November 2022).

### **2.3.3 Performance prediction and early stopping**

Performance prediction is an important step to reduce the amount of computation required for neural architecture search and hyperparameter optimization (Baker et al. November 8, 2017).

### 3 Methods

The scope of the study is limited to 5 performance metrics of 3 different ML models trained and tested on 3 different datasets.

This master's thesis asks the following research questions:

- *RQ1*: How can we measure the performance of real-world ML systems?
- *RQ2*: How can these performance metrics be effectively tuned?

TODO This is a methods chapter

TODO no tool comparison

TODO different models, different datasets

TODO different resources (memory, time, accuracy)

## 4 Results

TODO This is a results chapter

## **5 Discussion**

### **5.1 Research Questions revisited**

### **5.2 Limitations and Validity**

### **5.3 Related Work**

To find relevant related work both reverse snowballing and forward snowballing is used on a set of MLOps papers previously known to the author.

Benchmarking ML systems (Cardoso Silva et al. December 2020)

### **5.4 Future Work**

TODO This is a discussion chapter

## **6 Conclusions**

TODO This is a conclusions chapter

## Bibliography

Baker, Bowen, Otkrist Gupta, Ramesh Raskar, and Nikhil Naik. November 8, 2017. *Accelerating Neural Architecture Search Using Performance Prediction*, arXiv:1705.10823, November 8, 2017. Visited on January 25, 2023. <https://doi.org/10.48550/arXiv.1705.10823>. arXiv: 1705.10823 [cs]. <http://arxiv.org/abs/1705.10823>.

Breck, Eric, Shanqing Cai, Eric Nielsen, Michael Salib, and D. Sculley. 2017. “The ML Test Score: A Rubric for ML Production Readiness and Technical Debt Reduction”. In *Proceedings of IEEE Big Data*.

Brunnert, Andreas, Andre van Hoorn, Felix Willnecker, Alexandru Danciu, Wilhelm Hasselbring, Christoph Heger, Nikolas Herbst, et al. August 18, 2015. *Performance-Oriented DevOps: A Research Agenda*, arXiv:1508.04752, August 18, 2015. Visited on January 17, 2023. <https://doi.org/10.48550/arXiv.1508.04752>. arXiv: 1508.04752 [cs]. <http://arxiv.org/abs/1508.04752>.

Cardoso Silva, Lucas, Fernando Rezende Zagatti, Bruno Silva Sette, Lucas Nildaimon dos Santos Silva, Daniel Lucrédio, Diego Furtado Silva, and Helena de Medeiros Caseli. December 2020. “Benchmarking Machine Learning Solutions in Production”. In *2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA)*, 626–633. 2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA). <https://doi.org/10.1109/ICMLA51294.2020.00104>.

Finzer, William. 2013. “The Data Science Education Dilemma”. *Technology Innovations in Statistics Education* 7 (2). Visited on January 17, 2023. <https://doi.org/10.5070/T572013891>. <https://escholarship.org/uc/item/7gv0q9dc>.

Mishra, Alok, and Ziadoon Otaiwi. November 1, 2020. “DevOps and Software Quality: A Systematic Mapping”. *Computer Science Review* 38 (November 1, 2020): 100308. ISSN: 1574-0137, visited on January 17, 2023. <https://doi.org/10.1016/j.cosrev.2020.100308>. <https://www.sciencedirect.com/science/article/pii/S1574013720304081>.

Shankar, Shreya, Rolando Garcia, Joseph M. Hellerstein, and Aditya G. Parameswaran. September 16, 2022. *Operationalizing Machine Learning: An Interview Study*, arXiv:2209.09125, September 16, 2022. Visited on December 7, 2022. <https://doi.org/10.48550/arXiv.2209.09125>. arXiv: 2209.09125 [cs]. <http://arxiv.org/abs/2209.09125>.

Waller, Jan, Nils C. Ehmke, and Wilhelm Hasselbring. April 3, 2015. “Including Performance Benchmarks into Continuous Integration to Enable DevOps”. *ACM SIGSOFT Software Engineering Notes* 40, number 2 (April 3, 2015): 1–4. ISSN: 0163-5948, visited on January 17, 2023. <https://doi.org/10.1145/2735399.2735416>. <https://doi.org/10.1145/2735399.2735416>.

Yang, Li, and Abdallah Shami. November 2020. “On Hyperparameter Optimization of Machine Learning Algorithms: Theory and Practice”. *Neurocomputing* 415 (): 295–316. ISSN: 09252312, visited on January 25, 2023. <https://doi.org/10.1016/j.neucom.2020.07.061>. arXiv: 2007.15745 [cs, stat]. <http://arxiv.org/abs/2007.15745>.

———. November 2022. “IoT Data Analytics in Dynamic Environments: From An Automated Machine Learning Perspective”. *Engineering Applications of Artificial Intelligence* 116 (): 105366. ISSN: 09521976, visited on January 25, 2023. <https://doi.org/10.1016/j.engappai.2022.105366>. arXiv: 2209.08018 [cs, eess]. <http://arxiv.org/abs/2209.08018>.