

On the Benefit of Virtualization: Strategies for Flexible Server Allocation

Dushyant Arora, Anja Feldmann, Gregor Schaffrath, Stefan Schmid
Deutsche Telekom Laboratories / TU Berlin, Germany
{darora, anja, grsch, stefan}@net.t-labs.tu-berlin.de

Abstract—Virtualization technology facilitates a dynamic, demand-driven allocation and migration of servers. This paper studies how the flexibility offered by network virtualization can be used to improve Quality-of-Service parameters such as latency, while taking into account allocation costs. A generic use case is considered where both the overall demand issued for a certain service (for example, an SAP application in the cloud, or a gaming application) as well as the origins of the requests change over time (e.g., due to time zone effects or due to user mobility), and we present online and optimal offline strategies to compute the number and location of the servers implementing this service. These algorithms also allow us to study the fundamental benefits of dynamic resource allocation compared to static systems. Our simulation results confirm our expectations that the gain of flexible server allocation is particularly high in scenarios with moderate dynamics.

I. INTRODUCTION

Virtualization is an intriguing paradigm which loosens the ties between services and physical infrastructure. The gained flexibility promises faster innovations, enabling a more diverse Internet and ensuring coexistence of heterogeneous virtual network (VNet) architectures on top of the shared substrate. Moreover, the dynamic and demand driven allocation of resources may yield a “greener Internet” without sacrificing (or, in the presence of the corresponding migration technology: with improved!) quality-of-service (QoS) / quality-of-experience (QoE).

This paper studies flexible offline and online allocation strategies that achieve a good tradeoff between efficient resource usage and QoS. We attend to a rather general use case which describes a system providing a certain service to a dynamic set of users. The users’ access pattern can change over time (e.g., due to time zone effects or due to user mobility), and hence, in order to ensure a low latency, servers may adaptively be migrated closer to the current origins of the requests; moreover, in peak hours, it can be worthwhile to allocate additional resources.

This use case captures many different scenarios. For instance, imagine an SAP application in the cloud which is accessed by different users going online and offline over time, resulting in a temporal change of the demand characteristics. Or, consider a mobile provider which offers a gaming application to a set of mobile users updating their location over time, and where access latency is of prime concern.

A. Our Contributions

This paper studies strategies for flexible server allocation and migration in a generic use case capturing various scenarios, ranging from business applications such as SAP services in the cloud, to entertainment applications such as mobile gaming. We present algorithms which ensure a low access latency by adapting the resources over time, while taking into account the corresponding costs (communication cost, service interruption cost, allocation cost, migration cost, and cost of running the servers). These algorithms come in two flavors, exploring the two extremal perspectives: online algorithms where decisions are done without any information of future requests,

and offline algorithms where the (e.g., periodic) demand is known ahead of time. While optimal solutions are computationally hard, there exist efficient adaptations, which are discussed as well.

Our algorithms allow us to quantify the cost-benefit tradeoffs of dynamic resource allocation, and shed light on fundamental questions such as the use of migration compared to solutions with static resources. For example, our simulations show that the overall cost can be much higher if resources are static, in particular if the demand dynamics is moderate.

II. MODEL

Our work is motivated by the recent advances in the field of network virtualization. We are about to develop a prototype architecture, and refer the reader to [25] for more information on this project. This section first provides some basic background and subsequently identifies the major cost factors in our use case.

A. Architecture and Use Case

The virtualization architecture proposed in [25] distinguishes the following roles: The *(Physical) Infrastructure Provider (PIP)*, which owns and manages an underlying physical infrastructure (called “substrate”); the *Virtual Network Provider (VNP)*, which provides bit-pipes and end-to-end connectivity to end-users; and the *Service Provider (SP)*, which offers application, data and content services to end-users.

As a generic use case, we consider a service provider offering a service to mobile users which can benefit from the flexibility of network and service virtualization. The goal of the service provider is to minimize the round-trip-time of its service users to the servers, by triggering migrations depending, e.g., on (latency) measurements. Concretely, VNP and/or PIPs will react on the SP-side changes of the requirements on the paths between server and access points, and re-embed the servers accordingly.

B. Graph Model and Access Cost

Formally, we consider a substrate network $G = (V, E)$ managed by a substrate provider (PIP), where $v \in V$ are the substrate nodes and $e = (u, v) \in E$, with $u, v \in V$, are the substrate links; we will refer to the total number of substrate nodes by $n = |V|$. Each substrate node v has a certain strength $\omega(v)$ associated with it (number of CPU cores, memory size, bus speed, etc.). A link is characterized by a bandwidth capacity $\omega(e)$ and a latency $\lambda(e)$. In addition to the substrate network, there is a set T of external terminals (the mobile thin clients or the users) that access G by issuing requests for a given virtualized service hosted on a set of virtual servers S on G . We will assume that a service is offered redundantly by up to $k = |S|$ servers.

In order for the clients in T to access the servers S , a fixed subset of nodes $A \subseteq V$ serve as *Access Points* where clients in T can connect to G . Due to the request dynamics, the popularity of access points

can change frequently, which may trigger the migration algorithm. We define σ_t to be the multi-set of requests at time t where each element is a tuple $(a \in A, S \in \mathcal{S})$ specifying the access point and the requested service S . (For ease of notation, when clear from the context, we will sometimes simply write $v \in \sigma_t$ to denote the multi-set of access points used by the different requests.) Our main objective is to shed light onto the trade-off between the access costs Cost_{acc} of the mobile users to the service (delay of requests), the server migration cost Cost_{mig} , and the cost Cost_{run} for running the servers: While moving the servers closer to the requester may reduce the access costs and hence improve the quality of service, it also entails the overhead of migration; moreover, the more active servers, the more resources needed (processing requests, CPU, storage, etc.).

In this paper, a simplified model is considered where the cost of accessing the server, Cost_{acc} , is given by the request latency, i.e., the sum of the requests' latencies to the corresponding servers (e.g., along the shortest paths on the substrate network), plus the latency due to the server's load, which, for server v and at time t , is given by $\text{load}(v, t) = f(\omega(v), \eta(v, t))$, a function of the node strength $\omega(v)$ and the number of requests arriving at the servers hosted by v at time t , $\eta(v, t)$. For example, a simple model where the load increases linearly would be $\text{load}(v, t) = \eta(v, t)/\omega(v)$:

$$\text{Cost}_{\text{acc}}(t) = \sum_{r_t \in \sigma_t} \text{delay}(r_t) + \sum_{v \in V} \text{load}(v, t).$$

We will assume that requests are routed to the server of minimal access costs.

C. Server Model and Migration

Each of the (at most k) servers can assume three different states: *not in use*, *inactive*, and *active*. If a server is not in use, there are no costs. An inactive server comes at a certain cost R_i per time: this cost includes storing the application software (e.g., the game) plus certain maintenance costs. The running costs of an active server R_a are larger, as they also include CPU costs, maintaining state in the RAM, or bandwidth costs. In order to startup a server which is not in use, a fixed creation cost c is assumed. For instance, this cost captures the installation of the Linux box and the template (copy if already on disk or download from an NFS share), configuration of the template (e.g., setting up IP addresses manually or via DHCP), starting the server etc. Finally, we assume that the cost of changing from inactive to active state is negligible.

Also the cost of migration depends on many different factors. While the operating systems is typically subject to replication (copy Linux box from disk), the virtual server's configuration and data/state component must be transmitted over the network. Besides the cost for the bulk data transfer of the server state, there are opportunistic costs that depend on whether the system supports live migration or not, possibly some requests need to be routed to other servers during migration, there can be periods of service interruption, etc. We will consider a simplified scenario where migration cost is described by a constant β and where inactive servers are not migrated.

How does β relate to c ? It again depends on the scenario. For example, $c \gg \beta$ for systems which support live migration (almost no opportunistic or outage costs during migration), where there is an NFS share and only the server state is migrated, and where new servers need to be configured manually. On the other hand, for example $c \ll \beta$ in systems where configuration is simple and where migration happens over multiple provider domains. For the formal description and analysis of the algorithms we will focus on the more interesting case that $\beta < c$: if $\beta \geq c$, migration is never beneficial,

and the problem boils down on when and where to create and delete servers; our algorithms could be easily adapted for these situations as well.

Also note that in our model, migrating a server from node v to an empty node v' costs β and that subsequently, node v is empty. It is not possible to maintain a (for example inactive) copy of the server at v "for free"; rather, this would require to set up a new server which costs c , as the template needs reconfiguration (e.g., new unique network addresses are needed).

However, let us emphasize that our approach is general enough to be adopted in many alternative scenarios where the cost models are slightly different, e.g., where server copies can be kept during migration and c denotes the cost of investing into an additional server.

In order to clarify our model, we give three examples.

Example 1: Assume we have three active servers located at nodes v_1, v_2, v_3 . When adding an additional server at some node v_4 , either: (1) if there is no inactive server, an additional server has to be created at v_4 which costs c (we assume $\beta < c$); (2) if node v_4 already hosts an inactive server, this server is simply activated and the cost is zero; (3) otherwise, if there is an inactive server at some other location v_5 , this inactive server is migrated to v_4 which costs β (note that there will be no server at v_5 anymore).

Example 2: Assume we have three servers hosted by nodes v_1, v_2, v_3 . In order to change to a configuration where three servers are located at nodes v_1, v_2 , and v_4 , either (1) if node v_4 is an inactive node, the cost is zero; (2) if an inactive server at v_5 is migrated to v_4 , the cost is β , and the server from v_3 can become inactive (in our algorithm, it will be added to a cache of inactive servers)—there will be no server at v_5 anymore; (3) if the active server at v_3 is migrated to v_4 , the cost is β as well and there is no server at v_3 anymore.

Example 3: Assume we have three servers hosted by nodes v_1, v_2, v_3 . When removing one server, i.e., when changing to a configuration with two servers at nodes v_1 and v_3 , we do not incur any costs and the server at v_2 becomes inactive (i.e., is added to the cache of inactive servers in our algorithms).

D. Request Model

We next discuss the model for the terminal dynamics (e.g., due to user mobility). One approach could be to assume arbitrary request sets σ_t , where σ_t is completely independent of σ_{t-1} . However, for certain applications it may be more realistic to assume that the requests move "slowly" between the access points. We can distinguish between two different sources of request dynamics: *time zone effects* (users from different countries access a service at different times of the day) and *user mobility*. Note that while users typically travel between different cities or countries at a limited speed, these geographical movements may not translate to the topology of the substrate network. Thus, rather than modeling the users to travel along the links of G , we may consider on/off models where a user appears at some access point $a_1 \in A$ at time t , remains there for a certain period Δt , before moving to another arbitrary node $a_2 \in A$ at time $t + \Delta t$. Often, it is reasonable to assume some form of correlation between the individual users' movement. For example, in an urban area, workers commute downtown in the morning and return to suburbs in the evening.

E. Online and Offline Algorithms

Typically, resources need to be allocated dynamically (i.e., *online*) in virtual networks, without knowledge on the demand or request sequence in advance. In order to focus on the main properties and trade-offs involved in the the dynamic allocation and migration problem, we assume a simplified online framework (see also [4]).

We assume a synchronized setting where time proceeds in time slots (or *rounds*). In each round t , a set of σ_t terminal requests arrive in a worst-case and online fashion at an arbitrary set of access nodes A . Thus the embedding problem is equivalent to the following synchronous game, where an online algorithm ALG has to decide on the server allocation and migration strategy in each round t , without knowing about the future access requests. Concretely, in each round $t \geq 0$:

1. The requests σ_t arrive at some access nodes A .
2. The online algorithm ALG pays the requests' access costs $\text{Cost}_{\text{acc}}(t)$ to the corresponding servers.
3. The online algorithm ALG decides where in G to allocate new or remove existing servers, which servers should be active and which inactive, and where to migrate the servers S . Accordingly it incurs running costs $\text{Cost}_{\text{run}}(t)$ as well as migration costs $\text{Cost}_{\text{mig}}(t)$.

Observe that as we assume that the requests during one time slot are much cheaper than a migration operation, our results also apply for a scenario where the last two steps are reordered (see also [4]).

To evaluate the efficiency of an online algorithm, its performance is often compared to the performance of a (sometimes hypothetical) optimal offline algorithm for the given request sequence. The ratio of the two costs is called the *competitive ratio*.

III. ONLINE ALGORITHMS

This section presents strategies to allocate and migrate resources in an online fashion—without knowing the future request pattern—depending on the observed origins of the requests and the load. A natural idea is to pursue a server configuration approach which gives a general class of online algorithms.

Definition 3.1 (Configuration): A configuration γ describes, for each server, whether it is *not in use*, *inactive*, or *active*. In case of inactive and active servers, γ specifies where—i.e., on which node—the server is located.

In some sense, the single server algorithm proposed in [4] can be regarded as a special case of this idea. Generalizing the algorithm of [4] gives the following algorithm ONCONF:

ONCONF uses a counter $C(\gamma)$ for each configuration γ . Time is divided into *epochs*. In each epoch ONCONF monitors, for each configuration γ , the cost of serving all requests from this epoch by servers kept in configuration γ , including the access costs (latency plus induced load) of the requests, the server running costs, and possible creation costs. ONCONF stores this cost in $C(\gamma)$. The servers are kept in a given configuration $\hat{\gamma}$ until $C(\hat{\gamma})$ reaches $k \cdot c$. In this case, ONCONF changes to a configuration $\hat{\gamma}'$ chosen uniformly at random among configurations with the property $C(\gamma) < k \cdot c$. If there is no such configuration left, we do not migrate and the epoch ends in that round; the next epoch starts in the next round and the counters $C(\gamma)$ are reset to zero.

We can implement ONCONF in such a way that inactive servers are managed by a queue of constant size. Inactive servers in the queue are managed in a FIFO manner (older servers are replaced first); in addition an inactive server expires after x epochs for some constant parameter x . (This also means that in configurations γ in ONCONF, inactive servers are not included.)

Observe that during an epoch, ONCONF goes through at most $O(k \log n)$ many configuration changes (or epochs), as there are $\sum_{i=1}^k \binom{n}{i}$ many configurations. The cost per epoch is at most $k \cdot c$,

and it is easy to see that—under an overly pessimistic perspective—also an optimal offline algorithm has cost at least β per k epochs, so in competitive analysis parlance [4], we have a competitive ratio of at most $O(c/\beta \cdot k^3 \log n)$.

Clearly, ONCONF needs to be optimized in many respects. For instance, it can make sense to switch between “close” (with respect to costs) configurations only, or to deterministically switch to the configuration with the lowest counter. However, the main problem of ONCONF is different: due to the configuration complexity, the runtime is only acceptable for a small number of servers k . Therefore, rather than discussing possible optimizations, we concentrate on efficient variants for ONCONF.

A. Efficient Online Algorithms

There are several ways to speed up ONCONF such as clustering approaches where optimal configurations are only considered on a cluster granularity, or sampling approaches where, e.g., only k configurations are tracked, one for each possible number of current servers. In the following, we will focus on a sequential best-response variant of ONCONF called ONBR: updating one server after another greatly reduces the configuration complexity, while flexibility is maintained.

ONBR starts in an arbitrary configuration, e.g., hosting one server at the network center. Time is divided into epochs, and an epoch ends when the total cost accumulated during this epoch (including access cost and running cost) reaches a threshold θ . Then, ONBR changes to the cheapest (w.r.t. the passed epoch and including access, migration, running, and creation cost) configuration among the following: (1) γ (no change), (2) γ but where one server s is migrated to a different location ($O(n)$ options), (3) γ but where one server s becomes inactive ($O(k)$ options), (4) γ but where one inactive server s becomes active, or a new active server s is created ($O(n)$ options). Inactive servers are organized in a queue of constant size where the oldest server in the queue is the first to be replaced—i.e., this server will no longer be in use (in our simulations: size 3). In case a new server is created at an empty node, the oldest inactive server from the queue is migrated to the corresponding node. Inactive servers in the queue expire after x epochs for some parameter x ($x = 20$ in our simulation).

ONBR requires a good choice of the parameter θ in order to trade-off flexibility and optimality. An alternative intuitive approach is to add new servers when the access cost is higher than the total running cost (w.r.t. a certain threshold). This is automatized by the following algorithm ONTH.

ONTH starts in an arbitrary configuration, e.g., hosting one server at the network center. Time is divided into small and large epochs: a small epoch ends when we have accumulated a cost of $y \cdot \beta$ in a given configuration for some constant parameter y ($y = 2$ in our simulations), and a large epoch ends when the accumulated access cost is larger than the accumulated running cost (of the active servers); concretely, we will use the following condition: $\text{Cost}_{\text{acc}}/(k_{\text{cur}}+1) - \text{Cost}_{\text{run}} > c$, where k_{cur} denotes the current number of active servers. When a small epoch ends ONTH changes to the cheapest (w.r.t. the passed epoch and including access, migration, and running cost) configuration among the following: (1) γ (no change), (2) γ but where one server s is migrated to a different location (the

server at the migration origin becomes inactive), (3) γ but where one server s becomes inactive (at most k options). Inactive servers are organized in a first-in-first-out (FIFO) queue of constant size (in our simulations: size 3), i.e., inactive servers which fall out of the queue are no longer in use. Inactive servers in the queue expire after $y \cdot \beta$ rounds for some parameter x ($x = 20$ in our simulation). When a large epoch ends, a new server is activated at an optimal position with respect to the access cost of the latest large epoch.

Note that both ONBR and ONTH have the appealing property that in case of constant demand, they will eventually converge to a stable configuration.

IV. OFFLINE ALGORITHMS

While in the worst case, the decisions when and where to migrate servers typically needs to be done *online*, i.e., without the knowledge of future requests, there can be situations where it is interesting to study which migration pattern would have been good *at hindsight*. For example, if it is known that the requests follow a regular pattern (e.g., a periodic pattern per day or week), it can make sense to compute the migration strategy offline and apply it in the future. While Section III assumed an extreme standpoint and only discussed algorithms that do not have any knowledge of future requests, in the following, the other extreme standpoint is explored where future demand is completely known. Please note that there is also another reason why designing offline algorithms explicitly may be of interest, namely to compute competitive ratios (see [4]) in simulations.

A. Optimal Offline Algorithm

This section presents an optimal offline algorithm OPT for our resource allocation optimization problem. It turns out that offline strategies can be computed for many different scenarios, and we describe a very general algorithm here.

Algorithm OPT is based on dynamic programming techniques and also uses the concept of configurations (cf Definition 3.1). Recall that given a configuration γ , access costs Cost_{acc} , migration costs Cost_{mig} , and the running costs Cost_{run} over time can be computed.

OPT exploits the fact that the migration problem exhibits an optimal substructure property: Given that at time t , the k servers are in a configuration γ , then the most cost-efficient *path* (migrations, activation and deactivation of servers, creation, etc.) that leads to this configuration consists solely of optimal sub-paths. That is, if a cost minimizing path to configuration γ at time t leads over a configuration γ' at time $t' < t$, then there cannot be a cheaper migration sub-path that leads to γ' at time t' than the corresponding sub-path.

OPT essentially fills out a matrix $\text{opt}[\text{time}][\text{configuration}]$ where $\text{opt}[t][\gamma]$ contains the cost of the minimal path that leads to a configuration where the servers satisfy the requests of time t in a configuration γ . Recall from Definition 3.1 that a configuration γ describes for each virtual server s at which physical node v it is hosted and whether s is *not in use*, *inactive*, or *active*.

Assume that in the beginning, the system is located in configuration γ_0 . Thus, initially, $\text{opt}[0][\gamma] = \text{Cost}(\gamma_0 \rightarrow \gamma) + \text{Cost}_{\text{run}}(\gamma) + \left[\sum_{v \in \sigma_0} \text{Cost}_{\text{acc}}(v, \gamma) \right]$, where $\text{Cost}(\gamma_1 \rightarrow \gamma_2)$ denotes the cost of changing from configuration γ_1 to γ_2 (cost of migrations, creation costs, etc.), $\text{Cost}_{\text{run}}(\gamma)$ denotes the cost of running the inactive and active servers for one time unit in configuration γ , and $\sum_{v \in \sigma_0} \text{Cost}_{\text{acc}}(v, \gamma)$ denotes the access costs (request latency and server load) resulting from the requests of σ_0 accessing the active

servers in configuration γ . (W.l.o.g., we assume that the cost Cost_{acc} contains the first wireless hop from terminal to substrate network.)

For $t > 0$, we find the optimal values $\text{opt}[t][\gamma]$ by considering the optimal paths to any configuration γ' at time $t - 1$, and adding the migration cost from γ' to γ . That is, in order to find the optimal cost to arrive at a configuration with servers at γ at time t :

$$\min_{\gamma'} \left[\text{opt}[t-1][\gamma'] + \text{Cost}(\gamma' \rightarrow \gamma) + \text{Cost}_{\text{run}}(\gamma) + \sum_{v \in \sigma_t} \text{Cost}_{\text{acc}}(v, \gamma) \right]$$

where we assume that Cost_{acc} includes the first (wireless) hop of the request from the terminal to the substrate network, and where $\text{Cost}(\gamma \rightarrow \gamma) = 0 \quad \forall \gamma$.

Note that OPT is not an online algorithm, although $\text{opt}[t]$ does not depend on future requests: in order to reconstruct the optimal migration strategy at hindsight, the configuration of minimal cost after the last request is determined, and from there, the optimal path is given by recursively finding the optimal configuration at time $t - 1$ which led to the optimal configuration at time t .

B. Efficient Offline Algorithms

Unfortunately, the computational complexity of OPT is rather high for scenarios with many servers. We believe that it is difficult to significantly reduce the runtime. Again, clustering or sampling heuristics may be used to speed up the computations (which may come at a loss of allocation quality).

There is an interesting and natural adaption of the best response strategies of Section III: OFFBR is similar to ONBR, but rather than switching to the configuration of lowest cost w.r.t. the passed epoch, we switch to the configuration of lowest cost in the *upcoming epoch*! A similar transformation can be done from ONTH to OFFTH: we simply compute optimal strategies of small epochs at hindsight. This yields an acceptable runtime.

V. SIMULATIONS

We conducted several experiments and in the following, we report on our simulation results and main insights.

A. Set-up

We conducted experiments on both artificial *Erdős-Rényi graphs* random graphs (with connection probability 1%) as well as more realistic graphs taken from the *Rocketfuel project*¹ [26], [27] (including the corresponding latencies for the access cost). To simulate OPT, we constrain ourselves to line graphs. If not stated otherwise, we assume that link bandwidths are chosen at random (either T1 (1.544 Mbit/s) or T2 (6.312 Mbit/s)), that $\beta = 40$ and that $c = 400$; for experiments with $\beta > c$, we set $\beta = 400$ and $c = 40$.

As real traffic patterns are confidential and cannot be published, we consider two different simplified, artificial scenarios, reflecting the two main reasons for request dynamics: *time zone effects* (users from different countries access the service at different times of the day) and *user mobility*.

¹For maps and data, see <http://www.cs.washington.edu/research/networking/rocketfuel/>.

Time Zones Scenario: This scenario models an access pattern that can result from global daytime effects. We divide a day into T time periods. For each time t , $p\%$ of all requests originate from a node chosen uniformly at random from the substrate network (we assume that these locations are the same each day). (Recall that the substrate topology does not necessarily reflect the geographic situation, but note that the uniform random choice is still pessimistic.) The sojourn time of the requests at a given location is constant and given by a parameter τ . In addition, there is a background traffic: The remaining requests originate from nodes chosen uniformly at random from all access points.

Commuter Scenario: This scenario models an access pattern that can result from commuters traveling downtown for work in the morning and returning back to the suburbs in the evening. The scenario comes in two flavors: one with *static demand* and one with *dynamic demand*.

- 1) *Static Load:* We use a parameter T to model the *frequency* of the changes. At time $t \bmod T < T/2$, there are $2^{t \bmod T}$ requests originating from access points chosen uniformly at random around the center of the network. In the second half of the day, i.e., for $t \in [T/2, \dots, T]$, the pattern is reversed. Then a new day starts. The total number of requests per round is fixed to $2^{T/2}$. At time $t_i < T/2$, the requests originate from $p = 2^{t_i \bmod T}$ of all access points including the network center ($2^{T/2}/p$ requests per access point), until single requests originate from $2^{T/2}$ access points. Then, the same process is reversed until all $2^{T/2}$ requests originate from a single access point: the network center. We assume that the time period between t_i and t_{i+1} is given by a fixed parameter λ .
- 2) *Dynamic Load:* The total number of requests per round is not fixed to $2^{T/2}$. At time $t_i < T/2$, the requests originate from $p = 2^{t_i \bmod T}$ of all access points including the network center (one request per access point), until single requests originate from $2^{T/2}$ access points. Then, the same process is reversed until we have a single request originating from a single access point: the network center. We assume that the time period between t_i and t_{i+1} is fixed and we denote it by parameter λ .

B. Experiments

The main objective of our algorithms is to adapt to dynamically changing demands in an efficient manner. Where to allocate or migrate how many servers depends on the origins and the size of the requests, and also on how access cost increases as a function of load. As a motivation, consider the exemplary executions of ONTH depicted in Figure 1 for linear and quadratic load functions. It can be seen intuitively that ONTH reacts as desired in this example, reacting to higher loads (either due to higher demand or steeper load functions) by allocating more servers. Figure 2 shows the same execution for the static load variant. Initially, the system converges quickly to a certain number of servers. It can also be seen that the number of servers in this scenario is more or less independent of the number of access points from which the given requests originate from, and that a quadratic load model requires more servers.

In the following, after these motivating examples, the performance of our algorithms is studied more systematically by considering the dependencies of the performance on the various parameters. A first set

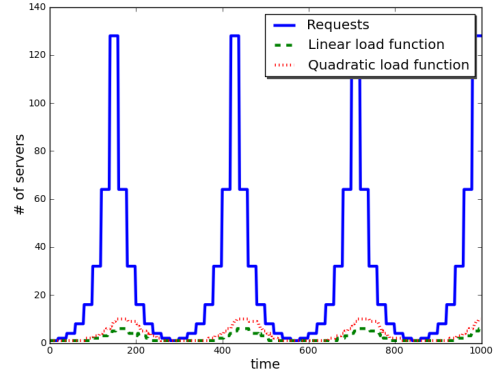


Fig. 1. Exemplary execution of ONTH in commuter scenario with dynamic load. In this setting, the runtime was 1000 rounds, $T = 14$, we considered a network of size 1000, and set $\lambda = 20$.

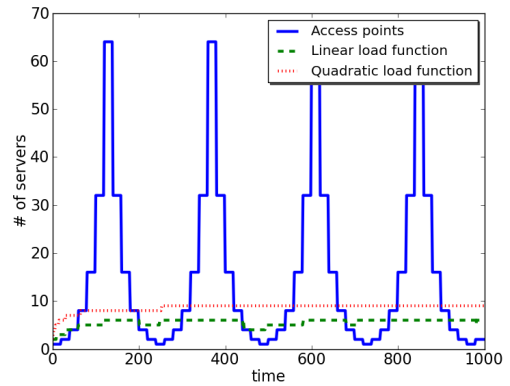


Fig. 2. Exemplary execution of ONTH in commuter scenario with static load. In this setting, the runtime was 1000 rounds, $T = 12$, we considered a network of size 500, and set $\lambda = 20$.

of experiments compares the performance of ONBR and ONTH. For ONBR, a threshold $2 \cdot c$ is considered. In addition, we experimented with a variant where the threshold also depends on the length ℓ of the preceding epoch, i.e., $2 \cdot c / \ell$ —in some sense, a shorter epoch denotes faster changes in the requests, and hence, the system should adapt more quickly. We will refer to the two variants of ONBR by “*fixed*” and “*dyn*”. Clearly, many other variants are possible, e.g., where the threshold depends on the variance (over time) of the access cost. However, since ONTH typically outperformed ONBR and as ONTH requires less parameter, we do not discuss the ONBR variants in more detail but will focus on ONTH. Moreover, note that for $\beta > c$, migration is never useful, and the three algorithms coincide; in this case, we will simply consider ONBR with fixed threshold $2c$ in this case.

Figure 3 compares the cost of the different algorithms in the commuter scenario with dynamic load and as a function of network size. We can see that ONTH has lower costs than the ONBR variants, although the cost increase slightly faster with the number of nodes in the substrate network. Figure 4 and Figure 5 show the same results for static load commuter scenario and the time zone scenario, respectively.

Of course, there are different costs involved in the different

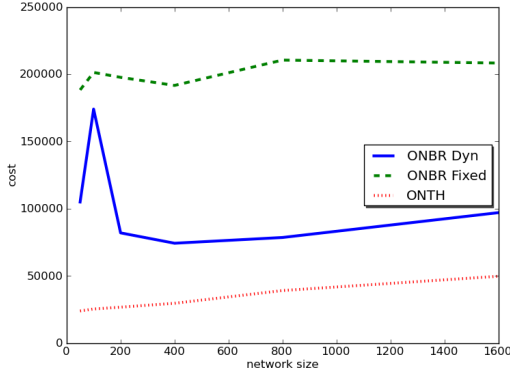


Fig. 3. Cost of different algorithms in commuter scenario with dynamic load as a function of network size. The runtime was 500 rounds, $\lambda = 10$, and we averaged over 5 runs. Note that T increases with network size in our model.

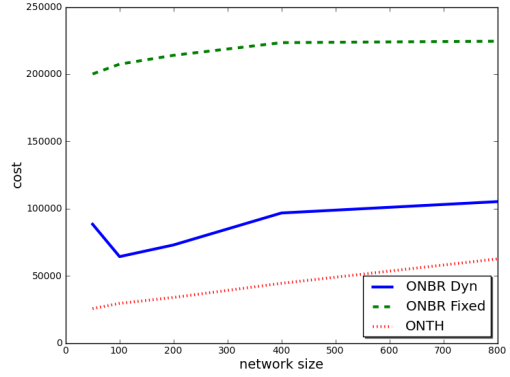


Fig. 5. Like Figure 3, but for time zone scenario.

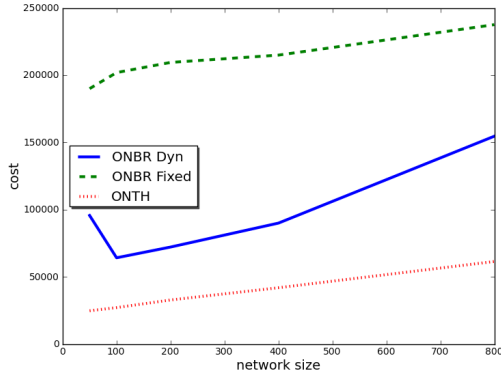


Fig. 4. Like Figure 3, but with static load.

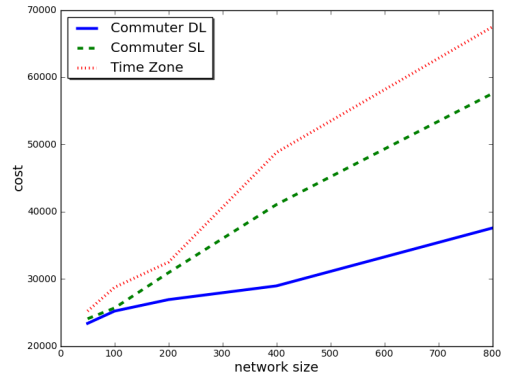


Fig. 6. Comparison of costs incurred by ONBR in different scenarios as a function of network size (runtime 500 rounds, $\lambda = 10$, $\beta = 400$, $c = 40$, and averaged over 5 runs.)

scenarios. Figure 6 shows how these costs relate to each other in the case of ONBR in a scenario where $\beta > c$. (Recall that for $\beta > c$, the three algorithms coincide, and we simply consider ONBR with fixed threshold $2c$ in this case.)

Figure 7 studies the cost of the different algorithms as a function of T in the commuter scenario with static load. Also here, ONTH always yields the best performance. (Note that the cost slightly increases with T , which is due to the fact that the request horizon is larger for larger T in our scenario.)

The total cost is more or less independent of λ , as shown in Figure 8, while ONTH is better by a factor of approximately two. Figure 9 presents the same results for a static load scenario, and Figure 10 presents the time zone scenario. In the latter, the total cost decreases slightly with λ , which is due to the fact that less migrations are needed for larger λ .

Our algorithms also allow us to get a glimpse² onto the price of online decision making, and more importantly, the question of when dynamic allocation and migration technology is most useful. Regarding the price of the lack of knowledge of future requests, Figure 11 shows the competitive ratio of ONTH as a function of λ . While the ratios are fairly low in all scenarios, the static load

²Unfortunately, the networks for which we can run these experiments are small.

commuter scenario reaches the highest peak for an intermediate λ .

In order to shed light onto the benefits of migration, as a reference point, we use the following static (but offline) algorithm OFFSTAT which we will compare to an optimal offline algorithm with migration.

For a given request sequence σ , OFFSTAT determines the optimal number of servers k_{opt} as follows. For each $i \in \{1, \dots, k\}$, we compute the cost of the following greedy static configuration for σ : one active server $j \in \{1, \dots, i\}$ after the other is placed greedily at the location which yields the lowest cost for σ , given the already placed servers $\{1, \dots, j-1\}$. k_{opt} is defined as the i with minimal cost. Figure 12 illustrates how OFFSTAT computes the number of servers.

Figure 13 shows the absolute costs of OFFSTAT and OPT as a function of λ in the commuter scenario with dynamic load. As expected, in less dynamic systems, the cost goes down, and the relative advantage of the allocation and migration flexibility declines. Figure 14 shows the same result for the $\beta > c$ scenario.

In contrast to the absolute costs, the relative costs capture the use of dynamic allocation more directly. Figure 15 plots the ratio of the total cost incurred by OFFSTAT by the total cost incurred by OPT

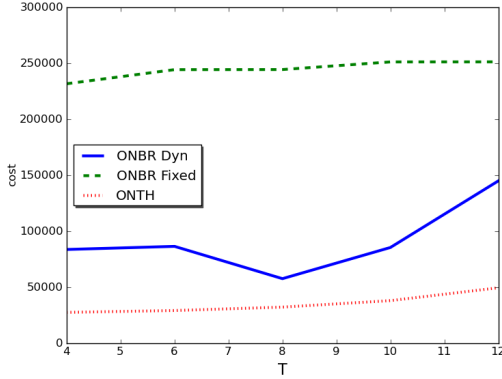


Fig. 7. Cost as a function of T for different strategies in a commuter scenario with static load (runtime 600, $\lambda = 20$, network size 1000, averaged over 10 runs.)

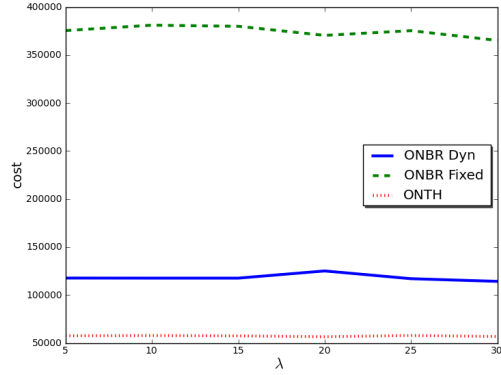


Fig. 9. Cost as a function of λ in commuter scenario with static load (runtime 900 rounds, $T = 10$, network size 200, averaged over 10 runs.)

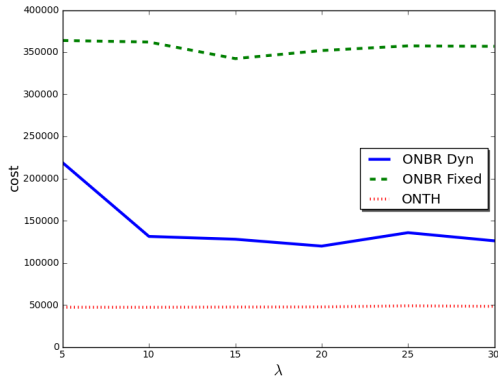


Fig. 8. Cost as a function of λ in commuter scenario with dynamic load (runtime 900 rounds, $T = 10$, network size 200, averaged over 10 runs.)

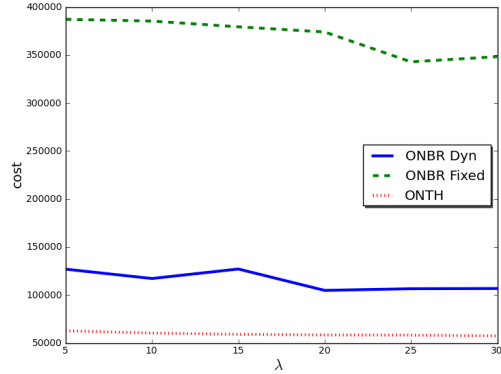


Fig. 10. Cost as a function of λ in time zone scenario with $p = 50\%$ (runtime 900 rounds, $T = 10$, network size 200, averaged over 10 runs.)

in the dynamic load commuter scenario as a function of λ . As can be seen, for very high dynamics as well as for very low dynamics, the flexibility of OPT is of limited benefit. However, for moderate dynamics, it is worthwhile for OPT to exploit the request patterns, and a better performance can be achieved (up to a factor of two). This result also meets our expectation. Interestingly, however, it turns out that OPT is relatively better if $\beta > c$, i.e., in scenarios where migration is never an option.

Also in a commuter scenario with static load (Figure 16), $\beta < c$ yields the lower ratio, fluctuating more or less constantly around 1.2 until it goes down to one for static access patterns. For $\beta > c$, the ratio goes up to almost two for intermediate λ values.

The dependency of the ratio between OFFSTAT and OPT costs on λ is more accentuated in the time zone scenario. Figure 17 shows that while a very high dynamic yields a moderate ratio, the ratio goes up quickly already for small λ , and then the use of migration declines more or less linearly with lower dynamics. Interestingly, the two variants $\beta < c$ and $\beta > c$ yield similar results in this time zone scenario. These results can be explained as follows: for the commuter scenarios, new servers have to be created rapidly as the requests fan out; migration does not help much, and the load on the servers plays a minor role. Thus, for $c < \beta$ we observe a better ratio between OFFSTAT and OPT (smaller cost c). For time zone on

the other hand, the requests move highly correlated, and there is not much of a difference whether new servers are created or existing servers are migrated. This explains why the ratio for $\beta < c$ is similar to the one for $\beta > c$.

We also conducted experiments studying the impact of T . Recall that for larger T , the request horizon becomes larger, and hence, we expect higher absolute costs as well as a higher benefit from migration. Figure 18 shows the ratio for a dynamic load commuter scenario, and Figure 19 shows the corresponding results in a static load scenario. The two experiments confirm our expectations, and also indicate that the $\beta > c$ variants typically benefit more from migration.

Finally, we briefly report on the results we obtained in the Rocketfuel network AS-7018 of ATT under the time zone scenario. ($c = 400$, $\beta = 40$, $R_a = 2.5$, $R_i = 0.5$, runtime 600 rounds, $\lambda = 20$, $p = 50\%$): the total cost of OFFSTAT was 26063.8129053. ONTH was a factor less than two higher (cost 44176.288923) while ONBR had costs 111470.296256.

VI. RELATED WORK

Our work is related to past and ongoing research in several fields. In the following, the literature in these fields is reviewed and discussed individually.

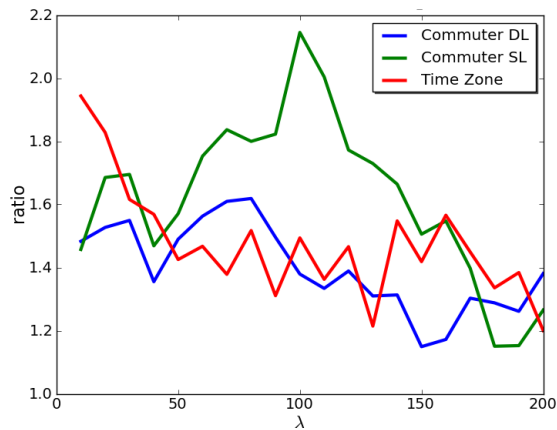


Fig. 11. Ratio of ONTH cost divided by OPT cost as a function of λ , runtime 200 rounds, in a network with five nodes, averaged over 10 runs.

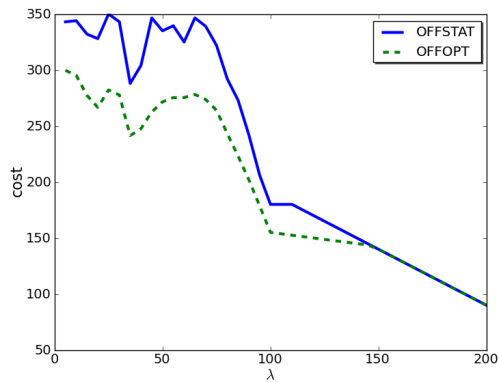


Fig. 13. The use of dynamic allocation in commuter scenario with dynamic load as a function of λ . The experiment ran for 200 rounds in a network of five nodes where $T = 4$, averaged over 10 runs.

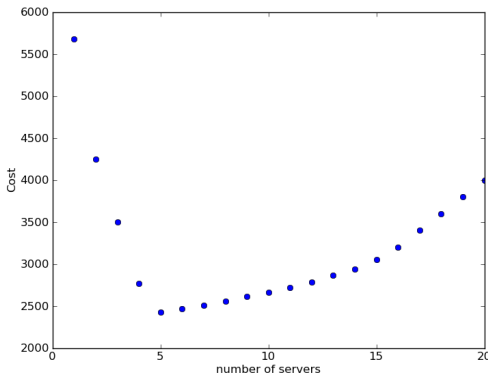


Fig. 12. OFFSTAT determines the best number of servers by minimizing the total cost.

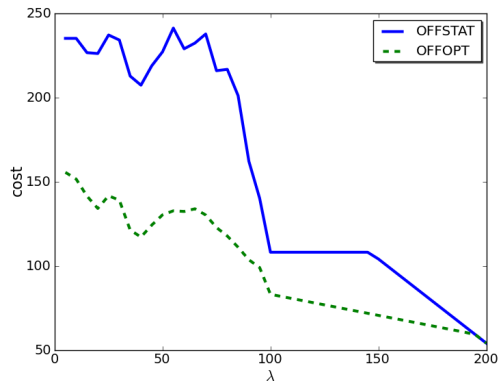


Fig. 14. Costs in dynamic load commuter scenario where $\beta = 400$ and $c = 40$ as a function of λ (runtime 200 rounds, network size five, $T = 4$, averaged over 10 runs).

Network Virtualization: It is expected that in the future, virtual networks will be allocated, maintained and managed like clouds, offering flexibility and elasticity of resources allocated for a limited time and driven by the demand; indeed, the algorithms in this paper are of applicable in cloud environments. Network virtualization has gained attention [28] because it enables the co-existence of innovation and reliability [25] and promises to overcome the “ossification” of the Internet [10]. For a more detailed survey on the subject, please refer to [9]. Virtualization allows to support a variety of network architectures and services over a shared substrate, that is, a *Substrate Network Provider (SNP)* provides a common substrate supporting a number of *Diversified Virtual Networks (DVN)*. OpenFlow [20] and VINI [2] are two examples that allow researchers to (simultaneously) evaluate protocols in a controllable and realistic environment. Trelis [3] provides such a software platform for hosting multiple virtual networks on shared commodity hardware and can be used for VINI. Network virtualization is also useful in data center architectures, see, e.g., [13].

Embedding: One major challenge in this context is the *embedding* [21] of VNets, that is, the question of how to efficiently and on-demand assign incoming service requests onto the topology. Due to its relevance, the embedding problem has been intensively studied in various settings, e.g., for an offline version of the embedding problem

see [18], for an embedding with only bandwidth constraints see [11], for heuristic approaches without admission control see [30], or for a simulated annealing approach see [24]. Since the general embedding problem is computationally hard, Yu et al. [19] advocate to rethink the design of the substrate network to simplify the embedding; for instance, they allow to split a virtual link over multiple paths and perform periodic path migrations. Lischka and Karl [17] present an embedding heuristic that uses backtracking and aims at embedding nodes and links concurrently for improved resource utilization. Such a concurrent mapping approach is also proposed in [8] with the help of a mixed integer program. Finally, several challenges of embeddings in wireless networks have been identified by Park and Kim [22].

Migration and Allocation: In contrast to the approaches discussed above we, in this paper, tackle the question of how to dynamically embed or migrate virtual servers [23] in order to efficiently satisfy connection requests arriving online at any of the network entry points, and thus use virtualization technology to improve the quality of service for mobile users. The relevance of this subproblem of the general embedding problem is underlined by Hao et al. [14] who show that under certain circumstances, migration of a Samba front-end server closer to the users can be beneficial even for bulk-data applications. Our paper builds upon the single server architecture

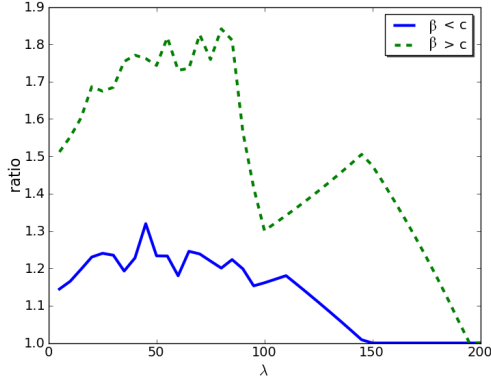


Fig. 15. Ratio of OFFSTAT and OPT costs in dynamic load commuter scenario as a function of λ , where runtime was 200 rounds, $T = 4$, network size 5, and averaged over 10 runs.

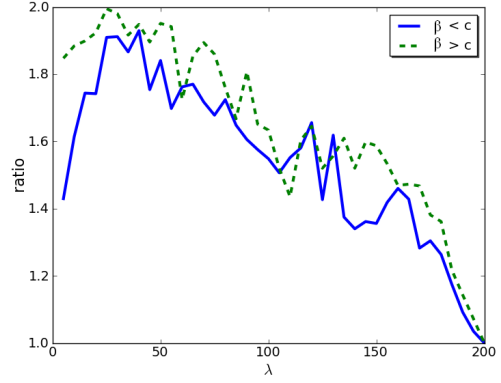


Fig. 17. Ratio between OFFSTAT and OPT cost as a function of λ in time zone scenario ($p = 50\%$). Runtime 200 rounds, three requests per round, network size five, and averaged over ten runs.

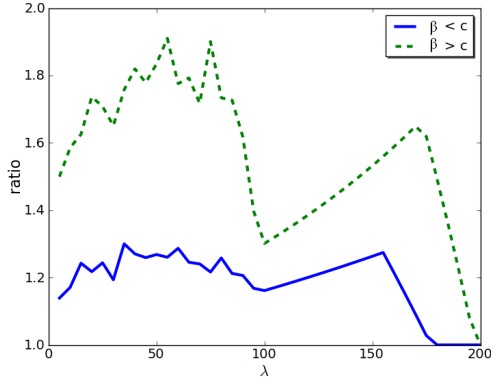


Fig. 16. Ratio of OFFSTAT and OPT costs in static load commuter scenario as a function of λ , where runtime was 200 rounds, $T = 4$, network size five, and averaged over 10 runs.

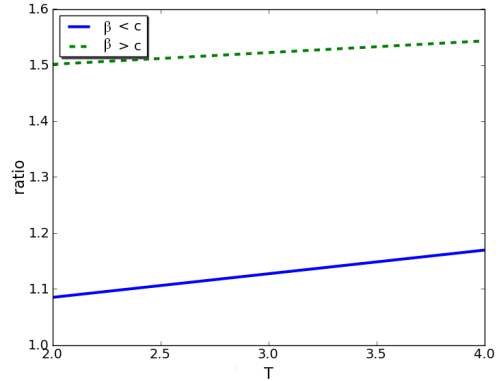


Fig. 18. Commuter scenario with dynamic load, where ratio of OFFSTAT and OPT costs are plotted as a function of T . Runtime 200 rounds, $\lambda = 10$, network size five, averaged over ten runs.

studied in [4], where a competitive online algorithm has been described to migrate a single server depending on the dynamics of mobile users. In contrast to [4] which attends to the question of *where* to migrate a server, we, in this paper, initiate the study of when to allocate *additional servers*; moreover, we generalize the migration model of [4] by taking into account the running costs of a server. Also, in contrast to [4] our model is aware of the load induced by the traffic arriving at a node, and we present solutions for very general load functions; for example, this implies that for scenarios where the cost highly depends on the number of requests, relatively more servers will be allocated to balance the cost. Note that how to dynamically allocate and expand/release resources has been studied in different contexts before as well. For instance, [15] describe a distributed virtual resource provisioning and embedding algorithm based on an autonomous agent framework, with a focus on fault-tolerance.

Online Algorithms: Due to the dynamic nature of virtual networks, the field of online algorithms and competitive analysis offers many tools that are useful to design strategies with performance guarantees under uncertainty of the request pattern. For instance, in the field of *facility location*, researchers aim at computing optimal facility locations that minimize building costs and access costs (see, e.g., [12] for an online algorithm). As in our model, the uncapacitated

online facility location problem allows to create (but not shut down again!) facilities when needed, and hence the set of “resources” is dynamic. In contrast to the classic facility location problems, in our model an additional server does not only come at a certain creation cost, but also entails running costs; moreover, our model incorporates a notion of mobility of requests, and servers can be migrated and shut down again. In [16], a heuristic algorithm is proposed for a variant of a facility location problem which allows for facility migration; this algorithm uses neighborhood-limited topology and demand information to compute optimal facility locations in a distributed manner. In contrast to our work, the setting is different and migration cost is measured in terms of hop count. In the field of *k-server problems* (e.g., [6]), an online algorithm must control the movement of a set of k servers, represented as points in a metric space, and handle requests that are also in the form of points in the space. As each request arrives, the algorithm must determine which server to move to the requested point. The goal of the algorithm is to reduce the total distance that all servers traverse. In contrast, in our model the server can be accessed remotely, that is, there is no need for the server to move to the request’s position. The *page migration problem* (e.g., [1]) occurs in managing a globally addressed shared memory in a multiprocessor system. Each physical page of

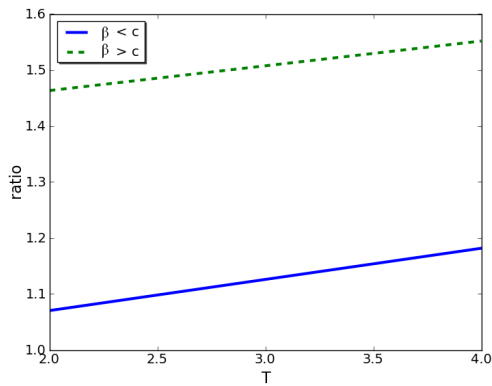


Fig. 19. As Figure 18 but for static load.

memory is located at a given processor, and memory references to that page by other processors are charged a cost equal to the network distance. At times, the page may migrate between processors, at a cost equal to the distance times a page size factor. The problem is to schedule movements on-line so as to minimize the total cost of memory references. In contrast to these page migration models, we differentiate between access costs that are determined by latency and migration costs that are determined by network bandwidth. Most of the models discussed are instances of so-called *metrical task systems* [6], [7]) for which there is, e.g., an asymptotically optimal deterministic $\Theta(n)$ -competitive algorithm, where n is the state space; or a randomized $O(\log^2 n \cdot \log \log n)$ -competitive algorithm given that the state space fulfills the triangle inequality: this algorithm uses a (well separated) tree approximation for the general metric space (in a preprocessing step) and subsequently solves the problem on this distorted space; unfortunately, both algorithmic parts are rather complex. As pointed out in [4], there is an intriguing relationship between server migration and *online function tracking* [5], [29]. In online function tracking, an entity Alice needs to keep an entity Bob (approximately) informed about a dynamically changing function, without sending too many updates. The online function tracking problem can be transformed into a chain network where the function values are represented by the nodes on the chain, and a sequence of value changes corresponds to a request pattern on the chain.

REFERENCES

- [1] Y. Bartal. Distributed paging. In *Dagstuhl Workshop on On-line Algorithms*, pages 97–117, 1996.
- [2] A. Bavier, N. Feamster, M. Huang, L. Peterson, and J. Rexford. In vini veritas: Realistic and controlled network experimentation. In *ACM SIGCOMM*, 2006.
- [3] S. Bhatia, M. Motiwala, W. Muhlbauer, V. Valancius, A. Bavier, N. Feamster, L. Peterson, and J. Rexford. Hosting virtual networks on commodity hardware. Technical Report GT-CS-07-10, Georgia Tech, 2008.
- [4] M. Bienkowski, A. Feldmann, D. Jurca, W. Kellerer, G. Schaffrath, S. Schmid, and J. Widmer. Competitive analysis for service migration in vnets. In *Proc. 2nd ACM SIGCOMM Workshop on Virtualized Infrastructure Systems and Architectures (VISA)*, 2010.
- [5] M. Bienkowski and S. Schmid. Online function tracking with generalized penalties. In *Proc. 12th Scandinavian Symposium and Workshops on Algorithm Theory (SWAT)*, 2010.
- [6] A. Borodin and R. El-Yaniv. *Online computation and competitive analysis*. Cambridge University Press, 1998.
- [7] A. Borodin, N. Linial, and M. E. Saks. An optimal on-line algorithm for metrical task system. *J. ACM*, 39(4):745–763, 1992.
- [8] K. Chowdhury, M. R. Rahman, and R. Boutaba. Virtual network embedding with coordinated node and link mapping. In *Proc. INFOCOM*, 2009.
- [9] M. K. Chowdhury and R. Boutaba. A survey of network virtualization. *Elsevier Computer Networks*, 54(5), 2010.
- [10] D. D. Clark, J. Wroclawski, K. R. Sollins, and R. Braden. Tussle in cyberspace: Defining tomorrow’s Internet. In *Proc. SIGCOMM*, 2002.
- [11] J. Fan and M. H. Ammar. Dynamic topology configuration in service overlay networks: A study of reconfiguration policies. In *Proc. INFOCOM*, 2006.
- [12] D. Fotakis. On the competitive ratio for online facility location. In *Proc. 30th International Conference on Automata, Languages and Programming (ICALP)*, also appeared in *Algorithmica* 50(1), pp. 1–57, 2008, pages 637–652, 2003.
- [13] C. Guo, G. Lu, H. Wang, S. Yang, C. Kong, P. Sun, W. Wu, and Y. Zhang. Secondnet: A data center network virtualization architecture with bandwidth guarantees. In *Proc. ACM CONEXT*, 2010.
- [14] F. Hao, T. V. Lakshman, S. Mukherjee, and H. Song. Enhancing dynamic cloud-based services using network virtualization. *SIGCOMM Comput. Commun. Rev.*, 40(1):67–74, 2010.
- [15] I. Houidi, W. Louati, D. Zeghlache, P. Papadimitriou, , and L. Mathy. Adaptive virtual network provisioning. In *Proc. 2nd ACM SIGCOMM Workshop on Virtualized Infrastructure Systems and Architectures (VISA)*, 2010.
- [16] N. Laoutaris, G. Smaragdakis, K. Oikonomou, I. Stavrakakis, and A. Bestavros. Distributed placement of service facilities in large-scale networks. In *IEEE INFOCOM*, 2007.
- [17] J. Lischka and H. Karl. A virtual network mapping algorithm based on subgraph isomorphism detection. In *Proc. VISA*, pages 81–88, 2009.
- [18] J. Lu and J. Turner. Efficient mapping of virtual networks onto a shared substrate. In *Technical Report, WUCSE-2006-35, Washington University*, 2006.
- [19] J. R. M. Yu, Y. Yi and M. Chiang. Rethinking virtual network embedding: Substrate support for path splitting and migration. *ACM SIGCOMM Computer Communication Review*, 38(2):17–29, Apr 2008.
- [20] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner. Openflow: Enabling innovation in campus networks. *SIGCOMM Comput. Commun. Rev.*, 38(2):69–74, 2008.
- [21] B. Monien and H. Sudborough. Embedding one interconnection network in another. In *Computational Graph Theory*, 1990.
- [22] K.-m. Park and C.-k. Kim. A framework for virtual network embedding in wireless networks. In *Proc. 4th International Conference on Future Internet Technologies (CFI)*, pages 5–7, 2009.
- [23] R. Potter and A. Nakao. Mobitopolo: A portable infrastructure to facilitate flexible deployment and migration of distributed applications with virtual topologies. In *Proc. 1st ACM Workshop on Virtualized Infrastructure Systems and Architectures (VISA)*, pages 19–28, 2009.
- [24] R. Ricci, C. Alfeld, and J. Lepreau. A solver for the network testbed mapping problem. *SIGCOMM Comput. Commun. Rev.*, 33(2):65–81, 2003.
- [25] G. Schaffrath, C. Werle, P. Papadimitriou, A. Feldmann, R. Bless, A. Greenhalgh, A. Wundsam, M. Kind, O. Maennel, and L. Mathy. Network virtualization architecture: Proposal and initial prototype. In *Proc. VISA*, 2009.
- [26] N. Spring, R. Mahajan, and T. Anderson. Quantifying the causes of path inflation. In *Proc. SIGCOMM*, 2003.
- [27] N. Spring, R. Mahajan, D. Wetherall, and T. Anderson. Measuring isp topologies with rocketfuel. *IEEE/ACM Trans. Netw.*, 12(1):2–16, 2004.
- [28] US National Science Foundation. Global environment for network innovations (geni). <http://www.geni.net/>, 2006.
- [29] K. Yi and Q. Zhang. Multi-dimensional online tracking. In *Proc. of the 19th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 1098–1107, 2009.
- [30] Y. Zhu and M. H. Ammar. Algorithms for assigning substrate network resources to virtual network components. In *Proc. INFOCOM*, 2006.