

Gelly-Scheduling: Distributed Graph Processing for Service Placement in Community Networks

This document merely serves the purpose of timely dissemination. Copyrights belong to original holders.

Miguel E. Coimbra
INESC-ID/IST, Universidade de Lisboa
Lisbon, Portugal
miguel.e.coimbra@tecnico.ulisboa.pt

Mennan Selimi
University of Cambridge
Cambridge, UK
ms2382@cam.ac.uk

Alexandre P. Francisco
INESC-ID/IST, Universidade de Lisboa
Lisbon, Portugal
aplf@inesc-id.pt

Felix Freitag
Universitat Politècnica de Catalunya
Barcelona, Spain
felix@ac.upc.edu

Luís Veiga
INESC-ID/IST, Universidade de Lisboa
Lisbon, Portugal
luis.veiga@inesc-id.pt

Abstract

Community networks (CNs) have seen an increase in the last fifteen years. Their members contact nodes which operate Internet proxies, web servers, user file storage and video streaming services, to name a few. Detecting communities of nodes with properties (such as co-location) and assessing node eligibility for service placement is thus a key-factor in optimizing the experience of users. We present a novel solution for the problem of service placement as a two-phase approach, based on: 1) community finding using a scalable graph label propagation technique and 2) a decentralized election procedure to address the multi-objective challenge of optimizing service placement in CNs. Herein we: *i)* highlight the applicability of leader election heuristics which are important for service placement in community networks and scheduler-dependent scenarios; *ii)* present a parallel and distributed solution designed as a scalable alternative for the problem of service placement, which has mostly seen computational approaches based on centralization and sequential execution.

1 Introduction

Community networks (CNs) are owned and managed by volunteers and offer various services to their members. Seamless computing and service sharing in CNs have gained momentum due to the emerging technology of CN micro-clouds. One such network is guifi.net, located in the Catalonia region of Spain. It is a successful example of this paradigm. Guifi.net is defined as an open, free and neutral CN built by its members pooling resources. Guifi.net was born in 2004, and until today, has grown into a network of more than 34,000 operational nodes. Previous work on guifi.net classified services into network-oriented and user-oriented. For these two types in the Catalonia region, the three most prevalent occurrences were [15]: *a) network-oriented services (558 in this region)* – network graph-servers (39.24%), DNS servers (35.48%) and NTP servers (17.20%); *b) user-oriented services (514 in this region)* – proxy servers for Internet access (53.50%), web pages (11.08%)

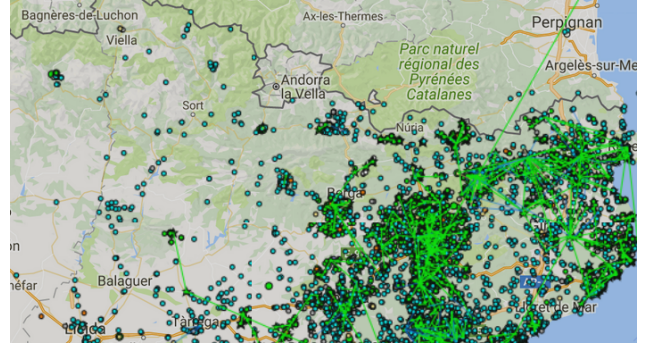


Figure 1. Depiction of guifi.net’s Osona region.

and communication applications such as VoIP, audio, video and instant messaging (9.33%). Nodes in guifi.net are exclusive to specific geographical zones (there are no overlays) such as what is depicted in Figure 1. There are special-purpose nodes called graph-servers, which are responsible for performing network measurements between nodes and have an API for querying node states [14]. These graph-servers comprise a distributed hierarchical monitoring system which records the network’s link data traffic properties. Guifi.net is thus a relevant testbed for developing and validating techniques to enhance service placement and system scheduling by exploring their requirement of leader election. In turn, these may be extrapolated to more complex scenarios, such as placement in P2P networks (typically irregular), industrial contexts and IoT scenarios. A simple web proxy would most likely have node latency as its most relevant parameter. On the other hand, a mission-critical quality-of-service proxy could place the focus on node availability. Heuristics may encompass network features such as topology, as well as domain attributes (such as availability and quality of specific resources). While one may intuitively define one heuristic as absolute, this could produce scenarios which are locally optimal but globally undesirable. What if the node with the highest availability happens

to be on the outer rims of the network? Aspects of network topology are as relevant for system efficiency as the service-level heuristics which traditionally guide leader election for placements.

Our objective is to devise an efficient, scalable solution which is easy to fine-tune regarding domain-specific attributes, and that provides seamless scalability for increasing network size and number of services. For this, we propose a platform that enables incremental processing in a scenario where information continuously arrives: changes in network, node and service quality are continuously monitored. Our solution is a two-phased approach which optimizes the definition of communities (**Phase One**) and election of leaders (**Phase Two**) in a community communications network. The paper is organized as follows. Section 2 explains the two main phases of our algorithm. Section 3 details our evaluation methodology and obtained results. Section 4 highlights relevant studies on community networks and service placement. Section 5 summarizes our contribution’s highlights.

2 Gelly-Scheduling Service Placement

The challenges inherent to service placement for large scale geo-distributed networks (such as community networks) are usually addressed in the literature with a batch-oriented non-scalable approach. The typical approach consists of performing a search (exhaustive or via heuristics) in a centralized computing unit. All the information about network links and nodes is centrally and sequentially processed, in order to determine the best network configurations as far as service placement is concerned. While the unit responsible for this search may benefit from hardware improvements, they are merely a form of vertical scaling (which is limited). This approach does not prioritize reaction to changes in the network and its nodes, in order to make service placement more dynamic in a context of continuous monitoring. It also doesn’t scale in the context of larger networks. We present a novel method capable of both achieving scale-out processing for optimizing community network topology as well as electing service placement targets within communities in a decentralized approach. We employ community detection as a parallel technique which enables the partitioning of the problem space to optimize node placement in communities. This allows for an efficient leader election to execute concurrently (each community being responsible for its leader) and in parallel within each community. This work aims to improve service placement for networks in a way that users and processing tasks are balanced regarding bandwidth restrictions and data sources.

Phase One: Community Finding. We use two definitions of community: *default* – the zone-based node distributions, provided in the dataset as-is (insignificant preprocessing is performed in this case); *custom* – a state-of-the-art label propagation technique [11] applied for detecting communities. We build an undirected graph $G = (V, E)$ by defining a set of n

nodes V and a set of m edges E such that an edge $e \in E$ will be created if and only if there is a corresponding link element between two working devices (each belonging to a working node) in the dataset. Single-leaf nodes were discarded as part of preprocessing. The goal of this phase is to rapidly partition the problem space into a configuration that promotes scalability of computation and efficient resource usage. We provide the pseudo-code for the most relevant actions of Phase One in Algorithm 1, where C is an upper bound on the number of iterations to execute (a default limit of $C = 10$ iterations is common in the literature for convergence [1]). Phase One thus becomes an important instrument in efficiently defining groups of network nodes by employing a state-of-the-art technique in community detection. These groups aid the optimization process of service placement, effectively serving as a useful blueprint for Phase Two of our algorithm. The two phases form a technique to harness current platforms and infrastructures to tackle service placement. Conceptually, there is a top-tier master node which is responsible for: 1) querying the graph-servers for all of the network’s node information; 2) executing Phase One of our algorithm to obtain a definition of communities; 3) informing each node of its community’s composition. This is depicted in Figure 2, where in the middle there is a centralized entity consisting of one or (potentially many) more computational workers. It initially queries graph-servers (or whatever network visibility mechanisms are in place) to obtain a snapshot of the network’s nodes. Then it executes Phase One of our algorithm, decomposing the network into communities. A major computational advantage of Phase One is that this master can be a single machine or a set of workers in a cluster, effectively scaling with the computational capability

Algorithm 1 Phase One: Community Finding

```

1: INPUT:  $G = (V, E), C = 10$ 
2: OUTPUT:  $Z \rightarrow$  Set of graphs representing communities
3: for all  $v \in V$  do
4:    $v.generateUniqueLabel()$ 
5: end for
6:  $G' \leftarrow G.setUndirectedEdges()$ 
7:  $i \leftarrow 1$ 
8: for  $i < C$  do
9:   for all  $v \in V$  do
10:     $M \leftarrow v.getInboundMessages()$ 
11:     $L \leftarrow M.getMostFrequentLabels()$ 
12:     $v.updateLabel(L.filterHighestLabel())$ 
13:   end for
14:    $i \leftarrow i + 1$ 
15:   if not  $G'.labelsChanged()$  then
16:     break
17:   end if
18: end for
19: return  $Z \leftarrow G'.groupByLabels()$ 

```

Table 1. Frequency of per-node device count categories. The most frequent services are Internet proxies, a consequence of guifi.net existing as an alternative to the standard ISP model.

Nodes	23,468 (100%)	β_3
i) Strong	337 (1.436%)	1
ii) Medium	1,666 (7.099%)	0.5
iii) Weak	21,465 (91.465%)	0.1

available to this top-tier master. Each community member is then informed of the elements of its own community: required to proceed to Phase Two.

Phase Two: Leader Election. Phase Two receives a set of communities and elects a leader for each one. This election phase is self-contained for each community, in the sense that a distributed implementation of this phase can be carried out concurrently with respect to communities and in parallel within each community with our graph-based approach. The right-side of Figure 2 illustrates this. There may be more than one connected component in geographical zones of guifi.net. Due to this, for every community network G , only the nodes belonging to the largest connected component of G are used to choose a leader for service placement. This election consists of Phase Two of our algorithm and is detailed in Algorithm 2. This phase serves the purpose of identifying the best node for service placement. Leadership is attributed through a scoring, where the score of each node i lies in defining a linear combination of two sets of heuristics. One set is based on system-centric values: availability β_1 and latency β_2 as defined by graph-servers [2], as well as computational class β_3 as per Table 1 and defined as part of this work; the other is calculated as part of this algorithm and consists of betweenness α_1 and closeness α_2 centralities. We defined heuristic β_3 as a score in three computational categories for nodes: i) server-type nodes which typically have stronger computational power to support more demanding services; ii) non-server nodes with more than one device; iii) non-server nodes with a single device. Table 1 shows the representation of β_3 for each category for the data we analyzed. The values we attribute to β_3 were selected arbitrarily to represent computational power of a given node. This categorization serves the purpose of approximating realistic tiers of computational capabilities for nodes in the network – information which, as far as the authors know, is not readily-available in the guifi.net CNML dataset. Thus, as an example, the initial score of a node i will be defined as:

$$s_i = w_1 \alpha_1 + w_2 \alpha_2 + w_3 \beta_1 + w_4 \beta_2 + w_5 \beta_3 \quad (1)$$

Table 2 details the specifics of each heuristic, namely their meaning and how they are obtained. Notation-wise, i is the node to be scored while u and v represent arbitrary nodes in the community graph G with n nodes, $\sigma_{u,v}$ is the number shortest paths from u to v , $\sigma_{u,v}(i)$ is the number of those that pass through i , and $d(i, v)$ is the geodesic distance between i

Algorithm 2 Phase Two: Leader Election

```

1: INPUT:  $G=(V, E), W$  ▷ Heuristic weight array
2: OUTPUT:  $R$  ▷ Decreasing-order ranks
3:  $\alpha \leftarrow [], \beta \leftarrow []$ 
4: for all  $v \in V$  do
5:    $\alpha[v] \leftarrow v.calculateAlphas()$ 
6:    $\beta[v] \leftarrow v.getBetas()$ 
7:    $v.setScore(W * [\alpha[v] \ \beta[v] ])$ 
8: end for
9:  $R \leftarrow G.getVertices().orderByScore()$ 
10: return  $R$ 

```

and v . Phase Two was designed under two types of evaluation based on configuration of heuristics: **Absolute heuristics** – in this case, leader selection is guided exclusively by exactly one of the heuristics. We analyze the impact of each individual heuristic, setting the weights of others to zero. **Combined heuristics** – we consider a linear combination of two heuristics. We set unbalanced weights in order to better determine the more significant contributions, in the sense that for two heuristics m_1 and m_2 , we may define the node score to be $v_s = (1 - f)m_1 + fm_2$, or the reverse. If one heuristic weights in for 60% of the score, the other will account for the remaining 40%.

Table 2. Algorithm’s heuristic symbols and meanings.

	Betweenness Centrality
α_1	$\sum_{u \neq i \neq v} \frac{\sigma_{uv}(i)}{\sigma_{uv}}$, fraction of shortest paths from u to v , for all nodes u and v , passing through node i .
	Closeness Centrality [9]
α_2	$(n - 1) / \sum_v d(i, v)$, where $d(v, i)$ is the geodesic distance from node i to node v .
	Availability
β_1	Percentage of ping responses received by a graph-server (%) over a specific time period.
	Latency
β_2	Ping response timing, measured by a graph-server (ms) over a specific time period.
	Computational Class
β_3	Defined by the number of devices handled by the node, as well as its role.

2.1 Implementation

Scores of heuristics α_1 and α_2 were obtained for each community G using the Python NetworkX library, for use in Phase Two. Overall, the time to calculate them is negligible when compared to the total amount of time required to compute Phase One plus Phase Two. There are common aspects to generating samples for bandwidth and round-trip time, but each was based on different statistical artifices.

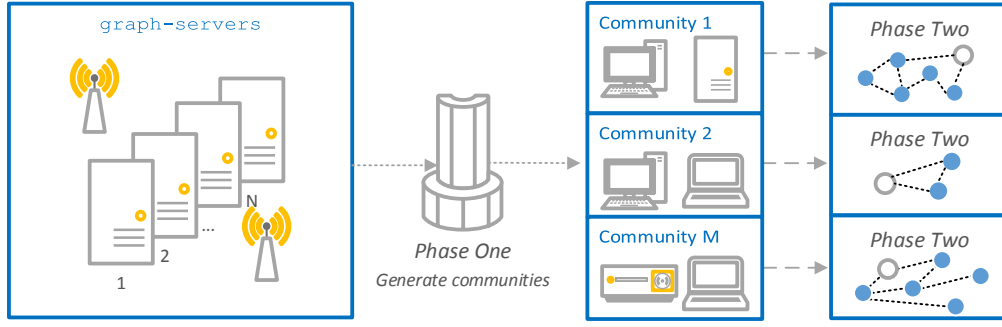


Figure 2. The graph-servers on the left send the network heuristics to the master node; the master node in the middle decides on the community configuration through Phase One; communities concurrently elect an internal leader during Phase Two.

Bandwidth. Let $BW \sim K(k, h, \xi, \alpha)$ represent the empirical bandwidth distribution. K stands for the four-parameter Kappa distribution [5], where k and h denote the shape of the distribution, ξ denotes its location and α is a scaling factor. These four parameters were estimated using L-moment statistics, namely through the `lmoms` function which computes the sample L-moments and the `parkap` function which estimates the four parameters of K based on the sample L-moments. Both functions are part of the R `lmomco` library. The four-parameter Kappa distribution is used for simulating additional samples based on the empirical distribution made by using the `rkappa4` function of the R `FAdist` library for random generation purposes.

Round-trip time. Let $RTT \sim GEV(\mu, \sigma, \xi)$ represent the empirical round-trip time distribution. GEV is the generalized extreme value distribution. It has four parameters: μ , which is the location of the distribution, σ which represents the scale and ξ which represents the shape of GEV (influencing the behavior of the distribution tail). The previously-referenced `lmoms` function was used as well, with `pargev` now being the function (also present in library `lmomco`) responsible for estimating the GEV parameters based on the sample L-moments. Additional round-trip time samples were simulated using the `rgev` function. GEV exists as a family of continuous probability distributions, stemming from extreme value theory [3].

The method of L-moments is used to understand insights of analyzed data and to estimate distributions [6, 7] using efficient techniques [4]. Figure 3 shows the \log_{10} plot of the bandwidth, while Figure 4 shows the same for round-trip time.

3 Experimental Evaluation

There are 23,391 nodes identified as working and, for the whole guifi.net, there are 878 nodes defined as servers. This implies that, at most, 3.75% of the working nodes could actually be sustaining full fledged services. We believe guifi.net, while it is in fact an open community network, has a type of topology which allows for extrapolating results into other sorts

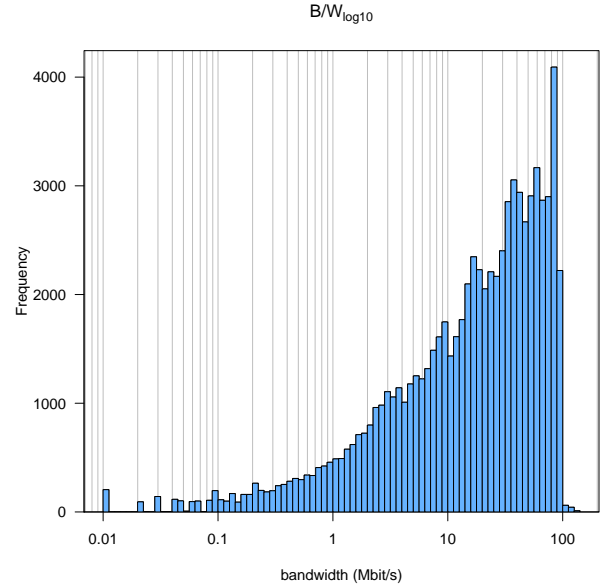


Figure 3. Bandwidth (B/W) in logarithmic scale.

of networks. This claim is made based on previous research work in the literature [14], which both analyzed the impact of prioritizing different heuristics on the computational and network resources available [16] and studied practical issues with micro-service architectures [13]. We used available statistical processing tools to attempt to fit several distributions and compare them. For the bandwidth, we present plots of the distribution fitting and Empirical Cumulative Distribution Function in Figures 5 and 6. In the same order, we also present the aforementioned plots for round-trip time in Figures 7 and 8.

We then modeled the ECDF of both network properties with the use of the `lmomco` and `FAdist` libraries in R.

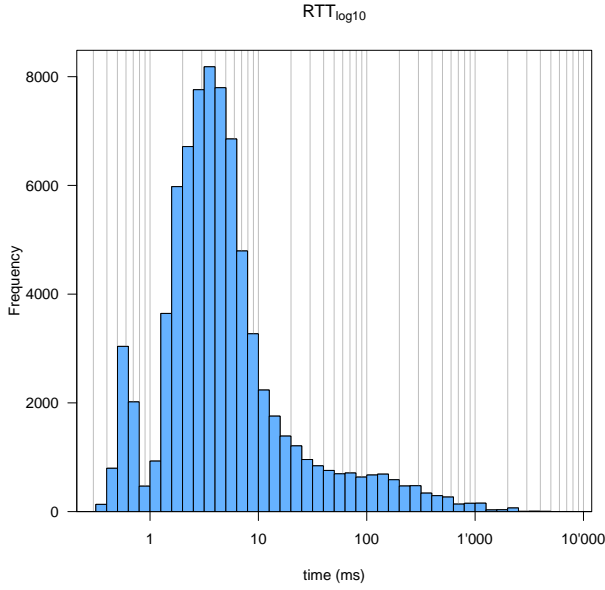


Figure 4. Round-trip time (RTT) in logarithmic scale.

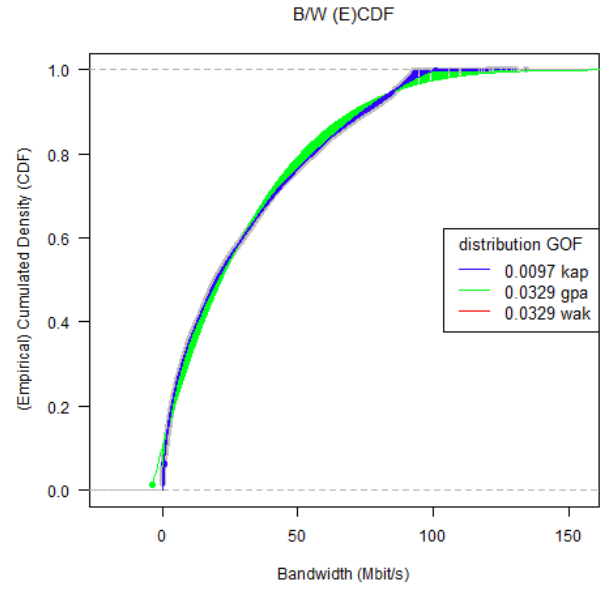


Figure 6. Bandwidth (B/W) Top 3 Empirical Cumulative Distribution Function (ECDF).

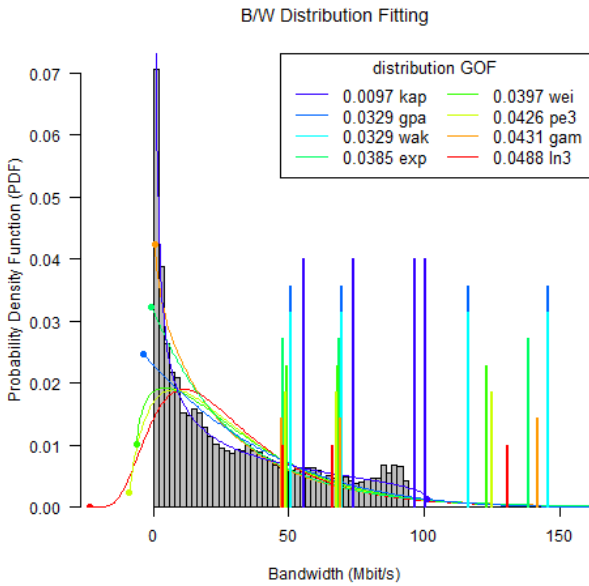


Figure 5. Bandwidth (B/W) observation comparison and goodness-of-fit of different candidate distributions. Lower goodness-of-fit is better.

Network Characteristics. Part of guifi.net exists as an instance of the *Quick Mesh Project* (QMP¹), a system for easily deploying MESH/MANET networks using Wi-Fi technology. QMP is an urban mesh network in Barcelona and it is a subset

¹<http://qmp.cat/>

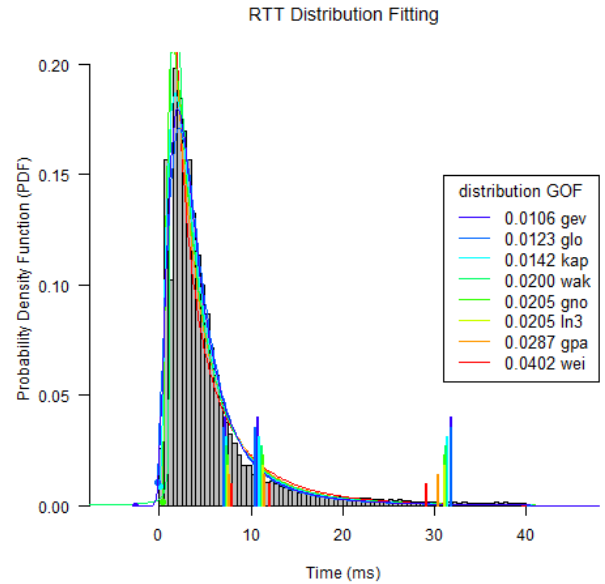


Figure 7. Round-trip time (RTT) observation comparison and goodness-of-fit of different candidate distributions. Lower goodness-of-fit is better.

of the guifi.net community network sometimes called Sants-UPC network. It was designed for use in scenarios such as free community networks, of which guifi.net is a rich example [13]. We use measurements of round-trip time (RTT) and bandwidth

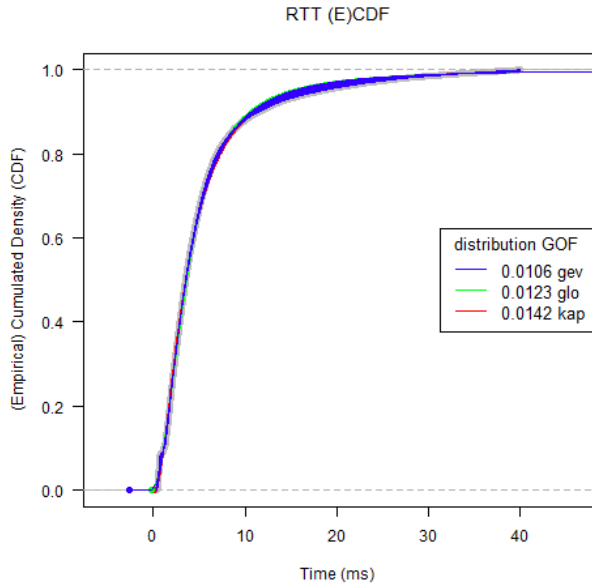


Figure 8. Round-trip time (RTT) Top 3 Empirical Cumulative Distribution Function (ECDF).

(B/W) from the Sants-UPC wireless mesh QMP instance to establish a model of these telecommunication heuristics for the remainder of the network. It would be through a hierarchy of graph-server nodes that one would acquire a view of all the nodes in the network. However, due to privacy and maintenance issues, many of these graph-server types fail to provide any type of information about queried nodes. Due to this, we employed a one-week snapshot of this seventy-node QMP instance to establish ground-truth relevance for our work. The measurements were taken for seven days from the 1st to the 8th of March 2017, with a snapshot taken every hour [2]. The measurement period and frequency produced enough samples for evaluating guifi.net in light of the results of our method. We remark that node links in QMP and guifi.net, in general, are not symmetrical: the bandwidth and round-trip time from node u to node v isn't necessarily the same from v to u .

Firstly, although QMP has a more uniform set of nodes compared to guifi.net, it is also subject to the same behavioral user factors which influence the whole network [19]. This means that it may be considered as a representative sampling of guifi.net. Secondly, the obtained number of samples is high enough to enable us to apply statistical techniques to define empirical models of bandwidth and round-trip time. This allows us to fit different distributions to the measurements and evaluate the resulting goodness-of-fit (GOF) values. Selecting the most fitting distributions, we then synthesize their parameters in order to generate functions to produce artificial observations. Qualitatively, these simulated values are representative of the behavior of the QMP network (and thus of guifi.net) and were used to populate the bulk of our

dataset (guifi.net snapshot of January, 2017) nodes, which were missing data.

Phase One: Network Impact. We present in Figure 9 the distribution of maximum and average degree versus the size of the communities. The left side pertains guifi.net zone-based communities (from the dataset as-is), while the right side is related to the configuration of network node groups obtained with Phase One of our algorithm. We derive from this that our algorithm produces groupings with a tendency for greater node inter-connectivity.

Moving on, we further evaluate this derivation by producing a visualization of the average number of hops-to-leader for each community versus community size. Figure 10 presents this with respect to natural geographical zones of guifi.net in the left, with our algorithm's results on the right side. Our algorithm led to an overall reduction in the number of hops, in particular for smaller and more frequent communities. Figure 11 highlights interesting tendencies with regard to the impact of absolute heuristic weights and their influence on the average number of hops. In particular, we achieve this by isolating the range of community sizes to a maximum size of 250 members. Plotting these ranges over a logarithmic scale, it can be seen that the contained communities exhibit a lower number of hops. This tendency is particularly manifested with heuristics α_1 and β_3 (betweenness centrality and computational class of the node, respectively). We extrapolate from this finding that the fixed-region geographical definition of guifi.net may be too rigid and that it may in fact provide a user experience which is probably below-optimal regarding typical services offered in CNMCs. Usage of the Phase One technique shows promise with respect to optimizing the length of the path taken from each community's node to the community leader, a sure benefit for many services.

Phase Two: Leader Election Results. It is relevant to note that after Phase Two of our algorithm, the application of heuristics over the propagation-based node sets (right side) yielded more outliers than the geographical zones (left side). While there were more outliers in the results of Phase One of our algorithm, lower values were achieved when compared to the geographical node groups. We presented obtained results evaluated under different criteria. Our focus is not on producing a one-size-fits-all hierarchy of heuristics: other real-world scenarios upon which to test our algorithm will have specific objective functions, bound by application needs. The results are promising as they highlight that our algorithm is a valid alternative to traditional computational approaches to optimizing responsibility assignment to network nodes. We present Figures 12 and 13, which depict the number of average hops-to-leader in decreasing order. Orthogonally to node group definitions, the tendencies in the influence of the heuristics remain valid, with the same patterns appearing for each of the cases. It is interesting to note that, for the right side (based on Phase One of our algorithm), heuristics α_2 and β_3 produced greater differences between them. Accounting for

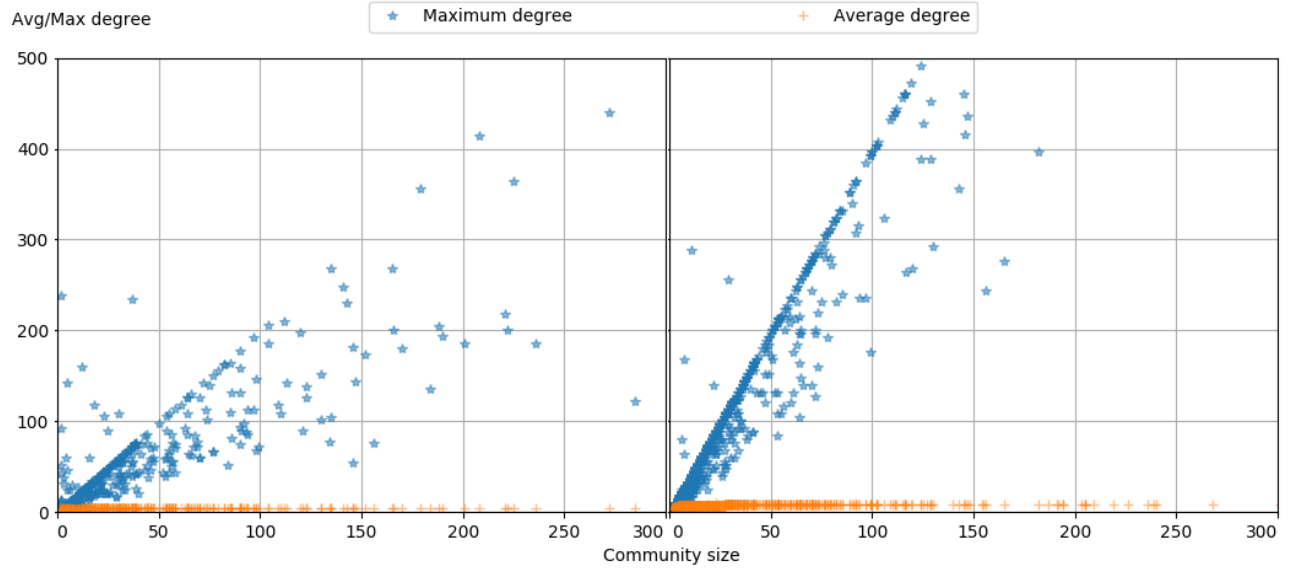


Figure 9. Plot of maximum and average degree distributions for each community. The left image is the geographical configuration of node sets, while the right side is based on Phase One of our algorithm.

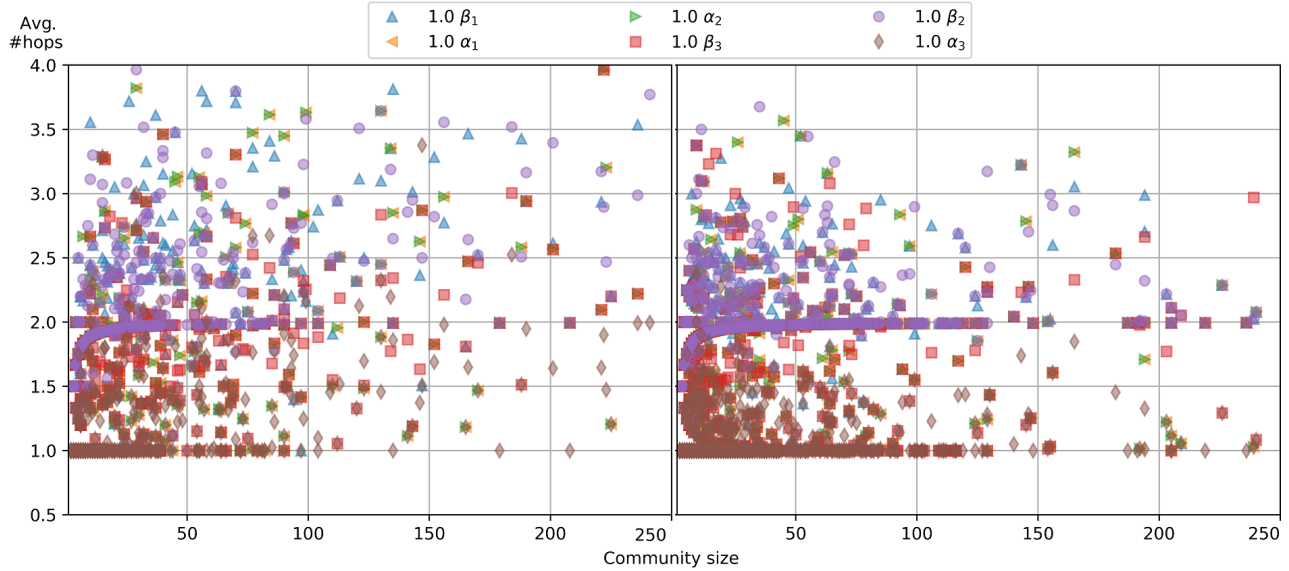


Figure 10. Average number of hops-to-leader plotted against each community’s size. The left image is the geographical configuration of node sets, while the right side is based on Phase One of our algorithm.

the computational class of nodes in the case of the right side led to a lower number of hops-to-leader compared to simply electing leaders based on centrality.

SLA Assessment. We also evaluate the quality of leaders in the context of the sampling performed for the QMP network. Namely, we modeled round-trip time (RTT) in milliseconds and bandwidth B/W in Mbit/s distributions based on around 70,000 samples (of bandwidth and round-trip time) obtained from QMP in guifi.net over a period of seven days. These two

features are relevant to types of SLAs inherent to services such as (RTT) web caching, web content requests, NoSQL cloud storage as well as (B/W) streaming and file download services. From these two features, we modeled their distribution and simulated their values for all of the guifi.net network snapshot mentioned earlier. Figures 14, 15, 16, and 17 were produced using the Python statsmodel package [12], which has a set of utilities to automate statistical processing tasks.

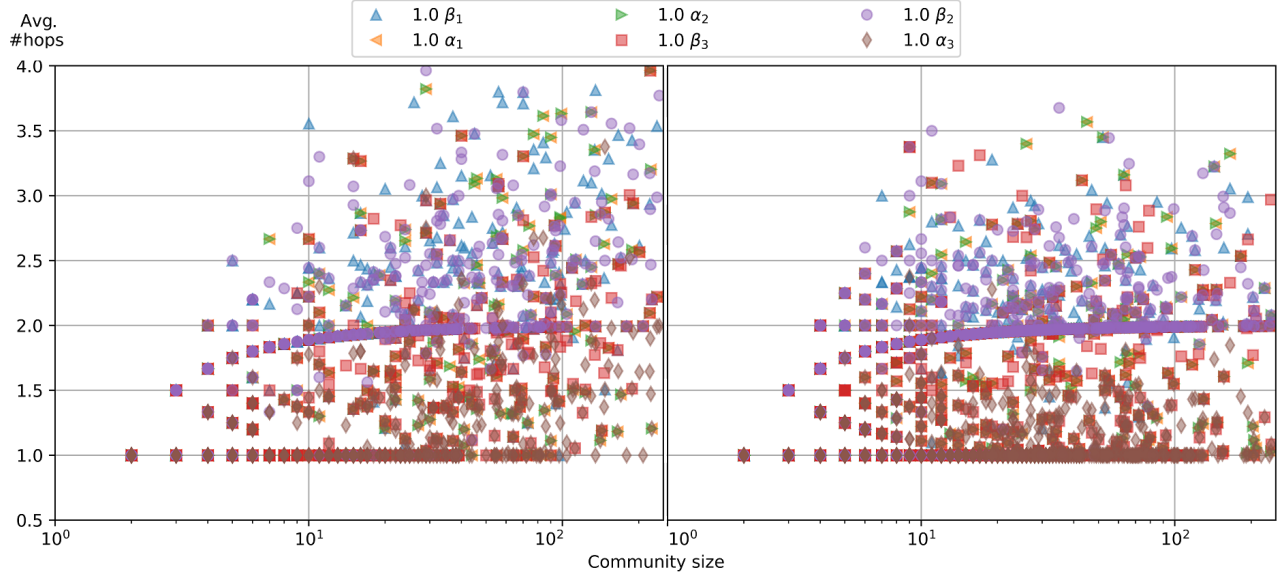


Figure 11. Average number of hops-to-leader plotted against a logarithmic scale of each community’s size. The left image is the geographical configuration of node sets, while the right side is based on Phase One of our algorithm.

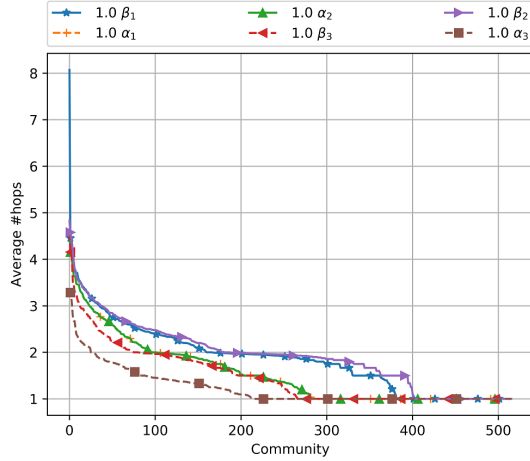


Figure 12. Average number of hops-to-leader against community in decreasing order for the original geographical configuration.

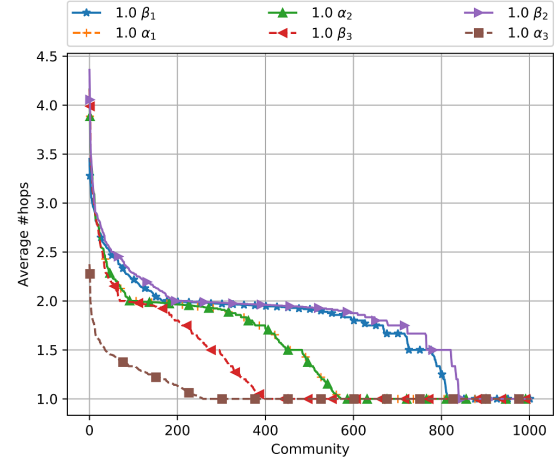


Figure 13. Average number of hops-to-leader against community in decreasing order for Phase One’s communities.

Figures 14 and 15 show the empirical cumulative distribution function of the bandwidth for the original guifi.net zones and for communities produced by Phase One of our algorithm, respectively. Interestingly, the bandwidth interval for [10; 30] Mbit/s in Figure 14 shows that Phase Two of our algorithm (the stage of leader election within a community – in this test case there is a one-to-one mapping between guifi.net zones and communities) fared better by using singular heuristics for electing the leader. That is, heuristics β_3 (computational class) and α_1 (betweenness centrality) produced, on average, more efficient bandwidth paths from a community’s nodes to

their leader. This tendency was also reproduced in the execution of Phase Two of our algorithm after Phase One (custom communities generated by Algorithm 1 instead of one-to-one mapping to guifi.net’s original zones), which can be seen in Figure 15.

For the round-trip time, the same two heuristic weights fared better than the others as well. All combinations seemingly max out (in terms of cumulative distribution) at around 125 milliseconds. However, up to about 100 milliseconds, heuristics β_3 (computational class) and α_1 (betweenness centrality) produced lower round-trip time. This occurs for the zone-based communities in Figure 16 and also the communities

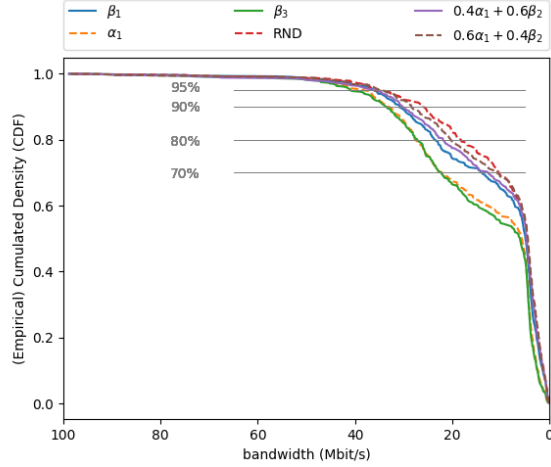


Figure 14. Zone-based bandwidth ECDF.

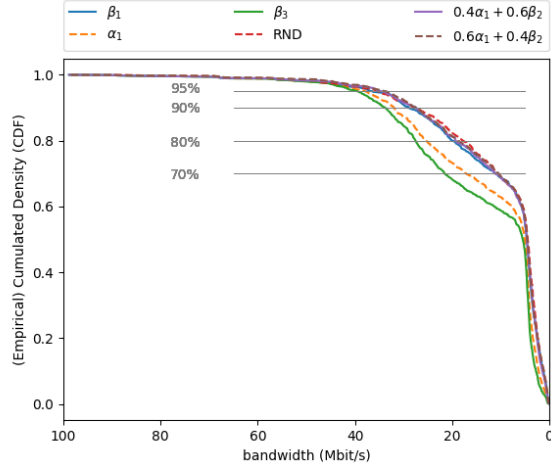


Figure 15. Community-based bandwidth ECDF.

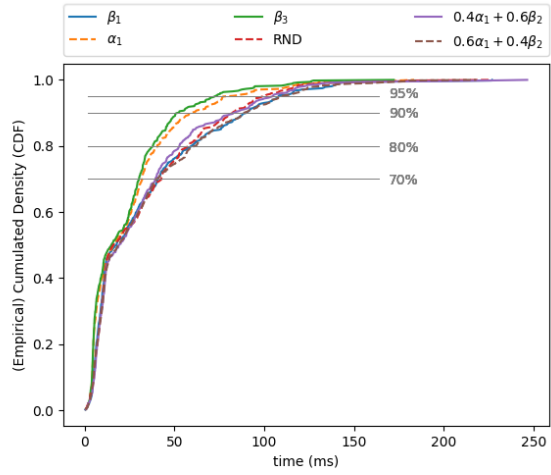


Figure 16. Zone-based round-trip time ECDF.

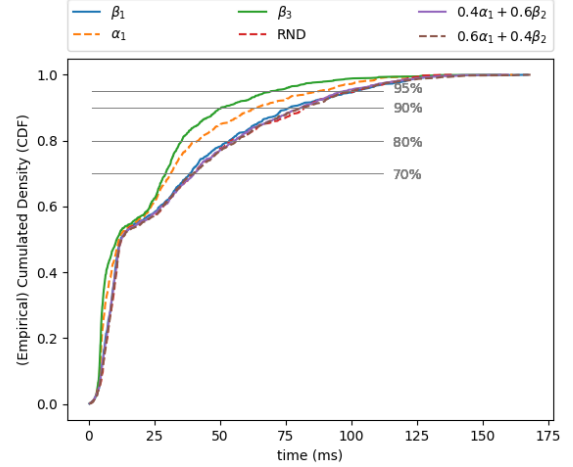


Figure 17. Community-based round-trip time ECDF.

resulting from Phase One of our algorithm, as illustrated in Figure 17. Curiously, we observe, as far as round-trip time is concerned, that the random leader election yielded practically the same results as the combined usage of betweenness centrality α_1 and latency β_2 (alternating between $0.4\alpha_1 + 0.6\beta_2$ and $0.6\alpha_1 + 0.4\beta_2$). We did not perform an exhaustive analysis of all possible combinations of heuristics and their weights. The combinations we present herein are relevant in terms of what the heuristics represent. Bandwidth samples were modeled as a four-parameter Kappa distribution, while round-trip time was modeled as a generalized extreme value (GEV) distribution. It is relevant to say that the empirical distributions of these two metrics exhibited a considerable degree of independence. In fact, $\text{corr}(BW, R) = -0.134$ for the sampled values, which means they appear to be only slightly inversely related. We assumed them to be independent with respect to results.

Summary. The method we present is inherently parallel and distributed, a break from traditionally-centralized often exhaustive optimization-driven solutions, opening possibilities for scalability. Phase Two of our algorithm was designed to be distributed with the purpose of executing concurrently among all communities. This implies that the computational time of this phase has an upper bound associated to the slowest-computing community. As far as the authors are aware, this work is the first that attempts to optimize service placement by defining communities using an analysis based purely on network theory and distributed graph processing. The guifi.net telecommunications network is one upon which different research projects have been executed [2, 19].

4 Related Work

Herein we go over alternative approaches to Phase One and Phase Two of our solution as a whole. We note that our work is novel, as far as we know, in the sense that it combines these two multidisciplinary phases, whose literature we analyze.

Community Networks. Different studies on guifi.net have drawn several insights: the network is not homogeneous – rural areas have topology properties different from those of metropolitan areas, such as density; the topology observed in rural areas is not scale-free (degree distribution does not fit a power law) due to the high number of terminals connected to some nodes; removing terminal nodes (with degree one) from the graphs in rural areas, however, reveals a scale-free *core-network* as in [19]. On the one hand, it is necessary to be aware of the challenges inherent to service allocation in different types of networks in the context of distributed systems. On the other hand, we highlight the existence of community detection techniques (in network theory) as a novel approach to these challenges. In recent years, metrics have been proposed for evaluating the quality of calculated communities have emerged: the most notorious one being that of modularity. However, focusing exclusively on modularity incurs community resolution penalties with smaller communities often not being detected. Considering this and focusing on scalability, other methods in the literature which do not use domain-specific heuristics were devised, such as the class of label propagation algorithms [8, 11]. These algorithms are inherently parallel and work well in practice for real world networks [1].

Service Placement. Typically, by monitoring all the physical and virtual resources on a system, service placement aims to balance load through the allocation, migration and replication of tasks. This can take place in cloud data-centers and in wireless networks that power a significant part of CNs. Most of the work in the data center environment, including distributed data centers, is not applicable to our case because we have a strong heterogeneity given by the limited capacity of nodes and links, as well as asymmetric quality of wireless links. The authors in [20] introduce a service allocation algorithm that provides near-optimal overlay allocations without the need to verify the whole solution space. They use static data from the network to identify node traits and minimize the coordination and overlay cost along a network. The work in [10] analyzes network topology and service dependencies, and combined with set of system constraints determines the placement of services within the wireless network. The authors use a multi-layer model to represent a service-based system embedded in a network topology and then apply an optimization algorithm to this model to find where best to place or reposition the services as the network topology and workload on the services changes.

In distributed micro-cloud environment (i.e., similar to our case), the work of Elmroth [18] takes into account rapid user mobility and resource cost when placing applications in Mobile Cloud Networks (MCN). A recent work of Tantawi [17] uses biased statistical sampling methods for cloud workload placement. Regarding the service placement through migration, the authors in [21] study the dynamic service migration problem in mobile edge-clouds that host cloud-based services at the

network edge. They formulate a sequential decision making problem for service migration using the framework of Markov Decision Process (MDP) and illustrate the effectiveness of their approach by simulation using real-world mobility traces of taxis in San Francisco. As a whole, mostly, service placement approaches are predominantly based on resource (CPU, memory) and node availability, and when they are network-aware, they are able just to employ static network information or at most process historical network data for availability predictions. Moreover, they are batch-oriented and execute sequentially in centralized settings and therefore cannot scale to larger network sizes, number of services, or greater network dynamism. Our approach is the first, to the best of our knowledge, that is dynamic, parallel and distributed, and therefore able scale seamlessly, by employing distributed graph processing systems, such as the Gelly library of Apache Flink. Thus, we are able to continually monitor service quality and perform service placement decisions continually/incrementally based on data gathered from the network (e.g., graph-servers in guifi.net).

5 Conclusion

In this paper, we presented a novel take on the processing steps that underlie service placement, a multi-objective problem. Compared to traditional system techniques (which, as far as we know, have not seen developments regarding parallel implementations and scalability with network size), our algorithm is expressed purely over state-of-the-art graph techniques which have inherent parallelism. This makes our algorithm a very competitive alternative, able to scale for networks which are orders of magnitude greater, when compared to other traditional techniques in the field.

Acknowledgments

This work was partly supported by the Portuguese government through FCT – Fundação para a Ciência e Tecnologia, under projects PTDC/EEI-SCR/6945/2014 and UID/CEC/500021/2013, by the ERDF through COMPETE 2020 Programme, within project POCI-01-0145-FEDER-016883, by the European H2020 project LightKone (H2020-732505), and by the Spanish government under contract TIN2016-77836-C2-2-R.

References

- [1] Paolo Boldi, Marco Rosa, Massimo Santini, and Sebastiano Vigna. 2011. Layered Label Propagation: A Multiresolution Coordinate-free Ordering for Compressing Social Networks. In *Proceedings of the 20th International Conference on World Wide Web (WWW '11)*. ACM, New York, NY, USA, 587–596. DOI : <https://doi.org/10.1145/1963405.1963488>
- [2] L. Cerdà-Alabern. 2012. On the topology characterization of Guifi.net. In *2012 IEEE 8th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*. 389–396. DOI : <https://doi.org/10.1109/WiMOB.2012.6379103>
- [3] Stuart Coles, Joanna Bawa, Lesley Trenner, and Pat Dorazio. 2001. *An introduction to statistical modeling of extreme values*. Vol. 208. Springer.
- [4] JRM Hosking. 2000. FORTRAN routines for use with the method of L-moments, Version 3.04. *IBM Research* (2000).

- [5] Jonathan RM Hosking. 1994. The four-parameter kappa distribution. *IBM Journal of Research and Development* 38, 3 (1994), 251–258.
- [6] J. R. M. Hosking. 1990. L-Moments: Analysis and Estimation of Distributions Using Linear Combinations of Order Statistics. *Journal of the Royal Statistical Society. Series B (Methodological)* 52, 1 (1990), 105–124. <http://www.jstor.org/stable/2345653>
- [7] Jonathan Richard Morley Hosking and James R Wallis. 2005. *Regional frequency analysis: an approach based on L-moments*. Cambridge University Press.
- [8] Ian X. Y. Leung, Pan Hui, Pietro Liò, and Jon Crowcroft. 2009. Towards real-time community detection in large networks. *Phys. Rev. E* 79 (Jun 2009), 066107. Issue 6. DOI : <https://doi.org/10.1103/PhysRevE.79.066107>
- [9] Mark Newman. 2010. *Networks: An Introduction*. Oxford University Press, Inc., New York, NY, USA.
- [10] Petr Novotny, Rahul Uргаonkar, Alexander L Wolf, and Bongjun Ko. 2015. Dynamic placement of composite software services in hybrid wireless networks. In *Military Communications Conference, MILCOM 2015-2015 IEEE*. IEEE, 1052–1057.
- [11] Usha Nandini Raghavan, Réka Albert, and Soundar Kumara. 2007. Near linear time algorithm to detect community structures in large-scale networks. *Phys. Rev. E* 76 (Sep 2007), 036106. Issue 3. DOI : <https://doi.org/10.1103/PhysRevE.76.036106>
- [12] Skipper Seabold and Josef Perktold. 2010. Statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference*.
- [13] Mennan Selimi, Llorenç Cerdà-Alabern, Marc Sanchez Artigas, Felix Freitag, and Luís Veiga. 2017. Practical Service Placement Approach for Microservices Architecture. In *IEEE/ACM 17th International Symposium On Cluster, Cloud And Grid (CCGRID 2017)*. ACM/IEEE.
- [14] Mennan Selimi, Felix Freitag, Llorenç Cerdà-Alabern, and Luís Veiga. 2016. Performance Evaluation of a Distributed Storage Service in Community Network Clouds. *Concurrency and Computation: Practice and Experience* 28, 11 (Aug. 2016), 3131–3148.
- [15] Mennan Selimi, Amin M Khan, Emmanouil Dimogerontakis, Felix Freitag, and Roger Pueyo Centelles. 2015. Cloud services in the Guifi. net community network. *Computer Networks* 93 (2015), 373–388.
- [16] Mennan Selimi, Davide Vega, Felix Freitag, and Luís Veiga. 2016. Towards Network-Aware Service Placement in Community Network Micro-Clouds. In *Euro-Par 2016: Parallel Processing - 22nd International Conference on Parallel and Distributed Computing*. Springer.
- [17] A. N. Tantawi. 2016. Solution Biasing for Optimized Cloud Workload Placement. In *2016 IEEE International Conference on Autonomic Computing (ICAC)*. 105–110. DOI : <https://doi.org/10.1109/ICAC.2016.34>
- [18] William Tärneberg, Amardeep Mehta, Eddie Wadbro, Johan Tordsson, Johan Eker, Maria Kihl, and Erik Elmroth. 2017. Dynamic application placement in the Mobile Cloud Network. *Future Generation Computer Systems* 70 (2017), 163 – 177. DOI : <https://doi.org/10.1016/j.future.2016.06.021>
- [19] D. Vega, L. Cerdà-Alabern, L. Navarro, and R. Meseguer. 2012. Topology patterns of a community network: Guifi.net. In *2012 IEEE 8th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*. 612–619. DOI : <https://doi.org/10.1109/WiMOB.2012.6379139>
- [20] Davide Vega, Roc Meseguer, Guillem Cabrera, and Joan Manuel Marqués. 2014. Exploring local service allocation in community networks. In *Wireless and Mobile Computing, Networking and Communications (WiMob), 2014 IEEE 10th International Conference on*. IEEE, 273–280.
- [21] Shiqiang Wang, Rahul Uргаonkar, Ting He, Kevin Chan, Murtaza Zafer, and Kin K. Leung. 2017. Dynamic Service Placement for Mobile Micro-Clouds with Predicted Future Costs. *IEEE Transactions on Parallel and Distributed Systems* 28, 4 (April 2017), 1002–1016. DOI : <https://doi.org/10.1109/TPDS.2016.2604814>