



**Czech
Technical
University
in Prague**

F3

**Faculty of Electrical Engineering
Department of Computer Science**

Unassisted project report

Lukáš Forst

**Supervisor: Ondřej Vaněk, Ph.D.
January 2019**

Contents

1 Introduction	1
1.1 Problem definition	2
1.1.1 Formal definition	2
1.1.2 Load Balancer Requirements .	2
1.2 Motivation to solve it	2
2 Technical Background	3
2.1 Optimization Algorithms	3
2.1.1 Linear Programming	3
2.1.2 Heuristic algorithms	3
2.1.3 Selected algorithm	4
3 State of the art	5
3.1 Load Balancing.....	5
Bibliography	7



Chapter 1

Introduction

Optimization algorithms and solutions build on them are widely used in current manufacturing industry to reduce production costs. With more and more production automatization, optimization algorithms can manage and schedule whole factories with maximum available efficiency.

Complexity of optimization problems could be huge and therefore performance requirements are sometimes not easily satisfiable. Using one powerful instance of optimization algorithm in cloud seems like a solution for problems with smaller complexity, but what if we have multiple huge problems where each is performance demanding? Of course, we can create multiple instances, but that would be expensive and not well manageable and scalable since adding another instances manually requires some time and it is not much flexible. Another disadvantage of this approach is the fact, that optimization algorithm is not running 100% of time and thus resources allocated by this algorithm are unused while other algorithm instances could be potentially overwhelmed. Also paying for unused hardware is wasting money and optimization algorithms are supposed to save money.

Now imagine having two completely different problems that each requires its own application which visualises data and optimization algorithm to compute some kind of plan, this algorithm can be generic enough to operate on both domains with same code base, but it requires a lot of performance resources. If we use monolithic architecture of both applications, we would have same code in two applications, but what is even worse, we would need two powerful machines to run our applications. As previously mentioned, these two machines would not be using their power whole time and would be mainly idle. What if one application runs only few minutes a day, but needs that power to complete tasks in time? A lot of resources would be wasted if it has its own server, but using not powerful server would lead to increasing duration of ongoing tasks which is something we do not want.

In this paper I would like to introduce **load balancer** specifically developed for optimization algorithms which could potentially minimize resources wasting and increase performance using correct utilization distribution across multiple instances of optimization algorithms.

■ 1.1 Problem definition

The problem with implementation of optimization algorithms in applications is that their performance requirements are quite high and are fully utilized only while working. Optimization algorithm is not running all the time and for that reason hardware resources are mainly unused. These unused resources could be potentially used by another instance of algorithm or can be shutdown completely to reduce hosting costs.

■ 1.1.1 Formal definition

Load balancer take in account following scheduling job properties.

■ 1.1.2 Load Balancer Requirements

■ 1.2 Motivation to solve it

Chapter 2

Technical Background

2.1 Optimization Algorithms

This work does not contain any own algorithm implementation for generic optimization problems, instead I would like to use pre-prepared and already implemented optimization solver. First we must specify which kind of approach we would like to choose. We have many options, how to represent and then solve scheduling problem such as

- Linear programming - mixed integers programming
- Constraint programming - heuristics algorithms

2.1.1 Linear Programming

Some general information about linear programming goes here

Advantages of linear programming approach

Disadvantages of linear programming approach

Existing solutions

GLPK Google Optimization Kit

2.1.2 Heuristic algorithms

General information about heuristic

Advantages of heuristic algorithms

The main advantage of heuristic algorithms is that they offer a quick solution for problem they are solving.

■ Disadvantages of heuristic algorithms

The main downside of HA is the fact, that they can't guarantee that found solution is the optimal one.

■ Existing solutions

I would like to mention two implementations of heuristics algorithms - OptaPlanner and TASP.

■ OptaPlanner

OptaPlanner is an open source generic heuristics based constraint solver. It is optimized to

■ TASP

Task and Asset Scheduling Platform is a lightweight framework developed by Blindspot Solutions designed to solve a large variety of optimization and scheduling problems from the area of logistics, workforce management, manufacturing, planning and others. It contains a modular, efficient planning engine utilizing latest optimization algorithms. TASP is delivered as a software library to be used through its API in applications which require powerful scheduling capabilities.

It is written in Kotlin and runs on JVM.

■ 2.1.3 Selected algorithm

I decided to use one linear solver and one heuristic algorithm to test load balancing server. This will provide us heterogeneous environment for distinguish optimization tasks as well as different demands on performance.

While choosing suitable solvers I was looking mainly at possibility running on JVM and their API as well as at their suitability for my paper. For final testing I selected **GLPK** as linear solver, mainly because it is widely used linear programming kit and because of it's convenient Java interface.

As a representative of heuristics algorithms I selected **TASP** because of it's great scalability, Kotlin interface and because I have already worked with it and I'm familiar with multiple TASP implementations.

Chapter 3

State of the art

3.1 Load Balancing

There will be some info about how should server balance itself.

- prioritisation - mainly done by priority queues
- handover
- instance sizing
- algorithms - following are methods used in network balancing -> probably can't be used because we need to manage scheduling which is heavy on computer resources like CPU/RAM/IO
 - The Least Connection Method
 - The Round Robin Method
 - The Least Response Time Method



Bibliography