

A Local Facility Location Algorithm for Large-Scale Distributed Systems

Denis Krivitski and Assaf Schuster

Technion – Israel Institute of Technology
(`{denisk,assaf}@cs.technion.ac.il`)

Ran Wolff

University of Maryland, Baltimore County
(`ranw@cs.umbc.edu`)

Abstract. In the *facility location problem* (FLP) we are given a set of facilities and a set of clients, each of which is to be served by one facility. The goal is to decide which subset of facilities to open, such that the clients will be served at a minimal cost.

In this paper we investigate the FLP in a setting where the cost depends on data known only to peer nodes. This setting typifies modern distributed systems: peer-to-peer file sharing networks, grid systems, and wireless sensor networks. All of them need to perform network organization, data placement, collective power management, and other tasks of this kind.

We propose a local and efficient algorithm that solves FLP in these settings. The algorithm presented here is extremely scalable, entirely decentralized, requires no routing capabilities, and is resilient to failures and changes in the data throughout its execution.

Keywords: local,distributed,data-mining,large scale

1. Introduction

The facility location problem (FLP) deals with finding an optimal subset of facilities that will be open to serve clients. The set of open facilities should minimize the cost function—the cost incurred by serving each client and that incurred by opening each facility. This clients-facilities metaphor can be used to model many practical optimization problems occurring in large scale distributed systems. Examples of such systems are peer-to-peer networks, grid computing, and wireless sensor networks.

Consider a scenario that may occur in large grids. These are used for the execution of user jobs on thousands of computers. Grid systems have many resources that are shared by many, or all, of the jobs. These include the job queue, the resource collector, and possibly large reference data files that are used by many jobs. All these resources can be replicated so as to allow better scalability, reliability, and response time. However, such replication does not come for free. It requires both

costly synchronization among replicas and the allocation of valuable resources (e.g., computers with large disks that can accommodate the reference tables). Hence, the problem of selecting resources for the different replicas, so that the system achieves best overall throughput, can also be represented as a facility location problem.

A more straight forward use of facility location may be found in data-grid systems. Data-grids are concerned with controlled sharing and management of large amounts of distributed data. In many cases, for example in AstroGrid (Page, 2001), data stored in the system needs to be mined. In addition to being distributed among thousands of computers, the data stored in those systems is constantly updated. FLP, being a close relative to clustering, can be used to extract knowledge from this data.

A more traditional application for decentralized facility location can be found in battery operated wireless sensor networks. In these networks the major challenge is to find an efficient way by which data gathered by the sensors can be routed to a control station using as little energy as possible. In order to save energy, sensors would not transmit such data directly, but rather relay it from one sensor to its neighbors until it reaches the control station. Assume a network enhanced with a few dozen high powered relay stations capable of communicating directly with the control station. Obviously these could dramatically increase the energy efficiency of the sensors. If the routers themselves are resource limited (e.g., the number of communication channels is bounded, or these relays are battery operated as well), then it would make sense to select a handful of them to be active, and shut down the rest. The optimal selection would take into account the bandwidth requirement and the battery limitation of each sensor. Thus, it is a facility location problem where the energy cost of a solution is the sum of costs for all sensors and all active relays.

The facility location problem (FLP) has been extensively studied in the last decade. Like many other optimization problems, optimal facility location is NP-Hard ((KLEINBERG et al., 1998; Jain and Vazirani, 1999)). Thus, the problem is often subjected to a hill-climbing heuristic ((Korupolu et al., 1998; Guha and Khuller, 1998; Charikar and Guha, 1999)). Hill-climbing is a simple and effective heuristic search technique in which it is assumed that a reasonable local optimum can be reasonably approached if on each search step the algorithm greedily chooses the direction that maximally decreases cost. The strength of the method is in its simplicity. It has been extensively tried for optimum search in exponential domains in problems such as genetic algorithms, clustering, etc. A rather surprising result by Arya et al. ((Arya et al., 2001)) states that, for FLP, hill climbing achieves a constant factor

approximation of 3, which means that the cost of the local optimum computed by the hill-climbing heuristic is not worse than thrice the cost of the globally optimal solution.

To the best of our knowledge, FLP has never been studied specifically in a distributed setting. Nevertheless, it is easy to see how FLP can be parallelized in the shared memory model. In addition, a related problem, the k -means clustering, which is also solved by a hill-climbing heuristic, was widely addressed in the parallel settings ((Dhillon and Modha, 2002; Foti et al., 2000; Forman and Zhang, 2000)). We note, however, that all previous work on distributed clustering assumes tight cooperation and synchronization between the peer nodes containing the data and a central node that collects the sufficient statistics needed in each step of the hill-climbing heuristic. Such central control is not practical in large-scale networks because it imposes large bandwidth requirements and is prone to errors even in the case of single failures. Even more importantly, central control is unscalable in the presence of dynamically changing data because any such change must be reported to the center, for fear it might alter the result.

In contrast, the most important features which qualify an algorithm for a large-scale distributed system are the following: the ability to efficiently scale-up (there are peer-to-peer systems today that consist of millions of peers), the ability to perform in a router-less network (critical for wireless sensor networks), and the ability to calculate the result in-network rather than to collect all of the data to a central processor (which would quickly exhaust bandwidth in both sensor and peer-to-peer networks (Gupta and Kumar, 2000)). Most important of all, because the data in a large-scale system usually changes before the computation is complete, it is crucial that the algorithm efficiently prune redundant messages and computation, as long as the data changes do not affect the global output. All these features typify *local* algorithms.

A local algorithm is one in which the complexity of computing the result does not directly depends on the number of participants. Instead, each node usually computes the result using information gathered from just few nearby neighbors. Because communication is restricted to neighbors, a local algorithm does not require message routing, performs all computation in-network, and in many cases is able to locally overcome failures and minor changes in the input (provided that these will not change its output). Local algorithms have been studied mainly in the context of graph related problems ((Awerbuch et al., 1989; Kutten and Patt-Shamir, 1997; Kutten and Peleg, 1995; Awerbuch et al., 1991; Kuhn et al., 2004; Linial, 1992; Naor and Stockmeyer, 1993)). Most recently, it has been demonstrated in (Wolff and Schuster, 2003) that local algorithms can be devised for complex data analysis

tasks, specifically, data mining of association in distributed transactional databases. The algorithm presented in (Wolff and Schuster, 2003) features local pruning of false propositions (candidates), in-network mining, asynchronous execution, as well as resilience to changes in the data and to partial failure.

In this work we develop a local algorithm that solves a specific version of FLP in which uncapacitated resources (i.e., those which can serve any number of clients) can be placed in any of the m possible locations. Initiating our algorithm from a fixed resource location, we show that the computation needed to agree on single hill-climbing step – shutting down an active resource or opening a new one at a free location – can be reduced to a group of majority votes. We then use a variation of the local majority voting algorithm presented in (Wolff and Schuster, 2003) to develop an algorithm which locally computes the exact same solution a hill-climbing algorithm would compute, had it been given the entire data. Our algorithm demonstrates that whenever the cost of a step is summed across the different sensors, peers, or resources, a hill-climbing heuristic can be computed using a local, in-network algorithm.

In a series of experiments employing networks of up to one thousand simulated processors, we prove that our algorithm has good locality, incurs reasonable communication costs, and quickly converges to the correct answer whenever the input stabilizes. We further show that when faced with constant data updates, the vast majority of sensors continue to compute the optimal solution. Most importantly, the algorithm is extremely robust to sporadic changes in the data. So long as these do not change the global result, they are pruned locally by the network.

The rest of this paper is organized as follows. We first describe our notations and formally define the problem. Then, in Section 3, we give our version for the majority voting algorithm originally described in (Wolff and Schuster, 2003). Section 4 describes a local facility location algorithm as an example of a hill climbing algorithm. In Section 5 preliminary experimental results are described. Section 6 ends the paper with some conclusions and open research problems.

2. Notations, Assumptions, and Problem Definition

A large number N of nodes are given, which can communicate with one another by sending messages. We assume that communication among neighboring nodes is reliable and ordered. This assumption can be enforced by using standard numbering, ordering and retransmission

mechanisms. For brevity, in this paper, we assume an undirected communication tree. As shown in (Yitzhak Birk, 2004), such a tree can be efficiently constructed and maintained using variations of Bellman-Ford algorithms ((Ford and Fulkerson, 1962; Jaffe and Moss, 1982)). Finally, we assume fail-stop failure and that when a node is disconnected or reconnected its neighbors are informed.

The input to the *facility location problem* consists of an *input points* database, where each point represent a client, $DB = \{p_1, p_2, \dots, p_n\}$, a set M of m possible *locations*, a cost function $d : DB \times M \rightarrow \mathbb{R}^+$, and a *configuration* (a set of open facilities) cost function $D : 2^M \rightarrow \mathbb{R}^+$. The task of a *facility location* algorithm is to find a set of *open facilities* $C \subseteq M$, such that the total cost of C and the cumulative distance of points from their nearest facility in C , $D(C) + \sum_{p_i \in DB} \min_{c \in C} d(p_i, c)$, is minimized.

To relate these definitions to the sensor networks example in the introduction, consider a database that includes a list of events that occurred in the last hour. Each event has a heuristic estimate of its importance. Furthermore, each sensor evaluates its hop distance from every relay and multiplies this by the heuristic importance of each event to produce its cost. Finally, the cost of each configuration is the number of active relays. Given this input, a facility location algorithm will compute the best combination of relays such that the most important events need not travel far before they reach the nearest relay, and not too many relays are active. The less important events, we assume, will be suppressed either in the sensor which produced them, or in-network by other sensors.

An *anytime* facility location algorithm is one which, at any given time during its operation, outputs a set of open facilities such that the cost of this ad hoc output improves with time until the optimal solution is found. A *distributed* facility location algorithm would compute the same result even though DB is partitioned into N mutually exclusive databases $\{DB^1, \dots, DB^N\}$, each of which is stored in a separate node, which is allowed to communicate with other nodes by passing messages. A *local* facility location algorithm is a distributed algorithm whose performance does not depend on N but rather corresponds to the difficulty of the problem instance at hand.

The *hill-climbing* heuristic for facility location begins from an initial set of open facilities (henceforth, *initial configuration*). Then it selects a single facility and a single empty location such that by doing one of the following: 1) moving the selected facility to this free location, 2) closing the selected facility, or 3) opening a new facility at the empty location, the cost of the solution is reduced to the largest possible

degree. If such a step exists, the algorithm changes the configuration accordingly and iterates. If any configuration which can be stepped into by closing, opening, or moving just one facility has a higher cost than the current configuration, the algorithm terminates and outputs the current configuration as the solution.

This paper presents a local, anytime algorithm which computes the hill-climbing heuristic for facility location. Note that this algorithm can easily be applied to any other hill-climbing problem. To do this, it is enough to describe the start point and the mechanism by which the next possible steps are created and their cost (or gain) evaluated. The rest of the algorithm remains the same.

3. Local Majority Voting

Our facility location algorithm reduces the problem to a large number of majority votes. In this section, we briefly describe a variation of the local majority voting algorithm from (Wolff and Schuster, 2003) which we use as the main building block for the algorithm. The algorithm assumes that messages sent between neighbors are reliable and ordered, and that node failure is reported to the node's neighbors. These assumptions can easily be enforced using standard numbering and retransmission, ordering, and heart-beat mechanisms. The algorithm requires no maximum time guarantee on message transfer or failure detection.

Given a set of nodes V , where each $u \in V$ contains a zero-one poll with c^u votes, s^u of which are one, and given the required majority $0 < \lambda < 1$, the objective of the algorithm is to decide whether $\sum_u s^u / \sum_u c^u \geq \lambda$. Equivalently, the algorithm can compute whether $\Delta = \sum_u s^u - \lambda \sum_u c^u$ is positive or negative. We call Δ the number of excess votes.

The following local algorithm decides whether $\Delta \geq 0$. Each node $u \in V$ computes the number of excess votes in its own poll $\delta^u = s^u - \lambda c^u$. Further, it stores the number of excess votes it reported to each neighbor v in δ^{uv} , and stores the number of excess votes reported to it by v in δ^{vu} . Node u computes the total number of excess votes it knows of, Δ^u , as the sum of its own excess votes and those reported to it by the set G^u of its neighbors: $\Delta^u = \delta^u + \sum_{v \in G^u} \delta^{vu}$. It also computes the number of excess votes it agreed on with every neighbor: $v \in G^u$, $\Delta^{uv} = \delta^{uv} + \delta^{vu}$. When u chooses to inform v about a change in the number of excess votes it knows of, u sets δ^{uv} to $\Delta^u - \delta^{vu}$, which results

in Δ^{uv} being equal to Δ^u . It then sends δ^{uv} to v . When u receives a message from v containing some δ , it sets δ^{vu} to δ —thus updating both Δ^{uv} and Δ^u . Finally, node u outputs that the majority is of ones if $\Delta^u \geq 0$, and of zeros otherwise.

The crux of the local majority voting algorithm is in determining when u must send a message to a neighbor v . More precisely, the question is when can sending a message be avoided, despite the fact that the local knowledge has changed. In the algorithm presented here, node u would send a message to a neighbor v in two cases: when u is initialized and when the condition $(\Delta^{uv} \geq 0 \wedge \Delta^{uv} > \Delta^u) \vee (\Delta^{uv} < 0 \wedge \Delta^{uv} < \Delta^u)$ evaluates true. Note that u must evaluate this condition upon receiving a message from a neighbor v (since this event updates Δ^u and the respective Δ^{uv}), when its input bit switches values, and when an edge connected to it fails (because then Δ^u is computed over a smaller set of edges and may change as a result). This means the algorithm is event driven and requires no form of synchronization.

The analysis in (Wolff and Schuster, 2003) reveals that the good performance of the above algorithm, in terms of message load and convergence time, stems directly from its locality. The average (as well as the worst) node would terminate after it has collected data from just a small number of nearby neighbors—its environment. The size of this environment depends on the difference, in the nearby surroundings of the node, between the average vote and majority threshold. If the two differ by as much as five percent, then we can expect the size of the environment to be limited to a few dozen nodes. It should be noted that in the worst case, that of a complete tie, all votes must be counted and the algorithm becomes global.

In order to use the local majority voting algorithm for facility location we modify it slightly. We add the ability to suspend and reactivate the vote using corresponding events. A node whose voting has been suspended will continue to receive messages and to modify the corresponding local variable, but will not send any messages. When the vote is activated, the node will always check if it is required to send a message as a result of the information received while in suspended state. Furthermore, we allow a bias towards one vote, which is equivalent to starting the vote with γ additional zeros (γ can be positive or negative). This is done by changing the condition for sending messages to $(\Delta^{uv} \geq \gamma \wedge \Delta^{uv} > \Delta^u) \vee (\Delta^{uv} < \gamma \wedge \Delta^{uv} < \Delta^u)$ and outputting a majority of ones if $\Delta^u \geq \gamma$ and of zeros otherwise. The pseudo-code of the modified algorithm is given in Algorithm 1.

Algorithm 1: Local Majority Vote

Input of node u : the local poll c^u , the local support s^u , and the set of neighbors G^u

Global constants: the majority threshold λ , the bias γ

Local variables: $\forall v \in G^u : \delta^{uv}, \delta^{vu}, active^u$.

Definitions: $\delta^u = s^u - \lambda c^u$, $\Delta^u = \delta^u + \sum_{v \in G^u} \delta^{vu}$,
 $\Delta^{uv} = \delta^{uv} + \delta^{vu}$

Initialization: $active^u = true$
 $\forall v \in G^u: \delta^{uv} = \delta^{vu} = 0$, SendMessage(v)

On activate: set $active^u \leftarrow true$

On suspend: set $active^u \leftarrow false$

On receive-message δ from $v \in G^u$: $\delta^{vu} \leftarrow \delta$

On notification of failure of $v \in G^u$: $G^u \leftarrow G^u \setminus \{v\}$

On notification of a new neighbor v : $G^u \leftarrow G^u \cup \{v\}$

On any of the above events and on change in δ^u :
 For all $v \in G^u$, if $(\Delta^{uv} \geq \gamma \wedge \Delta^{uv} > \Delta^u) \vee (\Delta^{uv} < \gamma \wedge \Delta^{uv} < \Delta^u)$
 then
 – SendMessage(v)

Procedure SendMessage(v):
 If $active^u = true$ then
 – $\delta^{uv} \leftarrow (\Delta^u - \delta^{vu})$, Send $\langle \delta^{uv} \rangle$ to v

Output of u :
 if $\Delta^u \geq \gamma$ then *positive* else *negative*

4. Majority Based Facility Location

The local facility location algorithm, which we now present, is based upon three fundamental ideas: The first is to have every node speculatively perform hill-climbing steps without waiting for a conclusive decision as to which step is globally optimal. Having performed such steps, the node continues to validate whether they agree with the globally correct ones. If there is no agreement, then these speculative steps are undone and better ones are chosen. The second idea is to choose the optimal step not by computing the cost of each step directly, but rather by voting. For each pair of possible steps, each node votes for the step it considers less costly. The third idea is a pruning technique by which many of these votes can be avoided altogether; avoiding unnecessary votes is essential because, as we further explain below, computing votes among each pair of optional steps might be arbitrarily more complicated than finding the best next step.

Key to our algorithm is the observation that the kernel problem of a hill-climbing facility location algorithm – choosing the step that reduces the cost of the solution as much as possible – is reducible to majority voting. We use this observation, together with the efficient local majority voting algorithm described in Section 3 to devise a local algorithm that computes the best among the set of possible configurations (ones reachable by moving just a single facility to a free location) and the current configuration. If the current configuration is the best possible one then it is a local minima and the algorithm makes no further steps. Otherwise, the algorithm steps to this best possible configuration and reiterates the computation.

4.1. SPECULATIVE COMPUTATION OF AN AD HOC SOLUTION

Most parallel data mining algorithms use synchronization to validate that their outcome represents the global data $\bigcup_u DB^u$. We find this approach impractical for large-scale distributed systems—specifically if one assumes that the data may change with time, making the global data impossible to determine. Instead, when performing local hill-climbing, we let each node proceed uphill whenever it has computed the best step according to the data it currently possesses. Then, we use local majority voting (as we describe next) to make sure that nodes which have taken erroneous steps will eventually be corrected. In the event that a node is corrected, a computation associated with configurations that were wrongly chosen is put on hold. These configurations are put aside in a designated cache in case additional data, accumulated later, will prove them correct after all.

We term the sequence of steps selected by node u at a given point in time its *path* through the space of possible configurations and denote it $R^u = \langle C_1^u, C_2^u, \dots, C_l^u \rangle$. C_1^u is always chosen to be the first location in M . C_l^u is the ad hoc solution of node u . u refrains from developing another configuration following C_l^u when no possible successor step has lower cost.

Since the computation of all of the configurations along every node's path is concurrent, messages sent by the algorithm contain a *context*—the configuration to which they relate. Since the computation is also speculative, it may well happen that two nodes u and v intermediately have different paths R^u and R^v . Whenever u receives a message in the context of some configuration $C \notin R^u$, this message is considered *out of context*. It is not accepted by u but rather is stored in u 's out-of-context message queue. Whenever a new configuration C enters R^u , u scans the out-of-context queue and accepts messages relating to C in the order by which they were received.

4.2. LOCALLY COMPUTING THE BEST POSSIBLE STEP

For each configuration $C_a^u \in R^u$, node u computes the best possible step as follows. First, it generates the set of possible successor configurations $Next(C_a^u)$, such that each member of $Next(C_a^u)$ adds one more location to C_a^u , removes one of C_a^u 's locations, or replaces one location in C_a^u with a location from $M \setminus C_a^u$. Next, for every $C_i, C_j \in Next(C_a^u)$, where $i < j$, node u initiates a majority vote $Majority_{C_a^u}^u \langle i, j \rangle$ which compares their costs and eventually outputs *negative* if the global cost of C_i is lower than that of C_j (as we explain below). Correctness of the majority vote process guarantees that the best configuration $C_{i_{best}} \in Next(C_a^u)$ will eventually have *negative* output for $Majority_{C_a^u}^u \langle i_{best}, j \rangle$ for all $j > i_{best}$, and *positive* output of $Majority_{C_a^u}^u \langle j, i_{best} \rangle$ for all $j < i_{best}$. Hence, the algorithm will speculatively choose C_i as the next configuration, whenever all votes indicate that C_i is better.

To determine which of two configurations has the better cost using a majority vote, we initialize $Majority_C^u \langle i, j \rangle$ with the following inputs: $s^u = \sum_{p \in DB^u} cost(p, C_i) - cost(p, C_j)$, where $cost(p, C) = \min_{f \in C} \{d(p, f)\}$, $c^u = 0$, $\lambda = 0$. Note that, as shown in (Wolff and Schuster, 2003), s^u and c^u can be set to arbitrary numbers and not just to zero or one. Further note that for every C_i, C_j the following equality holds:

$$\begin{aligned} \sum_{p \in DB} cost(p, C_i) - \sum_{p \in DB} cost(p, C_j) = \\ \sum_u \sum_{p \in DB^u} [cost(p, C_i) - cost(p, C_j)] \end{aligned}$$

Additionally, we set the bias of the vote to the difference in costs between the two configurations, $\gamma = D(C_j) - D(C_i)$. Hence, if the vote comparing the costs of C_i and C_j determines that $\Delta^u \langle i, j \rangle \geq \gamma$, then the cost of C_i is proven to be larger than the cost of C_j .

Note that since every majority vote is performed using the local algorithm described in Sec. 3, the entire computation is also local. Eventual correctness of the result and the ability to handle changes in DB^u or G^u also follow immediately from the corresponding features of the majority voting algorithm.

4.3. PRUNING THE SET OF COMPARISONS

In the above subsections we have shown how it is possible to reduce facility location to a set of majority votes. However, the reduction overshoots the objective of the algorithm. While a facility location algorithm only requires that the *best* possible successor configuration be

calculated given a certain configuration, the reduction above actually computes a *full order* on the possible successor configurations. This is problematic because for some inputs computing a full order may be arbitrarily more difficult (and thus less local) than computing only the best option. For instance, the algorithm may invest a lot of messages in deciding which of the two configurations is better, even though none of them is the best.

To overcome this problem we augment the algorithm with a pruning technique that limits the progress of comparisons such that only a small number of them actually take place. The technique we adopt is based on pivoting. First, each configuration is a candidate to be the least costly. We choose an arbitrary candidate configuration (w.l.g., the first one) as a pivot and compare all of the other candidates to it. Then, we choose the configurations which are indicated to be less costly than the pivot to be the next set of candidates and select one of them (again, the one with the lowest index) as the new pivot. As soon as the next candidate set becomes empty, the development process is stopped, and the last pivot is the least costly configuration.

Formally, given a configuration C and the set of successor configurations $Next(C)$, we define for every node u : $Pivot_i^u(C) = \min S_i^u(C)$ for $1 \geq i \geq k$, and $S_1^u(C) = \{1..|Next(C)|\}$. We define $S_i^u(C) = \{j \in S_{i-1}^u(C) \mid Majority_C^u(j, Pivot_{i-1}^u(C)).out = negative\}$ developing pivots until $S_k^u(C) = \emptyset$. Eventual correctness of all the majority votes assure that $Pivot_{k-1}^u(C)$ is the configuration with the lowest cost.

4.4. PSEUDOCODE OF THE ALGORITHM

The pseudocode of the algorithm is given in Algorithm 2. It relies on an underlying majority voting algorithm. The facility location algorithm tunnels messages to and from majority votes, removing and adding context information on the way.

5. Experiments

To evaluate the algorithm's performance, we ran it on simulated networks of up to one thousand nodes using databases of varying sizes. Our experiments test for two main properties of the algorithm: its ongoing behavior when the data is constantly altered and its dependency on the different operational parameters.

We are interested in three main metrics: the percentage of nodes which compute the exact solution at any point in time, the relative costs of the solutions computed by nodes which output a wrong solution, and

Algorithm 2: Local Facility Location

Global Constants: the set M of m possible locations

Input of processor u :

a database $DB^u = \{p_1^u, p_2^u, \dots\}$, a set of neighbors G^u , and
 the distances of points from the possible locations
 $d : DB^u \times M \rightarrow \mathbb{R}^+$

Local variables:

A vector $R^u = \langle C_1^u, \dots, C_l^u \rangle$ of configurations, where $C_i^u \subseteq M$
 A message queue $OutOfContext^u$
 For each $C \in R^u$: a vector of pairs
 $MV(C) = \langle (S_1^u(C), Pivot_1^u(C)), \dots, (S_k^u(C), Pivot_k^u(C)) \rangle$
 a set of majority votes referred to as $Majority_C^u \langle i, j \rangle$

Definitions:

$Next(C) = Swap(C) \cup Add(C) \cup Remove(C) \cup C$
 $Swap(C) = \{C \setminus \{f\} \cup \{f'\} \mid f \in C, f' \in M \setminus C\}$
 $Add(C) = \{C \cup \{f'\} \mid f' \in M \setminus C\}$
 $Remove(C) = \{C \setminus \{f\} \mid f \in C\}$

For some set of configurations N , $N[i]$ is i 's element of N
 according to lexicographic order.

For any $C_i \in R^u$ the following holds:

$C_1^u = \{M[1]\}$
 $C_{i+1}^u = Next(C_i^u)[Pivot_{k-1}^u(C_i^u)]$ for $1 \leq i < l$
 l is the minimal index for which $C_l^u = C_{l-1}^u$
 $S_1^u(C) = \{1..|Next(C)|\}$ —Set of indexes
 $\forall i > 1 : S_i^u(C) = \{j \in S_{i-1}^u(C) \mid$
 $Majority_C^u \langle j, Pivot_{i-1}^u(C) \rangle.out = negative\}$,
 k is the minimal index for which $S_k^u(C) = \emptyset$
 For $1 \leq i < k : Pivot_i^u(C) = \min S_i^u(C)$
 $ActiveSet^u = \{\langle C, i, j \rangle \mid C \in R^u, \exists q : i \in S_q^u(C), j =$
 $Pivot_q^u(C), i \neq j\}$

Init of $Majority_C^u \langle i, j \rangle$:
if $Majority_C^u \langle i, j \rangle$ exists then

 └ activate $Majority_C^u \langle i, j \rangle$
else

 └ create $Majority_C^u \langle i, j \rangle$ with inputs

 $s = \sum_{p \in DB^u} cost(p, N[i]) - cost(p, N[j])$, for $N = Next(C)$,

 $cost(p, C) = \min_{f \in C} \{d(p, f)\}$
 $c = 0, \lambda = 0, G^u = G^u$
 $\gamma = D(Next(C)[j]) - D(Next(C)[i])$

 └ tunnel to $Majority_C^u \langle i, j \rangle$ all messages in $OutOfContext^u$
 └ directed to it

Algorithm 2: Local Facility Location (cont.)

Initialization:

- | $OutOfContext^u \leftarrow \emptyset$
- | **while** R^u or $MV(C)$ for $C \in R^u$ changes **do**
 - | $\forall \langle C, i, j \rangle \in ActiveSet^u$ init $Majority_C^u \langle i, j \rangle$
 - | update R^u and $\forall C \in R^u$ update $MV(C)$
- | **On MessageSend** $\{\delta\}$ from $Majority_C^u \langle i, j \rangle$ to v :
 - | send message $\{C, i, j, \delta\}$ to v
- | **On message** $\{C, i, j, \delta\}$ from $v \in G^u$:
 - | **if** $Majority_C^u \langle i, j \rangle$ exist **then**
 - | tunnel message $\{\delta\}$ to $Majority_C^u \langle i, j \rangle$
 - | **else**
 - | enqueue $\{C, i, j, \delta\}$ in $OutOfContext^u$
- | **On change in** G^u :
 - | **foreach** existing $Majority_C^u \langle i, j \rangle$ **do**
 - | call on change in G^u for $Majority_C^u \langle i, j \rangle$
- | **On change in output of** $Majority_C^u \langle i, j \rangle$:
 - | **repeat**
 - | $OldActive \leftarrow ActiveSet^u$
 - | update R^u and $MV(C)$ for $C \in R^u$
 - | **foreach** $\langle C, i, j \rangle \in ActiveSet^u \setminus OldActive$ **do**
 - | init $Majority_C^u \langle i, j \rangle$
 - | **foreach** $\langle C, i, j \rangle \in OldActive \setminus ActiveSet^u$ **do**
 - | suspend $Majority_C^u \langle i, j \rangle$
 - | **until** $OldActive \neq ActiveSet^u$

Output of processor u : C_l^u

the communication cost for computing the solution. Ideally, most of the nodes will compute the exact solution, or else will compute a solution that is not much costlier than the exact one. Finally, all this will be done using very few messages.

The operational parameters we find most crucial are the size of the system (N), the number of data points in every local database (n/N), and the network topology. To test for the dependency of the algorithm's performance on these parameters, we ran batch mode experiments in which the data did not change during the execution. This provided a controlled experiment in which each parameter could be tested on its own.

We used a synthetic database created using the method described in (Ester et al., 1998). The data points (which represent clients) were

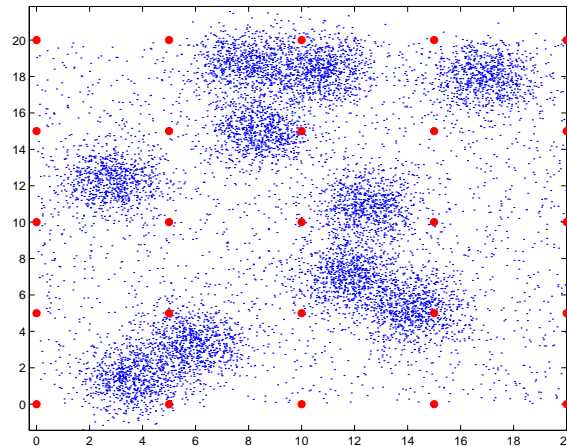


Figure 1. A typical database

generated in the 2D real domain $[0, 20] \times [0, 20]$. The source of the data points was a mixture of ten Gaussians with a random mean and variance of one, and 20% random uniform noise. Possible facility locations were placed on an equally spaced grid covering the points' domain. Figure 1 depicts a typical database with the large (red) dots signifying possible locations. Using these settings, each static experiment was run ten times with different data each time. In the batch experiments, databases of between one hundred and one thousand points (depending on the experiment) were generated for each node. In the dynamic experiments, once every few simulator clock cycles (again, depending on the experiment) a node was randomly selected, and a fraction of the points in its database were replaced with new points, sampled from the same distribution.

Finally, because the behavior of distributed algorithms may depend on network topology, we repeated our experiments for two different topologies: An Internet-like topology generated by a state-of-the-art BRITE (Medina et al., 2001) simulator and a de Bruijn topology (Kaashoek and Karger, 2003) that simulates a network with a fixed expansion rate.

5.1. ONGOING OPERATION

In this experiment we tried to realistically simulate a typical working scenario of the algorithm, in which the distribution of the data is stationary, but the data is continuously updated over time with new samples. To simulate dynamic data that retains a stationary distribution, we randomly select five percent of the nodes every five simulator cycles (about 35 times in an average edge delay). We replace 10% of

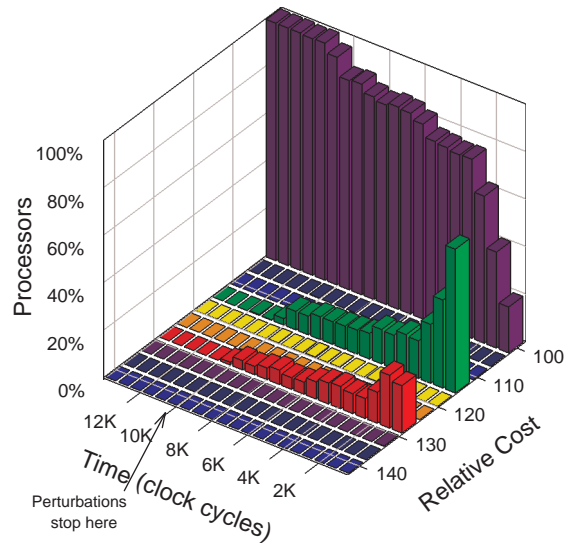


Figure 2. Evolution of solutions costs distribution in dynamic experiment is shown. At each time, a histogram of solutions costs is plotted. The cost of the optimal solution is 100. At time interval 0 - 687 simulator cycles (the first time slice in the graph), 60% of nodes output a solution of cost 113 to 118, and 20% of nodes output a solution of cost 127 to 131. Leaving only 20% of the nodes at the optimal solution. As time passes more nodes converge to the optimal solutions. At time interval 0 - 10,000, node's data undergoes continuous perturbations. From time 3000 to 10,000 the system is at steady state, where perturbations prevent from all nodes to converge to the optimal solution. As perturbations stop, the system rapidly converges to the optimum. Simulation parameters: Internet topology, $N = 512$, $n/N = 1000$, and $m = 25$.

each selected node's data points with new points selected from the same distribution. We keep changing the database this way for 10,000 simulator cycles and then stop the changes in order to let the algorithm converge to the exact result.

As the results in Figure 2 show, during the period of data changes, more than 80% of the nodes manage to quickly compute the optimal solution. Of the nodes that compute a different result, most compute one that is about 15% more costly than the optimal (note that costs here are normalized). When at cycle 10,000 the changes stop, all of the nodes immediately converge to the correct result. Similar results were computed for different amounts of data perturbation.

5.2. COMMUNICATION COST

We evaluated the communication cost of the algorithm by counting the number of messages sent by the average node during an entire bulk

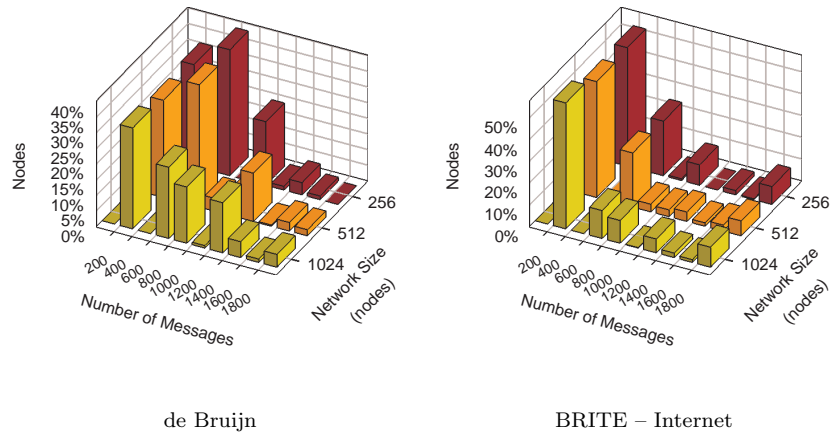


Figure 3. Distribution of number of messages sent by nodes for three network sizes and two topologies. Last bars represent all nodes which sent more than 1800 messages. The graph shows that as the network size grows, messages distribution remains essentially the same.

run. We were especially interested in the scale-up of the algorithm (i.e., the effect of increasing N on the message cost) and in the effect of different network topologies on the message cost. We ran experiments with 256, 512, and 1024 nodes using both a BRITE generated Internet topology and a de Bruijn topology. For each combination we generated ten different databases. We let the algorithm run through, counted the number of messages sent by each node, and then reported the averaged histogram for each topology and N .

The result, as reported in Figure 3, shows some interesting trends. First, about half the nodes use 200 to 400 messages throughout the algorithm. This seems a very reasonable number, considering that each of the algorithm's messages takes a few dozens bytes. Being small, these messages could be buffered. Second, the algorithm scales-up well, with no real increase in costs. Last, the algorithm requires fewer messages in an Internet topology, which can be explained by its superior mixing power.

5.3. LOCALITY

The next set of experiments measures the size of each node's environment directly. Node u 's environment is the set of neighboring nodes from which u gathers data. The size of the environment is important because the algorithm's performance strongly depends on it. For each network topology we again ran multiple experiments with different N . In each experiment, we counted the number of majority votes which

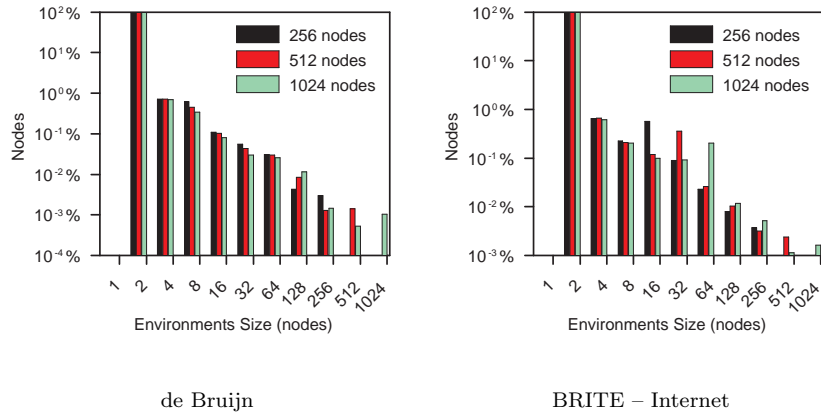


Figure 5. A histogram of maximum environment sizes is depicted for three network sizes and two network topologies (de Bruijn and BRITE). Each bar represents the maximum number of nodes over 10 experiments. These graphs show that even in the worst case analysis, the vast majority of nodes have small environments. Simulation parameters: $m = 25$, $n/N = 1000$.

5.4. ROBUSTNESS OF OUR RESULTS

One question that is often asked about average (or worst case) results is how robust they are. In Figure 6 we depict, for one of the experiments above, the range of results that were calculated. As can be seen, the difference between the minimal and maximal number of messages is not significant. The difference between the minimal and maximal environment size for the average node can be quite large because in certain instances of the problem there happen to be no ties at all.

5.5. SIZE OF DATABASE

In the last experiment, we measured the effect of the size of the local database (n/N) on the locality and message load of the algorithm. n/N is dependent mostly on the characteristics of the domain, e.g., the number of files on the disk of an average e-Mule peer or the buffer size of a sensor. We expect that as this number increases, the performance of the algorithm will improve because local statistics will become more accurate and more representative of the global ones.

We varied n/N from one thousand points down to one hundred. Figure 7 depicts the number of messages sent by each node and Figure ?? depicts the sizes of the average and largest environments. As can be seen, the same trends we discussed above persist when n/N equals 1000, 500, and 250. For n/N equals 100, the number of messages and

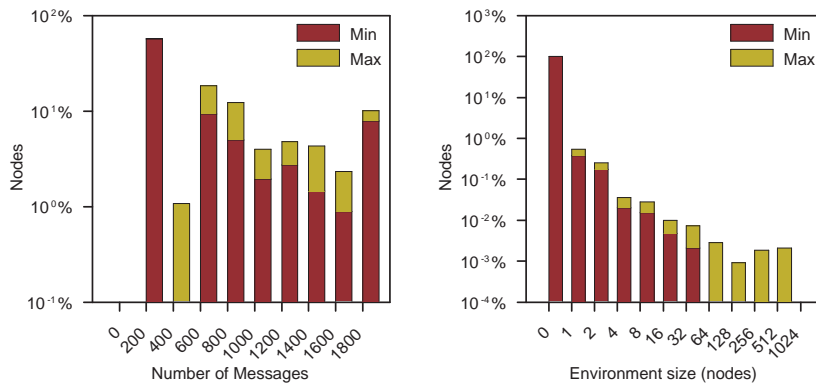


Figure 6. Comparison of best case and worst case results when using different databases. The small differences suggest that our results are robust. Simulation parameters: $m = 25$, $n/N = 1000$, $N = 1024$, BRITE-Internet topology.

average environment size grow significantly. However, since these experiments exhibit local behavior as well, we assume that, had we run our experiments with larger N values, the trend would have been visible for smaller n/N as well. Unfortunately, the performance of our simulator restricts us at this phase to $N = 1024$.¹

6. Conclusions and Further Research

We have described a new facility location algorithm suitable for large scale distributed systems. The characteristics which qualify the algorithm for systems of this type are its message efficiency, its strong local pruning, and its ability to efficiently sustain failures and changes in the input. All these qualities stem from the algorithm's *local* nature.

Besides its immediate value, the algorithm serves to demonstrate that various data mining problems can be solved in large scale distributed settings through reduction to basic primitives like majority-vote. These primitives can later be solved by efficient local algorithms. We believe that in-network data mining may very well become one of the key techniques for accessing the output of these systems.

¹ The slight discrepancies between the results for $n/N = 1000$ in Figure 7 and Figures 3 and 4 are due to different m values used in those experiments.

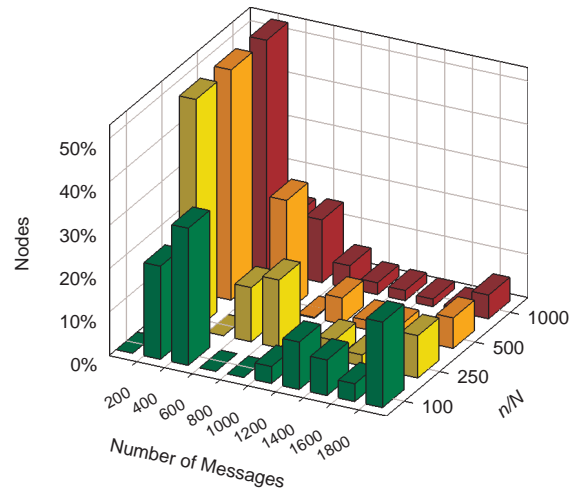


Figure 7. Effect of n/N (the number of clients on each node) on the number of messages. The graph shows four message distributions for four different n/N values. As n/N grows the algorithm sends less messages. Simulation parameters: $m = 20$, $N = 512$, BRITE—Internet topology.

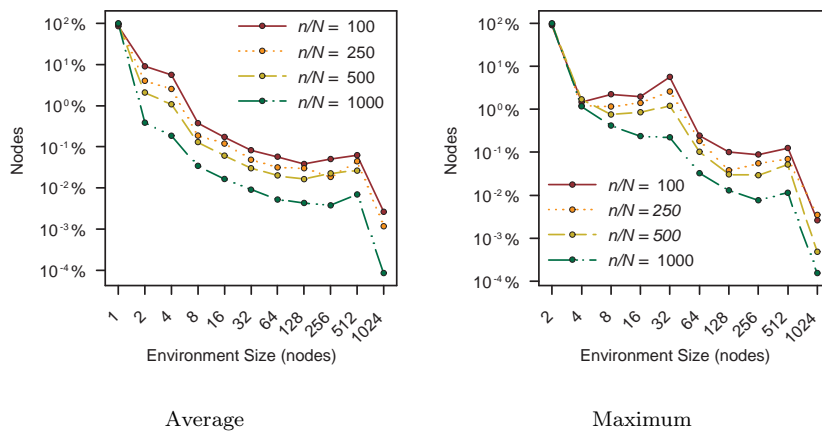


Figure 8. Effect of n/N (the number of clients on each node) on the environment size. Two graphs show the maximum and average number of messages histograms for 4 different n/N values. As n/N grows, environments (both maximum and average) become smaller, and therefore, the algorithm becomes more and more local. Simulation parameters: $m = 20$, $N = 512$, BRITE—Internet topology.

References

- Arya, V., N. Garg, R. Khandekar, K. Munagala, and V. Pandit: 2001, 'Local search heuristic for k-median and facility location problems'. In: *STOC '01: Proceedings of the thirty-third annual ACM symposium on Theory of computing*. New York, NY, USA, pp. 21–29, ACM Press.
- Awerbuch, B., A. Bar-Noy, N. Linial, and D. Peleg: 1989, 'Compact distributed data structures for adaptive routing'. In: *STOC '89: Proceedings of the twenty-first annual ACM symposium on Theory of computing*. New York, NY, USA, pp. 479–489, ACM Press.
- Awerbuch, B., B. Patt-Shamir, and G. Varghese: 1991, 'Self-stabilization by local checking and correction (extended abstract)'. In: *Proceedings of the 32nd annual symposium on Foundations of computer science*. Los Alamitos, CA, USA, pp. 268–277, IEEE Computer Society Press.
- Charikar, M. and S. Guha: 1999, 'Improved Combinatorial Algorithms for the Facility Location and k -Median Problems'. In: *IEEE Symposium on Foundations of Computer Science*. pp. 378–388.
- Dhillon, I. S. and D. S. Modha: 2002, 'A Data-Clustering Algorithm on Distributed Memory Multiprocessors'. In: *Large-Scale Parallel Data Mining*. pp. 245–260.
- Ester, M., H. Kriegel, J. Sander, M. Wimmer, and X. Xu: 1998, 'Incremental Clustering for Mining in a Data Warehousing Environment'. In: *VLDB*. pp. 323–333.
- Ford, L. and D. Fulkerson: 1962, *Flows in Networks*. Princeton University Press.
- Forman, G. and B. Zhang: 2000, 'Distributed data clustering can be efficient and exact'. *SIGKDD Explor. Newsl.* **2**(2), 34–38.
- Foti, D., D. Lipari, C. Pizzuti, and D. Talia: 2000, 'Scalable Parallel Clustering for Data Mining on Multicomputers'. In: *IPDPS '00: Proceedings of the 15 IPDPS 2000 Workshops on Parallel and Distributed Processing*. London, UK, pp. 390–398, Springer-Verlag.
- Guha and Khuller: 1998, 'Greedy Strikes Back: Improved Facility Location Algorithms'. In: *SODA: ACM-SIAM Symposium on Discrete Algorithms (A Conference on Theoretical and Experimental Analysis of Discrete Algorithms)*.
- Gupta, P. and P. R. Kumar: 2000, 'The capacity of wireless networks'. *IEEE Transactions on Information Theory* **46**(2), 388 – 404.
- Jaffe, J. and F. Moss: 1982, 'A Responsive Routing Algorithm for Computer Networks'. *IEEE Transactions on Communications* pp. 1758–1762.
- Jain, K. and V. V. Vazirani: 1999, 'Primal-Dual Approximation Algorithms for Metric Facility Location and k -Median Problems'. In: *IEEE Symposium on Foundations of Computer Science*. pp. 2–13.
- Kaashoek, F. and D. Karger: 2003, 'Koorde: A simple degree-optimal distributed hash table'. *Peer-to-Peer Systems II: Second International Workshop*.
- KLEINBERG, J., C. PAPANITRIOU, and P. RAGHAVAN: 1998, 'A Microeconomic View of Data Mining'. *Data Mining and Knowledge Discovery*.
- Korupolu, M. R., C. G. Plaxton, and R. Rajaraman: 1998, 'Analysis of a local search heuristic for facility location problems'. In: *Proceedings of the ninth annual ACM-SIAM symposium on Discrete algorithms*. pp. 1–10, Society for Industrial and Applied Mathematics.
- Kuhn, F., T. Moscibroda, and R. Wattenhofer: 2004, 'What cannot be computed locally!'. In: *PODC '04: Proceedings of the twenty-third annual ACM symposium on Principles of distributed computing*. New York, NY, USA, pp. 300–309, ACM Press.

- Kutten, S. and B. Patt-Shamir: 1997, 'Time-adaptive self stabilization'. In: *PODC '97: Proceedings of the sixteenth annual ACM symposium on Principles of distributed computing*. New York, NY, USA, pp. 149–158, ACM Press.
- Kutten, S. and D. Peleg: 1995, 'Fault-local distributed mending (extended abstract)'. In: *PODC '95: Proceedings of the fourteenth annual ACM symposium on Principles of distributed computing*. New York, NY, USA, pp. 20–27, ACM Press.
- Linial, N.: 1992, 'Locality in distributed graph algorithms'. *SIAM J. Comput.* **21**(1), 193–201.
- Medina, A., A. Lakhina, I. Matta, and J. Byers: 2001, 'BRITE: Universal Topology Generation from a User's Perspective'. Technical report, Boston, MA, USA.
- Naor, M. and L. Stockmeyer: 1993, 'What can be computed locally?'. In: *STOC '93: Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*. New York, NY, USA, pp. 184–193, ACM Press.
- Page, C.: 2001, 'Astrogrid and data mining'. In: *Proc. SPIE Vol. 4477, p. 53-60, Astronomical Data Analysis, Jean-Luc Starck; Fionn D. Murtagh; Eds.* pp. 53–60.
- Wolff, R. and A. Schuster: 2003, 'Association Rule Mining in Peer-to-Peer Systems'. In: *ICDM '03: Proceedings of the Third IEEE International Conference on Data Mining*. Washington, DC, USA, p. 363, IEEE Computer Society.
- Yitzhak Birk, Liran Liss, A. S. R. W.: 2004, 'A Local Algorithm for Ad Hoc Majority Voting via Charge Fusion'. In: *DISC*.