

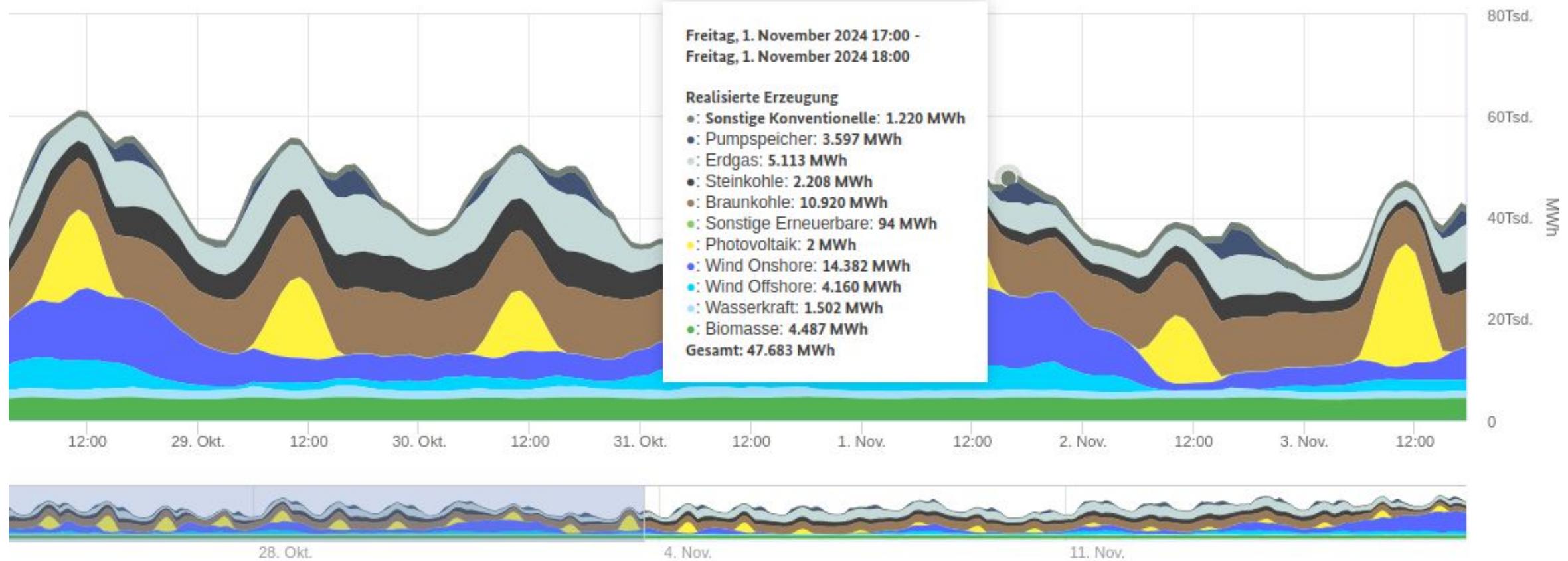
Masterarbeit Verteidigung

A Parameterizable Test Bed for Carbon Aware Job Scheduling

**Design IT.
Create Knowledge.**

www.hpi.de





Stromproduktion aus versch. Quellen, smard.de/home/marktdaten

Ein immer größerer Anteil des Stroms in Deutschland wird durch Erneuerbare Energien gedeckt. Das zeigen die neuen Zahlen zum 1. Halbjahr 2024. Ein Überblick zum Ausbau von Solar- und Windenergie.

HPI

bundesregierung.de/breg-de/aktuelles/ausbau-ern
euerbare-energien-2225808

Windräder auf See stehen häufiger still - vor allem vor SH

Stand: 25.04.2024 05:00 Uhr

Offshore-Windräder könnten deutlich mehr Strom nach SH liefern. Die Netzbetreiber schalten sie aber verstärkt ab, weil Leitungen überlastet sind. Das erhöht den Strompreis.

nrd.de/nachrichten/schleswig-holstein/Windraeder-auf-See-stehen-haeufiger-still-vor-allem-vor-SH,offshore924.html



Wusstest du schon: Wenn der Wind stark weht oder viel Sonne scheint, ist Energie besonders grün und günstig. Warum nicht dann deinen Strom verbrauchen? Tibber macht es dir einfach.

tibber.com

GLOBALE ERWÄRMUNG

Jeder einzelne der vergangenen zwölf Monate war weltweit der deutlich wärmste seiner Art

05. Juni 2024, 20:10 Uhr

mdr.de/wissen/umwelt-klima/waermste-zwoelf-monate-weltweit-waermster-meteorologischer-fruehling-deutschland-100.html

Strukturierte Literaturrecherche zum State of the Art

- Schlüsselwörter zu Nachhaltigkeit / Energieeffizienz und zu Rechenzentren



Google Scholar search results for "energy efficiency" "scheduling".

Artikel
Ungefähr 444.000 Ergebnisse (0,05 Sek.)

Beliebige Zeit
Seit 2024
Seit 2023
Seit 2020
Zeitraum wählen...

Scheduling in centralized cognitive radio networks for **energy efficiency**
S Bayhan, F Alagoz - IEEE Transactions on Vehicular ..., 2012 - ieeexplore.ieee.org
... We formulate the **scheduling** problem in CRNs as an **energy efficiency** maximization ... Next, we study the **scheduling** problem not only from an **energy efficiency** perspective but from a ...
☆ Speichern 99 Zitieren Zitiert von: 145 Ähnliche Artikel Alle 11 Versionen

Nach Relevanz sortieren
Nach Datum sortieren

Resource-conscious **scheduling** for **energy efficiency** on multicore processors
A Merkel, J Stoess, F Bellosta - ... of the 5th European conference on ..., 2010 - dl.acm.org
... and co-**scheduling** policies that improve performance and **energy efficiency** by combining

- jeweils 5 Abstracts lesen, Kategorisieren

Strukturierte Literaturrecherche zum State of the Art

- 145 Paper erkundet
- davon 99 "rot / irrelevant", 31 "orange / später", und 15 "grün / komplett lesen, sehr relevant"

| keyword 1 | keyword 2 | title |
|-------------------|-----------------|---|
| energy efficiency | datacenter | Beyond Energy-Efficiency: Evaluating Green Datacenter Applications for Energy-Agility |
| energy efficiency | datacenter | Energy Efficiency Dilemma: P2P-cloud vs. Datacenter |
| energy efficiency | datacenter | Energy Efficiency for Data Center and Cloud Computing: A Literature Review |
| energy efficiency | datacenter | Evaluation of the Heat and Energy Performance of a Datacenter Using a New Efficiency Index: Energy Usage Effectiveness Design – EUED. |
| energy efficiency | datacenter | Viability of datacenter cooling systems for energy efficiency in temperate or subtropical regions: Case study |
| energy efficiency | scheduling | Resource-conscious scheduling for energy efficiency on multicore processors |
| energy efficiency | scheduling | A review of energy-efficient scheduling in intelligent production systems |
| energy efficiency | scheduling | Integrating process optimization with energy-efficiency scheduling to save energy for paper mills |
| energy efficiency | scheduling | A genetic algorithm for energy-efficiency in job-shop scheduling |
| energy efficiency | scheduling | Optimal Scheduling of Demand Side Load Management of Smart Grid Considering Energy Efficiency |
| energy efficiency | compute cluster | GreenHDFS: Towards An Energy-Conserving, Storage-Efficient, Hybrid Hadoop Compute Cluster |
| energy efficiency | compute cluster | An Energy Efficiency Optimization and Control Model for Hadoop Clusters |
| energy efficiency | compute cluster | eeClust: Energy-Efficient Cluster Computing |
| energy efficiency | compute cluster | Energy efficient computing, clusters, grids and clouds: A taxonomy and survey |
| energy efficiency | compute cluster | An Empirical Study of Hadoop's Energy Efficiency on a HPC Cluster |
| energy efficiency | cloud | Managing overloaded hosts for energy-efficiency in cloud data centers |
| energy efficiency | cloud | Energy Efficiency for Cloud Computing Applications: A Survey on the Recent Trends and Future Scopes |
| energy efficiency | cloud | Adaptive Computational Solutions to Energy Efficiency in Cloud Computing Environment Using VM Consolidation |
| energy efficiency | cloud | A review of energy efficiency evaluation technologies in cloud data centers |
| energy efficiency | cloud | Energy efficiency in cloud computing data centers: a survey on software technologies |
| energy efficiency | load balancing | Energy efficiency of load balancing for data-parallel applications in heterogeneous systems |
| energy efficiency | load balancing | Energy Efficiency and Load Balancing in MANET: A Survey |

State of the Art - Papers

| Paper | Suspend & Resume Scheduling |
|--|-----------------------------|
| Let's Wait Awhile: How Temporal Workload Shifting Can Reduce Carbon Emissions in the Cloud | ✓ |
| Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions | ✓ |
| The War of the Efficiencies: Understanding the Tension between Carbon and Energy Optimization | ✓ |
| On the Limitations of Carbon-Aware Temporal and Spatial Workload Shifting in the Cloud | ✓ |
| The Online Pause and Resume Problem: Optimal Algorithms and An Application to Carbon-Aware Load Shifting | ✓ |

The thumbnail shows the title page of the paper. At the top, there are three circular icons: 'Check for updates' (blue), 'Citefactor Evaluate' (red), 'Citefactor Available' (green), and 'Ultra Reproducible' (blue). Below the title, the authors' names are listed: Walid A. Hanafy, Roozbeh Bostandoost, Noman Bashir, David Irwin, Mohammad Hajiesmaili, and Prashant Shenoy. The university affiliation is University of Massachusetts Amherst, USA. The abstract section is titled 'ABSTRACT'.

21

The War of the Efficiencies: Understanding the Tension between Carbon and Energy Optimization

Walid A. Hanafy, Roozbeh Bostandoost, Noman Bashir, David Irwin, Mohammad Hajiesmaili, and Prashant Shenoy
University of Massachusetts Amherst, USA

ABSTRACT
Motivating the investigation of carbon-aware scheduling in the cloud, this paper identifies the tension between energy efficiency and carbon reduction. We propose a two-dimensional framework for understanding the trade-off between the two metrics. The first dimension is the temporal dimension, where we distinguish between shifting workloads and shifting resources. The second dimension is the spatial dimension, where we distinguish between shifting workloads and shifting resources. We show that shifting workloads is more effective than shifting resources for reducing carbon emissions. We also show that shifting workloads is more effective than shifting resources for reducing energy consumption. We further show that shifting workloads is more effective than shifting resources for reducing both carbon emissions and energy consumption. We also show that shifting workloads is more effective than shifting resources for reducing both carbon emissions and energy consumption. We also show that shifting workloads is more effective than shifting resources for reducing both carbon emissions and energy consumption.

Probleme in Prior Work

- Ein Beispielhafter Machine Learning Workload

RoBERTa

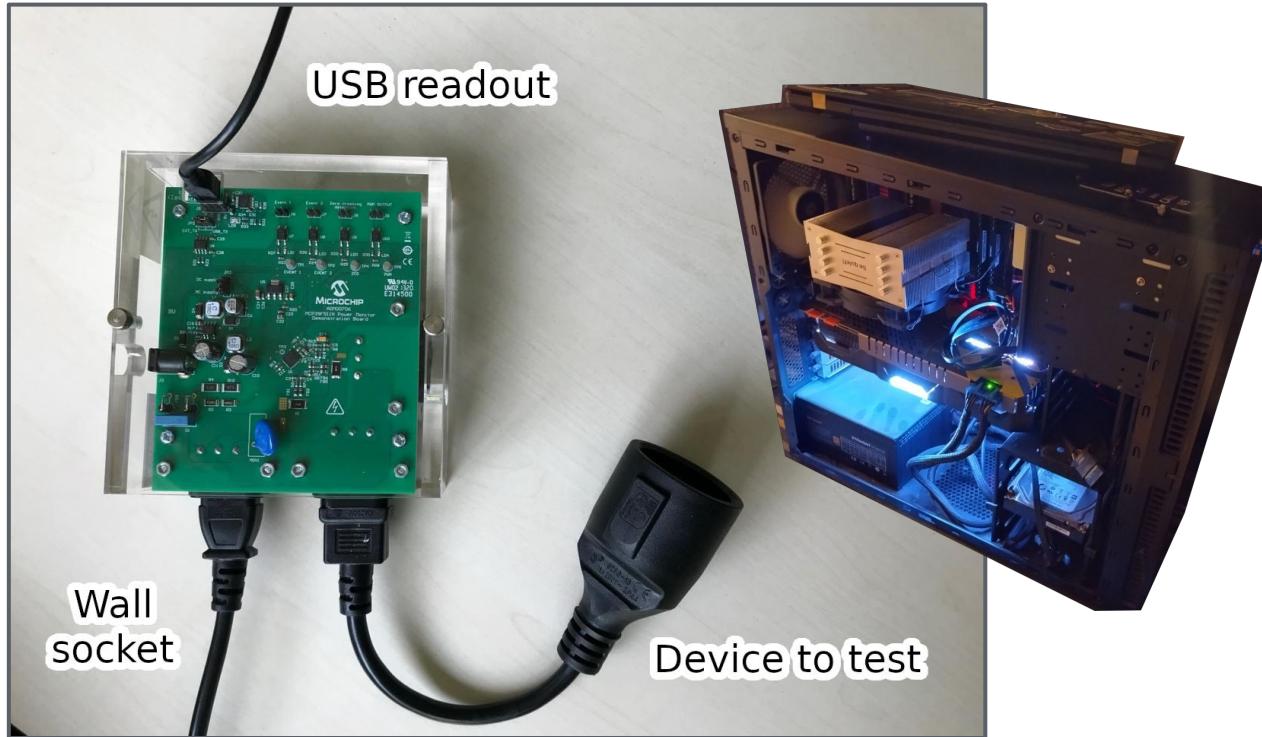
All model pages [roberta](#)  Hugging Face Spaces Paper page 1907.11692

Overview

The RoBERTa model was proposed in [RoBERTa: A Robustly Optimized BERT Pretraining Approach](#) by Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov. It is based on Google's BERT model released in 2018.

It builds on BERT and modifies key hyperparameters, removing the next-sentence pretraining objective and training with much larger mini-batches and learning rates.

Leistungsmessungen auf Machine Learning - Wie?



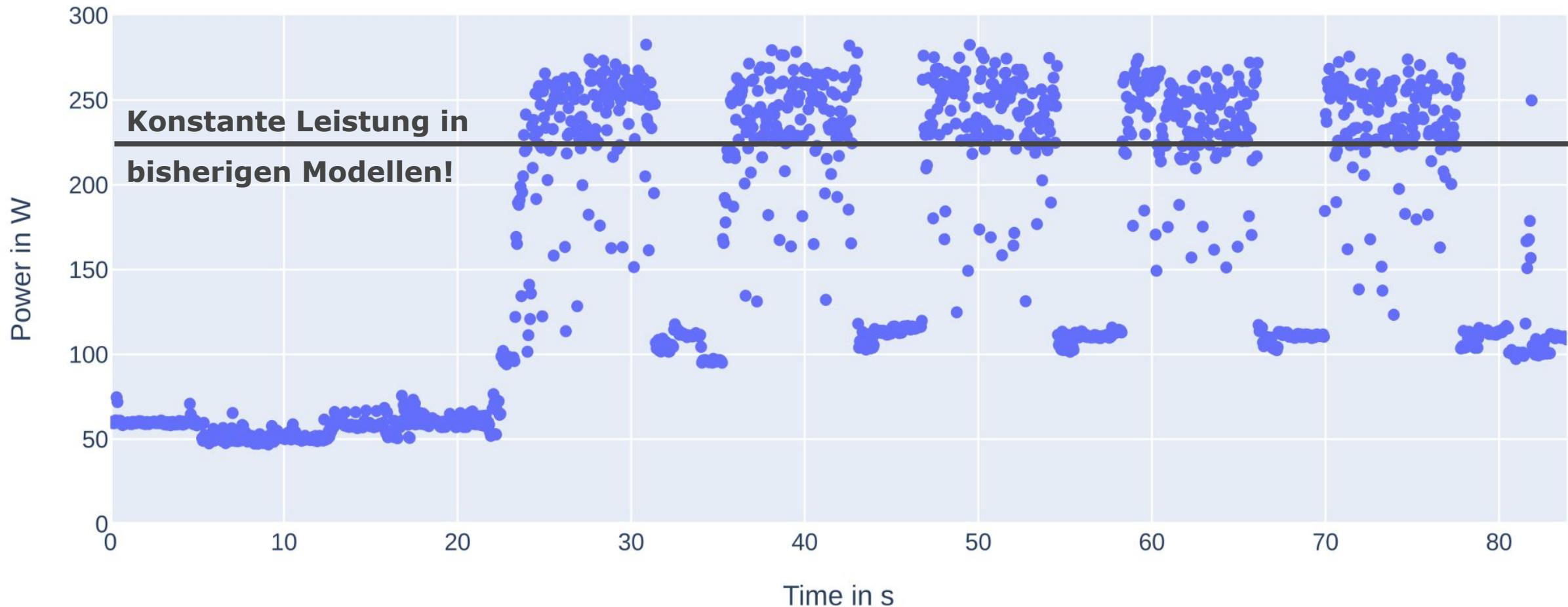
Microchip MCP39F511N Power Monitor

```
#!/bin/bash
```

```
# definiere Pfad zu Messprogramm pinpoint1
# setze CPU Frequenz
# setze Anzahl der Wiederholungen
# definiere Experimente
# Isoliere einzelne Runs voneinander
    # sleep
    # setze State zurück
echo "Done!"
```

¹ github.com/osmhpi/pinpoint

Leistungsmessungen von RoBERTa



State of the Art - Papers

| Paper | Suspend & Resume Scheduling | Covers Resume Overhead costs |
|--|-----------------------------|------------------------------|
| Let's Wait Awhile: How Temporal Workload Shifting Can Reduce Carbon Emissions in the Cloud | ✓ | X |
| On the Limitations of Carbon-Aware Temporal and Spatial Workload Shifting in the Cloud | ✓ | X |
| The War of the Efficiencies: Understanding the Tension between Carbon and Energy Optimization | ✓ | ~ |
| Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions | ✓ | X |
| The Online Pause and Resume Problem: Optimal Algorithms and An Application to Carbon-Aware Load Shifting | ✓ | ✓ |



Hasso
Plattner
Institut

Digital Engineering · Universität Potsdam



A Parameterizable Test Bed for Carbon Aware Job Scheduling

Ein parametrisierbares Testbed für kohlenstoffbewusste Jobplanung

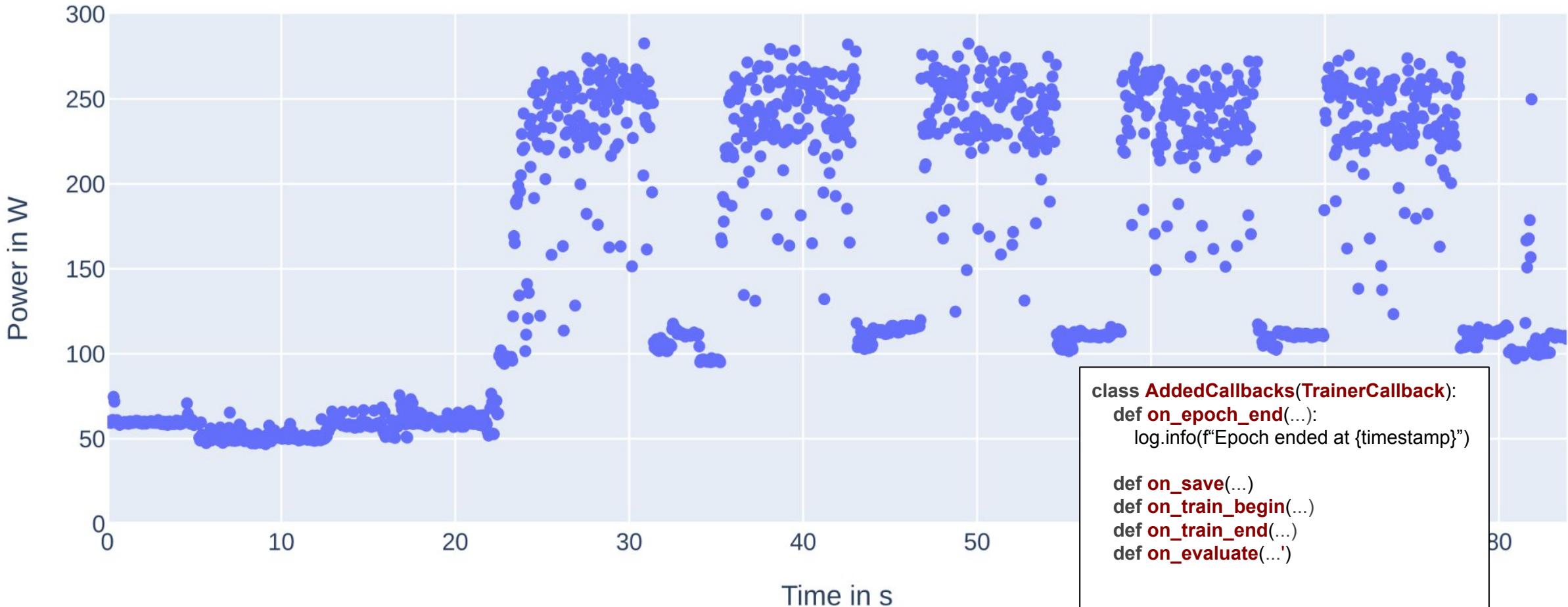
Vincent Opitz

Arbeit zur Erlangung des Grades “Master of Science” der Digital-Engineering-Fakultät der Universität Potsdam

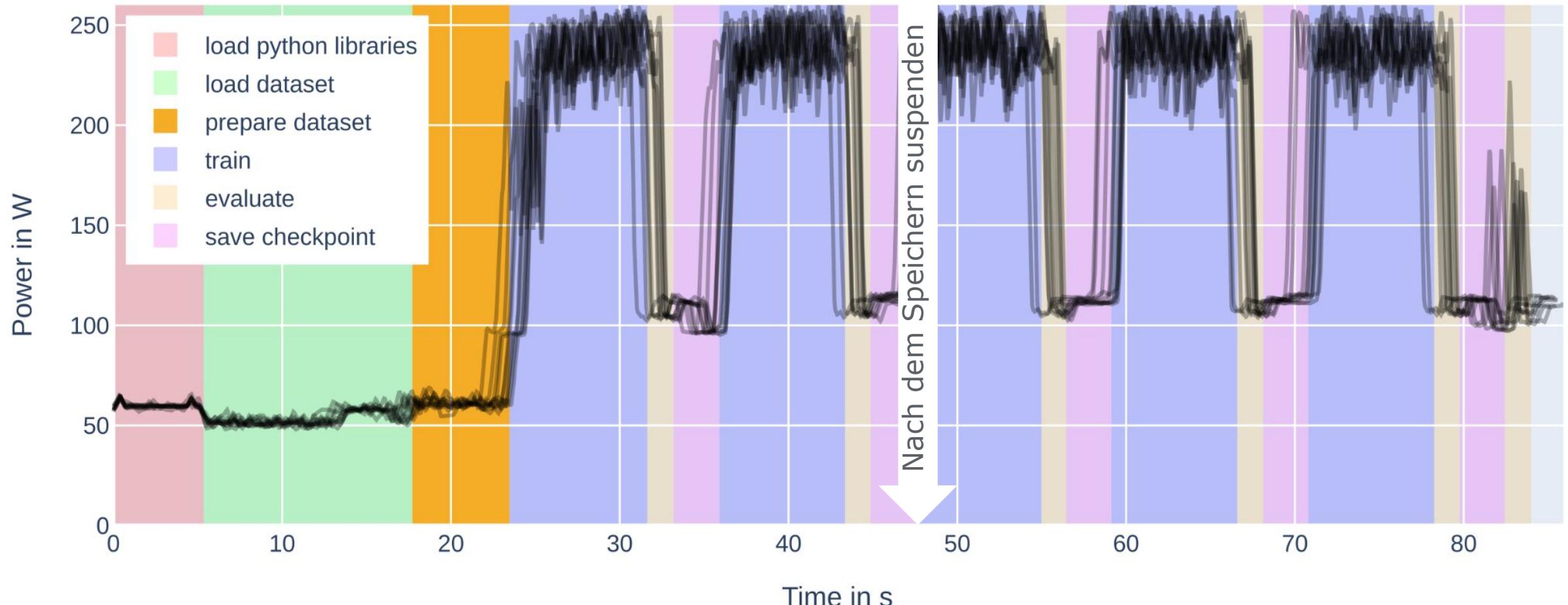
Contributions / Ablauf

- 1. Erstellen eines neuen Workload Modells**
welches Resumption Overhead und Phasen beinhaltet
- 2. Neue Scheduling Algorithmen**
darauf basierend implementieren in einem Testbed
- 3. Evaluation gegen das bestehende Modell**

Leistungsmessungen von RoBERTa



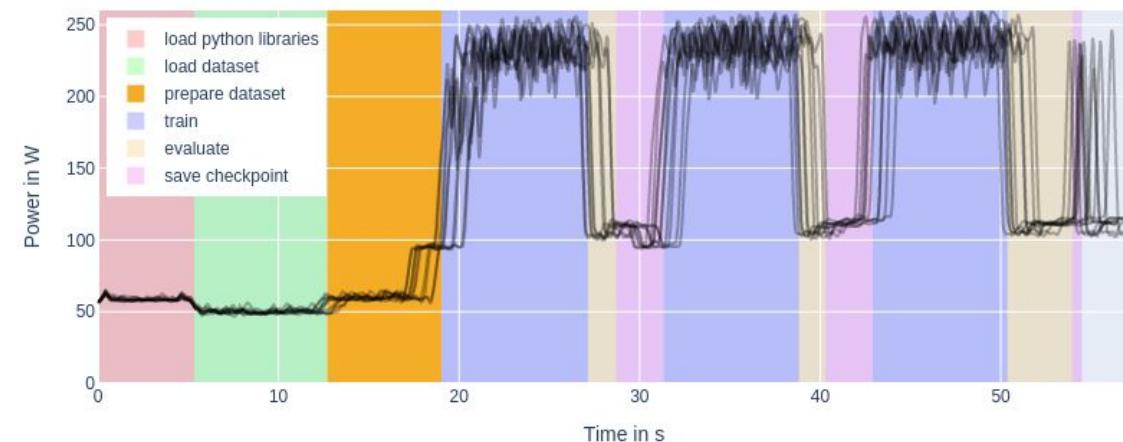
Leistungsmessungen von RoBERTa



Leistungsmessungen von RoBERTa

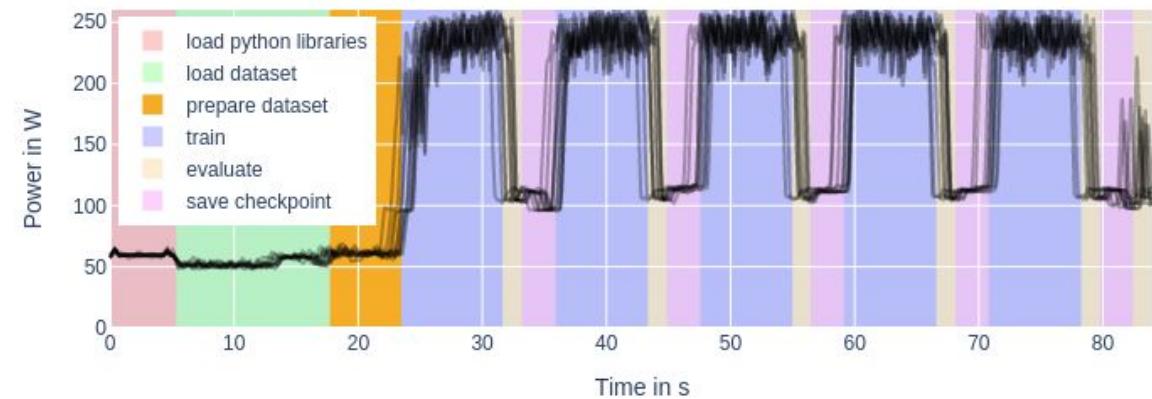
- Resumee nach der 2. Epoche

- Das Programm durchläuft erneut Startphasen
- Die verbleibenden 3 Epochen werden ausgeführt



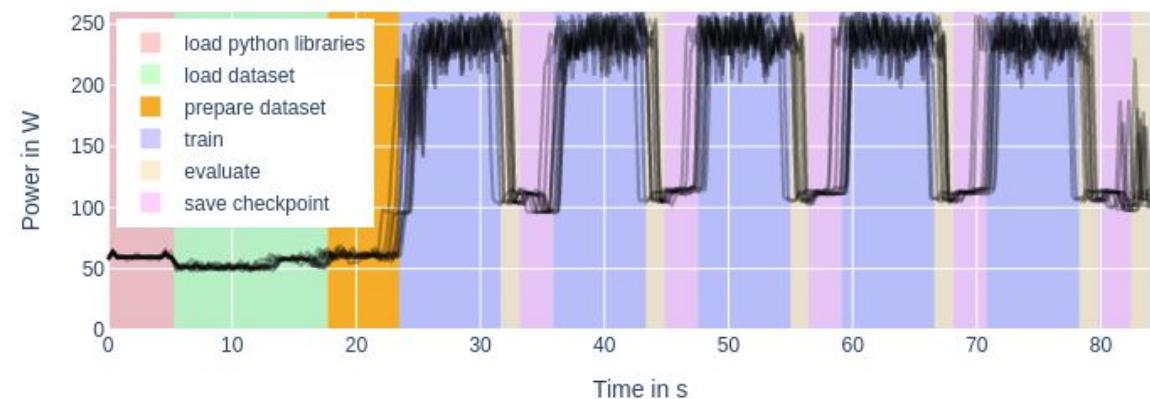
Leistungsmessungen auf Machine Learning - Zusammenfassung

- Es gibt (wiederkehrende) Phasen welche sich in ihrer Leistungsaufnahme unterscheiden
- Das Stoppen und Wiederaufnehmen von einem Job führt zu Overhead



Ein neues Workload Modell - STAWP

- Es gibt (wiederkehrende) Phasen welche sich in ihrer Leistungsaufnahme unterscheiden
- Das Stoppen und Wiederaufnehmen von einem Job führt zu Overhead



```
class Stawp(TypedDict):
    startup: List[Phase]
    work: List[Phase]
```

```
class Phase(TypedDict):
    name: str
    duration: float
    power: float
```

```
1.  {
2.  'startup': [
3.      {'name': 'load libraries', 'duration': 5.34, 'power': 59.9},
4.      {'name': 'load dataset', 'duration': 12.36, 'power': 53.77},
5.      {'name': 'prepare dataset', 'duration': 5.75, 'power': 63.17},
6.  ],
7.  'work': [
8.      {'name': 'train', 'duration': 8.17, 'power': 221.93},
9.      {'name': 'evaluate', 'duration': 1.54, 'power': 134.0},
10.     {'name': 'save', 'duration': 2.72, 'power': 105.1},
11.     ...
12. ]}
```

RoBERTa via STAWP modellieren

```

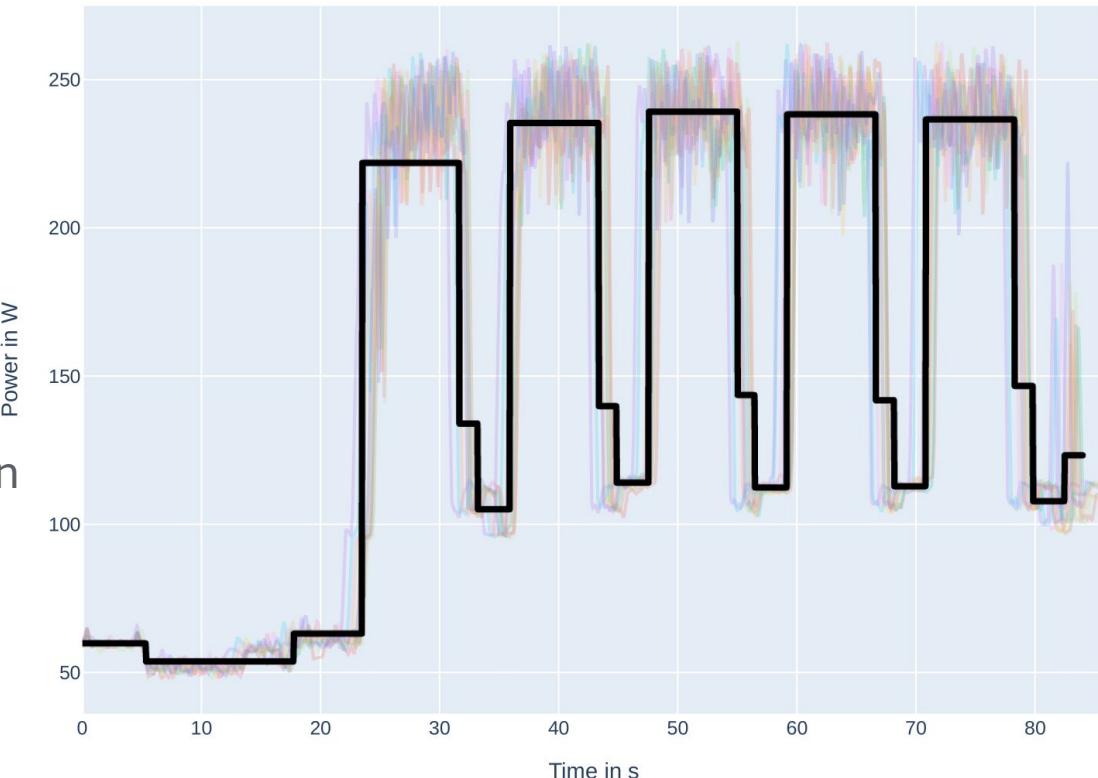
1.  {
2.    'startup': [
3.      {'name': 'load libraries', 'duration': 5.34, 'power': 59.9},
4.      {'name': 'load dataset', 'duration': 12.36, 'power': 53.77},
5.      {'name': 'prepare dataset', 'duration': 5.75, 'power': 63.17},
6.    ],
7.    'work': [
8.      {'name': 'train', 'duration': 8.17, 'power': 221.93},
9.      {'name': 'evaluate', 'duration': 1.54, 'power': 134.0},
10.     {'name': 'save', 'duration': 2.72, 'power': 105.1},
11.   ...
12. ]
}
```

```

class Stawp(TypedDict):
    startup: List[Phase]
    work: List[Phase]

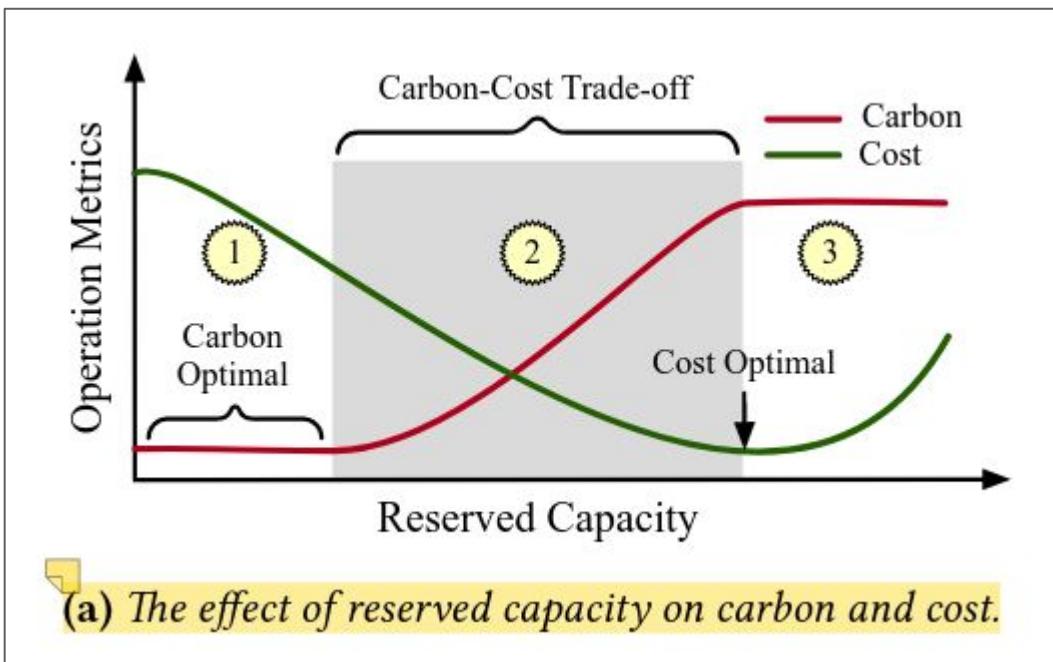
class Phase(TypedDict):
    name: str
    duration: float
    power: float

```



Implementation

- Hanafy et al. "Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions" (Mai 2024)



Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions

Walid A. Hanafy
University of Massachusetts Amherst USA

Qianlin Liang
University of Massachusetts Amherst USA

Noman Bashir
Massachusetts Institute of Technology USA

Abel Souza
University of Massachusetts Amherst USA

David Irwin
University of Massachusetts Amherst USA

Prashant Shenoy
University of Massachusetts Amherst USA

April 27-May 1, 2024, La Jolla, CA, USA. ACM, New York, NY, USA.
18 pages. <https://doi.org/10.1145/3620666.3651374>

1 Introduction
The increasing demand for computing and accelerating growth in cloud datacenter capacity have long raised concerns about their sustainability, environmental impact, and resulting carbon footprint [23]. Although cloud operators previously addressed such concerns by increasing datacenters' energy efficiency, or work done per unit of energy consumed, via software and hardware optimizations, recent work highlights that continuing improvements to energy efficiency are likely to see diminishing returns [12]. For example, large datacenters already operate near the optimal PUE value of 1.0, so there is little room to further improve PUE. Thus, to continue reducing their carbon emissions, cloud providers are increasingly employing *supply-side* optimizations, such as purchasing renewable energy from solar and wind farms to offset datacenter demand [22, 41]. However, eliminating all carbon emissions using supply-side techniques alone can be very expensive [6, 15].
Researchers have also proposed *demand-side* techniques to decrease computing's operational carbon emissions. These techniques utilize 1) visibility into grid energy's carbon intensity (in g CO₂/eq.kWh) and 2) application-level flexibility to *modulate* demand based on variations in energy's carbon intensity [6, 21, 31, 36, 44]. For example, prior work on Wait Awhile [44] and Ecosvisor [35] utilize batch workloads' *temporal flexibility* to optimize carbon by executing jobs when energy's carbon intensity is low and pausing execution when carbon intensity is high. Importantly, while state-of-the-art techniques consider carbon-performance trade-offs, they ignore the cost implications of carbon-aware optimizations. Specifically, carbon-aware scheduling exploits batch jobs' temporal flexibility to delay their execution until energy's carbon intensity is low. However, this delay also increases the completion times of batch jobs. Thus, carbon-aware scheduling exhibits carbon-performance trade-off: decreasing carbon emissions generally results in longer completion times.

In addition to the performance trade-off above, there is also a cost trade-off when using temporal shifting to optimize carbon emissions. This cost trade-off manifests in

ACM Concepts: • Computer systems organization → Cloud computing • Social and professional topics → Sustainability.

Keywords: Sustainable Computing, Cloud Computing

ACM Reference Format:
Walid A. Hanafy, Qianlin Liang, Noman Bashir, Abel Souza, David Irwin, and Prashant Shenoy. 2024. Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions. In *29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3 (ASPLOS '24)*.

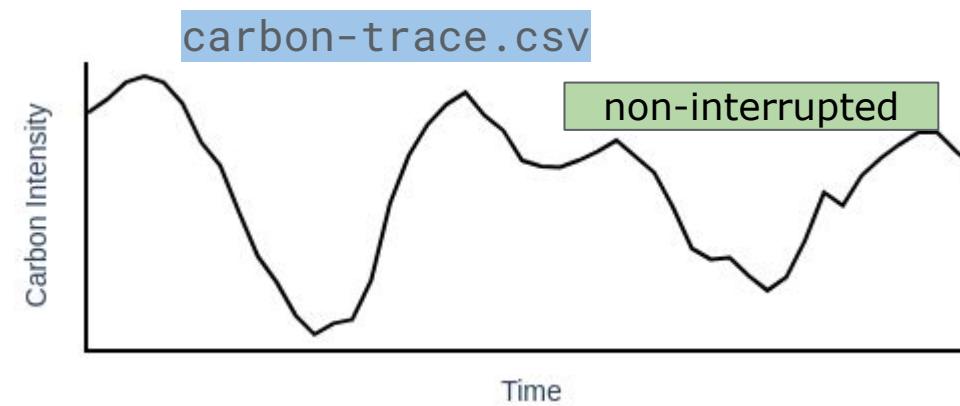
This work is licensed under a Creative Commons Attribution International 4.0 License.
ASPLOS '24, April 27-May 1, 2024, La Jolla, CA, USA
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0386-7/24/04.
<https://doi.org/10.1145/3620666.3651374>

Implementation - GAIA bisher

```
1. python3 gaia/run.py
2. --task-trace mytrace.csv
3. --carbon-trace carbon-trace.csv
4. --scheduling-policy gaia-non-interrupted
5. --carbon-policy oracle
6. --waiting-times 48
```

mytrace.csv

```
1. arrival_time, length
2. 0.0, 402166.0
3. 419.0, 263214.0
4. 113110.0, 166311.0
5. 113979.0, 221466.0
```

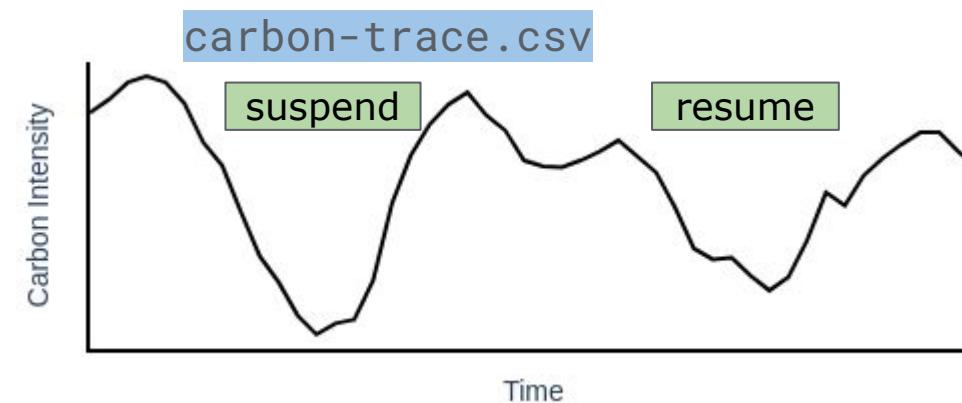


Implementation - GAIA bisher

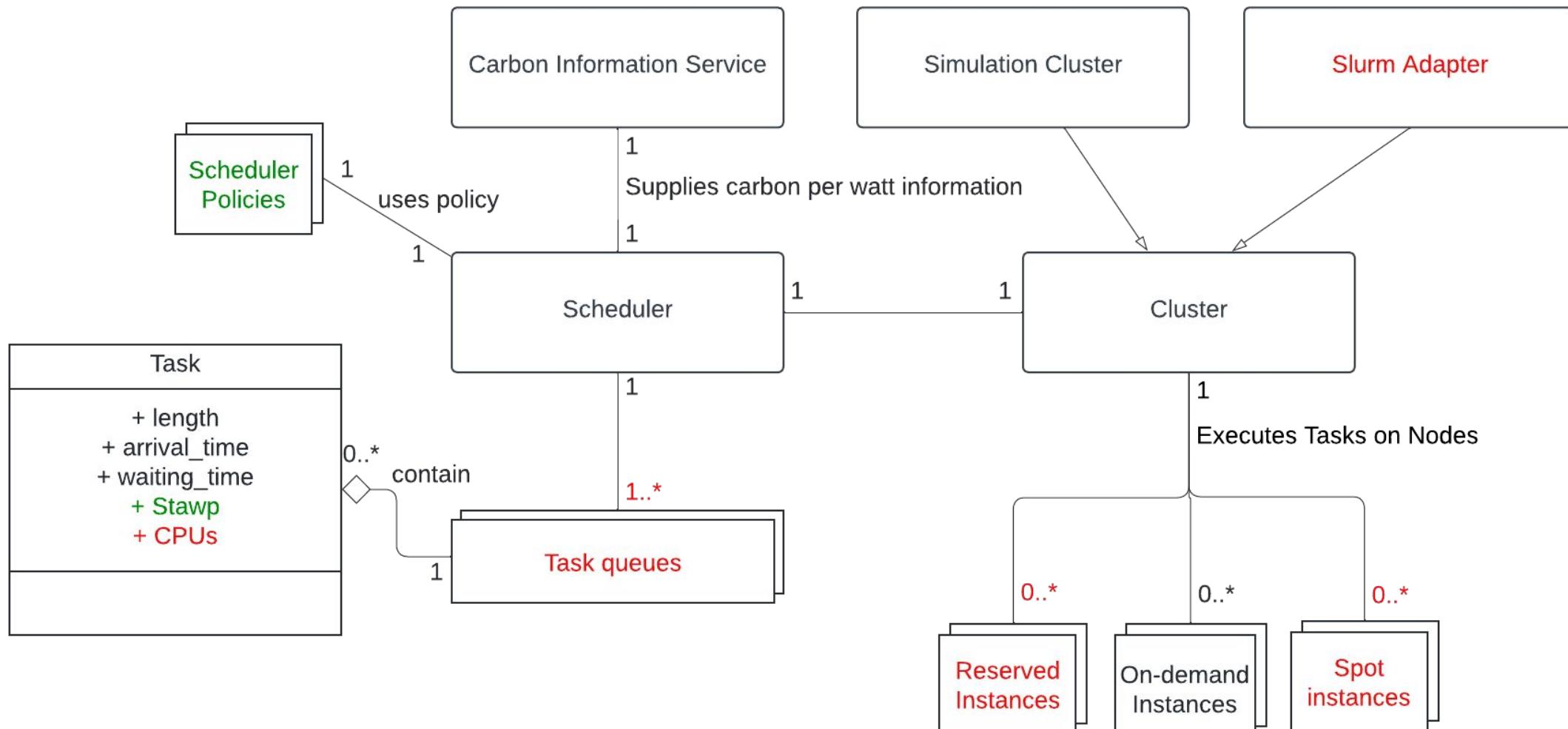
```
1. python3 gaia/run.py
2. --task-trace mytrace.csv
3. --carbon-trace carbon-trace.csv
4. --scheduling-policy gaia-suspend-resume
5. --carbon-policy oracle
6. --waiting-times 48
```

mytrace.csv

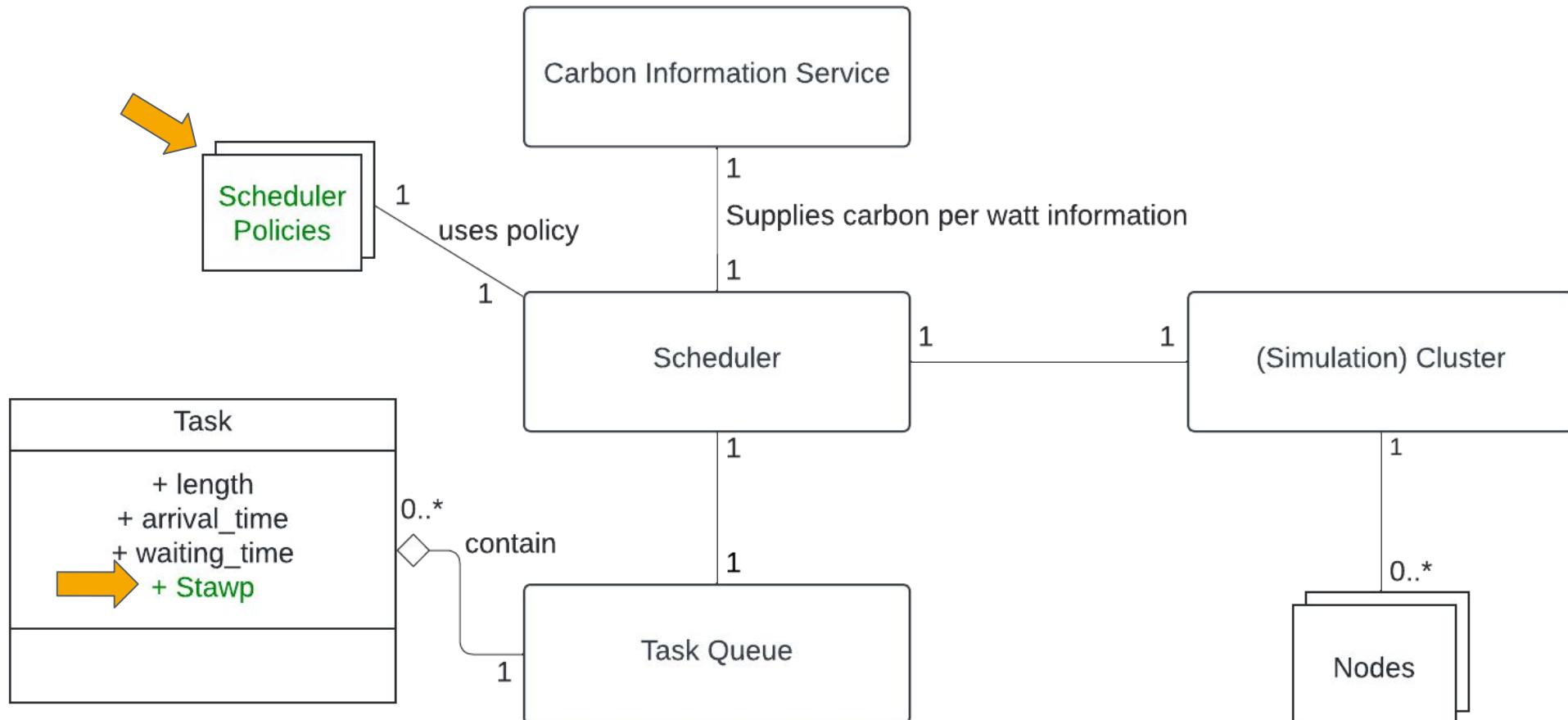
```
1. arrival_time,      length
2. 0.0,                402166.0
3. 419.0,              263214.0
4. 113110.0,           166311.0
5. 113979.0,            221466.0
```



Erweitern von GAIA zu meinem Testbed Carbs



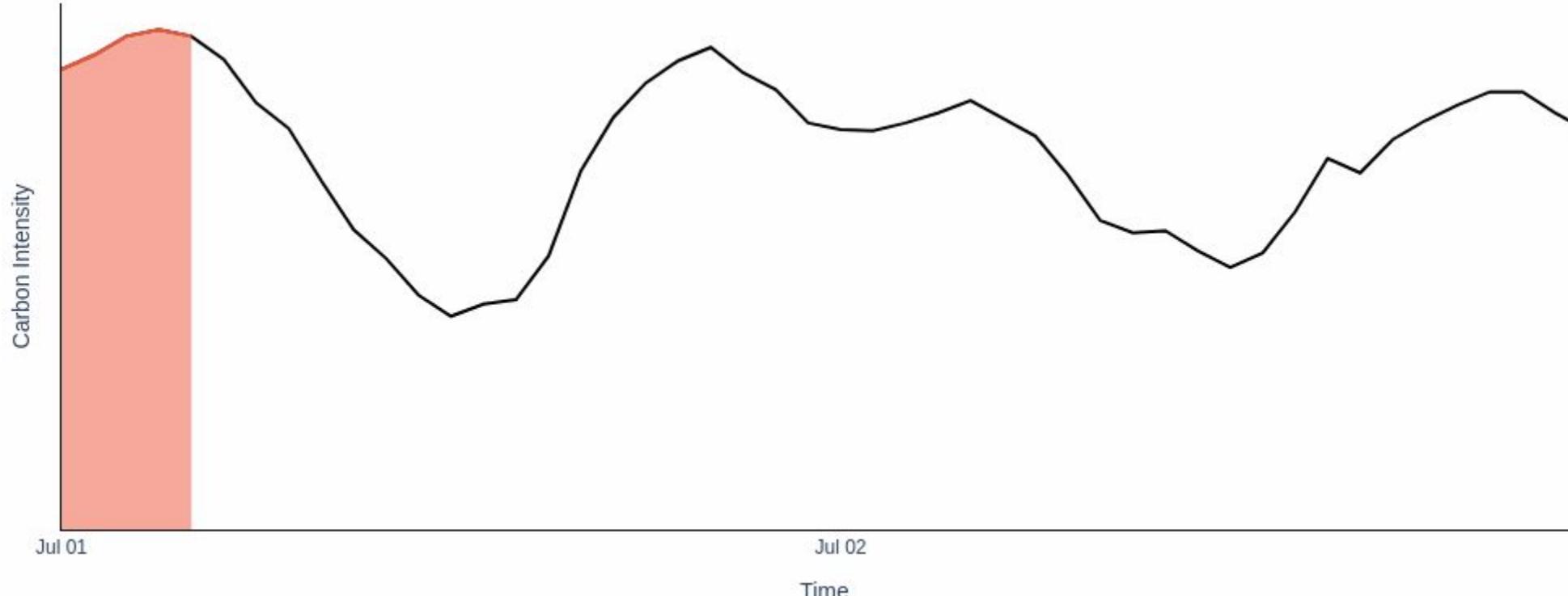
Erweitern von GAIA zu meinem Testbed Carbs



1. `python3 src/run.py`
2. `--scheduling-policy gaia-non-interrupted`
3. `--waiting-times 24`

Carbon-Aware Scheduling, ohne Suspend, in GAIA

$$\text{carbonEmissions} = \sum_{t = start}^{start + joblength} \text{carbonIntensity}(t) \cdot power$$



1. python3 src/run.py
2. --scheduling-policy carbs-non-interrupted
3. --waiting-times 24

Carbon-Aware Scheduling, ohne Suspend, in Carbs mit STAWP

$$\text{carbonEmissions} = \sum_{t = start}^{start + joblength} \text{carbonIntensity}(t) \cdot \text{stawp}(t)$$

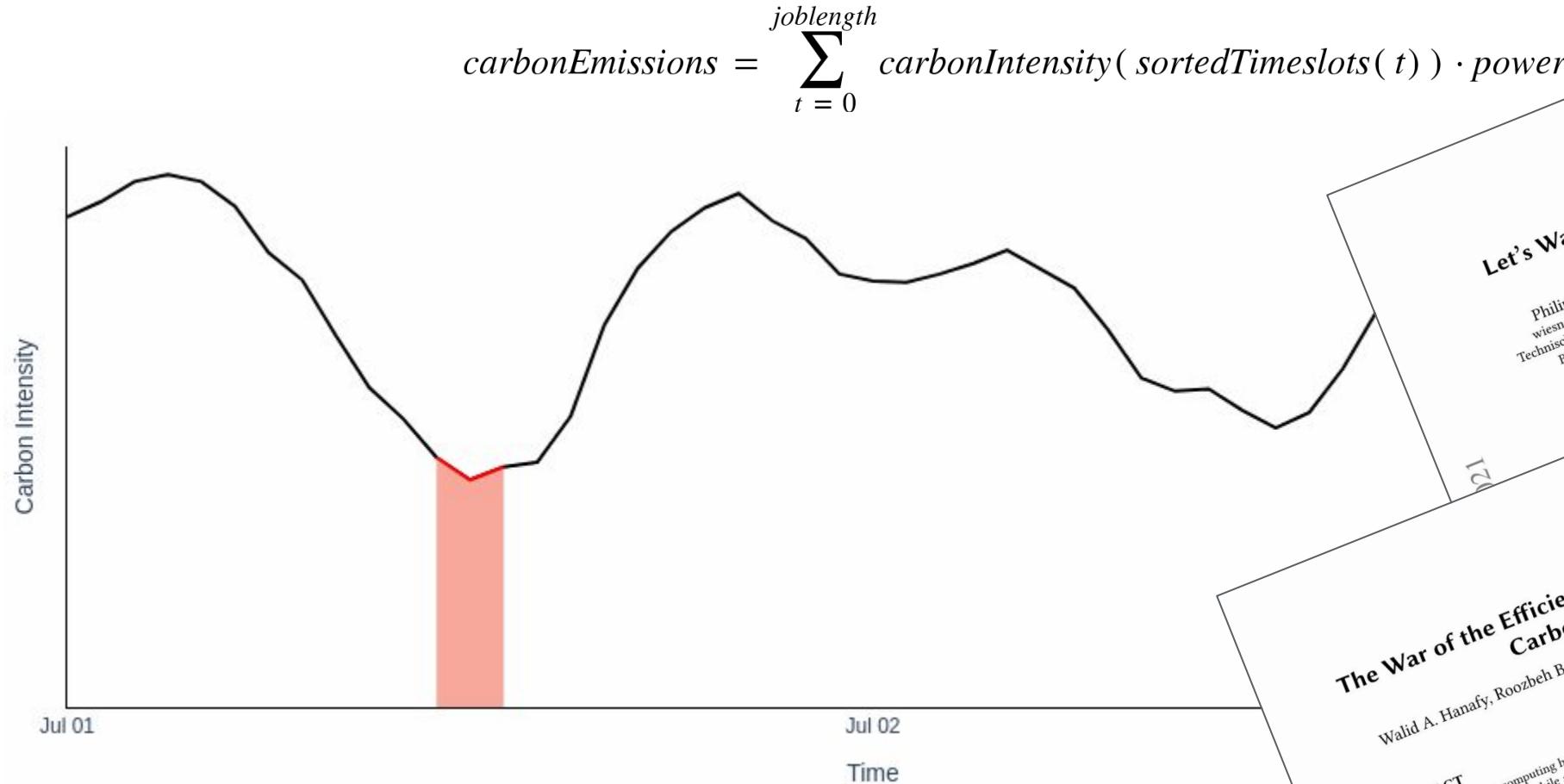
stawp(t) := Treppenfunktion aus

```
class Stawp(TypedDict):
    startup: List[Phase]
    work: List[Phase]
```

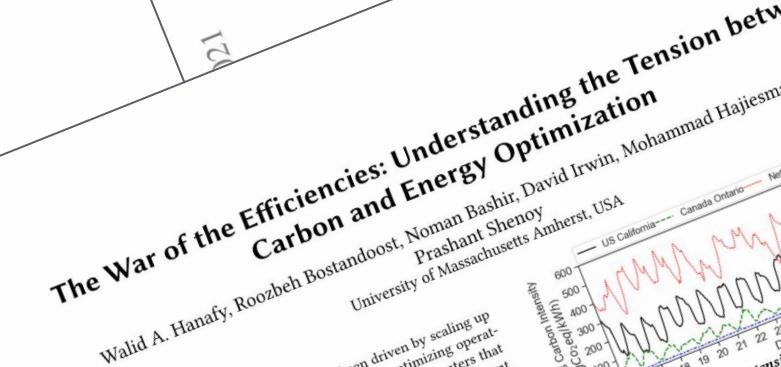
```
class Phase(TypedDict):
    name: str
    duration: float
    power: float
```

1. `python3 src/run.py`
2. `--scheduling-policy carbs-suspend-resume`
3. `--waiting-times 24`

Carbon-Aware Scheduling, mit Suspend, mit GAIA



21



ABSTRACT

Major innovations in computing have been driven by scaling up computing infrastructure, while aggressively optimizing operating costs. The result is a network of worldwide datacenters that consume a large amount of energy, mostly in an energy-inefficient manner. Since the electric grid powering these datacenters provided a simple, and opaque abstraction of an unlimited and reliable power supply, the computing industry remained largely oblivious to the carbon intensity of the electricity it uses. Much like the rest of the world, it generally treated the carbon intensity of the electricity as an afterthought – it was mostly true for a fossil fuel-driven grid. As the objective of increasing energy-efficiency – which has generally been the goal of the green IT approach. However, as the electricity cost per unit of energy – has generally been increasing, so has the cost of operation.

1 INTRODUCTION

The demand for computing has been expected to accelerate even further, as computing demand has not responded to the growth in demand so far [3]. The growth in computing demand can be kept in check by reducing the energy consumption per unit of work.

Philip Wiesner
wiesner@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Ilya Behnke
i.behnke@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Kordian Gontarska
kordian.gontarska@hpi.de
HPI, University of Potsdam
Potsdam, Germany

Lauritz Thamsen
lauritz.thamsen@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Let's Wait Awhile: How Temporal Workload Shifting Reduce Carbon Emissions in the Cloud

Philip Wiesner
wiesner@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Ilya Behnke
i.behnke@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Kordian Gontarska
kordian.gontarska@hpi.de
HPI, University of Potsdam
Potsdam, Germany

Lauritz Thamsen
lauritz.thamsen@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Lauritz Thamsen
lauritz.thamsen@tu-berlin.de
Technische Universität Berlin
Berlin, Germany

Carbon-Aware Scheduling, mit Suspend, in Carbs mit STAWP

Mit Linear Programming !

Linear programs are problems that can be expressed in [standard form](#) as

$$\begin{array}{ll} \text{Find a vector} & \mathbf{x} \\ \text{that maximizes} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & A\mathbf{x} \leq \mathbf{b} \\ \text{and} & \mathbf{x} \geq \mathbf{0}. \end{array}$$

Here the components of \mathbf{x} are the variables to be determined, \mathbf{c} and \mathbf{b} are given [vectors](#), and A is a given [matrix](#). The function whose value is to be maximized ($\mathbf{x} \mapsto \mathbf{c}^T \mathbf{x}$ in this case) is called the [objective function](#). The constraints $A\mathbf{x} \leq \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$ specify a [convex polytope](#) over which the objective function is to be optimized.

(Wikipedia)

Carbon-Aware Scheduling, mit Suspend, mit STAWP

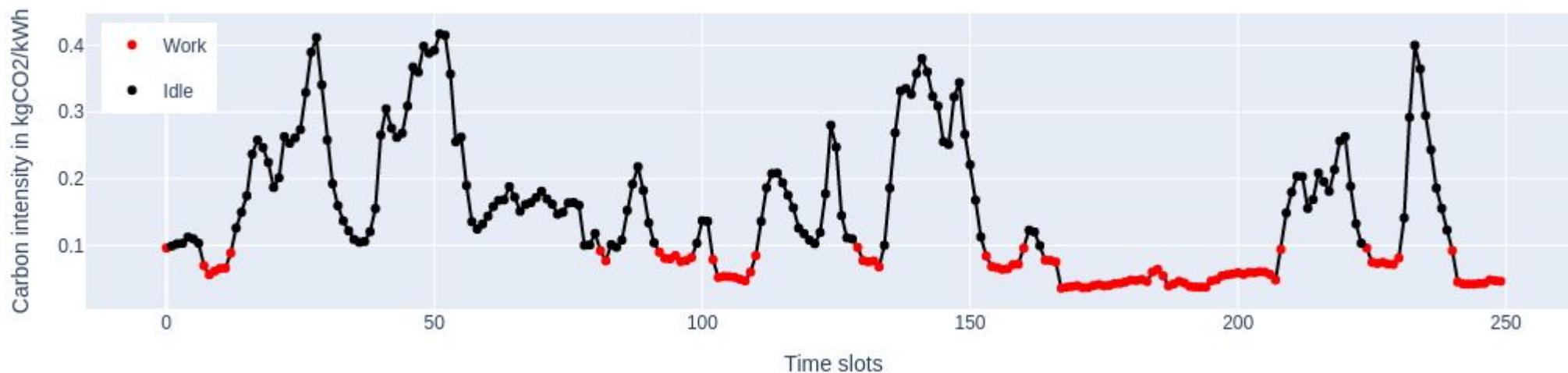
Step 1 - den Job ausführen

```

1. import pulp
2. prob = LpProblem("CarbonAwareScheduling", pulp.LpMinimize)
3. work = LpVariable.dicts("work", for t in range(DEADLINE), cat="Binary"))
4. prob += lpSum([work[t] * carbon_cost[t] for t in range(DEADLINE)])
5. prob += lpSum(work[t] for t in range(DEADLINE)) == WORK_LENGTH
6. solver = pulp.Gurobi_CMD(timeLimit)
7. prob.solve(solver)

```

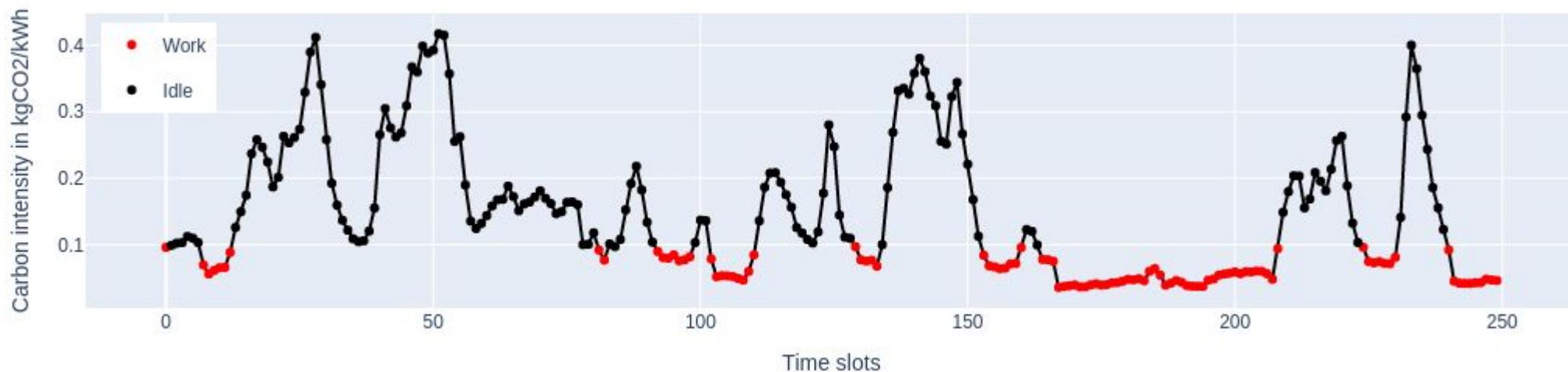
Find a vector \mathbf{x}
 that maximizes $\mathbf{c}^T \mathbf{x}$
 subject to $\mathbf{Ax} \leq \mathbf{b}$



Carbon-Aware Scheduling, mit Suspend, mit STAWP

Step 2 - Overhead bei resume

```
1. for t in range(DEADLINE - 1):
2.     prob += startup_finished[t] >= work[t + 1] - work[t]
3.     prob += startup_finished[t] + work[t] <= 1
```



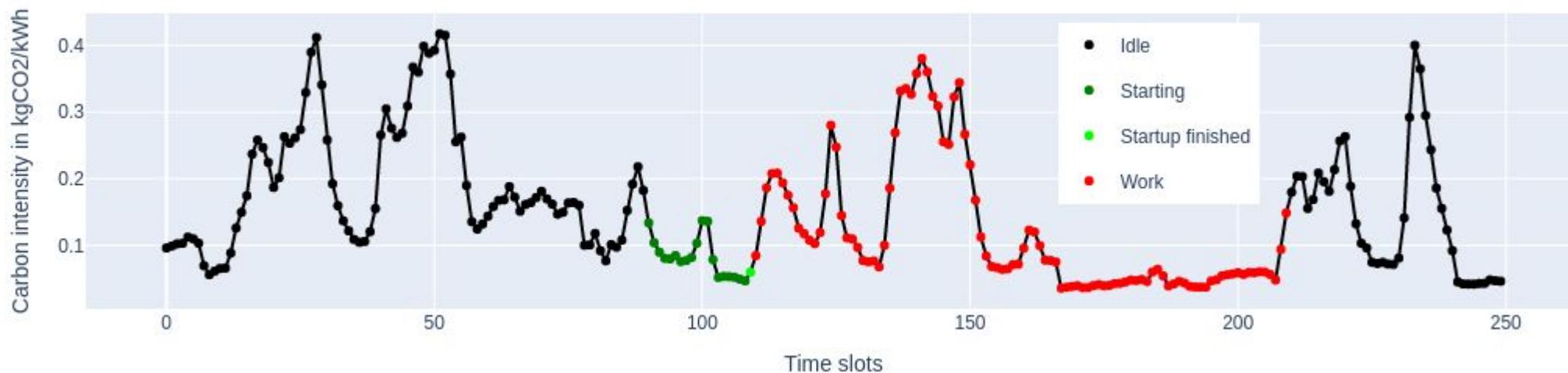
Carbon-Aware Scheduling, mit Suspend, mit STAWP

Step 2 - Overhead bei resume

```

1. for t in range(DEADLINE - 1):
2.     prob += startup_finished[t] >= work[t + 1] - work[t]
3.     prob += startup_finished[t] + work[t] <= 1
4.     prob += starting[t] + work[t] <= 1
5.
6. for i in range(STARTUP_LENGTH - 1, DEADLINE):
7.     prob += pulp.lpSum([starting[i - j] for j in range(STARTUP_LENGTH)]) >= STARTUP_LENGTH * startup_finished[i]

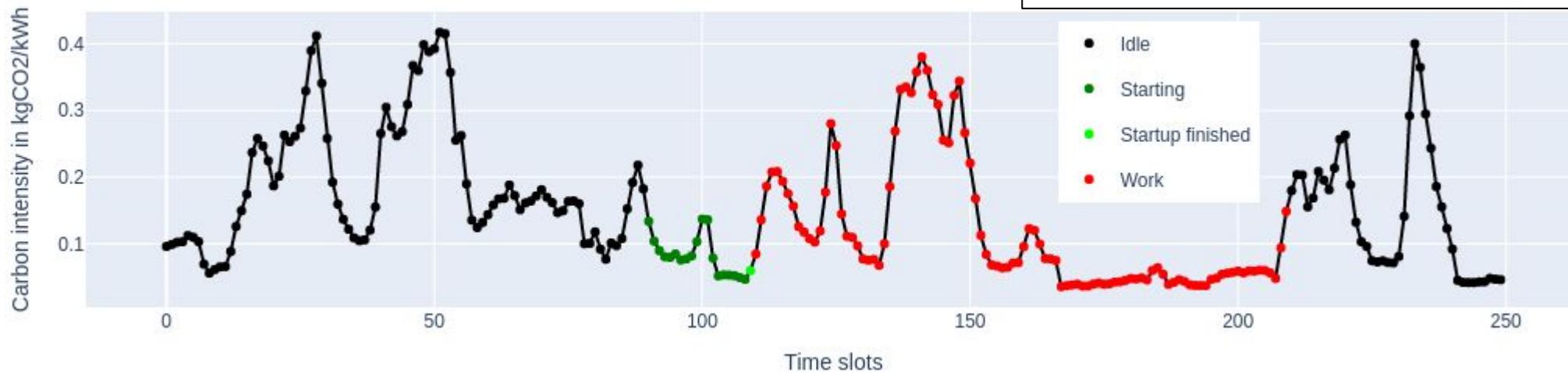
```



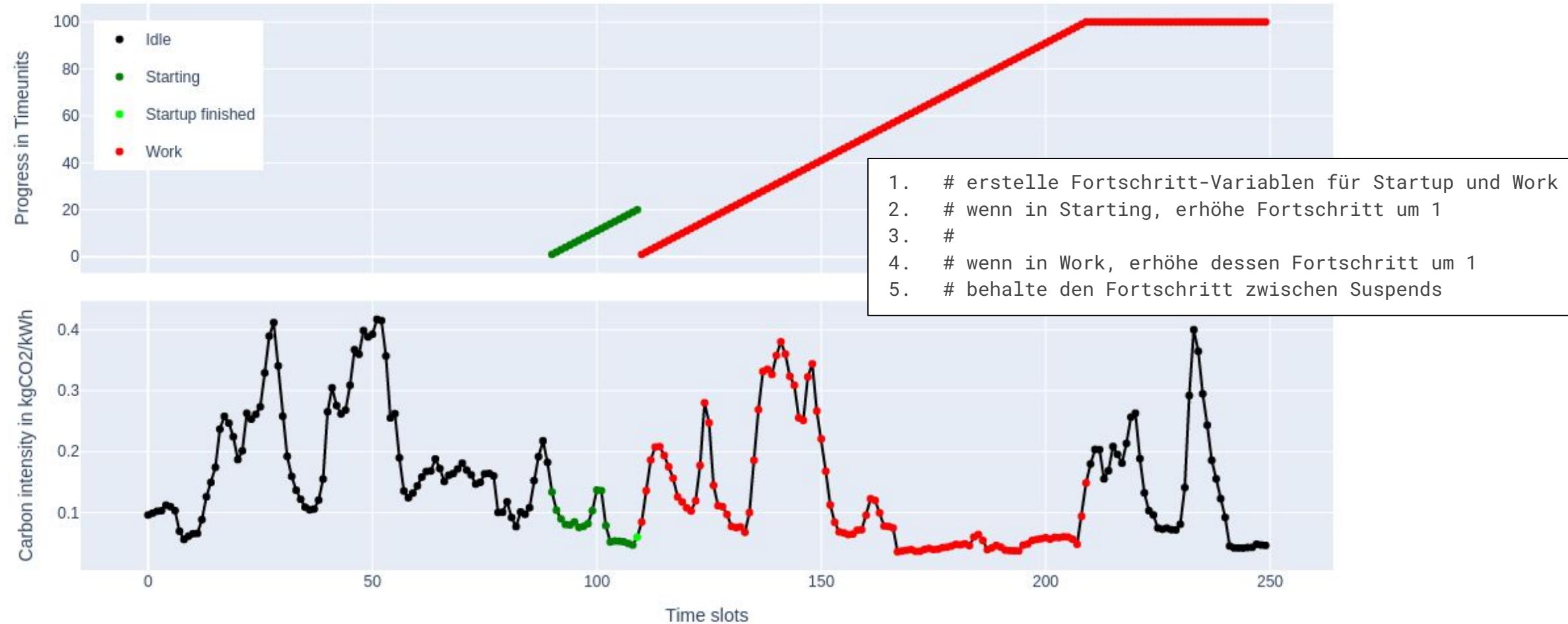
Carbon-Aware Scheduling, mit Suspend, mit STAWP

Step 3 - Zeit innerhalb des Jobs

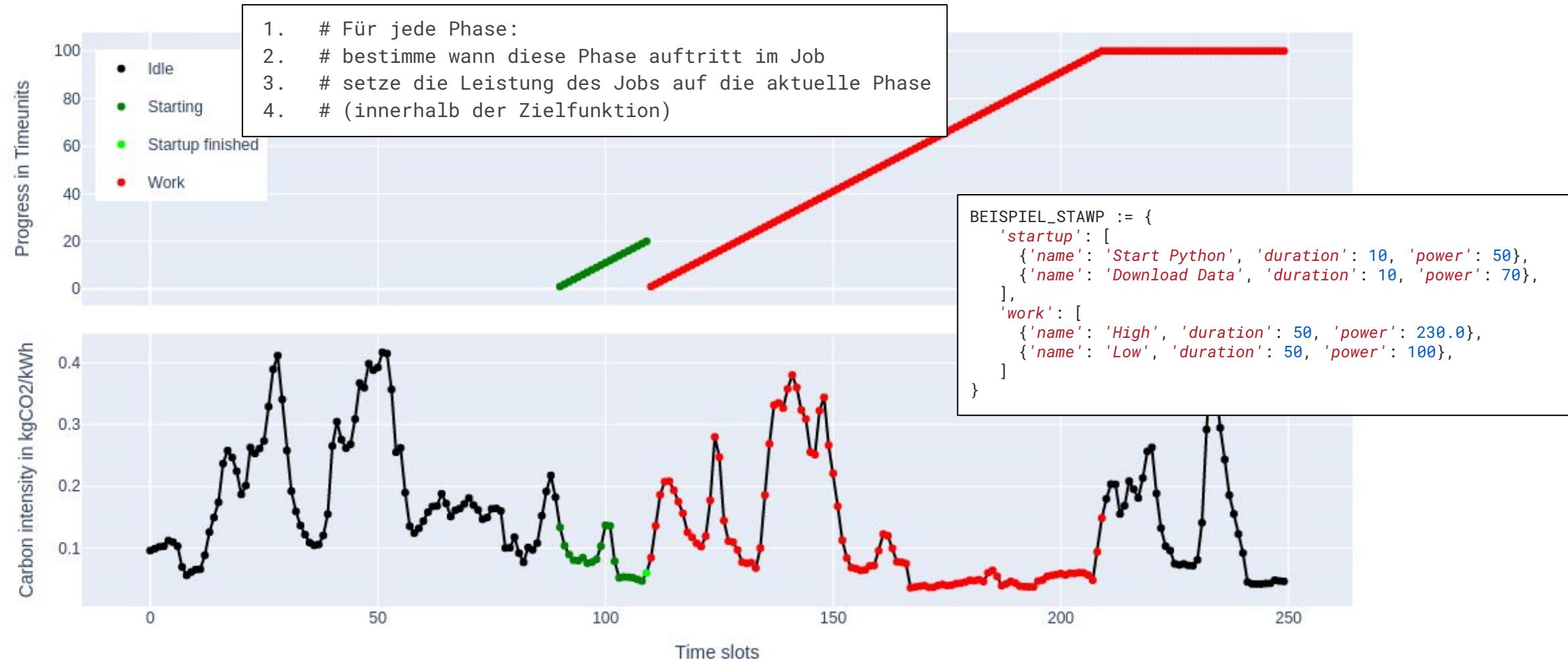
1. # erstelle Fortschritt-Variablen für Startup und Work
2. # wenn in Starting, erhöhe Fortschritt um 1
3. #
4. # wenn in Work, erhöhe dessen Fortschritt um 1
5. # behalte den Fortschritt zwischen Suspends



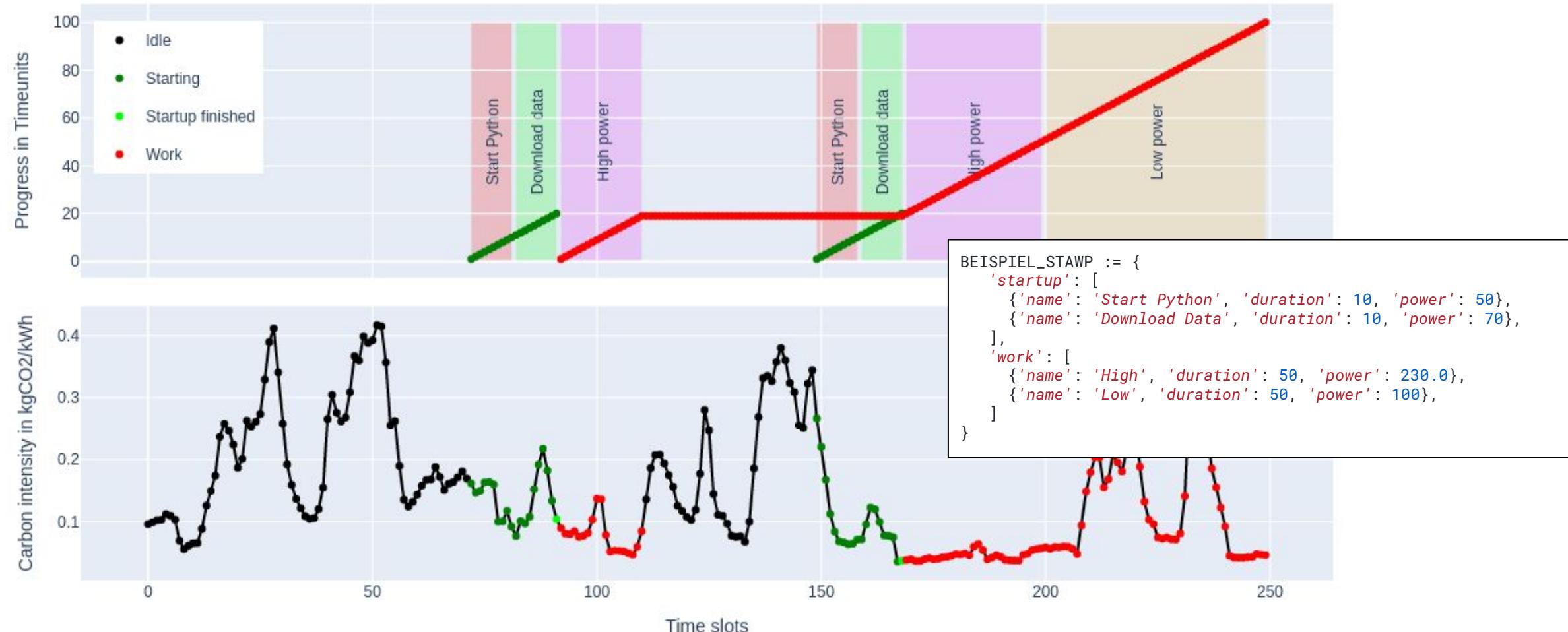
Step 3 - Zeit innerhalb des Jobs



Step 4 - Phasen abh. von der Zeit



Step 4 - Phasen abh. von der Zeit



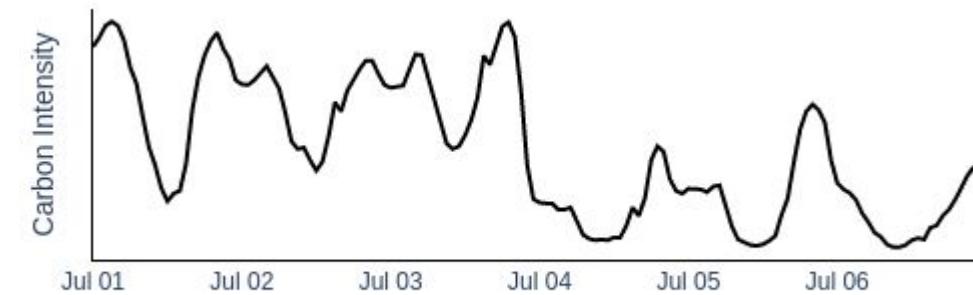
Evaluation - Setup

- Vergleich von STAWP gegen das bisherige State-of-the-Art Modell

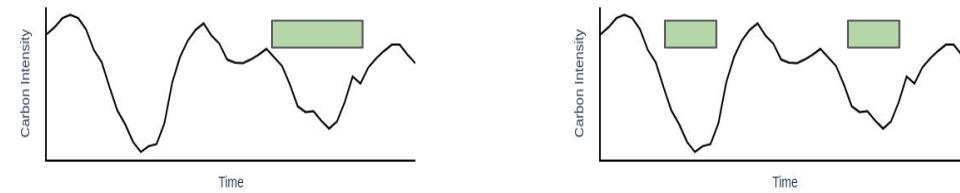
| | |
|------------------|-------------------------------------|
| Startup Länge | kein Startup, 5, 10, und 30 minuten |
| Startup Leistung | 100 W, 200 W |
| Job Länge | 1, 2, 4 , 8, und 16 Stunden |
| Phasen Länge | kurz oder lang |
| Benutze Phasen ? | ja oder nein |



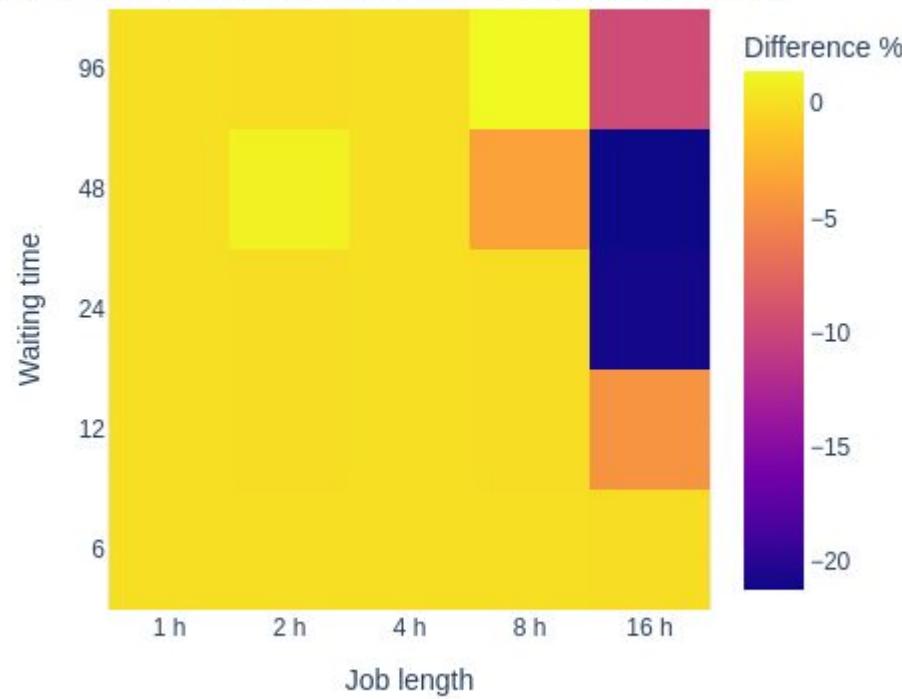
| | |
|----------------|------------------------------|
| Zeitspanne | Erste Juli Woche 2024 |
| max. Wartezeit | 4h, 12h, 1, 2, und 4 Tage |
| Scheduler | suspend & resume, oder nicht |



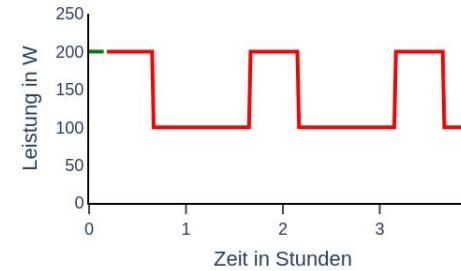
Evaluation - Vergleich von



Difference in Emissions between
Non-interruptable and suspend-resume Scheduling



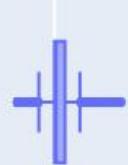
Evaluation - Vergleich von



vs.



non-interrupted



Explored 173519 nodes (14943839 simplex iterations) in 1200.31 seconds (501.96 work units)
Thread count was 256 (of 256 available processors)

Solution count 10: 0.523778 0.523778 0.523778 ... 0.542611

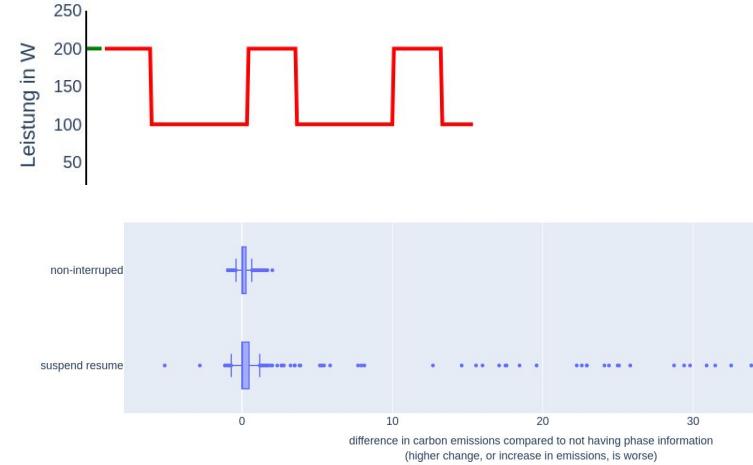
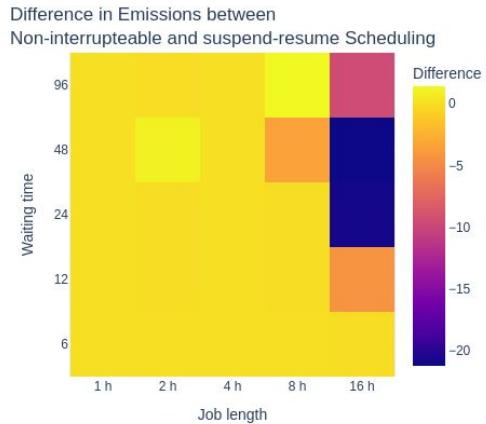
Time limit reached

suspend resume

0 10 20 30 40

difference in carbon emissions compared to not having phase information
(higher change, or increase in emissions, is worse)

Diskussion und Future Work



- Beobachtungen aus Prior Work traten auch auf.**
- Lange Jobs & Hohe Wartezeiten lohnen sich**
- Nur bis zu 16 Stunden und auf einem Carbon Trace untersucht**

Agenda

Workload Modell

Scheduler

Evaluation

Diskussion

- Effekte der Leistungsheterogenität nicht klar, Anzahl der Phasen skaliert schlecht**
- Hier müssen andere, echte Workloads untersucht werden**

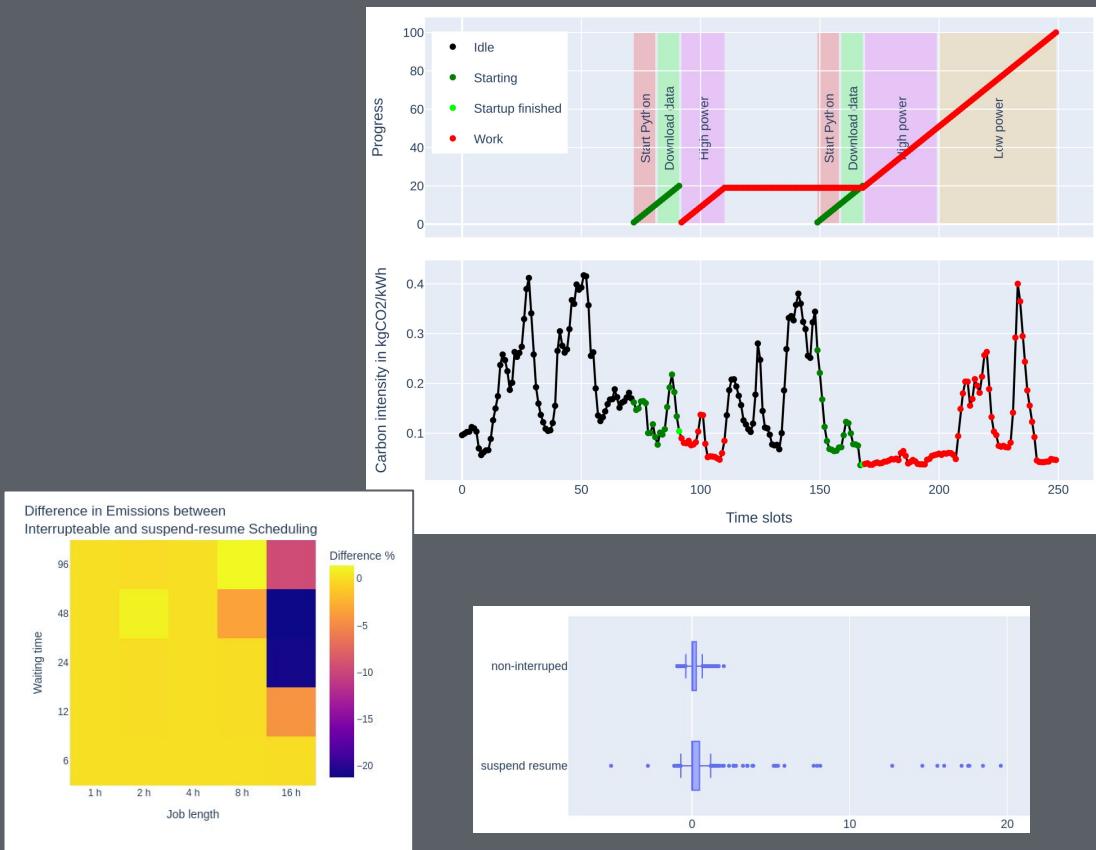
- Annahme, dass Jobs immer suspended werden können**
- Könnte Scheduling beeinflussen, oder für lange Jobs vernachlässigbar sein?**

```
class Stawp(TypedDict):
    startup: List[Phase]
    work: List[Phase]
```

```
class Phase(TypedDict):
    name: str
    duration: float
    power: float
```

```
class Stawp(TypedDict):
    startup: List[Phase]
    work: List[Phase]

class Phase(TypedDict):
    name: str
    duration: float
    power: float
```



Zusammenfassung

- 1. Stawp als neues Workload Modell vorgestellt**
enthält Startup (Overhead), Work, und Phasen
- 2. Neue Scheduling Algorithmen implementiert**
in meinem Testbed Carbs als Erweiterung von GAIA
- 3. Gegen das bestehende Modell evaluiert**

| | |
|------------------|--------------------------------------|
| Operating system | Ubuntu 24.04 LTS |
| Kernel | Linux 6.8.0-39-generic |
| CPU | AMD Ryzen 5 1600X Six-Core Processor |
| Memory | 16GiB |
| GPU | GeForce GTX 1070 |

Table 4.1: Environment parameters of the power measurements

Listing 4.1: Used operating system information

```
MINFREQ=$(cpupower frequency-info --hwlimits | sed -n '1d;p' \
    | awk '{print $1}')
MAXFREQ=$(cpupower frequency-info --hwlimits | sed -n '1d;p' \
    | awk '{print $1}')

cpupower frequency-set --min ${MAXFREQ} &>/dev/null
cpupower frequency-set --max ${MAXFREQ} &>/dev/null

# ... conduct experiments

cpupower frequency-set --min ${MINFREQ} &>/dev/null
cpupower frequency-set --max ${MAXFREQ} &>/dev/null
```

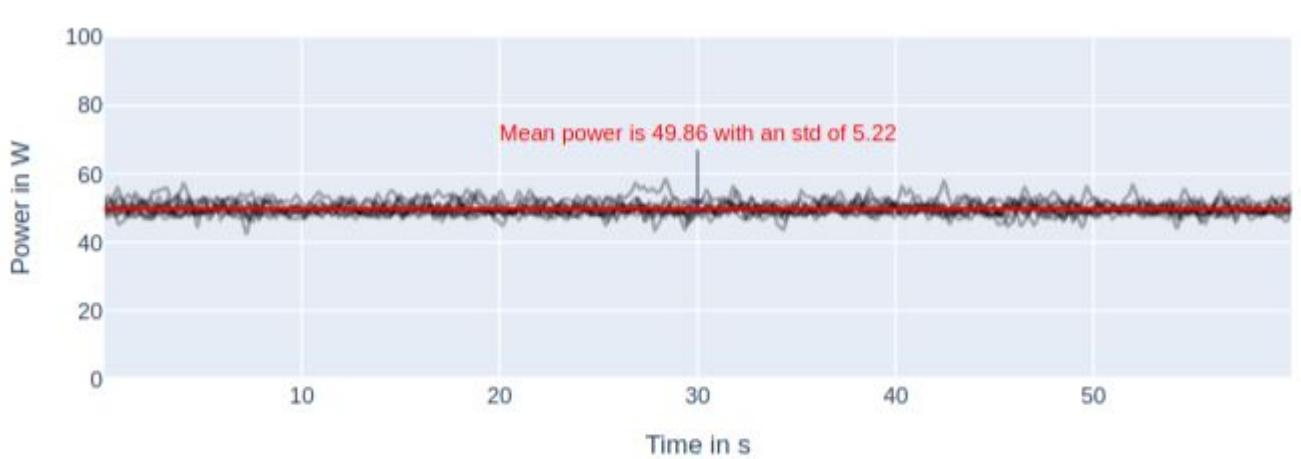


Figure 4.2: In the baseline experiment of measuring the system at rest, the average power draw is about 50 W. The black lines are Gaussian smoothed trend lines with $\sigma = 2$.

| Experiment | average energy cost | standard deviation |
|--|---------------------|--------------------|
| 1, whole run | 12.97 kJ | 0.04 kJ |
| 2, suspend after checkpoint and resume | 13.94 kJ | 0.1 kJ |
| 3, abort early and resume | 15.72 kJ | 0.07 kJ |

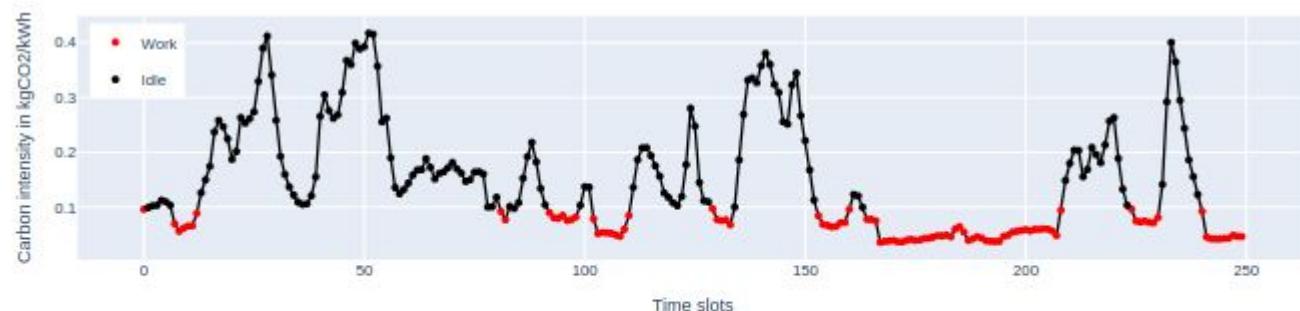
Table 4.2: Energy costs of the different experiments. Unlike the no-overhead assumption in prior work, we can show that suspending a workload does carry overhead.

Listing 6.5: LP Implementation for basic scheduling

```

1 prob = LpProblem("CarbonAwareScheduling", pulp.LpMinimize)
2 work = LpVariable.dicts("work",
3     (t for t in range(DEADLINE)), cat="Binary")
4
5 prob += lpSum([work[t] * carbon_cost[t] for t in range(DEADLINE)])
6 prob += lpSum(work[t] for t in range(DEADLINE)) == WORK_LENGTH
7
8 solver = pulp.Gurobi_CMD(timeLimit=timelimit)
9 prob.solve(solver)

```

**Figure 6.1:** Result using the first iteration of the LP scheduler. Each time slot, with its corresponding carbon intensity, is assigned to be worked or not. As this iteration has no notion of startup-costs, only the lowest time slots are used with many suspends in between.

Listing 6.6: LP Implementation for overhead

```
1 for t in range(DEADLINE - 1):
2     prob += startup_finished[t] >= work[t + 1] - work[t]
3     prob += startup_finished[t] + work[t] <= 1
4     prob += starting[t] + work[t] <= 1
5
6 for i in range(STARTUP_LENGTH - 1, DEADLINE):
7     prob += pulp.lpSum([starting[i - j] for j in range(STARTUP_LENGTH)])
8     >= STARTUP_LENGTH * startup_finished[i]
```

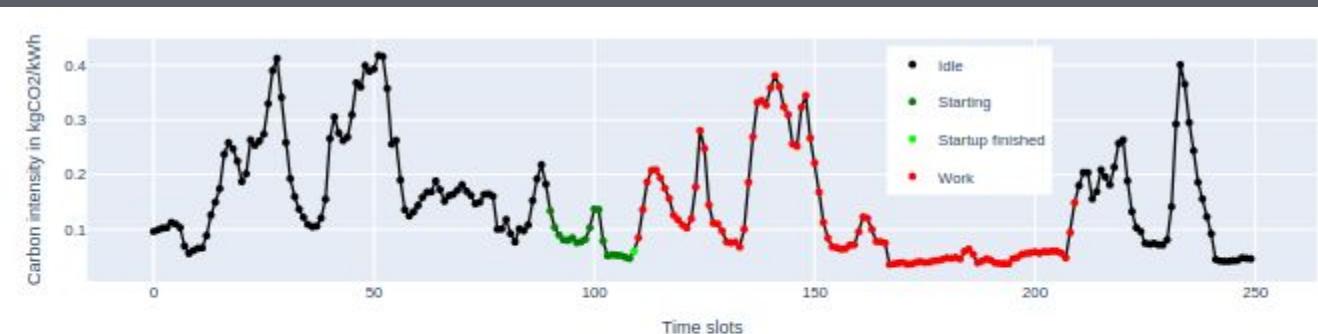


Figure 6.2: The schedule of the second iteration. In addition to time slots being assigned for work, a startup phase must also be scheduled before work can begin. The light-green mark indicates a helper variable, used for this mechanic. In comparison to the previous figure, there is now one long, connected, block where the job is worked on, as resuming a job carries an extra cost.

Listing 6.7: Progress Variables in LP

```
1 # define "work_progress" and "startup_progress" as
2 # DEADLINE-many integer variables
3
4 M = DEADLINE * 2
5 for t in range(DEADLINE-1):
6     if t > 0:
7         prob += startup_progress[t] >= startup_progress[t-1] + 1 - (1 -
8             starting[t]) * M
9         prob += startup_progress[t] <= startup_progress[t-1] + 1 + (1 -
10            starting[t]) * M
11        prob += startup_progress[t] <= starting[t] * M
12        prob += work_progress[0] == 0
13        if t > 0:
14            prob += work_progress[t] == work_progress[t-1] + work[t]
```

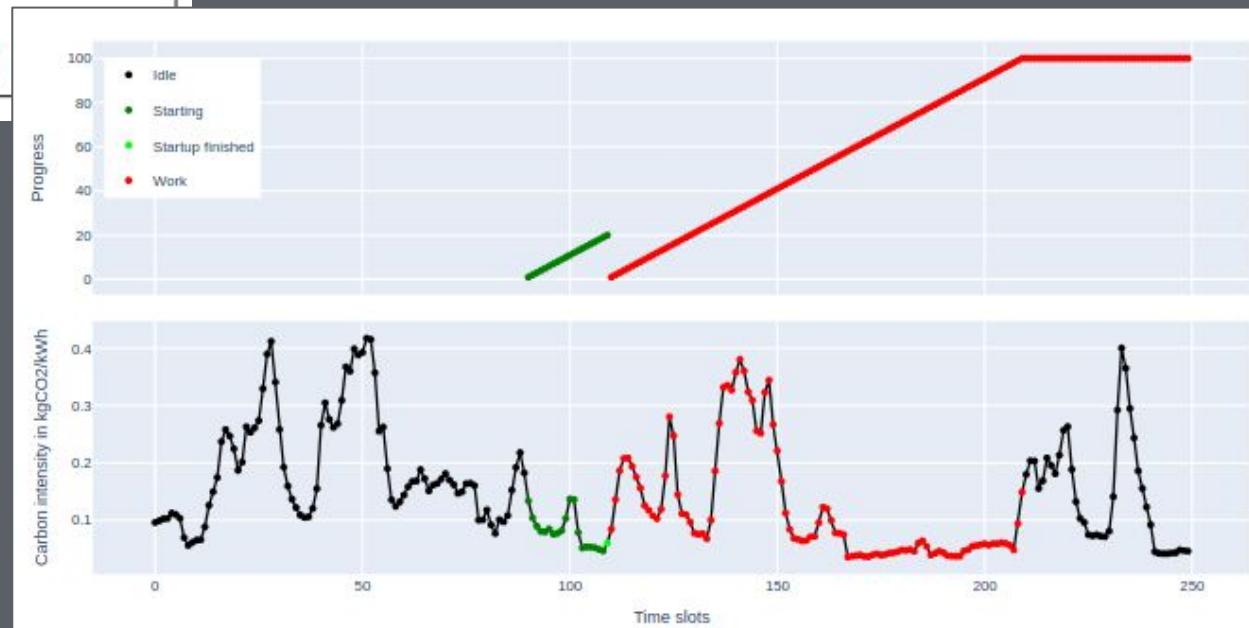


Figure 6.3: Visualisation of the result set after adding progress indicators. For either starting up or working, count up a respective integer variable. These do no yet impact the schedule itself, compared to Figure 6.2. Progress is kept even when idling.

Listing 6.8: Phase detection in LP

```

1 # for each phase
2 # let phase_indicator, phase_indicator_upper, phase_indicator_lower be
   DEADLINE-many boolean variables
3 # set lower_bound to be the minimum progress this phase can occur in
4 # set upper_bound to be the maximum similarly
5
6 # do the following for each phase
7 for t in range(DEADLINE):
8     prob += progress[t] - lower_bound <= M*phase_indicator_lower[t]
9     prob += lower_bound - progress[t] <= M*(1-phase_indicator_lower[t])
10
11    prob += upper_bound - progress[t] <= M*phase_indicator_upper[t]
12    prob += progress[t] - upper_bound <= M*(1-phase_indicator_upper[t])
13
14    prob += phase_active[t] >= phase_indicator_lower[t] +
15        phase_indicator_upper[t] - 1
16    prob += phase_active[t] <= phase_indicator_lower[t]
17    prob += phase_active[t] <= phase_indicator_upper[t]

```

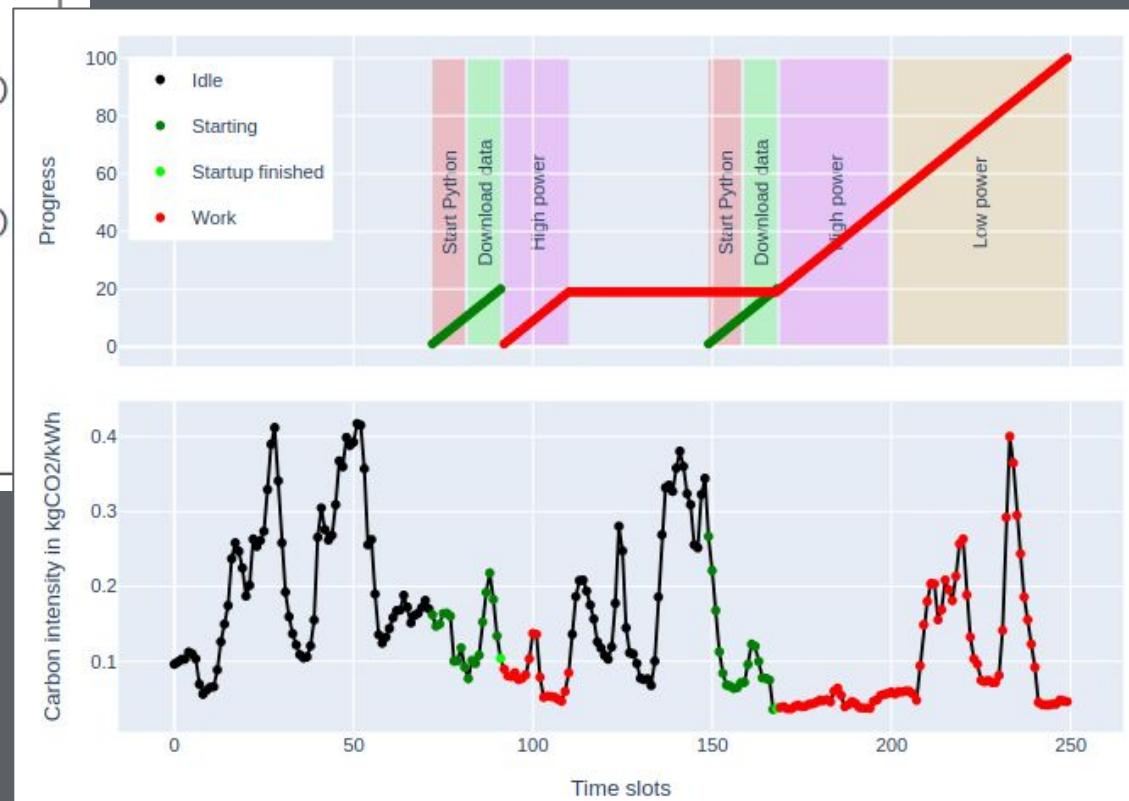


Figure 6.4: The final scheduling including differing power levels according to phases. In this example, there is a high-powered phase at 230 W and a low powered phase at 100 W. Unlike before, it is now optimal to suspend & resume the job to schedule the high-powered phase on the least carbon intensive time slots (170 to 200)

```
#!/bin/bash
export GUROBI_HOME="/home/vincent.opitz/master-thesis/gurobi1103/linux64"
export LD_LIBRARY_PATH=$GUROBI_HOME/lib
export PATH=$GUROBI_HOME/bin:$PATH
export GRB_LICENSE_FILE="/home/vincent.opitz/master-thesis/gurobi1103/gurobi.lic"
python3 src/run.py
  --scheduling-policy carbon \
  --carbon-trace DE-hourly-09-02-to-09-15 \
  --task-trace evaluation_jobs \
  --dynamic-power-draw \
  --dynamic-power-draw-type periodic-phases \
  --dynamic-power-draw-phases "[{'startup':[{'name': 'startup', 'duration': 0, 'power': 100}], 'work':[{'name': 'work', 'duration': 12000, 'power': 100}], 'idle': 0}]"
  --carbon-policy oracle \
  --start-index 0 \
  --w 12 \
  --filename results/simulation/evaluation_jobs/carbon_periodic-phases_0_0_100_12 \
  --repeat\
```

bash Script von job-1 der Evaluation

| arrival_time | length | cpus | Yield |
|--------------|---------|------|-------|
| 0 3600.0 1 | 12000.0 | 1 | 0 |
| 0 7200.0 1 | 12000.0 | 1 | 0 |
| 0 14400.0 1 | 12000.0 | 1 | 0 |
| 0 28800.0 1 | 12000.0 | 1 | 0 |
| 0 57600.0 1 | 12000.0 | 1 | 0 |

evaluation_jobs.csv

```
index=0
```

```
mkdir -p jobs  
mkdir -p results/simulation/evaluation_jobs  
  
rm -r jobs/job-*
```

```
# now we use a bunch of nested loops to create the cross product between these scenarios
```

```
for scheduling_policy in "${scheduling_policies[@]}"; do  
    for work_type in "${work_types[@]}"; do  
        for work_phase_index in "${!work_phases[@]}"; do  
            for startup_length in "${startup_lengths[@]}"; do  
                for startup_power_level in "${startup_power_levels[@]}"; do  
                    for waiting_time in "${waiting_times[@]}"; do  
                        ((index++))
```

```
                        work_phase=${work_phases[$work_phase_index]}
```

```
                        filename="results/simulation/evaluation_jobs/${scheduling_policy}_${work_type}_${work_phase_index}_${startup_length}_${scheduling_policies[@]}  
                        phases="{'startup':[{'name': 'startup', 'duration': $startup_length, 'power': $startup_power_level}], 'work':$work_phase}"
```

```
                        if [ ! -f "$filename" ]; then
```

```
                            echo Creating $filename
```

```
                            # create a .sh script with the job for later submission
```

```
                            echo '#!/bin/bash > jobs/job-${index}.sh'
```

```
                            echo 'export GUROBI_HOME="/home/vincent.opitz/master-thesis/gurobill03/linux64"' >> jobs/job-${index}.sh
```

```
                            echo 'export LD_LIBRARY_PATH=$GUROBI_HOME/lib' >> jobs/job-${index}.sh
```

```
                            echo 'export PATH=$GUROBI_HOME/bin:$PATH' >> jobs/job-${index}.sh
```

```
                            echo 'export GRB_LICENSE_FILE="/home/vincent.opitz/master-thesis/gurobill03/gurobi.lic"' >> jobs/job-${index}.sh
```

```
                            echo python3 src/run.py \
```

```
                                --scheduling-policy $scheduling_policy \
```

```
work_types=("periodic-phases" "constant-from-periodic-phases")  
work_phases=( \  
    # balanced  
    "[{'name': 'high', 'power': 200, 'duration': $long},  
    # long high  
    "[{'name': 'high', 'power': 200, 'duration': $long},  
    # short high  
    "[{'name': 'high', 'power': 200, 'duration': $short}  
)  
  
startup_lengths=(0 300 600 1800)  
startup_power_levels=(100 200)  
waiting_times=(6 12 24 48 96)  
scheduling_policies=("carbon" "suspend-resume")
```

Script für die Erstellung der Jobs der Evaluation

```
#!/bin/bash
```

```
chmod +x run_all_even_jobs.sh  
chmod +x run_all_odd_jobs.sh
```

```
# parallelization could be better, but the gurobi license only allows two parallel sessions
```

```
sbatch -A polze -p magic \  
    --container-image=python \  
    --container-name=test-2 \  
    --container-writable \  
    --mem=128G \  
    --cpus-per-task=128 \  
    --time=24:0:0 \  
    --comment="even" \  
    --output=slurmlogs/output_%j.txt \  
    --error=slurmlogs/error_%j.txt \  
    --constraint=ARCH:X86 \  
    --container-mounts=/hpi/fs00/home/vincent.opitz:/home/vincent.opitz \  
    --container-workdir=/home/vincent.opitz/master-thesis/GAIA run_all_even_jobs.sh
```

```
sbatch -A polze -p magic \  
    --container-image=python \  
    --container-name=test-2 \  
    --container-writable \  
    --mem=128G \  
    --cpus-per-task=128 \  
    --time=24:0:0 \  
    --comment="odd" \  
    --output=slurmlogs/output_%j.txt \  
    --error=slurmlogs/error_%j.txt \  
    --constraint=ARCH:X86 \  
    --container-mounts=/hpi/fs00/home/vincent.opitz:/home/vincent.opitz \  
    --container-workdir=/home/vincent.opitz/master-thesis/GAIA run_all_odd_jobs.sh |
```

```
#!/bin/bash
```

```
find jobs -name "job-*.*" | while read -r script; do  
    index=$(echo "$script" | grep -o -E '[0-9]+')  
  
    if (( index % 2 == 0 )); then  
        echo "Executing: $script"  
        ./"$script"  
    fi  
done
```

```
#!/bin/bash
```

```
find jobs -name "job-*.*" | while read -r script; do  
    index=$(echo "$script" | grep -o -E '[0-9]+')  
  
    if (( index % 2 == 1 )); then  
        echo "Executing: $script"  
        ./"$script"  
    fi  
done
```

Listing 7.1: Submitting the evaluation to SCORE Lab's Slurm

```
1 sbatch -A polze -p magic \
2   --container-image=python \
3   --container-name=test \
4   --container-writable \
5   --mem=128G \
6   --cpus-per-task=128 \
7   --time=24:0:0 \
8   --output=slurmlogs/output_%j.txt \
9   --error=slurmlogs/error_%j.txt \
10  --constraint=ARCH:X86 \
11  --container-mounts=/hpi/fs00/home/vincent.opitz:/home/vincent.opitz \
12  --container-workdir=/home/vincent.opitz/master-thesis/carbs \
13  run_all_even_jobs.sh
```

| Parameter | <i>p</i> value |
|------------------------|----------------|
| Job length | 0.00 |
| Waiting time | 0.00 |
| Phases | 0.0001 |
| Scheduling Algorithm | 0.18 |
| Startup power level | 0.66 |
| Phases averaged or not | 0.74 |
| Startup length | 0.99 |

Table 7.2: Results of the ANOVA for total carbon emissions per scheduler. The lower the p-value, the higher the impact on carbon emissions.

```

carbs > slurmlogs > output_1215025.txt
 501 H280997 76721          0.5237778  0.43812 16.4% 214 1084s
 503 287337 76914 cutoff 56   0.52378  0.44016 16.0% 212 1091s
 504 292331 76475 0.48840 59 435   0.52378  0.44121 15.8% 210 1095s
 505 305445 74293 0.51606 63 364   0.52378  0.44432 15.2% 205 1102s
 506 311380 73262 0.51210 55 592   0.52378  0.44604 14.8% 203 1107s
 507 318109 72136 0.51135 61 520   0.52378  0.44797 14.5% 201 1110s
 508 331886 69693 0.47927 52 832   0.52378  0.45052 14.0% 197 1119s
 509 338987 68893 cutoff 59   0.52378  0.45179 13.7% 195 1126s
 510 350442 65975 cutoff 56   0.52378  0.45433 13.3% 192 1131s
 511 358194 64206 cutoff 54   0.52378  0.45554 13.0% 190 1136s
 512 366285 62221 0.50677 55 589   0.52378  0.45736 12.7% 188 1141s
 513 375823 59977 cutoff 60   0.52378  0.45799 12.6% 186 1146s
 514 384769 57922 cutoff 61   0.52378  0.45990 12.2% 184 1151s
 515 393709 55782 cutoff 60   0.52378  0.46131 11.9% 182 1156s
 516 402240 53577 cutoff 58   0.52378  0.46275 11.7% 180 1162s
 517 412712 53450 cutoff 57   0.52378  0.46387 11.4% 178 1168s
 518 422946 54350 0.51404 54 423   0.52378  0.46437 11.3% 176 1177s
 519 437245 54955 cutoff 57   0.52378  0.46712 10.8% 174 1184s
 520 448501 55311 cutoff 59   0.52378  0.46861 10.5% 172 1190s
 521 460202 55356 0.50235 64 409   0.52378  0.46962 10.3% 171 1197s
 522 472395 55240 0.50530 56 500   0.52378  0.47130 10.0% 169 1200s
 523
 524 Cutting planes:
 525 Learned: 478
 526 Gomory: 22
 527 Cover: 8048
 528 Implied bound: 80
 529 Clique: 3376
 530 MIR: 423
 531 StrongCG: 73
 532 Flow cover: 1515
 533 Inf proof: 13
 534 Zero half: 647
 535 Mod-K: 2
 536 Network: 8
 537 RLT: 262
 538 Relax-and-lift: 610
 539 BQP: 20
 540
 541 Explored 476681 nodes (80360387 simplex iterations) in 1200.84 seconds (767.94 work units)
 542 Thread count was 256 (of 256 available processors)
 543
 544 Solution count 10: 0.523778 0.523778 0.523778 ... 0.573111
 545
 546 Time limit reached
 547 Best objective 5.237777778e-01, best bound 4.712998343346e-01, gap 10.0191%
 548

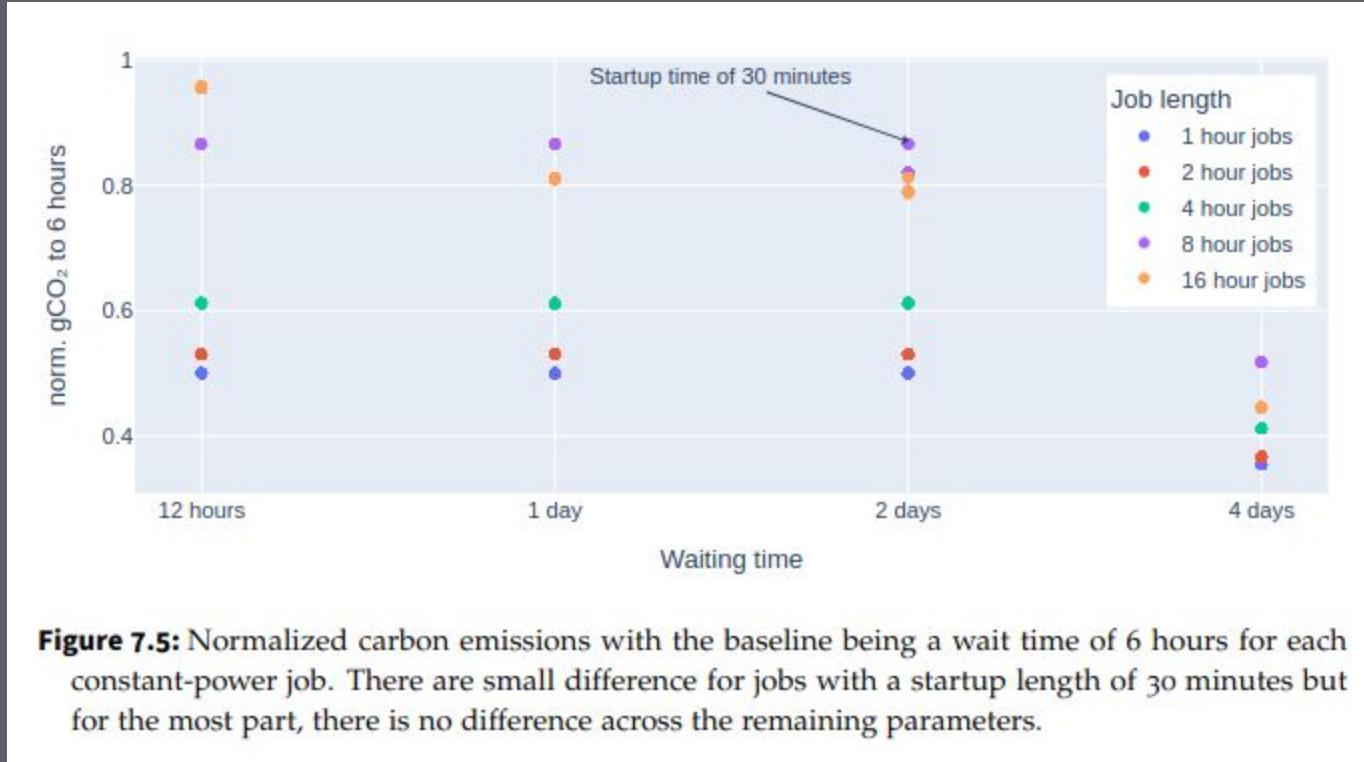
```

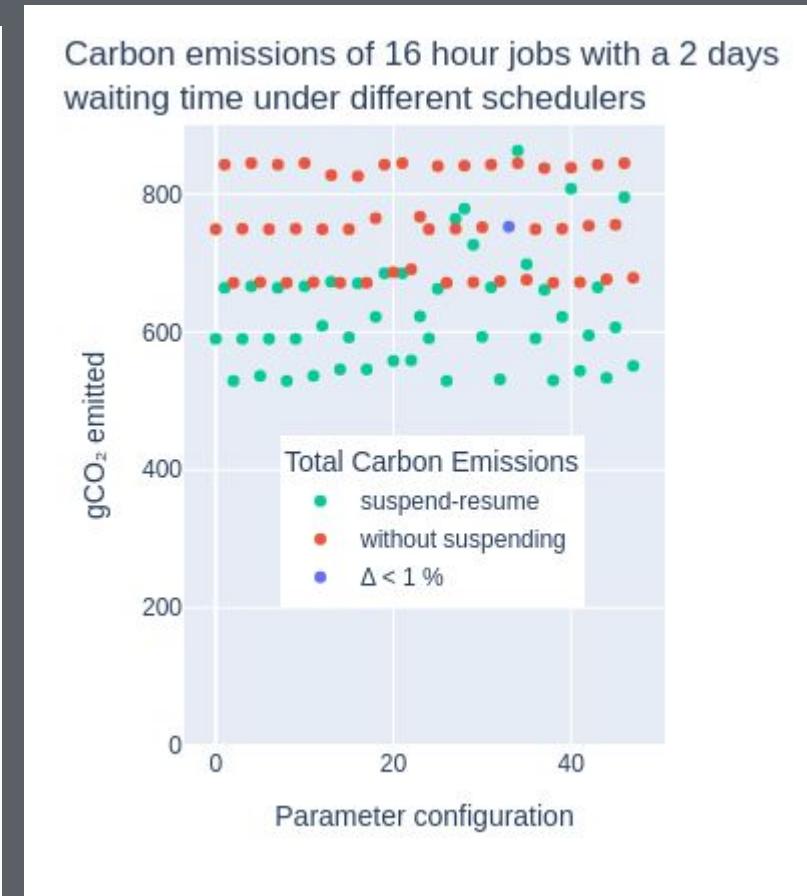
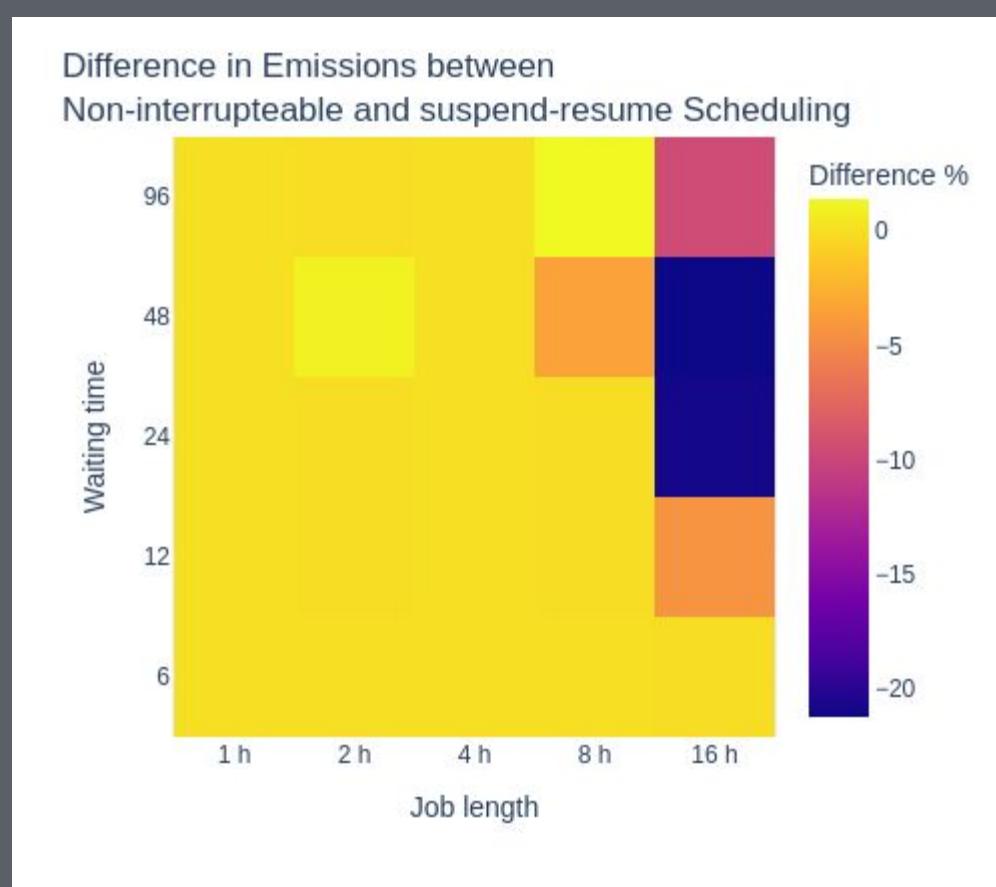
| Top 16 moments | |
|----------------|---------|
| 1 | gap |
| 2 | 0.0000 |
| 3 | 0.0000 |
| 4 | 0.0000 |
| 5 | 10.0191 |
| 6 | 91.3230 |
| 7 | 0.0000 |
| 8 | 0.0000 |
| 9 | 0.0000 |
| 10 | 0.0000 |
| 11 | 0.0000 |
| 12 | 0.0000 |
| 13 | 0.0000 |
| 14 | 0.0000 |
| 15 | 0.0000 |
| 16 | 0.0000 |

Filtern der Logs nach der Gap Value
\$ logs | grep "Best Object" | awk ...

Ausschnitt der Evaluations Logs

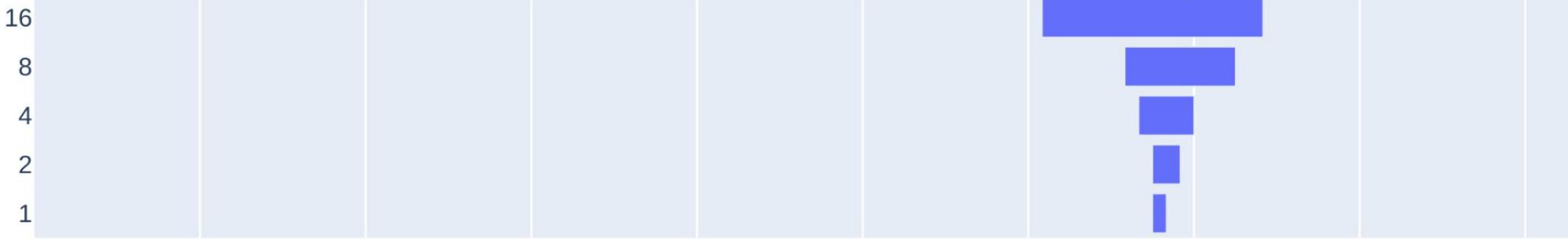






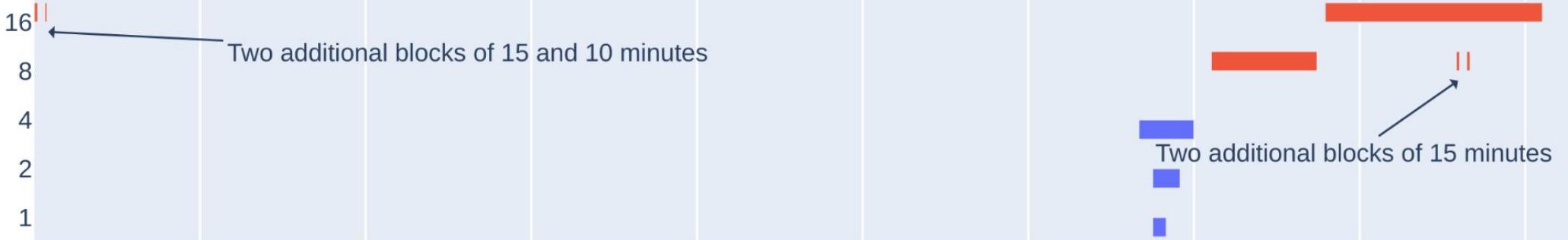
Scheduling without suspend & resume

Job length



Scheduling with suspend & resume

Job length

gCO₂/kWh