

Power-aware replica placement in tree networks with multiple servers per client

Guillaume Aupy, Anne Benoit, Matthieu Journault, Yves Robert

► **To cite this version:**

Guillaume Aupy, Anne Benoit, Matthieu Journault, Yves Robert. Power-aware replica placement in tree networks with multiple servers per client. [Research Report] RR-8474, INRIA. 2014. <hal-00949252>

HAL Id: hal-00949252

<https://hal.inria.fr/hal-00949252>

Submitted on 19 Feb 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Power-aware replica placement in tree networks with multiple servers per client

Guillaume Aupy, Anne Benoit, Matthieu Journault, Yves Robert

**RESEARCH
REPORT**

N° 8474

February 2014

Project-Team ROMA



Power-aware replica placement in tree networks with multiple servers per client

Guillaume Aupy*, Anne Benoit*[†], Matthieu Journault*, Yves Robert*^{‡†}

Project-Team ROMA

Research Report n° 8474 — February 2014 — 27 pages

Abstract: In this paper, we revisit the well-studied problem of replica placement in tree networks. Rather than minimizing the number of servers needed to serve all client requests, we aim at minimizing the total power consumed by these servers. In addition, we use the most general (and powerful) server assignment policy, where the requests of a client can be served by multiple servers located in the (unique) path from this client to the root of the tree. We consider multi-modal servers that can operate at a set of discrete speeds, using the dynamic voltage and frequency scaling (DVFS) technique. The optimization problem is to determine an optimal location of the servers in the tree, as well as the speed at which each server is operated. A major result is the NP-completeness of this problem, to be contrasted with the minimization of the number of servers, which has polynomial complexity. Another important contribution is the formulation of a Mixed Integer Linear Program (MILP) for the problem, together with the design of several polynomial-time heuristics. We assess the efficiency of these heuristics by simulation. For mid-size instances (up to 30 nodes in the tree), we evaluate their absolute performance by comparison with the optimal solution (obtained via the MILP). The most efficient heuristics provide satisfactory results, within 20% of the optimal solution.

Key-words: replica placement, tree network, multiple servers, power-aware algorithm, complexity, discrete speeds.

* Ecole Normale Supérieure de Lyon, LIP, CNRS, INRIA, Université de Lyon, France

† Institut Universitaire de France, France

‡ University of Tennessee Knoxville, USA

**RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES**

Inovallée
655 avenue de l'Europe Montbonnot
38334 Saint Ismier Cedex

Placement de répliques dans des réseaux en forme d'arbre sous contraintes de puissance avec plusieurs serveurs par client

Résumé : Dans ce papier, nous revisitons le problème de placement de répliques dans des réseaux en forme d'arbre. Plutôt que de minimiser le nombre de serveurs requis pour satisfaire toutes les requêtes des clients, nous visons à minimiser la puissance totale consommée par ces serveurs. Nous utilisons la politique d'assignation la plus générale où les requêtes d'un client peuvent être servies sur plusieurs serveurs, tous placés sur le chemin entre ce client et la racine de l'arbre. Nous considérons des serveurs multi-modaux qui peuvent travailler à plusieurs vitesses (ensemble discret de vitesses), en employant la technique du DVFS (Dynamic Voltage and Frequency Scaling). Le problème d'optimisation revient à choisir un placement optimal des serveurs sur les nœuds de l'arbre, ainsi qu'à déterminer leurs vitesses d'exécution. Alors que le problème de minimisation du nombre de serveurs pouvait être résolu en temps polynomial, nous montrons que ce nouveau problème est NP-complet. Nous formulons toutefois un programme linéaire en nombres entiers permettant de résoudre ce problème de façon optimale (mais en un temps exponentiel), ainsi que plusieurs heuristiques polynomiales. Par de multiples simulations, nous avons étudié ces heuristiques, en les comparant pour des instances de tailles moyenne (jusqu'à 30 nœuds) à la solution optimale obtenue par le programme linéaire. Les heuristiques les plus efficaces sont en moyenne à un facteur de 30% de la solution optimale.

Mots-clés : placement de répliques, arbres, serveurs multiples, puissance, complexité, vitesses discrètes.

1 Introduction

In this paper, we revisit the well-studied problem of replica placement in tree networks. Replica placement in tree networks is an important problem [7, 19, 2], with a broad spectrum of applications, such as electronic, ISP, or VOD service delivery (see [11, 7, 13] and additional references in [19]). The problem is the following: one is given a tree-shaped network where clients are periodically issuing requests to be satisfied by servers. The clients are known (both their position in the tree and their number of requests per time unit), while the number and location of the replicas (also called servers) are to be determined. Clients are leaves of the tree, and requests can be served by one or several internal nodes. Note that the distribution tree (clients and nodes) is fixed in the approach.

Initially, there is no replica; when a node is equipped with a replica, it can process a number of requests, up to its capacity limit. Nodes equipped with a replica, also called servers, can only serve clients located in their subtree (so that the root, if equipped with a replica, can serve any client); this restriction is usually adopted to enforce the hierarchical nature of the target application platforms, where a node has knowledge only of its parent and children in the tree. More precisely, there are three classical policies to serve the requests of a client [2]: (i) *Closest*: All requests of a client must be served by the first server located in the path from this client to the root; (ii) *Single*: All requests of a client must be served by a single server, located anywhere in the path from this client to the root; and (iii) *Multiple*: The requests of a client can be served by several servers, all located in the path from this client to the root. For instance in the *Multiple* policy, half the requests of a client can be served by one server, and the other half by another server located higher in this path. In this paper, we study the *Multiple* policy, because it is the most flexible, hence it will lead to the most efficient solution in terms of both the number of servers and total consumed power.

The classical optimization objective in the literature is the number of servers needed to serve all requests. However, minimizing the total power consumed by the servers has recently become a very important objective, both for economic and environmental reasons [15]. To help reduce power dissipation, multi-modal servers are used: each server has a *discrete* number of predefined speeds, which correspond to different voltages that the server can be subjected to. State-of-the-art processors can only be operated with a restricted number of voltage levels, hence with a few speeds [12, 10]. The power consumption is the sum of a static part (the cost for a server to be on and operated) and a dynamic part. This dynamic part is a strictly convex function of the server speed, so that the execution of a given amount of work costs more power if a server runs at a higher speed [10]. More precisely, a server operated at speed s dissipates s^3 watts [5, 4, 17]. Faster speeds allow servers to handle more requests per time unit, but at the price of a much higher (supra-linear) power consumption.

A major contribution of this paper is to show that minimizing power consumption is an NP-complete problem, even if the servers are already placed in the network (and without static power). This is to be contrasted with the polynomial complexity of minimizing the number of servers [2]. Another major contribution is the design of a set of heuristics to minimize power consumption. These heuristics work in two steps: (i) server placement and (ii) request assign-

ment. The placement step relies on an interesting theoretical result: given a fixed set of servers that should all be used, and assuming continuous speeds, it is possible to optimally assign the requests to these servers in polynomial time. We can therefore easily derive a greedy algorithm to place the servers in the continuous case, because for a given placement, we can directly compute the corresponding optimal power consumption. Of course, assuming continuous speeds is not realistic, but it is a handy simplification of the problem: with continuous speeds, once requests are assigned to servers, each server can operate just at the right speed, namely the sum of its requests, so that selecting the server speeds is immediate. With discrete speeds, the problem is more challenging and may well lead to re-assign the requests, for a given placement of servers. To see this, we start from the solution with continuous speeds (including the greedy placement and the optimal request assignment). Let r be the number of requests processed by a given server N in the solution with continuous speeds. With discrete speeds, we have to use the smallest speed s that is larger than r , thereby losing a lot of power if the difference $s - r$ is large. If it is the case, we can try and re-assign some requests to another server N' located upper in the path from N to the tree root. There would then remain only s' requests to be served by N , where s' is the largest speed that is smaller than r : this saves power locally by avoiding the large $s - r$ gap, but we have to re-assign $r - s'$ requests to another server, and this has a cost that should be balanced with the local gain. Such trade-off decisions are exactly those taken in the request assignment step of the heuristics.

To the best of our knowledge, this paper is the first to propose heuristics for power minimization with multiple servers, hence we cannot use any heuristics from the literature as reference. However, we have derived a Mixed Integer Linear Program (MILP) to compute the optimal solution to the power minimization problem. Using this linear program has (potentially) an exponential cost, but it enables us to assess the absolute performance of the heuristics, at least for small-size problems.

The rest of the paper is organized as follows. Section 2 surveys related work. Section 3 is devoted to a precise statement of the framework. Section 4 assesses the complexity of the power minimization problem, through an intricate NP-completeness proof. This section also provides the MILP to compute the optimal solution. Section 5 introduces several heuristics to solve the problem. The placement step is an incremental greedy procedure, whose evaluation is based on the optimal solution for request assignment with fixed servers, when assuming continuous speeds. Section 6 reports simulation results and comparisons of the heuristics, together with their absolute performance evaluation: the distance to the optimal solution is computed through the linear program for instances with up to 30 servers.

2 Related work

Many papers considering the replica placement problem deal with general graphs, while we focus in this work on tree networks. In the problem with a general graph network, it is already difficult to decide which spanning tree to use, in order to optimize some global objective function. A survey of work targeting performance issues can be found in [14]. Recently, some work start to tackle

energy-related problems. For instance, in [18], the authors discuss thermal and power-aware task scheduling and data placement heuristics, in the context of a Hadoop system. All problems are NP-hard, and there is no tree structure but rather a set of racks, and a set of data nodes per rack.

For tree networks, a large effort has been spent to optimize the performance of replica placements, assuming that the spanning tree was given, or that the network had a tree structure initially. Most work has focused on the *Closest* policy, where a client is requested to be served by the closest server on the path towards the root of the tree, see for instance [7, 13]. Kalpakis et al. [11] studied a variant with bi-directional links, and therefore the tree structure may not be respected anymore, and a client may be served by a node that is not its ancestor in the tree. While the problem with a tree structure has polynomial complexity, the bi-directional problem becomes NP-complete.

Following this line of work, we had investigated in our previous work [3] the complexity of the power-aware replica placement problem with the *Closest* policy, and proved that the problem becomes NP-complete when the objective is to minimize the total power consumption. We considered servers with several distinct possible speeds, and a server operating at a given speed consumes a power composed of a static part and a dynamic part proportional to the cube of the speed. We keep the same model in this paper, because it is a classical model extensively used when considering dynamic voltage and frequency scaling (DVFS) technique [5, 4, 17].

The *Multiple* policy is more flexible than *Closest* because it loosens placement rules: the requests of a client can be processed by several servers located anywhere in the path from the client to the root. Similarly to the *Closest* policy, the problem of minimizing the cost of the replica placement can be solved in polynomial time [2]. However, we are not aware on any other work aiming at optimizing the power consumption on tree networks for this *Multiple* policy.

3 Framework

This section is devoted to a precise statement of the framework. We start with a description of the replica placement problem. Then we detail the power consumption model. Finally, we state the objective function, and the corresponding optimization problems.

3.1 Replica placement

We consider a distribution tree whose nodes are partitioned into a set of clients \mathcal{C} , and a set of N nodes, \mathcal{N} . The clients are leaf nodes of the tree, while \mathcal{N} is the set of internal nodes. Each client $i \in \mathcal{C}$ (leaf of the tree) is sending r_i requests per time unit to a database object. Internal nodes equipped with a replica (also called *servers*) can process requests from clients in their subtree. If a server $j \in \mathcal{N}$ is operated at speed s_j , then it can process up to s_j requests per time unit. Both the r_i 's and the s_j 's are assumed to take rational values. Note that it would be easy to allow *client-nodes* that play both the rule of a client and of a node (possibly a server), by dividing such a node into an internal node and a leaf in the tree.

For each client $i \in \mathcal{C}$ and each node $j \in \mathcal{N}$, $r_{i,j}$ is the number of requests from client i processed by server j . We must have $\sum_{j \in \mathcal{N}} r_{i,j} = r_i$ for all $i \in \mathcal{C}$, i.e., all requests are processed. Furthermore, a server cannot process more requests than its assigned speed, i.e., $w_j = \sum_{i \in \mathcal{C}} r_{i,j} \leq s_j$ for all $j \in \mathcal{N}$, where w_j is the load of server j . The set of replicas is defined as $\mathcal{R} = \{j \in \mathcal{N} \mid \exists i \in \mathcal{C}, r_{i,j} > 0\}$.

3.2 Power consumption model

We (realistically) consider discrete speeds. Servers may operate only at a set $\{s_1, \dots, s_K\}$ of different (rational) speeds, depending upon the number of requests that they have to process per time unit. We assume that $0 \leq s_1 \leq \dots \leq s_K$, and therefore no server can handle more than s_K requests. A server with a load w will therefore operate at speed s_k , where $s_{k-1} < w \leq s_k$ (letting $s_0 = -1$ for the limit case). The power consumption of a server $j \in \mathcal{R}$ operated at speed $s(j)$ obeys the classical model,

$$\mathcal{P}(j) = \mathcal{P}_{\text{static}} + s(j)^3,$$

and the total power consumption $\mathcal{P}(\mathcal{R})$ of the solution is the sum of the power consumption of all server nodes:

$$\mathcal{P}(\mathcal{R}) = \sum_{j \in \mathcal{R}} \mathcal{P}(j) = \sum_{j \in \mathcal{R}} (\mathcal{P}_{\text{static}} + s(j)^3) = |\mathcal{R}| \times \mathcal{P}_{\text{static}} + \sum_{j \in \mathcal{R}} s(j)^3, \quad (1)$$

where $|\mathcal{R}|$ is the total number of servers in the solution.

3.3 Optimization problems

The main optimization problem is the DISCRETE problem: given a distribution tree (with a number of requests per client), decide where to place the servers, and how to distribute client requests among them (which can also be seen as assigning the speed of each server), in order to minimize the total power consumption.

We also consider the sub-problem where the servers are already placed in the tree, DISCRETE-PLACED. The goal is then only to decide how to distribute requests among servers, hence at which speed to operate each server, in order to minimize the total power consumption.

4 Complexity results

In this section, we establish the NP-completeness of DISCRETE and DISCRETE-PLACED (Section 4.1). Then in Section 4.2, we provide a Mixed Integer Linear Program (MILP) to solve DISCRETE.

4.1 NP-completeness

Theorem 1. *The DISCRETE and DISCRETE-PLACED problems are NP-complete, even with $\mathcal{P}_{\text{static}} = 0$.*

Proof. Consider the decision problem associated to DISCRETE: given a tree network with a set of requests, a set of server speeds, and a bound K on power consumption, can we find a placement of servers and an assignment of the requests to these servers so that the bound K on power consumption is met? The problem is clearly in NP: given the servers and the request assignments, we can compute the load of each server, select the smallest possible speed and compute the total power consumption in polynomial time.

To establish the completeness, we use a reduction from 2-Partition [9]. We consider an instance \mathcal{I}_1 of 2-Partition: given n strictly positive integers a_1, \dots, a_n , does there exist a subset I of $\{1, \dots, n\}$ such that $\sum_{i \in I} a_i = \sum_{i \notin I} a_i$? Let $S = \frac{1}{2} \sum_{i=1}^n a_i$. We further assume, without loss of generality, that S is an integer, and that $a_i < S$ for $1 \leq i \leq n$.

We build the following instance \mathcal{I}_2 of the DISCRETE problem:

- $\mathcal{P}_{\text{static}} = 0$
- We first define two constants (see below for an intuitive explanation):

$$X = 2n(7n^2S)^4, \quad (2)$$

$$M = 7n^2X^2. \quad (3)$$

- There are $K = 2n + 1$ possible speeds for servers:

$$\text{for } 1 \leq i \leq n, \begin{cases} s_i = iX \\ \tilde{s}_i = iX \left(1 + \frac{1}{3} \frac{Ma_i}{(iX)^3} - \frac{1}{9} \left(\frac{Ma_i}{(iX)^3} \right)^2 + \frac{5}{81} \left(\frac{Ma_i}{(iX)^3} \right)^3 \right) \end{cases} \quad (4)$$

$$\text{and } s_{n+1} = (n+1)X.$$

- The tree is a fork with $n+1$ internal nodes; node 0 is the root of the tree and node i is a child of node 0, for $1 \leq i \leq n$. Each internal node i has a unique leaf child with r_i requests, where $r_0 = (n+1)X - S$ and $r_i = s_i + a_i$ for $1 \leq i \leq n$.
- The bound on power consumption is $B = (n+1)^3X^3 + \sum_{i=1}^n ((iX)^3 + Ma_i) - MS + 0.5$.

Note that for all i , \tilde{s}_i is a rational number polynomial in the size of the instance. The size of \mathcal{I}_2 is polynomial in the size of \mathcal{I}_1 . Let us further define a new set of speeds:

$$\text{for } 1 \leq i \leq n, \quad s'_i = ((iX)^3 + Ma_i)^{1/3}. \quad (5)$$

The speeds s'_i are not part of the instance but they will be used for the proof. Intuitively, \tilde{s}_i is a close approximation of s'_i for $1 \leq i \leq n$.

The intuition behind the proof is that there is a server on each node, hence the proof works for both problems (DISCRETE and DISCRETE-PLACED). Furthermore, the root node will run at speed s_{n+1} , and it will process exactly $(n+1)X$ requests. Therefore, S requests can come from its children. Because of the speed configurations, server i has exactly two choices, either speed s_i or speed \tilde{s}_i . The constants X and M ensure that no other speed can be used,

because speeds are spaced apart enough. Either server i processes all its requests, and it must run at speed \tilde{s}_i , or a_i requests are forwarded to the root server and server i can then run at speed s_i . There can be at most S requests going up (server capacity limit at the root node), and because of the bound on the power consumption, there should be at least S requests going up, hence the 2-partition. Indeed, the power saved when moving from speed s'_i (and hence almost \tilde{s}_i) to speed s_i is linear in a_i : $s'_i{}^3 - s_i^3 = Ma_i$.

First we prove some preliminary properties on the \tilde{s}_i 's to ensure that the s'_i 's are a good approximation and that speeds are spaced apart enough so that each server is restricted to its set of two speeds.

Lemma 1. *The following properties hold true:*

1. for $1 \leq i \leq n$, $s'_i \leq \tilde{s}_i$;
2. for $1 \leq i \leq n$, $s_i + a_i = iX + a_i \leq \tilde{s}_i$;
3. for $1 \leq i \leq n$, $\tilde{s}_i \leq (i+1)X = s_{i+1}$;
4. $\sum_{i=1}^n s'_i{}^3 \leq \sum_{i=1}^n \tilde{s}_i^3 \leq \sum_{i=1}^n s'_i{}^3 + 0.5$;
5. for $1 \leq i \leq n$, $\tilde{s}_i \leq iX + (Ma_i)^{1/3}$.

Proof. Consider the following functions on $[0, 1]$: $f_1 : x \mapsto (1+x)^{1/3}$ and $f_2 : x \mapsto 1+x/3-x^2/9+5x^3/81$. We have $s'_i = iX f_1\left(\frac{Ma_i}{(iX)^3}\right)$ and $\tilde{s}_i = iX f_2\left(\frac{Ma_i}{(iX)^3}\right)$, and note that $\frac{Ma_i}{(iX)^3} \leq \frac{MS}{X^3} < 1$, because $M = 7n^2X^2$ and $X > 7n^2S$.

We prove the five properties:

1. ($\tilde{s}_i - s'_i \geq 0$) We can verify, using several differentiations (or a computer algebra software), that $f_2(x) - f_1(x) > 0$ for $x \in [0, 1]$.
2. ($iX + a_i \leq \tilde{s}_i$) Thanks to Property 1, $\tilde{s}_i > s'_i$, and showing that $s'_i \geq s_i + a_i$ gives the result:

$$\begin{aligned} & s'_i \geq s_i + a_i \\ \Leftrightarrow & (s_i^3 + Ma_i)^{1/3} \geq s_i + a_i \\ \Leftrightarrow & s_i^3 + Ma_i \geq s_i^3 + 3s_i^2a_i + 3s_ia_i^2 + a_i^3 \\ \Leftrightarrow & M \geq 3(iX)^2 + 3iXa_i + a_i^2 \end{aligned}$$

Because $X > a_i$ and $n \geq i \geq 1$, we have $M = 7n^2X^2 \geq 3(iX)^2 + 3iXa_i + a_i^2$, hence the result.

3. ($\tilde{s}_i \leq (i+1)X$) We have already seen that $\frac{Ma_i}{X^3} < 1$, and therefore $\tilde{s}_i \leq iX \left(1 + \frac{1}{3i} + \frac{5}{81i}\right) = \left(i + \frac{32}{81}\right)X \leq (i+1)X$.
4. ($\sum_{i=1}^n s'_i{}^3 \leq \sum_{i=1}^n \tilde{s}_i^3 \leq \sum_{i=1}^n s'_i{}^3 + 0.5$) We can verify, using several differentiations (or a computer algebra software), that $\forall x \in [0, 1]$,

$$f_2(x)^3 - f_1(x)^3 \leq x^4.$$

Hence, $\tilde{s}_i^3 - s_i'^3 = (iX)^3 \left(f_2 \left(\frac{Ma_i}{(iX)^3} \right)^3 - f_1 \left(\frac{Ma_i}{(iX)^3} \right)^3 \right) \leq \frac{(7n^2 a_i)^4}{i^9 X}$. Hence:

$$\begin{aligned} \sum_{i=1}^n \tilde{s}_i^3 &\leq \sum_{i=1}^n \left(s_i'^3 + \frac{(7n^2 a_i)^4}{i^9 X} \right) \\ &\leq \sum_{i=1}^n s_i'^3 + n \frac{(7n^2 a_i)^4}{X} \\ &\leq \sum_{i=1}^n s_i'^3 + n \frac{(7n^2 S)^4}{X} \\ &\leq \sum_{i=1}^n s_i'^3 + 0.5 \quad \text{because } X \geq 2n (7n^2 S)^4 \end{aligned}$$

5. ($\tilde{s}_i \leq iX + \frac{Ma_i}{3}$) We have $\tilde{s}_i - iX - (Ma_i)^{1/3} = \frac{Ma_i}{3(iX)^2} - (Ma_i)^{1/3} - \frac{iX}{9} \left(\frac{Ma_i}{(iX)^3} \right)^2 + \frac{5iX}{81} \left(\frac{Ma_i}{(iX)^3} \right)^3$. Then, first let us show that $-\frac{iX}{9} \left(\frac{Ma_i}{(iX)^3} \right)^2 + \frac{5iX}{81} \left(\frac{Ma_i}{(iX)^3} \right)^3 < 0$:

$$\begin{aligned} \frac{iX}{9} \left(\frac{Ma_i}{(iX)^3} \right)^2 &> \frac{5iX}{81} \left(\frac{Ma_i}{(iX)^3} \right)^3 \\ \frac{9}{5} (iX)^3 &> Ma_i \\ \frac{9}{5} &> \frac{7n^2 a_i}{i^3 X} \end{aligned}$$

This last equation is always true since $X > 7n^2 a_i$. Then let us show that $\frac{Ma_i}{3(iX)^2} - (Ma_i)^{1/3} < 0$:

$$\begin{aligned} (Ma_i)^{1/3} &> \frac{Ma_i}{3(iX)^2} \\ 3i^2 X^2 &> (7n^2 X^2 a_i)^{2/3} \\ X^{1/3} &> \frac{7^{2/3} n^{4/3} a_i}{3i^2} \\ X &> \frac{49n^4 a_i^3}{27i^6} \end{aligned}$$

which is true since $X = 2n (7n^2 S)^4$.

This concludes the proof of Lemma 1. \square

With these properties, we are now ready to show the NP-hardness of the problem. Note that $s_1 \leq \tilde{s}_1 \leq s_2 \leq \tilde{s}_2 \leq \dots \leq s_n \leq \tilde{s}_n \leq s_{n+1}$.

Suppose first that instance \mathcal{I}_1 has a solution I . We assign the servers as follows: the root node is assigned a server at speed s_{n+1} ; for all $i \in I$, we assign to node i a server at speed s_i , and for all $i \notin I$, we assign to node i a server at speed \tilde{s}_i .

Because $\tilde{s}_i \geq s_i + a_i$ (Lemma 1), the server on node i can process all its requests, for $i \notin I$. On the contrary, for $i \in I$, the server can only process s_i requests, and therefore a_i requests are going up in the tree. Hence, a total of $\sum_{i \in I} a_i = S$ requests need to be processed by the root node, in addition to the $(n+1)X - S$ requests that arrive directly to this node (and the speed is s_{n+1} so the root server can process all its requests).

Finally we need to compute the power consumption, that is:

$$\begin{aligned}
E &= (n+1)^3 X^3 + \sum_{i \notin I} \tilde{s}_i^3 + \sum_{i \in I} s_i^3 \\
&= (n+1)^3 X^3 + \sum_{i \notin I} \tilde{s}_i^3 + \sum_{i \in I} (s_i^3 - M a_i) \\
&\leq (n+1)^3 X^3 + \sum_{i \notin I} \tilde{s}_i^3 + \sum_{i \in I} \tilde{s}_i^3 - M \sum_{i \in I} a_i \\
&= (n+1)^3 X^3 + \sum_{i=1}^n \tilde{s}_i^3 - MS \\
&\leq (n+1)^3 X^3 + \sum_{i=1}^n s_i^3 + 0.5 - MS = B
\end{aligned}$$

The power bound is respected, and therefore we have a solution to \mathcal{I}_2 .

Suppose now that \mathcal{I}_2 has a solution. We define the set I of nodes as follows:

$$i \in I \iff \text{there is a server on node } i \text{ operating at speed } s < \tilde{s}_i$$

We show that I is a solution to \mathcal{I}_1 . First let us show that for $i \in I$, the speed of the server on node i is exactly s_i . We prove this by contradiction: if this server has a lower speed s , then the number of requests from node i that goes up to the root node is at least

$$\begin{aligned}
iX + a_i - s &\geq iX + a_i - ((i-1)X + (Ma_{i-1})^{1/3}) \\
&\geq X + a_i - (Ma_{i-1})^{1/3},
\end{aligned}$$

because $s \leq \tilde{s}_{i-1}$, and $\tilde{s}_{i-1} \leq (i-1)X + (Ma_{i-1})^{1/3}$ (Lemma 1). This would lead to execute at the root more requests than the maximum $(n+1)X$. Indeed,

$$\begin{aligned}
(n+1)X - S + X + a_i - (Ma_{i-1})^{1/3} &> (n+1)X \\
\iff X &> (Ma_{i-1})^{1/3} + S - a_i \\
\iff X &> (a_{i-1}7n^2X^2)^{1/3} + S - a_i
\end{aligned}$$

We know that $S < X^{2/3}$, and that $X^{1/3} > (1+(7n^2)^{1/3})$ (we use $a_i < S$), hence the contradiction. Therefore we conclude that for all $i \in I$, the server on node i works at speed s_i .

Then for all $i \notin I$, there is a server on node i (otherwise the root node could not process the $iX + a_i$ upcoming requests with the largest speed $(n+1)X$), and its speed is at least \tilde{s}_i by definition of I . We have seen that $\tilde{s}_i \geq iX + a_i$, so every request can be processed directly on node i .

Consequently, the power consumption is at least $E \geq (n+1)^3 X^3 + \sum_{i \notin I} \tilde{s}_i^3 + \sum_{i \in I} s_i^3 \geq (n+1)^3 X^3 + \sum_{i=1}^n s_i^3 - M \sum_{i \in I} a_i$. Then we have, $M \sum_{i \in I} a_i \geq$

$MS - 0.5$ because $E \leq B$ (we have a solution of the problem). Because M and S are integers, it implies that $\sum_{i \in I} a_i \geq S$. Because of the constraint on the maximum speed at the root node, the quantity of requests coming from other nodes and processed at the root cannot be greater than S , and we have $\sum_{i \in I} a_i \leq S$. Indeed, for all $i \in I$, no more than iX requests can be executed on node i , and therefore there are at least a_i requests from node i executed at the root.

Finally, we conclude that $\sum_{i \in I} a_i = S$, and therefore that \mathcal{I}_1 has a solution. This concludes the proof for DISCRETE. The same instance can be used to prove the NP-completeness of DISCRETE-PLACED (which belongs to NP since DISCRETE does), because each internal node is equipped with a server in the solution. \square

4.2 Mixed integer linear program

Theorem 2. *The following MILP characterizes the DISCRETE problem, where the unknown variables are the $x_{j,k}$'s (Boolean variables) and the $y_{i,j}$'s (rational variables), for $j \in \mathcal{N}$, $1 \leq k \leq K$ and $i \in \mathcal{C}$:*

$$\begin{aligned} & \text{Minimize } \sum_{j \in \mathcal{N}} \sum_{1 \leq k \leq K} x_{j,k} (\mathcal{P}_{static} + s_k^3) && \text{subject to} \\ & (i) \sum_{j \in \mathcal{N}} y_{i,j} = r_i, && i \in \mathcal{C} \\ & (ii) \sum_{1 \leq k \leq K} x_{j,k} \leq 1, && j \in \mathcal{N} \\ & (iii) \sum_{i \in \mathcal{C}} y_{i,j} \leq \sum_{1 \leq k \leq K} x_{j,k} s_k, && j \in \mathcal{N} \end{aligned} \quad (6)$$

Proof. The constants are the r_i 's for $i \in \mathcal{C}$, and the s_k 's for $1 \leq k \leq K$, and we consider the following variables:

- $x_{j,k}$ is a boolean variable equal to 1 if j is a server operated at speed s_k , for $j \in \mathcal{N}$ and $1 \leq k \leq K$; $x_{j,k} = 0$ otherwise.
- $y_{i,j}$ is a rational variable equal to $r_{i,j}$, the number of requests of client $i \in \mathcal{C}$ processed by server $j \in \mathcal{N}$; if j is not an ancestor of i in the tree, we directly set $y_{i,j} = 0$.

Then the constraints are:

- For all $i \in \mathcal{C}$, all requests of client i must be processed: $\forall i \in \mathcal{C}, \sum_{j \in \mathcal{N}} y_{i,j} = r_i$;
- Each server is assigned at most one speed: $\forall j \in \mathcal{N}, \sum_{1 \leq k \leq K} x_{j,k} \leq 1$; note that a node j is equipped with a server if and only if $\sum_{1 \leq k \leq K} x_{j,k} = 1$;
- The processing capacity of any server cannot be exceeded: $\forall j \in \mathcal{N}, \sum_{i \in \mathcal{C}} y_{i,j} \leq \sum_{1 \leq k \leq K} x_{j,k} s_k$.

Finally, we minimize the total power consumption. Overall, there are $|\mathcal{C}| + 2|\mathcal{N}|$ constraints and $|\mathcal{N}| \cdot (|K| + |\mathcal{C}|)$ variables in this MILP. \square

5 Heuristics

In this section, we propose some polynomial-time heuristics for the DISCRETE problem. We start by outlining the general principles that have guided their design before exposing the details for each heuristic.

5.1 General principles

As already mentioned in Section 1, the heuristics work in two steps: (i) server placement and (ii) request assignment. The placement step of the heuristics relies on the following result, whose proof can be found in Appendix A:

Proposition 1. *Given a fixed set of servers deployed on a tree of size $s = |\mathcal{C}| + |\mathcal{N}|$ and assuming continuous speeds, the optimal assignment ALG-CONT-PLACED of requests to servers that uses all these servers and minimizes power consumption can be determined in time $O(s^2)$.*

The placement step works incrementally: to compute a solution with k servers, the heuristic starts from a solution with $k - 1$ servers, and then greedily tests the addition of one additional server. It uses Proposition 1 to compute the optimal assignment of requests with this additional server, computes the corresponding power, and returns the best solution over all possible choices for the additional server.

The placement step assumes continuous speeds, hence the loads assigned by ALG-CONT-PLACED to each server do not take the set of actual speeds into account. The second step of the heuristics consists in determining a discrete speed for each server, which usually leads to re-assigning some requests, as explained in Section 1. While the first step of the heuristics is common to all heuristics, we outline below three different methods to perform this request assignment step.

5.2 List of heuristics

We provide here three different heuristics to determine the actual speed of each server. In the first heuristic, GREEDY, we assign the smallest speed equal to or greater than the load given by ALG-CONT-PLACED to each server. While simple, GREEDY provides a $\left(1 + \frac{\max_i(s_{i+1} - s_i)}{s_{\min}}\right)^3$ -approximation for the problem with the placement given by ALG-CONT-PLACED, where s_{\min} is the smallest speed available. Finally, we point out that if there is no speed greater than the value determined by ALG-CONT-PLACED for some server, then there does not exist a solution for this (given) placement (see Appendix B for further details).

The next two heuristics, SPEED and EXCESS, improve the GREEDY heuristic by trying to modify the load of each server, via request re-assignment. The goal is to decrease the speed of some servers. More precisely, in the procedure, which is called EQUILIBRATE and detailed in Algorithm 1, if a server is not loaded up to its full capacity (meaning its load is equal to its capacity), then the heuristics take some load out of its children until this server reaches its capacity (see the loop line 15). The capacity of a server is defined as the maximum between its actual speed and the maximum speed of its children (see line 12), hence transferring even more load to this server if one of its children has a higher

Algorithm 1: Procedure EQUILIBRATE

```

1 procedure EQUILIBRATE( $T$ )
2 begin
3   Apply the GREEDY heuristic on  $T$ ;
4   return sub-EQUILIBRATE ( $T$ , root)
5
6
7 procedure sub-EQUILIBRATE( $T$ ,  $N$ ) /* here  $N$  is a node of  $T$  */
8 begin
9   Let  $L$  the list of children of  $N$  sorted as favored by the heuristic;
10  Let  $(i_N, w_N)$  be the indices of the speed and load of  $N$ ;
11  Let  $i_{N'}$  be the index of the maximum speed of the children of  $N$ ;
12  if  $s_{i_{N'}} > s_{i_N}$  then
13     $i_N \leftarrow i_{N'}$ ;
14  remaining  $\leftarrow s_{i_N} - w_N$ ;
15  while remaining  $\neq 0$  do
16    Let  $N'$  be the first element of  $L$ ;
17    Let  $(i_{N'}, w_{N'})$  be the indices of the speed and load of  $N'$ ;
18    temp  $\leftarrow \min(\text{remaining}, w_{N'} - s_{i_{N'}})$ ;
19     $w_{N'} \leftarrow w_{N'} - \text{temp}$ ;
20    remaining  $\leftarrow \text{remaining} - \text{temp}$ ;
21     $w_N \leftarrow w_N + \text{temp}$ ;
22    Update the index of the speed of  $N'$ ;
23    Update the order of  $L$  to account for the new load and speed of
     $N'$ ;
24  for  $N' \in \text{children of } N$  do
25     $T \leftarrow \text{sub-EQUILIBRATE}(T, N')$ ;
26  return  $T$ 

```

speed (and thus we should be able to reduce the speed of at least one child). This may happen if we have decreased the speed of the current node in a previous step of the algorithm, but not the speed of its children.

The main difference between the two heuristics SPEED and EXCESS lies in the selection of the children whose load is taken out:

- In the SPEED heuristic, we favor the children whose servers have the largest speeds. To break ties if two children of a given server have the same speed, we favor the one with the smallest load. The idea is that the gain in power will be more important if we can decrease the execution speed of a server with a large speed (favor large speeds); and if there is a tie, there is more chance to decrease the speed of a server if its load is small.
- On the contrary, in the EXCESS heuristic, we favor children with small *excess*. The excess of a server is defined as the difference between its load and the largest speed below it. The idea is that we will be able to decrease the speed of more servers if we favor small excess. Finally, when two children have the same excess, we favor the one with the largest load.

Recall that K is the number of speeds and $s = |\mathcal{C}| + |\mathcal{N}|$. The complexity of the GREEDY heuristic is $O(s^2)$, the most costly part being the call to ALG-CONT-PLACED. For the EQUILIBRATE procedure, and hence for the SPEED and EXCESS heuristics, the complexity becomes $O(Ks^2 \log s)$. First, it takes $s \log s$ operations to sort the children of each of the s nodes in a pre-treatment phase. Then, the procedure sub-EQUILIBRATE is called s times, and it has a cost $O(Ks \log s)$: in the worst case, we visit each child, for each possible speed, to fill the parent at full capacity, and we have to include the cost of updating the list of sorted children, that can be done in $O(\log s)$.

6 Simulations

In this section, we report extensive simulations to assess the performance of the heuristics presented in Section 5. All the source code, together with scripts to obtain additional results, are publicly available [1]. The heuristics have been coded using the programming language OCaml, while the MILP computing the optimal solution is generated using the C language and solved using IBM Cplex [8].

In order to evaluate the heuristics, we have generated more than 100 random trees for each simulation. To simplify the generation, each internal node in the tree has a unique client leaf, which is assigned a random rational number of requests between 0 and 100. For processor speeds, unless stated otherwise, we use five speeds spaced as those of the Intel Xscale, following [6, 16]: we suppose that the largest speed can process 150 requests, and the ratio of the different speeds to the largest speed is then: (0.15, 0.4, 0.6, 0.8, 1). In [6, 16], the static power is equal to the power consumed in the lowest speed, which here corresponds to $(0.15 * 150)^3 \approx 11,000$.

We have conducted four different sets of simulations to assess the impact of the number of nodes, of static power, of the number of available speeds and of the total load of requests.

Impact of the number of nodes. In the first set of simulations, we study the impact of the number of nodes on power consumption: in Figure 1, we plot the ratio of the power returned by the heuristics over the power of the optimal solution, with a static power of 1,000, 5,000, 10,000, 20,000 and 100,000 respectively. Note that Figure 1f is different from the others and provides results at larger scale: there we plot the ratio of the power returned by SPEED and EXCESS over the power consumption of GREEDY, with a static power of 50,000, but for a larger number of nodes (up to 200 nodes). While one could expect that the performance of the heuristics would decrease for larger trees, it seems that SPEED and EXCESS reach a tangent after ≈ 15 nodes, and that on average the maximum waste is between 20 and 25%. Furthermore, when the static power is higher (100,000), this maximum waste is even below 20%. This observed tangent is very likely correlated to the set of speeds and to the static power. This tangent makes sense in practice if we assume that the first step (the placement step) of the heuristics is not too far from the optimal solution because the GREEDY heuristic is an approximation algorithm. SPEED and EXCESS are just improvements of the GREEDY heuristic. It would be interesting to see how much

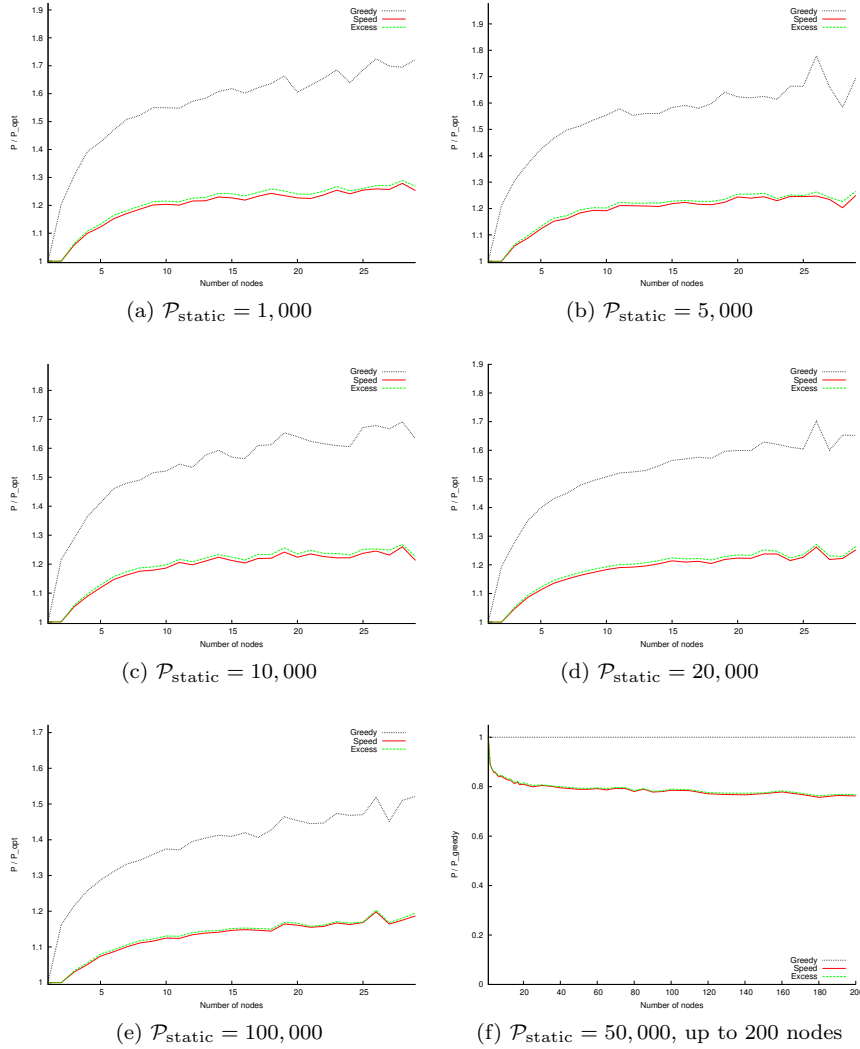


Figure 1 – Study of the impact of the number of nodes, for random requests between 0 and 100, average on 100 tests.

they improve the approximation factor, though probably complicated. Note that in Figure 1f, we see that SPEED and EXCESS are still constantly better than GREEDY even with a larger number of nodes, with a power consumption of around 80% of the power consumption for GREEDY.

Impact of the static power. In the second set of simulations, we have studied the impact of the static power on total power consumption. In Figure 2, we plot the ratio of the power returned by the heuristics over the power of the optimal solution, with a static power varying between 0 and 200,000 for trees of 20 nodes. Note that the higher the static power, the better the results. Indeed, at some point, what matters most is the number of servers placed, and not the

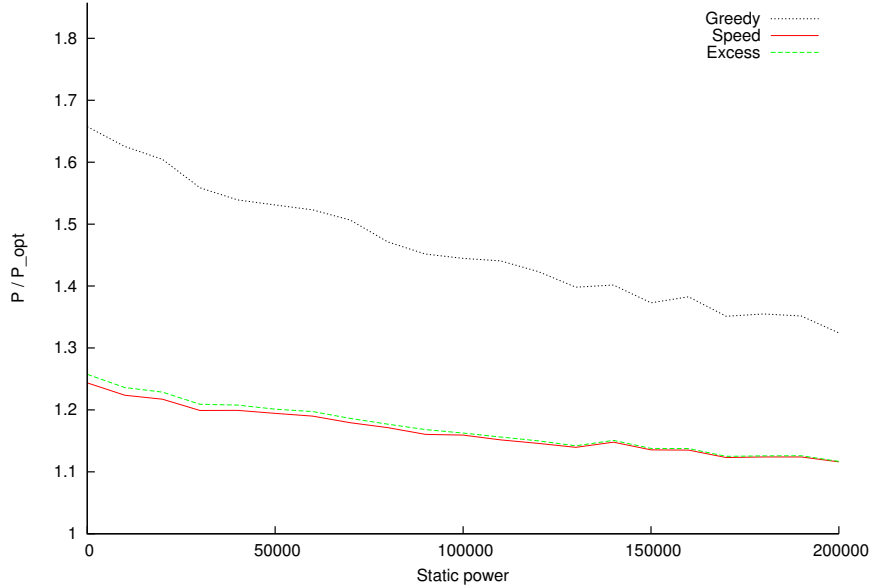


Figure 2 – Study of the impact of the static power, for trees of 20 nodes, for random requests between 0 and 100, average on 100 tests.

allocation of requests, hence GREEDY gets closer to the optimal solution as well.

Impact of the number of speeds. In this third set of simulations, we have studied the impact of the number of speeds on power consumption. For this set of simulations, we do not use Intel speeds anymore, but instead speeds that are equally distributed between 0 and 150. In Figure 3, we plot the ratio of the power returned by the heuristics over the power of the optimal solution, with a static power of 50,000, for trees of 20 nodes, with the number of speeds varying from one (150) and ten (15, 30, 45, 60, 75, 90, 105, 120, 135, 150). When there is only one speed, obviously the results are as good as they can be and only depend on the allocation heuristic. Then starting from three speeds, we observe that the more speeds, the better the results. This was expected since the more speeds we have, the closer we can get to the optimal solution computed by ALG-CONT-PLACED, and the better the results. The fact that the results are better with two speeds than three can be explained by the fact that with only two speeds, it is still easier to find the optimal speed (hence a lower ratio than with three speeds), but a mistake is very expensive, hence a result that is not as good as with four speeds. A final remark: when speeds are equally distributed, the (proven) approximation ratio of the GREEDY heuristic is 8. However in Figure 3 we see that the ratio never goes above 1.8.

Impact of the total load of the tree. Finally, in the last set of simulations, we have studied the impact of the total load of the tree on the power consumption. In Figure 4, we plot the ratio of the power returned by the heuristics over the power of the optimal solution, for trees of 20 nodes, and for a static power of 1,000, 5,000, 10,000, 20,000 and 100,000 respectively. Overall, the total number

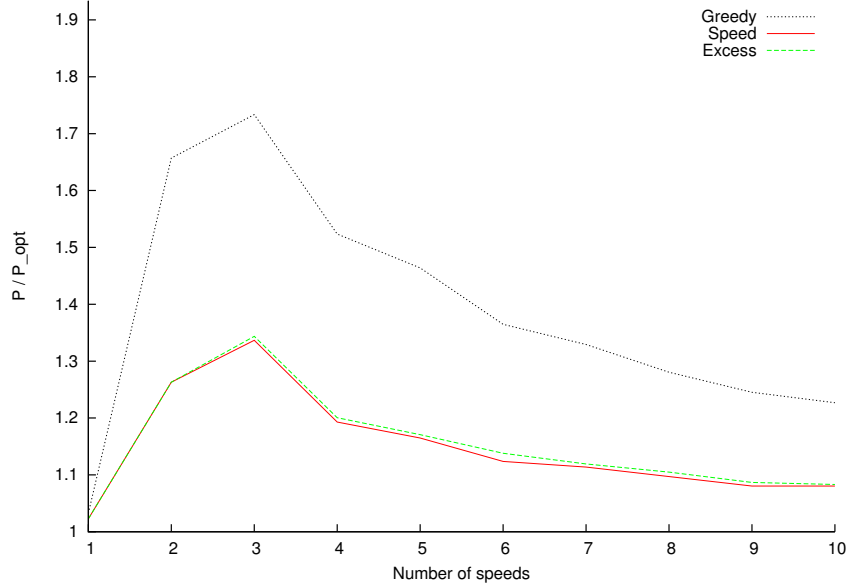


Figure 3 – Study of the impact of the number of speeds, for trees of 20 nodes, static power of 50,000, for random requests between 0 and 100, average on 100 tests.

of requests does not impact the performance of the heuristics: the average ratio stays constant with the total load of the tree, for all values of static power.

Summary of simulation results. To conclude on the different studies, the first observation is somewhat expected: there is a huge gap between the GREEDY heuristic, and the SPEED and EXCESS heuristics: there is a degradation w.r.t. the optimal of 50 to 70% when using the GREEDY heuristic with 10 to 30 nodes, while it is only approximately 20% (or less with larger static power) when using the SPEED or EXCESS heuristic. The difference between the SPEED and EXCESS heuristics is negligible, although it should be noted that on average, the SPEED heuristic performs slightly ($\approx 1\%$) better than the EXCESS heuristic. Furthermore, it seems that what matters most for the competitiveness of the heuristics is the set of speeds and the static power $\mathcal{P}_{\text{static}}$. In particular, the number of speeds is very important: the closer the speeds are to each other, the better the results. Above a certain number of nodes (≈ 15), the ratio of the results of the heuristics over the optimal solution seems to reach a threshold (independently of the load and the static power), but the value of this threshold depends on the set of speeds and on the static power. Higher static power lowers the value of the threshold: at some point, what matters most is the number of servers, even if they are all at maximum speed. Similarly, the smaller the gap between two consecutive speeds, then the closer we can get to the optimal solution computed by ALG-CONT-PLACED, and the better the results.

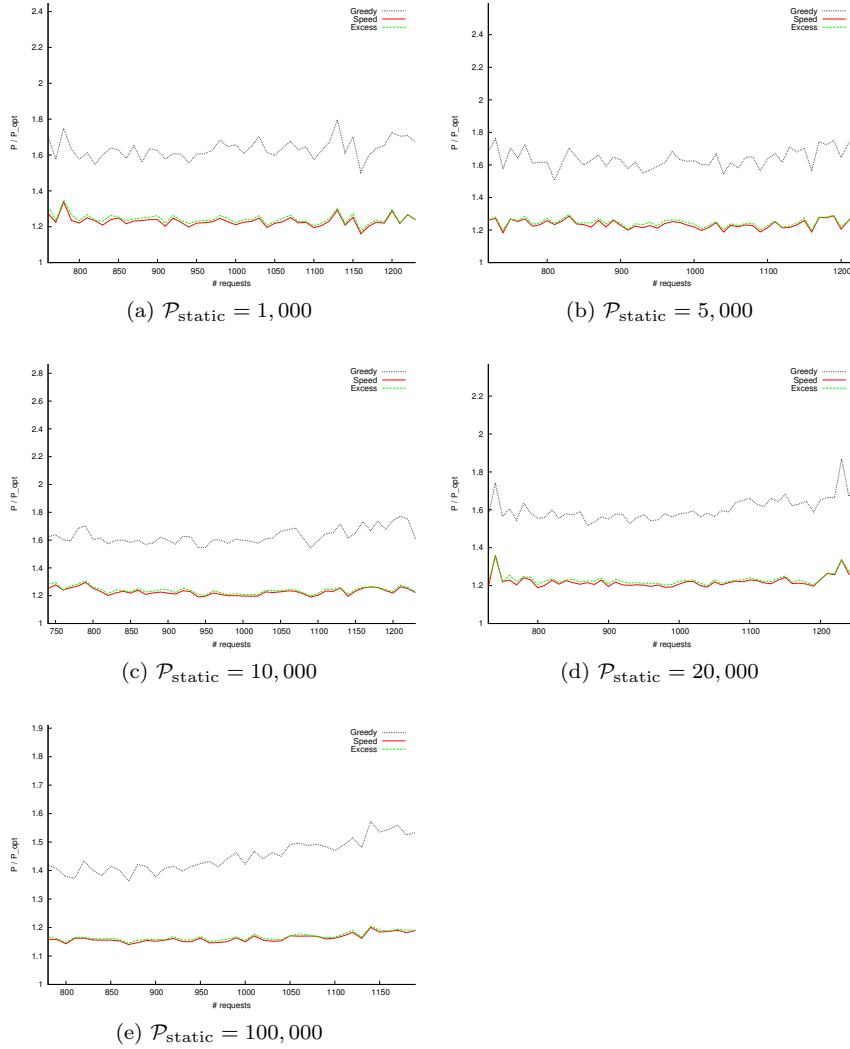


Figure 4 – Study of the impact of the total load, for trees of 20 nodes, for random requests between 0 and 100, average on 100 tests.

7 Conclusion

In this paper, we have revisited the well-known replica problem in tree networks under power constraints, in the most flexible scenario where requests of a client can be split between multiple servers. While the problem of minimizing the number of servers has polynomial complexity, we have proved that the problem of minimizing the power consumption is NP-complete, even if the servers are already placed in the tree. We assume that the server speeds can be modified using dynamic voltage and frequency scaling, depending upon the number of requests to be processed, and that a set of discrete speeds is available. Therefore, the core of the difficulty lies in assigning requests to servers in order to optimize

the speeds given to each server. Building upon the optimal solution with already placed servers and continuous speeds, we have designed efficient polynomial-time heuristics to solve the general optimization problem (deciding where to place servers and how to assign requests).

In order to assess the performance of the heuristics, we have also provided a mixed integer linear program (MILP) that returns the optimal solution of the problem for small instances (up to 30 nodes in the tree). The heuristics are always quite close to the optimal solution, and the sophisticated versions that readjust the request assignment to better fit server speeds prove to be valuable improvements of the basic greedy solution.

For future work, it would be very interesting to prove a competitive ratio for the heuristics that we have designed. However, this is quite a challenging work for arbitrary trees, and one may try to design approximation algorithms only for special tree structures, e.g. binary trees.

Acknowledgements

This work was supported in part by the French Research Agency (ANR) through the *Rescue* project. Anne Benoit and Yves Robert are with Institut Universitaire de France.

References

- [1] Guillaume Aupy and Matthieu Journault. Source code for the simulations. <http://github.com/Gaupy/replica>.
- [2] Anne Benoit, Veronika Rehn-Sonigo, and Yves Robert. Replica placement and access policies in tree networks. *IEEE Trans. Parallel and Distributed Systems*, 19(12):1614–1627, 2008.
- [3] Anne Benoit, Paul Renaud-Goud, and Yves Robert. Power-aware replica placement and update strategies in tree networks. In *Proceedings of IPDPS'2011, the 25th IEEE Int. Parallel and Distributed Processing Symposium*, May 2011.
- [4] Anantha P. Chandrakasan and Amit Sinha. Jouletrack: A web based tool for software energy profiling. In *Design Automation Conference*, pages 220–225. IEEE, 2001.
- [5] J.-J. Chen and T.-W. Kuo. Multiprocessor energy-efficient scheduling for real-time tasks. In *Proceedings of Int. Conf. on Parallel Proc. (ICPP)*, pages 13–20. IEEE, 2005.
- [6] Jian-Jia Chen. Expected energy consumption minimization in DVS systems with discrete frequencies. In *Proc. of SAC'08, Symp. on Applied Computing*, pages 1720–1725, 2008.
- [7] I. Cidon, S. Kutten, and R. Soffer. Optimal allocation of electronic content. *Computer Networks*, 40:205–218, 2002.

-
- [8] Cplex. ILOG CPLEX: High-performance software for mathematical programming and optimization. <http://www.ilog.com/products/cplex/>.
 - [9] M. R. Garey and D. S. Johnson. *Computers and Intractability, a Guide to the Theory of NP-Completeness*. W.H. Freeman and Company, 1979.
 - [10] Y. Hotta, M. Sato, H. Kimura, S. Matsuoka, T. Boku, and D. Takahashi. Profile-based optimization of power performance by using dynamic voltage scaling on a PC cluster. In *Proceedings of IPDPS, the IEEE Int. Parallel and Distributed Processing Symposium*, 2006.
 - [11] K. Kalpakis, K. Dasgupta, and O. Wolfson. Optimal placement of replicas in trees with read, write, and storage costs. *IEEE Trans. Parallel and Distributed Systems*, 12(6):628–637, 2001.
 - [12] Michael Larabel. Intel EIST SpeedStep.
 - [13] Pangfeng Liu, Yi-Fang Lin, and Jan-Jan Wu. Optimal placement of replicas in data grid environments with locality assurance. In *Int. Conf. on Parallel and Distr. Syst.*, 2006.
 - [14] Thanasis Loukopoulos, Ishfaq Ahmad, and Dimitris Papadias. An overview of data replication on the Internet. In *Proc. Int. Symp. on Parallel Architectures, Algorithms and Networks ISPAN'02*. IEEE Computer Society Press, 2002.
 - [15] M. P. Mills. The internet begins with coal. *Environment and Climate News*, 1999.
 - [16] Linwei Niu. Energy Efficient Scheduling for Real-Time Embedded Systems with QoS Guarantee. In *Proc. of RTCSA, the 16th Int. Conf. on Embedded and Real-Time Computing Systems and App.*, pages 163–172, August 2010.
 - [17] Kirk Pruhs, Rob van Stee, and Patchrawat Uthaisombut. Speed scaling of tasks with precedence constraints. *Theory of Computing Systems*, 43:67–80, 2008.
 - [18] Bing Shi and Ankur Srivastava. Thermal and Power-Aware Task Scheduling and Data Placement for Storage Centric Datacenters. In *Handbook of Energy-Aware and Green Computing, edited by S. Ranka and I. Ahmad, Vol. 1*. CRC Press, 2012.
 - [19] Jan-Jan Wu, Yi-Fang Lin, and Pangfeng Liu. Optimal replica placement in hierarchical Data Grids with locality assurance. *J. Parallel and Distributed Computing*, 68(12):1517–1538, 2008.

A Proof of Proposition 1

First let us define formally the problem with continuous speeds and already placed servers that we are solving in Proposition 1:

Definition 1 (CONTINUOUS-PLACED). Given a distribution tree (with a number of requests per client), with servers turned on and already placed on the tree, decide how to distribute client requests among them (which can also be seen as assigning the speed, equal to the number of requests, of each server), in order to minimize the total power consumption.

Note that in the following, because placed servers are all assumed turned on, we do not need to account for $\mathcal{P}_{\text{static}}$.

Definition 2. Let T be a distribution tree. Let N and N' be two nodes of T :

- N' is a *server-child* of N if there is a path from N' to N (i.e., N is an ancestor of N' in the tree), and there is a server on N' . Furthermore, we call it a *direct server-child* if it is a server-child and there are no other servers on the path from N' to N . Symmetrically, N is a server-parent of N' .
- Given a solution where N' is a child of N , and such that there is a server on N' that computes r' requests, and a server on N that computes r requests, then a *transfer of $t \in [0, r']$ requests* from N' to N is a solution where N' computes $r' - t$ requests and N computes $r + t$ requests.

Given a solution to CONTINUOUS-PLACED we consider the two following hypothesis:

Hypothesis 1. The number of requests processed by a server is never smaller than the number of requests processed by each of its server-children.

Hypothesis 2. If a server does not process all requests that are available in its subtree (i.e., some of these requests are processed higher in the tree), then the number of requests that it processes is equal to the number of requests processed by its direct server-parent.

Lemma 2. *There exists an optimal solution to CONTINUOUS-PLACED that satisfies hypothesis 1.*

Proof. Consider an optimal solution to CONTINUOUS-PLACED that does not satisfy hypothesis 1. Then there exists a server N_1 processing w_1 requests, while one of its server-children N_2 is processing $w_2 > w_1$ requests. By transferring $w_2 - w_1$ requests from N_2 to N_1 (which is possible because N_1 is a parent of N_2), we can construct a solution such that N_1 executes w_2 requests and N_2 executes w_1 requests. This solution has the same power consumption as the previous one.

Similarly we can do this in the new solution, on all nodes that violate hypothesis 1, until the solution does not violate anymore hypothesis 1. \square

Lemma 3. *If a solution to CONTINUOUS-PLACED is optimal, then it satisfies hypothesis 2.*

Proof. Let us consider a solution that does not satisfy hypothesis 2. Then there exists a server N_2 processing w_2 requests, and its direct parent-server N_1 is processing $w_1 \neq w_2$ requests. Moreover, N_2 leaves l requests to be processed higher in the tree (either by N_1 or by an ancestor of N_1).

1. if $w_1 < w_2$, then transferring $\frac{w_2-w_1}{2}$ additional requests from N_2 to N_1 is still a valid solution since N_1 is a parent of N_2 , and the power consumption of this solution is better ($2\left(\frac{w_1+w_2}{2}\right)^3 < w_1^3 + w_2^3$).
2. if $w_2 < w_1$, then it would be better to process $l' = \min\left(l, \frac{w_1-w_2}{2}\right)$ more requests at node N_2 . This is possible because these requests are processed higher in the tree and we can always exchange them with some requests of N_1 if N_1 is not processing them. In the new solution, N_2 processes $w_2 + l'$ requests, while N_1 processes $w_1 - l'$ requests. Because $l' > 0$, and by convexity, this new solution has a better power consumption.

Hence the solution is not optimal, and all optimal solutions satisfy hypothesis 2. \square

Theorem 3. *If a solution to CONTINUOUS-PLACED satisfies hypothesis 1 and 2, then it is optimal.*

To prove this theorem, let us show that

1. there is a unique solution that satisfies both hypotheses;
2. this solution is optimal.

Proposition 2. *Given two optimal solutions to CONTINUOUS-PLACED that satisfy hypothesis 1 and 2, consider a server N that computes w_1 (resp. w_2) requests in the first (resp. second) solution, and that leaves t_1 (resp. t_2) requests to be processed higher in the tree.*

- *If $w_1 > w_2$ and $t_1 \geq t_2$, then there exists a direct server-child N' of N that computes w'_1 (resp. w'_2) requests in the first (resp. second) solution, and that transfers t'_1 (resp. t'_2) requests to N , such that $w'_1 > w'_2$ and $t'_1 > t'_2$.*
- *If $w_1 \geq w_2$ and $t_1 > t_2$, then there exists a direct server-child N' of N that computes w'_1 (resp. w'_2) requests in the first (resp. second) solution, and that transfers t'_1 (resp. t'_2) requests to N , such that $w'_1 \geq w'_2$ and $t'_1 > t'_2$.*

Proof. In order to prove this proposition, note that we know that $w_i + t_i$ is equal to the sum of the requests coming directly into N , and of all requests transferred from N 's server-children. Then, because $w_1 + t_1 > w_2 + t_2$ in both cases, there exists a direct server-child of N such that $t'_1 > t'_2$. Then necessarily $t'_1 > 0$, and with hypothesis 2, $w'_1 = w_1$. Finally, with hypothesis 1, $w'_2 \leq w_2$, hence the results: if $w_2 < w_1$, then $w'_2 \leq w_2 < w_1 = w'_1$, and if $w_2 \leq w_1$, then $w'_2 \leq w_2 \leq w_1 = w'_1$. \square

Proof of Theorem 3. Uniqueness of the optimal solution: To show that there is a unique optimal solution, we proceed by induction. If there is a unique server, then there is a unique solution where this server processes all requests.

Let us now consider a tree with n servers. Let us consider two optimal solutions that satisfy hypothesis 1 and 2. Finally, let us consider w_1 (resp. w_2)

the number of requests computed by the server at the root of the tree in the first (resp. second) solution. The transfer at the root of the tree is necessary null, i.e., $t_1 = t_2 = 0$ at the root. Note that if there is no server at the root of the tree, then we can divide the tree into subtrees with a server at their roots, and by induction the two solutions are identical. If there is only one subtree with a server at its root, then we conduct the same reasoning on this subtree, and the transfer is necessary null because there are no servers upper in the tree. N denotes the server at the root of the tree.

Then, there are five different possibilities:

- If $w_1 > w_2$, with Proposition 2, we can construct an infinite sequence of servers s starting from N such that $w_1^{(s)} > w_2^{(s)}$ and $t_1^{(s)} \geq t_2^{(s)}$, because initially $t_1 = t_2$ (no transfer from the root node). Hence, we have a contradiction since there are only n servers.
- The case with $w_2 > w_1$ is symmetrical.
- If $w_1 = w_2$, and if there exists a direct server-child s of N such that $t_1^{(s)} > t_2^{(s)}$, then with hypothesis 2, $w_1^{(s)} = w_1$, and with hypothesis 1, $w_2^{(s)} \leq w_2$. Finally, with Proposition 2, we can construct an infinite sequence of servers \bar{s} starting from s such that $w_1^{(\bar{s})} \geq w_2^{(\bar{s})}$ and $t_1^{(\bar{s})} > t_2^{(\bar{s})}$. Hence we have a contradiction since there are only n servers.
- The case with $w_1 = w_2$, and there exists a direct server-child s of N such that $t_1^{(s)} < t_2^{(s)}$, is symmetrical.
- Finally, if $w_1 = w_2$, and for all direct server-child s of N , $t_1^{(s)} = t_2^{(s)}$, then we look at the sub-problem for all sub-trees of N , and by induction hypothesis, the two solutions are equal.

Finally, we have shown that two solutions to CONTINUOUS-PLACED that satisfy both hypothesis are identical.

Existence of a solution: With Lemma 2, we know that there exists an optimal solution to CONTINUOUS-PLACED that satisfies hypothesis 1. With Lemma 3, we know that this solution satisfies hypothesis 2. Therefore, it is exactly the solution that satisfies both hypothesis.

Finally, we have shown that if a solution satisfies both hypothesis, it is optimal. \square

We now introduce the algorithm ALG-CONT-PLACED that computes a distribution of the requests over a set of placed servers. First, note that if there are requests coming in a node such that there are no servers placed on that node or on any of the ancestors of that node, then there is no solution to the problem CONTINUOUS-PLACED: these requests cannot be satisfied.

Definition 3 (server-subtrees). For a given node N of a tree T , a *server-subtree* of N is a subtree of T below N such that

1. there is a server at the root N' of the subtree;
2. there are no servers on any nodes of the path of T from N to N' .

Hence we can consider trees such that there is a server on the root node (otherwise, either there is no solution to CONTINUOUS-PLACED, or we can divide a tree into the disjoint server-subtrees of the root of the original tree).

Algorithm 2: Procedure ALG-CONT-PLACED

```

1 procedure ALG-CONT-PLACED( $T$ )
2 begin
3    $\lfloor$  return sub-ALG-CONT-PLACED ( $T$ , root)
4
5 procedure sub-ALG-CONT-PLACED( $T$ ,  $N$ ) /* here  $N$  is a node of  $T$  */
6 begin
7   for  $s \in \text{Child}(N)$  do
8      $\lfloor$   $T_s =$  sub-ALG-CONT-PLACED ( $T$ ,  $s$ );
9   Let  $R$  be the number of requests that are sent by clients on internal
   nodes that are not in server-subtrees of  $N$ ;
10  Let  $S_N = R$  be the speed of the server on node  $N$ ;
11  while there exists a server-child of  $N$  with a speed greater than  $S_N$  do
12    Let  $S_{\max}$  be the maximum speed of a server-child of  $N$ ;
13    Let  $S_2$  be the second maximum speed of a server-child of  $N$ 
    ( $S_2 \neq S_{\max}$ ),  $S_2 = 0$  if there is no second maximum speed;
14    Let  $N_1, \dots, N_k$  be the  $k$  server-children of  $N$  that work at
    speed  $S_{\max}$ ;
15    Let  $S' = \frac{kS_{\max} + R}{k+1}$ ;
16    if  $S' > S_2$  then
17       $N_1, \dots, N_k$  now work at speed  $S'$ , and transfer  $S_{\max} - S'$ 
      additional requests to  $N$ ;
18       $R \leftarrow S'$  /* We will exit the while loop at the next
      step */
19    else
20       $N_1, \dots, N_k$  now work at speed  $S_2$ , and transfer  $S_{\max} - S_2$ 
      additional requests to  $N$ ;
21       $R \leftarrow R + k(S_{\max} - S_2)$  /* Note that  $R < S_2$  */
22  return  $T$ 

```

Theorem 4. *The solution returned by Algorithm 2 satisfies hypothesis 1 and 2, therefore it is optimal.*

Proof. It is easy to verify that both hypothesis are verified by the result of the algorithm. Given a node N ,

- They are verified at the beginning of the execution of the algorithm for all server subtrees of N .
- At all time during the execution of the “while” loop, they are verified for all the servers below N .
- The “while” loop ensures that both hypothesis are verified for N at the end of the execution.

□

Note that the proofs also work if the only set of speeds allowed is $[s_{\min}, s_{\max}]$. We detail in Algorithm 3 what ALG-CONT-PLACED would become.

Algorithm 3: Procedure ALG-CONT-PLACED where allowed speeds range from s_{\min} to s_{\max} .

```

1 procedure ALG-CONT-PLACED-BOUNDED( $T$ )
2 begin
3   Let  $T = \text{sub-ALG-CONT-PLACED-BOUNDED}(T, \text{root})$ ;
4   Let  $S$  be the speed of the root of  $T$ ;
5   if  $S > s_{\max}$  then
6     | return There are no solutions
7   else
8     | return  $T$ 
9
10 procedure sub-ALG-CONT-PLACED-BOUNDED( $T, N$ ) /* here  $N$  is a node
    of  $T$  */
11 begin
12   for  $s \in \text{Child}(N)$  do
13     |  $T_s = \text{sub-ALG-CONT-PLACED-BOUNDED}(T, s)$ ;
14   Let  $R$  be the number of requests that are sent by clients on internal
    nodes that are not in server-subtrees of  $N$ ;
15   Let  $S_N = R$  be the speed of the server on node  $N$ ;
16   while there exists a server-child of  $N$  with a speed greater than  $S_N$  do
17     | Let  $S_{\max}$  be the maximum speed of a server-child of  $N$ ;
18     | Let  $S_2$  be the second maximum speed of a server-child of  $N$ 
    ( $S_2 \neq S_{\max}$ ),  $S_2 = s_{\min}$  if there is no second maximum speed;
19     | Let  $N_1, \dots, N_k$  be the  $k$  server-children of  $N$  that work at
    speed  $S_{\max}$ ;
20     | Let  $S' = \frac{kS_{\max} + R}{k+1}$ ;
21     | if  $S' > S_2$  then
22       |  $N_1, \dots, N_k$  now work at speed  $S'$ , and transfer  $S_{\max} - S'$ 
    additional requests to  $N$ ;
23       |  $R \leftarrow S'$  /* We will exit the while loop at the next
    step */
24     | else
25       |  $N_1, \dots, N_k$  now work at speed  $S_2$ , and transfer  $S_{\max} - S_2$ 
    additional requests to  $N$ ;
26       |  $R \leftarrow R + k(S_{\max} - S_2)$  /* Note that  $R < S_2$  */
27   return  $T$ 

```

B Theoretical results on heuristics

In this section, we detail some theoretical results, on which the heuristics are based. Recall that the set of discrete speeds is $\{s_1, \dots, s_K\}$.

Theorem 5. *Given a distribution tree (with a number of requests per client), with servers turned on and already placed on the tree, if the maximum speed given by the solution of ALG-CONT-PLACED is greater than the maximum discrete speed, then there is no solution.*

Proof. Assume a distribution tree with servers turned on and already placed on the tree. Let us call \mathcal{S} the solution given by ALG-CONT-PLACED, such that one server works at a speed s greater than the maximum discrete speed (s_K). Since the result given by ALG-CONT-PLACED follows hypothesis 1 (Theorem 4), then in particular the server placed at the root of the tree works at speed s .

Let us assume that there exists a solution \mathcal{S}' to DISCRETE-PLACED. Then the root of the tree is processed at a speed lower than or equal to s_K . Let us consider the $s - s_K$ requests that are processed on the root in \mathcal{S} , but are not processed on the root in \mathcal{S}' . Then, in \mathcal{S}' they are processed on a server below the root:

- Either this server works at speed s in \mathcal{S} , but then we can consider the subtree rooted in this server and reach a contradiction;
- Either this server works at a speed lower than s in \mathcal{S} , and because \mathcal{S} follows hypothesis 2 (Theorem 4), it does not transfer any requests higher the tree. Hence necessarily, it cannot process some of the $s - s_K$ requests that are not processed on the root in \mathcal{S}' .

Finally, necessarily, the $s - s_K$ requests that are not processed on the root are processed on a node that works at speed $s > s_K$ in \mathcal{S} , and we have a contradiction.

In conclusion, if the maximum speed given by the solution of ALG-CONT-PLACED is greater than the maximum discrete speed, there is no solution to DISCRETE-PLACED. \square

Theorem 6. *Given a distribution tree (with a number of requests per client), with servers turned on and already placed on the tree, GREEDY is a $(1 + \frac{\max_{1 < i < K} (s_{i+1} - s_i)}{s_1})^3$ -approximation to the DISCRETE-PLACED problem.*

Proof. First, note that the optimal solution to CONTINUOUS-PLACED, where the set of speeds is bounded in the interval $[s_1, s_K]$, is a lower-bound on the optimal solution to DISCRETE-PLACED. Indeed, any solution to DISCRETE-PLACED is a solution to CONTINUOUS-PLACED.

Then, for each server N , let us call $s_N^{\{c\}}$ the speed given by the ALG-CONT-PLACED-BOUNDED with bounded speeds (Algorithm 3), and $s_N^{\{g\}}$ the speed given by the GREEDY algorithm. There exists i such that $s_N^{\{g\}} = s_{i+1}$ (with $0 \leq i < K$ and $s_0 = 0$), and by definition of $s_N^{\{g\}}$,

$$s_N^{\{g\}} \geq s_N^{\{c\}} > s_i = s_N^{\{g\}} - (s_{i+1} - s_i).$$

Furthermore, if $s_N^{\{g\}} = s_1$, then $s_N^{\{c\}} = s_N^{\{g\}}$. Therefore, if $i > 1$, then $s_N^{\{g\}} < s_N^{\{c\}} \left(1 + \frac{s_{i+1} - s_i}{s_N^{\{c\}}}\right)$, and in all cases,

$$s_N^{\{g\}} < s_N^{\{c\}} \left(1 + \frac{\max_{1 < i < K} (s_{i+1} - s_i)}{s_1}\right).$$

Finally, we can bound the total energy consumption $E^{\{g\}}$ of the GREEDY solution:

$$\begin{aligned} E^{\{g\}} &= \sum_N s_N^{\{g\}3} \leq \sum_N s_N^{\{c\}3} \left(1 + \frac{\max_{1 < i < K} (s_{i+1} - s_i)}{s_1}\right)^3 \\ &\leq E^{\{opt\}} \left(1 + \frac{\max_{1 < i < K} (s_{i+1} - s_i)}{s_1}\right)^3, \end{aligned}$$

where $E^{\{opt\}}$ is the optimal energy consumption. Indeed, $\sum_N s_N^{\{c\}3}$ corresponds to the lower bound of the optimal energy consumption in the continuous case. This concludes the proof. \square



**RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES**

Inovallée
655 avenue de l'Europe Montbonnot
38334 Saint Ismier Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399