

Carbon-Aware Computing for Data Centers with Probabilistic Performance Guarantees

Sophie Hall, Francesco Micheli, Giuseppe Belgioioso, Ana Radovanović, Florian Dörfler

arXiv:2410.21510v2 [eess.SY] 30 Oct 2024

Abstract—Data centers are significant contributors to carbon emissions and can strain power systems due to their high electricity consumption. To mitigate this impact and to participate in demand response programs, cloud computing companies strive to balance and optimize operations across their global fleets by making strategic decisions about when and where to place compute jobs for execution. In this paper, we introduce a load shaping scheme which reacts to time-varying grid signals by leveraging both temporal and spatial flexibility of compute jobs to provide risk-aware management guidelines and job placement with provable performance guarantees based on distributionally robust optimization. Our approach divides the problem into two key components: (i) day-ahead planning, which generates an optimal scheduling strategy based on historical load data, and (ii) real-time job placement and (time) scheduling, which dynamically tracks the optimal strategy generated in (i). We validate our method in simulation using normalized load profiles from randomly selected Google clusters, incorporating time-varying grid signals. We can demonstrate significant reductions in carbon cost and peak power with our approach compared to myopic greedy policies, while maintaining computational efficiency and abiding to system constraints.

Index Terms—Data-driven distributionally robust optimization, job scheduling, data center optimization, demand shift

I. INTRODUCTION

The number of hyperscale data centers (DCs) doubled between 2015 and 2021 [1]. In 2022, energy demand from DCs and data transmission networks accounted for about 1.5% of the global electricity demand (240–340 TWh) and for 1% of the energy-related greenhouse gas emissions [2]. Despite the rapid growth in the number of DCs, their global power demand and related greenhouse gas emissions have remained nearly constant between 2015 and 2020, largely due to significant improvements in operational efficiency of individual data centers [3]. However, efficiency gains are stagnating, and with the rising demand of AI technologies, ever more capacity is needed. As a result, global electricity demand from data centers is expected to triple between 2020 and 2030 [3], placing substantial strain on local grid infrastructure [4]. Nevertheless, well-coordinated networks of DCs can also provide flexibility to the power system by participating in demand response (DR) programs and offering ancillary services [5].

S. Hall, F. Micheli and F. Dörfler are with the Automatic Control Laboratory, ETH Zürich, 8092 Zürich, Switzerland (e-mail: {shall, frmicheli, doerfler}@ethz.ch), G. Belgioioso is with KTH Royal Institute of Technology, 11428 Stockholm, Sweden (e-mail: giubel@kth.se), and A. Radovanović with Google Inc., Mountain View, CA 94043 USA (e-mail: {anaradovanovic}@google.com). This work was supported by the NCCR Automation, a National Centre of Competence in Research, funded by the Swiss National Science Foundation (grant number 51NF40_225155), and by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation.

Cloud computing companies usually own large data center fleets spread across the globe. For example, in 2021, Microsoft, Amazon, and Google collectively owned more than 50% of all hyperscale DCs [6]. As efficiency gains at individual DCs begin to plateau and flexibility provision to the grid is rewarded, these companies are looking for ways to jointly optimize operations across their entire fleet to reduce operational costs and carbon footprint to reach their sustainability goals [3], [7]–[10], such as Google’s goal of net-zero emissions by 2030 [11]. Moreover, optimizing their global fleet will allow cloud computing companies to participate in DR schemes, unlocking new revenue streams and supporting the power system as a result. Several studies in Europe [12], [13] and the U.S. [14] have explored the potential for grid stabilization through their participation in DR programs. Google has recently performed DR pilot studies across Asia, Europe and the U.S. [15].

It is well-known that certain types of compute jobs, such as offline data processing, model training and simulation pipelines, are temporally flexible, while they consume significant computing resources. At the same time, large global balancing systems serving user requests have the potential to become carbon-aware and grid-aware through smart coordination. A key enabler for cloud computing companies to optimize operations and respond to time-varying grid signals is to exploit the spatial and temporal flexibility of different classes of compute jobs - meaning that they can be executed in different DC locations and/or delayed [16]–[18].

Whether reacting to DR signals or internally optimizing a global DC fleet, its time-varying operation is a challenging problem for multiple reasons: (i) it requires thousands of placement decisions every second, with compute jobs arriving continuously throughout the day; (ii) data on individual compute jobs is not accessible before submission; and (iii) job placement within clusters (with each data center containing multiple clusters) is managed by low-level operating systems, which are difficult to model and modify by high-level planners.

The literature covers various aspects of the problem but has key limitations. Specifically, most works (i) do not exploit temporal and spatial flexibility jointly [19]–[22]; (ii) rely on either optimization-based planning approaches, which are not tractable in real-time operation [21], [23], or myopic greedy policies, which are suboptimal as they do not take into account future grid signals or compute load predictions [16], [24], [25]; and (iii) either assume exact knowledge of incoming compute jobs, which is unrealistic, or lack specifics on real-time job placement [18], [21].

The authors in [21] propose geographic load shifting by

solving a bilevel program based on locational marginal carbon intensities. While promising, marginal carbon intensity data is not readily available in practice and it is unclear how the scheme would be applied in real-time operation. In [18], a bilevel program combines power system planning with a latency-minimizing task shifting problem, approximating the solution through linear regression. However, for a practical implementation, considering the stochasticity of compute load and giving robust service guarantees, would be essential.

In this work, we build on the model recently presented by Google [26], which leverages temporal flexibility of compute jobs to minimize the expected carbon footprint and daily peak power over the entire fleet of DCs. In this model, temporal shifting of compute jobs is indirectly regulated by introducing *Virtual Capacity Curves* (VCCs), which are hourly limits that artificially cap the computing resources available to the real-time job placement algorithm in each cluster.

If cloud computing companies want to participate in DR schemes, they will have to redistribute compute jobs while maintaining a guaranteed high level of service. However, as compute load of incoming jobs is highly stochastic, any redistribution and scheduling scheme must ensure robustness with respect to job queuing latency. Distributionally robust optimization (DRO) is an optimization technique which inherently offers robustness against unexpected compute load profiles, rare events, and distribution shifts that might be caused by changes in the job submission pattern.

We propose a distributionally-robust day-ahead planning problem coupled with a real-time placement algorithm which is driven by data, offers robustness guarantees, and considers various degrees of flexibility of jobs. Specifically, our approach is unique in the following ways:

- 1) Temporal and spatial flexibility: Our approach accommodates jobs with any flexibility, from 0 to 24h delay tolerance and local to global execution range, whereas in [26] only the temporal flexibility was considered.
- 2) Fully data-driven: We directly approximate the probability distribution of compute loads from data within the planning problem.
- 3) Tuning robustness: The probabilistic performance guarantees are tunable within the distributionally-robust optimization problem allowing to trade-off reliability and performance. All types of jobs are treated equally and performance guarantees hold across compute job classes.
- 4) Co-design of VCCs and the scheduling strategy: Deciding when and where jobs are being sent while limiting cluster capacity allows to fully exploit the spatial and temporal correlations in the data and perform preemptive peak planning such that spatially flexible loads go to underused clusters, thus reducing job queuing latency.

The proposed approach provides rigorous guarantees based on DRO theory while abiding to the key constraints discussed in [26] and from our discussions with Google. We validate our approach using a sample of normalized load profiles from randomly selected Google clusters and show that it outperforms commonly used greedy policies in terms of peak power and carbon cost while providing theoretical guarantees.

The rest of the paper is structured as follows. In Section II, we introduce the model and formulate the stochastic optimization problem. In Section III, we derive a distributionally-robust load schedule, approximating the ambiguity set from historical data. Section IV discusses real-time job placement. Section V presents simulation results and Section VI concludes the paper.

A. Basic Notation

We denote the set of the first K positive integers as $\mathbb{Z}_K := \{1, \dots, K\}$, and the set of positive integers from t to T , with $t < T$, as $\mathbb{Z}_{[t:T]} := \{t, \dots, T\}$. The operator $[\cdot]_+ := \max\{\cdot, 0\}$ is the projection on the positive orthant. Given some variables $v_{t,d}$, with $t \in \mathcal{T}$ and $d \in \mathcal{D}$, we denote the stacked vector as $v := \text{col}(v_{t,d})_{t \in \mathcal{T}, d \in \mathcal{D}} = [v_{1,1}, \dots, v_{T,D}]^\top$. Table I summarizes all parameters, sets, and variables used.

Indices, Sets and Parameters	
$k \in \mathcal{H} := \{1, \dots, K\}$ [hrs]	Submission time
$t \in \mathcal{T} := \{1, \dots, T\}$ [hrs]	Execution time
$d \in \mathcal{D} := \{1, \dots, D\}$	Data center cluster index
$c \in \mathcal{C} := \{1, \dots, C\}$	Job class
$\mathcal{D}_c \subseteq \mathcal{D}$	Set of clusters at which jobs of class c can be allocated
h_c [hrs]	Time flexibility of class c
$Y \in \mathbb{R}_{\geq 0}^{K \times C \times T \times D}$, $Y_{k,c,t,d} \in \mathbb{R}_{\geq 0}$,	Scheduling strategy as tensor and individual entry
$y \in \mathbb{R}_{\geq 0}^{KCTD}$,	Scheduling strategy as vector
$v_{t,d}$	Virtual Capacity Curve
$\bar{v}_{t,d}$	Real (cluster) machine capacity
$\rho_{t,d}^{\text{carb}}$	Carbon cost metric
ρ_d^{in}	Infrastructure cost metric

Table I: Modeling notation

B. Preliminaries

In the following, we introduce relevant concepts from distributionally robust optimization (DRO) [27]. The conditional value at risk (CVaR) is a popular risk measure that provides the average of the tail end of the loss distribution. Intuitively, it represents the expected average of the worst-case outcomes above a certain threshold $\beta \in [0, 1]$. It is defined as follows.

Definition 1. For a random variable $\omega \in \Omega \subset \mathbb{R}^r$ with distribution \mathbb{P}_ω and a function $\phi : \mathbb{R}^r \mapsto \mathbb{R}$, the CVaR of level β is defined as

$$\text{CVaR}_{1-\beta}^{\omega \sim \mathbb{P}_\omega}(\phi(\omega)) := \inf_{q \in \mathbb{R}} [\beta^{-1} \mathbb{E}^{\omega \sim \mathbb{P}_m} [[\phi(\omega) + q]_+] - q]. \quad (1)$$

The Wasserstein metric quantifies the minimum cost required to transport one distribution \mathbb{Q}_1 into another \mathbb{Q}_2 .

Definition 2 (Wasserstein distance [28]). Consider distributions $\mathbb{Q}_1, \mathbb{Q}_2 \in \mathcal{M}(\mathcal{Y})$ where $\mathcal{M}(\mathcal{Y})$ is the set of all probability distributions \mathbb{Q} supported on $\mathcal{Y} \subseteq \mathbb{R}^{nT}$ such that $\mathbb{E}[|\mathcal{Y}|] < \infty$. The Wasserstein metric $d_W : \mathcal{M}(\mathcal{Y}) \times \mathcal{M}(\mathcal{Y}) \rightarrow \mathbb{R}_{\geq 0}$ between the distributions \mathbb{Q}_1 and \mathbb{Q}_2 is defined as

$$d_W(\mathbb{Q}_1, \mathbb{Q}_2) := \inf \left\{ \int_{\mathcal{Y}^2} \|\mathbf{y}_1 - \mathbf{y}_2\| \Pi(d\mathbf{y}_1, d\mathbf{y}_2) \right\}, \quad (2)$$

where Π is the joint distribution of \mathbf{y}_1 and \mathbf{y}_2 with marginals \mathbb{Q}_1 and \mathbb{Q}_2 , respectively.

This distance metric d_w is also referred to as the “type-1 Wasserstein distance”.

II. PROBLEM SETTING

We consider the problem of dynamically allocating an aggregate compute load, denoted by s , among clusters $\mathcal{D} := \{1, \dots, D\}$ geographically distributed across locations, or data centers. Each data center contains multiple clusters. The compute load comprises resource usage across individual compute jobs that get placed in clusters over a 24-hour (planning) horizon, at one-hour intervals k , with $k \in \mathcal{H} := \{1, \dots, K\}$, and leave the system upon job completion or cancellation. We consider compute jobs that tolerate delayed execution, e.g., data processing pipelines, log analysis, large-scale simulations, and nightly builds [29]. These jobs vary in terms of their temporal and spatial flexibility. To categorize them, we define a set $\mathcal{C} := \{1, \dots, C\}$ of flexibility classes, wherein each class $c \in \mathcal{C}$ is associated with:

- (i) a temporal flexibility $h_c \in \mathbb{Z}_{\geq 0}$, namely, the maximum delay (in hours) the compute jobs can tolerate,
- (ii) a spatial range, encoded via $\mathcal{D}_c \subseteq \mathcal{D}$ consisting of all clusters at which jobs of class c are allocable¹.

Crucially, the aggregate load profile of each class c is unknown at the beginning of the planning horizon. Hence, we treat it as a random vector s_c drawn from an underlying probability distribution, denoted by \mathbb{P} . Each element $s_{c,k}$, with $k \in \mathcal{H}$, within this random vector represents the aggregate load of compute jobs of class c submitted at time k for processing.

Remark 1. *Classifying the spatial and temporal flexibility of compute jobs is challenging and an active area of research. Choosing the number of classes involves a trade-off between per-class prediction accuracy, computational complexity, and performance gains with increasing class granularity.*

Our goal is to design a scheduling strategy that is practical for real-world deployment, and leverages the temporal and spatial flexibility of compute jobs to ensure the DC network operates reliably (i.e., respecting capacity limits) and efficiently (i.e., quickly allocating jobs while minimizing carbon footprint), despite the uncertainty of next day’s compute load.

We represent the scheduling strategy for these compute loads with a tensor $Y \in \mathbb{R}^{K \times C \times T \times D}$, where each entry $Y_{k,c,t,d} \geq 0$, describes the fraction of the aggregate load $s_{c,k}$ of class c , submitted at time k , that is allocated for processing at time t to cluster d . Notably, $t \in \mathcal{T} := \{1, \dots, T\}$ with $T = K + \max_{c \in \mathcal{C}} h_c$ represents the planning horizon length for job execution. This horizon length must be larger than the submission horizon because jobs of class c submitted at time K still have a time flexibility of h_c .

As the scheduling strategy assigns fractions of compute load to multiple clusters and times, it is essential to ensure that these

¹The spatial range is a function of the compute job recurrence, network consumption, data dependence and more.

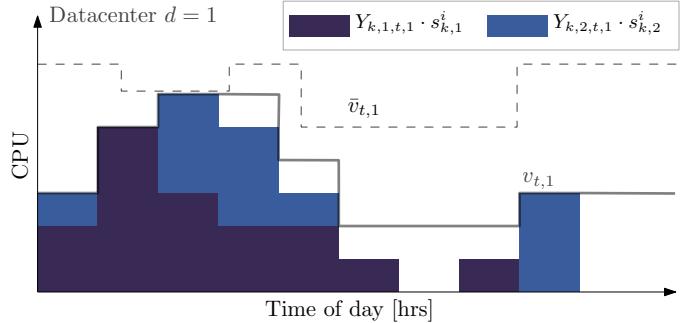


Figure 1: An example of the aggregate load schedule for compute loads from two flexibility classes $c \in \{1, 2\}$ at a single cluster $d = 1$ over 24 hours. The VCC (solid line) $v_{t,1}$ limits the allocable load at each time interval t , while the true capacity $\bar{v}_{t,1}$ (dashed line) is obtained by subtracting the inflexible load² from the (cluster) machine capacity.

fractions sum up to one. This guarantees that the entire load is accounted for, as captured by the following constraint:

$$\sum_{d \in \mathcal{D}} \sum_{t \in \mathcal{T}} Y_{k,c,t,d} \geq 1, \quad \forall k \in \mathcal{H}, c \in \mathcal{C}. \quad (3)$$

Since each class c has limited spatial range and temporal flexibility, the fraction of its load scheduled for clusters outside the spatial range, $d \in \mathcal{D} \setminus \mathcal{D}_c$, and all hours before the submission time k and after the delay tolerance $k + h_c$, must equal zero. This requirement is formalized as follows:

$$Y_{k,c,t,d} = 0, \quad \forall k \in \mathcal{H}, c \in \mathcal{C}, t \in \mathcal{T}, d \in \mathcal{D} \setminus \mathcal{D}_c, \quad (4a)$$

$$Y_{k,c,t,d} = 0, \quad \forall k \in \mathcal{H}, c \in \mathcal{C}, t \in \mathbb{Z}_{[k:k+h_c]}, d \in \mathcal{D}. \quad (4b)$$

The operational capacities $v_{t,d}$ of each cluster d , known as virtual capacity curve (VCC), determine the largest load allocable at a cluster. Inherently, each $v_{t,d}$ is limited by the true capacity $\bar{v}_{t,d}$, which takes into account the (cluster) machine capacity and the predicted inflexible load² yielding the following constraint

$$0 \leq v_{t,d} \leq \bar{v}_{t,d}, \quad \forall t \in \mathcal{T}, d \in \mathcal{D}. \quad (5)$$

In the current Google model [26], the VCCs were computed independently of the job scheduling strategy. In this paper, we treat $v = \text{col}(v_{t,d})_{d \in \mathcal{D}, t \in \mathcal{T}}$ as an additional decision variable and propose the co-design of VCCs and scheduling strategy Y . Intuitively, the allocated aggregate load per cluster d and time t must be lower than the VCC limit v . Since the compute loads $s_{k,c}$ are random variables, we enforce this constraint in a probabilistic sense, using the following Conditional Value-at-Risk (CVaR) constraint:

$$\text{CVaR}_{1-\beta}^{s \sim \mathbb{P}} \left[\max_{\substack{t \in \mathcal{T}, \\ d \in \mathcal{D}}} \left(\sum_{c \in \mathcal{C}} \sum_{k \in \mathcal{H}} Y_{k,c,t,d} \cdot s_{k,c} - v_{t,d} \right) \right] \leq 0. \quad (6)$$

Loosely speaking, this constraint ensures that, for each cluster d and time step t , a large percentile (tuned by β) of random

²Inflexible load consists of compute jobs that cannot be shifted in time and across-locations.

load profiles $s_{c,k}$ allocated using the strategy $Y_{k,c,t,d}$ satisfies the capacity constraint $v_{t,d}$. The CVaR operator in (6) is formally introduced in Definition 1.

The objective is a piece-wise linear function of compute usage and can incorporate different time-varying grid signals (carbon intensities, carbon free energy score, electricity prices, DR signals, etc.) as well as infrastructure related costs (peak power, stand-by machine cost etc.). Here, we use a similar objective function to Google's real implementation in [26], considering carbon footprint and daily peak power as a function of VCCs:

$$f(v) = \sum_{t \in \mathcal{T}} \sum_{d \in \mathcal{D}} \rho_{t,d}^{\text{carb}} v_{t,d} + \sum_{d \in \mathcal{D}} \rho_d^{\text{in}} \|\text{col}(v_{t,d})_{t \in \mathcal{T}}\|_\infty, \quad (7)$$

where $\rho_{t,d}^{\text{carb}}$ is a metric for the generated carbon footprint [30] and ρ_d^{in} is associated with infrastructure costs driven by a cluster's peak power consumption. The carbon impact of utilizing computing power is modeled using a linear relationship, as described in [16], [21]. As predictions for time-varying grid signals can be inaccurate or change rapidly, so re-optimizing throughout the day is necessary for better performance and to participate in DR schemes.

Overall, the problem of co-designing the optimal scheduling strategy and VCCs is a stochastic program of the form

$$\min_{y,v} f(v) \quad (8a)$$

$$\text{s.t. } \text{CVaR}_{1-\beta}^{s \sim \mathbb{P}}[F(y, v, s)] \leq 0, \quad (8b)$$

$$(3), (4), (5), \quad \forall k, c, t, d, \quad (8c)$$

where $y = \text{col}(y_r)_{r=0}^{KCTD}$, $y_r = Y_{k,c,t,d}$ with $r = K(c-1) + k + KC(t-1) + KCT(d-1)$ is a vectorized version of the load schedule, Y , which is more suitable as an optimization variable. The constraint (8b) is a reformulation of (6) in terms of y with F defined as

$$F(y, v, s) := \max_{t \in \mathcal{T}, d \in \mathcal{D}} y^\top A_{td} s + v^\top b_{td}, \quad (9)$$

where the matrix $A_{td} \in \mathbb{R}^{KCTD \times KC}$ has ones in its $(kctd, kc)$ entries and zeroes elsewhere, while the vector $b_{td} \in \mathbb{R}^{TD}$ has its td entries set to -1 and zero elsewhere. An example of a load schedule is shown in Figure 1.

Remark 2. The stochastic problem in (8a) could be formulated only in terms of scheduling strategy y and the random load profile s . The epigraph reformulation we chose through constraint (6) gives an explicit CVaR bound $v_{t,d}$ for each cluster d and time t for the random aggregate load profile $s_{k,c}$. As done by Google [26], the virtual capacities $v_{t,d}$ can be enforced as hard constraints on the executable load in each cluster. \square

Unfortunately, the stochastic problem (8) cannot be solved directly, since the underlying unknown probability distribution \mathbb{P} of the aggregate load s is unknown. Instead, we only have access to historical data, provided as 24-hour aggregate load samples for each flexibility class s_c^i . In the next section, we show how we leverage data-driven distributionally-robust optimization (DRO) [28] to solve (8), essentially building a robust forecast into the optimization problem using the finite set of historical data.

III. DISTRIBUTIONALLY-ROBUST SCHEDULING

Assuming that the per-class load profiles s_c^i in the dataset are independent samples drawn from the unknown distribution \mathbb{P} , we can derive an empirical compute load distribution as

$$\hat{\mathbb{P}} = \frac{1}{N} \sum_{i=1}^N \delta_{s^i}, \quad (10)$$

where δ_{s^i} denotes the Dirac distribution centred at the vector $s^i \in \mathbb{R}^{KC}$, namely, the i^{th} sample of per-class aggregate load obtained from the historical data. This empirical aggregate load distribution offers an approximate representation of the distribution of loads expected for the next day.

One could solve problem (8) with respect to the empirical distribution $\hat{\mathbb{P}}$ in place of the true distribution \mathbb{P} , obtaining the so called Sample Average Approximation (SAA) [31]. While simple to implement, the SAA relies on the fact that the distribution $\hat{\mathbb{P}}$ built from the historical data can be used as an accurate forecast of the distribution of the future aggregate load. When this is not the case, the SAA may lead to poor out-of-sample performance and cannot provide tight probabilistic guarantees if a limited amount of data is available or if there is a mismatch between the approximate and true aggregate load distribution.

In practice, the probability distribution of the aggregate load can exhibit significant variations across different clusters and through seasonal changes [32]. Moreover, the load on specific days, such as public holidays or close to project deadlines, can deviate substantially from historical patterns, making predictions particularly challenging [32]. Consequently, the empirical distribution (10) may not fully capture the nuances of future load distributions. This discrepancy leads to a potential shift, the infamous distribution shift, between the empirical distribution, which informs the derivation of the optimal load schedule, and the actual future load distribution encountered.

To address this uncertainty, we formulate a distributionally-robust version of (8) where the optimization is carried out against the worst-case distribution in a “neighboorhood” of the empirical distribution $\hat{\mathbb{P}}$, commonly known as *ambiguity set*. Formally, we define this ambiguity set as

$$\mathcal{B}^\varepsilon := \left\{ \mathbb{Q} \in \mathcal{P}_1(\mathcal{S}) \mid d_W(\hat{\mathbb{P}}, \mathbb{Q}) \leq \varepsilon \right\}, \quad (11)$$

where $\mathcal{P}_1(\mathcal{S})$ is the set of Borel probability measures with finite first moment, and $d_W(\hat{\mathbb{P}}, \mathbb{Q})$ is the so-called Wasserstein distance, given in Definition 2, between the probability distributions \mathbb{Q} and $\hat{\mathbb{P}}$ [33]. Intuitively, (11) describes the set of distributions that are within a radius ε , as measured by the Wasserstein metric d_W of the empirical distribution $\hat{\mathbb{P}}$. Ambiguity sets based on the Wasserstein distance are expressive and particularly well-suited for modeling and robustifying against so-called black-swan events, namely, rare and unpredictable outlier events with extreme impact³ [34].

The radius ε of the ambiguity set in (11) shall be chosen such that the ambiguity set is large enough to contain the true (unknown) future load distribution. It can be regarded as a

³For instance, an extremely large compute load that, under classical placement schemes, would lead to execution delays for time-sensitive jobs.

tuning knob that allows to trade-off between performance (cost reduction) and probabilistic constraint satisfaction. Namely, for $\varepsilon = 0$ we recover the SAA formulation that is not robust against distribution shifts whereas for large ε the solution to (12) will be robust but possibly conservative. In practice, a suitable radius ε can be obtained by analyzing the available data, e.g., by cross-validation.

With these definitions in place, we can formulate the distributionally-robust version of (8) as

$$\min_{y,v} f(v) \quad (12a)$$

$$\text{s.t. } \sup_{\mathbb{Q} \in \mathcal{B}^\varepsilon} \text{CVaR}_{1-\beta}^{s \sim \mathbb{Q}} [F(y, v, s)] \leq 0, \quad (12b)$$

$$(3), (4), (5), \quad \forall t, d, c, k. \quad (12c)$$

This approach enables us to achieve robust (probabilistic) constraint satisfaction while providing a probabilistic guarantee on the realized cost.

Remarkably, the worst-case constraint over the set of distribution in the ambiguity set that appears in problem (12) admits a tractable reformulation as a linear program (LP) that depends on the observed samples $s^i \in \mathcal{S}$ with \mathcal{S} being an a priori known support set for \mathbb{P} defined here as

$$\mathcal{S} := \{s \in \mathbb{R}^{KC} \mid Gs \leq h\}. \quad (13)$$

The matrix $G \in \mathbb{R}^{g \times KC}$ and vector $h \in \mathbb{R}^g$ encode prior information on the random aggregate load s , for example, the fact that job volumes can only be positive and that there exists a finite upper bound on the total aggregate load. A solution to the distributionally robust problem (12), i.e., the optimal schedule y^* and VCC v^* , can thus be obtained by solving an LP as shown in the following proposition.

Proposition 1. Assume that $s \in \mathcal{S}$ then, the DRO problem (12) can be reformulated as an LP⁴

$$\begin{aligned} & \min_{y,v,q} f(v) \\ & \text{s.t. } \lambda\varepsilon + \frac{1}{N} \sum_{i=1}^N p^i \leq q\beta \\ & \quad [v^\top b_{td} + q + (y^\top A_{td} - \eta_{itd}^\top G) s^i + \eta_{itd}^\top h]_+ \leq p^i \quad (14) \\ & \quad \|y^\top A_{td} - \eta_{itd}^\top G\|_\infty \leq \lambda \\ & \quad q \in \mathbb{R}, \eta_{itd} \geq 0, \lambda \geq 0 \\ & \quad (3), (4), (5), \\ & \quad \forall t \in \mathcal{T}, d \in \mathcal{D}, c \in \mathcal{C}, k \in \mathcal{H}, i \in \{1, \dots, N\} \end{aligned}$$

where $q \in \mathbb{R}$, $\lambda \in \mathbb{R}$, $p^i \in \mathbb{R}$, and $\eta_{itd} \in \mathbb{R}^g$ are auxiliary variables.

Proof: The proof follows directly from Proposition V.1. in [35] and is omitted here due to space limitations. ■

Next, we show that any solution to the DRO reformulation is a feasible point of problem (8).

Proposition 2. Assume that ε is such that $\mathbb{P} \in \mathcal{B}^\varepsilon$, and that (y^*, v^*) are component of a solution $(y^*, v^*, q^*, p^{i*}, \lambda^*, \nu^*)$ to

⁴The infinity norm term in the cost function (7) can be reformulated in LP form, but we omit it here to focus on the constraint reformulation.

the LP (14) in Proposition 1. Then, y^* and v^* are a feasible solution to the original stochastic optimization problem (8).

Proof: The proof follows directly from the definition of the worst-case distribution through the ambiguity set in (11). If (14) is feasible for all distributions in the set \mathcal{B}^ε , then it is also feasible for true unknown distribution \mathbb{P} as the ambiguity set contains \mathbb{P} by assumption. ■

This statement shows that if there exists a solution to the LP in Proposition 1, then the constraint satisfaction is guaranteed for all distributions inside the ambiguity set \mathcal{B}^ε . In other words, for any feasible scheduling strategy (thus, also for the optimal one), given a new aggregate load profile s^{new} drawn from a distribution inside the ambiguity set, the following holds:

- (i) With high probability⁵, the total compute load deployed on cluster d at time t will not exceed the planned operational capacity $v_{t,d}^*$;
- (ii) With high probability, the carbon footprint and peak load cost of the data center network over the 24-hour horizon is upper-bounded by $f(v^*)$.

IV. REAL-TIME JOB PLACEMENT

In the real application of job placement, the aggregate load of each class c consists of discrete compute jobs submitted continuously over the planning horizon. Thus, we can not perfectly implement the scheduling strategy Y^* computed in (14). Instead, in this section we present a real-time placement scheme which uses Y^* as a reference signal and tries to track it with every incoming job.

In fact, in real-time operation, the scheduler needs to make hundreds of thousands of placement decisions per second [36] making it imperative that the placement algorithm is computationally inexpensive, allowing jobs to flow immediately to the available compute resources. To achieve this, we designed our scheme to operate on two time-scales:

- 1) Day-ahead planning: DRO problem takes into account past aggregate load data, future load predictions, as well as the constraints of jobs and data centers. The output of the DRO problem is the optimal scheduling strategy Y^* which acts as a reference signal for job placement and the VCC curve for every cluster $v_{t,d}$, $\forall t \in \mathcal{T}, \forall d \in \mathcal{D}$.
- 2) Real-time job placement: Places every incoming job with the goal of tracking the optimal scheduling strategy Y^* .

The approach is illustrated in Figure 2. The time-scale separation of the overall control architecture allows to have a complex planning problem with rigorous theoretical guarantees while having a simple real-time placement which can handle a continuous flow of incoming compute jobs.

We point out that in the DRO problem used for day-ahead planning it is assumed that the continuous aggregate load profile can be split into arbitrary small fractions given by the scheduling strategy $Y_{k,c,t,d}$. However, in reality the aggregate load $s_{k,c}$ of class c at time k is the sum of compute volumes of discrete jobs $s_{k,c} = \sum_j u_c^{j,\text{vol}}$. Thus, in practice, the extent to which the optimal scheduling strategy Y^* can be approximated with discrete jobs depends on their size,

⁵As specified by the CVaR in (12b), tuned by adjusting parameter β .

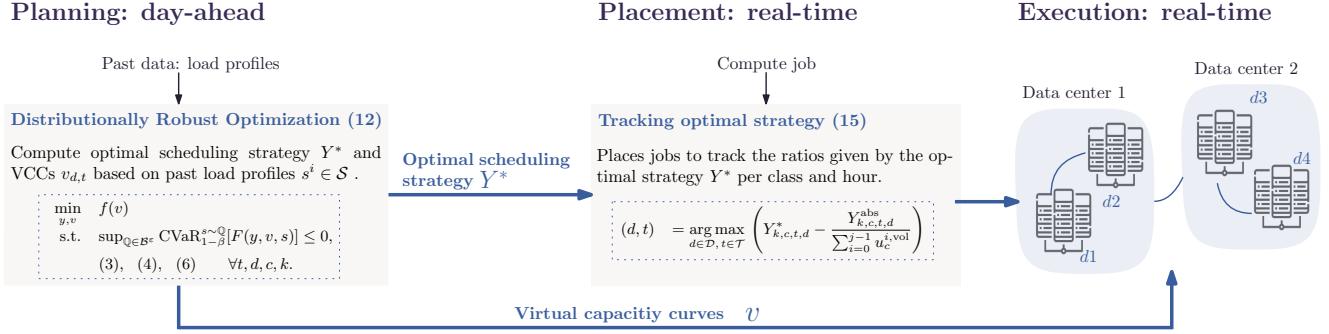


Figure 2: Schematic of the two layered control approach separated into day-ahead planning and real-time execution.

number and runtime. More details on application specific considerations and extensions are given in Subsection V-A.

In real-time operation, a job u of class c arrives, and we want to place it in cluster d and time slot t which is furthest away from fulfilling the optimal load fraction Y^* . Thus, for the j^{th} job of volume $u_c^{j,\text{vol}}$ arriving in hour k , we choose the placement tuple (d, t) for which the difference between Y^* and the current load ratios $\frac{Y_{k,c,t,d}^{\text{abs}}}{\sum_{i=0}^{j-1} u_c^{i,\text{vol}}}$ is maximal, expressed in the following:

$$(d, t) = \arg \max_{d \in \mathcal{D}, t \in \mathcal{T}} \left(\underbrace{Y_{k,c,t,d}^*}_{\text{Optimal fraction}} - \underbrace{\frac{Y_{k,c,t,d}^{\text{abs}}}{\sum_{i=0}^{j-1} u_c^{i,\text{vol}}}}_{\text{Current fraction}} \right) \quad (15)$$

where $Y_{k,c,t,d}^* \in \mathbb{R}$ is the optimal scheduling strategy for jobs of class c submitted at time k and $u_c^{i,\text{vol}}$ is the compute volume of job i . The tensor $Y_{k,c,t,d}^{\text{abs}} \in \mathbb{R}$ is the absolute compute volume already placed for execution in cluster d at time t . It is initialized at the beginning of every hour k before the first job arrives, i.e., when $j = 0$, $Y_{k,c,t,d}^{\text{abs}} = 0$. In practice, a job placed at a cluster d for execution at time t is stored in a cloud until the optimal time of execution, i.e., until $k = t$. A cluster-level operating system (known as Borg at Google [36], [37]) handles placements to available virtual machines.

V. IMPLEMENTATION AND NUMERICS

A. Application specific considerations and extensions

- 1) **Job runtimes:** We currently assume job runtimes of one time step. To handle longer runtimes, we can introduce additional job classes and constraints to ensure time continuity of the scheduling strategy.
- 2) **Cross-resource requirements:** Compute jobs require a variety of resources on machines (CPU, memory, disk, etc.). The DRO problem (12) can be formulated taking into account job's usage across all resource dimensions by considering vectorized VCC limits per cluster.
- 3) **Capacity constraints:** The VCCs $v_{t,d}$ sent to clusters can be used in two ways: (i) As information for the cluster-level operating system (known as Borg at Google [36], [37]) on the CVAR bound of the daily load profile; (ii) As a hard capacity constraint for the placement algorithm. The implementation by Google in [26] incorporates both (i) and (ii). We choose approach (i) for the simulations in

Section V.B-F as it aligns with the theoretical derivations in Section II-III. For Section V.G we use approach (ii) which is closer to the practical implementation.

- 4) **Data handling:** Samples s^i used to define (10) and to solve (14) can be (i) past per-class aggregate load profiles, which is simple but ignores correlations like day-of-week effects, or (ii) samples can be obtained from a calibrated stochastic predictor using large datasets, capturing side information and correlations.
- 5) **Computational complexity:** The resulting DRO problem (14) is an LP, thus large amounts of data and different classes can be included without jeopardizing solvability within one hour. Nevertheless, a data pre-processing pipeline would be essential to handle the massive amounts of data, their anomaly detection and missing data.
- 6) **Receding-horizon implementation and feedback:** As load predictions and grid signals may update throughout the day, and the state of the network may change, recomputing the DRO problem (12) every hour, including the system state, and applying it in a receding-horizon manner is expected to yield performance gains in practice.

B. Simulation set-up and processing of data

We consider a network of computing clusters $\mathcal{D} = \mathbb{Z}_4$ consisting of 2 data centers with two clusters each, as illustrated in the right-hand side of Figure 2. For exposition purposes, we consider a simplified setup in which the maximum delay time is $h_c = 10$ hrs for all compute jobs. However, we stress that our approach generally accommodates any number of time flexibility classes. This setup results in a total of seven job classes $\mathcal{C} = \mathbb{Z}_7$: four cluster-bound classes $\mathcal{D}_c = d, \forall c = d \in \{1, 2, 3, 4\}$, i.e., class one is bound to cluster one, two data center-bound classes $\mathcal{D}_5 = \{1, 2\}$, $\mathcal{D}_6 = \{3, 4\}$, and one globally flexible class, $\mathcal{D}_7 = \mathcal{D}$.

We use 60 normalized daily compute usage shapes from randomly selected clusters in the Google fleet. The LP in (14) is solved in Python using Gurobi [38] with 60 training load samples (80% of the total dataset). We use the compute load shapes from the training data set directly as individual samples s^i in the LP given in (14). The numerical case studies are conducted with 15 previously unobserved validation samples (20% of the dataset). The aggregate compute load curves after allocation are generated by multiplying the optimal scheduling

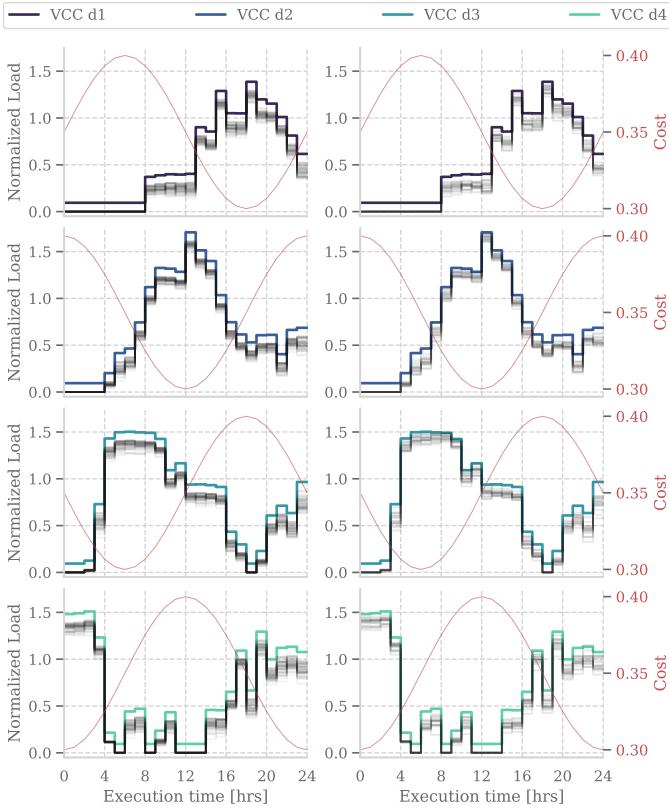


Figure 3: Comparing load profiles under the optimal schedule $Y_{k,c,t,d}^* \cdot s_{k,c}^i, \forall t \in \mathcal{T}, d \in \mathcal{D}$ for 60 training scenarios $s_{\text{train}}^i, i \in \mathbb{Z}_{60}$ (left) and 15 validation scenarios $s_{\text{val}}^i, i \in \mathbb{Z}_{15}$ (right).

strategy Y^* with the unobserved training samples. The cost function parameters in (7) are generated using sinusoids with a phase shift for different clusters to emulate variations in time-varying grid signals by location. All simulations are solved for CVaR level $\beta = 0.2$ and ambiguity set radius $\varepsilon = 8 \cdot 10^{-3}$ if not indicated otherwise.

C. Day-ahead: Training and validation scenarios

In Figure 3 we plot the output of the DRO day-ahead planning problem (12). Specifically, we show how the optimal load strategy $Y_{k,c,t,d}^*$ would distribute the aggregate load of the training and validation scenarios s^i across time $t \in \mathcal{T}$ and clusters $d \in \mathcal{D}$. We observe that even for unobserved aggregate load profiles in the validation set, load is shifted successfully and stays almost always within VCC bounds as given by the probabilistic constraints. We make the following observations:

- The robust scheduling strategy, $Y_{k,c,t,d}^*$, leverages the spatial range and temporal flexibility of loads to minimize the overall cost. For example, minimal load is executed during the first 8 hours of the day at location $d = 1$, where costs are considerably higher.
- For all the training scenarios, the optimal robust schedule results in loads slightly below the VCC capacity, $v_{t,d}$, at all times and across all clusters.
- For some of the validation scenarios, the optimal robust strategy results in load profiles that slightly exceed the

Policy	DRO	Greedy
Total cost increase: Mean [%] (STD [%])	$\varepsilon = 8 \cdot 10^{-3}: 2.57 (1.05)$	14.05 (0.63)
	$\varepsilon = 5 \cdot 10^{-2}: 8.86 (1.00)$	

Table II: Comparing the cost of the DRO and greedy policy to the optimal cost with perfect forecast.

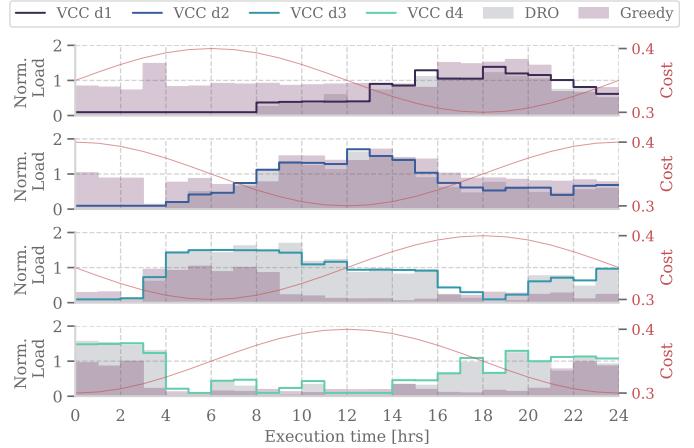


Figure 4: Comparing the running load of the DRO and greedy policy for discrete jobs over one day.

VCC limits, $v_{t,d}$, for example for hours 8 to 10 at cluster 3. This is expected due to the probabilistic nature of the CVaR constraint in (6).

D. Greedy vs planning

In this study, we compare the optimal scheduling strategy resulting from our proposed DRO approach with a basic greedy approach that places every incoming job at hour k to the “cheapest” cluster within the job classes’ spatial range $d \in \mathcal{D}_c$ and with available (cluster) machine capacity. The load profiles over one day for both policies are shown in Figure 4. Table II shows the cost across all 15 validation scenarios compared to an ideal placement that has a perfect forecast of the next day’s aggregate load $s_{k,c}$. Our analysis reveals several key findings:

- The DRO scheduling strategy shifts job execution across clusters and time to exploit low cost hours. For example, in the first 10 hours the compute load executed in $d \in \{1, 2\}$ is nearly zero as load is shifted to later hours.
- The costs of the DRO strategy are just 2.57 % higher than those of the policy with perfect forecast for $\varepsilon = 8 \cdot 10^{-3}$ and 8.86 % higher for $\varepsilon = 5 \cdot 10^{-2}$. For both robustness levels, it significantly reduces cost compared to the greedy policy, as shown in Table II.

This study shows that the DRO approach significantly reduces operational costs compared to a simplistic greedy policy while also giving guarantees on constraint satisfaction.

E. Exploiting spatial range and temporal flexibility

In this study, we aim to demonstrate that the DRO scheduling strategy utilizes the spatial range and temporal flexibility

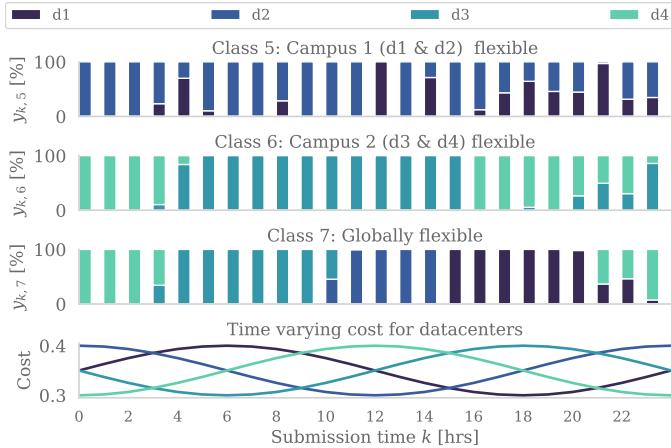


Figure 5: Plotting the optimal schedule Y^* for spatially flexible classes $c \in \mathbb{Z}_{[5:7]}$, showing the percentage of aggregate load submitted at hour k that should be sent to clusters $d \in \mathbb{Z}_4$.

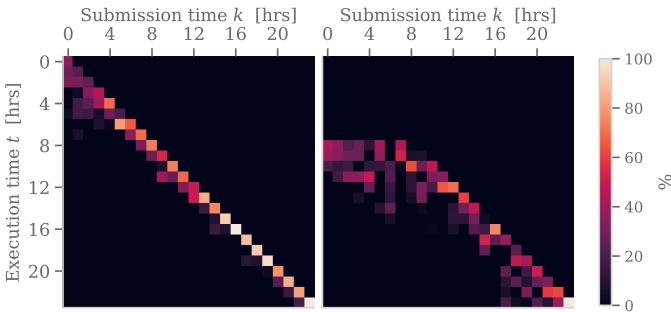


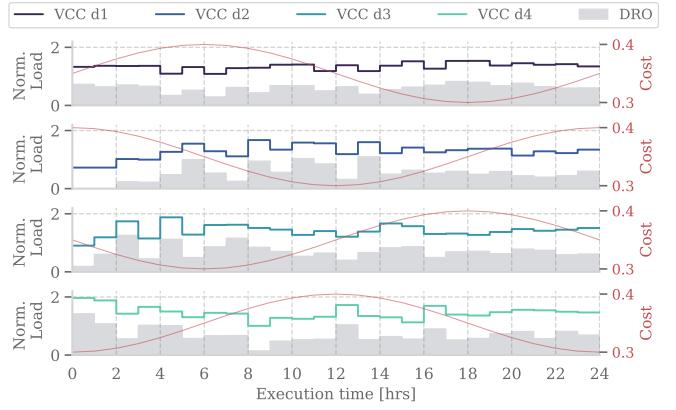
Figure 6: Two heatmaps of $Y_{k,c,t,d=1}^*$ show the aggregate load percentages submitted at k to be executed at t , with $\beta = 0.02$ (left) and $\beta = 0.2$ (right). With increased flexibility ($\beta = 0.2$) hours with high cost at $d = 1$ are avoided, e.g., hours 0 to 8.

of jobs to shift them to less expensive hours and clusters. In Figure 5, we present the optimal load fractions for all spatially flexible classes: DC-flexible classes $c = 4$ and $c = 5$, and the globally flexible class $c = 6$. The load distribution for these classes aligns with cluster-level price curves to avoid high-cost periods. For example, in the globally flexible class $c = 7$, all load shifts to $d = 1$ in the early hours to minimize costs.

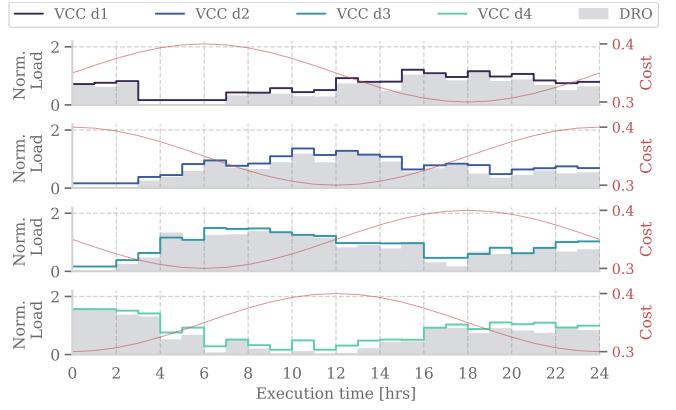
In Figure 6, we illustrate how the DRO scheme exploits temporal flexibility with the two time axes, submission and execution time. Without the second time axis (i.e., if $t = k$), the heatmap would show only diagonal entries, indicating 100% execution at submission time, thereby eliminating any temporal flexibility and enforcing that all loads are processed immediately. Figure 6 also highlights the effect of the tuning parameter β . For $\beta = 0.05$, much of the load is executed immediately upon submission to meet the delay constraint of $h_c = 10$ hours. In contrast, for $\beta = 0.7$, the probabilistic constraint in (6) is relaxed, allowing compute loads to be pushed to later hours to minimize costs.

F. Influence of CVaR constraint parameters β and ε

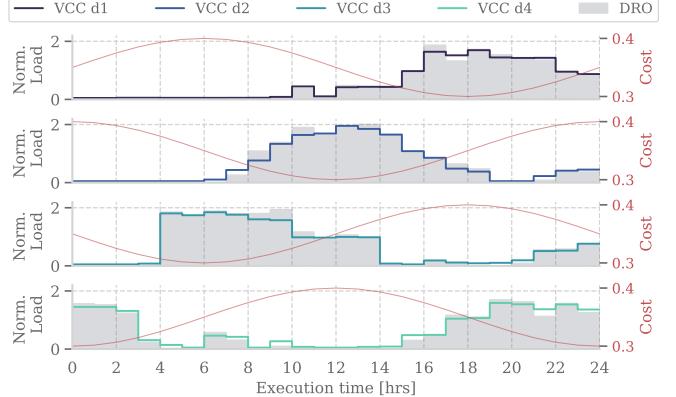
In this section, we highlight the role of the CVaR level β and DRO radius ε influencing the CVaR constraint in (6).



(a) With $\beta = 0.02$ and 0 VCC violations.



(b) With $\beta = 0.1$ and 4 VCC violations.



(c) With $\beta = 0.5$ and 33 VCC violations.

Figure 7: Comparison of VCCs and the realised load distribution for discrete jobs submitted over the day for three different β values.

Figure 7 shows the load profiles of discrete jobs over one day for increasing values of the CVaR level β . We observe that VCCs $v_{d,t}$ tighten with increasing β , corresponding to a relaxation of the probabilistic constraint against violating the VCC limits. While a value of $\beta = 0.5$ allows for significant shifting of compute load in time, thereby reducing operational costs, it may result in (cluster) level machine capacity being violated and some jobs not being executed due to insufficient capacity. However, note that in practice, the DRO would be implemented in a receding-horizon fashion, in which case it

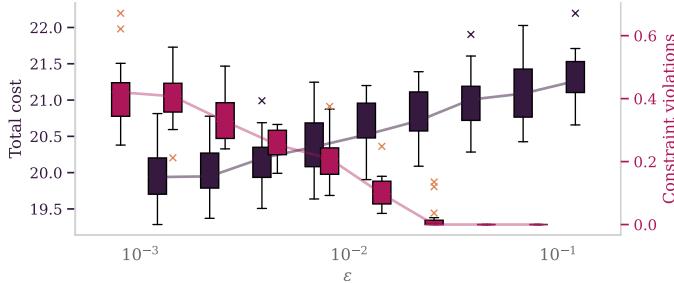


Figure 8: A boxplot showing the trade-off between cost and max VCC constraint violations (6) for increasing ε .

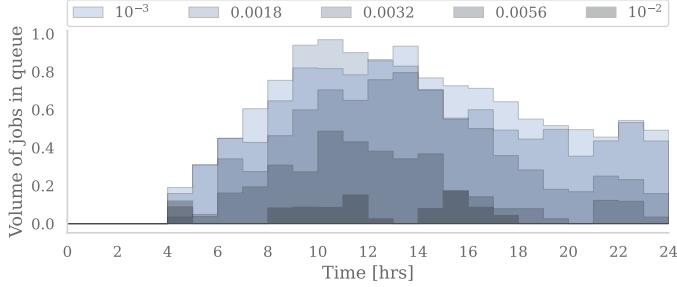


Figure 9: The queue length at cluster $d = 3$ for increasing ε .

is not necessary for all jobs to be executed by the end of the finite horizon.

Figure 8 shows how ε , the ambiguity set radius in (11), is used as a tuning knob to trade-off between performance and robustness against future realizations of the compute loads. The problem reduces to SAA with $\varepsilon = 0$. A choice of $\varepsilon = 10^{-2}$ achieves a good balance between performance and VCC constraint violations across all 15 validation scenarios.

G. Practical implementation

In the real application presented by Google in [26], VCC curves serve as a hard constraint for resource availability and the cluster-level operating system. When the load pushed to a specific cluster exceeds its planned capacity, jobs are placed in a cluster-level queue and executed as resources become available. In this numerical study, we tested how our proposed real-time job placement scheme in (15) handles discrete jobs. Jobs have a one-hour runtime and volume size drawn from a per-class normal distribution $\mathcal{N}_c(u_c^{\text{vol},\text{mean}}, 0.1)$, with $u_c^{\text{vol},\text{mean}}$ being the per-class compute load divided by job count.

In Figure 9 we compare the resulting queue length for different radii of the ambiguity set ε . We make three key observations: (i) The queue length is shortest during hours 0–10, as jobs are executed immediately during low-cost periods; (ii) Queue lengths are considerably longer for lower ε values, as the VCC limits become tighter and it is more likely that the scenario realization is outside the ambiguity set defined by the past data; (iii) A receding-horizon implementation that adjusts the planning policy by integrating the feedback from the queue lengths would be necessary in a real application.

VI. CONCLUSION

We presented a scheme for spatial and temporal management of flexible computing jobs in large-scale data center fleets which is driven by data, computationally efficient for real-time operation and retains theoretical performance guarantees provided by the DRO solution. The day-ahead planning problem incorporates historical compute load data and future predictions on load and prices to compute a cost-efficient robust aggregate load schedule which is tracked in real-time. Using normalized load profiles from randomly selected Google clusters, we demonstrated that the DRO scheduling strategy outperforms commonly used greedy policies while offering a direct way to trade off performance versus constraint violations through tuning parameters of the CVaR constraint. Future work will investigate a receding-horizon implementation of the DRO policy with real job-level data.

REFERENCES

- [1] E. Masanet, A. Shehabi, N. Lei, S. Smith, and J. Koomey, “Recalibrating global data center energy-use estimates,” *Science*, vol. 367, pp. 984–986, Feb. 2020.
- [2] International Energy Agency, “Data centres and data transmission networks.” [online], July 2023.
- [3] The Goldman Sachs Group, “Generational Growth AI, data centers and the coming US power demand surge.” [online], 4 2024.
- [4] C. Hodgson, “Booming ai demand threatens global electricity supply,” *Financial Times*, 2024. Accessed: 2024-10-18.
- [5] Enel X, “How data centers support the power grid with ancillary services,” 2024. Accessed: 2024-10-18.
- [6] Synergy Research Group, “Microsoft, Amazon and Google Account for Over Half of Today’s 600 Hyperscale Data Centers.” <https://tinyurl.com/3tz73kv7>, Jan. 2021. Accessed: 2024-07-05.
- [7] B. Johnson, “Carbon-aware kubernetes: Reducing emissions with smart scaling,” Oct. 2020. Microsoft Developer Blog.
- [8] R. Ramachandran, “Announcing the public preview of azure compute fleet,” Microsoft, May 2024. Accessed: 2024-10-10.
- [9] H. D. Dixit and J. Tse, “Retinas: Real-time infrastructure accounting for sustainability,” Meta Engineering Blog, 2024. Accessed: 2024-10-10.
- [10] Verrus, “How verrus is powering the data future.” Verrus News, 2024. Accessed: 2024-10-10.
- [11] Google, “Net zero carbon: Operating sustainably,” 2024. Accessed: 2024-10-18.
- [12] F. Bovera, M. Delfanti, and F. Bellifemine, “Economic opportunities for demand response by data centers within the new italian ancillary service market,” in *2018 IEEE International Telecommunications Energy Conference (INTELEC)*, vol. 10, pp. 1–8, IEEE, Oct. 2018.
- [13] J. Hansson, “The potential of data centre participation in ancillary service markets in Sweden,” Master’s thesis, KTH, School of Industrial Engineering and Management (ITM), 2022.
- [14] A. Wierman, Z. Liu, I. Liu, and H. Mohsenian-Rad, “Opportunities and challenges for data center demand response,” in *International Green Computing Conference*, vol. 14, pp. 1–10, IEEE, Nov. 2014.
- [15] V. Mehra and R. Hasegawa, “Using demand response to reduce data center power consumption.” <https://tinyurl.com/msj84hcy>, 2024. Accessed: 2024-10-18.
- [16] M. Xu and R. Buyya, “Managing renewable energy and carbon footprint in multi-cloud computing environments,” *Journal of Parallel and Distributed Computing*, vol. 135, pp. 191–202, 2020.
- [17] M. Abu Sharh, A. Shami, and A. Ouda, “Optimal and suboptimal resource allocation techniques in cloud computing data centers,” *Journal of Cloud Computing*, vol. 6, Mar. 2017.
- [18] V. Dvorkin, “Agent coordination via contextual regression (agentconcur) for data center flexibility,” *IEEE Trans. Power Syst.*, pp. 1–11, 2024.
- [19] T. Chen, A. G. Marques, and G. B. Giannakis, “Dglb: Distributed stochastic geographical load balancing over cloud networks,” *IEEE Trans. Parallel Distrib. Syst.*, vol. 28, no. 7, pp. 1866–1880, 2017.
- [20] Z. Liu, M. Lin, A. Wierman, S. Low, and L. L. H. Andrew, “Greening geographical load balancing,” *IEEE ACM Transactions on Networking*, vol. 23, no. 2, pp. 657–671, 2015.

- [21] J. Lindberg, B. C. Lesieurte, and L. A. Roald, "Using geographic load shifting to reduce carbon emissions," *Electric Power Systems Research*, vol. 212, p. 108586, 2022.
- [22] D. Paul and W.-D. Zhong, "Price and renewable aware geographical load balancing technique for data centres," in *2013 9th International Conference on Information, Communications and Signal Processing*, pp. 1–5, 2013.
- [23] E. Breukelman, S. Hall, G. Belgioioso, and F. D'orfler, "Carbon-aware computing in a network of data centers: A hierarchical game-theoretic approach," in *2024 European Control Conference (ECC)*, pp. 798–803, IEEE, 2024.
- [24] R. Wang, Y. Lu, K. Zhu, J. Hao, P. Wang, and Y. Cao, "An optimal task placement strategy in geo-distributed data centers involving renewable energy," *IEEE Access*, vol. 6, pp. 61948–61958, 2018.
- [25] A. Khosravi, L. L. H. Andrew, and R. Buyya, "Dynamic vm placement method for minimizing energy and carbon cost in geographically distributed cloud data centers," *IEEE Trans. Sustain. Comput.*, vol. 2, no. 2, pp. 183–196, 2017.
- [26] A. Radovanović, R. Koningstein, I. Schneider, B. Chen, A. Duarte, B. Roy, D. Xiao, M. Haridasan, P. Hung, N. Care, S. Talukdar, E. Mullen, K. Smith, M. Cottman, and W. Cirne, "Carbon-Aware Computing for Datacenters," *IEEE Trans. Power Syst.*, vol. 38, pp. 1270–1280, mar 2023.
- [27] D. Kuhn, P. M. Esfahani, V. A. Nguyen, and S. Shafeezahe-Abadeh, "Wasserstein distributionally robust optimization: Theory and applications in machine learning," in *Operations research & management science in the age of analytics*, pp. 130–166, Informs, 2019.
- [28] P. M. Esfahani and D. Kuhn, "Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations," *Mathematical Programming*, vol. 171, pp. 115–166, jul 2017.
- [29] J. Dean and S. Ghemawat, "MapReduce:simplified data processing on large clusters," in *OSDI'04: Sixth Symposium on Operating System Design and Implementation*, (San Francisco, CA), pp. 137–150, 2004.
- [30] Electricity Maps, "Carbon Intensity Data." [online], 2024.
- [31] A. J. K. A. S. T. Homem-de-Mello, "The Sample Average Approximation Method for Stochastic Discrete Optimization," *SIAM Journal on Optimization*, vol. 12, pp. 479–502, jan 2002.
- [32] J. Subirats and J. Guitart, "Assessing and forecasting energy efficiency on cloud computing platforms," *Future Generation Computer Systems*, vol. 45, pp. 70–94, 2015.
- [33] C. Villani *et al.*, *Optimal transport: old and new*, vol. 338. Springer, 2009.
- [34] N. N. Taleb, *The black swan : the impact of the highly improbable*. New York Times Bestseller, New York: Random House Trade Paperbacks, 2nd ed., random trade pbk. ed. ed., 2010.
- [35] A. R. Hota, A. Cherukuri, and J. Lygeros, "Data-driven chance constrained optimization under wasserstein ambiguity sets," in *2019 American Control Conference (ACC)*, pp. 1501–1506, IEEE, 2019.
- [36] M. Tirmazi, A. Barker, N. Deng, M. E. Haque, Z. G. Qin, S. Hand, M. Harchol-Balter, and J. Wilkes, "Borg: the next generation," in *Proceedings of the fifteenth European conference on computer systems*, pp. 1–14, 2020.
- [37] A. Verma, L. Pedrosa, M. R. Korupolu, D. Oppenheimer, E. Tune, and J. Wilkes, "Large-scale cluster management at Google with Borg," in *Proceedings of the European Conference on Computer Systems (EuroSys)*, (Bordeaux, France), 2015.
- [38] Gurobi Optimization, LLC, "Gurobi Optimizer Reference Manual," 2023.



Francesco Micheli is a PhD student at the Automatic Control Laboratory at ETH Zurich, under the supervision of Prof. J. Lygeros. He received his B.Sc. and M.Sc. degrees in Mechanical Engineering from Politecnico di Milano, Italy, in 2017 and 2018, respectively. His research focuses on safe learning and control, distributionally robust optimization, and robotics.



Giuseppe Belgioioso (Member, IEEE) is an Assistant Professor at the Division of Decision and Control Systems at KTH Royal Institute of Technology, Sweden. He received the Bachelor's degree in Information Engineering in 2012 and the Master's degree (cum laude) in Control Systems in 2015, both at the University of Padova, Italy. In 2020, he obtained the Ph.D. degree in Automatic Control at Eindhoven University of Technology (TU/e), The Netherlands. From 2021 to 2024, he was first a Postdoctoral researcher and then Senior Scientist at the Automatic Control Laboratory, ETH Zürich, Switzerland. His research lies at the intersection of optimization, game theory, and automatic control with applications in complex systems, such as power grids and traffic networks.



Ana Radovanović has been a research scientist at Google since early 2008, after she earned her PhD Degree in Electrical Engineering from Columbia University (2005) and worked for 3 years as a Research Staff Member in the Mathematical Sciences Department at IBM TJ Watson Research Center. For more than 10 years, Ana Radovanović has focused all her research efforts at Google on building innovative technologies and business models with two goals in mind: (i) to deliver more reliable, affordable and clean electricity to everyone in the world, and (ii) to help Google become a thought leader in decarbonizing the electricity grid. Nowadays, Ana is widely recognized as a technical lead and research entrepreneur. She is a Senior Staff Research Scientist, serving as a Technical Lead for Energy Analytics and Carbon Aware Computing at Google



Sophie Hall is a PhD student at the Automatic Control Laboratory at ETH Zurich. She received the Bachelor's degree in Mechanical Engineering from the University of Surrey, UK, and Nanyang Technological University, Singapore. She completed her Master's degree at ETH Zurich in Biomedical Engineering focusing on modelling and control. She was a finalist for the IFAC NMPC 2024 Young Authors Award. Her research interests revolve around game theory, model predictive control and real-time optimization with applications in network systems such as energy and supply chains.



Florian Dörlfer is a Full Professor at the Automatic Control Laboratory at ETH Zurich. He received his Ph.D. degree in Mechanical Engineering from the University of California at Santa Barbara in 2013, and a Diplom degree in Engineering Cybernetics from the University of Stuttgart in 2008. From 2013 to 2014 he was an Assistant Professor at the University of California Los Angeles. He has been serving as the Associate Head of the ETH Zurich Department of Information Technology and Electrical Engineering from 2021 until 2022. His research interests are centered around control, optimization, and system theory with applications in network systems, in particular electric power grids.