



# Leveraging Low-Carbon Energy for Flexible Cloud and Edge Workloads

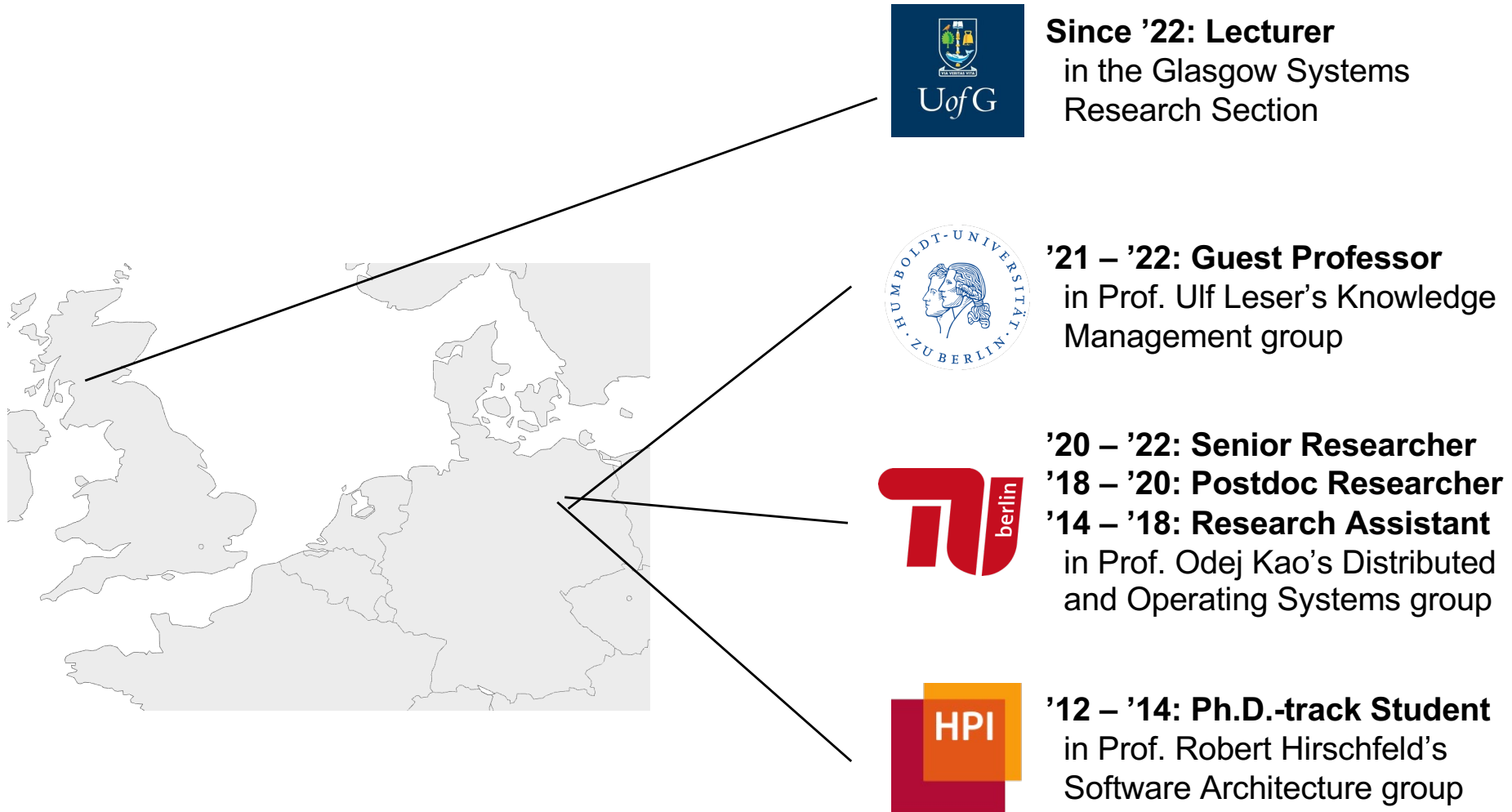
HPI Research Symposium 2024

Dr Lauritz Thamsen

University of Glasgow

<https://lauritzthamsen.org>

# My Path from HPI to Glasgow



# Adaptive Resource Management

---



## Adaptive Resource Allocation

Allocate compute resources to meet specific performance objectives and constraints

e.g. Cluster'21, ICFEC'21, IC2E'21 & '22, EuroPar'22, BigData'22 & '23, FGCS'24



## Scheduling & Dynamic Scaling

Adjust resource configurations at runtime as workloads change or components fail

e.g. CCPE'20, Middleware'21, SPE'21, IPCCC'23, CCGrid'23, e-Energy'24



## Automatic System Tuning

Tune system configurations using monitoring data, profiling, and performance models

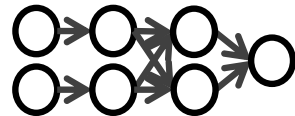
e.g. BigData'19 & '20, IC2E'22, ICWS'22, ICPE'24

# Data-Intensive Systems

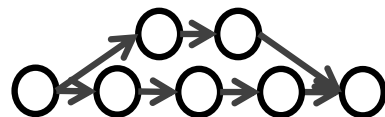
---

Many data-intensive applications run on top of scalable and fault-tolerant distributed processing systems

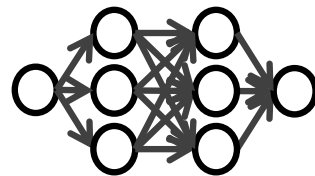
Distributed Batch /  
Stream Processing



Scientific  
Workflows



Machine  
Learning



# Computing's Growing Footprint

---

- **Data centers already consume > 1% of the globally produced energy**, a share that is projected to rise sharply over the next decades
- More and more large-scale, long-running, resource-intensive **data processing jobs** (e.g. Big Data, AI/ML, and IoT)
- Emissions depend on the **energy consumption**, yet also the **specific sources of energy**

# Carbon-Conscious Computing

---

- Objective: **Reducing the carbon footprint** of large-scale data processing applications on today's diverse distributed computing infrastructure
1. Compute when and where low-carbon energy is going to be available
  2. Allocate resources for high resource utilization and highly utilize allocated resources
  3. Save computation and communication through distributed and dynamic architectures

---

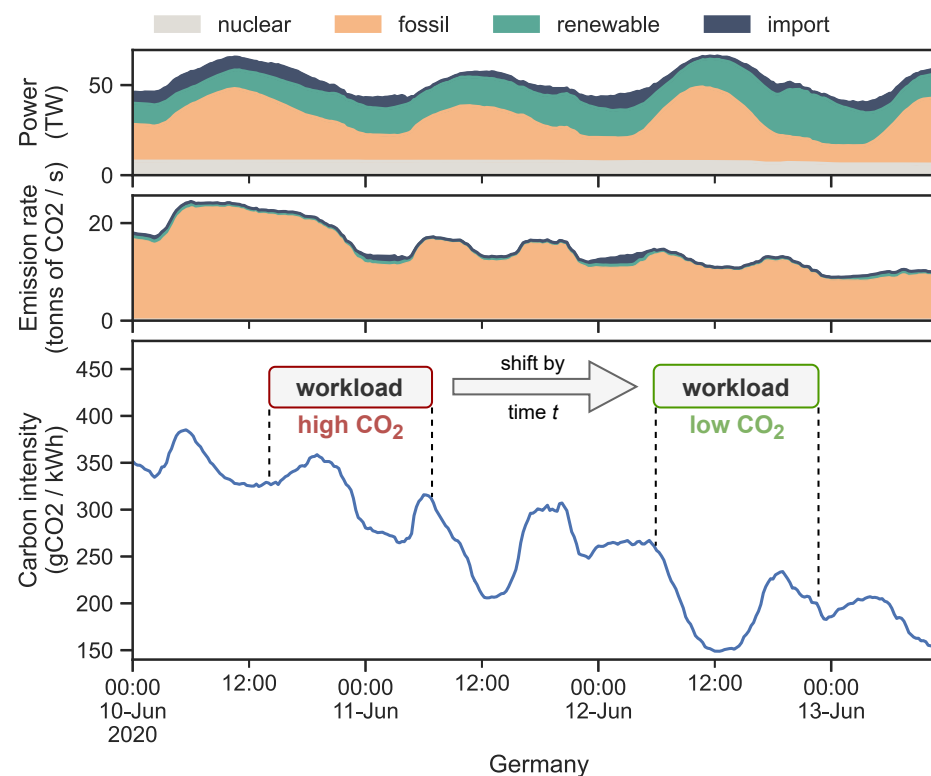
# Carbon-Aware Cloud Workload Shifting

---

Let's Wait Awhile: How Temporal Workload Shifting Can Reduce Carbon Emissions in the Cloud.  
Wiesner, Behnke, Scheinert, Gontarska, Thamsen. ACM/IFIP Middleware 2021.

# Motivation

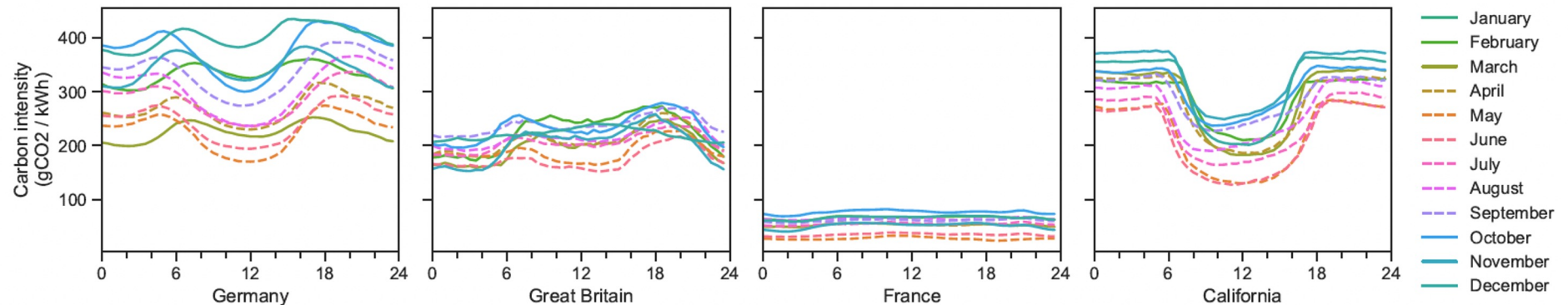
- Emissions of power grids are determined by the energy mix and demand
- Low-carbon objective:  
Compute when and where low-carbon energy is available





# Changing Carbon Intensity (1/2)

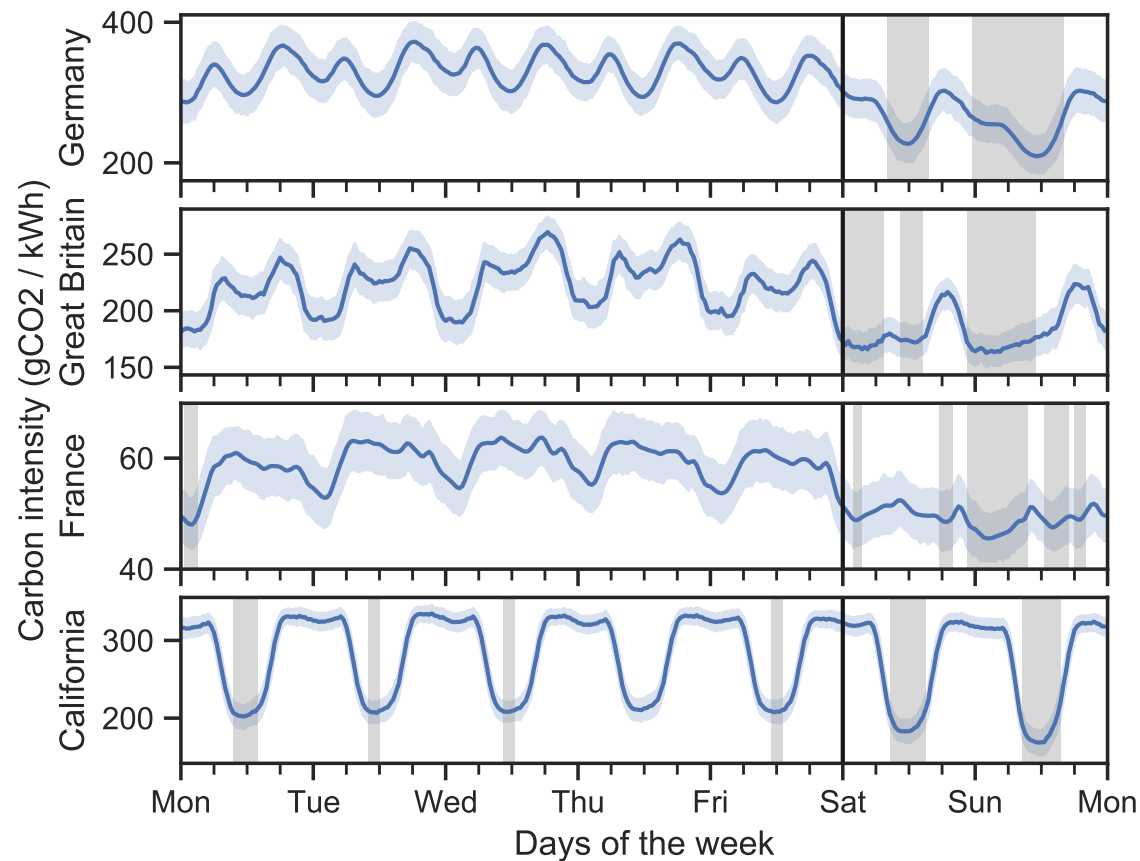
- What are the most promising times to shift work to?



- Average carbon intensity (== CO<sub>2</sub>-equiv. greenhouse gas emissions per kilowatt hour of energy) in 2020

# Changing Carbon Intensity (2/2)

- What are the most promising days to shift work to?



# Simulations

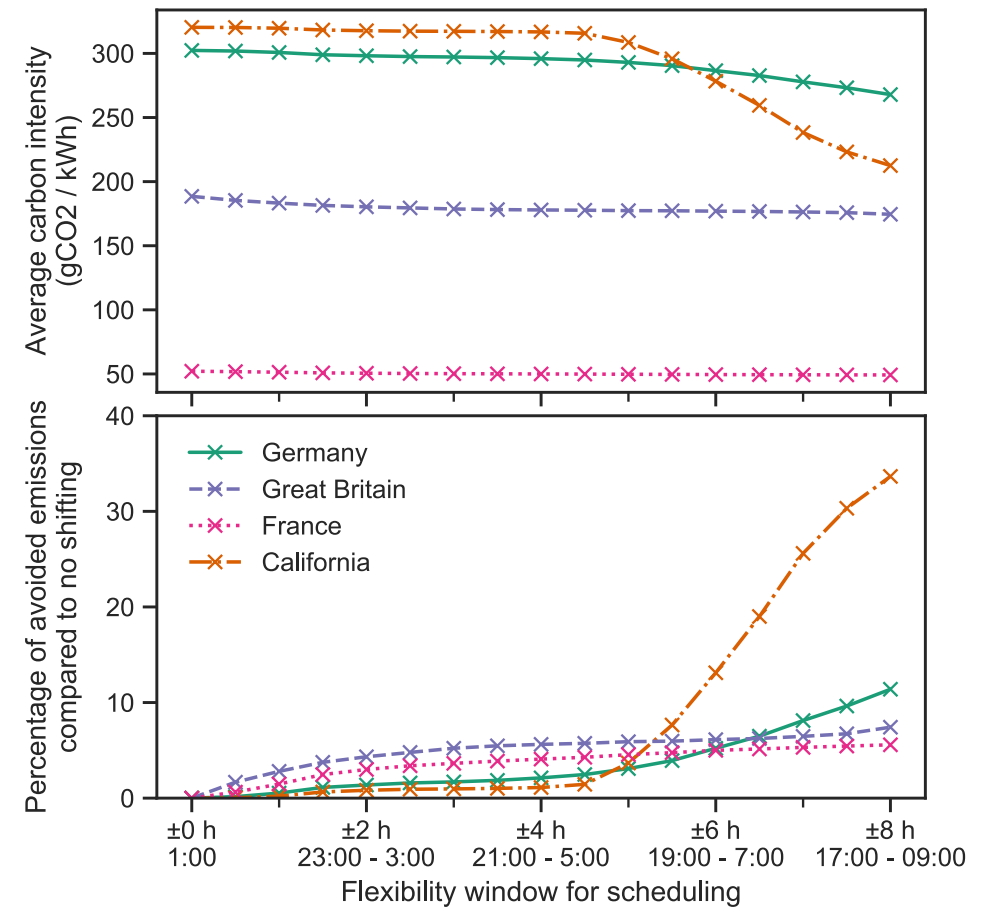
---



- Evaluation of two scenarios using our simulator (<https://github.com/dos-group/leaf>)
- Scenario 1 – Periodic Jobs: Nightly builds, integration tests, recurring generation of business reports, ...
- Scenario 2 – Ad Hoc Jobs: ML training jobs, data analysis pipelines, scientific simulations, ...

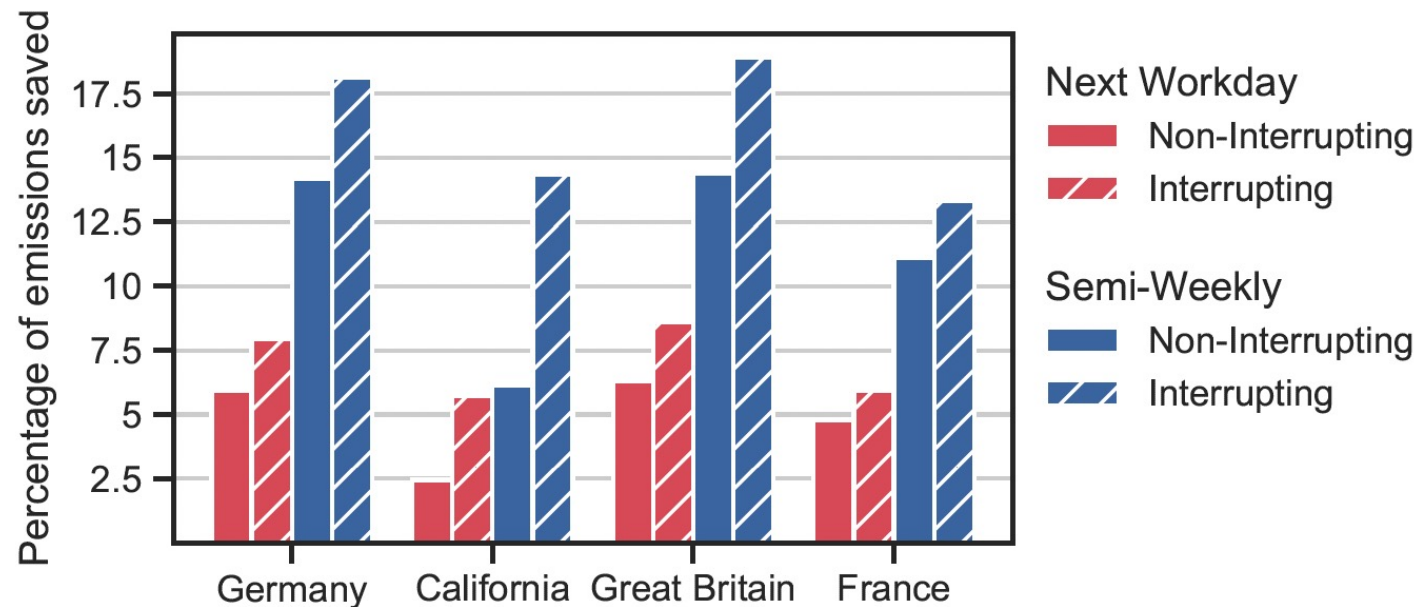
# Scenario 1: Periodic Jobs

- Baseline: All jobs scheduled at 1 am in the night
- Increasing the window by  $\pm 1$  h to allow scheduling between
  - 00:00 to 3:00 ( $\pm 1$  h)
  - 23:00 to 4:00 ( $\pm 2$  h)
  - ...
  - 17:00 to 9:00 ( $\pm 8$  h)



# Scenario 2: Large Ad Hoc Jobs

- Based on an NVIDIA research project, which ran 3387 ML training jobs using 145.76 GPU years and 325 MWh
- Baseline: Instant scheduling of jobs that arrive randomly during working hours



---

# Edge Computing & FL on Renewable Excess Energy

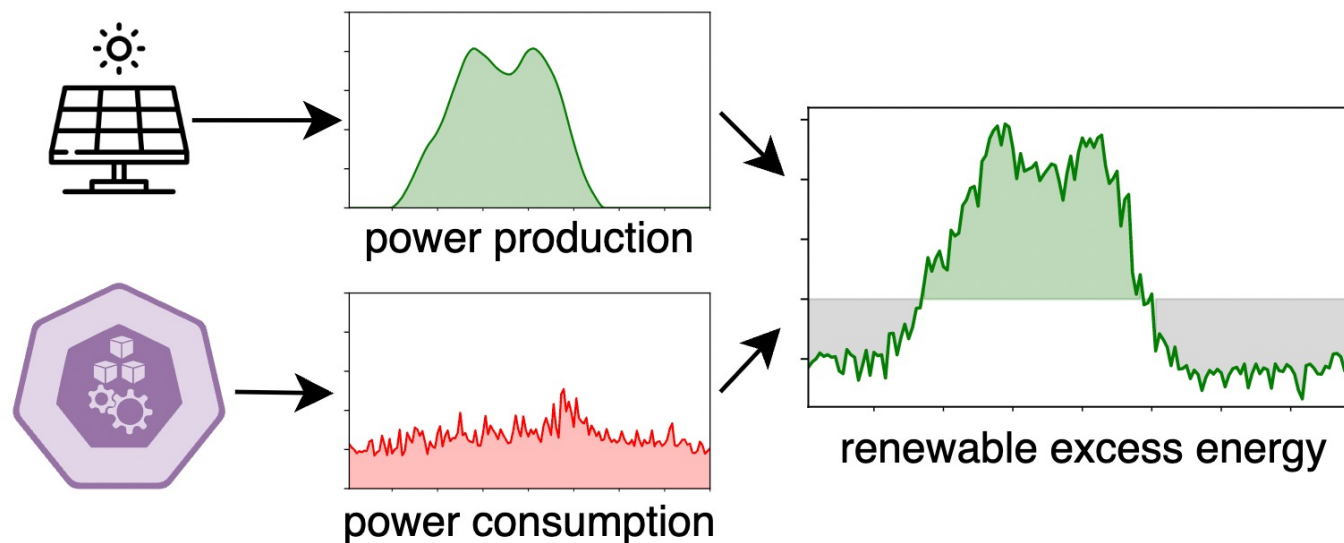
---

Cucumber: Renewable-Aware Admission Control for Delay-Tolerant Cloud and Edge Workloads.  
Wiesner, Scheinert, Wittkopp, Thamsen, Kao. EuroPar 2022.

FedZero: Leveraging Renewable Excess Energy in Federated Learning. Wiesner, Khalili, Grinwald, Pratik  
Agrawal, Thamsen, Kao. ACM e-Energy 2024.

# Renewable Excess Energy

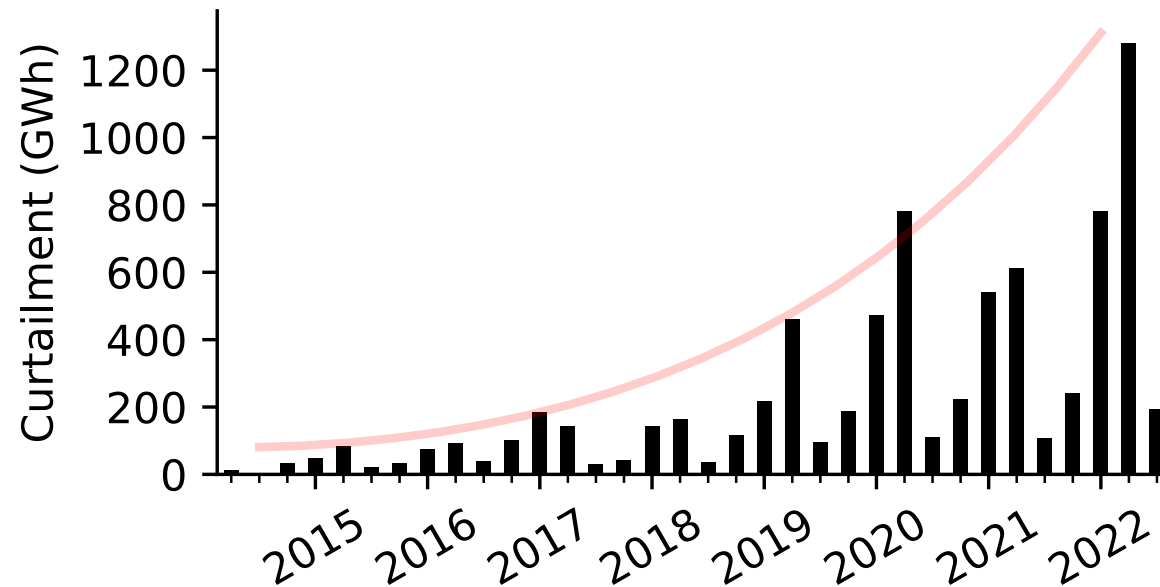
The output from renewables such as solar and wind varies, and there can be more energy than demand



# Solar Energy Curtailment in California

---

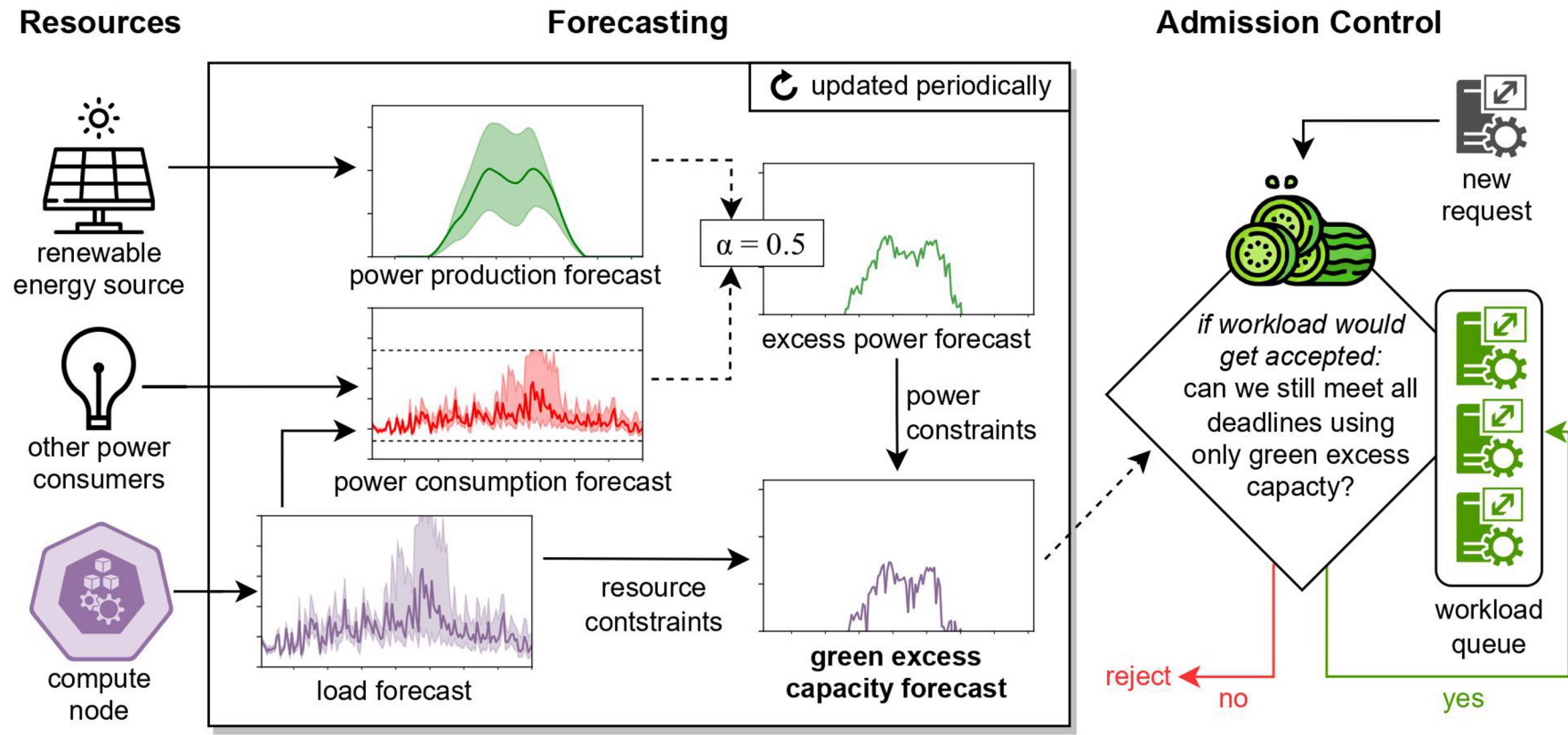
Around 7% of solar power is being curtailed already



Source: California Independent System Operator (CAISO)

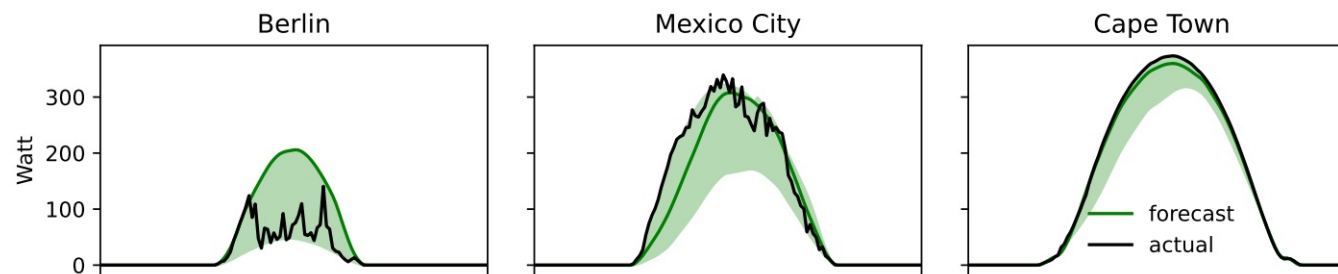


# “Cucumber” Overview

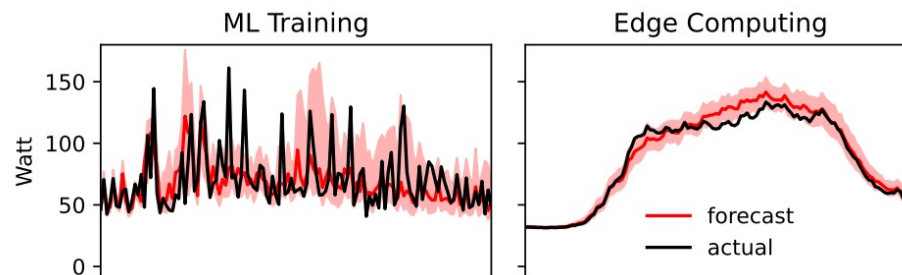


# Simulation Setup

Two weeks of solar production forecasts for 400W panels across three sites (using <https://solcast.com/>):

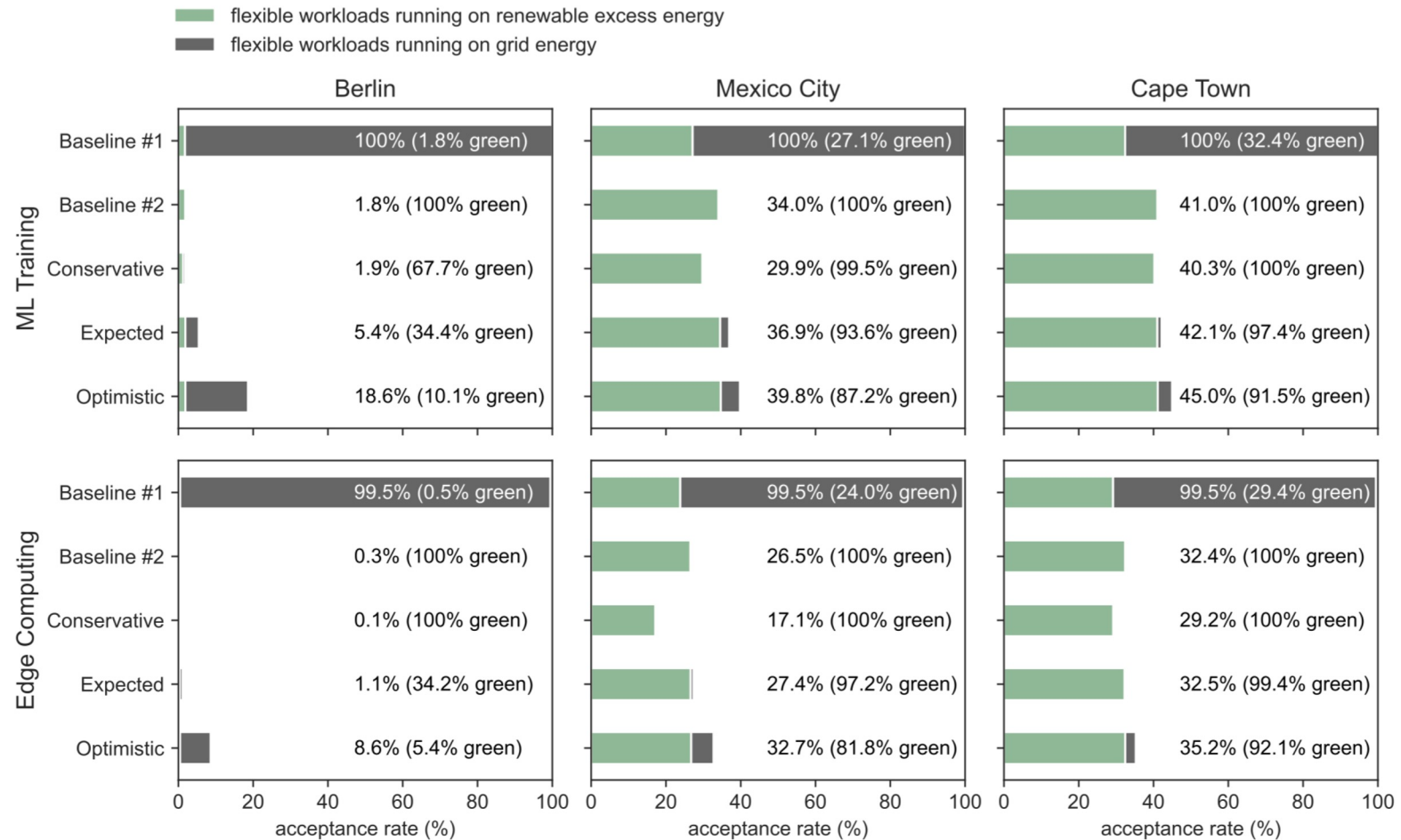


Two workload traces:



**ML Training** based on Alibaba GPU cluster traces (deadlines set to midnight)  
**Edge Computing** based on a NYC taxi trip dataset (deadlines derived from trip lengths)

# Simulation Results



# Where are the Flexible Low-Priority Jobs Coming From?

---

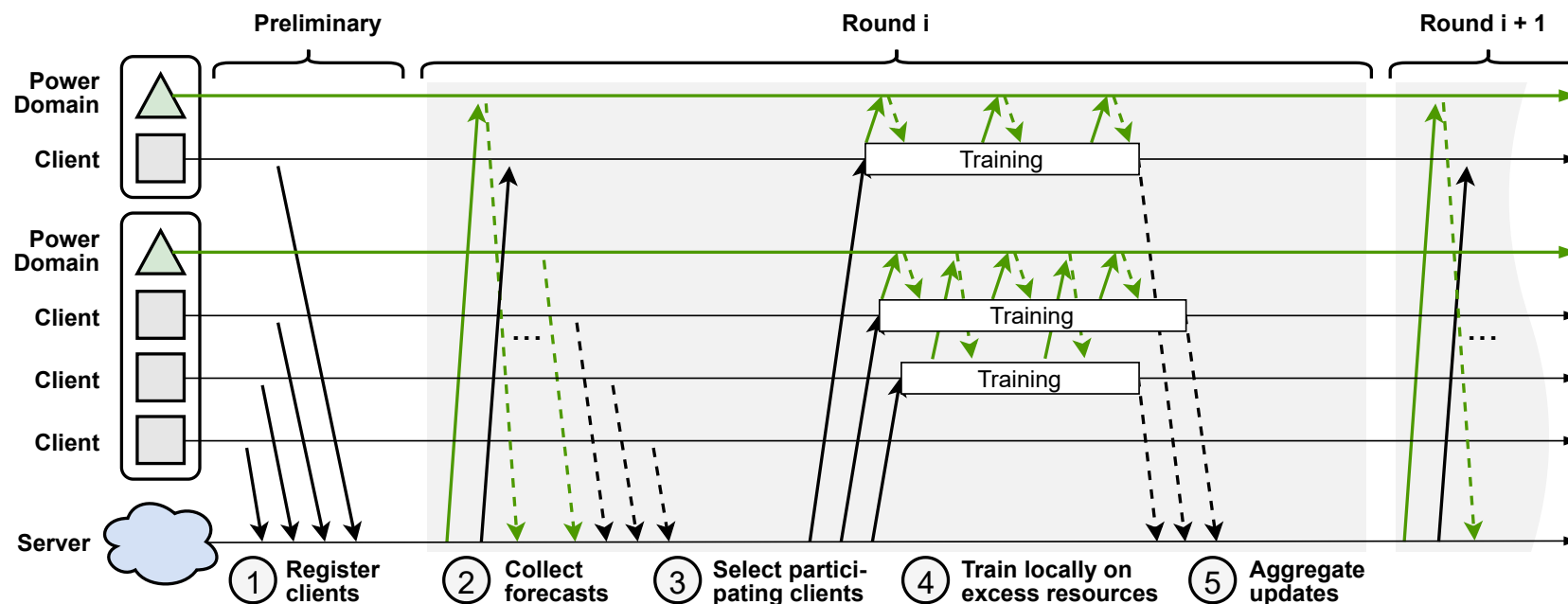
A recent trend in distributed ML is Federated Learning (FL), improving data privacy

FL is a promising candidate for carbon-aware computing:

- It constitutes energy-intensive batch jobs
- It is scheduled in geo-distributed environments
- It is rather flexibility (e.g. in terms of which clients participate in an iteration and iteration deadlines)

# FedZero Client Selection Protocol

Scalable client selection strategy for Zero-Carbon Federated Learning with fast convergence and fair client participation based on forecasts



# FedZero Overall Results

Dataset & model	Data distribution & aggregation strategy	Approach	Global			Co-located		
			Target accuracy	Time-to-accuracy	Energy-to-accuracy	Target accuracy	Time-to-accuracy	Energy-to-accuracy
CIFAR-10 ResNet-18	iid FedAvg	Constrained	83.26 %	7.0 d	72.1 kWh	84.04 %	6.6 d	101.5 kWh
		FedZero		4.6 d	74.5 kWh		5.6 d	109.7 kWh
		Unconstrained		1.5 d	85.1 kWh		2.2 d	128.1 kWh
	non-iid FedProx	Constrained	79.11 %	6.7 d	71.3 kWh	80.53 %	6.6 d	86.9 kWh
		FedZero		4.3 d	67.4 kWh		4.8 d	88.0 kWh
		Unconstrained		1.4 d	68.7 kWh		2.1 d	105.6 kWh
CIFAR-100 DenseNet-121	iid FedAvg	Constrained	57.85 %	6.7 d	78.6 kWh	58.90 %	6.7 d	101.2 kWh
		FedZero		4.4 d	82.1 kWh		4.5 d	94.3 kWh
		Unconstrained		1.5 d	89.0 kWh		2.0 d	119.6 kWh
	non-iid FedProx	Constrained	56.32 %	6.7 d	76.6 kWh	57.63 %	6.8 d	102.1 kWh
		FedZero		4.5 d	82.7 kWh		4.6 d	99.5 kWh
		Unconstrained		1.5 d	88.3 kWh		2.2 d	128.7 kWh
Shakespeare LSTM	non-iid FedProx	Constrained	52.14 %	5.7 d	79.3 kWh	52.57 %	6.7 d	77.2 kWh
		FedZero		1.4 d	27.7 kWh		2.2 d	37.8 kWh
		Unconstrained		1.4 d	46.7 kWh		1.9 d	65.0 kWh

---

# Summary

---

# Carbon-Aware Edge/Cloud Computing

---

- Fluctuations in grid carbon intensities can be leveraged to reduce the footprint of flexible workloads (e.g. by 5-20% in some cases)
- Renewable excess energy and spare compute capacity can drive down the operational footprint of e.g. ML and FL substantially
- Interesting open questions for realizing the potential:
  - Which signal to use? e.g. carbon intensity or curtailed energy?
  - How best to profile and predict application performance?
  - How to conduct experiments? e.g. (hybrid) simulations, real data?
- Contact: [lauritz.thamsen@glasgow.ac.uk](mailto:lauritz.thamsen@glasgow.ac.uk)