

# Sustainability for the Cloud Native world

Data Centers are generating over 2 Billion Tons of CO2 annually and by 2030 may consume up to 10% of the world’s energy

**Objective:**  
To reduce the world’s Data Center carbon footprint by making the best use of servers across multiple workloads.

**Approach:**  
1. Figure out which servers connect to which outlets  
2. Use this data to generate a heuristic for each workloads power use on each server  
3. Use the Heuristic to fill in empirical gaps in our very sparse knowledge then optimize a new distribution  
4. Run the new assignment of jobs and observe the results, improving the data set at hand and the heuristic model

**(A) Matching Phase:**

- We have Power Data Stream and Server Data Streams
- We don’t know which matches to which
- Use a Factor Graph model to match the Servers to Outlets

**(B) Extrapolation Phase:**

- We have some data about the power usage of Servers and Workloads
- Not every workload has been observed on every server type.
- Use a neural net to generate a heuristic Predicted Power for every Server Type/Workload that we don’t have observed data for

**(C) Optimization Phase:**

- We now have Real or heuristic estimates of power
- We are given a set of workloads to run on our servers.
- Use a CSP Solver to optimize the overall power usage of our systems

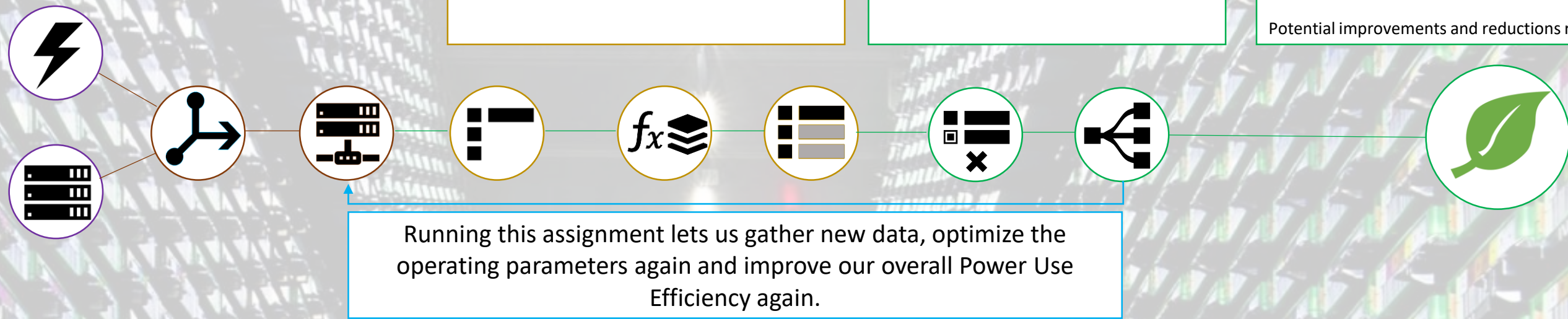
**Motivation:**

Data Centers account for an estimated 1-2% of global energy usage forecasted to grow to 10% by 2030  
Gartner estimates that ongoing power costs are increasing at least 10% per year

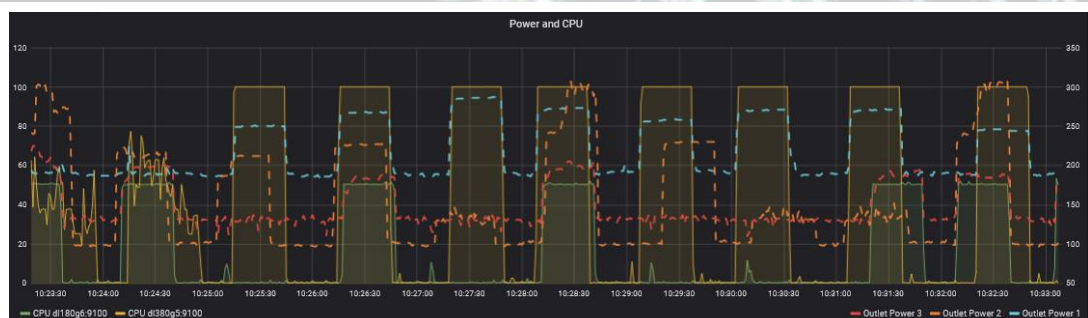
**Project Outcomes:**

**-12% Energy Used**  
**+2% Asset Efficiency**

Potential improvements and reductions relative to baselines

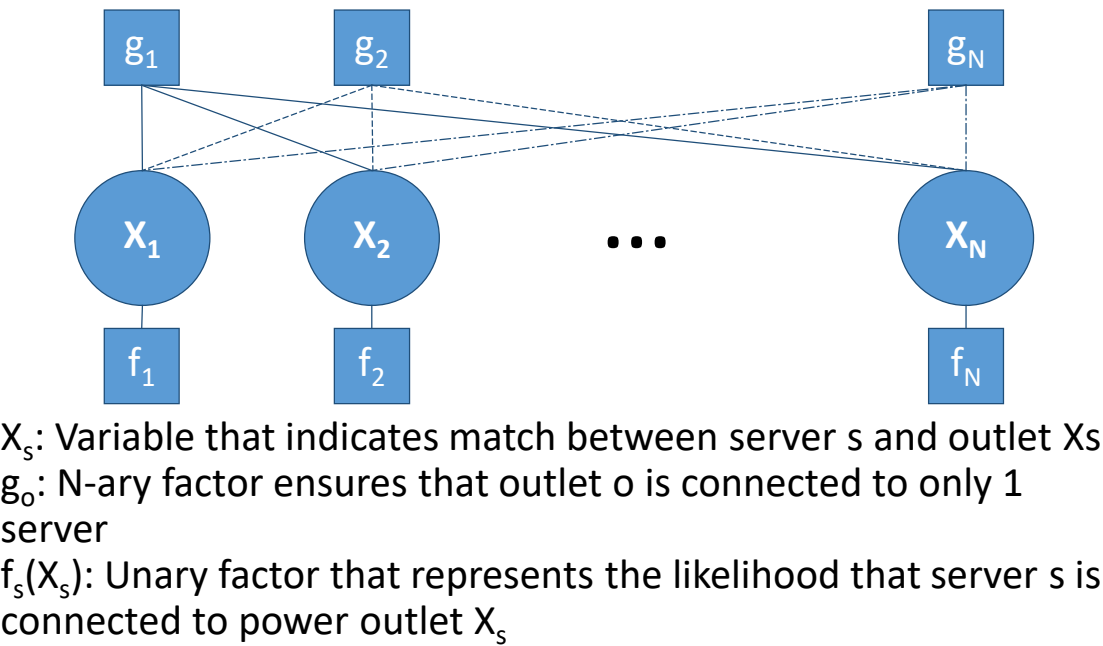


**(A) Matching: Data and Model**



Two CPU data streams and three Power Outlet streams

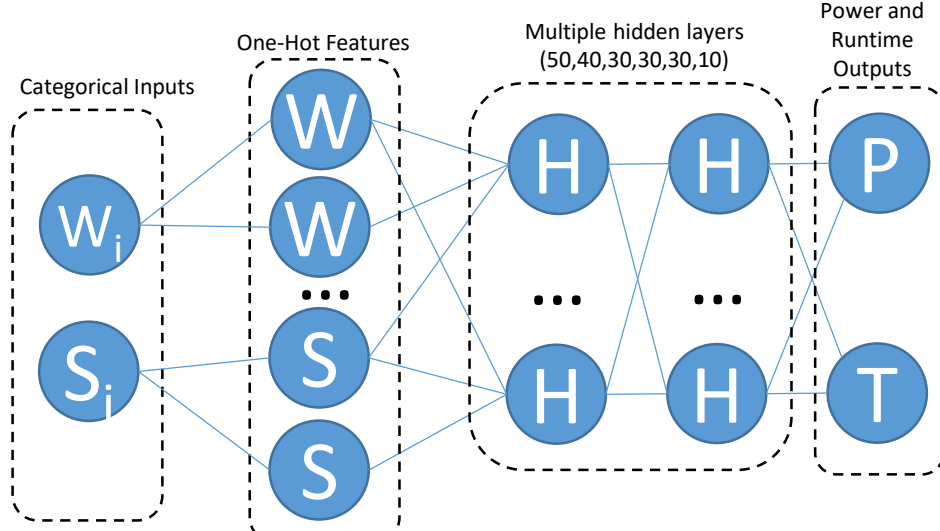
**Matching factor graph:**



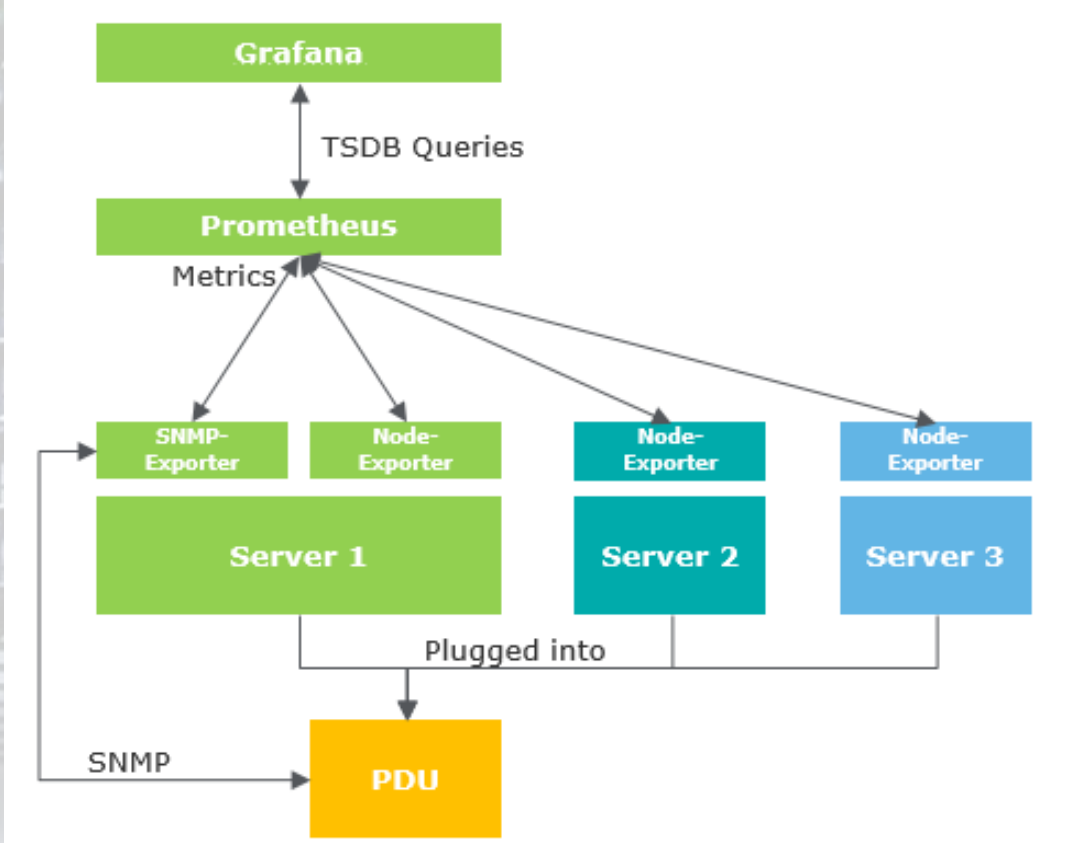
**(B) Extrapolation: Data and Model**

Server Type	Workload			
	MatrixProd	MemoryBM	Streaming	DiskBM
DI360G7	0.3Wh	0.65Wh	0.6Wh	0.3Wh
DI360G6	0.35Wh	0.7Wh	0.72Wh	0.4Wh
DI360G6 (SSD)	0.32Wh	0.7Wh	0.8Wh	0.2Wh

Illustrative **Measured** and **Extrapolated** Data (note we predict both power and runtime)



**Data Capture Setup**

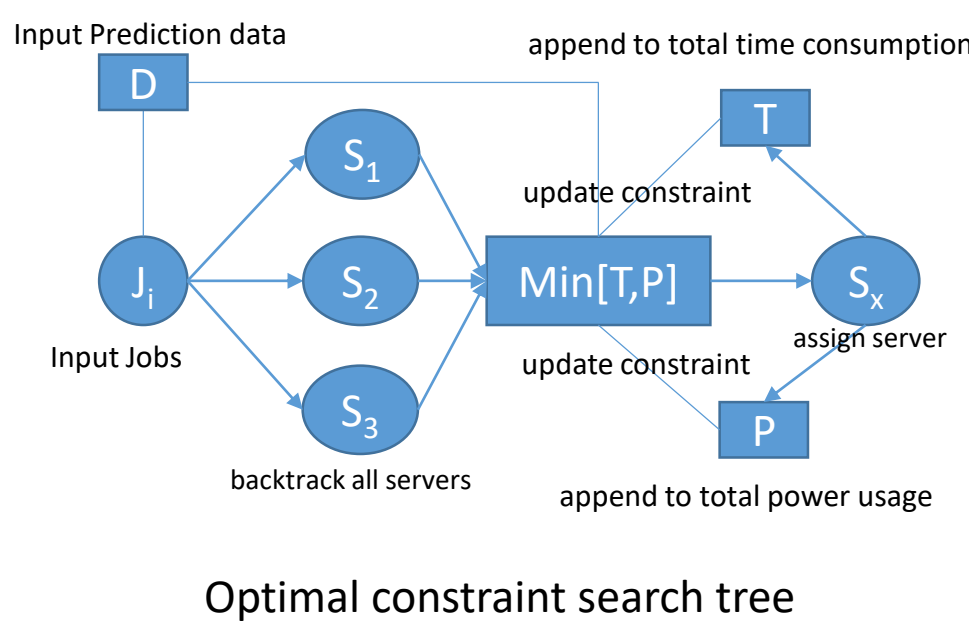


- 5 servers commonly used in cloud datacenters (HP ProLiant series) are deployed in a lab,
- Each with different configurations of CPU, Disks, Memory.
- Workloads are dispatched to each from a central server allowing us to run our experiments and monitor results

**(C) Optimization: Data and Model**

**Approach**  
We employed three strategies to explore the problem and solution space, assessing total power use and total runtime of the assignments generated  
1. Assignment to servers without regard to predicted values, distributing jobs with uniform distribution  
2. Assigning all jobs to the most efficient server and shutting down all others  
3. A CSP/Search Problem that takes into account the Runtime and Power predictions to generate an optimal distribution.

Note: Since we are optimizing for carbon emission reduction we entertained the greedy option which is equivalent to replacing all servers with higher efficiency servers.



We generated trials consisting of a collection of workloads to execute then assigned them using each model.

**Key Observations:**

- The most meaningful results are from the CSP and Uniform Random strategies.
- Since widely diverse workloads are considered the savings vary from Workload to Workload

**Samples:**

Workload (50 executions)	Baseline Duration	CSP Duration	Baseline Power Usage	CSP Power Usage
HDD	21min55	21min38	39.43 Wh	36.26 Wh
10kStream	29min23	29min27	169.68 Wh	142.83 Wh
50kMemory	32min02	31min59	129.99 Wh	127.73 Wh

**(A) Matching: Analysis and Results**

**Baseline algorithm:** Assign outlet with highest average power usage to server with highest average CPU usage. Repeat process until all outlets are matched to servers.

**Factor graph algorithm:** Find maximum weight assignment of factor graph:

$$\operatorname{argmax}_x = \prod_{s=1}^N f_s(X_s) \cdot \prod_{o=1}^N g_o(X)$$

**Results:** We tested the baseline and factor graph algorithm on a problem with 219 independent pairings\*

	Accuracy
Baseline algorithm	59%
Factor graph algorithm	97%

\*Generated from 3 servers and 3 power outlets over 73 time slots of 2 minutes.

**(B) Extrapolation: Analysis and Results**

Using Matched Data, our predictive model is able to predict both Power Usage and Runtime of a standardized workload using a Multi-Layer-Perceptron model.

	Power MAE %	Power RMSE	Time MAE%	Time RMSE
Baseline	47.70%	7,164	30.60%	23.2
Oracle	6.90%	512	4.20%	2.5
MCP	27.30%	2,891	22.60%	11.7

RMSE: Root Mean Squared Error  
MAE%: Mean Average Error Percentage

Power and Execution time are jointly estimated although improvements to time estimate have been seen using alternative Neural Net architectures in Keras. Generalization is a significant factor for us as well and the results have so far been very similar even when altering the nature and number of the workloads or introducing new servers

**(C) Optimization: Analysis and Results**

Our CSP model shows significant improvement on expected time consumption and power usage compared to random assignment. Power Usage shows the most improvement (aligned to our goals) while allowing for utilization of even the less efficient servers. Expected time consumption varies less among different servers but also shows some improvement, likely because there is a penalty associated with server overhead (OS, Fans etc.).

Strategy	Total Runtime	Total Power Usage
Random Assignment	2hr 41min	1.787 kWh
Uniform Random	1day 8hrs 39min	1.584 kWh
Final CSP Model	2hr 30min	1.587

**Total Improvement on Baseline:**

↓ -12.60% Total Power Use    ↓ -2.11% Total Runtime

**Summary:**

An innovative combination of AI approaches helps provide a practical solution for **reducing the Carbon Footprint of Datacenters**.

Throughout the solution we focus on developing a model that is implementable for most Data Center operators using the systems they have installed today and data that is already at their disposal.

As a result we can make a difference while providing operators a **net positive financial return** on investment due to **greater asset efficiency** and **lower power consumption**.

**References**

- Makris, Theodosios, "Measuring and Analyzing Energy Consumption of the Data Center". <https://pdfs.semanticscholar.org/a2da/e6059153e57a8bbd7f8eeca59c0a5d90404.pdf>.
- Whitehead, Beth, Andrews, Deborah, Shah, Amip, Maidment, Graeme, "Assessing the environmental impact of data centres part 1: Background, energy use and metrics". <https://www.sciencedirect.com/science/article/pii/S036013231400273X>
- Jones, Nicola, "How to stop data centers from gobbling up the world’s electricity". <https://www.nature.com/articles/d41586-018-06610-y>
- Pettey, Christy, "5 Steps to Maximize Data Center Efficiency". <https://www.gartner.com/smarterwithgartner/5-steps-to-maximize-data-center-efficiency/>

Developed for CS221 by:  
Danni Ma, Roy Justus, Yannick Meyer

