

# **Carbon-Aware Computing: How to Get Power Systems and Data Centers to Talk to Each Other**

**Vladimir Dvorkin**

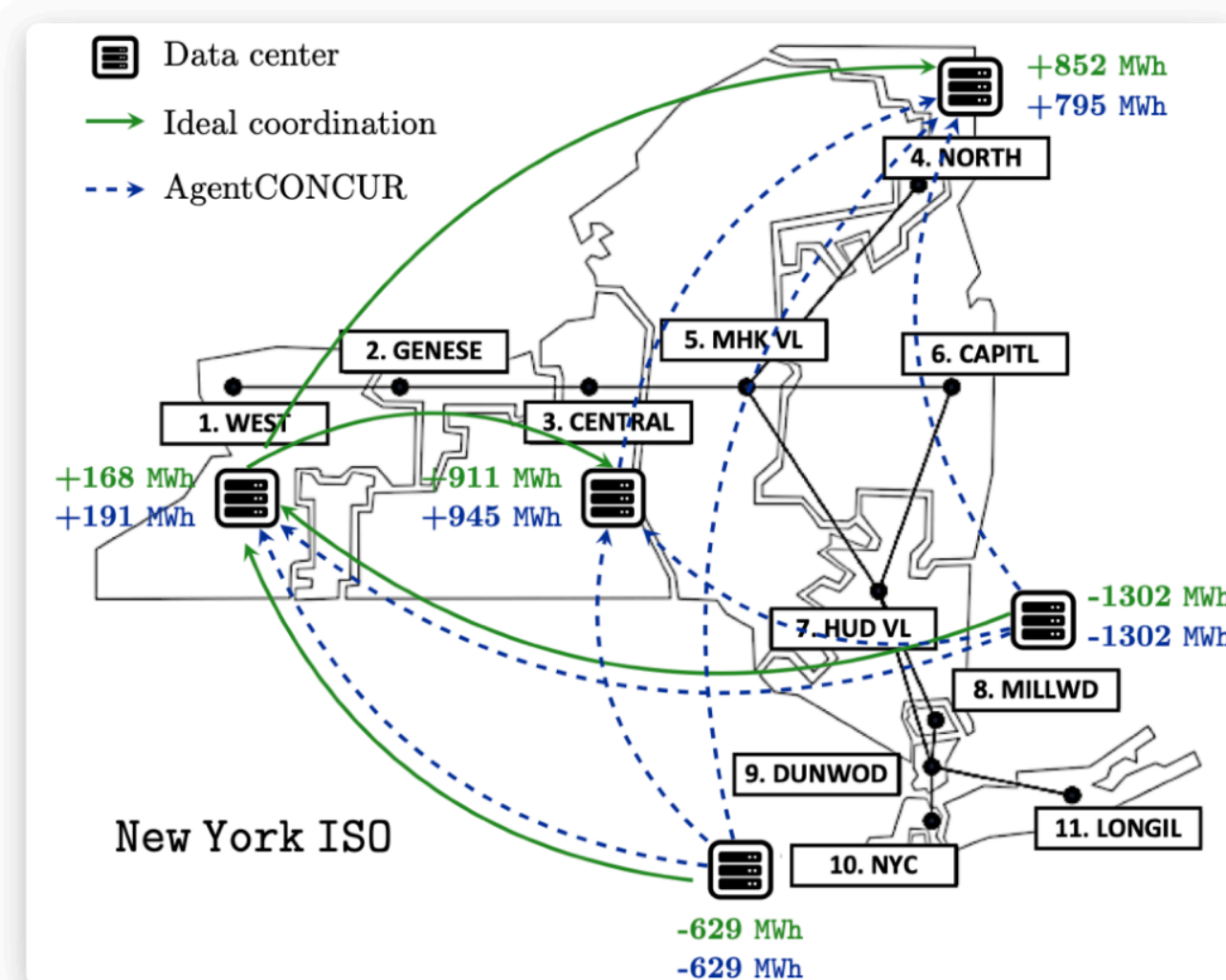
Department of Electrical Engineering and Computer Science

University of Michigan – Ann Arbor

IEEE CSS Day 2024

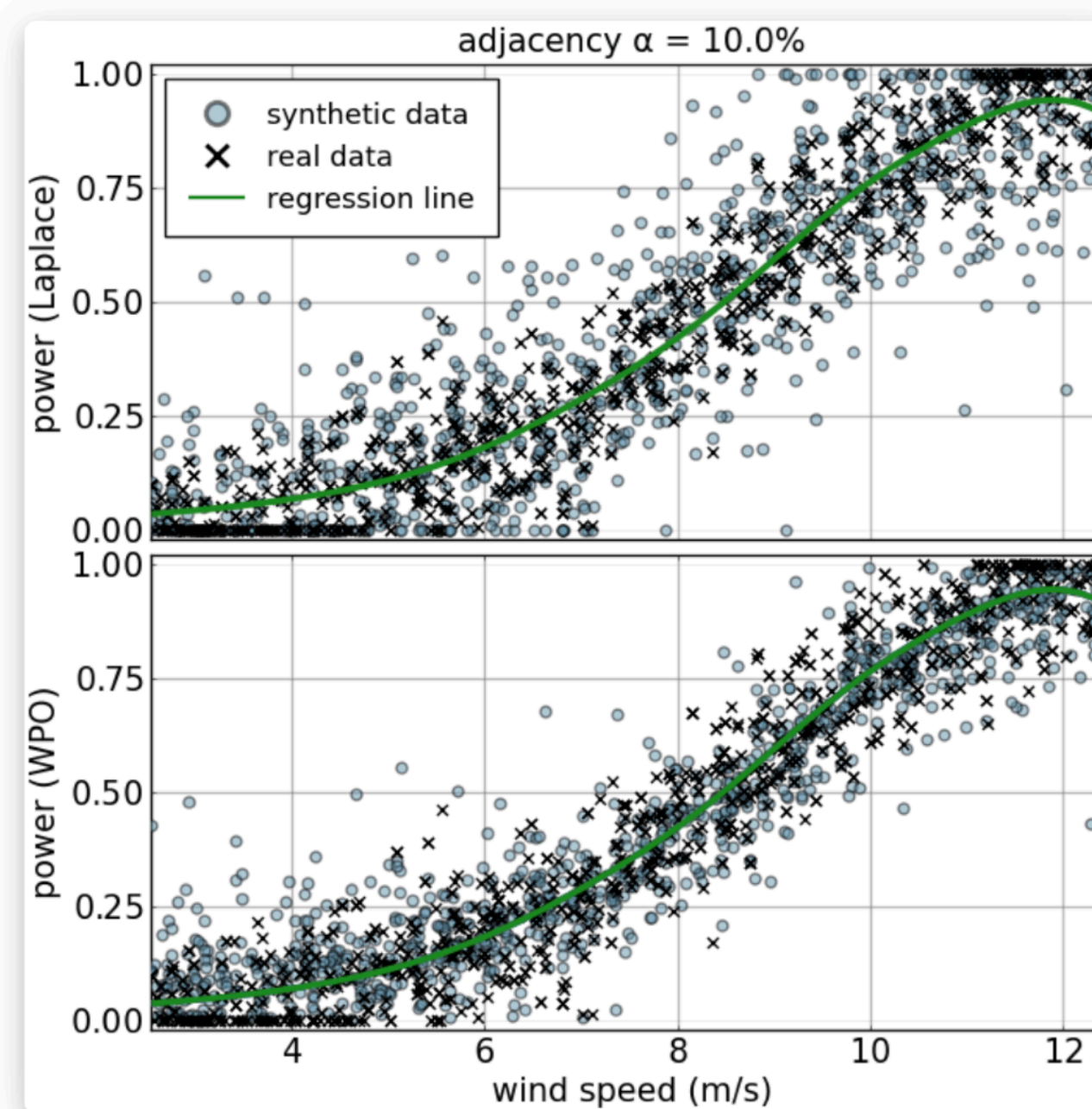
October 21, 2024

## Optimizing energy. Empowering society.



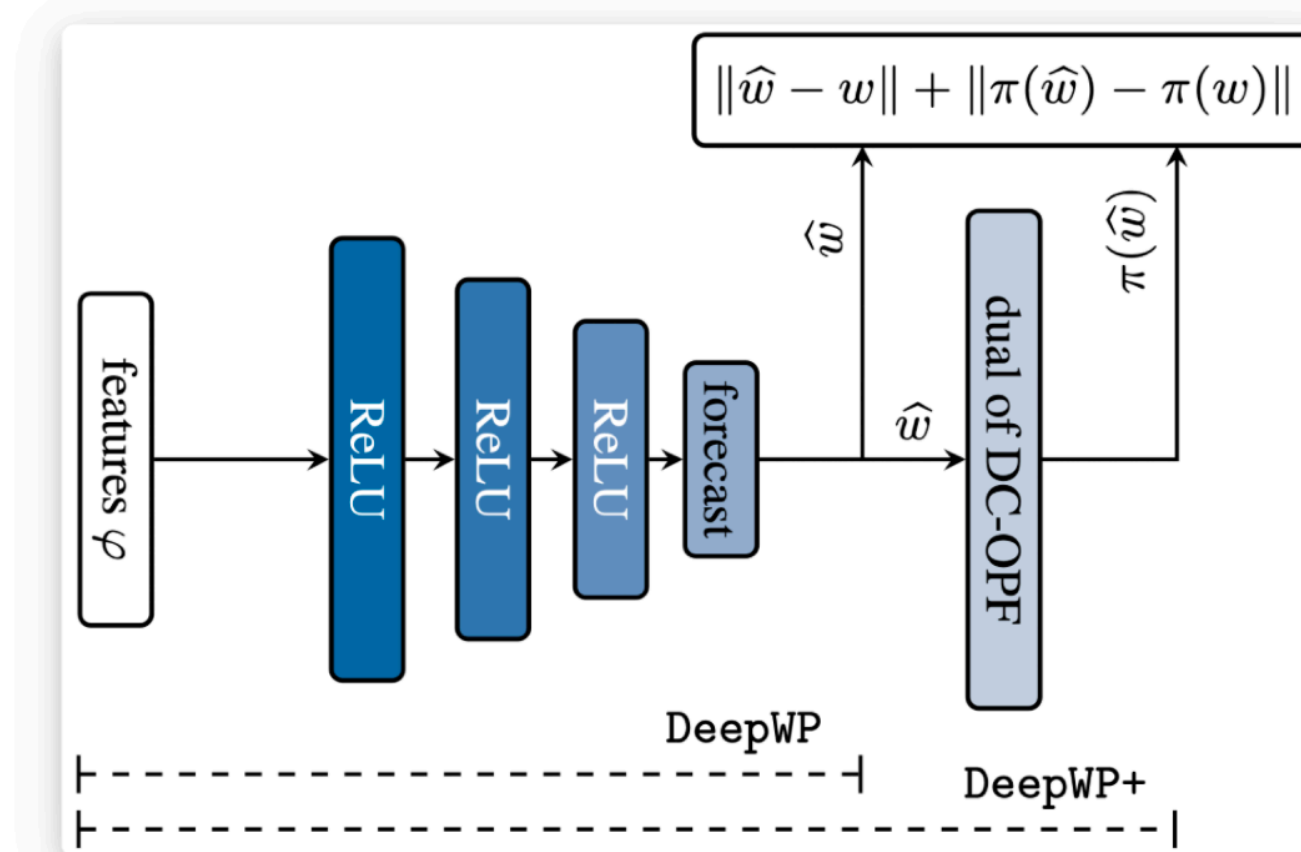
Agent coordination via contextual regression (AgentCONCUR) is a cost-effective protocol, both in terms of data and computation, designed to coordinate spatial shifts of data center electricity consumption to support power grid operations during peak hours.

[\[paper\]](#) [\[code\]](#)



Wind Power Obfuscation (WPO) is a noise-additive algorithm to create synthetic data from real measurements while maintaining higher quality than standard Laplace-based algorithms.

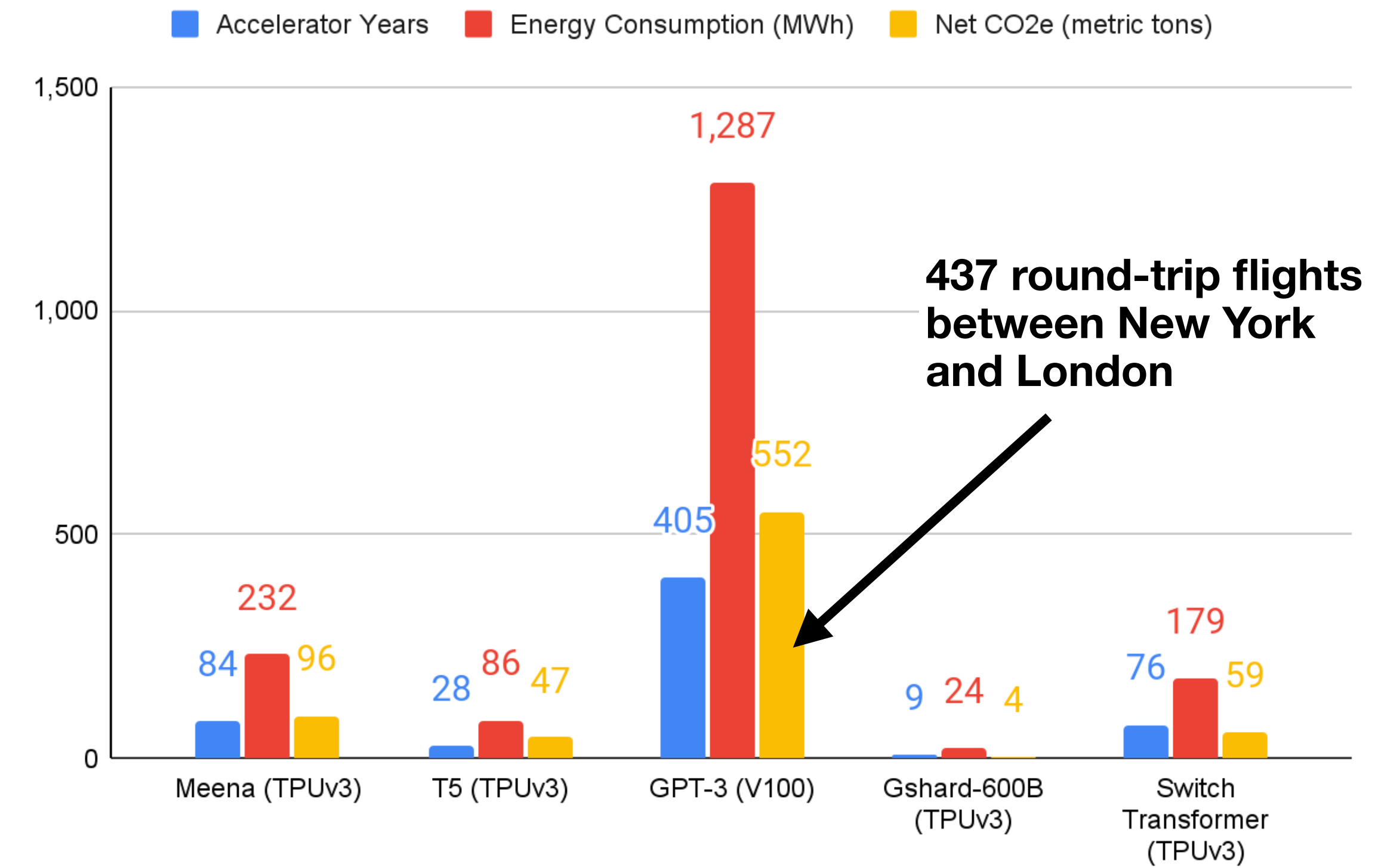
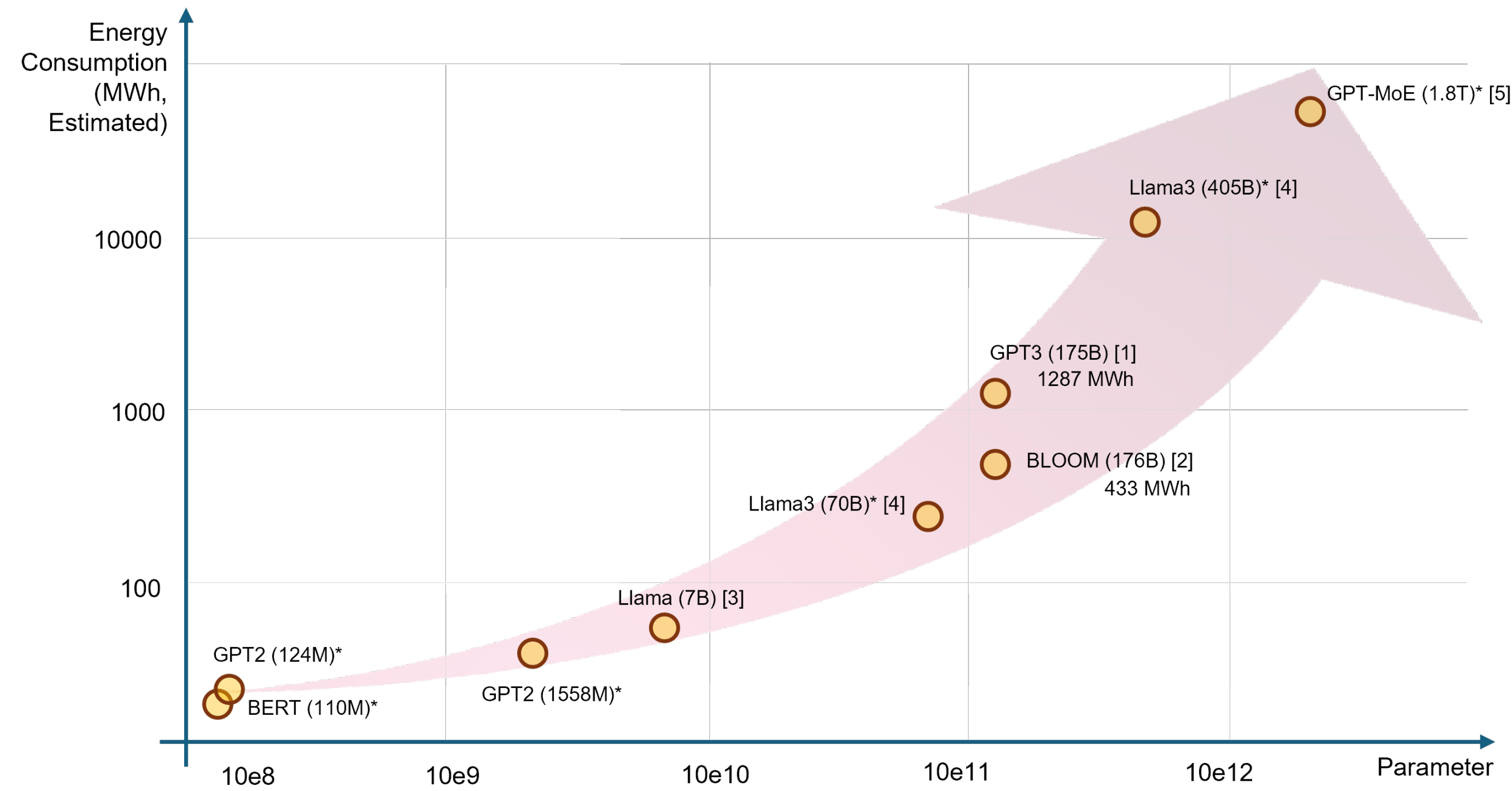
[\[paper\]](#) [\[code\]](#) [\[FERC talk\]](#)



DeepWP+ is a new deep learning architecture designed to enhance wind power predictions by addressing errors and unfairness in electricity prices. Unlike traditional architectures, DeepWP+ embeds a market-clearing optimization problem that guides predictions towards more accurate and fair outcomes.

[\[paper\]](#) [\[slides\]](#)

# AI is currently the fastest growing electricity demand

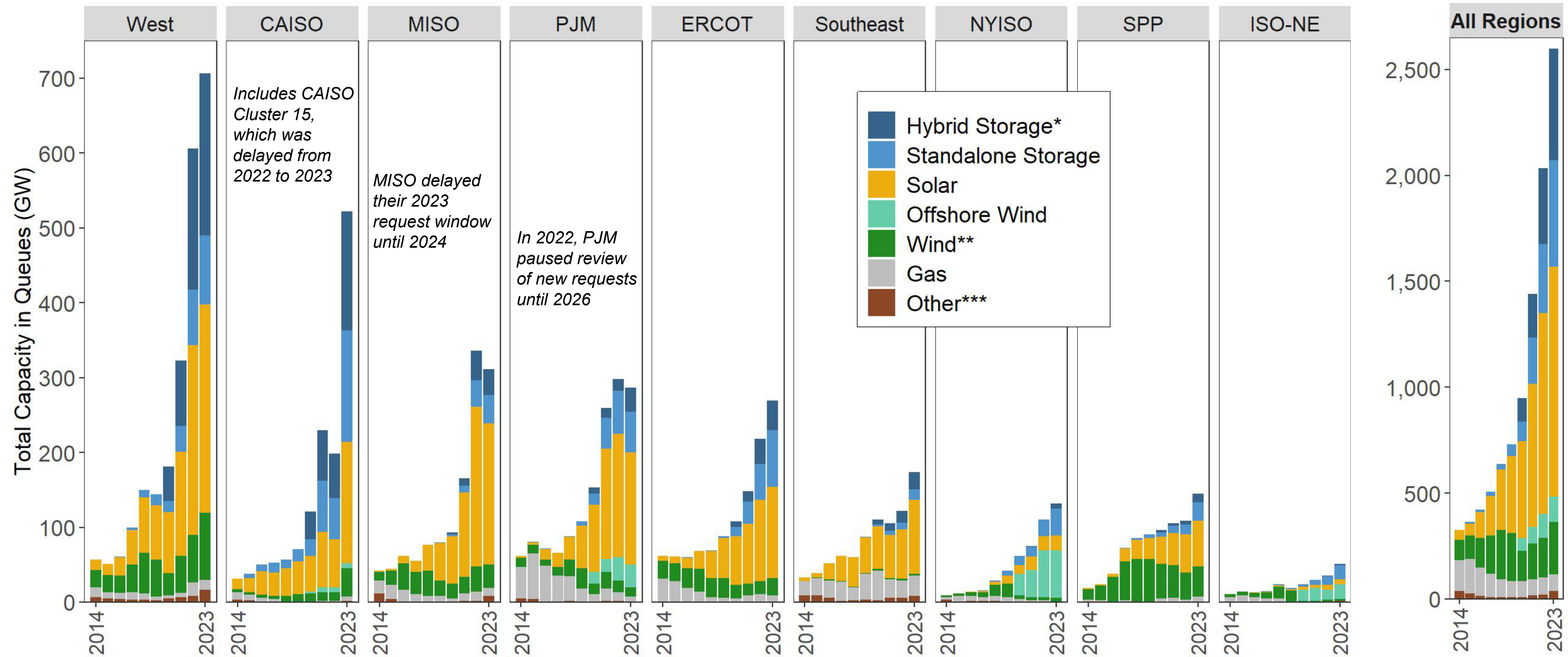


Li et al. The Unseen AI Disruptions for Power Grids: LLM-Induced Transients. 2024

Patterson et al. Carbon Emissions and Large Neural Network Training. 2021

- ▶ Data centers in the U.S. could consume as much electricity by 2030 as some entire industrialized economies
- ▶ Data centers have a significant environmental impact

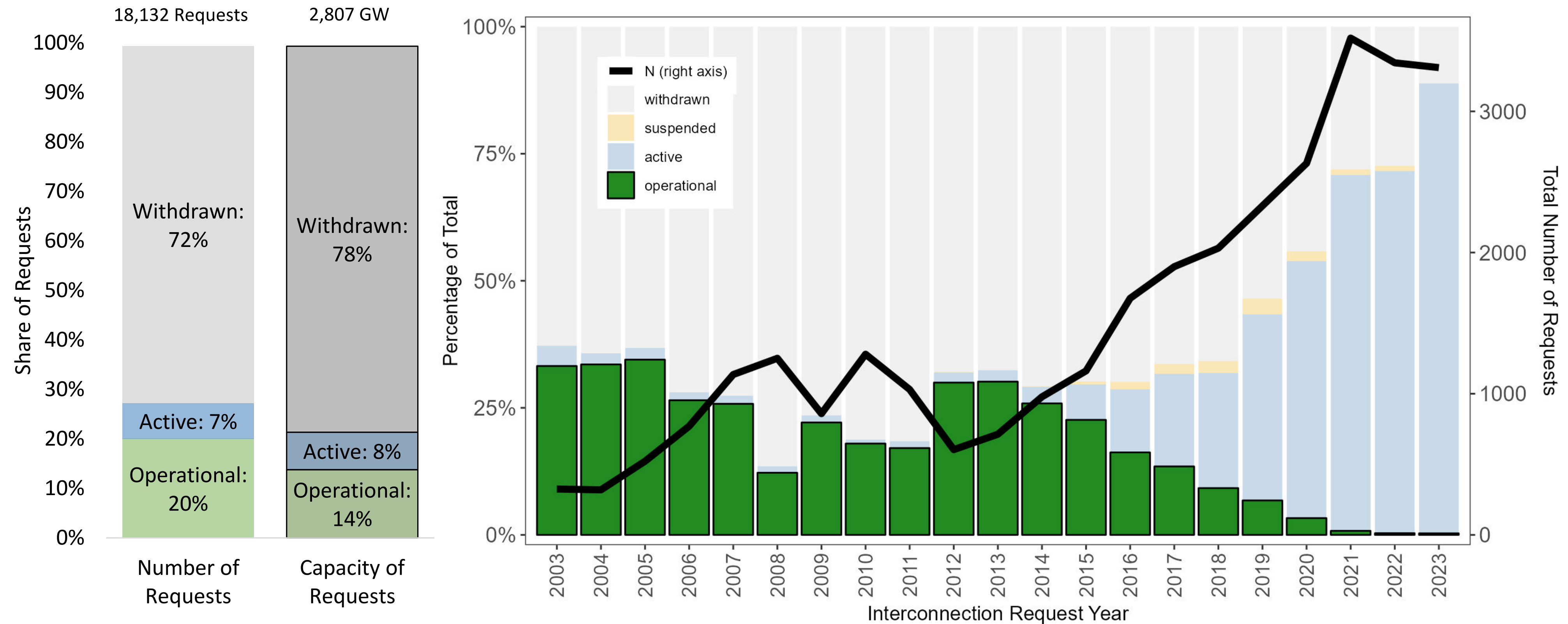
# AI electricity demand growth outpaces power grid development



Lawrence Berkeley NL. 2024 Edition Characteristics of Power Plants Seeking Transmission Interconnection As of the End of 2023.

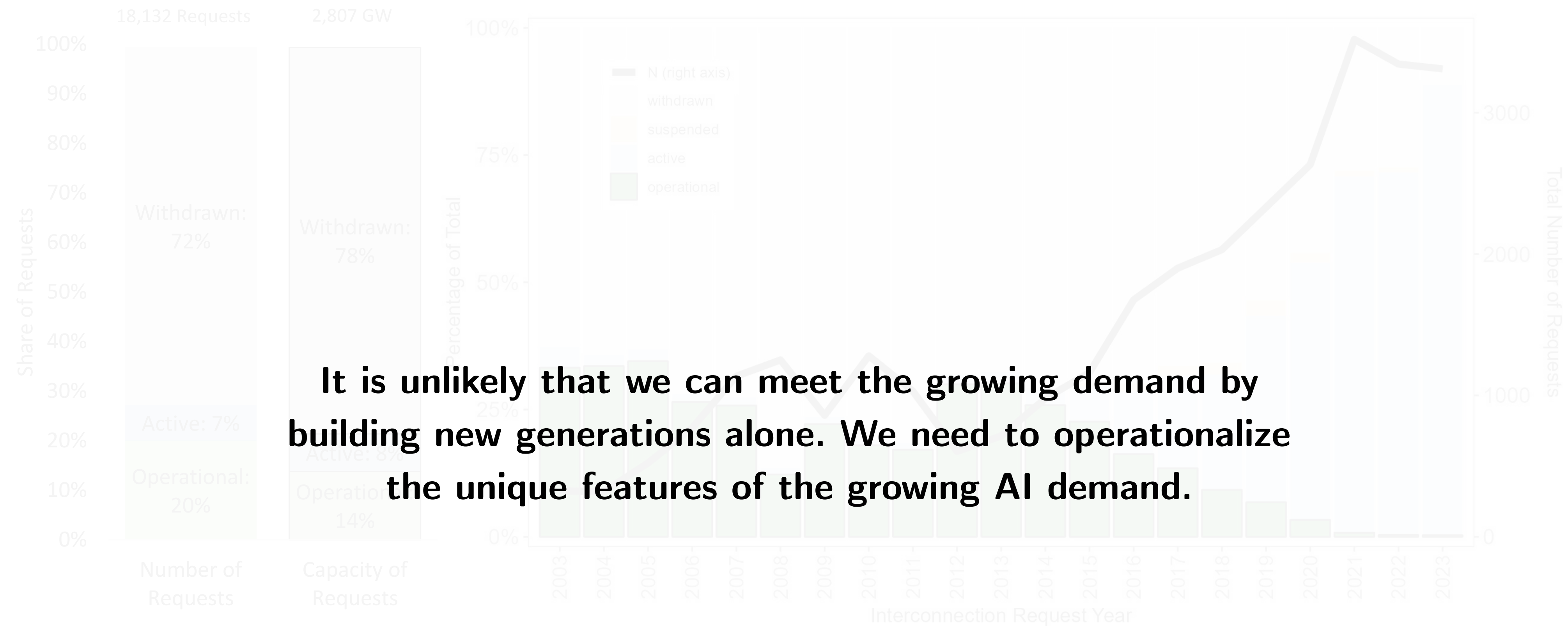
- ▶ Grid operators require generation projects to undergo a series of impact studies before they can be built
- ▶ The total capacity in the queues is growing year-over-year, but the impact studies do not keep up

# AI electricity demand growth outpaces power grid development



Lawrence Berkeley NL. 2024 Edition Characteristics of Power Plants Seeking Transmission Interconnection As of the End of 2023.

- ▶ The majority (> 70%) of interconnection requests are withdrawn
- ▶ Just 20% of requests (14% of capacity) submitted from 2000-2018 had been built as of the end of 2023
- ▶ Power grids expansion does not keep with the pace of AI electricity growth



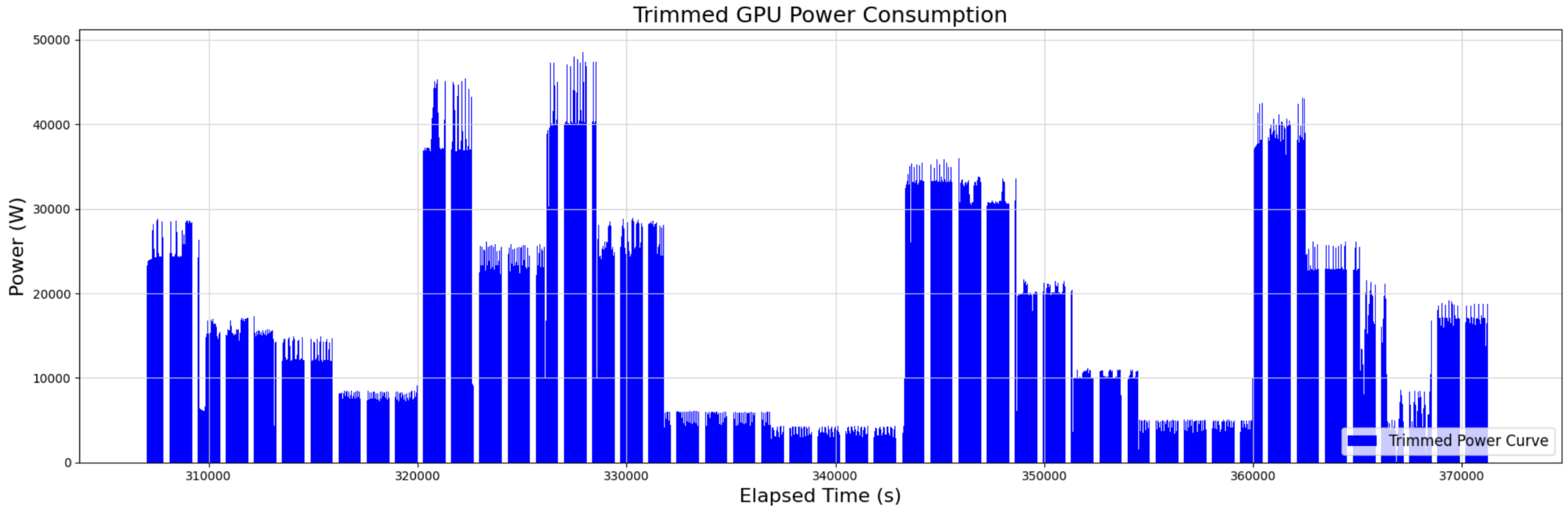
**It is unlikely that we can meet the growing demand by building new generations alone. We need to operationalize the unique features of the growing AI demand.**

Lawrence Berkeley NL. 2024 Edition Characteristics of Power Plants Seeking Transmission Interconnection As of the End of 2023.

- ▶ The majority (> 70%) of interconnection requests are withdrawn
- ▶ Just 20% of requests (14% of capacity) submitted from 2000-2018 had been built as of the end of 2023
- ▶ Power grids expansion does not keep with the pace of AI electricity growth

## Disruptive electricity consumption patterns

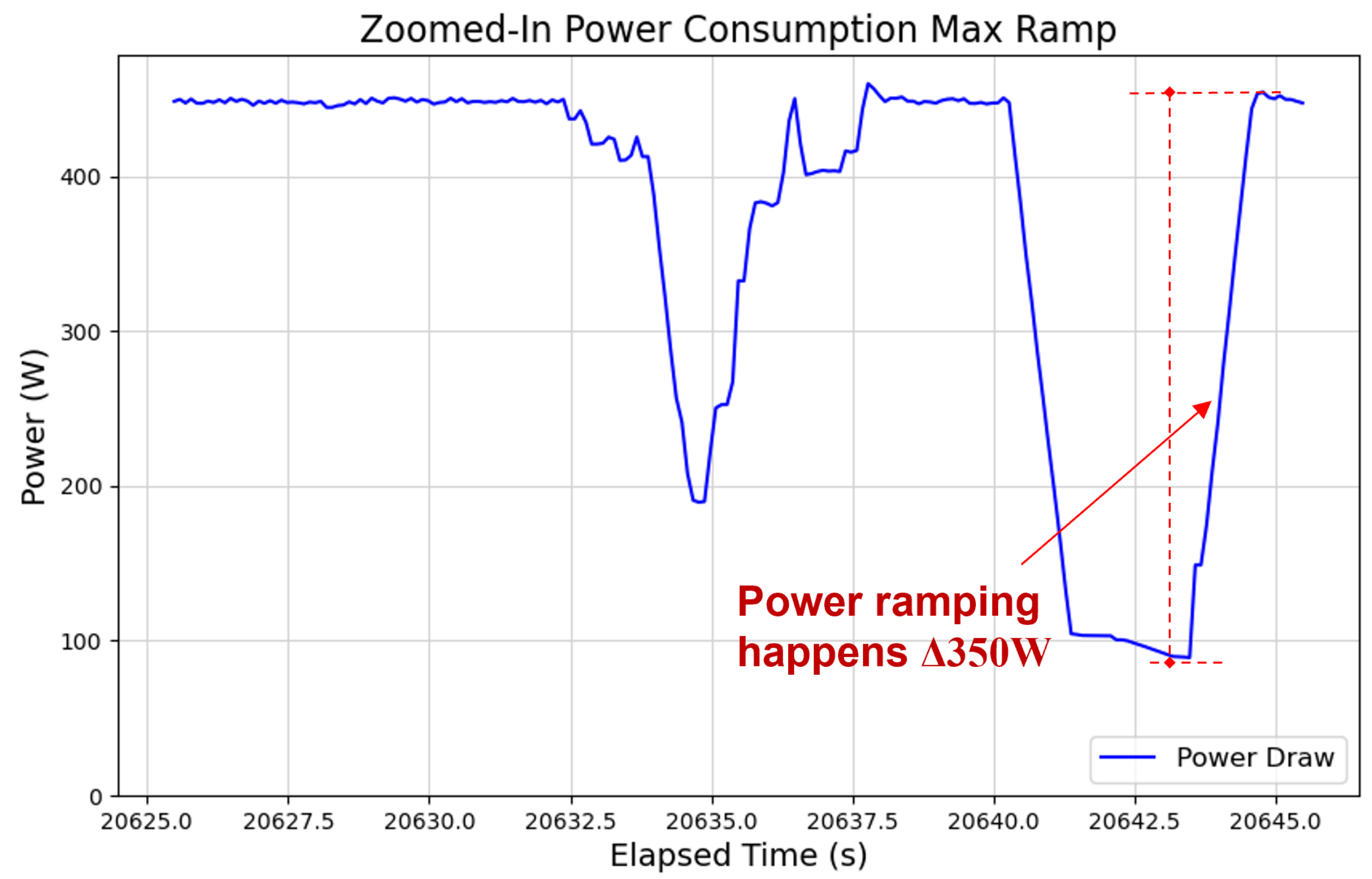
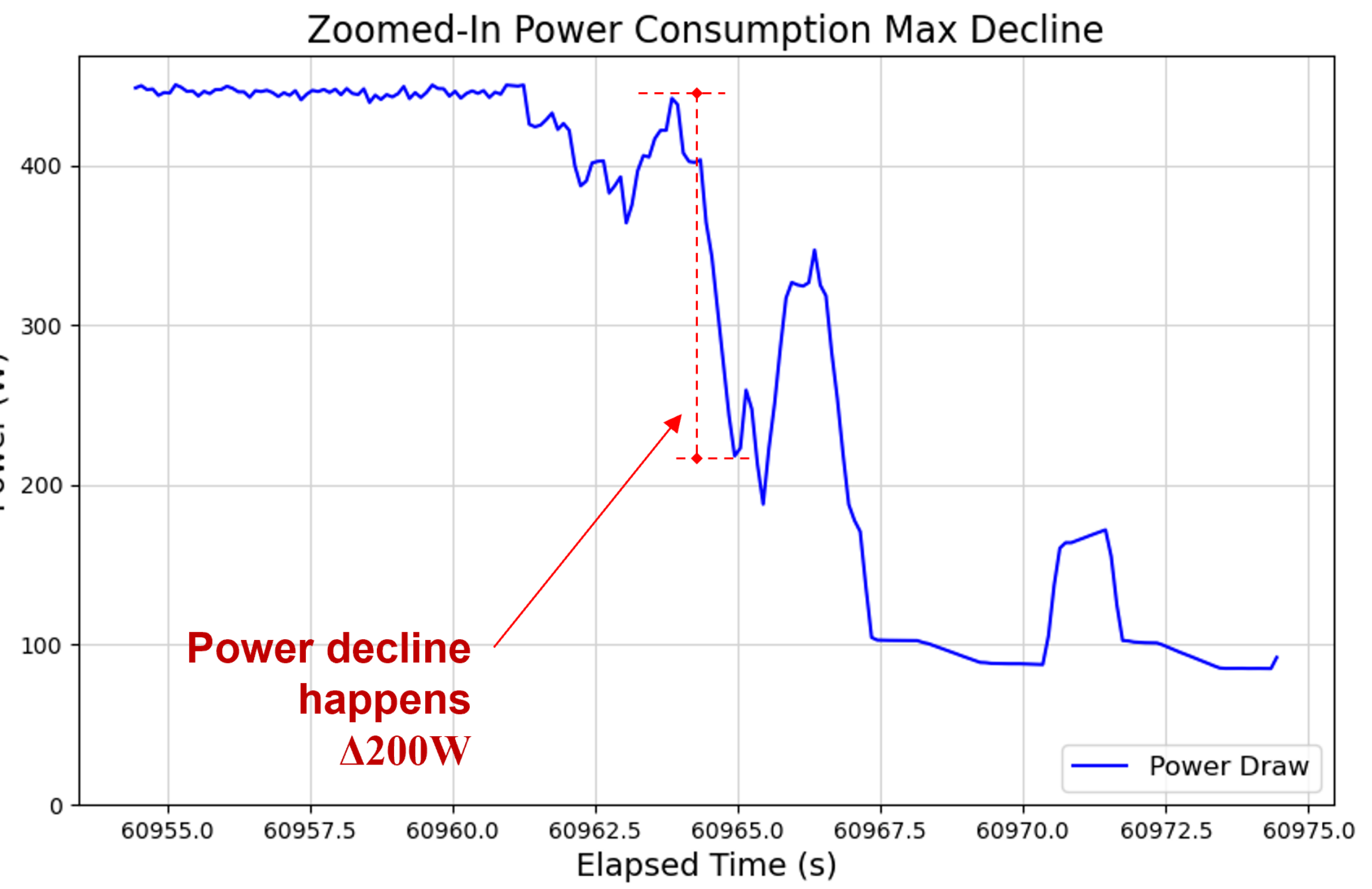
## Disruptive electricity consumption patterns



Power Consumption of BERT in MIT Supercloud Dataset



## Disruptive electricity consumption patterns

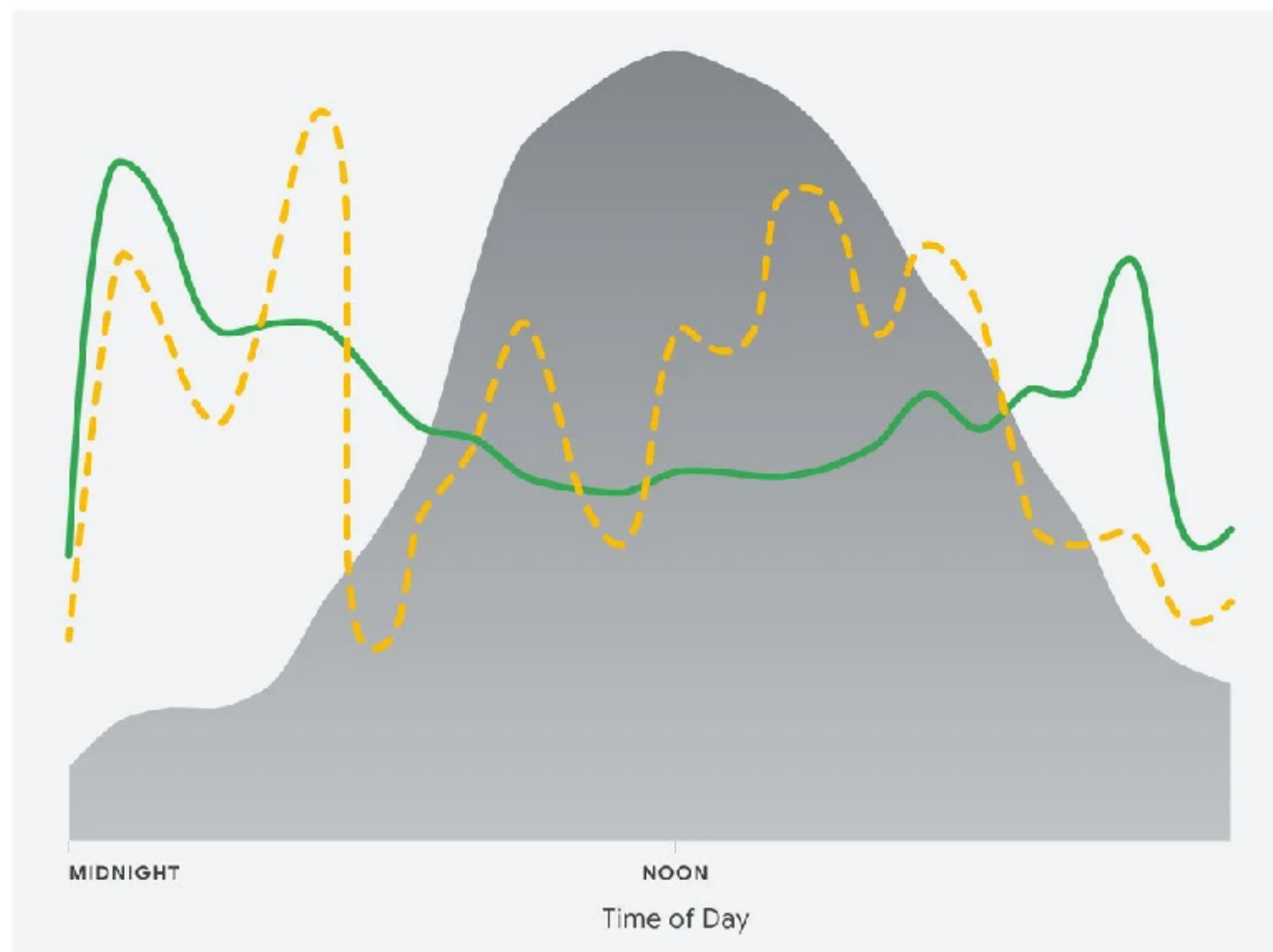


Power Transients of GPT-2 124M

## Temporal and **geospatial** load flexibility

### Baseline versus Carbon-aware Load

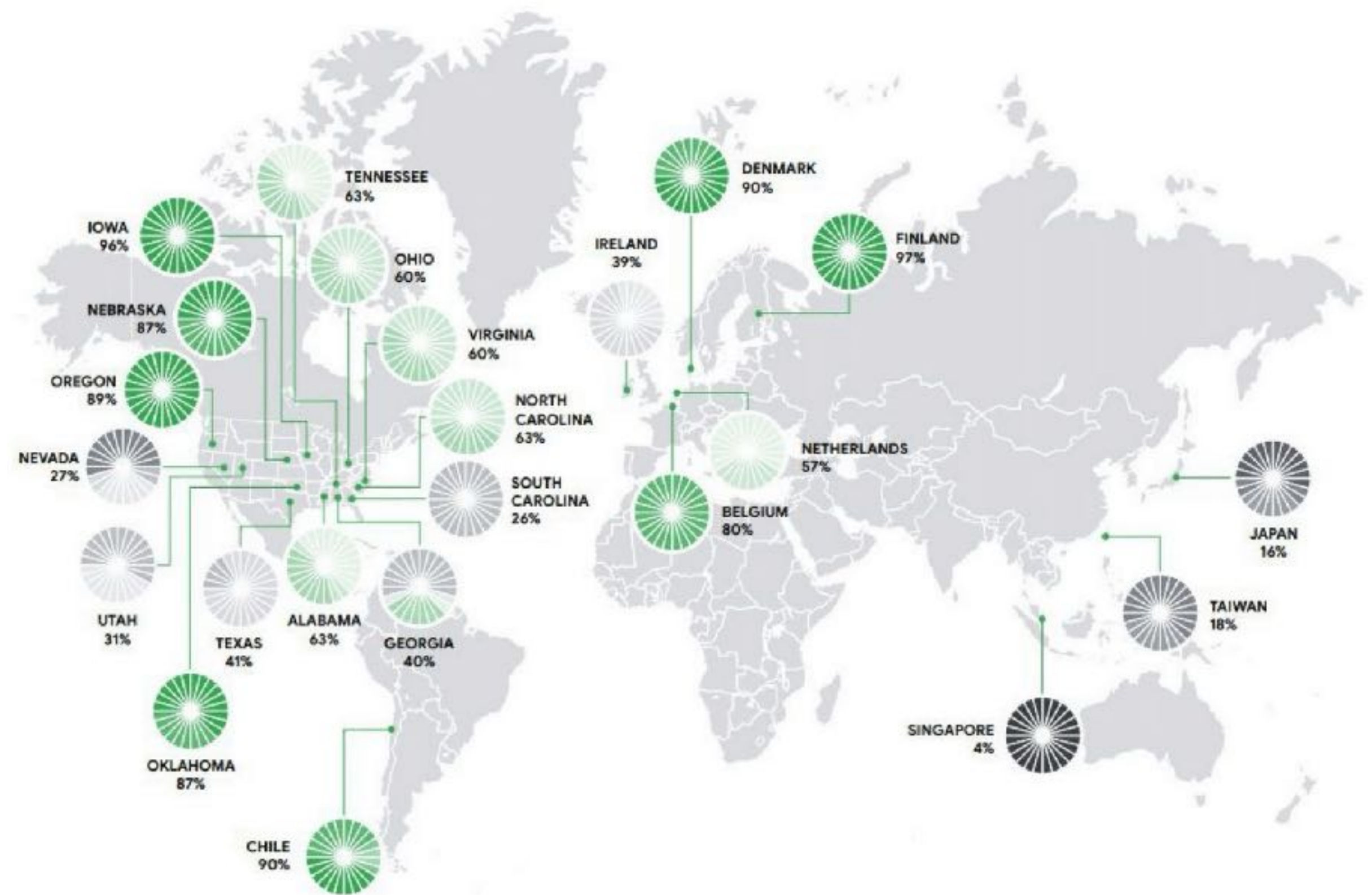
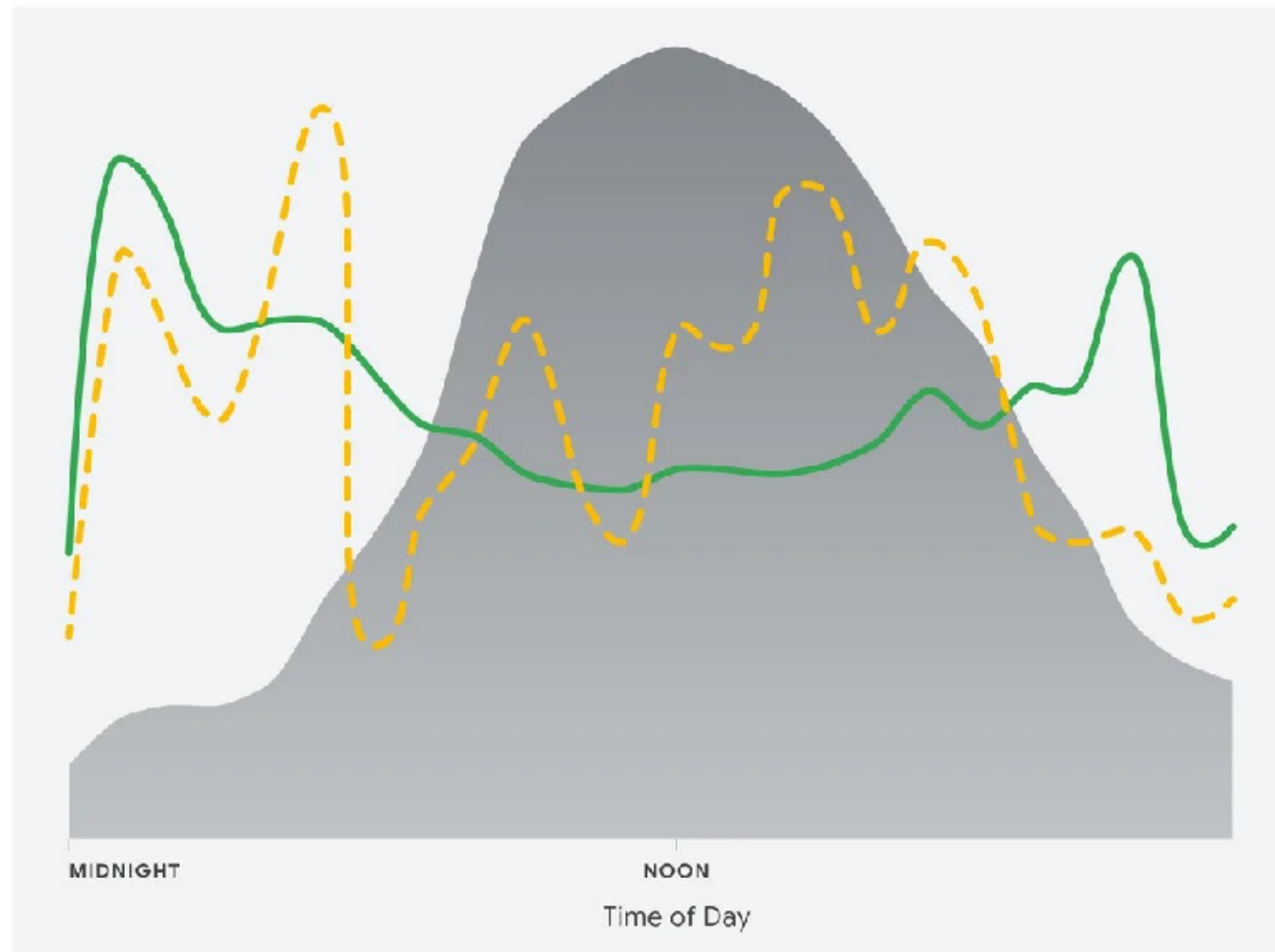
--- Baseline Load    — Carbon-aware Load    ● Carbon Intensity



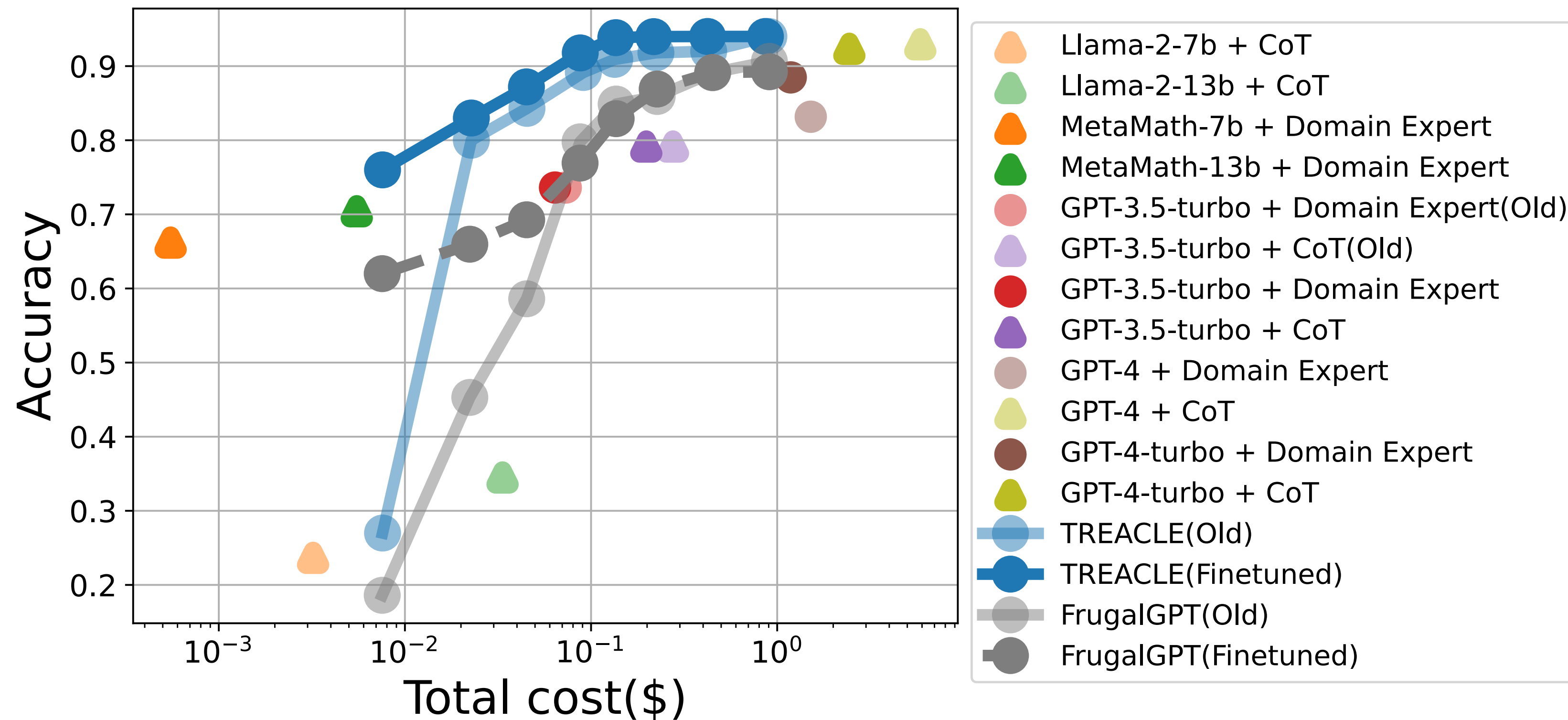
## Temporal and geospatial load flexibility

### Baseline versus Carbon-aware Load

--- Baseline Load    — Carbon-aware Load    ● Carbon Intensity

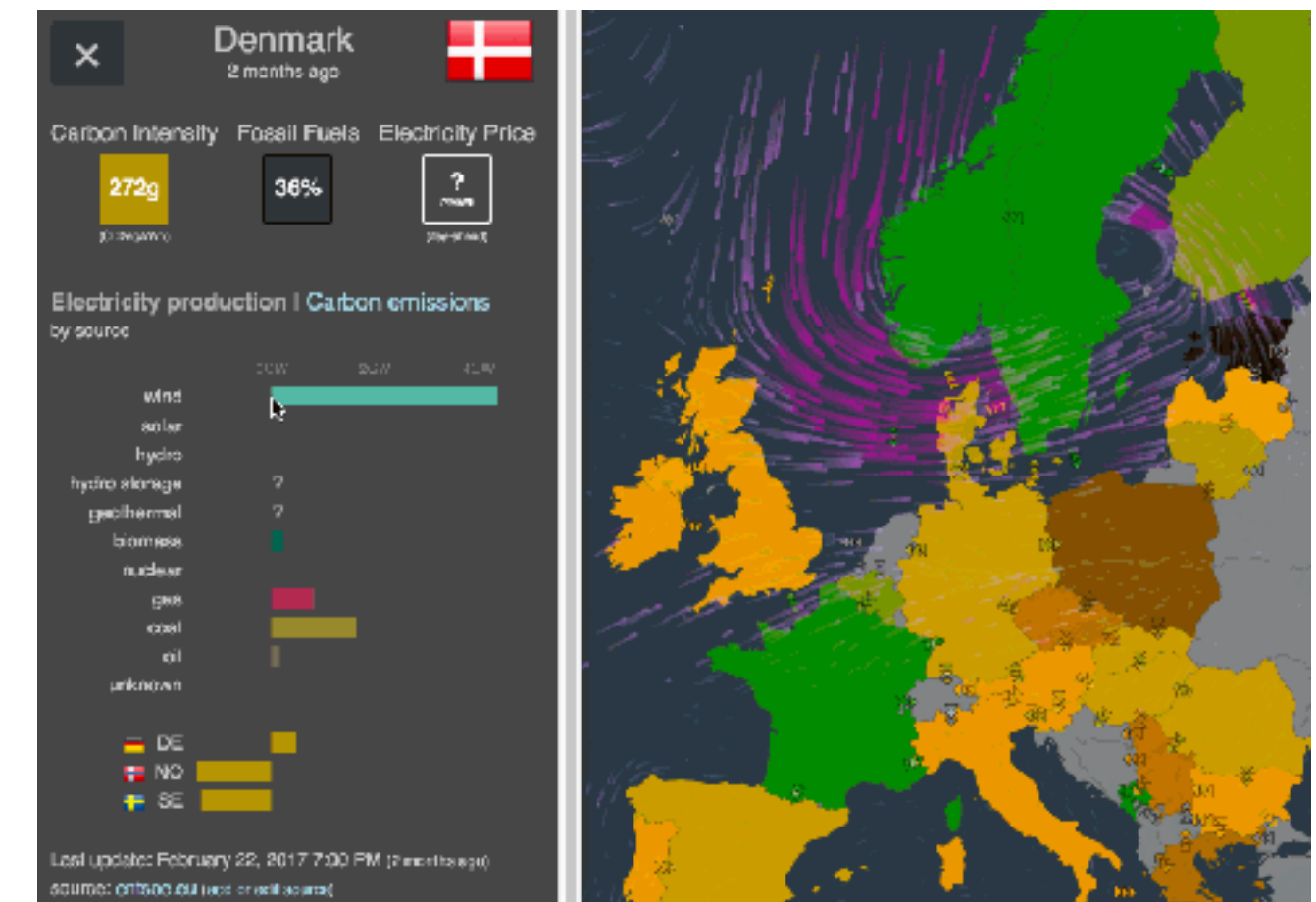


LLM service can be instantaneously delivered by multiple vendors:  
random job allocation in the grid, depending on user preferences.



**How to coordinate power grids and data centers,  
considering such a disruptive nature of AI loads?**

- ▶ Data center and grid co-design<sup>1</sup>:
  - ▶ Optimal sizing, siting and timing of data centers in the grid
  - ▶ Delivers benefits in the long run, but not in the short run
- ▶ New electricity market design<sup>2</sup>:
  - ▶ Definition of new market products is difficult (e.g., pricing geospatial shifts)
  - ▶ Mostly bilateral, out of market contracts. Microsoft & Three Mile Island
- ▶ Demand response participation<sup>3</sup>
  - ▶ Currently a working solution, yet with limited scalability
- ▶ Co-optimization of grid and data-center operations:
  - ▶ Ideal yet unattainable in practice solution
  - ▶ Significant **privacy** concerns and computational requirements



<https://app.electricitymaps.com/map>

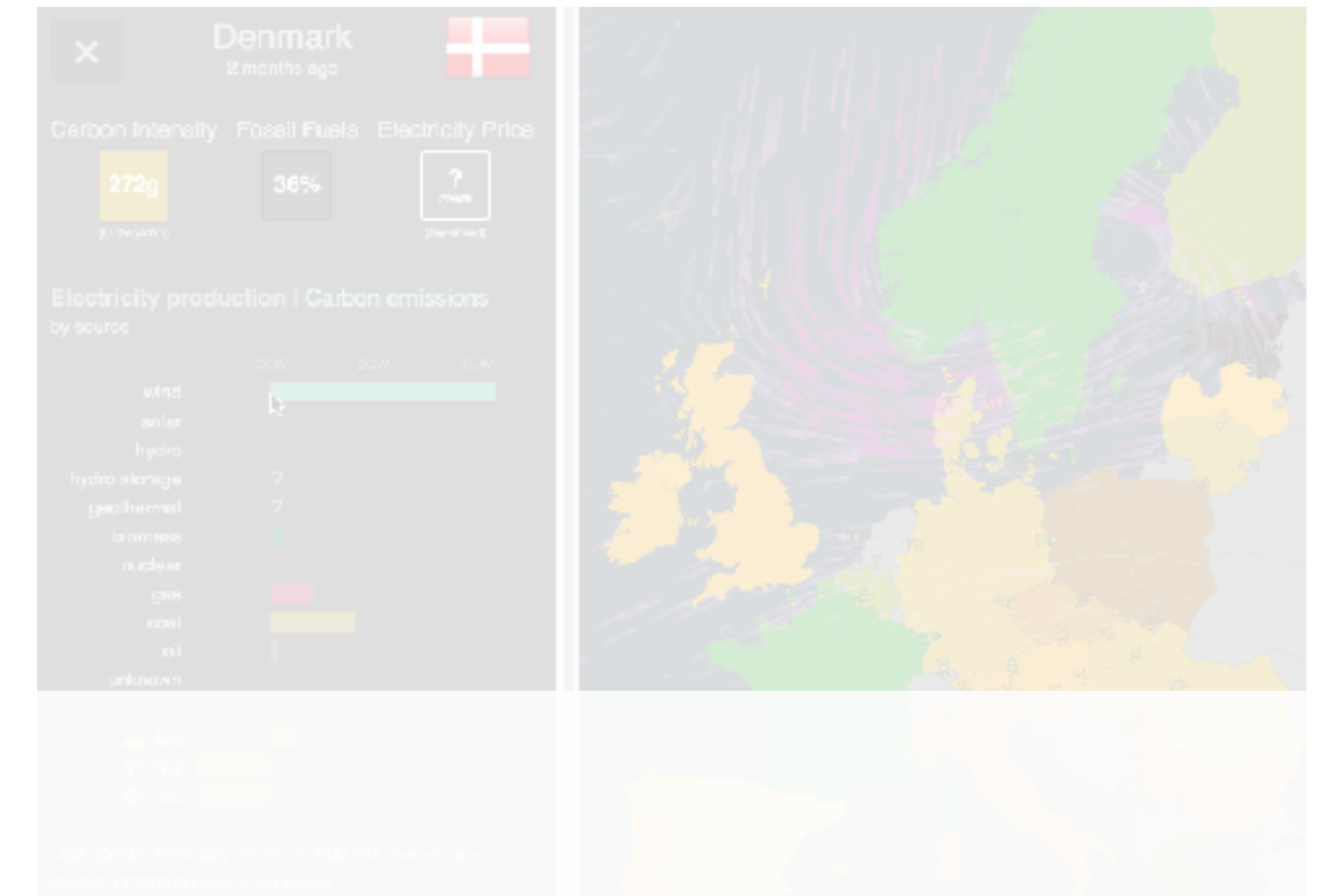


<sup>1</sup>Y. Abdennadher et al. Carbon Efficient Placement of Data Center Locations. 2022

<sup>2</sup>W. Zhang et al. Flexibility from networks of data centers: A market clearing formulation with virtual links. 2020

<sup>3</sup><https://cloud.google.com/blog/products/infrastructure/using-demand-response-to-reduce-data-center-power-consumption>

- ▶ Data center and grid co-design<sup>1</sup>:
  - ▶ Optimal sizing, siting and timing of data centers in the grid
  - ▶ Delivers benefits in the long run, but not in the short run
- ▶ New electricity market design<sup>2</sup>:
  - ▶ Definition of new market products is difficult (e.g., pricing geospatial shifts)
  - ▶ Mostly bilateral, out of market contracts. Microsoft & Three Mile Island
- ▶ Demand response participation<sup>3</sup>
  - ▶ Currently a working solution, yet with limited scalability
- ▶ Co-optimization of grid and data-center operations:
  - ▶ Ideal yet unattainable in practice solution
  - ▶ Significant **privacy** concerns and computational requirements



<https://app.electricitymaps.com/map>



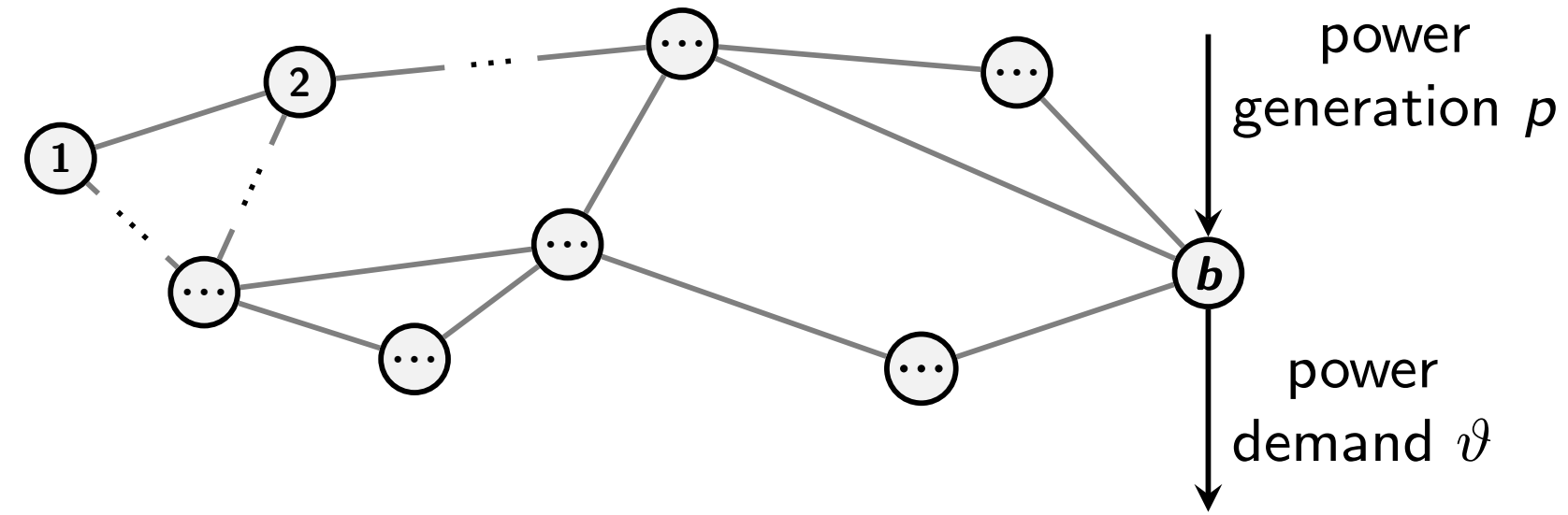
<sup>1</sup>Y. Abdennadher et al. Carbon Efficient Placement of Data Center Locations. 2022

<sup>2</sup>W. Zhang et al. Flexibility from networks of data centers: A market clearing formulation with virtual links. 2020

<sup>3</sup><https://cloud.google.com/blog/products/infrastructure/using-demand-response-to-reduce-data-center-power-consumption>

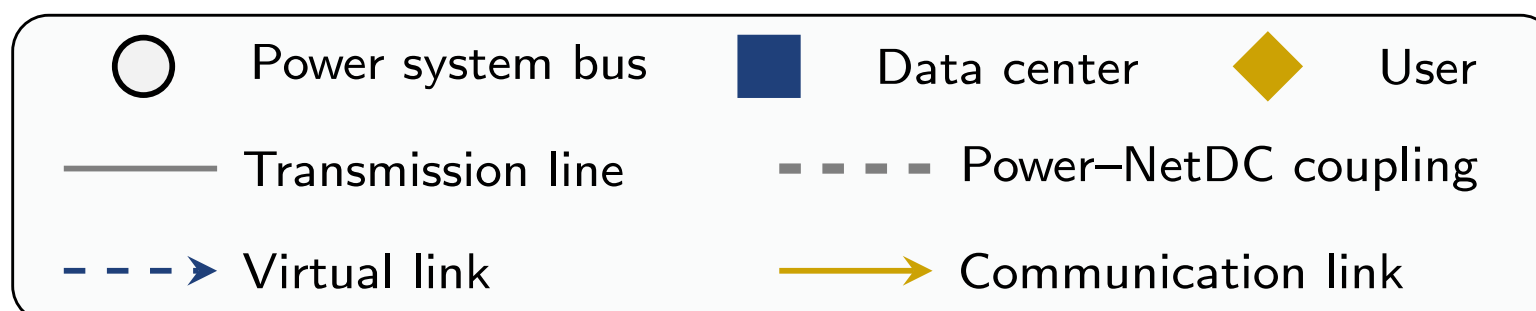
# **Co-optimization of grid and data-center operations**

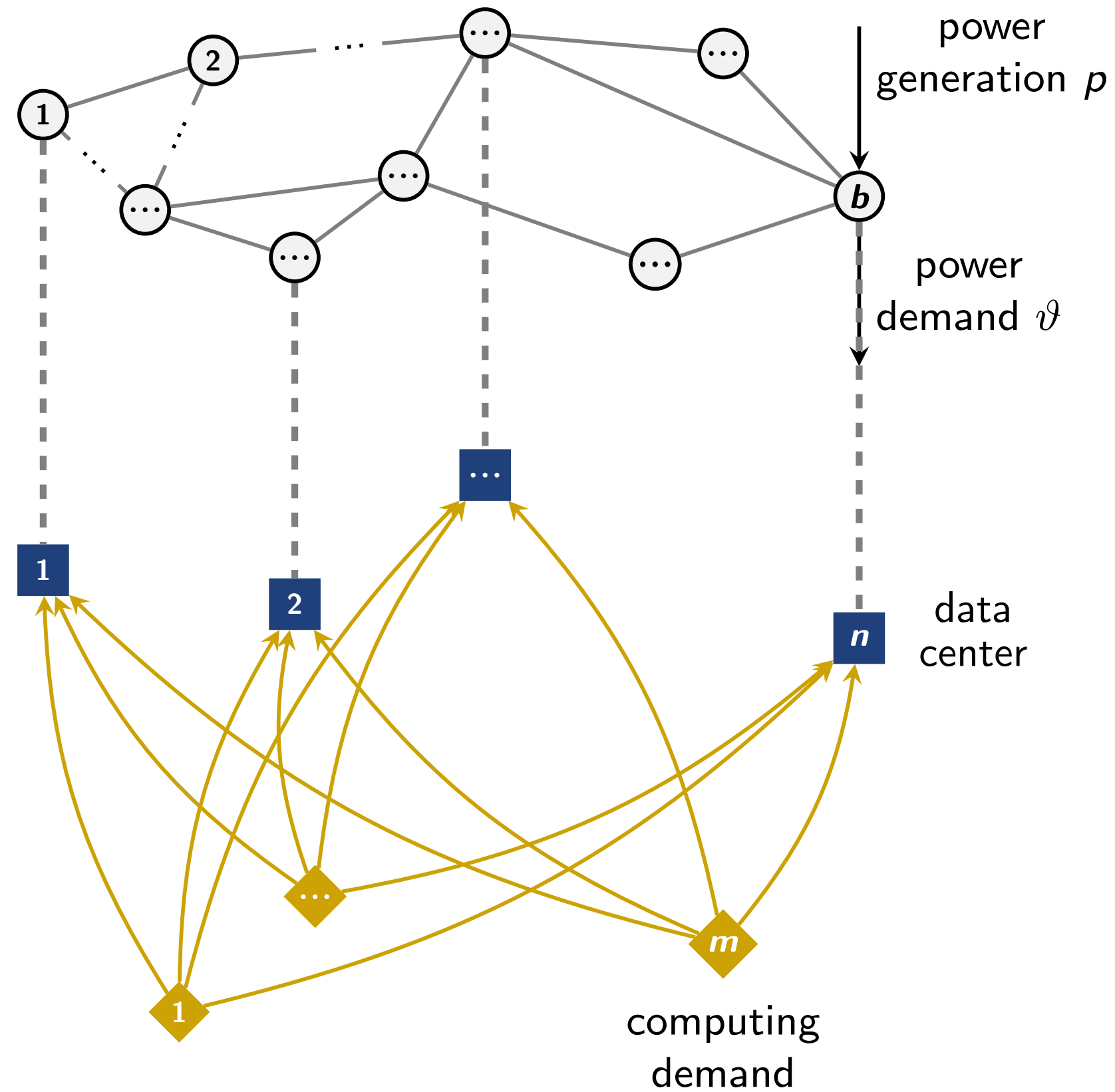




## Power grid optimization problem:

$$\begin{aligned}
 & \underset{p}{\text{minimize}} && c_{\text{pwr}}(p) && \triangleright \text{Dispatch cost} \\
 & \text{subject to} && p \in \mathcal{P}_{\text{pwr}}(v) && \triangleright \text{Grid feasibility}
 \end{aligned}$$



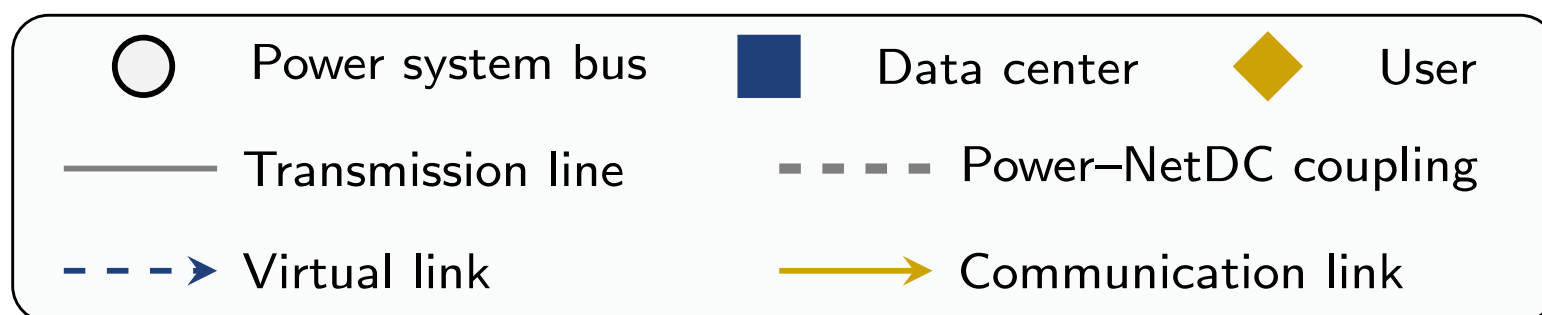
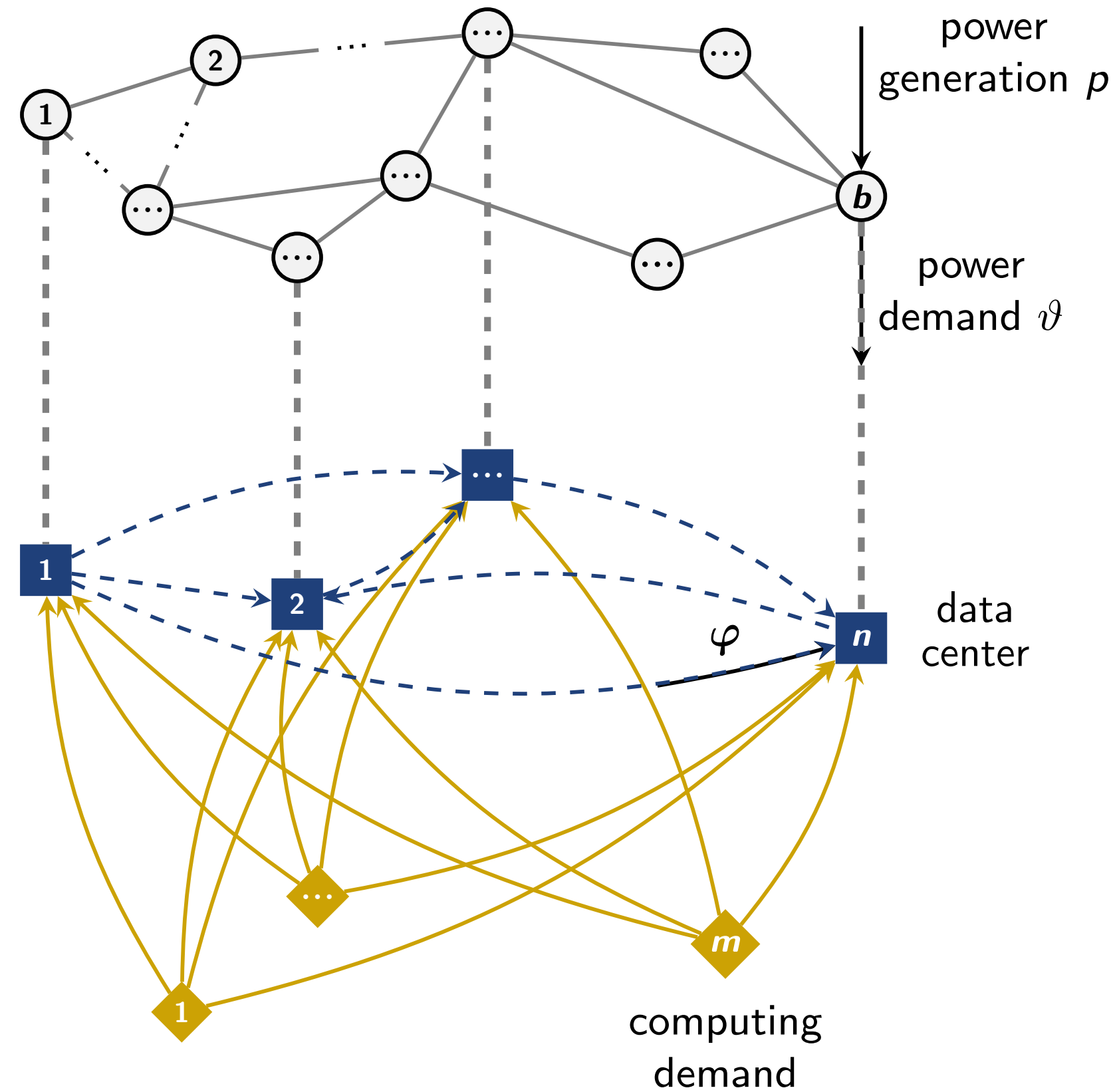


## Power grid optimization problem:

$$\begin{aligned} & \underset{p}{\text{minimize}} && c_{\text{pwr}}(p) && \triangleright \text{Dispatch cost} \\ & \text{subject to} && p \in \mathcal{P}_{\text{pwr}}(\vartheta) && \triangleright \text{Grid feasibility} \end{aligned}$$

## Data centers optimization problem:

$$\begin{aligned} & \underset{\vartheta}{\text{minimize}} && c_{\text{net-dc}}(\vartheta) && \triangleright \text{Latency loss} \\ & \text{subject to} && \vartheta \in \mathcal{W}_{\text{net-dc}}(\varphi) && \triangleright \text{NetDC feasibility} \end{aligned}$$

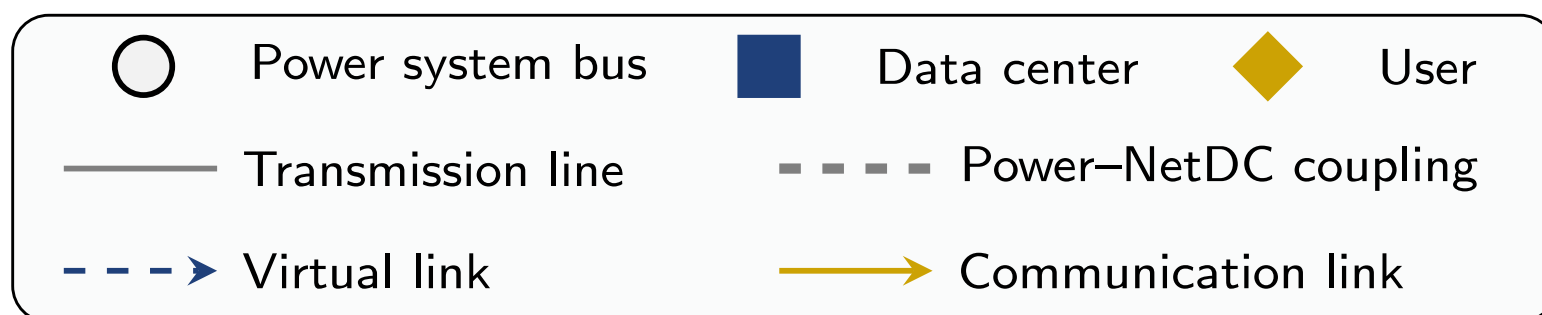
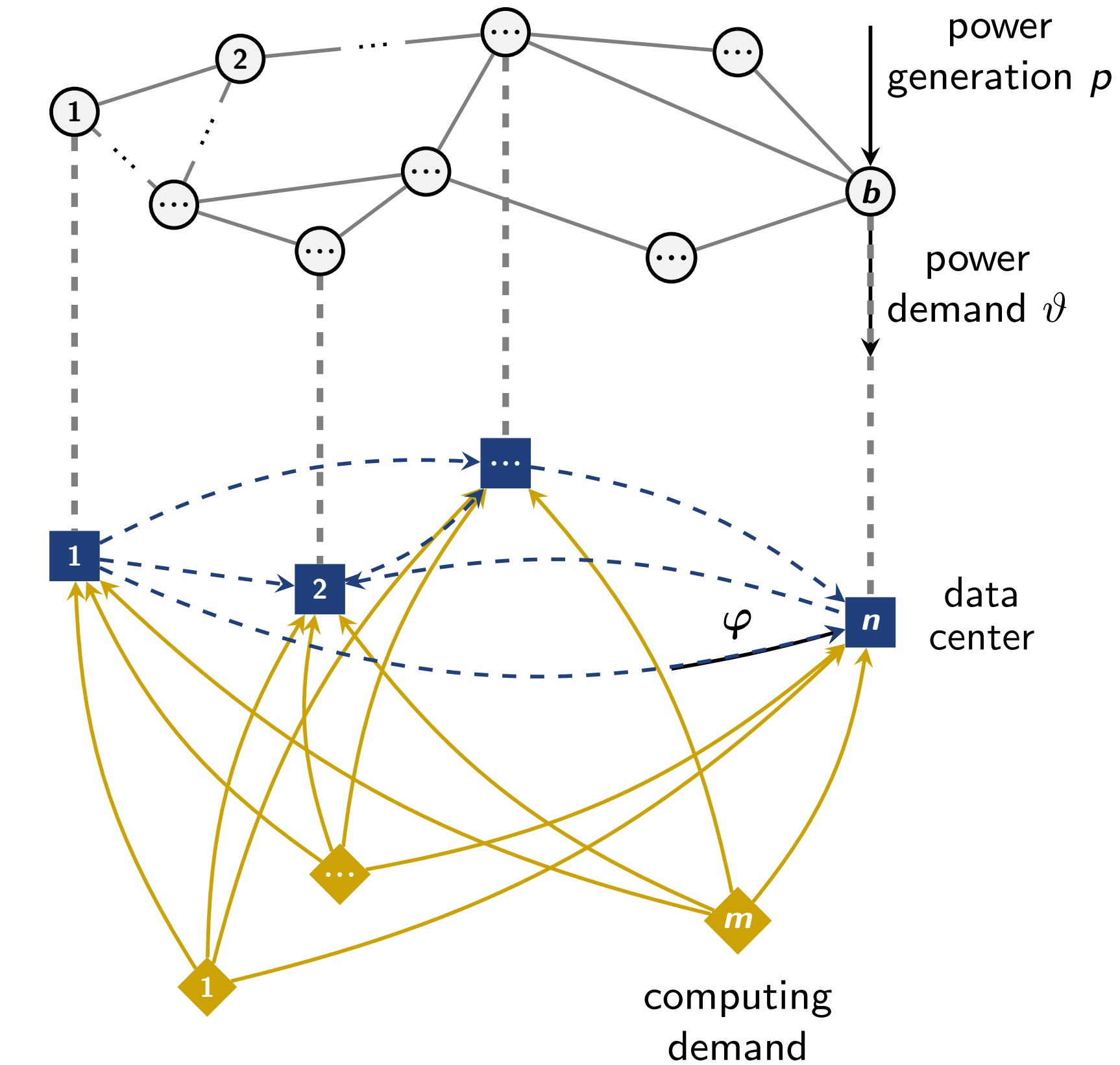


## Power grid optimization problem:

$$\begin{aligned} & \underset{p}{\text{minimize}} && c_{\text{pwr}}(p) && \triangleright \text{Dispatch cost} \\ & \text{subject to} && p \in \mathcal{P}_{\text{pwr}}(\vartheta) && \triangleright \text{Grid feasibility} \end{aligned}$$

## Data centers optimization problem:

$$\begin{aligned} & \underset{\vartheta}{\text{minimize}} && c_{\text{net-dc}}(\vartheta) && \triangleright \text{Latency loss} \\ & \text{subject to} && \vartheta \in \mathcal{W}_{\text{net-dc}}(\varphi) && \triangleright \text{NetDC feasibility} \end{aligned}$$



## Power Grid-NetDC coordination problem:

$$\begin{aligned}
 & \underset{\varphi, p}{\text{minimize}} && c_{\text{pwr}}(p) && \triangleright \text{Dispatch cost} \\
 & \text{subject to} && p \in \mathcal{P}_{\text{pwr}}(\vartheta) && \triangleright \text{Grid feasibility} \\
 & && \text{el. demand} && \\
 & && \leftarrow && \\
 & && \underset{\vartheta}{\text{minimize}} && c_{\text{net-dc}}(\vartheta) && \triangleright \text{Latency loss} \\
 & && \text{subject to} && \vartheta \in \mathcal{W}_{\text{net-dc}}(\varphi) && \triangleright \text{NetDC feasibility} \\
 & && \uparrow && \text{computing task shift} && 
 \end{aligned}$$

Bilevel structure: the grid operator minimizes dispatch costs by optimizing task shifts that reshape data-center electricity demand.

# **Power-NetDC coordination: From optimization to regression**

- ▶ Solving the bilevel problem in real-time is extremely challenging:
  - ▶ Power and NetDC needs to be exchanged in real time (privacy barriers)
  - ▶ Large-scale bilevel optimization in real time (computationally intractable)
- ▶ Instead, we consider a contextual regression policy for real-time coordination

$$\phi(x) = \beta_0 + \beta_1 x$$

where  $\beta = (\beta_0, \beta_1)$  are regression weights and  $x$  is the vector of real-time coordination features (e.g, locational marginal prices, zonal electricity demand, renewable power generation)

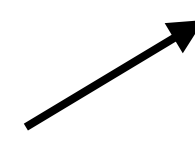
- ▶ The baseline training of the regression policy:
  1. Collect a dataset  $\{(x_1, \varphi_1^*), \dots, (x_q, \varphi_q^*)\}$  of  $q$  coordination scenarios.

Each scenario  $i$  includes contextual features  $x_i$  and the optimal solution  $\varphi_i^*$  to the coordination problem

2. Train a contextual regression to map coordination features into the optimal task shifts in real-time

minimize  $\frac{1}{q} \sum_{i=1}^q \|\beta_0 + \beta_1 x_i - \varphi_i^*\|_2^2$ 
← *Minimum prediction loss, but no guarantees on policy cost-optimality and feasibility*

feature selection



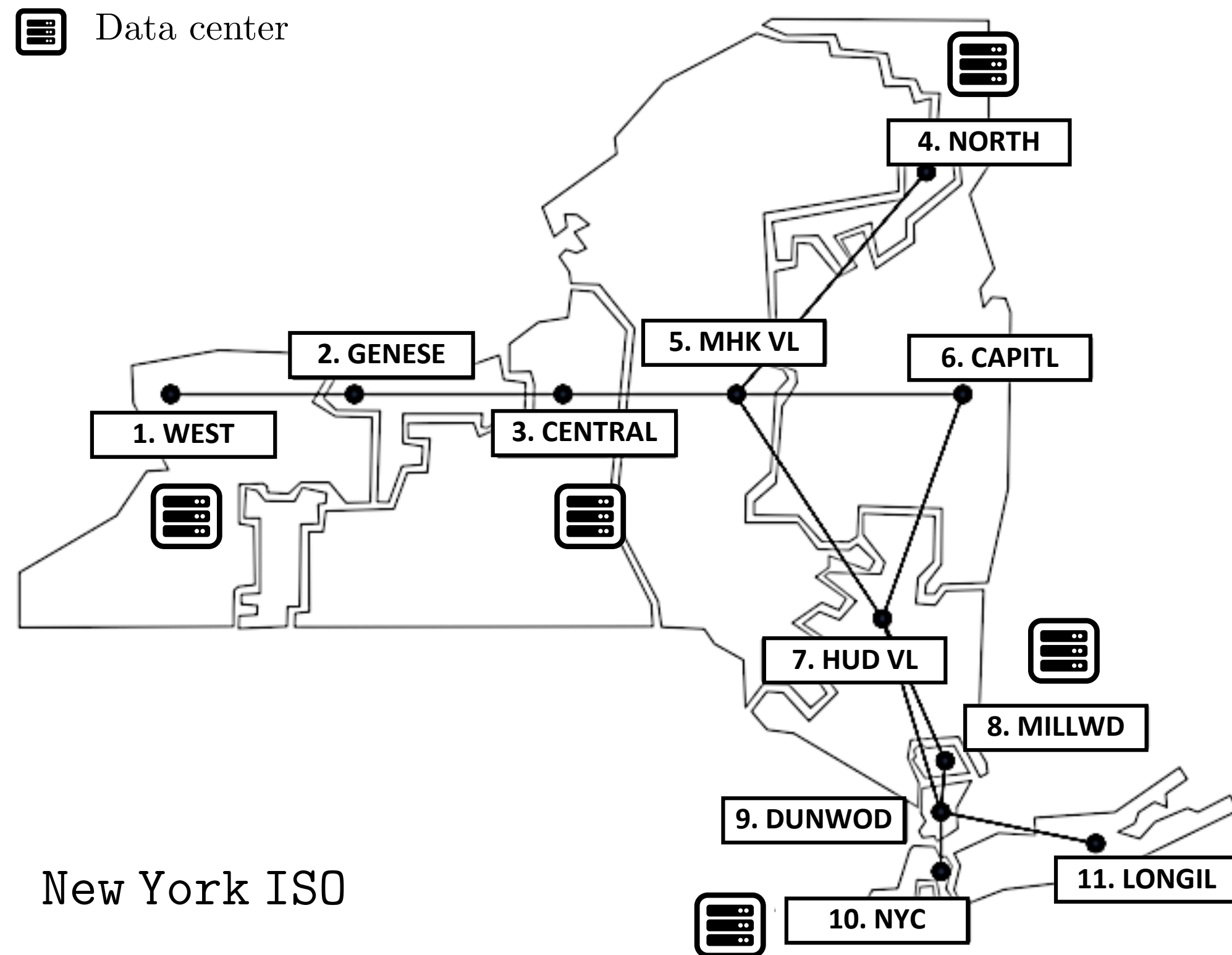
**AgentCONCUR protocol with performance guarantees**

$$\begin{aligned}
 & \underset{\beta, \varphi_1, \dots, \varphi_q, p_1, \dots, p_q}{\text{minimize}} && \frac{1}{q} \sum_{i=1}^q c_{\text{pwr}}(p_i) && \triangleright \text{Average dispatch cost} \\
 & \text{subject to} && p_i \in \mathcal{P}_{\text{pwr}}(\vartheta_i^*), \quad \forall i = 1, \dots, q && \triangleright \text{Grid equations for each scenario} \\
 & && \varphi_i = \beta_0 + \beta_1 x_i, \quad \|\beta\|_1 \leq \varepsilon, \quad \forall i = 1, \dots, q && \triangleright \text{Coupling contextual regression} \\
 & && \vartheta_i^* \in \underset{\vartheta_i}{\text{minimize}} c_{\text{net-dc}}(\vartheta_i) && \triangleright \text{Latency loss} \\
 & && \text{subject to } \vartheta_i \in \mathcal{W}_{\text{net-dc}}(\varphi_i), \quad \forall i = 1, \dots, q && \triangleright \text{NetDC feasibility}
 \end{aligned}$$

- ▶ The task shifts are restricted to the affine policy of contextual features.
- ▶ Optimization anticipates how the affine restriction affects the average OPF costs.
- ▶ Feasibility guarantees on the training set  $\rightarrow$  also holds on the testing set



# **Numerical Experiments on the New York Independent System Operator's System**



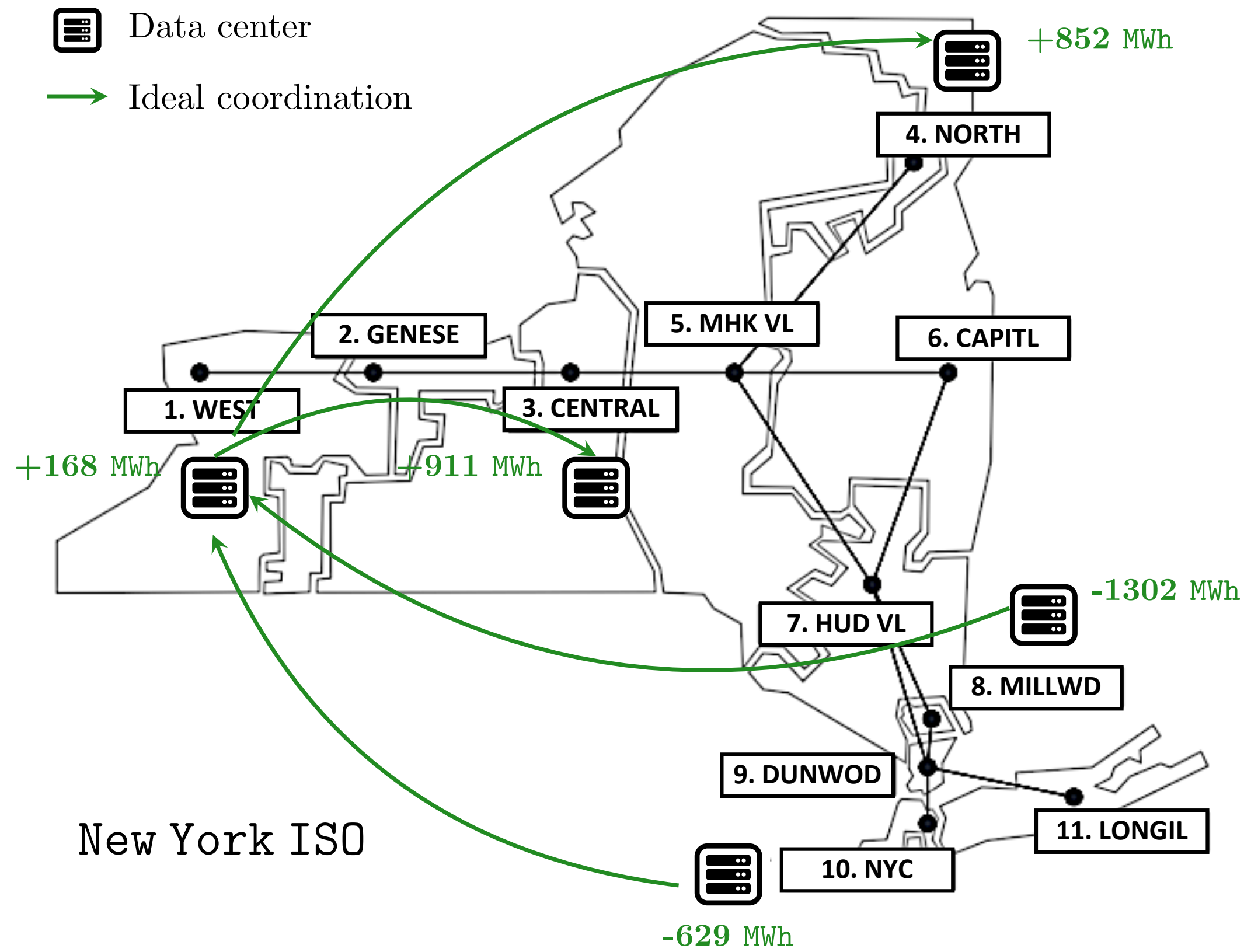
## Data inputs:

- ▶ 11-zone aggregation of the New York ISO
- ▶ Network of 5 data centers (10 virtual links)
- ▶ Varying demand from 5% to 30% of the peak load

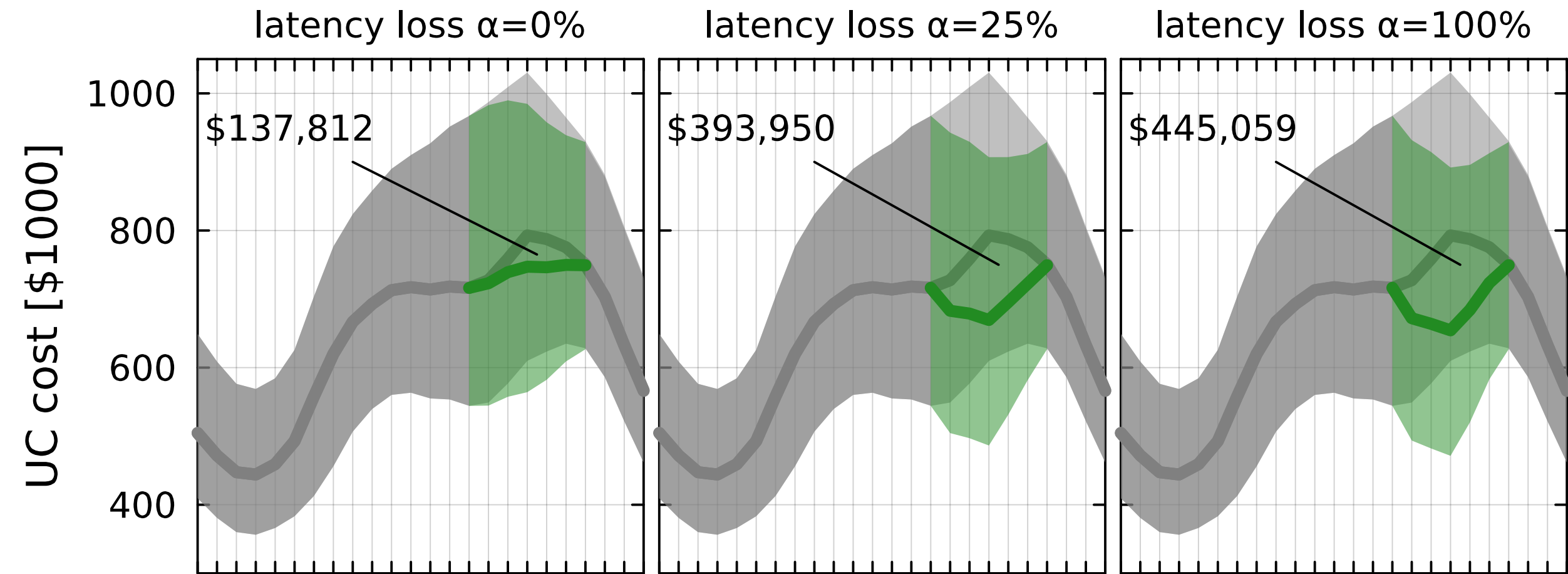
## We study two coordination settings:

- ▶ Ideal day-ahead coordination with optimization
- ▶ Real-time coordination with contextual regression

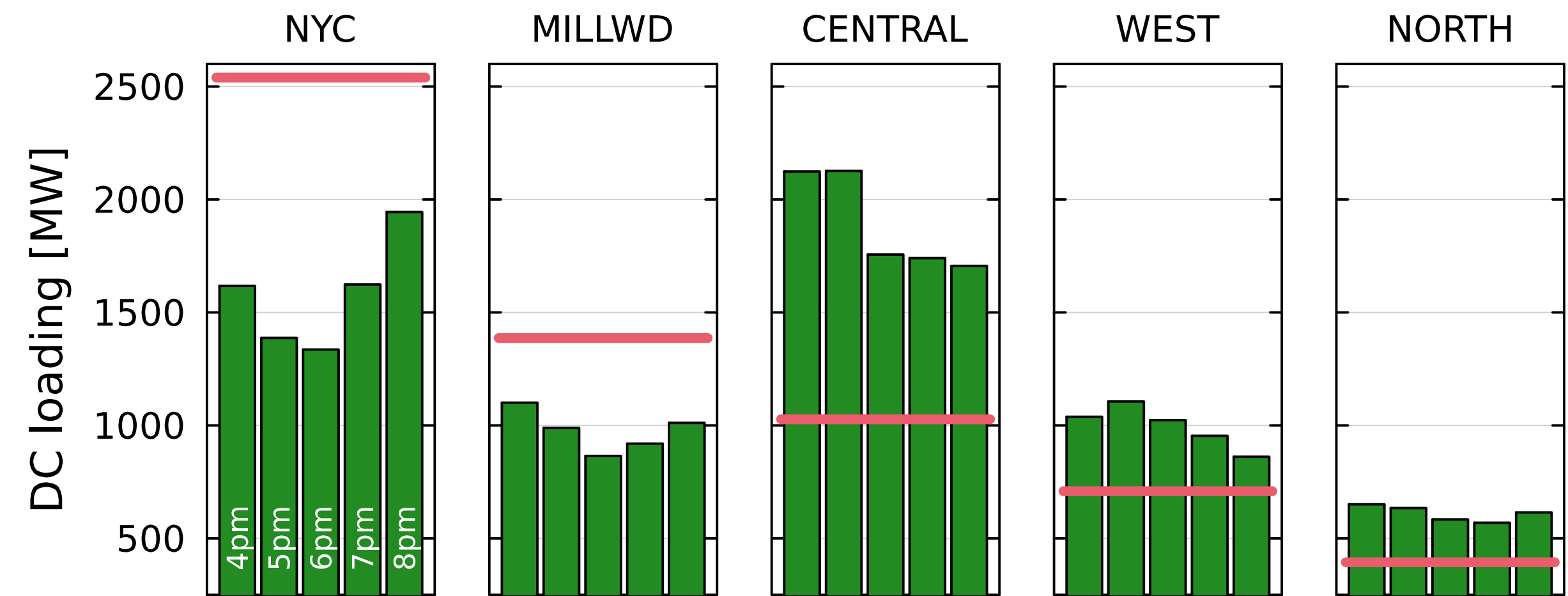
# NYISO: Ideal coordination at the day-ahead stage



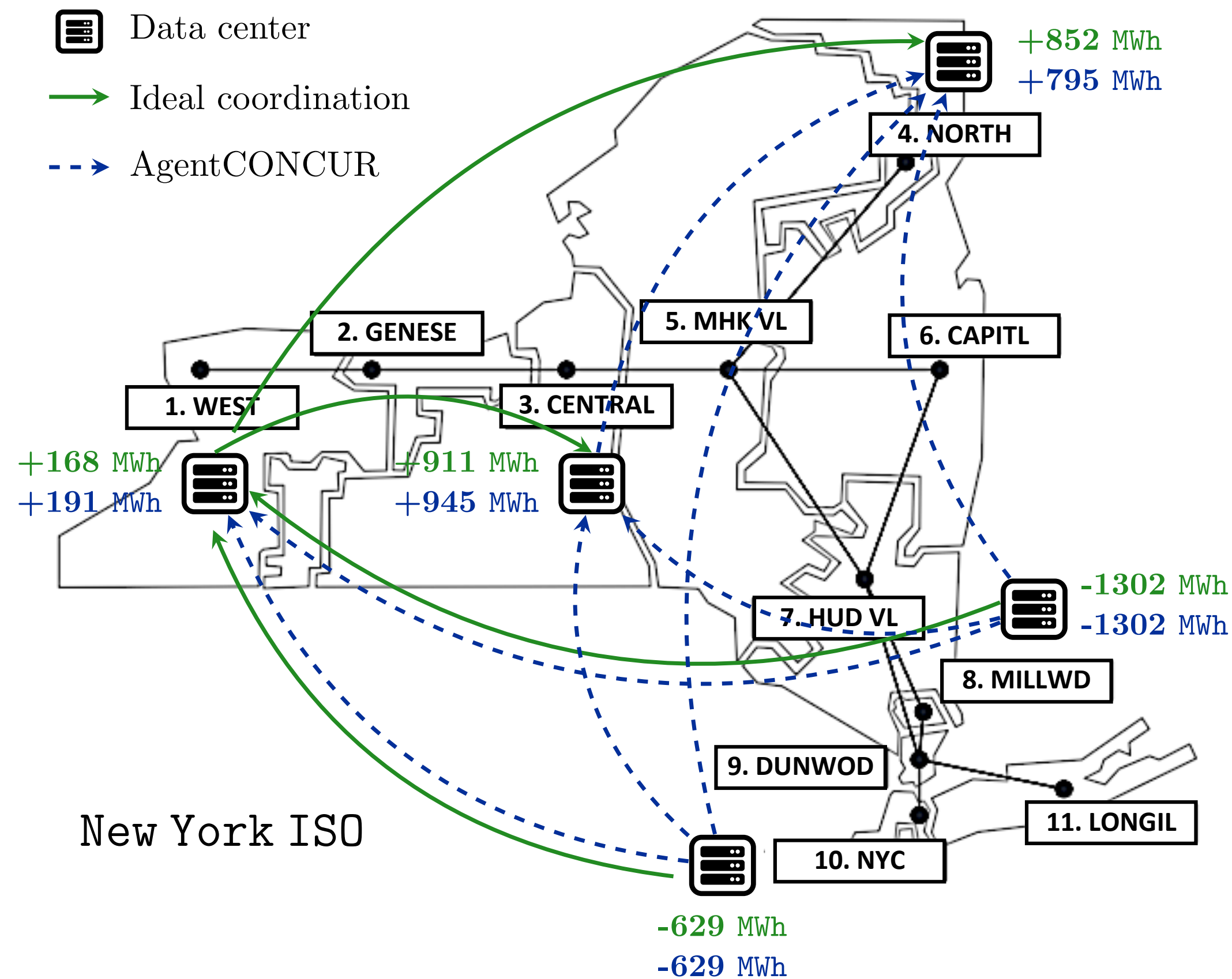
Example of task shifts across the network of data centers



Relaxation of latency constraints  $\Rightarrow$  greater generation cost savings  
 Unit commitment constraints prevent unlocking the whole NetDC flexibility



The flat (in red) loading profile is re-distributed in space and time



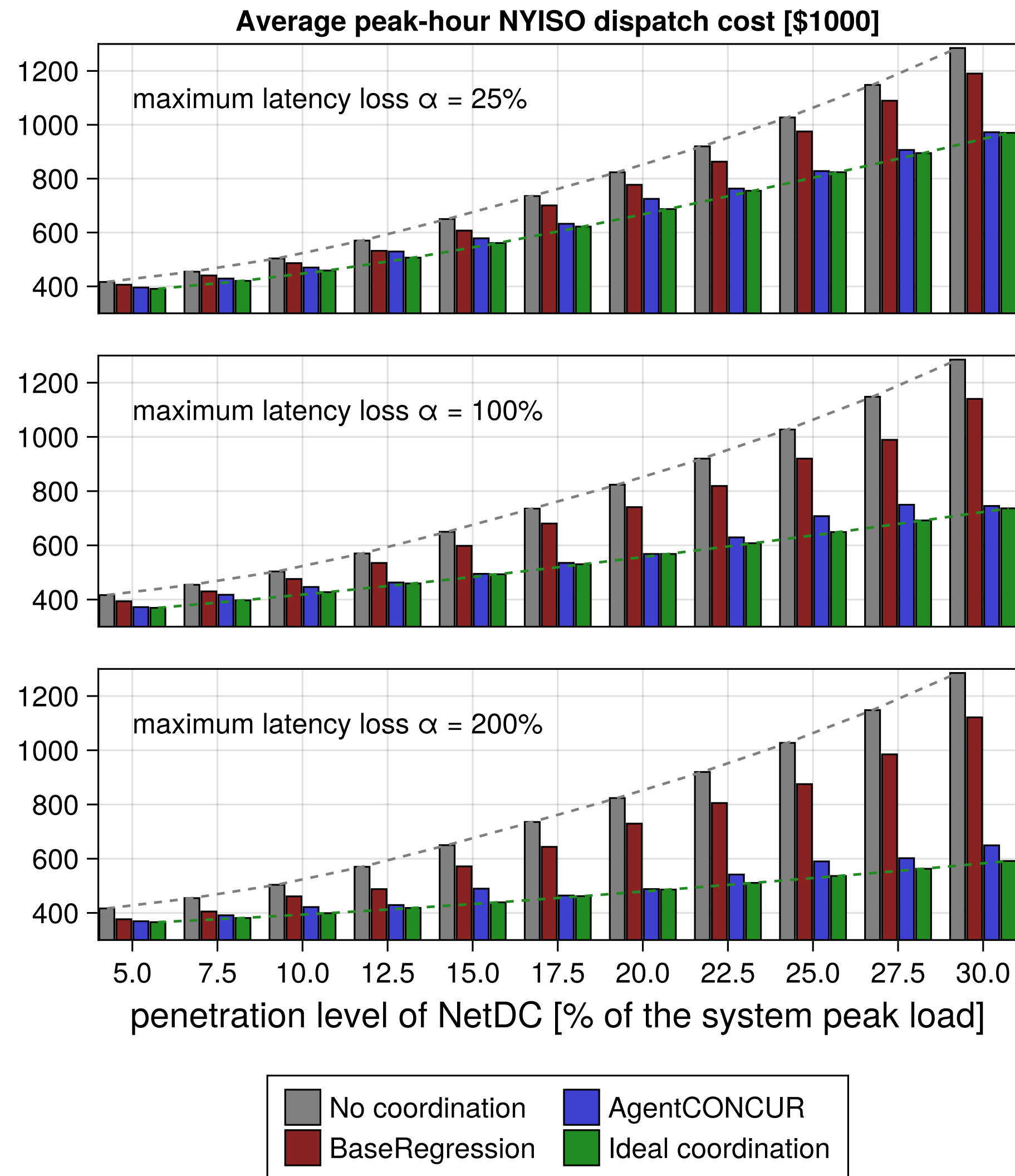
Ideal coordination versus the AgentCONCUR solution

## Contextual features from NYISO website

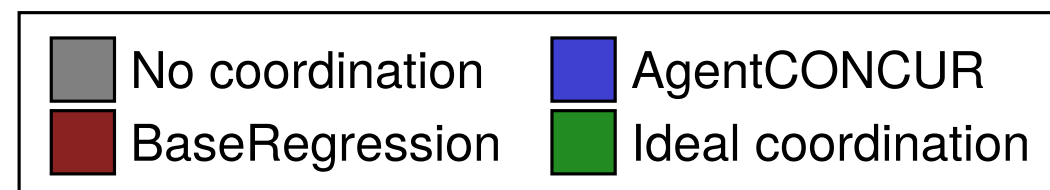
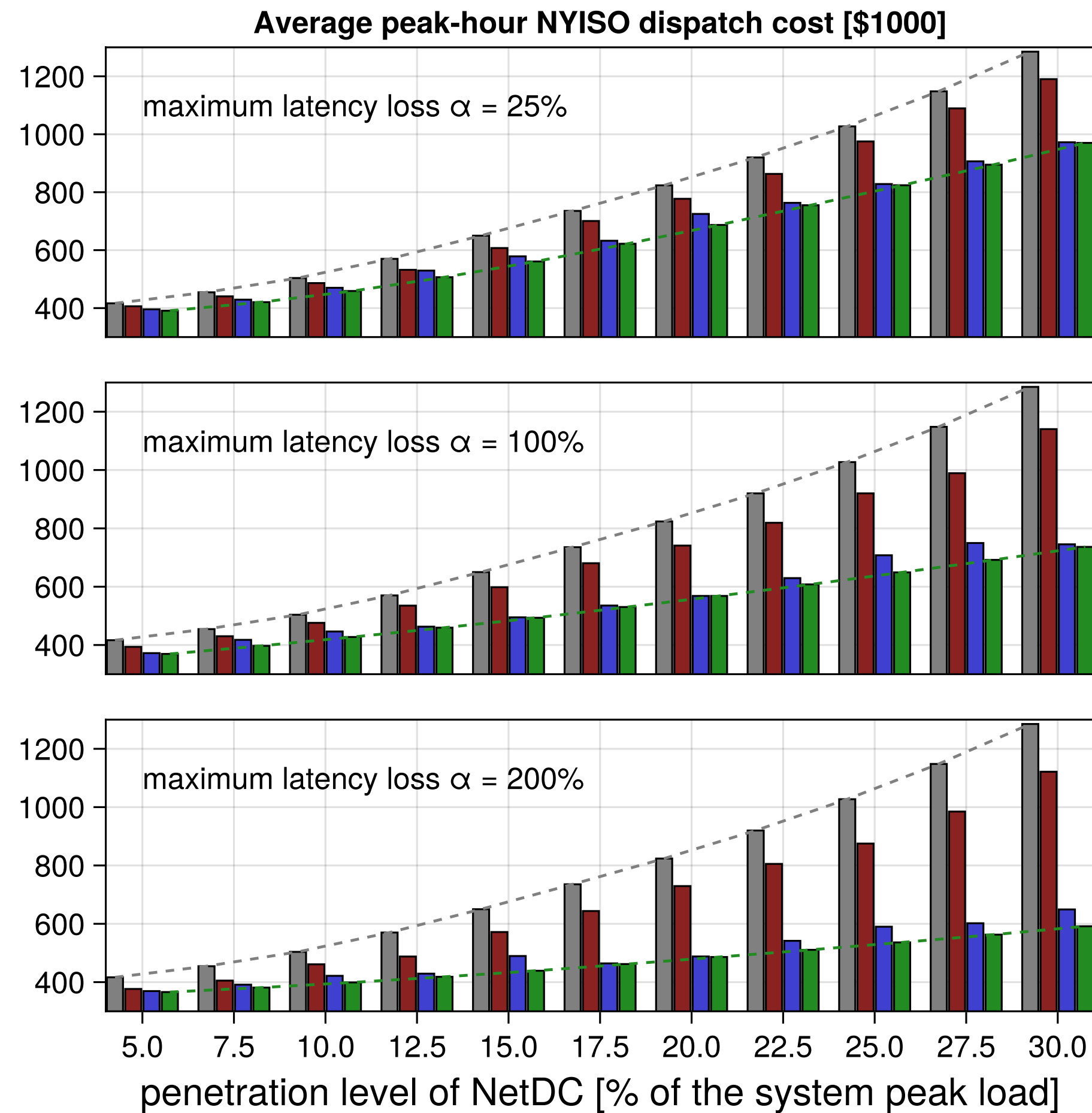
- ▶ Zonal real-time electricity demand ( $d$ );
- ▶ Zonal electricity prices ( $\lambda$ );
- ▶ Zonal renewable power generation ( $r$ );
- ▶ Power flows between aggregation zones ( $f$ ).

## Coordination policy to be optimized offline:

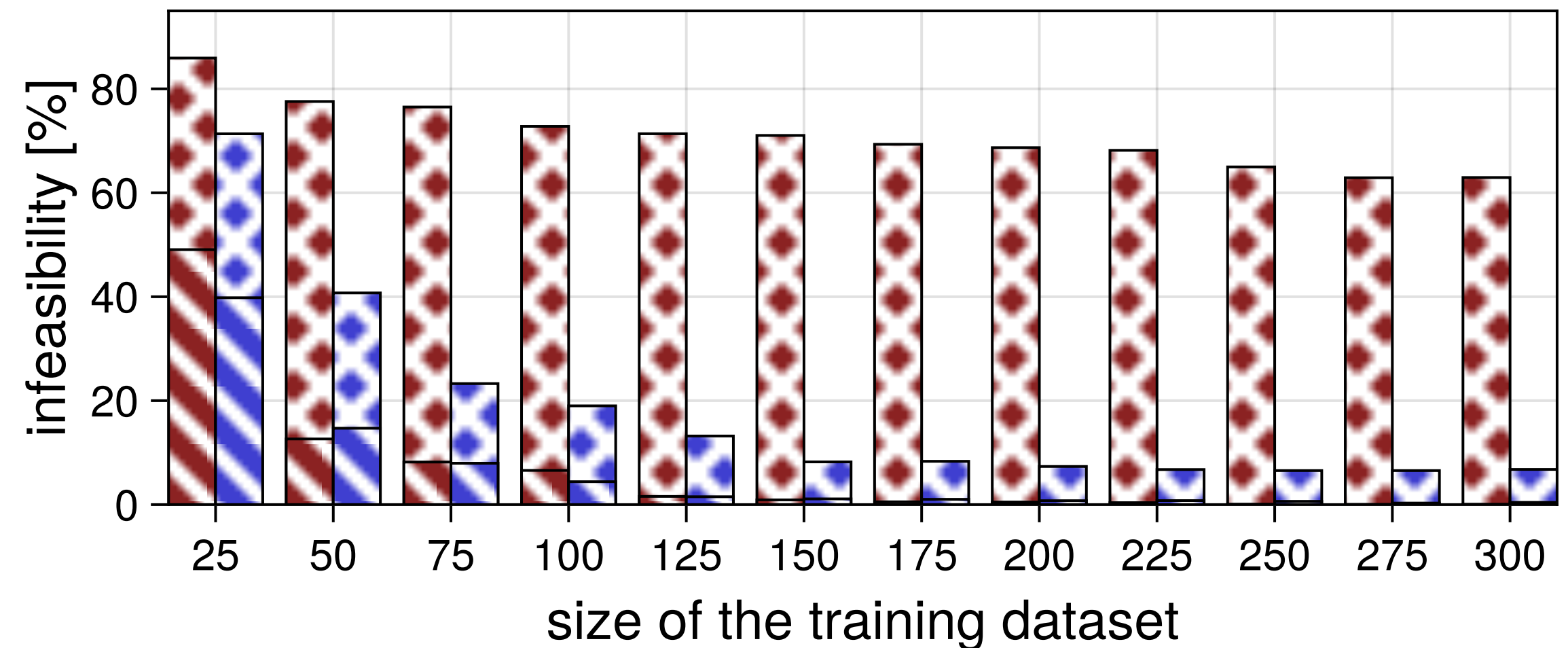
$$\phi \triangleq \beta_0 + \beta_1^d \begin{bmatrix} d_1 \\ \vdots \\ d_{11} \end{bmatrix} + \beta_1^\lambda \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_{11} \end{bmatrix} + \beta_1^r \begin{bmatrix} r_1 \\ \vdots \\ r_{11} \end{bmatrix} + \beta_1^f \begin{bmatrix} f_1 \\ \vdots \\ f_{12} \end{bmatrix}$$

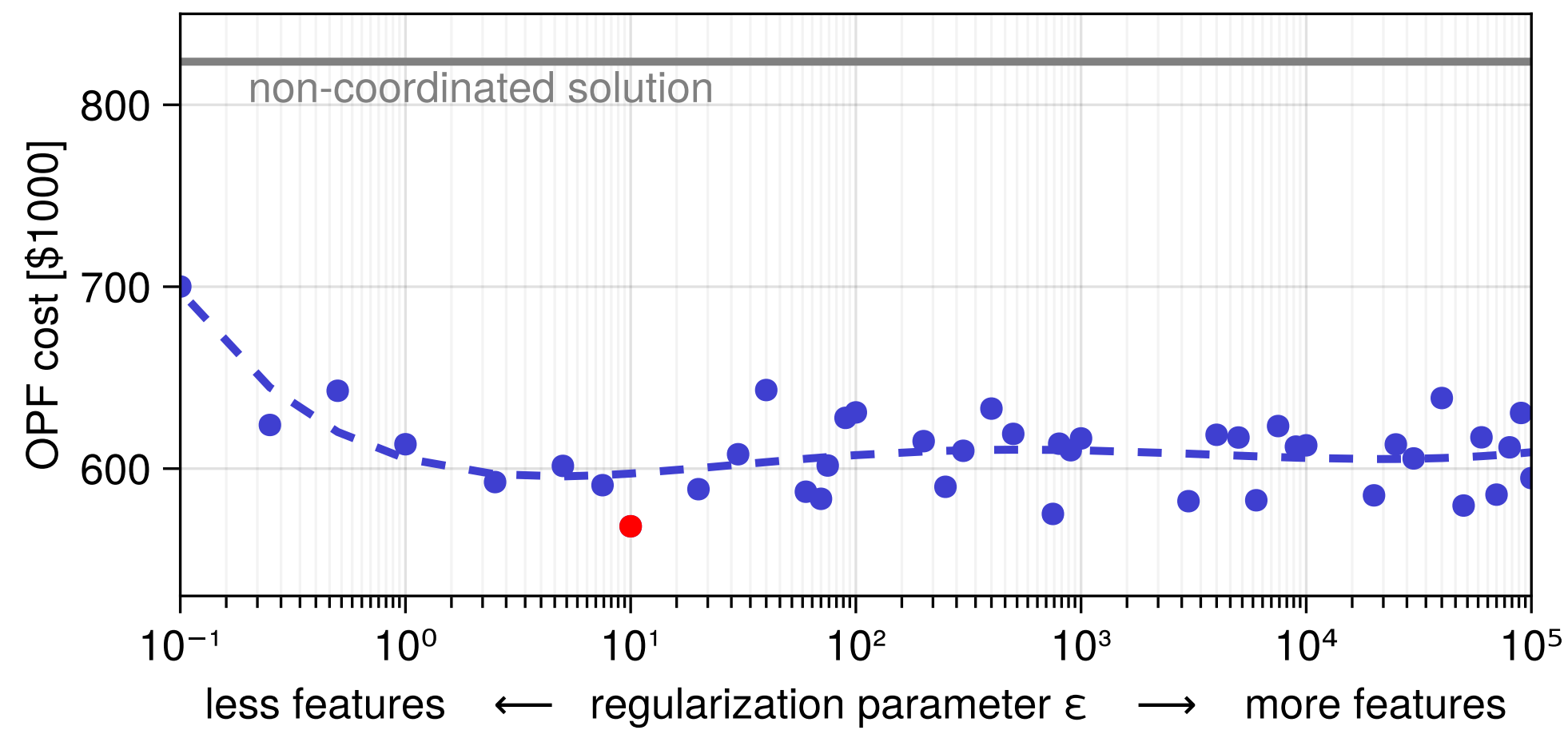


- ▶ Non-coordinated solution  $\Rightarrow$  quadratic cost growth
- ▶ **Ideal coordination**  $\Rightarrow$  more linear cost growth
- ▶ **Baseline regression** is often infeasible  $\Rightarrow$  small savings
- ▶ Feasibility guarantees of the **AgentCONCUR** protocol  $\Rightarrow$  efficient approximation of the ideal coordination



- ▶ Non-coordinated solution  $\Rightarrow$  quadratic cost growth
- ▶ **Ideal coordination**  $\Rightarrow$  more linear cost growth
- ▶ **Baseline regression** is often infeasible  $\Rightarrow$  small savings
- ▶ Feasibility guarantees of the **AgentCONCUR** protocol  $\Rightarrow$  efficient approximation of the ideal coordination





- ▶ Feature selection by means of  $\ell_1$ -regularization
- ▶  $\ell_1$ -regularization also ensures coordination robustness
- ▶ Can we organize coordination using just one feature?







## Policy optimization still requires sensitive data

$$\begin{aligned}
 & \underset{\beta, \varphi_1, \dots, \varphi_q, p_1, \dots, p_q}{\text{minimize}} && \frac{1}{q} \sum_{i=1}^q c_{\text{pwr}}(p_i) && \leftarrow \text{cost function?} \\
 & \text{subject to} && p_i \in \mathcal{P}_{\text{pwr}}(\vartheta_i^*), \quad \forall i = 1, \dots, q && \leftarrow \text{network/gen/load data?} \\
 & && \varphi_i = \beta_0 + \beta_1 x_i, \quad \|\beta\|_1 \leq \varepsilon, \quad \forall i = 1, \dots, q \\
 & && \vartheta_i^* \in \underset{\vartheta_i}{\text{minimize}} \quad c_{\text{net-dc}}(\vartheta_i) \\
 & && \text{subject to } \vartheta_i \in \mathcal{W}_{\text{net-dc}}(\varphi_i), \quad \forall i = 1, \dots, q
 \end{aligned}$$

- How to enable the exchange of sensitive training datasets?

## Policy optimization still requires sensitive data

$$\begin{aligned}
 & \underset{\beta, \varphi_1, \dots, \varphi_q, p_1, \dots, p_q}{\text{minimize}} && \frac{1}{q} \sum_{i=1}^q c_{\text{pwr}}(p_i) && \leftarrow \text{cost function?} \\
 & \text{subject to} && p_i \in \mathcal{P}_{\text{pwr}}(\vartheta_i^*), \quad \forall i = 1, \dots, q && \leftarrow \text{network/gen/load data?} \\
 & && \varphi_i = \beta_0 + \beta_1 x_i, \quad \|\beta\|_1 \leq \varepsilon, \quad \forall i = 1, \dots, q \\
 & && \vartheta_i^* \in \underset{\vartheta_i}{\text{minimize}} \quad c_{\text{net-dc}}(\vartheta_i) \\
 & && \text{subject to } \vartheta_i \in \mathcal{W}_{\text{net-dc}}(\varphi_i), \quad \forall i = 1, \dots, q
 \end{aligned}$$

- ▶ How to enable the exchange of sensitive training datasets?
- ▶ Differential privacy provides an answer

## Differentially Private Algorithms for Synthetic Power System Datasets

Vladimir Dvorkin, Jr., *Member, IEEE*, and Audun Botterud, *Member, IEEE*

- ▶ We build significantly less generation capacity than what we need to accommodate the growing AI demand
- ▶ We need to leverage the unique (geospatial) flexibility of datacenters to accommodate the loads
- ▶ Legacy optimization-based solutions to coordination are not the option
  - ▶ Lack of real-time communication interfaces
  - ▶ Privacy barriers for information exchange
  - ▶ Computationally intractable problem
- ▶ We developed a contextual regression mechanism (AgentCONCUR) to:
  - ▶ Minimize real-time communication requirements (contextual features)
  - ▶ Enable coordination at minimum data exchange (e.g., feature selection)
  - ▶ Computationally tractable real-time computations (millisecond)

IEEE TRANSACTIONS ON POWER SYSTEMS, AUGUST 2024

1

## Agent Coordination via Contextual Regression (AgentCONCUR) for Data Center Flexibility

Vladimir Dvorkin, *Member, IEEE*