



**Nuno Morais**

Master of Science

## **Edge DeMon Edge Decentralized Monitoring**

Dissertação para obtenção do Grau de Mestre em  
**Engenharia Informática**

Orientador: João Leitão, Assistant Professor,  
NOVA University of Lisbon

Júri

Presidente: Name of the committee chairperson  
Arguente: Name of a rapporteur  
Vogal: Yet another member of the committee



FACULDADE DE  
CIÊNCIAS E TECNOLOGIA  
UNIVERSIDADE NOVA DE LISBOA

**December, 2019**



## RESUMO

---

Lorem ipsum em Português.

**Palavras-chave:** Palavras-chave (em Português) ...

---



## ABSTRACT

---



# ÍNDICE

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Motivation . . . . .	2
1.3	Expected Contributions . . . . .	3
1.4	Document structure . . . . .	3
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	Offloading computation to the Edge . . . . .	6
2.1.1	Edge Computing . . . . .	6
2.1.2	Discussion . . . . .	9
2.2	Topology Management . . . . .	10
2.2.1	Evaluating topologies . . . . .	10
2.2.2	Taxonomy of overlay networks . . . . .	11
2.2.3	Unstructured Overlays . . . . .	11
2.2.4	Structured overlays . . . . .	13
2.2.5	Discussion . . . . .	14
2.3	Resource Location and Discovery . . . . .	14
2.3.1	Query taxonomy . . . . .	14
2.3.2	Query dissemination . . . . .	15
2.3.3	Centralized Resource Location . . . . .	16
2.3.4	Decentralized Resource Location . . . . .	16
2.3.5	Discussion . . . . .	19
2.4	Resource monitoring . . . . .	19
2.4.1	Device monitoring . . . . .	19
2.4.2	End-to-end link monitoring . . . . .	20
2.4.3	Aggregation techniques . . . . .	21
2.4.4	Relevant aggregation protocols . . . . .	23
2.4.5	Monitoring systems . . . . .	24
2.5	Resource management . . . . .	26
2.5.1	Resource Management Taxonomy . . . . .	26
2.5.2	Relevant Resource Monitoring Systems . . . . .	28
2.5.3	Discussion . . . . .	31

## ÍNDICE

---

<b>3 Planning</b>	<b>33</b>
3.1 Proposed Solution . . . . .	33
3.2 Scheduling . . . . .	34
<b>Bibliografia</b>	<b>35</b>



## INTRODUCTION

### 1.1 Context

Nowadays, the Cloud Computing paradigm is the standard for development, deployment and management of services, most software present in our everyday life such as Google Apps, Amazon, Twitter, among many others is deployed on some form of cloud service. This paradigm has proven to have massive economic benefits that make it very likely to remain permanent in future of the computing landscape. Cloud Computing provides the illusion of unlimited computing power, which and has revolutionized the way developers, users and businesses rationalize about building and deploying applications [1].

However, the rise in popularity of mobile applications and IoT applications differs from the centralized model proposed by the Cloud Computing paradigm. With recent advances in the IoT industry, it is safe to assume that in the future almost all consumer electronics will play a role in producing and consuming data. However, when computation resides in the data center (DC), far from the source of the data, problems arise: from the physical space needed to contain all the infrastructure, the increasing amount of bandwidth needed to support the information exchange from the DC to the client, the latency in communication from the client to the DC as well as the security aspects that arise from offloading data storage and computation.

The aforementioned aspects have directed us into a post-cloud era where a new computing paradigm emerged, Edge Computing. Edge computing takes into consideration all the computing and network resources which act as an "edge" along the path between the data source and the DC. It addresses the increasing need for supporting interaction between cloud computing systems and mobile or IoT applications [41], and allows the emergence of novel edge-enabled applications (e.g. traffic management, smart city management, mobile games, among others).

Additionally, systems that require real-time processing of data may be not be feasible with cloud computing, Google's self-driving car generates 1 Gigabyte every second [39], while a Boeing 787 will create around 5 gigabytes of data per second [12]. If this data were to be processed in real-time (e.g. towards self-driving), it would be infeasible to transport it to cloud and back.

## 1.2 Motivation

When accounting for all the devices that are external to the DC, we are met by a huge increase in heterogeneity of devices: from Data Centers to private servers, desktops and mobile devices to 5G towers and ISP servers, among others. Contrary to the cloud, edge environments tend to be highly dynamic, devices have constrained computational power and their connections are often limited in capacity and reliability.

Developing an efficient general compute platform for edge environments is still an open challenge in Edge Computing. A crucial requirement towards this is performing efficient **resource management**, which consists of keeping track of the tasks to perform and manage the utilization of computational resources of each device. General compute platforms are extensively used in Cloud systems (e.g. Mesos [18], Yarn [45]), however, those solutions are tailored towards small numbers of homogenous resource-heavy devices, which mismatches the edge environment.

To overcome this, a resource management solution for the edge must be capable of managing very large numbers of heterogenous devices and tasks. An infrastructure capable of performing the aforementioned tasks successfully has strong applicability in today's world (e.g. Cloud providers, Smart Cities, among others).

A particularly hard task in resource management is **scheduling**, which consists in distributing tasks among nodes in the system, ideally, the task distribution must promote a balanced resource usage among nodes in the system. One of the popular solutions is transporting all the monitoring data towards a centralized point and redistribute the tasks among nodes. However, this presents a centralized point of failure and a point of contention in a large scale system.

Alternatively, nodes with only partial knowledge of the system (ideally even without sending any additional messages) must be able to autonomously offload tasks towards neighbors. However, offloading tasks of varied complexity in heterogeneous nodes is not an easy task. The accuracy and freshness of the **monitoring information** each peer has dictates how efficiently they can offload tasks such that the system remains balanced, and applications running on it maintain (and possibly improve) quality of service. All of this while adapting to environment changes.

Given this, peers must integrate a robust decentralized **resource monitoring system** which tracks device and service metrics. This system must federate peers such that they leverage on heterogeneity to build a hierarchical infrastructure, which combines naturally with the device taxonomy. Finally, resource monitoring systems must provide

### 1.3 Expected Contributions

The expected contributions, as will be further detailed in section 3.1, derive from the aforementioned challenges:

- Devise a decentralized monitoring infrastructure for edge devices which employs a topology tailored towards edge environments, and a combination of monitoring and aggregation techniques in a natural way which promotes load-balancing and networking locality.
- Evaluate the designed infrastructure through simulation (e.g. iFogSim or PeerSim) and compare the performance with similar systems.
- Design an experimental scenario to test the system under various scenarios. For this, we must design or adapt an existing scheduler and test the feasibility of the solution by deploying applications and collecting metrics about the overlay, the potential reduction in cost and the quality of service.

### 1.4 Document structure

This document is structured in the following manner:

**Chapter 2** focuses on the related work, first section covers edge computing in further detail, next we cover P2P systems and the different types of topology management protocols. Third section studies the different types of resource location architectures, namely how to efficiently find a specific peer in the system and common techniques towards performing efficient searches over networks composed by a large number of devices. Fourth section covers aggregation techniques and popular implementations of these systems **era para meter isto em monitoring?**. Finally, last section covers popular resource monitoring systems.

**Chapter 3** further explains the proposed contribution and proposes the work plan for the remainder of the thesis.



## RELATED WORK

The goal of this chapter is to further identify the challenges that arise from creating a device monitoring infrastructure which enables efficient scheduling/offloading of tasks across edge environments, and to study related work towards the identified challenges.

Computing platforms (e.g. Mesos [18], Yarn [45]) are usually tailored towards cloud computing systems. As such, they are usually centralized and tailored towards clusters of homogeneous devices which have very different monitoring requirements as those in edge environments. Given this, we identify a set of requirements which we believe to be important towards monitoring systems in general, and critical in an edge environment:

1. *Scalability*, a scalable monitoring system can handle a remarkable number of monitored devices and services. This property is paramount towards edge computing environments because of the necessity of managing a wide variety of parameters which needs to be monitored across a largely decentralized and fault-prone environment. Existing cloud monitoring tools fail to distribute the monitoring load, which in turn leads to a single point of failure and lack of scalability.
2. *Non-intrusiveness*, following the edge computing viewpoint of performing simple lightweight tasks, the monitoring solution should follow this approach. Given this, a monitoring tool which adopts minimal processing, consumes low memory and generates low amounts of traffic is essential.
3. *Interoperability*, edge computing aims at a cooperative deployment of applications interconnected over both cloud and edge infrastructures. Ideally, companies should be able to deploy components across a wide range of infrastructures owned by different providers.

4. *Robustness*, the monitoring system must be robust in the face of failures or network partitions. Especially in the edge environment, where devices are scattered and restricted in capacity.
5. *Live Migration Support*, with the current offers of virtualization technologies, such as VMs or containers, there is a wide variety of migration tools for resource management systems. For edge-enabled applications, live migration is a highly desirable feature, as it is a crucial tool towards enabling the self-adaptation to the environment (e.g. dealing with failures, or scaling systems).

Now that we have identified the main requirements of such a system, we study related work towards solving these challenges. First, we study the taxonomy of resources in the edge of the network, and discuss how they can be employed towards the proposed solution (Section 2.1).

After understanding the environment taxonomy, we need to federate devices and sub-networks in an efficient abstraction layer (to ensure interoperability), it is paramount that such layer provides efficient communication patterns, traffic locality and load balancing (according to device heterogeneity) (Section 2.2).

Furthermore, nodes must employ the aforementioned topology in order to efficiently find resources (e.g. other peers, services or even computing power) in the system, and to keep track of which tasks are running and where. For this, we cover different types of queries employed in resource location systems, how to disseminate them, and finally study related resource location systems in the state of the art (Section 2.3).

Next, we study how to **monitor** the state of the devices and the tasks running on them, we study popular techniques of monitoring device and task status, and how to aggregate those values (Section 2.4). Lastly, we briefly study how the monitoring data can be used to distribute and offload tasks (Section 2.5).

## 2.1 Offloading computation to the Edge

In this chapter we provide context about related work towards decentralizing computation from the Cloud. First, we define how we characterize Edge Computing and discuss how it is related to edge-related paradigms. Following, we study the taxonomy of the environment and focus on which computations each device can perform.

### 2.1.1 Edge Computing

As previously mentioned, edge computing calls for the processing of data in the edge of the network, specifically servicing IoT devices and performing computations on behalf of cloud services [40].

As previously mentioned, edge computing has the potential of enabling novel edge-enabled applications along with providing many improvements to existing systems:

- *Reducing the amount of traffic*, if the workload of a client can be computed in the cloud, then the latency can be greatly improved vs when doing the work in the cloud. For example, in the case of a smart home which computes the average home temperature in all divisions, almost all the data collected by the sensors can be aggregated (averaged) in a home gateway, effectively summarizing the data from a set of values to one single value. Reducing the data size improves the transmission reliability and saves bandwidth for other edge devices.
- *Latency* is one of the most important aspects towards the performance of applications / services, especially time-critical applications (e.g. traffic monitoring) have special latency needs. Latency consists in the time it takes for a packet to travel the network from the origin endpoint to the target endpoint. Edge computing solves this by offloading the computation towards the nearest edge device, for example a cellphone offloading facial recognition to the nearest 5G cell tower.
- Facilitating new approaches of *load-balancing*, given that there is a massive increase of devices in the edge of the network, if all devices cooperate towards a common goal, they can achieve tasks similar to those of more complex systems (e.g. sensor networks performing sql queries [29]). For example, edge computing solves service load-balancing by offloading its computation towards the nearest available neighbor which has more capacity, or by offloading a portion of its tasks.
- *Minimizing energy consumption*, edge computing has the potential of performing some of the aforementioned tasks, however, instead of optimizing towards latency or traffic locality, devices take into consideration their energy percentage. For example, a device which has very low latency but also low battery in a time-critical service can choose to offload any unnecessary computations towards neighbors, as a cost-saving measure. Additionally, given that devices are closer to their targets, communications are faster and consequently consume less energy.

Many approaches have already leveraged on some form of Edge computing in the past. **Cloudlets** [46] are an extension of the cloud computing paradigm beyond the DC, which consists of deploying resource rich computers near the vicinity of users that provide cloud functionality.

A limitation of cloudlets is that because they are specialized computers, they cannot not guarantee low-latency ubiquitous service provision, and consequently cannot satisfy QoS of large hotspots of users. Cloudlets have become a trending subject, and have been employed towards resource management, Big Data analytics, security, among others.

**Content Distribution networks** [33] (CDNs) emerged to address the overwhelming utilization of network bandwidth and server capacity that arose with bandwidth-intensive content (e.g. streaming HD video). In short, CDNs consist of specialized high bandwidth servers strategically located at the edge of the network, these servers replicate content

from a certain origin and serve it at reduced latencies, effectively decentralizing the content delivery.

Additionally, many paradigms have emerged which propose to solve similar problems to the Edge Computing paradigm. **Fog Computing** [4] proposes to provide compute, storage and networking services between end devices and traditional cloud computing data centers, typically, but not exclusively located at the edge of the network. Fog computing is interchangeable with our vision of all devices acting as an "edge" contributing towards computing tasks.

**Osmotic Computing** [47] envisions the automatic deployment and management of inter-connected microservices on both edge and cloud infrastructures. Osmotic computing envisions edge devices employing an orchestration technique similar to the process of "osmosis". Translated, this consists in dynamically detecting and resolving resource contention via coordinated microservice deployments, furthermore, this paradigm is focused towards ensuring and maintaining quality of service.

**Multi-access edge computing** [30] (MEC) formerly known as mobile-edge cloud computing, is a network architecture that proposes to provide fast-interactive responses for mobile applications. It solves this by employing the network edge (e.g. base stations and access points) to provide compute resources for latency-critical mobile applications. MEC is a subset of our edge computing vision, although with a higher focus on communications technology and how to offload the computation from mobile to the cloud and not vice-versa.

#### 2.1.1.1 Edge Environment Taxonomy

Similar to [26], we classify edge device according to 3 main attributes: **capacity** refers to computational, storage and connectivity capabilities of the device, **availability** consists in the probability of a device being reachable, and finally, **domain** characterizes the way in which a device may be employed towards applications. In short, if the device can support the whole *applicational domain* or only the activities of a single user (*user domain*).

Tabela 2.1: Taxonomy of the edge environment

Level	Category	Availability	Capacity	Level	Category	Availability	Capacity
L0	Cloud Data Centers	High	High	L4	Priv. Servers & Desktops	Medium	Medium
L1	ISP, Edge & Private DCs	High	High	L5	Laptops	Low	Medium
L2	5G Towers	High	Medium	L6	Mobile devices	Low	Low
L3	Networking devices	High	Low	L7	Actuators & Sensors	Varied	Low

Table 2.1 shows the categories of edge devices, we assign levels to categories as a function of the distance from the cloud infrastructure. Coincidentally, the levels are correlated to the number of devices and their computational power, where higher levels tend to have more devices that are closer to the origin of the data and have lower computational power.

**Levels 0 and 1** *cloud and edge DCs* offer pools of computational and storage resources, that can dynamically scale to support the operation of edge-enabled applications. Both of these options have high availability and large amounts of storage and computational



power, as such, there is no limitations on which type of computations these devices can perform.

**Level 2** is composed of *5G cell towers*, which serve as access points for mobile devices, **Level 3** also consists of *networking devices*, although with lower capacity than those in level 2. Devices in both levels have high availability, and they can contribute to the applicational domain, however in a limited fashion (e.g. coordinate resource management, host a microservice, or just act as a gateway for mobile devices). Devices in this level can easily improve the management of the network, for example, by manipulating data flows.

**Level 4** Consists of *private servers and desktops*, it is the first level where devices belong to the user domain. Devices in this level have medium capacity and availability, and can perform a varied amount of tasks on behalf of devices in higher levels (e.g. compute on behalf of smartphones, act as logical gateways or just cache data).

**Level 5** consists of *laptops*, which are also on the user domain, and can perform a role similar to devices in level 4, although with lower availability and capacity. The main differentiating factor with devices in level 4 is that laptops are battery-powered, which means that energy consumption must also be taken into account whenever monitoring and computing on these devices.

**Level 6** consists of *tablets and mobile devices*, devices in this level act as producers and consumers of data and belong to the user domain. Because because of their low capacity, availability and short battery life, mobile devices are limited in how they can perform computations and contribute towards edge applications. Aside from caching user data, common usages are filtering or aggregation of data generated from devices in level 7.

Finally, **level 7** consists of *actuators, sensors and things*, these devices are the most limited in their capacity, and varied availability. *Things* act both as data producers and consumers towards edge-enabled applications. They enable limited forms of computation in the form of aggregation and filtering.

### 2.1.2 Discussion

Intuitively, the lower the level the harder it is to employ devices towards applications. Devices in levels 6 and 7 are especially restricted due to having lower availability and computational power, however, these can still be used in specific scenarios (e.g. an application with very low computational overhead but real-time latency requirements).

Devices in levels 0-5 are potential candidates towards building the resource monitoring system we intend to create. Given the low availability of devices in higher levels, they are not very suitable, as they could not contribute much towards the system and would incur churn. This effect can be circumvented by employing devices in other levels as gateways for mobile devices and *things*.

## 2.2 Topology Management

As previously mentioned, a challenge towards solving the proposed solution is to federate all peers in an abstraction layer that allows intercommunication and efficient resource discovery. Given that this is a classic P2P problem, this section provides context about P2P systems and the taxonomy of overlay networks.

In P2P, participants contribute to the system with a portion of their resources, so that the overall system can accomplish tasks which would otherwise be impossible for a single peer to solve. However, due to memory and communication overhead, it is undesirable that all nodes in a P2P system collaborate with all other peers (unless in specific scenarios which we will further elaborate).

Instead, peers select a subset of peers in the system to establish neighboring relations. These neighboring relations are usually constructed on top of underlying links from an already existing network (called an underlay). Intuitively, the accumulation of neighboring relations on top of the underlay network is what constitutes the **overlay network**.

### 2.2.1 Evaluating topologies

If we look at neighboring connections as edges in a graph, we can define a set of metrics to measure graph-related metrics to measure overlay performance:

1. **Connectivity.** A connected graph is one where there is at least one path from each node to all other nodes in the system. The absence of this property means that there are nodes in the system that are isolated, thus will not be able to cooperate towards the overall behavior of the system. This property is usually measured as a percentage, corresponding to the largest portion of the system that is connected. Intuitively, a connected overlay has 100% connectivity.
2. **Degree Distribution.** The degree of a node consists in the number of arcs that are connected to it. Depending on the type of system, the connections may be directed or undirected, in a directed graph there is a distinction between **in-degree** and **out-degree** of a node. Intuitively, nodes with a high in-degree have higher reachability in the system, and nodes with 0 in-degree cannot be reached. In flat overlays, where load distribution is desired, degree distribution should be as similar as possible in all nodes. By contrast, in hierarchical overlays, designs take advantage of device heterogeneity to differentiate between peers and promote scalability.
3. **Average Shortest Path.** A path is composed by the edges of the graph that a message would have to cross to get from one node to other. The average shortest path consists in the average of all those paths, to promote efficient communication patterns, is desirable that this value is as low as possible.
4. **Clustering Coefficient.** The clustering coefficient provides a measure of the density of neighboring relations across the neighbors of a given node. It consists in the

number of a node's neighbors divided by the maximum number of links between those neighbors. Similar to the average shortest path, the clustering coefficient of an overlay consists in the average of the clustering coefficient of all the peers. A high value of clustering coefficients will result in a higher number of redundant messages, and by consequence, additional localized traffic. Finally, areas of an overlay with a higher clustering coefficient tend to be more easily isolated from the remaining system.

5. **Overlay Cost.** If we assume that a link in the overlay has a *cost*, then the overlay cost is the sum of all the links that form the overlay. Link cost can derive from overlay metrics (numeric distance, XOR distance, etc), or external metrics such as latency.

### 2.2.2 Taxonomy of overlay networks

There are two main approaches to build decentralized P2P systems, which are split in two categories: structured and unstructured. Within those categories we have identified the following sub-categories: flat and hierarchical.

### 2.2.3 Unstructured Overlays

**Unstructured overlays** usually impose little to no rules in neighboring relations, peers may pick random peers to be their neighbors, or alternatively employ strategies to "rank" neighbors and selectively pick the "best".

A key factor of unstructured overlays is their low maintenance cost, given that nodes can easily create neighboring relations and replace failed ones. Consequently, this is the type of overlay which offers better resilience to churn [43] (participants concurrently entering and leaving the system).

#### 2.2.3.1 Flat unstructured Overlays

A flat unstructured overlay is an overlay where peers contribute evenly towards a common system goal. These overlays attempt to have even degree distributions while providing good connectivity. A prime example of a flat unstructured overlay which is highly resilient to churn and catastrophic failures is *Hyparview* [24].

*Hyparview* (Hybrid Partial View) gets its name from maintaining two exclusive views: the *active* and *passive* view, that are distinguished by their maintenance strategy.

The *passive view* is a larger view which consists of a random set of peers in the system, this view is maintained by a simple gossip protocol which periodically gossips a message to a random peer in the active view. This message contains a subset of the neighbors of the sending node and a time-to-live (TTL), the message is forwarded in the system until the TTL expires. In contrast, the *active view* is a smaller view (around  $\log(n)$ ) created during the bootstrap of the protocol, and actively maintained by monitoring peers with a TCP connection (effectively making the active view connections bidirectional and act

as a failure detector). Whenever peers from the active view fail, nodes attempt to replace them with nodes in the passive view.

Hyparview achieves high reliability even in the face of high percentage of high node failures (up to 80-90% of all nodes). This is highly desirable in edge environments. However, due to battery constraints, many devices cannot actively maintain TCP connections, which limits the applicability of Hyparview towards levels 0-4 of the taxonomy described in 2.1.1.1. Hyparview is often used as a *peer sampling service* for other protocols which rely on the connections from the active view to collaborate (e.g. PlumTree [25] and X-BOT [0]).

T-MAN [20] is protocol based on a gossiping scheme, which proposes to build a wide range of overlay networks from scratch (e.g. ring, mesh, tree, etc.). To achieve this, T-MAN takes an overlay as an input to the protocol, this overlay is represented by *ranking method*. The ranking method sorts a set of nodes according to a given metric, where first nodes are the most "preferable" of the list. The resulting protocol is scalable and fast, with convergence times that grow as the logarithm of the network size. Furthermore, it is completely decentralized and extremely robust. Limitations that arise from using T-Man is that it does not ensure stability of in-degree of nodes during the optimization of the overlay.

X-BOT [0] X-BOT is a protocol which constructs an unstructured overlay network where neighboring relations are biased towards a certain metric. X-BOT does so while preserving key properties of the overlay such as the node degree and consequently, the overlay connectivity.

Neighboring connections are attributed a metric according to an *oracle*, which consists in a component which exports a function that takes a pair of peers in the system and attributes a cost to that neighboring connection. An oracle may take into account device latency, ISP distribution, stretch, among others.

The rationale X-BOT is as follows, nodes maintain active and passive views similar to Hyparview [24]. Then, nodes periodically trigger optimization rounds where nodes attempt to swap one neighbor from their active active view with another neighbor which ranks higher according to the oracle.

An important factor (which contrasts with [20]) is that nodes in the system maintain unbiased neighbors such that the overlay maintains low average path length and low clustering coefficient.

### 2.2.3.2 Unstructured Hierarchical Overlays

Unstructured hierarchical overlays are characterized as overlays where peers have different tasks in the system, which easily accommodates device heterogeneity while potentially increasing the performance of the system.

An example is **super-peers**, which are peers that have increased capacity and stability, that are commonly assigned towards disseminating queries throughout the system or

caching file locations.

should i do trees here?

Super-peers have proven their effectiveness in reducing the number of peers that have to exchange messages, which by consequence raises system scalability. This is the approach taken by Gia [6] to improve the scalability of Gnutella [16]. However, this approach is inefficient when finding rare resources in the system.

**Overnesia** [27] is a protocol which creates virtual super-peers by maintaining fully connected groups of nodes with the same identifier known by all elements in the group. Nodes join the system by sending a JOIN request to a bootstrap node which triggers a random walk, the requesting node joins the group where random walk finishes (either because it finds an underpopulated group or because the TTL expires).

Then, in order to promote intra-group membership consistency, nodes employ an anti-entropy mechanism where they periodically exchange messages containing their own view of the group.

When a group detects that its size has become too large, it triggers a dividing procedure which splits the groups in two halves. This division reduces the costs of replicating data among its members, but also promotes good load-balancing of the system.

When group size has fallen below a certain threshold, Overnesia triggers a collapse procedure, where each node takes the initiative to relocate itself to another group, resulting in the collapse of the group.

Finally, nodes perform walks along the overlay in order to establish inter-group links which in turn use to perform efficient broadcasts.

#### 2.2.4 Structured overlays

**Structured overlays** enforce stronger rules towards neighbor selection (generally based on identifiers of peers). As a result, the overlay generally converges to a topology known a priori, where the target topologies are tailored towards applicational requirements.

A canonical example of a structured overlay is a distributed hash table (DHT). DHTs offer efficient routing capabilities over the identifier space (usually routing procedures take a logarithmic number of steps), DHTs have been extensively used to support many large scale services (publish-subscribe, file sharing, among others) and are especially used in Cloud-based environments.

##### 2.2.4.1 Flat structured overlays

In a flat DHT, peers use consistent hashing functions to select random identifiers which are uniformly distributed over the identifier space. We will cover distributed hash tables with a bigger emphasis in section 2.3

#### 2.2.4.2 Hierarchical structured overlays

The second approach to building a hierarchical topology is to employ a DHT, hierarchical DHTS usually form contained DHTS within other DHTS (e.g. a ring within a ring). This offers several important advantages over a flat DHT: first, lookups take less hops and messages to reach the target, second, organizing nodes in disjoint groups allows traffic locality if groups of nodes are close in the underlay, finally, churn events within a group stay contained within that group. However, many of these systems either employ more memory to accommodate the many levels hierarchical DHT, or tradeoff reliability (by shortening the number of connections) for memory and communication efficiency.

#### 2.2.5 Discussion

Unstructured overlays are an attractive option to federate large amounts of devices in heavily dynamic environments. They provide a low clustering coefficient and good connectivity even in the face of churn, which makes them appealing towards the edge environments.

Conversely, structured overlays enable efficient routing procedures with very low message overhead, which makes them suitable for resource location systems. However, given the strict neighboring rules, nodes cannot replace neighbors easily, which hinders the fault-tolerance of these types of topologies.

### 2.3 Resource Location and Discovery

Given that the main challenge to solve is to provide a platform which enables monitoring of edge devices and tasks running on them, it is imperative that peers are able to find the resources they need (e.g. services, peers or computing power) to meet their requests. For this, peers need implement a **resource location system**.

Resource location systems are one of the most common applications of the P2P paradigm. In these systems, a participant provided with a resource descriptor is able to query peers and obtain an answer to the location (or absence) of that resource in the system within a reasonable amount of time.

#### 2.3.1 Query taxonomy

1. **Exact Match queries** specify the resource to search by the value of a specific attribute (for example, a hash of the value).
2. **Keyword queries** employ one or more keywords (or tags) combined with logical operators to describe resources (e.g. "pop", "rock", "pop and rock"...). These queries return a list of resources and peers that own a resource whose description matches the keyword(s).

3. **Range queries** retrieve all resources whose value is contained in a given interval (e.g. "movies with 100 to 300 minutes of duration"). These queries are especially applied in databases.
4. **Arbitrary queries** are queries that aim to find a set of nodes or resources that satisfy one or more arbitrary conditions, a possible example of an arbitrary query is looking for a set resources with a certain size or format.

### 2.3.2 Query dissemination

**Disseminating** the previously mentioned queries in an efficient manner through the overlay is a challenge in P2P resource location systems, there have been devised many **dissemination strategies** whose applicability depends on the applicational requirements and system capabilities. There are two main types of dissemination of queries, **flooding** and **random walks**.

#### 2.3.2.1 Flooding

When **flooding**, peers eagerly forward queries to other peers in the system, the objective of flooding is to contact a certain number of distinct peers in the system that may have the desired resource.

One approach is **complete flooding** which consists in contacting every node in the system, this guarantees that if the resource exists, it will be found (this is the only way to provide exact resource location in a decentralized resource location system), however, complete flooding is not scalable and incurs lots of message redundancy.

**Flooding with limited horizon** minimizes the message redundancy overhead by attaching a TTL to messages that limits the number of times a message can be retransmitted. However, there is a trade-off for efficiency: flooding with limited horizon does not provide exact resource location. There are many other dissemination techniques, often tailored towards specific application requirements.

#### 2.3.2.2 Random Walks

**Random Walks** are a dissemination strategy that attempts to minimize the communication overhead that accompanies flooding. Instead of flooding, a random walk consists of a message with a TTL that is randomly forwarded one peer at a time throughout the network. Walks may also take a biased paths in the system based on information accumulated by peers, through aggregation, this is called a **random guided walk**. Random guided walks forward queries to neighbors that are more likely to have answers [8].

A common approach to bias walks is to use bloom filters [44], which are space-efficient probabilistic data structures that support set membership queries. There are many other techniques of performing guided walks, often tailored for application needs.

### 2.3.3 Centralized Resource Location

Throughout the years 3 popular architectures emerged that are common towards indexing resources in a distributed system:

**Centralized architectures** rely on one (or a group of) centralized peers that index all resources in the system. This type of architecture greatly reduces the complexity of systems, as peers only need to contact a subset of nodes to locate resources. However, their scalability is limited, due to the centralized point of failure.

It is important to notice that in a centralized architecture, while the indexation of resources is centralized, the resource access may still be distributed (e.g. a centralized server provides the addresses of peers who have the files, and files are obtained in a pure P2P fashion). A system which employs this architecture with success is BitTorrent [7].

Because centralized architectures have limited scalability, purely centralized architectures cannot be applied in large scale Edge environments. However, there are many ways that a hybrid architecture can be applied to Edge computing: since the failure rate of a single DC is low, if we assume a system composed by multiple DCs, they may act as a reliable failover for whenever edge devices are partitioned or fail.

### 2.3.4 Decentralized Resource Location

#### 2.3.4.1 Distributed Hash tables

**Distributed Hash Tables** (DHTs) contrast with centralized servers, where the index distribution is split among peers in the system. In a DHT, peers are assigned uniformly distributed IDs using hash functions, then, peers employ a global coordination mechanism that restricts their neighboring relations (usually called routing tables) such that the resulting overlays commonly consist in low-diameter geometric structures like rings, hypercubes, among others.

Peers maintain routing tables to forward messages in the system, such that they can contact any other participant in a bounded number of steps, where the bound (usually logarithmic) is dictated by the topology. Finally, using the same hash functions to map resources (files, multimedia, messages, among others) to the peer identifier space, and assigning a key-space interval to each peer, peers can store and find any resource in a bounded number of steps (**exact resource location**).

One particular type of DHT that is commonly employed in small to medium sized storage solutions is the One-Hop DHT, nodes in a one-hop DHT have **Full membership** of the system and consequently, can perform lookups in  $O(1)$  time and message complexity. Facebook's Cassandra [23] and Amazon's Dynamo [10] are widely used implementations of one-hop DHTs. However, full membership solutions have scalability problems due to the required memory and message volume necessary to maintain the full membership information up-to-date. Finally, maintaining full membership is costly in the presence



churn (participants entering and leaving the system concurrently). The accumulation of these factors make full-membership solutions impractical in Edge environments.

Given this, the usual approach to building a DHT is through **partial membership** systems, which rely on some membership mechanism that restricts neighboring relations that are used to perform communication. DHTs with partial membership are attractive because they provide exact resource location while maintaining very little membership information (typically 1% of the peers), .

There have been attempts to apply DHTs towards Edge Computing. Common limitations that arise from this is that the common flat design that goes against the device heterogeneity of Edge Environments. Furthermore, given that devices in those environments have lower computational power and weaker connectivity, devices in the edge may even be a bottleneck to the system.

Following, we present some popular implementations of relevant DHT's along with a discussion on their applicability towards Edge environments:

**Chord** [42] is a distributed lookup protocol that addresses the need to locate the node that stores a particular data item, it specifies how to find the locations of keys, how nodes recover from failures, and how nodes join the system. Chord assigns each node and key an  $m$ -bit identifier that is uniformly distributed in the id space (peers receive roughly the same number of keys). Peers are ordered by identifier in a clockwise circle, then, any key  $k$  is assigned to the first peer whose identifier is equal or follows  $k$  in the identifier space.

Chord implements a system of "shortcuts" called the **finger table**. The finger table contains at most  $m$  entries, each  $i^{th}$  entry of this table corresponds to the first peer that succeeds a certain peer  $n$  by  $2^{i^{th}}$  in the circle. This means that whenever the finger table is up-to-date, lookups only take logarithmic time to finish.

Chord, although provides the best trade-off between bandwidth and lookup latency [28], however, chord presents some limitations: peers do not learn routing information from incoming requests and links have no correlation to latency or traffic locality.

Chord is a basis for lots of work: Cyclone [3] is a hierarchical version of Chord provides that constructs a hierarchy by splitting the ID space into a PREFIX and SUFIX. The PREFIX provides intra-cluster identity, whereas the SUFIX is used towards creating clusters of nodes. Routing procedures are executed in lower rings and move up the hierarchy.

Hieras [50] uses a binning scheme according the underlay topology to group peers into smaller rings. The lower the ring, the smaller the average link latency. Routing is similar to Cyclone.

Crescendo [13] splits the ID range into domains (similar to DNS), where nodes in leaf-domains form Chord rings, then nodes merge rings by applying rules such that rings in different domains can communicate. The resulting routing table and the routing procedures in Crescendo are similar to chord.

**Pastry** [37] is a DHT that assigns a 128-bit node identifier (nodeId) to each peer in the system. The nodes are randomly generated thus uniformly distributed in the 128-bit nodeId space.

Nodes store values whose keys are also distributed in the nodeId space, and Key-value pairs are stored among nodes which are numerically closest to the key.

The routing procedure is as follows: in each routing step, messages are forwarded to nodes whose nodeId shares a prefix that is at least one bit closer to the key. If there are no nodes available, Pastry routes messages towards the numerically closest nodeId.

This routing technique accomplishes routing in  $O(\log N)$ , where  $N$  is the number of Pastry nodes in the system. This protocol has been widely used and tested in applications such as Scribe [38] and PAST [11].

Limitations from using Pastry arise from the use of a numeric distance function towards the end of the routing process, which creates discontinuities at some node ID values, and complicates attempts at formal analysis of worst case behavior.

**Kademlia** [32] is a DHT with provable consistency and performance in a fault-prone environment. Kademlia nodes are assigned 160-bit identifiers uniformly distributed in the ID space.

Peers route queries and locate nodes by employing a novel **XOR-based distance** function which is symmetric and unidirectional. Each node in Kademlia is a router where its routing tables consist of shortcuts to peers whose XOR distance is between  $2^i$  by  $2^{i+1}$  in the ID space. Intuitively, and similar to Pastry, "closer" nodes are those that share a longer common prefix.

The main benefits that Kademlia draws from this approach are: nodes learn routing information from receiving messages, there is a single routing algorithm for the whole routing process (unlike Pastry) which eases formal analysis of worst-case behavior.

Finally, Kademlia exploits the fact that node failures are inversely related to uptime by prioritizing nodes that are already present in the routing table.

**Kelips** [17] exploits increased memory usage and constant background communication to achieve  $O(1)$  lookup time and message complexity.

Kelips nodes are split in  $k$  affinity groups split in the intervals  $[0, k-1]$  of the ID space, thus, with  $n$  nodes in the system, each affinity group contains  $\frac{n}{k}$  peers. Each node stores a partial set of nodes contained in the same affinity group and a small set of nodes lying in foreign affinity groups.

Assuming a proportional number of files and peers in the system and a fixed view of nodes in foreign affinity groups, Kelips achieves  $O(1)$  time and message complexity in lookups at the cost of increased memory consumption ( $O(\sqrt{n})$ ). Due to this, system scalability is limited when compared to Pastry, Chord or Kademlia.

**Tapestry** [49] Is a DHT similar to pastry where messages are incrementally forwarded to the destination digit by digit (e.g.  $***8 \rightarrow **98 \rightarrow *598 \rightarrow 4598$ ). Lookups have  $\log_b(n)$  time complexity where  $b$  is the base of the ID space. A system with  $n$  nodes has a resulting topology composed of  $n$  spanning trees, where each node is the root of its own tree. Because nodes assume that the preceding digits all match the current node's suffix, it only needs to keep a constant size of entries at each route level. Thus, nodes contain entries for a fixed-sized neighbor map of size  $b \cdot \log(N)$ .

#### 2.3.4.2 Unstructured Overlays

SOSP-Net (Self-Organizing Super-Peer Network) [14] proposes to optimize super-peer networks by organizing super-peers in a topology that reflects the semantic similarity of peers with similar interests. Super-peers maintain cache references to files which were recently requested by weak peers, while the weak peers maintain caches containing super-peers that satisfied most of its requests. This yields good cache hit ratios, and load balancing. However, the original paper does not specify how super-peers communicate among themselves, nor the election process of super-peers.

#### 2.3.4.3 Hybrid approaches

**Curiata & Build One Get One Free**

#### 2.3.5 Discussion

### 2.4 Resource monitoring

Monitoring is key to the effective management of applications.

#### 2.4.1 Device monitoring

In order to adapt edge computing applications to changes in the environment and ensure that the above requirements can be met, it is necessary to tailor the monitoring system to support the whole spectrum of underlying infrastructures (section 2.1).

A modern approach towards deploying systems is to deploy them in loosely coupled independent components running some form of virtualization software, as it enables co-deployment of logical machines in the same physical node.

##### 2.4.1.1 VM monitoring

In virtual machines (VMs), all the physical resources can be virtualized (CPU, memory, disk and network). Multiple VMs can co-deployed in the same physical node and thus share resources among each other.

In order to have efficient resource utilization and prevent any problems in the virtualized resources, monitoring of VMs is critical. This can be best achieved by tracking the utilization of the virtualized resources, mainly usage of CPU, memory, storage and network.

- **CPU** usage tracks the amount of usage of the CPU as a percentage of all available CPU to the machine. It is never desirable that the CPU usage reaches 100%, as queues start filling up and it means that the device has run out of capacity for processing tasks.
- **Memory** indicates the amount of virtualized memory left.

- **Disk usage** tracks the amount of data read or written by a VM. Alternatively, it can indicate the percentage of used space.
- **Network usage** consists in the volume of traffic on a specific network interface of the VM, either external or internal traffic.

Although VMs are widely present in the cloud infrastructure, we believe their applicability towards edge scenarios is limited, due to the fact that they have significant start up time (having to start-up an OS).

#### 2.4.1.2 Container monitoring

Containers are the recent alternative to VMs, as previously mentioned, containers do not require an OS to boot up [], and the images that compose them are significantly smaller. Container-based virtualization can be compared to an OS running on bare-metal in terms of memory, CPU and disk usage, however at a cost of network utilization [34], which makes them an attractive options towards resource-constrained devices.

Due to their lightweight nature, it is possible to deploy container-based applications (e.g. microservices), which can perform fast migration across nodes in the edge environment in order to improve QoS. This flexibility towards the migration process is an efficient tool to deal with many challenges such as load balancing, scaling, resource reallocation and fault-tolerance.

There are many container-specific tools which provide statistics about a given container, most use REST APIs to expose this functionality to external entities.

Docker, which is a container provider, *citação?* has a built tool called *docker stats* which provides runtime metrics for a given container.

*Container Advisor* [15] (cAdvisor) is a service which analyzes and exposes resource usage and performance data from running containers. Simply put, cAdvisor consists of a daemon which collects, aggregates and exports information. The information it collects resource isolation parameters, historical resource usage and network statistics. cAdvisor includes native support for Docker containers and supports a wide variety of other container implementations.

#### 2.4.2 End-to-end link monitoring

Given that the edge infrastructure envisions cooperation from all devices in the path from the origin of the data to the DC, devices need to be interconnected across an underlying infrastructure which is continuously changing. This raises concerns about the network quality of links between devices across the system, especially if they are running time-critical services.

It is paramount to analyze how to monitor and improve link quality, for providing traffic locality, latency, among others. According to the literature [], the most popular metrics to analyze are:

1. Network throughput, which is the average rate of successful data transfer through a network connection.
2. Latency, which consists in how long a packet takes to travel across a link from one endpoint to another
3. Packet loss, which consists in how many packets are lost when traveling towards their destination.
4. Jitter is the variation in latency of sequential received packets.

### 2.4.3 Aggregation techniques

**Aggregation** is an essential building block towards monitoring distributed systems, it enables the determination of important system wide properties in a decentralized manner [22].

Aggregation consists in computing an aggregation function over a set of input values where each node has one input value. Common aggregation functions consist in sum, count, average, min, max.

Towards monitoring edge devices and tasks running on them, aggregation is paramount, examples of usages are (e.g. computing the average latency of the closest available service that meets a certain criteria; counting nearby available computing resources that can be used to offload services, or identify hotspots by aggregating the average system load in certain areas).

Given this, it is important to understand the taxonomy of aggregation functions. There are two properties of aggregation functions: *decomposability* and *duplicate sensitiveness*.

#### **Decomposability**

For some aggregation functions, we may need to involve all elements in the multiset, however, for memory and bandwidth issues, it is impractical to perform a centralized computation, hence, the aim is to employ *in-transit computation*. In order to enable this, it is required that the aggregation function is **decomposable**.

Intuitively, a decomposable aggregation function is one where a function may be composed defined as a composition of other functions. Decomposable functions may *self-decomposable*, which intuitively means that the aggregated value is the same for all possible combinations of all sub-multisets partitioned in the multiset. This happens whenever the applied function is commutative and associative (e.g. min, max, sum, count).

A canonical example of a decomposable function that is not self-decomposable is average, which consists in the sum of all pairs divided by the count of peers that contributed to the aggregation.

The second property of aggregation is **duplicate sensitiveness**, and it is related to whether a given value occurs several times in a multiset. Depending on the aggregation function used, the presence of repeated values may influence the result, it is said that a function is **duplicate sensitive** if the result of the aggregation function is influenced

	Decomposable		Non-Decomposable
	Self-decomposable		
Duplicate insensitive	Min, Max	Range	Distinct Count
Duplicate sensitive	Sum, Count	Average	Median, Mode

Tabela 2.2: popular aggregation functions in function of decomposability and duplicate sensitiveness

by the repeated values (e.g. SUM). Conversely, if the aggregation function is **duplicate insensitive**, it can be successfully repeated any number of times to the same multiset without affecting the result (e.g. MIN and MAX).

Table 2.2 classifies popular aggregation functions in function of decomposability and duplicate sensitiveness as found in [22]:

Building on the concepts of duplicate sensitiveness and decomposability, we show that aggregation functions present their own particularities which dictate their applicability in particular scenarios. For example, a Min or Max function may be easier to implement with a simpler algorithm, while Sum, Count and Average require extra considerations. This presents a limitation towards calculating exact aggregations in large scale systems, to circumvent this, some systems do not require obtaining exact aggregated values to perform near optimally (e.g. estimating the system size in order to select the optimal fanout for a gossip system only requires an estimation of the magnitude of the system).

#### 2.4.3.1 Aggregation techniques

Following, we present a taxonomy of aggregation techniques for edge devices, for each technique, we provide context and discuss its possible advantages and limitations in the edge environment.

#### 2.4.3.2 Hierarchical

**Tree-based** approaches leverage directly on the decomposability of aggregation functions. Aggregations from this class depend on the existence of a hierarchical communication structure, (e.g. a spanning tree) with one root (sink node). Aggregations take place by splitting inputs into groups and aggregating values bottom-up in the hierarchy.

Commonly, hierarchical aggregation systems have nodes whose roles are *aggregators* or *forwarders*, intuitively, aggregators compute the aggregation functions and forward results to forwarders who then retransmit the results to upper levels in the hierarchy. In the absence of faults, the correct final result is obtained in the sink node.

**Cluster-based** techniques rely on clustering the nodes in the network according to a certain criterion (e.g. latency, energy efficiency). In each cluster a representative is responsible for local aggregation and for transmitting the results to other nodes.

Hierarchical approaches, due to taking advantage of device heterogeneity, are attractive in edge environments. However, due to the low computational power of devices, not

all nodes may be able to handle the additional overhead of maintaining the hierarchical topology.

### 2.4.3.3 Averaging

Averaging aggregation consists in the continuous computation and exchanging of partial averages data among all active nodes in the aggregation process. In this type of systems, after a few rounds, all nodes usually converge to the correct value with high accuracy, as shown in [21].

This type of aggregation is attractive for gossip protocols, where nodes may employ varied gossip techniques to continuously share and update their values with random neighbors.

Algorithms from this category are also attractive to use in edge environments, because they are accurate while employing random unstructured overlays, which retain their fault-tolerance and resilience to churn.

**Sketches** are fixed-size data structures that hold a *sketch* of all network values. Multiple sketches are usually forwarded throughout the system, and nodes who forward sketches apply (usually commutative and associative) operations to update and merge them. [functioning and edge discussion](#)

**Digests** are an aggregation technique that gathers a representation of all system values, it supports complex aggregation functions such as Median and Mode. In short, algorithms employ a fixed-size data structures commonly composed of a set of values and associated counters) which compacts the data distribution (e.g. into a histogram). [edge discussion](#)

**Counting** algorithms target the same aggregation function: Count, algorithms from this class usually employ some randomized procedure to achieve a probabilistic approximation of the population size.

## 2.4.4 Relevant aggregation protocols

In this subsection we will analyze relevant aggregation protocols that illustrate some techniques discussed above.

### 2.4.4.1 TAG: Tiny AGgregation

**TAG: Tiny AGgregation**[29] is a service for aggregation in low-power, distributed, wireless sensor networks. TAG distributes queries in the network in a time and power-efficient manner by employing a hierarchical aggregation pattern. For each aggregation procedure, there is a *root* nodes which broadcasts a message to start the tree-building process, each message contains two fields: a level and a an ID. Whenever a node without an assigned level receives a tree-building message, it assigns its own level as the message level plus one, and its own parent as the message sender. Then, it reassigns the level and ID to its



own and forwards the message to other nodes. Then, whenever a node wishes to send a message to the root, it simply forwards the message bottom-up in the tree. The formed topology allows the computation of Count, Maximum, Minimum, Sum and Average. It is important to notice that the formed tree will be unbalanced as a function of the underlay latency and processing time.

#### 2.4.4.2 SingleTree

**SingleTree** [] //TODO

#### 2.4.4.3 MultipleTree

**MultipleTree** [] //TODO

#### 2.4.4.4 DECA

**DECA** [0] //TODO

### 2.4.5 Monitoring systems

#### 2.4.5.1 Astrolabe

**Astrolabe** [36] is a platform which aims at monitoring the dynamically changing state of a collection of distributed resources.

Astrolabe introduces a hierarchical architecture defined by zones, a zone is recursively defined to be either a host or a set of non-overlapping zones. Each zone (minus the root zone) has a local identifier, which is unique within the zone where it is contained, zones are globally identified by their *zone name*, which consists of the concatenation of all zone identifiers within the path from the root to the zone.

Associated with each zone there is a Management Information Base (MIB), which consists in a set of attributes from that zone. Zone attributes are not directly writable, instead, they are generated by aggregation functions contained in special entries in the MIB. Leaf zones are the exception to the aforementioned mechanism, leaf zones contain *virtual child zones* which are directly writable by devices within that virtual child zone.

The aggregation functions which produce the MIBs are contained in *aggregation function certificates* (AFCs), which contain a user-programmable SQL function, a timestamp and a digital signature. In addition to the aforementioned code, AFCs may contain other information. An *Information Request AFC*, in addition to the code, specifies which information to retrieve from each participating host, and how to summarize the retrieved information. Alternatively, we may have a *Configuration AFC* which specifies runtime parameters that applications may use for dynamic configuration.

Astrolabe employs gossip which provides an eventual consistency model: if updates cease to exist for a long enough time, all the elements of the system converge towards the same state. This is achieved by employing a gossip algorithm which selects another



agent at random and exchanges zone state with it. If the agents are within the same zone, they simply exchange information relative to their zone. Conversely, if agents are in different zones, they exchange information relative to the zone which is their least common ancestor.

Not all nodes gossip information, within each zone, a node is elected (the authors do not specify how) to perform gossip on behalf of that zone, additionally, nodes can represent nodes from other zones, in this case nodes run one instance of the gossip protocol per zone it represents. The number of represented zones is bounded by the number of levels in the Astrolabe tree.

Astrolabe relies on clock time from the devices enforce an order among MIBs. Originally, it compared timestamps contained in MIBs to impose an order among them, however, as the system scaled, this approach is became unreliable. To circumvent this, Astrolabe agents store the last MIB from each agent that gossiped with them, and only compare the timestamps of MIBs which originated from the same representative.

Given this, attribute updates may not be monotonic, this happens due to the fact that aggregations do not take into account the time where the leaf attribute was updated.

A agents' zone is defined by the system administrator, which is a potential limitation towards scalability, given that configuration errors have the potential to heavily reduce system latency and traffic locality. Additionally, the original authors state that the size of gossip messages scales with the branching factor, often exceeding the maximum size of a UDP packet.

#### 2.4.5.2 Ganglia

**Ganglia** is a scalable distributed monitoring system [31] for high performance computing systems, namely clusters and grids. In short, Ganglia groups nodes in clusters, in each cluster, there are representative cluster nodes which federate devices and aggregate internal cluster state.

Ganglia relies on IP multicast to perform intra-cluster aggregation, it is mainly designed to monitor infrastructure monitoring data about machines in a high-performance computing cluster. Then, it displays the time series data in a Web interface.

As previously mentioned, Ganglia is aimed towards monitoring high-performance grids. And its applicability is limited towards edge environments: (1) clusters are situated in stable environments, which contrasts with the edge environment; (2) it relies on IP multicast (which has been proven not to hold in a number of cases) to perform automatic discovery of nodes as they are added and removed within clusters; (3) has no mechanism to prevent network congestion

#### 2.4.5.3 SDIMS

[0] proposes a scalable distributed information management system (SDIMS) which leverages on distributed hash tables to create scalable aggregation trees. [completar...](#)

#### 2.4.5.4 Prometheus

**Prometheus** [35] is an open-source monitoring and alerting toolkit originally built at SoundCloud. Prometheus works well for recording any purely numeric time series. It supports machine-centric monitoring as well as monitoring of highly dynamic service-oriented architectures.

Prometheus is especially useful towards querying and collecting multi-dimensional data collections. Furthermore, it offers a platform towards configuring alerts, which are triggered whenever a certain criteria is met.

Prometheus allows a prometheus server to scrape selected time-series from another Prometheus server. This allows deploying scalable monitoring setups or scraping related metrics from another service's Prometheus into another. Prometheus federation is split in two categories, *hierarchical federation* and *cross-service federation*

In *hierarchical federation*, prometheus servers are organized into a topology which resembles a tree, where each server aggregates aggregated time series data from a larger number of subordinated servers. Alternatively, *cross-service federation* enables scraping selected data from another service's prometheus server to enable alerting and queries against both datasets within a single server.

Prometheus employs a variety of service discovery options for discovering and scraping targets, including Kubernetes, Consul, among others. A key feature which improves its flexibility is the option to use file-based discovery, which enables the use of customized service discovery. This is achieved by programmers by providing a custom file with list of targets (and target metadata) which prometheus will scrape.

Mais papers sobre monitoring, especialmente na edge

## 2.5 Resource management

It is paramount to understand how resource monitoring systems can be tailored in order to meet the need of resource monitoring systems. Currently, most popular resource management systems are tailored towards Cloud systems [45] [18] [0] [0].

We classify resource management systems as platforms which provide resources (e.g. computing power, memory, among other) to tenants (i.e. applications, frameworks, among others) such that applications can perform their tasks.

In this section we present an overview of the taxonomy of resource management architectures, along with relevant resource management systems used in Cloud and Edge environments. In each example, we analyze the type of monitoring system employed.

### 2.5.1 Resource Management Taxonomy

A resource management system aims at controlling the distribution of resources among application(s). We may classify resource management architectures according to their

*control* and *tenancy*. **Control** refers to how the resources are provisioned towards applications, whereas *tenancy* refers to how whether applications share the system resources or not.

### 2.5.1.1 Tenancy

As previously mentioned, *tenancy* refers to whether entities share system resources or not. The term *tenancy* in distributed systems refers to whether or not underlying hardware resources are shared among entities [19].

**Single tenancy** refers to an architecture in which a single instance of a software application and supporting infrastructure serves one customer. In single-tenancy architectures, a customer (tenant) has nearly full control over customization of software and infrastructure.

Intuitively, there are many advantages of having single-tenancy architectures, such as data security: Even if there is a data breach to one tenant with the same provider, another tenant would be safe from the breach, since the data is stored in a separate instance.

However, single tenancy goes against most distributed systems today, which are inherently *shared* by multiple tenants, both on private and public clouds [0].

**Multi-tenancy** consists in tenants sharing multiple resources across multiple processes and machines. This approach has clear advantages, as sharing the infrastructure leads to lower costs (e.g. electricity), and companies of all sizes like to share infrastructure and in order to achieve lower operational costs.

However, providing performance guarantees and isolation in multi-tenant systems is extremely hard, resource management systems must avoid mismatching the resource allocation, as tenant-generated requests compete with each other and with the system generated tasks. Furthermore, tenant workload can change in unpredictable ways depending on the input workload, the workload of other tenants in the system and the underlying topology. Given this, mismanaging the system resources may degrade tenant performance and in extreme cases even hinder the stability of the resource management system.

Furthermore, in the context of edge-enabled workloads, some may have critical latency or availability requirements, which are extremely hard to ensure for a single tenant.

### 2.5.1.2 Control

**Control** refers to how resource management systems allocates tasks. There are two alternatives towards allocating resources.

**Centralized control** consists of a centralized controller deciding the computations, networks and communication of devices. This means that there is a centralized component in the system which has a global view of the resource usage within both and across processes and machines.

Intuitively, given that a centralized component generates manages all the resources in the system, this component can enforce policies that achieve the desired performance guarantees or fairness goals by identifying and only throttling the tenants or system activities responsible for resource bottlenecks. [0] Additionally, centralized control ensures that the system has the potential to attribute the best possible workload scheduling which ensures the requirements of tenants while employing the available resources as efficiently as possible.

Many resource management systems rely on a centralized component having a global view of the system, especially in Cloud environments, where nodes are federated in large clusters of machines with similar computing power. Cloud systems commonly employ Zookeeper [0] to maintain the component available and to transfer state in case it fails.

However, we argue that centralized control is not feasible in an edge environment, as the number of devices in the system increases, so does the number of resources to track, and the harder it is for a centralized component to have an up-to-date global view of the system. This occurs either because the system is heavily decentralized which means that monitoring information takes a long time to propagate to the centralized component, or because there is not enough computing power or memory to handle all the information in a centralized component.

**Decentralized control** architectures are those where the the decision-making process is distributed across multiple components [19]. This topic has yet not been subject to much research, although it is of extreme relevance towards edge environments. For example, if the system is globally distributed, it may take too long for a centralized controller to identify hotspots in a certain zone and load-balance them.

One of the biggest challenges of controlling the resource management in a distributed fashion is ensuring that the components which perform resource assignments do not conflict with each other. Additionally, in a multi-tenant decentralized resource control system, tenants may request resources to different resource controllers in the system, and if they do not coordinate themselves, the application may be provisioned with too many (or too little) resources.

## 2.5.2 Relevant Resource Monitoring Systems

In this subsection we will cover popular resource monitoring systems found in the literature, and discuss their advantages and / or limitations.

### 2.5.2.1 Mesos

Mesos [18] is a multi-tenant centralized resource sharing platform which attempts to provide fine-grained resource sharing in the Data Center. The tenants for this platform are frameworks such as HDFS [0], MapReduce [0], among others, which in turn support multiple applications running in the DC.

In short, the Mesos resource sharing system consists of a *master* process which manages *slave* daemons running on each cluster node, in order to achieve fault-tolerance, Mesos employs Zookeeper [0] for fault tolerance (by electing a new master and transferring state to it if the active master crashes).

The master node implements fine-grained sharing of resources across frameworks by employing *resource offers*, which consist of lists of free resources distributed among slaves.

Resource offers provide flexibility for frameworks operating on Mesos, given that each framework has different communication patterns, task dependencies and data placement, a centralized scheduler would need to provide an extremely expressive API to capture all frameworks' requirements. Instead, the master makes decisions about how many resources to offer to each framework, and its decision-making process is based on an arbitrary organizational policy, such as fair sharing or priority.

Each framework that wishes to use Mesos must implement *scheduler* and an *executor*. The scheduler registers with the Mesos master to receive resource offers, and the executor is the process that is launched on slave nodes to run the framework's tasks.

A limitation of the Mesos resource sharing platform is that it is not scalable (the original authors mention the system scales up to 50000 slave daemons on 99 physical machines), which is not enough for an edge environment. Furthermore, the resource offer model forces frameworks to employ a specific programming model based on schedulers and executors, which may not be desirable by the framework author.

### 2.5.2.2 Borg

### 2.5.2.3 Yarn

Yarn (Yet Another Resource Negotiator) [45] is a compute platform which attempts decouple the programming model from the resource management infrastructure, and delegate many scheduling functions (e.g. task fault-tolerance) to per-application components.

The architecture of YARN is as follows: the system is composed of a per-cluster Resource Manager (RM), multiple Application Masters (AM) and Node Managers (NM).

The Resource Manager (RM) tracks resource usage and node liveness, enforces allocation invariants and arbitrates contention among tenants. It matches a global model of the cluster state against the digest of resource requirements reported by the applications, which allows the RM to tightly enforce global scheduling properties such as capacity or fairness.

The Application Master (AM) runs arbitrary user code, and can be written in any programming language, its duties in the system consist of managing the lifecycle aspects including dynamically increasing and decreasing resource consumption, managing the flow of execution and handling faults.

The Node Manager (NM) is the worker daemon in YARN. Its responsibilities in the system consist of managing container dependencies, monitor their execution and provide a set of services for containers. Node Managers

AMs send resource requests to the RM, these contain the number of containers to request, the resources per container, locality preferences and a priority level within the application. Resource requests are designed to capture the needs of applications, while at the same time removing specific application concerns (such as task dependencies) from the scheduler.

Similar to MESOS [18], the resource attribution process is done by centralized component, however, contrasting with MESOS which uses Zookeeper to maintain secondary masters, the RM is a single point of failure and contention in the system, which makes it not suitable towards edge environments.

#### 2.5.2.4 ENORM

Edge NNode Resource Management [0] (ENORM) framework aimed at employing edge resources towards applications by provisioning and auto-scaling edge node resources. It attempts to answer these two questions: "How to deploy a workload on the edge node?" and "How much of the workload can be deployed on the edge node?".

ENORM proposes a three-tier architecture: (1) the Cloud tier, where application servers are hosted; (2) the middle tier, where the edge nodes are situated; (3) the bottom tier, where user devices (e.g. smartphones, wearables, gadgets) are situated.

To enable the use edge nodes, ENORM deploys a *cloud server manager* on each application server which communicates with potential edge nodes requesting computing services, deploys partitioned servers on the edge nodes, and finally receives updates from the edge nodes to update the global view of the application server on the cloud.

Each edge node is composed by the following 5 components: (1) *Resource Allocator* which provides computing as a service, this is the basic service of the edge node, and cannot be compromised, which means that it has priority over offloaded computations; (2) *Edge Manager* deals with the requests that are obtained by the server manager in the cloud, and once a request is successful, it initializes a container and allocates the necessary ports for communication; (3) *Monitor* which periodically reports a number of metrics related to the application running on the edge server (e.g. communication latency and computing latency); (4) *Auto-scaler* dynamically allocates and de-allocates hardware resources to the container executing application servers. It is responsible for scaling the process such that the application can accommodate a larger number of users. (5) *Application Edge Server* the partitioned server which is hosted on the edge node.

ENORM authors test the designed system using an online game based on Pokemon GO (iPokemon)[0], the framework partitions the game server and sends user data relevant to the geographical location of the edge node. Users from the relevant geographical zone connect to the edge server and are serviced as if they were connected to the data center. ENORM observes up to 95% application latency reduction, and a reduction of data transferred between edge nodes and the cloud of up to 95%.

Limitations from this framework are similar to Mesos and Yarn, specifically lack of

fault-tolerance and scalability, which arise from a centralized component handling all the edge nodes at the same time.

### 2.5.3 Discussion

Although resource management systems have been present for many years, these are often tailored towards small scale environments composed by high-capacity devices in stable environments. This environment contrasts with the edge of the network, where devices are extremely numerous, decentralized, and heterogenous.

Due to their low capacity, devices in the edge of the network are very susceptible to workload changes, for example, a 5G tower which is hosting a service cannot handle a drastic increase in the number of users it is servicing. In this scenario, we argue that in order to maintain QoS, devices must autonomously take the decision to scale horizontally or vertically in order to adapt to the hotspot of users.

Additionally, with the recent emergence of IoT and the increase of devices in the edge of the network, we argue that a centralized controller would not be able to process the resource management of a large scale system, especially if applications have special requirements such as real-time latency.





### 3.1 Proposed Solution

To achieve this, we propose to create a new novel algorithm which employs a hierarchical topology that resembles the device distribution of the Edge Infrastructure. This topology is created by assigning a level to each device and leveraging on gossip mechanisms to build a structure resembling a FAT-tree [].

The levels of the tree will be determined by **...undecided...** and will the tree be used to employ efficient aggregation and search algorithms. Each level of the tree will be composed by many devices that form groups among themselves, the topology of the groups **...undecided...**

The purpose of this algorithm is to allow:

1. Efficient resource monitoring to deploy services on.
2. Offloading computation from the cloud to the Edge and vice-versa through elastic management of deployed services.
3. Service discovery enabled by efficiently searching over large amount of devices
4. Federate large amount of heterogeneous devices and use heterogeneity as an advantage for building the topology.

We plan to research existing protocols (both for topology management and aggregation) and enumerate their trade-offs along with how they behave across different environments. Then, employ a combination of different techniques in a unique way to monitor the topology.

## **3.2 Scheduling**

## BIBLIOGRAFIA

- [1] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica e et al. “A View of Cloud Computing”. Em: *Commun. ACM* 53.4 (abr. de 2010), 50–58. ISSN: 0001-0782. DOI: [10.1145/1721654.1721672](https://doi.org/10.1145/1721654.1721672). URL: <https://doi.org/10.1145/1721654.1721672>.
- [0] M. S. Artigas, P. García e A. F. Skarmeta. “DECA: A hierarchical framework for DE-Centralized aggregation in DHTs”. Em: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 4269 LNCS (2006), pp. 246–257. ISSN: 16113349. DOI: [10.1007/11907466\\_23](https://doi.org/10.1007/11907466_23).
- [3] M. S. Artigas, P. G. López, J. P. Ahulló e A. F. Skarmeta. “Cyclone: A novel design schema for hierarchical DHTs”. Em: *Proceedings - Fifth IEEE International Conference on Peer-to-Peer Computing, P2P 2005*. Vol. 2005. 2005, pp. 49–56. ISBN: 0769523765. DOI: [10.1109/P2P.2005.5](https://doi.org/10.1109/P2P.2005.5).
- [4] F. Bonomi, R. Milito, J. Zhu e S. Addepalli. “Fog computing and its role in the internet of things”. Em: *Proceedings of the first edition of the MCC workshop on Mobile cloud computing*. 2012, pp. 13–16.
- [0] D. Borthakur et al. “HDFS architecture guide”. Em: *Hadoop Apache Project* 53.1-13 (2008), p. 2.
- [0] *Catch Pokémon in the Real World with Pokémon GO!* URL: <https://www.pokemongo.com/en-us/>.
- [6] Y. Chawathe, S. Ratnasamy, L. Breslau, N. Lanham e S. Shenker. “Making Gnutella-like P2P Systems Scalable”. Em: *Computer Communication Review* 33.4 (2003), pp. 407–418. ISSN: 01464833. DOI: [10.1145/863997.864000](https://doi.org/10.1145/863997.864000).
- [7] B. Cohen. “Incentives build robustness in BitTorrent”. Em: *Workshop on Economics of Peer-to-Peer systems*. Vol. 6. 2003, pp. 68–72.
- [8] A. Crespo e H. Garcia-Molina. “Routing indices for peer-to-peer systems”. Em: *Proceedings 22nd International Conference on Distributed Computing Systems*. 2002, pp. 23–32. DOI: [10.1109/ICDCS.2002.1022239](https://doi.org/10.1109/ICDCS.2002.1022239).
- [0] J. Dean e S. Ghemawat. “MapReduce: simplified data processing on large clusters”. Em: *Communications of the ACM* 51.1 (2008), pp. 107–113.

- [10] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall e W. Vogels. “Dynamo: amazon’s highly available key-value store”. Em: *ACM SIGOPS operating systems review*. Vol. 41. 6. ACM. 2007, pp. 205–220.
- [11] P. Druschel e A. Rowstron. “PAST: a large-scale, persistent peer-to-peer storage utility”. Em: *Proceedings Eighth Workshop on Hot Topics in Operating Systems*. 2001, pp. 75–80. DOI: [10.1109/HOTOS.2001.990064](https://doi.org/10.1109/HOTOS.2001.990064).
- [12] M. Finnegan. *Boeing 787s to create half a terabyte of data per flight, says Virgin Atlantic*. 2013. URL: <https://www.computerworld.com/article/3417915/boeing-787s-to-create-half-a-terabyte-of-data-per-flight--says-virgin-atlantic.html>.
- [13] P. Ganesan, K. Gummadi e H. Garcia-Molina. “Canon in G major: Designing DHTs with hierarchical structure”. Em: *Proceedings - International Conference on Distributed Computing Systems* 24 (2004), pp. 263–272. DOI: [10.1109/icdcs.2004.1281591](https://doi.org/10.1109/icdcs.2004.1281591).
- [14] P. Garbacki, D. H. Epema e M. Van Steen. “Optimizing peer relationships in a super-peer network”. Em: *27th International Conference on Distributed Computing Systems (ICDCS’07)*. IEEE. 2007, pp. 31–31.
- [15] Google. *google/cadvisor*. 2020. URL: <https://github.com/google/cadvisor>.
- [16] *Gtk-Gnutella*. 2019. URL: <https://sourceforge.net/projects/gtk-gnutella/>.
- [17] I. Gupta, K. Birman, P. Linga, A. Demers e R. Van Renesse. “Kelips: Building an efficient and stable P2P DHT through increased memory and background overhead”. Em: *International Workshop on Peer-to-Peer Systems*. Springer. 2003, pp. 160–169.
- [18] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. H. Katz, S. Shenker e I. Stoica. “Mesos: A platform for fine-grained resource sharing in the data center.” Em: *NSDI*. Vol. 11. 2011. 2011, pp. 22–22.
- [19] C. H. Hong e B. Varghese. “Resource management in fog/Edge computing: A survey on architectures, infrastructure, and algorithms”. Em: *ACM Computing Surveys* 52.5 (2019). ISSN: 15577341. DOI: [10.1145/3326066](https://doi.org/10.1145/3326066).
- [0] P. Hunt, M. Konar, F. P. Junqueira e B. Reed. “ZooKeeper: Wait-free Coordination for Internet-scale Systems.” Em: *USENIX annual technical conference*. Vol. 8. 9. 2010.
- [20] M. Jelasity e O. Babaoglu. “T-Man: Gossip-based overlay topology management”. Em: *International Workshop on Engineering Self-Organising Applications*. Springer. 2005, pp. 1–15.
- [21] M. Jelasity, A. Montresor e O. Babaoglu. “Gossip-Based Aggregation in Large Dynamic Networks”. Em: *ACM Transactions on Computer Systems* 23 (ago. de 2005), pp. 219–252. DOI: [10.1145/1082469.1082470](https://doi.org/10.1145/1082469.1082470).

- [22] P. Jesus, C. Baquero e P. S. Almeida. “A Survey of Distributed Data Aggregation Algorithms”. Em: *CoRR* abs/1110.0725 (2011). arXiv: [1110.0725](https://arxiv.org/abs/1110.0725). URL: <http://arxiv.org/abs/1110.0725>.
- [23] A. Lakshman e P. Malik. “Cassandra: a decentralized structured storage system”. Em: *ACM SIGOPS Operating Systems Review* 44.2 (2010), pp. 35–40.
- [24] J. Leitaó, J. Pereira e L. Rodrigues. “HyParView: A Membership Protocol for Reliable Gossip-Based Broadcast”. Em: *37th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN’07)*. 2007, pp. 419–429. DOI: [10.1109/DSN.2007.56](https://doi.org/10.1109/DSN.2007.56).
- [25] J. Leitaó, J. Pereira e L. Rodrigues. “Epidemic broadcast trees”. Em: *2007 26th IEEE International Symposium on Reliable Distributed Systems (SRDS 2007)*. IEEE. 2007, pp. 301–310.
- [0] J. Leitão, J. P. Marques, J. Pereira e L. Rodrigues. “X-bot: A protocol for resilient optimization of unstructured overlay networks”. Em: *IEEE Transactions on Parallel and Distributed Systems* 23.11 (2012), pp. 2175–2188.
- [26] J. Leitão, P. Á. Costa, M. C. Gomes e N. Preguiça. “Towards Enabling Novel Edge-Enabled Applications”. Em: 732505 (2018). arXiv: [1805.06989](https://arxiv.org/abs/1805.06989). URL: <http://arxiv.org/abs/1805.06989>.
- [27] J. C. A. Leitão e L. E. T. Rodrigues. “Overnesia: a resilient overlay network for virtual super-peers”. Em: *2014 IEEE 33rd International Symposium on Reliable Distributed Systems*. IEEE. 2014, pp. 281–290.
- [28] J. Li, J. Stribling, T. Gil, R. Morris e M. Kaashoek. “Comparing the Performance of Distributed Hash Tables Under Churn”. Em: mar. de 2004. DOI: [10.1007/978-3-540-30183-7\\_9](https://doi.org/10.1007/978-3-540-30183-7_9).
- [0] J. Mace, P. Bodik, R. Fonseca e M. Musuvathi. “Retro: Targeted resource management in multi-tenant distributed systems”. Em: *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*. 2015, pp. 589–603.
- [29] S. Madden, M. J. Franklin, J. Hellerstein e W. Hong. “{TAG}: {Tiny} {AGgregate} Queries in Ad-Hoc Sensor Networks”. Em: *Proceedings of the {USENIX} Symposium on Operating Systems Design and Implementation* (2002), pp. 131–146. URL: [xxx](#).
- [30] Y. Mao, C. You, J. Zhang, K. Huang e K. Letaief. “A Survey on Mobile Edge Computing: The Communication Perspective”. Em: *IEEE Communications Surveys & Tutorials* PP (ago. de 2017), pp. 1–1. DOI: [10.1109/COMST.2017.2745201](https://doi.org/10.1109/COMST.2017.2745201).
- [31] M. L. Massie, B. N. Chun e D. E. Culler. “The ganglia distributed monitoring system: design, implementation, and experience”. Em: *Parallel Computing* 30.7 (2004), pp. 817–840.

- [32] P. Maymounkov e D. Mazieres. “Kademlia: A peer-to-peer information system based on the xor metric”. Em: *International Workshop on Peer-to-Peer Systems*. Springer. 2002, pp. 53–65.
- [33] G. Peng. “CDN: Content distribution network”. Em: *arXiv preprint cs/0411069* (2004).
- [34] E. Preeth, F. J. P. Mulerickal, B. Paul e Y. Sastri. “Evaluation of Docker containers based on hardware utilization”. Em: *2015 International Conference on Control Communication & Computing India (ICCC)*. IEEE. 2015, pp. 697–700.
- [35] Prometheus. *From metrics to insight*. URL: <https://prometheus.io/>.
- [36] R. V. A. N. Renesse, K. P. Birman e W. Vogels. “Astrolabe: A Robust and Scalable Technology for Distributed System Monitoring, Management, and Data Mining”. Em: *ACM Transactions on Computer Systems* 21.2 (2003), pp. 164–206.
- [37] A. Rowstron e P. Druschel. “Pastry: Scalable, decentralized object location, and routing for large-scale peer-to-peer systems”. Em: *IFIP/ACM International Conference on Distributed Systems Platforms and Open Distributed Processing*. Springer. 2001, pp. 329–350.
- [38] A. Rowstron, A.-M. Kermarrec, M. Castro e P. Druschel. “Scribe: The Design of a Large-Scale Event Notification Infrastructure”. Em: *Networked Group Communication*. Ed. por J. Crowcroft e M. Hofmann. Berlin, Heidelberg: Springer Berlin Heidelberg, 2001, pp. 30–43. ISBN: 978-3-540-45546-2.
- [39] *Self-driving Cars Will Create 2 Petabytes Of Data, What Are The Big Data Opportunities For The Car Industry?* URL: <https://datafloq.com/read/self-driving-cars-create-2-petabytes-data-annually/172>.
- [40] W. Shi, J. Cao, Q. Zhang, Y. Li e L. Xu. “Edge Computing: Vision and Challenges”. Em: *IEEE Internet of Things Journal* 3.5 (2016), pp. 637–646. ISSN: 2372-2541. DOI: [10.1109/JIOT.2016.2579198](https://doi.org/10.1109/JIOT.2016.2579198).
- [41] W. Shi, J. Cao, Q. Zhang, Y. Li e L. Xu. “Edge Computing: Vision and Challenges”. Em: *IEEE Internet of Things Journal* 3 (out. de 2016), pp. 1–1. DOI: [10.1109/JIOT.2016.2579198](https://doi.org/10.1109/JIOT.2016.2579198).
- [42] I. Stoica, R. Morris, D. Liben-Nowell, D. R. Karger, M. F. Kaashoek, F. Dabek e H. Balakrishnan. “Chord: a scalable peer-to-peer lookup protocol for internet applications”. Em: *IEEE/ACM Transactions on Networking (TON)* 11.1 (2003), pp. 17–32.
- [43] D. Stutzbach e R. Rejaie. “Understanding churn in peer-to-peer networks”. Em: *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*. ACM. 2006, pp. 189–202.

- 
- [44] S. Tarkoma, C. E. Rothenberg e E. Lagerspetz. “Theory and Practice of Bloom Filters for Distributed Systems”. Em: *IEEE Communications Surveys Tutorials* 14.1 (2012), pp. 131–155. ISSN: 2373-745X. DOI: [10.1109/SURV.2011.031611.00024](https://doi.org/10.1109/SURV.2011.031611.00024).
- [45] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, B. Saha, C. Curino, O. O’Malley, S. Radia, B. Reed e E. Baldeschwieler. “Apache Hadoop YARN: yet another resource negotiator”. Em: SOCC ’13. 2013.
- [46] T. Verbelen, P. Simoens, F. De Turck e B. Dhoedt. “Cloudlets: Bringing the Cloud to the Mobile User”. Em: *Proceedings of the Third ACM Workshop on Mobile Cloud Computing and Services*. MCS ’12. Low Wood Bay, Lake District, UK: Association for Computing Machinery, 2012, 29–36. ISBN: 9781450313193. DOI: [10.1145/2307849.2307858](https://doi.org/10.1145/2307849.2307858). URL: <https://doi.org/10.1145/2307849.2307858>.
- [0] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune e J. Wilkes. “Large-scale cluster management at Google with Borg”. Em: *Proceedings of the Tenth European Conference on Computer Systems*. 2015, pp. 1–17.
- [47] M. Villari, M. Fazio, S. Dustdar, O. Rana e R. Ranjan. “Osmotic computing: A new paradigm for edge/cloud integration”. Em: *IEEE Cloud Computing* 3.6 (2016), pp. 76–83.
- [0] N. Wang, B. Varghese, M. Matthaiou e D. S. Nikolopoulos. “ENORM: A framework for edge node resource management”. Em: *IEEE transactions on services computing* (2017).
- [0] P. Yalagandula e M. Dahlin. “A Scalable Distributed Information Management System”. Em: *SIGCOMM Comput. Commun. Rev.* 34.4 (ago. de 2004), 379–390. ISSN: 0146-4833. DOI: [10.1145/1030194.1015509](https://doi.org/10.1145/1030194.1015509). URL: <https://doi.org/10.1145/1030194.1015509>.
- [49] B. Zhao, L. Huang, J. Stribling, S. Rhea, A. Joseph e J. Kubiatowicz. “Tapestry: A Resilient Global-Scale Overlay for Service Deployment”. Em: *IEEE Journal on Selected Areas in Communications* 22 (jul. de 2003). DOI: [10.1109/JSAC.2003.818784](https://doi.org/10.1109/JSAC.2003.818784).
- [50] Zhiyong Xu, Rui Min e Yiming Hu. “HIERAS: a DHT based hierarchical P2P routing algorithm”. Em: *2003 International Conference on Parallel Processing, 2003. Proceedings*. 2003, pp. 187–194. DOI: [10.1109/ICPP.2003.1240580](https://doi.org/10.1109/ICPP.2003.1240580).

