

Greening Geographical Load Balancing

Zhenhua Liu, Minghong Lin,
Adam Wierman, Steven H. Low
CMS, California Institute of Technology, Pasadena
{zliu2,mhlin,adamw,slow}@caltech.edu

Lachlan L. H. Andrew
Faculty of ICT
Swinburne University of Technology, Australia
landrew@swin.edu.au

ABSTRACT

Energy expenditure has become a significant fraction of data center operating costs. Recently, “geographical load balancing” has been suggested to reduce energy cost by exploiting the electricity price differences across regions. However, this reduction of cost can paradoxically increase total energy use.

This paper explores whether the geographical diversity of Internet-scale systems can additionally be used to provide environmental gains. Specifically, we explore whether geographical load balancing can encourage use of “green” renewable energy and reduce use of “brown” fossil fuel energy. We make two contributions. First, we derive two distributed algorithms for achieving optimal geographical load balancing. Second, we show that if electricity is dynamically priced in proportion to the instantaneous fraction of the total energy that is brown, then geographical load balancing provides significant reductions in brown energy use. However, the benefits depend strongly on the degree to which systems accept dynamic energy pricing and the form of pricing used.

Categories and Subject Descriptors

C.2.4 [Computer-Communication Networks]: Distributed Systems

General Terms

Algorithms, Performance

1. INTRODUCTION

Increasingly, web services are provided by massive, geographically diverse “Internet-scale” distributed systems, some having several data centers each with hundreds of thousands of servers. Such data centers require many megawatts of electricity and so companies like Google and Microsoft pay tens of millions of dollars annually for electricity [31].

The enormous, and growing energy demands of data centers have motivated research both in academia and industry on reducing energy usage, for both economic and environmental reasons. Engineering advances in cooling, virtualization, multi-core servers, DC power, etc. have led to significant improvements in the Power Usage Effectiveness (PUE) of data centers; see [6, 37, 19, 21]. Such work focuses on reducing the *energy use* of data centers and their components.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGMETRICS’11, June 7–11, 2011, San Jose, California, USA.

Copyright 2011 ACM 978-1-4503-0262-3/11/06 ...\$10.00.

A different stream of research has focused on exploiting the geographical diversity of Internet-scale systems to reduce the *energy cost*. Specifically, a system with clusters at tens or hundreds of locations around the world can dynamically route requests/jobs to clusters based on proximity to the user, load, and local electricity price. Thus, dynamic geographical load balancing can balance the revenue lost due to increased delay against the electricity costs at each location.

In recent years, many papers have illustrated the potential of geographical load balancing to provide significant cost savings for data centers, e.g., [24, 28, 31, 32, 34, 39] and the references therein. The goal of the current paper is different. Our goal is to explore the social impact of geographical load balancing systems. In particular, geographical load balancing aims to reduce energy costs, but this can come at the expense of increased total energy usage: by routing to a data center farther from the request source to use cheaper energy, the data center may need to complete the job faster, and so use more service capacity, and thus energy, than if the request was served closer to the source.

In contrast to this negative consequence, geographical load balancing also provides a huge opportunity for environmental benefit as the penetration of green, renewable energy sources increases. Specifically, an enormous challenge facing the electric grid is that of incorporating intermittent, unpredictable renewable sources such as wind and solar. Because generation supplied to the grid must be balanced by demand (i) instantaneously and (ii) locally (due to transmission losses), renewable sources pose a significant challenge. A key technique for handling the unpredictability of renewable sources is *demand-response*, which entails the grid adjusting the demand by changing the electricity price [2]. However, demand response entails a *local* customer curtailing use. In contrast, the demand of Internet-scale systems is flexible geographically; thus traffic can be routed to different regions to “follow the renewables”, providing demand-response without service interruption. Since data centers represent a significant and growing fraction of total electricity consumption, and the IT infrastructure is already in place, geographical load balancing has the potential to provide an extremely inexpensive approach for enabling large scale, global demand-response.

The key to realizing the environmental benefits above is for data centers to move from the fixed price contracts that are now typical toward some degree of dynamic pricing, with lower prices when green energy is available. The demand response markets currently in place provide a natural way for this transition to occur, and there is already evidence of some data centers participating in such markets [2].

The contribution of this paper is twofold. (1) We develop distributed algorithms for geographical load balancing with provable optimality guarantees. (2) We use the proposed algorithms to explore the feasibility and consequences of using geographical load balancing for demand response in the grid.

Contribution (1): To derive distributed geographical load balancing algorithms we use a simple but general model, described in detail in Section 2. In it, each data center minimizes its cost, which is a linear combination of an energy cost and the lost revenue due to the delay of requests (which includes both network propagation delay and load-dependent queueing delay within a data center). The geographical load balancing algorithm must then dynamically define both how requests should be routed to data centers and how to allocate capacity in each data center (i.e., how many servers are kept in active/energy-saving states).

In Section 3, we characterize the optimal geographical load balancing solutions and show that they have practically appealing properties, such as sparse routing tables. Then, in Section 4, we use the previous characterization to give two distributed algorithms which provably compute the optimal routing and provisioning decisions, and which require different types of coordination of computation. Finally, we evaluate the distributed algorithms in a trace-driven numeric simulation of a realistic, distributed, Internet-scale system (Section 5). The results show that a cost saving of over 40% during light-traffic periods is possible.

Contribution (2): In Section 6 we evaluate the feasibility and benefits of using geographical load balancing to facilitate the integration of renewable sources into the grid. We do this using a trace-driven numeric simulation of a realistic, distributed Internet-scale system in combination with models for the availability of wind and solar energy over time.

When the data center incentive is aligned with the social objective or reducing brown energy by dynamically pricing electricity proportionally to the fraction of the total energy coming from brown sources, we show that “follow the renewables” routing ensues (see Figure 5), causing significant social benefit. In contrast, we also determine the wasted brown energy when prices are static, or are dynamic but do not align data center and social objectives.

2. MODEL AND NOTATION

We now introduce the workload and data center models, followed by the geographical load balancing problem.

2.1 The workload model

We consider a discrete-time model whose timeslot matches the timescale at which routing decisions and capacity provisioning decisions can be updated. There is a (possibly long) interval of interest $t \in \{1, \dots, T\}$. There are $|J|$ geographically concentrated sources of requests, i.e., “cities”, and the mean arrival rate from source j at time t is $L_j(t)$. Job inter-arrival times are assumed to be much shorter than a timeslot, so that provisioning can be based on the average arrival rate during a slot. In practice, T could be a month and a slot length could be 1 hour. Our analytic results make no assumptions on $L_j(t)$; however to provide realistic estimates we use real-world traces to define $L_j(t)$ in Sections 5 and 6.

2.2 The data center cost model

We model an Internet-scale system as a collection of $|N|$ geographically diverse data centers, where data center i is modeled as a collection of M_i homogeneous servers. The model focuses on two key control decisions of geographical load balancing: (i) determining $\lambda_{ij}(t)$, the amount of traffic routed from source j to data center i ; and (ii) determining $m_i(t) \in \{0, \dots, M_i\}$, the number of active servers at data center i . The system seeks to choose $\lambda_{ij}(t)$ and $m_i(t)$ in order to minimize cost during $[1, T]$. Depending on the system design these decisions may be centralized or decentralized. Algorithms for these decisions are the focus of Section 4.

Our model for data center costs focuses on the server costs of the data center.¹ We model costs by combining the *energy cost* and the *delay cost* (in terms of lost revenue). Note that, to simplify the model, we do not include the switching costs associated with cycling servers in and out of power-saving modes; however the approach of [24] provides a natural way to incorporate such costs if desired.

Energy cost. To capture the geographic diversity and variation over time of energy costs, we let $g_i(t, m_i, \lambda_i)$ denote the energy cost for data center i during timeslot t given m_i active servers and arrival rate λ_i . For every fixed t , we assume that $g_i(t, m_i, \lambda_i)$ is continuously differentiable in both m_i and λ_i , strictly increasing in m_i , non-decreasing in λ_i , and convex in m_i . This formulation is quite general, and captures, for example, the common charging plan of a fixed price per kWh plus an additional “demand charge” for the peak of the average power used over a sliding 15 minute window [27]. Additionally, it can capture a wide range of models for server power consumption, e.g., energy costs as an affine function of the load, see [14], or as a polynomial function of the speed, see [40, 5].

Defining $\lambda_i(t) = \sum_{j \in J} \lambda_{ij}(t)$, the total energy cost of data center i during timeslot t , $\mathcal{E}_i(t)$, is simply

$$\mathcal{E}_i(t) = g_i(t, m_i(t), \lambda_i(t)). \quad (1)$$

Delay cost. The delay cost captures the lost revenue incurred because of the delay experienced by the requests. To model this, we define $r(d)$ as the lost revenue associated with a job experiencing delay d . We assume that $r(d)$ is strictly increasing and convex in d .

To model the delay, we consider its two components: the network delay experienced while the request is outside of the data center and the queueing delay experienced while the request is at the data center.

To model the *network delay*, we let $d_{ij}(t)$ denote the network delay experienced by a request from source j to data center i during timeslot t . We make no requirements on the structure of the $d_{ij}(t)$.

To model the *queueing delay*, we let $f_i(m_i, \lambda_i)$ denote the queueing delay at data center i given m_i active servers and an arrival rate of λ_i . We assume that f_i is strictly decreasing in m_i , strictly increasing in λ_i , and strictly convex in both m_i and λ_i . Further, for stability, we must have that $\lambda_i = 0$ or $\lambda_i < m_i \mu_i$, where μ_i is the service rate of a server at data center i . Thus, we define $f_i(m_i, \lambda_i) = \infty$ for $\lambda_i \geq m_i \mu_i$. Elsewhere, we assume f_i is finite, continuous and differentiable. Note that these assumptions are satisfied by most standard queueing formulae, e.g., the mean delay under M/GI/1 Processor Sharing (PS) queue and the 95th percentile of delay under the M/M/1. Further, the convexity of f_i in m_i models the law of diminishing returns for parallelism.

Combining the above gives the following model for the total delay cost $\mathcal{D}_i(t)$ at data center i during timeslot t :

$$\mathcal{D}_i(t) = \sum_{j \in J} \lambda_{ij}(t) r(f_i(m_i(t), \lambda_i(t)) + d_{ij}(t)). \quad (2)$$

2.3 The geographical load balancing problem

Given the cost models above, the goal of geographical load balancing is to choose the routing policy $\lambda_{ij}(t)$ and the number of active servers in each data center $m_i(t)$ at each time t in order to minimize the total cost during $[1, T]$. This is captured by the following optimization problem:

$$\min_{\mathbf{m}(t), \boldsymbol{\lambda}(t)} \sum_{t=1}^T \sum_{i \in N} (\mathcal{E}_i(t) + \mathcal{D}_i(t)) \quad (3a)$$

¹Minimizing server energy consumption also reduces cooling and power distribution costs.

$$\text{s.t. } \sum_{i \in N} \lambda_{ij}(t) = L_j(t), \quad \forall j \in J \quad (3b)$$

$$\lambda_{ij}(t) \geq 0, \quad \forall i \in N, \forall j \in J \quad (3c)$$

$$0 \leq m_i(t) \leq M_i, \quad \forall i \in N \quad (3d)$$

$$m_i(t) \in \mathbb{N}, \quad \forall i \in N \quad (3e)$$

To simplify (3), note that Internet data centers typically contain thousands of active servers. So, we can relax the integer constraint in (3) and round the resulting solution with minimal increase in cost. Also, because this model neglects the cost of turning servers on or off, the optimization decouples into independent sub-problems for each timeslot t . For the analysis we consider only a single interval and omit the explicit time dependence.² Thus (3) becomes

$$\min_{\mathbf{m}, \boldsymbol{\lambda}} \sum_{i \in N} g_i(m_i, \lambda_i) + \sum_{i \in N} \sum_{j \in J} \lambda_{ij} r(d_{ij} + f_i(m_i, \lambda_i)) \quad (4a)$$

$$\text{s.t. } \sum_{i \in N} \lambda_{ij} = L_j, \quad \forall j \in J \quad (4b)$$

$$\lambda_{ij} \geq 0, \quad \forall i \in N, \forall j \in J \quad (4c)$$

$$0 \leq m_i \leq M_i, \quad \forall i \in N, \quad (4d)$$

We refer to this formulation as GLB. Note that GLB is jointly convex in λ_{ij} and m_i and can be efficiently solved centrally. However, a distributed solution algorithm is usually required, such as those derived in Section 4.

In contrast to prior work studying geographical load balancing, it is important to observe that this paper is the first, to our knowledge, to incorporate jointly optimizing the total energy cost and the end-to-end user delay with consideration of both price diversity and network delay diversity.

GLB provides a general framework for studying geographical load balancing. However, the model ignores many aspects of data center design, e.g., reliability and availability, which are central to data center service level agreements. Such issues are beyond the scope of this paper; however our designs merge nicely with proposals such as [36] for these goals.

The GLB model is too broad for some of our analytic results and thus we often use two restricted versions.

Linear lost revenue. This model uses a lost revenue function $r(d) = \beta d$, for constant β . Though it is difficult to choose a “universal” form for the lost revenue associated with delay, there is evidence that it is linear within the range of interest for sites such as Google, Bing, and Shopzilla [13]. GLB then simplifies to

$$\min_{\mathbf{m}, \boldsymbol{\lambda}} \sum_{i \in N} g_i(m_i, \lambda_i) + \beta \left(\sum_{i \in N} \lambda_i f_i(m_i, \lambda_i) + \sum_{i \in N} \sum_{j \in J} d_{ij} \lambda_{ij} \right) \quad (5)$$

subject to (4b)–(4d). We call this optimization GLB-LIN.

Queueing-based delay. We occasionally specify the form of f and g using queueing models. This provides increased intuition about the distributed algorithms presented.

If the workload is perfectly parallelizable, and arrivals are Poisson, then $f_i(m_i, \lambda_i)$ is the average delay of m_i parallel queues, each with arrival rate λ_i/m_i . Moreover, if each queue is an M/GI/1 Processor Sharing (PS) queue, then $f_i(m_i, \lambda_i) = 1/(\mu_i - \lambda_i/m_i)$. We also assume $g_i(m_i, \lambda_i) = p_i m_i$, which implies that the increase in energy cost per timeslot for being in an active state, rather than a low-power state, is m_i regardless of λ_i .

Under these restrictions, the GLB formulation becomes:

$$\min_{\mathbf{m}, \boldsymbol{\lambda}} \sum_{i \in N} p_i m_i + \beta \sum_{j \in J} \sum_{i \in N} \lambda_{ij} \left(\frac{1}{\mu_i - \lambda_i/m_i} + d_{ij} \right) \quad (6a)$$

²Time-dependence of L_j and prices is re-introduced for, and central to, the numeric results in Sections 5 and 6.

subject to (4b)–(4d) and the additional constraint

$$\lambda_i \leq m_i \mu_i \quad \forall i \in N. \quad (6b)$$

We refer to this optimization as GLB-Q.

Additional Notation. Throughout the paper we use $|S|$ to denote the cardinality of a set S and bold symbols to denote vectors or tuples. In particular, $\boldsymbol{\lambda}_j = (\lambda_{ij})_{i \in N}$ denotes the tuple of λ_{ij} from source j , and $\boldsymbol{\lambda}_{-j} = (\lambda_{ik})_{i \in N, k \in J \setminus \{j\}}$ denotes the tuples of the remaining λ_{ik} , which forms a matrix. Similarly $\mathbf{m} = (m_i)_{i \in N}$ and $\boldsymbol{\lambda} = (\lambda_{ij})_{i \in N, j \in J}$.

We also need the following in discussing the algorithms. Define $F_i(m_i, \lambda_i) = g_i(m_i, \lambda_i) + \beta \lambda_i f_i(m_i, \lambda_i)$, and define $F(\mathbf{m}, \boldsymbol{\lambda}) = \sum_{i \in N} F_i(m_i, \lambda_i) + \sum_{i,j} \lambda_{ij} d_{ij}$. Further, let $\hat{m}_i(\lambda_i)$ be the unconstrained optimal m_i at data center i given fixed λ_i , i.e., the unique solution to $\partial F_i(m_i, \lambda_i)/\partial m_i = 0$.

2.4 Practical considerations

Our model assumes there exist mechanisms for dynamically (i) provisioning capacity of data centers, and (ii) adapting the routing of requests from sources to data centers.

With respect to (i), many dynamic server provisioning techniques are being explored by both academics and industry, e.g., [4, 11, 16, 38]. With respect to (ii), there are also a variety of protocol-level mechanisms employed for data center selection today. They include, (a) dynamically generated DNS responses, (b) HTTP redirection, and (c) using persistent HTTP proxies to tunnel requests. Each of these has been evaluated thoroughly, e.g., [12, 25, 30], and though DNS has drawbacks it remains the preferred mechanism for many industry leaders such as Akamai, possibly due to the added latency due to HTTP redirection and tunneling [29]. Within the GLB model, we have implicitly assumed that there exists a proxy/DNS server co-located with each source.

Our model also assumes that the network delays, d_{ij} can be estimated, which has been studied extensively, including work on reducing the overhead of such measurements, e.g., [35], and mapping and synthetic coordinate approaches, e.g., [22, 26]. We discuss the sensitivity of our algorithms to error in these estimates in Section 5.

3. CHARACTERIZING THE OPTIMA

We now provide characterizations of the optimal solutions to GLB, which are important for proving convergence of the distributed algorithms of Section 4. They are also necessary because, a priori, one might worry that the optimal solution requires a very complex routing structure, which would be impractical; or that the set of optimal solutions is very fragmented, which would slow convergence in practice. The results here show that such worries are unwarranted.

Uniqueness of optimal solution.

To begin, note that GLB has at least one optimal solution. This can be seen by applying Weierstrass theorem [7], since the objective function is continuous and the feasible set is compact subset of \mathbb{R}^n . Although the optimal solution is generally not unique, there are natural aggregate quantities unique over the set of optimal solutions, which is a convex set. These are the focus of this section.

A first result is that for the GLB-LIN formulation, under weak conditions on f_i and g_i , we have that λ_i is common across all optimal solutions. Thus, the input to the data center provisioning optimization is unique.

Theorem 1. *Consider the GLB-LIN formulation. Suppose that for all i , $F_i(m_i, \lambda_i)$ is jointly convex in λ_i and m_i , and continuously differentiable in λ_i . Further, suppose that $\hat{m}_i(\lambda_i)$ is strictly convex. Then, for each i , λ_i is common for all optimal solutions.*

The proof is in the Appendix. Theorem 1 implies that the server arrival rates at each data center, i.e., λ_i/m_i , are common among all optimal solutions.

Though the conditions on F_i and \hat{m}_i are weak, they do not hold for GLB-Q. In that case, $\hat{m}_i(\lambda_i)$ is linear, and thus not strictly convex. Although the λ_i are not common across all optimal solutions in this setting, the server arrival rates remain common across all optimal solutions.

Theorem 2. *For each data center i , the server arrival rates, λ_i/m_i , are common across all optimal solutions to GLB-Q.*

Sparsity of routing.

It would be impractical if the optimal solutions to GLB required that traffic from each source was divided up among (nearly) all of the data centers. In general, each λ_{ij} could be non-zero, yielding $|N| \times |J|$ flows of traffic from sources to data centers, which would lead to significant scaling issues. Luckily, there is guaranteed to exist an optimal solution with extremely sparse routing. Specifically:

Theorem 3. *There exists an optimal solution to GLB with at most $(|N| + |J| - 1)$ of the λ_{ij} strictly positive.*

Though Theorem 3 does not guarantee that every optimal solution is sparse, the proof is constructive. Thus, it provides an approach which allows one to transform an optimal solution into a sparse optimal solution.

The following result further highlights the sparsity of the routing: any source will route to at most one data center that is not fully active, i.e., where there exists at least a server in power-saving mode.

Theorem 4. *Consider GLB-Q where power costs p_i are drawn from an arbitrary continuous distribution. If any source $j \in J$ has its traffic split between multiple data centers $N' \subseteq N$ in an optimal solution, then, with probability 1, at most one data center $i \in N'$ has $m_i < M_i$.*

4. ALGORITHMS

We now focus on GLB-Q and present two distributed algorithms that solve it, and prove their convergence.

Since GLB-Q is convex, it can be efficiently solved centrally if all necessary information can be collected at a single point, as may be possible if all the proxies and data centers were owned by the same system. However there is a strong case for Internet-scale systems to outsource route selection [39]. To meet this need, the algorithms presented below are decentralized and allow each data center and proxy to optimize based on partial information.

These algorithms seek to fill a notable hole in the growing literature on algorithms for geographical load balancing. Specifically, they have provable optimality guarantees for a performance objective that includes both energy and delay, where route decisions are made using both energy price and network propagation delay information. The most closely related work [32] investigates the total electricity cost for data centers in a multi-electricity-market environment. It contains the queuing delay inside the data center (assumed to be an $M/M/1$ queue) but neglects the end-to-end user delay. Conversely, [39] uses a simple, efficient algorithm to coordinate the “replica-selection” decisions, but assumes the capacity at each data center is fixed. Other related works, e.g., [32, 34, 28], either do not provide provable guarantees or ignore diverse network delays and/or prices.

Algorithm 1: Gauss-Seidel iteration

Algorithm 1 is motivated by the observation that GLB-Q is separable in m_i , and, less obviously, also separable in

$\lambda_j := (\lambda_{ij}, i \in N)$. This allows all data centers as a group and each proxy j to iteratively solve for optimal \mathbf{m} and λ_j in a distributed manner, and communicate their intermediate results to each other. Though distributed, Algorithm 1 requires each proxy to solve an optimization problem.

To highlight the separation between data centers and proxies, we reformulate GLB-Q as:

$$\min_{\lambda_j \in \Lambda_j} \min_{m_i \in \mathcal{M}_i} \sum_{i \in N} \left(p_i m_i + \frac{\beta \lambda_i}{\mu_i - \lambda_i/m_i} \right) + \beta \sum_{i,j} \lambda_{ij} d_{ij} \quad (7)$$

$$\mathcal{M}_i := [0, M_i] \quad \Lambda_j := \{ \lambda_j \mid \lambda_j \geq 0, \sum_{i \in N} \lambda_{ij} = L_j \} \quad (8)$$

Since the objective and constraints \mathcal{M}_i and Λ_j are separable, this can be solved separately by data centers i and proxies j .

The iterations of the algorithm are indexed by τ , and are assumed to be fast relative to the timeslots t . Each iteration τ is divided into $|J| + 1$ phases. In phase 0, all data centers i concurrently calculate $m_i(\tau + 1)$ based on their own arrival rates $\lambda_i(\tau)$, by minimizing (7) over their own variables m_i :

$$\min_{m_i \in \mathcal{M}_i} \left(p_i m_i + \frac{\beta \lambda_i(\tau)}{\mu_i - \lambda_i(\tau)/m_i} \right) \quad (9)$$

In phase j of iteration τ , proxy j minimizes (7) over its own variable by setting $\lambda_j(\tau + 1)$ as the best response to $\mathbf{m}(\tau + 1)$ and the most recent values of $\lambda_{-j} := (\lambda_k, k \neq j)$. This works because proxy j depends on λ_{-j} only through their aggregate arrival rates at the data centers:

$$\lambda_i(\tau, j) := \sum_{l < j} \lambda_{il}(\tau + 1) + \sum_{l > j} \lambda_{il}(\tau) \quad (10)$$

To compute $\lambda_i(\tau, j)$, proxy j need not obtain individual $\lambda_{il}(\tau)$ or $\lambda_{il}(\tau + 1)$ from other proxies l . Instead, every data center i measures its local arrival rate $\lambda_i(\tau, j) + \lambda_{ij}(\tau)$ in every phase j of the iteration τ and sends this to proxy j at the beginning of phase j . Then proxy j obtains $\lambda_i(\tau, j)$ by subtracting its own $\lambda_{ij}(\tau)$ from the value received from data center i . When there are fewer data centers than proxies, this has less overhead than direct messaging.

In summary, the algorithm is as follows (noting that the minimization (9) has a closed form). Here, $[x]^a := \min\{x, a\}$.

Algorithm 1. *Starting from a feasible initial allocation $\lambda(0)$ and the associated $\mathbf{m}(\lambda(0))$, let*

$$m_i(\tau + 1) := \left[\left(1 + \frac{1}{\sqrt{p_i/\beta}} \right) \cdot \frac{\lambda_i(\tau)}{\mu_i} \right]^{M_i} \quad (11)$$

$$\lambda_j(\tau + 1) := \arg \min_{\lambda_j \in \Lambda_j} \sum_{i \in N} \frac{\lambda_i(\tau, j) + \lambda_{ij}}{\mu_i - (\lambda_i(\tau, j) + \lambda_{ij})/m_i(\tau + 1)} + \sum_{i \in N} \lambda_{ij} d_{ij}. \quad (12)$$

Since GLB-Q generally has multiple optimal λ_j^* , Algorithm 1 is not guaranteed to converge to one optimal solution, i.e., for each proxy j , the allocation $\lambda_{ij}(\tau)$ of job j to data centers i may oscillate among multiple optimal allocations. However, both the optimal cost and the optimal per-server arrival rates to data centers will converge.

Theorem 5. *Let $(\mathbf{m}(\tau), \lambda(\tau))$ be a sequence generated by Algorithm 1 when applied to GLB-Q. Then*

- (i) *Every limit point of $(\mathbf{m}(\tau), \lambda(\tau))$ is optimal.*
- (ii) *$F(\mathbf{m}(\tau), \lambda(\tau))$ converges to the optimal value.*
- (iii) *The per-server arrival rates $(\lambda_i(\tau)/m_i(\tau), i \in N)$ to data centers converge to their unique optimal values.*

The proof of Theorem 5 follows from the fact that Algorithm 1 is a modified Gauss-Seidel iteration. This is also the reason for the requirement that the proxies update sequentially. The details of the proof are in Appendix B.

Algorithm 1 assumes that there is a common clock to synchronize all actions. In practice, updates will likely be asynchronous, with data centers and proxies updating with different frequencies using possibly outdated information. The algorithm generalizes easily to this setting, though the convergence proof is more difficult.

The convergence rate of Algorithm 1 in a realistic scenario is illustrated numerically in Section 5.

Algorithm 2: Distributed gradient projection

Algorithm 2 reduces the computational load on the proxies. In each iteration, instead of each proxy solving a constrained minimization (12) as in Algorithm 1, Algorithm 2 takes a single step in a descent direction. Also, while the proxies compute their $\lambda_j(\tau+1)$ sequentially in $|J|$ phases in Algorithm 1, they perform their updates all at once in Algorithm 2.

To achieve this, rewrite GLB-Q as

$$\min_{\lambda_j \in \Lambda_j} \sum_j F_j(\lambda) \quad (13)$$

where $F(\lambda)$ is the result of minimization of (7) over $m_i \in \mathcal{M}_i$ given λ_i . As explained in the definition of Algorithm 1, this minimization is easy: if we denote the solution by (cf. (11)):

$$m_i(\lambda_i) := \left[\left(1 + \frac{1}{\sqrt{p_i/\beta}} \right) \cdot \frac{\lambda_i}{\mu_i} \right]^{M_i} \quad (14)$$

then

$$F(\lambda) := \sum_{i \in N} \left(p_i m_i(\lambda_i) + \frac{\beta \lambda_i}{\mu_i - \lambda_i / m_i(\lambda_i)} \right) + \beta \sum_{i,j} \lambda_{ij} d_{ij}.$$

We now sketch the two key ideas behind Algorithm 2. The first is the standard gradient projection idea: move in the steepest descent direction

$$-\nabla F_j(\lambda) := - \left(\frac{\partial F(\lambda)}{\partial \lambda_{1j}}, \dots, \frac{\partial F(\lambda)}{\partial \lambda_{|N|j}} \right)$$

and then project the new point into the feasible set $\prod_j \Lambda_j$. The standard gradient projection algorithm will converge if $\nabla F(\lambda)$ is Lipschitz over our feasible set $\prod_j \Lambda_j$. This condition, however, does not hold for our F because of the term $\beta \lambda_i / (\mu_i - \lambda_i / m_i)$. The second idea is to construct a compact and convex subset Λ of the feasible set $\prod_j \Lambda_j$ with the following properties: (i) if the algorithm starts in Λ , it stays in Λ ; (ii) Λ contains all optimal allocations; (iii) $\nabla F(\lambda)$ is Lipschitz over Λ . The algorithm then projects into Λ in each iteration instead of $\prod_j \Lambda_j$. This guarantees convergence.

Specifically, fix a feasible initial allocation $\lambda(0) \in \prod_j \Lambda_j$ and let $\phi := F(\lambda(0))$ be the initial objective value. Define

$$\Lambda := \Lambda(\phi) := \prod_j \Lambda_j \cap \left\{ \lambda \mid \lambda_i \leq \frac{\phi M_i \mu_i}{\phi + \beta M_i}, \forall i \right\}. \quad (15)$$

Even though the Λ defined in (15) indeed has the desired properties (see Appendix B), the projection into Λ requires coordination of all proxies and is thus impractical. In order for each proxy j to perform its update in a decentralized manner, we define proxy j 's own constraint subset:

$$\hat{\Lambda}_j(\tau) := \Lambda_j \cap \left\{ \lambda_j \mid \lambda_i(\tau, -j) + \lambda_{ij} \leq \frac{\phi M_i \mu_i}{\phi + \beta M_i}, \forall i \right\}$$

where $\lambda_i(\tau, -j) := \sum_{l \neq j} \lambda_{il}(\tau)$ is the arrival rate to data center i , excluding arrivals from proxy j . Even though $\hat{\Lambda}_j(\tau)$ involves $\lambda_i(\tau, -j)$ for all i , proxy j can easily calculate these quantities from the measured arrival rates $\lambda_i(\tau)$ it is told by data centers i , as done in Algorithm 1 (cf. (10) and the discussion thereafter), and does not need to communicate with other proxies. Hence, given $\lambda_i(\tau, -j)$ from data centers i , each proxy can project into $\hat{\Lambda}_j(\tau)$ to compute the next iterate $\lambda_j(\tau+1)$ without the need to coordinate with other proxies.³ Moreover, if $\lambda(0) \in \Lambda$ then $\lambda(\tau) \in \Lambda$ for all iterations τ . In summary, Algorithm 2 is as follows.

Algorithm 2. Starting from a feasible initial allocation $\lambda(0)$ and the associated $\mathbf{m}(\lambda(0))$, each proxy j computes, in each iteration τ :

$$\mathbf{z}_j(\tau+1) := [\lambda_j(\tau) - \gamma_j (\nabla F_j(\lambda(\tau)))]_{\hat{\Lambda}_j(\tau)} \quad (16)$$

$$\lambda_j(\tau+1) := \frac{|J|-1}{|J|} \lambda_j(\tau) + \frac{1}{|J|} \mathbf{z}_j(\tau+1) \quad (17)$$

where $\gamma_j > 0$ is a stepsize and $\nabla F_j(\lambda(\tau))$ is given by

$$\frac{\partial F(\lambda(\tau))}{\partial \lambda_{ij}} = \beta \left(d_{ij} + \frac{\mu_i}{(\mu_i - \lambda_i(\tau) / m_i(\lambda_i(\tau)))^2} \right).$$

Implicit in the description is the requirement that all data centers i compute $m_i(\lambda_i(\tau))$ according to (14) in each iteration τ . Each data center i measures the local arrival rate $\lambda_i(\tau)$, calculates $m_i(\lambda_i(\tau))$, and broadcasts these values to all proxies at the beginning of iteration $\tau+1$ for the proxies to compute their $\lambda_j(\tau+1)$.

Algorithm 2 has the same convergence property as Algorithm 1, provided the stepsize is small enough.

Theorem 6. Let $(\mathbf{m}(\tau), \lambda(\tau))$ be a sequence generated by Algorithm 2 when applied to GLB-Q. If, for all j , $0 < \gamma_j < \min_{i \in N} \beta^2 \mu_i^2 M_i^4 / (|J|(\phi + \beta M_i)^3)$, then

- (i) Every limit point of $(\mathbf{m}(\tau), \lambda(\tau))$ is optimal.
- (ii) $F(\mathbf{m}(\tau), \lambda(\tau))$ converges to the optimal value.
- (iii) The per-server arrival rates $(\lambda_i(\tau) / m_i(\tau), i \in N)$ to data centers converge to their unique optimal values.

Theorem 6 is proved in Appendix B. The key novelty of the proof is (i) handling the fact that the objective is not Lipschitz and (ii) allowing distributed computation of the projection. The bound on γ_j in Theorem 6 is more conservative than necessary for large systems. Hence, a larger stepsize can be chosen to accelerate convergence. The convergence rate is illustrated in a realistic setting in Section 5.

5. CASE STUDY

The remainder of the paper evaluates the algorithms presented in the previous section under a realistic workload. This section considers the data center perspective (i.e., cost minimization) and Section 6 considers the social perspective (i.e., brown energy usage).

5.1 Experimental setup

We aim to use realistic parameters in the experimental setup and provide conservative estimates of the cost savings resulting from optimal geographical load balancing. The setup models an Internet-scale system such as Google within the United States.

³The projection to the nearest point in $\hat{\Lambda}_j(\tau)$ is defined by $[\lambda]_{\hat{\Lambda}_j(\tau)} := \arg \min_{y \in \hat{\Lambda}_j(\tau)} \|y - \lambda\|_2$.

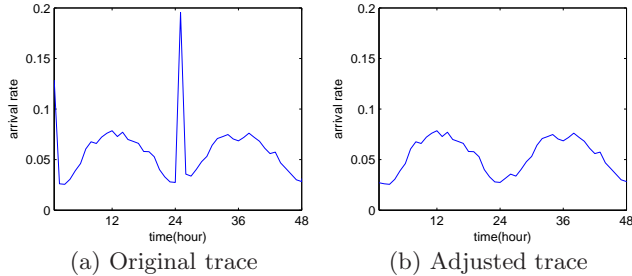


Figure 1: Hotmail trace used in numerical results.

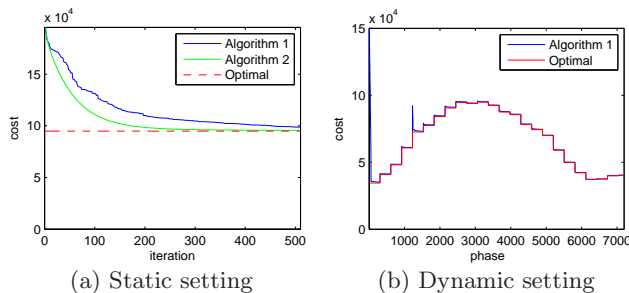


Figure 2: Convergence of Algorithm 1 and 2.

Workload description.

To build our workload, we start with a trace of traffic from Hotmail, a large Internet service running on tens of thousands of servers. The trace represents the I/O activity from 8 servers over a 48-hour period, starting at midnight (PDT) on August 4, 2008, averaged over 10 minute intervals. The trace has strong diurnal behavior and has a fairly small peak-to-mean ratio of 1.64. Results for this small peak-to-mean ratio provide a lower bound on the cost savings under workloads with larger peak-to-mean ratios. As illustrated in Figure 1(a), the Hotmail trace contains significant nightly activity due to maintenance processes; however the data center is provisioned for the peak foreground traffic. This creates a dilemma about whether to include the maintenance activity or not. We have performed experiments with both, but report only the results with the spike removed (as illustrated in Figure 1(b)) because this leads to a more conservative estimate of the cost savings.

Building on this trace, we construct our workload by placing a source at the geographical center of each mainland US state, co-located with a proxy or DNS server (as described in Section 2.4). The trace is shifted according to the time-zone of each state, and scaled by the size of the population in the state that has an Internet connection [1].

Data center description.

To model an Internet-scale system, we have 14 data centers, one at the geographic center of each state known to have Google data centers [17]: California, Washington, Oregon, Illinois, Georgia, Virginia, Texas, Florida, North Carolina, and South Carolina.

We merge the data centers in each state and set M_i proportional to the number of data centers in that state, while keeping $\sum_{i \in N} M_i \mu_i$ twice the total peak workload, $\max_t \sum_{j \in J} L_j(t)$. The network delays, d_{ij} , between sources and data centers are taken to be proportional to the distances between the centers of the two states and comparable to queuing delays. This lower bound on the network delay ignores delay due to congestion or indirect routes.

Cost function parameters.

To model the costs of the system, we use the GLB-Q formulation. We set $\mu_i = 1$ for all i , so that the servers at each location are equivalent. We assume the energy consumption of an active server in one timeslot is normalized to 1. We set constant electricity prices using the industrial electricity price of each state in May 2010 [18]. Specifically, the price (cents per kWh) is 10.41 in California; 3.73 in Washington; 5.87 in Oregon, 7.48 in Illinois; 5.86 in Georgia; 6.67 in Virginia; 6.44 in Texas; 8.60 in Florida; 6.03 in North Carolina; and 5.49 in South Carolina. In this section, we set $\beta = 1$; however Figure 3 illustrates the impact of varying β .

Algorithm benchmarks.

To provide benchmarks for the performance of the algorithms presented here, we consider three baselines, which are approximations of common approaches used in Internet-scale systems. They also allow implicit comparisons with prior work such as [32]. The approaches use different amounts of information to perform the cost minimization. Note that each approach must use queuing delay (or capacity information); otherwise the routing may lead to instability.

Baseline 1 uses network delays but ignores energy price when minimizing its costs. This demonstrates the impact of price-aware routing. It also shows the importance of dynamic capacity provisioning, since without using energy cost in the optimization, every data center will keep every server active.

Baseline 2 uses energy prices but ignores network delay. This illustrates the impact of location aware routing on the data center costs. Further, it allows us to understand the performance improvement of Algorithms 1 and 2 compared to those such as [32, 34] that neglect network delays in their formulations.

Baseline 3 uses neither network delay information nor energy price information when performing its cost minimization. Thus, the traffic is routed so as to balance the delays within the data centers. Though naive, designs such as this are still used by systems today; see [3].

5.2 Performance evaluation

The evaluation of our algorithms and the cost savings due to optimal geographic load balancing will be organized around the following topics.

Convergence.

We start by considering the convergence of each of the distributed algorithms. Figure 2(a) illustrates the convergence of each of the algorithms in a static setting for $t = 11\text{am}$, where load and electricity prices are fixed and each phase in Algorithm 1 is considered as an iteration. It validates the convergence analysis for both algorithms. Note here Algorithm 2 used a step size $\gamma = 10$; this is much larger than that used in the convergence analysis, which is quite conservative, and there is no sign of causing lack of convergence.

To demonstrate the convergence in a dynamic setting, Figure 2(b) shows Algorithm 1's response to the first day of the Hotmail trace, with loads averaged over one-hour intervals for brevity. One iteration (51 phases) is performed every 10 minutes. This figure shows that even the slower algorithm, Algorithm 1, converges fast enough to provide near-optimal cost. Hence, the remaining plots show only the optimal solution.

Energy versus delay tradeoff.

The optimization objective we have chosen to model the data center costs imposes a particular tradeoff between the delay and the energy costs, β . It is important to understand the impact of this factor. Figure 3 illustrates how the delay

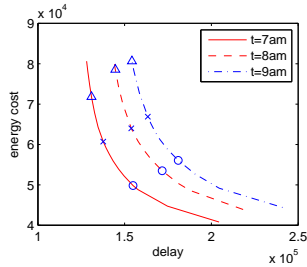


Figure 3: Pareto frontier of the GLB-Q formulation as a function of β for three different times (and thus arrival rates), PDT. Circles, x-marks, and triangles correspond to $\beta = 2.5$, 1, and 0.4, respectively.

and energy cost trade off under the optimal solution as β changes. Thus, the plot shows the Pareto frontier for the GLB-Q formulation. The figure highlights that there is a smooth convex frontier with a mild ‘knee’.

Cost savings.

To evaluate the cost savings of geographical load balancing, Figure 4 compares the optimal costs to those incurred under the three baseline strategies described in the experimental setup. The overall cost, shown in Figures 4(a) and 4(b), is significantly lower under the optimal solution than all of the baselines (nearly 40% during times of light traffic). Recall that Baseline 2 is the state of the art, studied in recent papers such as [32, 34].

To understand where the benefits are coming from, let us consider separately the two components of cost: delay and energy. Figures 4(c) and 4(d) show that the optimal algorithm performs well with respect to both delay and energy costs individually. In particular, Baseline 1 provides a lower bound on the achievable delay costs, and the optimal algorithm nearly matches this lower bound. Similarly, Baseline 2 provides a natural bar for comparing the achievable energy cost. At periods of light traffic the optimal algorithm provides nearly the same energy cost as this baseline, and (perhaps surprisingly) during periods of heavy-traffic the optimal algorithm provides significantly lower energy costs. The explanation for this is that, when network delay is considered by the optimal algorithm, if all the close data centers have all servers active, a proxy might still route to them; however when network delay is not considered, a proxy is more likely to route to a data center that is not yet running at full capacity, thereby adding to the energy cost.

Sensitivity analysis.

Given that the algorithms all rely on estimates of the L_j and d_{ij} it is important to perform a sensitivity analysis to understand the impact of errors in these parameters on the achieved cost. We have performed such a sensitivity analysis but omit the details for brevity. The results show that even when the algorithms have very poor estimates of d_{ij} and L_j there is little effect on cost. Baseline 2 can be thought of as applying the optimal algorithm to very poor estimates of d_{ij} (namely $d_{ij} = 0$), and so the Figure 4(a) provides some illustration of the effect of estimation error.

6. SOCIAL IMPACT

We now shift focus from the cost savings of the data center operator to the social impact of geographical load balancing. We focus on the impact of geographical load balancing on the usage of “brown” non-renewable energy by Internet-scale systems, and how this impact depends on pricing.

Intuitively, geographical load balancing allows the traffic

to “follow the renewables”; thus providing increased usage of green energy and decreased brown energy usage. However, such benefits are only possible if data centers forgo static energy contracts for dynamic energy pricing (either through demand-response programs or real-time markets). The experiments in this section show that if dynamic pricing is done optimally, then geographical load balancing can provide significant social benefits.

6.1 Experimental setup

To explore the social impact of geographical load balancing, we use the setup described in Section 5. However, we add models for the availability of renewable energy, the pricing of renewable energy, and the social objective.

The availability of renewable energy.

To model the availability of renewable energy we use standard models of wind and solar from [15, 20]. Though simple, these models capture the average trends for both wind and solar accurately. Since these models are smoother than actual intermittent renewable sources, especially wind, they conservatively estimate the benefit due to following renewables.

We consider two settings (i) high wind penetration, where 90% of renewable energy comes from wind and (ii) high solar penetration, where 90% of renewable energy comes from solar. The availability given by these models is shown in Figure 5(a). Setting (i) is motivated by studies such as [18]. Setting (ii) is motivated by the possibility of on-site or locally contracted solar, which is increasingly common.

Building on these availability models, for each location we let $\alpha_i(t)$ denote the fraction of the energy that is from renewable sources at time t , and let $\bar{\alpha} = (|N|T)^{-1} \sum_{t=1}^T \sum_{i \in N} \alpha_i(t)$ be the “penetration” of renewable energy. We take $\bar{\alpha} = 0.30$, which is on the progressive side of the renewable targets among US states [10].

Finally, when measuring the brown/green energy usage of a data center at time t , we use simply $\sum_{i \in N} \alpha_i(t)m_i(t)$ as the green energy usage and $\sum_{i \in N} (1 - \alpha_i(t))m_i(t)$ as the brown energy usage. This models the fact that the grid cannot differentiate the source of the electricity provided.

Demand response and dynamic pricing.

Internet-scale systems have flexibility in energy usage that is not available to traditional energy consumers; thus they are well positioned to take advantage of demand-response and real-time markets to reduce both their energy costs and their brown energy consumption.

To provide a simple model of demand-response, we use time-varying prices $p_i(t)$ in each time-slot that depend on the availability of renewable resources $\alpha_i(t)$ in each location.

The way $p_i(t)$ is chosen as a function of $\alpha_i(t)$ will be of fundamental importance to the social impact of geographical load balancing. To highlight this, we consider a parameterized “differentiated pricing” model that uses a price p_b for brown energy and a price p_g for green energy. Specifically,

$$p_i(t) = p_b(1 - \alpha_i(t)) + p_g\alpha_i(t).$$

Note that $p_g = p_b$ corresponds to static pricing, and we show in the next section that $p_g = 0$ corresponds to socially optimal pricing. Our experiments vary $p_g \in [0, p_b]$.

The social objective.

To model the social impact of geographical load balancing we need to formulate a social objective. Like the GLB formulation, this must include a tradeoff between the energy usage and the delay users of the system experience, because

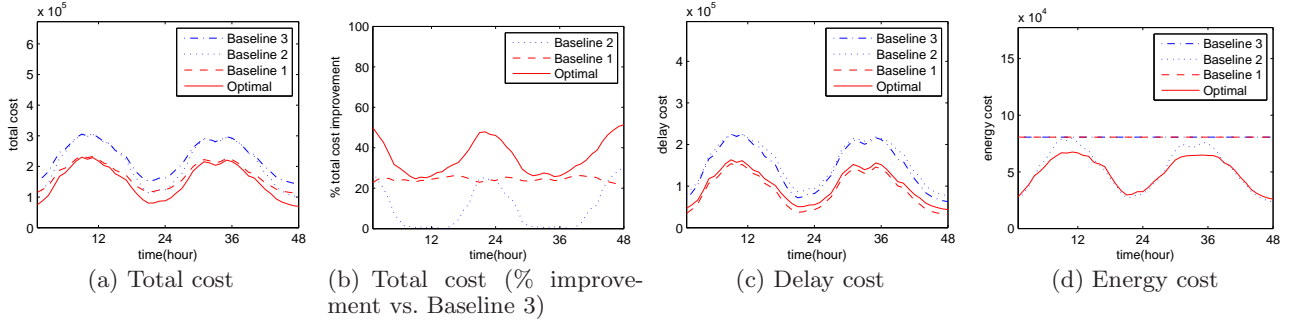


Figure 4: Impact of information used on the cost incurred by geographical load balancing.

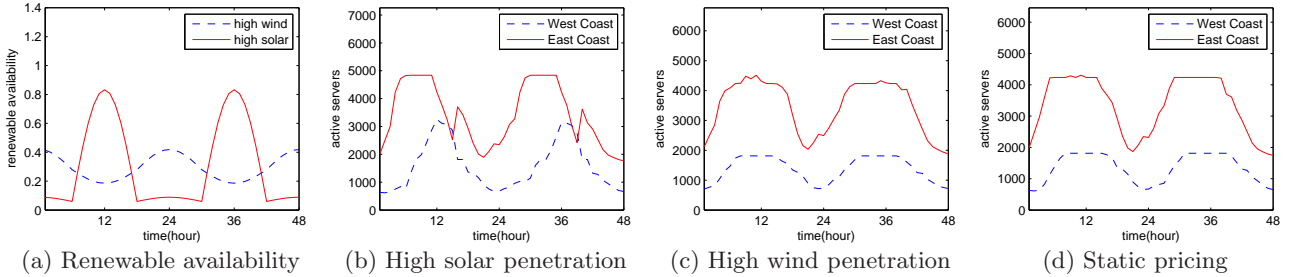


Figure 5: Geographical load balancing “following the renewables” under optimal pricing. (a) Availability of wind and solar. (b)–(d) Capacity provisioning of east coast and west coast data centers when there are renewables, high solar penetration, and high wind penetration, respectively.

purely minimizing brown energy use requires all $m_i = 0$. The key difference between the GLB formulation and the social formulation is that the *cost* of energy is no longer relevant. Instead, the environmental impact is important, and thus the brown energy usage should be minimized. This leads to the following simple model for the social objective:

$$\min_{\mathbf{m}(t), \lambda(t)} \sum_{t=1}^T \sum_{i \in N} \left((1 - \alpha_i(t)) \frac{\mathcal{E}_i(t)}{p_i(t)} + \tilde{\beta} \mathcal{D}_i(t) \right) \quad (18)$$

where $\mathcal{D}_i(t)$ is the delay cost defined in (2), $\mathcal{E}_i(t)$ is the energy cost defined in (1), and $\tilde{\beta}$ is the relative valuation of delay versus energy. Further, we have imposed that the energy cost follows from the pricing of $p_i(t)$ cents/kWh in timeslot t . Note that, though simple, our choice of $\mathcal{D}_i(t)$ to model the disutility of delay to users is reasonable because lost revenue captures the lack of use as a function of increased delay.

An immediate observation about the above social objective is that to align the data center and social goals, one needs to set $p_i(t) = (1 - \alpha_i(t))/\tilde{\beta}$, which corresponds to choosing $p_b = 1/\tilde{\beta}$ and $p_g = 0$ in the differentiated pricing model above. We refer to this as the “optimal” pricing model.

6.2 The importance of dynamic pricing

To begin our experiments, we illustrate that optimal pricing can lead geographical load balancing to “follow the renewables.” Figure 5 shows this in the case of high solar penetration and high wind penetration for $\tilde{\beta} = 0.1$. By comparing Figures 5(b) and 5(c) to Figure 5(d), which uses static pricing, the change in capacity provisioning, and thus energy usage, is evident. For example, Figure 5(b) shows a clear shift of service capacity from the east coast to the west coast as solar energy becomes highly available and then

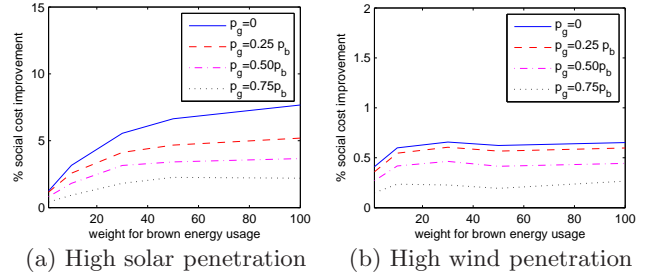


Figure 6: Reduction in social cost from dynamic pricing compared to static pricing as a function of the weight for brown energy usage, $1/\tilde{\beta}$, under (a) high solar penetration and (b) high wind penetration.

back when solar energy is less available. Similarly, Figure 5(c) shows a shift, though much smaller, of service capacity toward the evenings, when wind is more available. Though not explicit in the figures, this “follow the renewables” routing has the benefit of significantly reducing the brown energy usage since energy use is more correlated with the availability of renewables. Thus, geographical load balancing provides the opportunity to aid the incorporation of renewables into the grid.

Figure 5 assumed the optimal dynamic pricing, but currently data centers negotiate fixed price contracts. Although there are many reasons why grid operators will encourage data center operators to transfer to dynamic pricing over the coming years, this is likely to be a slow process. Thus, it is important to consider the impact of partial adoption of dynamic pricing in addition to full, optimal dynamic pricing.

Figure 6 focuses on this issue. To model the partial adoption of dynamic pricing, we can consider $p_g \in [0, p_b]$. Figure 6(a) shows that the benefits provided by dynamic pricing are moderate but significant, even at partial adoption (high p_g), when there is high solar penetration. Figure 6(b) suggests that there would be much less benefit if renewable sources were dominated by wind with only diurnal variation, because the availability of solar energy is much more correlated with the traffic peaks. Specifically, the three hour gap in time zones means that solar on the west coast can still help with the high traffic period of the east coast, but the peak average wind energy is at night. However, wind is vastly more bursty than this model predicts, and a system which responds to these bursts will still benefit significantly.

Another interesting observation about the Figure 6 is that the curves increase faster in the range when $\tilde{\beta}$ is large, which highlights that the social benefit of geographical load balancing becomes significant even when there is only moderate importance placed on energy. When p_g is higher than p_b , which is common currently, the cost increases, but we omit the results due to space constraints.

7. CONCLUDING REMARKS

This paper has focused on understanding algorithms for and social impacts of geographical load balancing in Internet-scaled systems. We have provided two distributed algorithms that provably compute the optimal routing and provisioning decisions for Internet-scale systems and we have evaluated these algorithms using trace-based numerical simulations. Further, we have studied the feasibility and benefits of providing demand response for the grid via geographical load balancing. Our experiments highlight that geographical load balancing can provide an effective tool for demand-response: when pricing is done carefully electricity providers can motivate Internet-scale systems to “follow the renewables” and route to areas where green energy is available. This both eases the incorporation of renewables into the grid and reduces brown energy consumption of Internet-scale systems.

There are a number of interesting directions for future work that are motivated by the studies in this paper. With respect to the design of distributed algorithms, one aspect that our model has ignored is the switching cost (in terms of delay and wear-and-tear) associated with switching servers into and out of power-saving modes. Our model also ignores issues related to reliability and availability, which are quite important in practice. With respect to the social impact of geographical load balancing, our results highlight the opportunity provided by geographical load balancing for demand response; however there are many issues left to be considered. For example, which demand response market should Internet-scale systems participate in to minimize costs? How can policy decisions such as cap-and-trade be used to provide the proper incentives for Internet-scale systems, such as [23]? Can Internet-scale systems use energy storage at data centers in order to magnify cost reductions when participating in demand response markets? Answering these questions will pave the way for greener geographic load balancing.

8. ACKNOWLEDGMENTS

This work was supported by NSF grants CCF 0830511, CNS 0911041, and CNS 0846025, DoE grant DE-EE0002890, ARO MURI grant W911NF-08-1-0233, Microsoft Research, Bell Labs, the Lee Center for Advanced Networking, and ARC grant FT0991594.

9. REFERENCES

[1] US Census Bureau, <http://www.census.gov>.

- [2] Server and data center energy efficiency, Final Report to Congress, U.S. Environmental Protection Agency, 2007.
- [3] V. K. Adhikari, S. Jain, and Z.-L. Zhang. YouTube traffic dynamics and its interplay with a tier-1 ISP: An ISP perspective. In *ACM IMC*, pages 431–443, 2010.
- [4] S. Albers. Energy-efficient algorithms. *Comm. of the ACM*, 53(5):86–96, 2010.
- [5] L. L. H. Andrew, M. Lin, and A. Wierman. Optimality, fairness and robustness in speed scaling designs. In *Proc. ACM Sigmetrics*, 2010.
- [6] A. Beloglazov, R. Buyya, Y. C. Lee, and A. Zomaya. A taxonomy and survey of energy-efficient data centers and cloud computing systems, Technical Report, 2010.
- [7] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, 1999.
- [8] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, 1989.
- [9] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [10] S. Carley. State renewable energy electricity policies: An empirical evaluation of effectiveness. *Energy Policy*, 37(8):3071–3081, Aug 2009.
- [11] Y. Chen, A. Das, W. Qin, A. Sivasubramaniam, Q. Wang, and N. Gautam. Managing server energy and operational costs in hosting centers. In *Proc. ACM Sigmetrics*, 2005.
- [12] M. Conti and C. Nazionale. Load distribution among replicated web servers: A QoS-based approach. In *Proc. ACM Worksh. Internet Server Performance*, 1999.
- [13] A. Croll and S. Power. How web speed affects online business KPIs. <http://www.watchingwebsites.com>, 2009.
- [14] X. Fan, W.-D. Weber, and L. A. Barroso. Power provisioning for a warehouse-sized computer. In *Proc. Int. Symp. Comp. Arch.*, 2007.
- [15] M. Fripp and R. H. Wiser. Effects of temporal wind patterns on the value of wind-generated electricity in California and the northwest. *IEEE Trans. Power Systems*, 23(2):477–485, May 2008.
- [16] A. Gandhi, M. Harchol-Balter, R. Das, and C. Lefurgy. Optimal power allocation in server farms. In *Proc. ACM Sigmetrics*, 2009.
- [17] <http://www.datacenterknowledge.com>, 2008.
- [18] <http://www.eia.doe.gov>.
- [19] S. Irani and K. R. Pruhs. Algorithmic problems in power management. *SIGACT News*, 36(2):63–76, 2005.
- [20] T. A. Kattakayam, S. Khan, and K. Srinivasan. Diurnal and environmental characterization of solar photovoltaic panels using a PC-AT add on plug in card. *Solar Energy Materials and Solar Cells*, 44(1):25–36, Oct 1996.
- [21] S. Kaxiras and M. Martonosi. *Computer Architecture Techniques for Power-Efficiency*. Morgan & Claypool, 2008.
- [22] R. Krishnan, H. V. Madhyastha, S. Srinivasan, S. Jain, A. Krishnamurthy, T. Anderson, and J. Gao. Moving beyond end-to-end path information to optimize CDN performance. In *Proc. ACM Sigcomm*, 2009.
- [23] K. Le, R. Bianchini, T. D. Nguyen, O. Bilgir, and M. Martonosi. Capping the brown energy consumption of internet services at low cost. In *Proc. IGCC*, 2010.
- [24] M. Lin, A. Wierman, L. L. H. Andrew, and E. Thereska. Dynamic right-sizing for power-proportional data centers. In *Proc. IEEE INFOCOM*, 2011.
- [25] Z. M. Mao, C. D. Cranor, F. Bouglis, M. Rabinovich, O. Spatscheck, and J. Wang. A precise and efficient evaluation of the proximity between web clients and their local DNS servers. In *USENIX*, pages 229–242, 2002.
- [26] E. Ng and H. Zhang. Predicting internet network distance with coordinates-based approaches. In *Proc. IEEE INFOCOM*, 2002.
- [27] S. Ong, P. Denholm, and E. Doris. The impacts of commercial electric utility rate structure elements on the economics of photovoltaic systems. Technical Report NREL/TP-6A2-46782, National Renewable Energy Laboratory, 2010.
- [28] E. Pakbaznia and M. Pedram. Minimizing data center cooling and server power costs. In *Proc. ISLPED*, 2009.
- [29] J. Pang, A. Akella, A. Shaikh, B. Krishnamurthy, and

- S. Seshan. On the responsiveness of DNS-based network control. In *Proc. IMC*, 2004.
- [30] M. Pathan, C. Vecchiola, and R. Buyya. Load and proximity aware request-redirection for dynamic load distribution in peering CDNs. In *Proc. OTM*, 2008.
- [31] A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag, and B. Maggs. Cutting the electric bill for internet-scale systems. In *Proc. ACM Sigcomm*, Aug. 2009.
- [32] L. Rao, X. Liu, L. Xie, and W. Liu. Minimizing electricity cost: Optimization of distributed internet data centers in a multi-electricity-market environment. In *INFOCOM*, 2010.
- [33] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [34] R. Stanojevic and R. Shorten. Distributed dynamic speed scaling. In *Proc. IEEE INFOCOM*, 2010.
- [35] W. Theilmann and K. Rothermel. Dynamic distance maps of the internet. In *Proc. IEEE INFOCOM*, 2001.
- [36] E. Thereska, A. Donnelly, and D. Narayanan. Sierra: a power-proportional, distributed storage system. Technical Report MSR-TR-2009-153, Microsoft Research, 2009.
- [37] O. S. Unsal and I. Koren. System-level power-aware design techniques in real-time systems. *Proc. IEEE*, 91(7):1055–1069, 2003.
- [38] R. Urgaonkar, U. C. Kozat, K. Igarashi, and M. J. Neely. Dynamic resource allocation and power management in virtualized data centers. In *IEEE NOMS*, Apr. 2010.
- [39] P. Wendell, J. W. Jiang, M. J. Freedman, and J. Rexford. Donar: decentralized server selection for cloud services. In *Proc. ACM Sigcomm*, pages 231–242, 2010.
- [40] A. Wierman, L. L. H. Andrew, and A. Tang. Power-aware speed scaling in processor sharing systems. In *Proc. IEEE INFOCOM*, 2009.

APPENDIX

A. PROOFS FOR SECTION 3

We now prove the results from Section 3, beginning with the illuminating Karush-Kuhn-Tucker (KKT) conditions.

A.1 Optimality conditions

As GLB-Q is convex and satisfies Slater's condition, the KKT conditions are necessary and sufficient for optimality [9]; for the other models they are merely necessary.

GLB-Q: Let $\underline{\omega}_i \geq 0$ and $\bar{\omega}_i \geq 0$ be Lagrange multipliers corresponding to (4d), and $\delta_{ij} \geq 0$, ν_j and σ_i be those for (4c), (4b) and (6b). The Lagrangian is then

$$\begin{aligned} \mathcal{L} = & \sum_{i \in N} m_i p_i + \beta \sum_{j \in J} \sum_{i \in N} \left(\frac{\lambda_{ij}}{\mu_i - \lambda_i/m_i} + \lambda_{ij} d_{ij} \right) \\ & - \sum_{i \in N} \sum_{j \in J} \delta_{ij} \lambda_{ij} + \sum_{j \in J} \nu_j \left(L_j - \sum_{i \in N} \lambda_{ij} \right) \\ & + \sum_{i \in N} (\bar{\omega}_i (m_i - M_i) - \underline{\omega}_i m_i) + \sum_{i \in N} \sigma_i (m_i \mu_i - \lambda_i) \end{aligned}$$

The KKT conditions of stationarity, primal and dual feasibility and complementary slackness are:

$$\beta \left(\frac{\mu_i}{(\mu_i - \lambda_i/m_i)^2} + d_{ij} \right) - \nu_j - \delta_{ij} - \sigma_i = 0 \quad (19)$$

$$\delta_{ij} \lambda_{ij} = 0; \quad \delta_{ij} \geq 0, \quad \lambda_{ij} \geq 0 \quad (20)$$

$$\sigma_i (m_i \mu_i - \lambda_i) = 0; \quad \sigma_i \geq 0, \quad m_i \mu_i - \lambda_i \geq 0 \quad (21)$$

$$\sum_{i \in N} \lambda_{ij} = L_j \quad (22)$$

$$p_i - \beta \left(\frac{\lambda_i/m_i}{\mu_i - \lambda_i/m_i} \right)^2 + \bar{\omega}_i - \underline{\omega}_i + \sigma_i \mu_i = 0 \quad (23)$$

$$\bar{\omega}_i (m_i - M_i) = 0; \quad \bar{\omega}_i \geq 0, \quad m_i \leq M_i \quad (24)$$

$$\underline{\omega}_i m_i = 0; \quad \underline{\omega}_i \geq 0, \quad m_i \geq 0. \quad (25)$$

The conditions (19)–(22) determine the sources' choice of λ_{ij} , and we claim they imply that source j will only send data to those data centers i which have minimum marginal cost $d_{ij} + (1 + \sqrt{p_i^*/\beta})^2/\mu_i$, where $p_i^* = p_i - \underline{\omega}_i + \bar{\omega}_i$. To see this, let $\bar{\lambda}_i = \lambda_i/m_i$. By (23), the marginal queueing delay of data centre i with respect to load λ_{ij} is $\mu_i/(\mu_i - \bar{\lambda}_i)^2 = (1 + \sqrt{p_i^*/\beta})^2/\mu_i$. Thus, from (19), at the optimal point,

$$d_{ij} + \frac{(1 + \sqrt{p_i^*/\beta})^2}{\mu_i} = d_{ij} + \frac{\mu_i}{(\mu_i - \bar{\lambda}_i)^2} = \frac{\nu_j + \delta_j}{\beta} \geq \frac{\nu_j}{\beta} \quad (26)$$

with equality if $\lambda_{ij} > 0$ by (20), establishing the claim.

Note that the solution to (19)–(22) for source j depends on λ_{ik} , $k \neq j$, only through m_i . Given λ_i , data center i finds m_i as the projection onto $[0, M_i]$ of the solution $\hat{m}_i = \lambda_i(1 + \sqrt{p_i/\beta})/(\mu_i \sqrt{p_i/\beta})$ of (23) with $\bar{\omega}_i = \underline{\omega}_i = 0$.

GLB-LIN again decouples into data centers finding m_i given λ_i , and sources finding λ_{ij} given the m_i . Feasibility and complementary slackness conditions (20), (22), (24) and (25) are as for GLB-Q; the stationarity conditions are:

$$\frac{\partial g_i(m_i, \lambda_i)}{\partial \lambda_i} + \beta \left(\frac{\partial (\lambda_i f_i(m_i, \lambda_i))}{\partial \lambda_i} + d_{ij} \right) - \nu_j - \delta_{ij} = 0 \quad (27)$$

$$\frac{\partial g_i(m_i, \lambda_i)}{\partial m_i} + \beta \lambda_i \frac{\partial f_i(m_i, \lambda_i)}{\partial m_i} + \bar{\omega}_i - \underline{\omega}_i = 0. \quad (28)$$

Note the feasibility constraint (6b) of GLB-Q is no longer required to ensure stability. In GLB-LIN, it is instead assumed that f is infinite when the load exceeds capacity.

The objective function is strictly convex in data center i 's decision variable m_i , and so there is a unique solution $\hat{m}_i(\lambda_i)$ to (28) for $\bar{\omega}_i = \underline{\omega}_i = 0$, and the optimal m_i given λ_i is the projection of this onto the interval $[0, M_i]$.

GLB in its general form has the same KKT conditions as GLB-LIN, with the stationary conditions replaced by

$$\begin{aligned} \frac{\partial g_i}{\partial \lambda_i} + r(f_i + d_{ij}) + \sum_{k \in J} \lambda_{ik} r'(f_i + d_{ik}) \frac{\partial f_i}{\partial \lambda_i} \\ - \nu_j - \delta_{ij} = 0 \quad (29) \end{aligned}$$

$$\frac{\partial g_i}{\partial m_i} + \sum_{j \in J} \lambda_{ij} r'(f_i + d_{ij}) \frac{\partial f_i}{\partial m_i} + \bar{\omega}_i - \underline{\omega}_i = 0 \quad (30)$$

where r' denotes the derivative of $r(\cdot)$.

GLB again decouples, since it is convex because $r(\cdot)$ is convex and increasing. However, now data center i 's problem depends on all λ_{ij} , rather than simply λ_i .

A.2 Characterizing the optima

Lemma 7 will help prove the results of Section 3.

Lemma 7. *Consider the GLB-LIN formulation. Suppose that for all i , $F_i(m_i, \lambda_i)$ is jointly convex in λ_i and m_i , and differentiable in λ_i where it is finite. If, for some i , the dual variable $\bar{\omega}_i > 0$ for an optimal solution, then $m_i = M_i$ for all optimal solutions. Conversely, if $m_i < M_i$ for an optimal solution, then $\bar{\omega}_i = 0$ for all optimal solutions.*

PROOF. Consider an optimal solution S with $i \in N$ such that $\bar{\omega}_i > 0$ and hence $m_i = M_i$. Let S' be some other optimal solution.

Since the cost function is jointly convex in λ_{ij} and m_i , any convex combination of S and S' must also be optimal. Let $m_i(s)$ denote the m_i value of a given solution s . Since $m_i(S) = M_i$, we have $\lambda_i > 0$ and so the optimality of S implies f_i is finite at S and hence differentiable. By (28) and the continuity of the partial derivative [33, Corollary 25.51], there is a neighborhood \mathcal{N} of S within which all optimal solutions have $\bar{\omega}_i > 0$, and hence $m_i(s) = M_i$ for all $s \in \mathcal{N}$.

Since $S + \epsilon(S' - S) \in \mathcal{N}$ for sufficiently small ϵ , the linearity of $m_i(s)$ implies $M_i = m_i(S + \epsilon(S' - S)) = m_i(S) + \epsilon(m_i(S') - m_i(S))$. Thus $m_i(S') = m_i(S) = M_i$. \square

PROOF OF THEOREM 1. Consider first the case where there exists an optimal solution with $m_i < M_i$. By Lemma 7, $\bar{\omega}_i = 0$ for all optimal solutions. Recall that $\hat{m}_i(\lambda_i)$, which defines the optimal m_i , is strictly convex. Thus, if different optimal solutions have different values of λ_i , then a convex combination of the two yielding (m'_i, λ'_i) would have $\hat{m}_i(\lambda'_i) < m'_i$, which contradicts the optimality of m'_i .

Next consider the case where all optimal solutions have $m_i = M_i$. In this case, consider two solutions S and S' that both have $m_i = M_i$. If λ_i is the same under both S and S' , we are done. Otherwise, let the set of convex combinations of S and S' be denoted $\{s(\lambda_i)\}$, where we have made explicit the parameterization by λ_i . The convexity of each F_k in m_k and λ_k implies that $F(s(\lambda_i)) - F_i(s(\lambda_i))$ is also convex, due to the fact that the parameterization is by definition affine. Further, since F_i is strictly convex in λ_i , this implies $F(s(\lambda_i))$ is strictly convex in λ_i , and hence has a unique optimal λ_i . \square

PROOF OF THEOREM 2. The proof when $m_i = M_i$ for all optimal solutions is identical to that of Theorem 1. Otherwise, when $m_i < M_i$ in an optimal solution, the definition of \hat{m} gives $\bar{\lambda}_i = \mu_i \sqrt{p_i/\beta_i} / (\sqrt{p_i/\beta_i} + 1)$ for all optimal solutions. \square

PROOF OF THEOREM 3. For each optimal solution S , consider an undirected bipartite graph G with a vertex representing each source and each data center and with an edge connecting i and j when $\lambda_{ij} > 0$. We will show that at least one of these graphs is acyclic. The theorem then follows since an acyclic graph with K nodes has at most $K - 1$ edges.

To prove that there exists one optimal solution with acyclic graph we will inductively reroute traffic in a way that removes cycles while preserving optimality. Suppose G contains a cycle. Let C be a minimal cycle, i.e., no strict subset of C is a cycle, and let C be directed.

Form a new solution $S(\xi)$ from S by adding ξ to λ_{ij} if $(i, j) \in C$, and subtracting ξ from λ_{ij} if $(j, i) \in C$. Note that this does not change the λ_i . To see that $S(\xi)$ maintains the optimal cost, first note that the change in the objective function of the GLB between S and $S(\xi)$ is equal to

$$\xi \left(\sum_{(j,i) \in C} r(d_{ij} + f_i(m_i, \lambda_i)) - \sum_{(i,j) \in C} r(d_{ij} + f_i(m_i, \lambda_i)) \right) \quad (31)$$

Next note that the multiplier $\delta_{ij} = 0$ since $\lambda_{ij} > 0$ at S . Further, the KKT condition (29) for stationarity in λ_{ij} can be written as $K_i + r(d_{ij} + f_i(m_i, \lambda_i)) - \nu_j = 0$, where K_i does not depend on the choice of j .

Since C is minimal, for each $(i, j) \in C$ where $i \in I$ and $j \in J$ there is exactly one (j', i) with $j' \in J$, and vice versa. Thus,

$$\begin{aligned} 0 &= \sum_{(j,i) \in C} (K_i + r(d_{ij} + f_i(m_i, \lambda_i)) - \nu_j) \\ &\quad - \sum_{(i,j) \in C} (K_i + r(d_{ij} + f_i(m_i, \lambda_i)) - \nu_j) \\ &= \sum_{(j,i) \in C} r(d_{ij} + f_i(m_i, \lambda_i)) - \sum_{(i,j) \in C} r(d_{ij} + f_i(m_i, \lambda_i)). \end{aligned}$$

Hence, by (31) the objective of $S(\xi)$ and S are the same.

To complete the proof, we let $(i^*, j^*) = \arg \min_{(i,j) \in C} \lambda_{ij}$. Then $S(\lambda_{i^*, j^*})$ has $\lambda_{i^*, j^*} = 0$. Thus, $S(\lambda_{i^*, j^*})$ has at least one fewer cycle, since it has broken C . Further, by construction, it is still optimal. \square

PROOF OF THEOREM 4. It is sufficient to show that, if $\lambda_{kj} \lambda_{k'j} > 0$ then either $m_k = M_k$ or $m_{k'} = M_{k'}$. Consider a case when $\lambda_{kj} \lambda_{k'j} > 0$.

For a generic i , define $c_i = (1 + \sqrt{p_i/\beta})^2 / \mu_i$ as the marginal cost (26) when the Lagrange multipliers $\bar{\omega}_i = \underline{\omega}_i = 0$. Since the p_i are chosen from a continuous distribution, we have that with probability 1

$$c_k - c_{k'} \neq d_{k'j} - d_{kj}. \quad (32)$$

However, (26) holds with equality if $\lambda_{ij} > 0$, and so $d_{kj} + (1 + \sqrt{p_k^*/\beta})^2 / \mu_k = d_{k'j} + (1 + \sqrt{p_{k'}^*/\beta})^2 / \mu_{k'}$. By the definition of c_i and (32), this implies either $p_k^* \neq p_k$ or $p_{k'}^* \neq p_{k'}$. Hence at least one of the Lagrange multipliers $\underline{\omega}_k, \bar{\omega}_k, \underline{\omega}_{k'}$ or $\bar{\omega}_{k'}$ must be non-zero. However, $\underline{\omega}_i > 0$ would imply $m_i = 0$ whence $\lambda_{ij} = 0$ by (21), which is false by hypothesis, and so either $\bar{\omega}_k$ or $\bar{\omega}_{k'}$ is non-zero, giving the result by (24). \square

B. PROOFS FOR SECTION 4

Algorithm 1

To prove Theorem 5 we apply a variant of Proposition 3.9 of Ch 3 in [8], which gives that if

- (i) $F(\mathbf{m}, \boldsymbol{\lambda})$ is continuously differentiable and convex in the convex feasible region (4b)–(4d);
- (ii) Every limit point of the sequence is feasible;
- (iii) Given the values of $\boldsymbol{\lambda}_{-j}$ and \mathbf{m} , there is a unique minimizer of F with respect to λ_j , and given $\boldsymbol{\lambda}$ there is a unique minimizer of F with respect to \mathbf{m} .

Then, every limit point of $(\mathbf{m}(\tau), \boldsymbol{\lambda}(\tau))_{\tau=1,2,\dots}$ is an optimal solution of GLB-Q.

This differs slightly from [8] in that the requirement that the feasible region be closed is replaced by the feasibility of all limit points, and the requirement of strict convexity with respect to each component is replaced by the existence of a unique minimizer. However, the proof is unchanged.

PROOF OF THEOREM 5. To apply the above to prove Theorem 5, we need to show that $F(\mathbf{m}, \boldsymbol{\lambda})$ satisfies the differentiability and continuity constraints under the GLB-Q model.

GLB-Q is continuously differentiable and, as noted in Appendix A.1, a convex problem. To see that every limit point is feasible, note that the only infeasible points in the closure of the feasible region are those with $m_i \mu_i = \lambda_i$. Since the objective approaches ∞ approaching that boundary, and Gauss-Seidel iterations always reduce the objective [8], these points cannot be limit points.

It remains to show the uniqueness of the minimum in \mathbf{m} and each λ_j . Since the cost is separable in the m_i , it is sufficient to show that this applies with respect to each m_i individually. If $\lambda_i = 0$, then the unique minimizer is $m_i = 0$. Otherwise

$$\frac{\partial^2 F(\mathbf{m}, \boldsymbol{\lambda})}{\partial m_i^2} = 2\beta\mu_i \frac{\lambda_i^2}{(m_i\mu_i - \lambda_i)^3}$$

which by (6b) is strictly positive. The Hessian of $F(\mathbf{m}, \boldsymbol{\lambda})$ with respect to $\boldsymbol{\lambda}_j$ is diagonal with i th element

$$2\beta\mu_i \frac{m_i^2}{(m_i\mu_i - \lambda_i)^3} > 0$$

which is positive definite except the points where some $m_i = 0$. However, if $m_i = 0$, the unique minimum is $\lambda_{ij} = 0$. Note we cannot have all $m_i = 0$. Except these points, $F(\mathbf{m}, \boldsymbol{\lambda})$ is strictly convex in $\boldsymbol{\lambda}_j$ given \mathbf{m} and $\boldsymbol{\lambda}_{-j}$. Therefore $\boldsymbol{\lambda}_j$ is unique given \mathbf{m} .

Part (ii) of Theorem 5 follows from part (i) and the continuity of $F(\mathbf{m}, \boldsymbol{\lambda})$. Part (iii) follows from part (i) and Theorem 2, which provides the uniqueness of optimal per-server arrival rates $(\lambda_i(\tau)/m_i(\tau), i \in N)$. \square

Algorithm 2

As discussed in the section on Algorithm 2, we will prove Theorem 6 in three steps. First, we will show that, starting from an initial feasible point $\boldsymbol{\lambda}(0)$, Algorithm 2 generates a sequence $\boldsymbol{\lambda}(\tau)$ that lies in the set $\Lambda := \Lambda(\phi)$ defined in (15), for $\tau = 0, 1, \dots$. Moreover, $\nabla F(\boldsymbol{\lambda})$ is Lipschitz over Λ . Finally, this implies that $F(\boldsymbol{\lambda}(\tau))$ moves in a descent direction that guarantees convergence.

Lemma 8. *Given an initial point $\boldsymbol{\lambda}(0) \in \prod_j \Lambda_j$, let $\phi := F(\boldsymbol{\lambda}(0))$. Then*

1. $\boldsymbol{\lambda}(0) \in \Lambda := \Lambda(\phi)$;
2. If $\boldsymbol{\lambda}^*$ is optimal then $\boldsymbol{\lambda}^* \in \Lambda$;
3. If $\boldsymbol{\lambda}(\tau) \in \Lambda$, then $\boldsymbol{\lambda}(\tau+1) \in \Lambda$.

PROOF. We claim $F(\boldsymbol{\lambda}) \leq \phi$ implies $\boldsymbol{\lambda} \in \Lambda$. This is true because $\phi \geq F(\boldsymbol{\lambda}) \geq \sum_k \frac{\beta \lambda_k}{\mu_k - \lambda_k / m_k(\lambda_k)} \geq \frac{\beta \lambda_i}{\mu_i - \lambda_i / m_i(\lambda_i)} \geq \frac{\beta \lambda_i}{\mu_i - \lambda_i / M_i}, \forall i$. Therefore $\lambda_i \leq \frac{\phi}{\phi + \beta M_i} M_i \mu_i, \forall i$. Consequently, the initial point $\boldsymbol{\lambda}(0) \in \Lambda$ and the optimal point $\boldsymbol{\lambda}^* \in \Lambda$ because $F(\boldsymbol{\lambda}^*) \leq F(\boldsymbol{\lambda})$.

Next we show that $\boldsymbol{\lambda}(\tau) \in \Lambda$ implies $\mathbf{Z}^j(\tau+1) \in \Lambda$, where $\mathbf{Z}^j(\tau+1)$ is $\boldsymbol{\lambda}(\tau)$ except $\lambda_j(\tau)$ is replaced by $\mathbf{z}_j(\tau)$. This holds because $Z_{ik}^j(\tau+1) = \lambda_{ik}(\tau) \geq 0, \forall k \neq j, \forall i$ and $\sum_i Z_{ik}^j(\tau+1) = \sum_i \lambda_{ik}(\tau) = L_k, \forall k \neq j$. From the definition of the projection on $\hat{\Lambda}_j(\tau)$, $Z_{ij}^j(\tau+1) \geq 0, \forall i$, $\sum_i Z_{ij}^j(\tau+1) = L_j$, and $\sum_k Z_{ik}^j(\tau+1) \leq \frac{\phi}{\phi + \beta M_i} M_i \mu_i, \forall i$. These together ensure $\mathbf{Z}^j(\tau+1) \in \Lambda$.

The update $\boldsymbol{\lambda}_j(\tau+1) = \frac{|J|-1}{|J|} \boldsymbol{\lambda}_j(\tau) + \frac{1}{|J|} \mathbf{z}_j(\tau), \forall j$ is equivalent to $\boldsymbol{\lambda}(\tau+1) = \frac{\sum_j \mathbf{Z}^j(\tau+1)}{|J|}$. Then from the convexity of Λ , we have $\boldsymbol{\lambda}(\tau+1) \in \Lambda$. \square

Let $F(\mathbf{M}, \boldsymbol{\lambda})$ be the total cost when all data centers use all servers, and $\nabla F(\mathbf{M}, \boldsymbol{\lambda})$ be the derivatives with respect to $\boldsymbol{\lambda}$. To prove that $\nabla F(\boldsymbol{\lambda})$ is Lipschitz over Λ , we need the following intermediate result. We omit the proof due to space constraint.

Lemma 9. *For all $\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b \in \Lambda$, we have*

$$\left\| \nabla F(\boldsymbol{\lambda}^b) - \nabla F(\boldsymbol{\lambda}^a) \right\|_2 \leq \left\| \nabla F(\mathbf{M}, \boldsymbol{\lambda}^b) - \nabla F(\mathbf{M}, \boldsymbol{\lambda}^a) \right\|_2.$$

Lemma 10. $\left\| \nabla F(\boldsymbol{\lambda}^b) - \nabla F(\boldsymbol{\lambda}^a) \right\|_2 \leq K \left\| \boldsymbol{\lambda}^b - \boldsymbol{\lambda}^a \right\|_2, \forall \boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b \in \Lambda$, where $K = |J| \max_i 2(\phi + \beta M_i)^3 / (\beta^2 M_i^4 \mu_i^2)$.

PROOF. Following Lemma 9, here we continue to show $\left\| \nabla F(\mathbf{M}, \boldsymbol{\lambda}^b) - \nabla F(\mathbf{M}, \boldsymbol{\lambda}^a) \right\|_2 \leq K \left\| \boldsymbol{\lambda}^b - \boldsymbol{\lambda}^a \right\|_2$.

The Hessian $\nabla^2 F(\mathbf{M}, \boldsymbol{\lambda})$ of $F(\mathbf{M}, \boldsymbol{\lambda})$ is given by

$$\nabla^2 F_{ij,kl}(\mathbf{M}, \boldsymbol{\lambda}) = \begin{cases} \frac{2\beta \mu_i / M_i}{(\mu_i - \lambda_i / M_i)^3} & \text{if } i = k \\ 0 & \text{otherwise.} \end{cases}$$

Then $\left\| \nabla^2 F(\mathbf{M}, \boldsymbol{\lambda}) \right\|_2 \leq \left\| \nabla^2 F(\mathbf{M}, \boldsymbol{\lambda}) \right\|_1 \left\| \nabla^2 F(\mathbf{M}, \boldsymbol{\lambda}) \right\|_\infty = \left\| \nabla^2 F(\mathbf{M}, \boldsymbol{\lambda}) \right\|_\infty^2$. The inequality is a property of norms and the equality is from the symmetry of $\nabla^2 F(\mathbf{M}, \boldsymbol{\lambda})$. Finally,

$$\begin{aligned} \left\| \nabla^2 F(\mathbf{M}, \boldsymbol{\lambda}) \right\|_\infty &= \max_{ij} \left\{ \sum_{kl} \nabla^2 F_{ij,kl}(\mathbf{M}, \boldsymbol{\lambda}) \right\} \\ &= \max_i \left\{ |J| \frac{2\beta \mu_i / M_i}{(\mu_i - \lambda_i / M_i)^3} \right\} \leq |J| \max_i \frac{2(\phi + \beta M_i)^3}{\beta^2 M_i^4 \mu_i^2}. \end{aligned}$$

In the last step we substitute λ_i by $\frac{\phi M_i \mu_i}{\phi + \beta M_i}$ because $\lambda_i \leq \frac{\phi}{\phi + \beta M_i} M_i \mu_i, \forall i$ and $\frac{2\mu_i / M_i}{(\mu_i - \lambda_i / M_i)^3}$ is increasing in λ_i . \square

Lemma 11. *When applying Algorithm 2 to GLB-Q,*

(a) $F(\boldsymbol{\lambda}(\tau+1)) \leq F(\boldsymbol{\lambda}(\tau)) - \left(\frac{1}{\gamma_m} - \frac{K}{2}\right) \left\| \boldsymbol{\lambda}(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2$, where $K = |J| \max_i 2(\phi + \beta M_i)^3 / (\beta^2 M_i^4 \mu_i^2)$, $\gamma_m = \max_j \gamma_j$, and $0 < \gamma_j < \min_i \beta^2 \mu_i^2 M_i^4 / (|J|(\phi + \beta M_i)^3), \forall j$.

(b) $\boldsymbol{\lambda}(\tau+1) = \boldsymbol{\lambda}(\tau)$ if and only if $\boldsymbol{\lambda}(\tau)$ minimizes $F(\boldsymbol{\lambda})$ over the set Λ .

(c) The mapping $T(\boldsymbol{\lambda}(\tau)) = \boldsymbol{\lambda}(\tau+1)$ is continuous.

PROOF. From the Lemma 10, we know

$$\left\| \nabla F(\boldsymbol{\lambda}^b) - \nabla F(\boldsymbol{\lambda}^a) \right\|_2 \leq K \left\| \boldsymbol{\lambda}^b - \boldsymbol{\lambda}^a \right\|_2, \forall \boldsymbol{\lambda}^a \in \Lambda, \forall \boldsymbol{\lambda}^b \in \Lambda$$

where $K = |J| \max_i 2(\phi + \beta M_i)^3 / (\beta^2 M_i^4 \mu_i^2)$.

Here $\mathbf{Z}^j(\tau+1) \in \Lambda, \boldsymbol{\lambda}(\tau) \in \Lambda$, therefore we have

$$\left\| \nabla F(\mathbf{Z}^j(\tau+1)) - \nabla F(\boldsymbol{\lambda}(\tau)) \right\|_2 \leq K \left\| \mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2.$$

From the convexity of $F(\boldsymbol{\lambda})$, we have

$$\begin{aligned} F(\boldsymbol{\lambda}(\tau+1)) &= F\left(\frac{\sum_j \mathbf{Z}^j(\tau+1)}{|J|}\right) \\ &\leq \frac{1}{|J|} \sum_j F(\mathbf{Z}^j(\tau+1)) \\ &\leq \frac{1}{|J|} \sum_j \left(F(\boldsymbol{\lambda}(\tau)) - \left(\frac{1}{\gamma_j} - \frac{K}{2}\right) \left\| \mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2 \right) \\ &= F(\boldsymbol{\lambda}(\tau)) - \sum_j \left(\frac{1}{\gamma_j} - \frac{K}{2} \right) \frac{\left\| \mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2}{|J|} \\ &\leq F(\boldsymbol{\lambda}(\tau)) - \left(\frac{1}{\gamma_m} - \frac{K}{2} \right) \frac{\sum_j \left\| \mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2}{|J|} \end{aligned}$$

where $K = |J| \max_i 2(\phi + \beta M_i)^3 / (\beta^2 M_i^4 \mu_i^2)$.

The first line is from the update rule of $\boldsymbol{\lambda}(\tau)$. The second line is from the convexity of $F(\boldsymbol{\lambda})$. The third line is from the property of gradient projection. The last line is from the definition of γ_m .

Then from the convexity of $\left\| \cdot \right\|_2^2$, we have

$$\begin{aligned} \frac{\sum_j \left\| \mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2}{|J|} &\geq \left\| \frac{\sum_j (\mathbf{Z}^j(\tau+1) - \boldsymbol{\lambda}(\tau))}{|J|} \right\|_2^2 \\ &= \left\| \frac{\sum_j \mathbf{Z}^j(\tau+1)}{|J|} - \boldsymbol{\lambda}(\tau) \right\|_2^2 = \left\| \boldsymbol{\lambda}(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2. \end{aligned}$$

Therefore we have

$$F(\boldsymbol{\lambda}(\tau+1)) \leq F(\boldsymbol{\lambda}(\tau)) - \left(\frac{1}{\gamma_m} - \frac{K}{2} \right) \left\| \boldsymbol{\lambda}(\tau+1) - \boldsymbol{\lambda}(\tau) \right\|_2^2.$$

(b) $\boldsymbol{\lambda}(\tau+1) = \boldsymbol{\lambda}(\tau)$ is equivalent to $\mathbf{Z}^j(\tau+1) = \boldsymbol{\lambda}_j(\tau), \forall j$. Moreover, if $\mathbf{Z}^j(\tau+1) = \boldsymbol{\lambda}_j(\tau), \forall j$, then from the definition of each gradient projection, we know it is optimal. Conversely, if $\boldsymbol{\lambda}(\tau)$ minimizes $F(\boldsymbol{\lambda}(\tau))$ over the set Λ , then the gradient projection always projects to the original point, hence $\mathbf{Z}^j(\tau+1) = \boldsymbol{\lambda}_j(\tau), \forall j$. See also [8, Ch 3 Prop. 3.3(b)] for reference.

(c) Since $F(\boldsymbol{\lambda})$ is continuously differentiable, the gradient mapping is continuous. The projection mapping is also continuous. T is the composition of the two and is therefore continuous. \square

PROOF OF THEOREM 6. Lemma 11 is parallel to that of Proposition 3.3 in Ch 3 of [8], and Theorem 6 here is parallel to Proposition 3.4 in Ch 3 of [8]. Therefore, the proof for Proposition 3.4 immediately applies to Theorem 6. We also have $F(\boldsymbol{\lambda})$ is convex in $\boldsymbol{\lambda}$, which completes the proof. \square