

Dear Author/Editor,

Here are the proofs of your chapter as well as the metadata sheets.

### Metadata

- Please carefully proof read the metadata, above all the names and address.
- In case there were no abstracts for this book submitted with the manuscript, the first 10-15 lines of the first paragraph were taken. In case you want to replace these default abstracts, please submit new abstracts with your proof corrections.

### Page proofs

- Please check the proofs and mark your corrections either by
  - entering your corrections online
  - or
  - opening the PDF file in Adobe Acrobat and inserting your corrections using the tool "Comment and Markup"
  - or
  - printing the file and marking corrections on hardcopy. Please mark all corrections in dark pen in the text and in the margin at least  $\frac{1}{4}$ " (6 mm) from the edge.
- You can upload your annotated PDF file or your corrected printout on our Proofing Website. In case you are not able to scan the printout , send us the corrected pages via fax.
- Please note that any changes at this stage are limited to typographical errors and serious errors of fact.
- If the figures were converted to black and white, please check that the quality of such figures is sufficient and that all references to color in any text discussing the figures is changed accordingly. If the quality of some figures is judged to be insufficient, please send an improved grayscale figure.

# Metadata of the chapter that will be visualized online

Book Title	Cloud Computing	
Chapter Title	Beyond the Clouds: How Should Next Generation Utility Computing Infrastructures Be Designed?	
Copyright	Springer International Publishing Switzerland 2014	
Corresponding Author	Prefix	
	Family name	<b>Bertier</b>
	Particle	
	Given name	<b>Marin</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Marin.Bertier@inria.fr
Author	Prefix	
	Family name	<b>Desprez</b>
	Particle	
	Given name	<b>Frédéric</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Frederic.Desprez@inria.fr
Author	Prefix	
	Family name	<b>Fedak</b>
	Particle	
	Given name	<b>Gilles</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Gilles.Fedak@inria.fr
Author	Prefix	
	Family name	<b>Lebre</b>
	Particle	
	Given name	<b>Adrien</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Adrien.Lebre@inria.fr
Author	Prefix	
	Family name	<b>Orgerie</b>
	Particle	
	Given name	<b>Anne-Cécile</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Anne-Cecile.Orgerie@inria.fr

Author	Prefix	
	Family name	<b>Pastor</b>
	Particle	
	Given name	<b>Jonathan</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Jonathan.Pastor@inria.fr
Author	Prefix	
	Family name	<b>Quesnel</b>
	Particle	
	Given name	<b>Flavien</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Flavien.Quesnel@inria.fr
Author	Prefix	
	Family name	<b>Rouzaud-Cornabas</b>
	Particle	
	Given name	<b>Jonathan</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Jonathan.Rouzaud-Cornabas@inria.fr
Author	Prefix	
	Family name	<b>Tedeschi</b>
	Particle	
	Given name	<b>Cédric</b>
	Suffix	
	Division	Inria
	Organization	Campus universitaire de Beaulieu
	Address	35042 Rennes, France
	Email	Cedric.Tedeschi@inria.fr
Abstract	To accommodate the ever-increasing demand for Utility Computing (UC) resources while taking into account both energy and economical issues, the current trend consists in building even larger data centers in a few strategic locations. Although, such an approach enables to cope with the actual demand while continuing to operate UC resources through centralized software system, it is far from delivering sustainable and efficient UC infrastructures. In this scenario, we claim that a disruptive change in UC infrastructures is required in the sense that UC resources should be managed differently, considering locality as a primary concern. To this aim, we propose to leverage any facilities available through the Internet in order to deliver widely distributed UC platforms that can better match the geographical dispersal of users as well as the unending resource demand. Critical to the emergence of such locality-based UC (LUC) platforms is the availability of appropriate operating mechanisms. We advocate the implementation of a unified system driving the use of resources at an unprecedented scale by turning a complex and diverse infrastructure into a collection of abstracted computing facilities that is both easy to operate and reliable. By deploying and using such a LUC Operating System on backbones, our ultimate vision is to make possible to host/operate a large part of the Internet by its internal structure itself: a scalable and nearly infinite set of resources delivered by any computing facilities forming the Internet, starting from the larger hubs operated by ISPs, governments, and academic institutions to any idle resources that may be provided by end users.	



# Chapter 14

## Beyond the Clouds: How Should Next Generation Utility Computing Infrastructures Be Designed?

Marin Bertier, Frédéric Desprez, Gilles Fedak, Adrien Lebre,  
Anne-Cécile Orgerie, Jonathan Pastor, Flavien Quesnel,  
Jonathan Rouzaud-Cornabas and Cédric Tedeschi

1   **Abstract** To accommodate the ever-increasing demand for Utility Computing (UC)  
2   resources while taking into account both energy and economical issues, the cur-  
3   rent trend consists in building even larger data centers in a few strategic locations.  
4   Although, such an approach enables to cope with the actual demand while continuing  
5   to operate UC resources through centralized software system, it is far from delivering  
6   sustainable and efficient UC infrastructures. In this scenario, we claim that a disrup-  
7   tive change in UC infrastructures is required in the sense that UC resources should be  
8   managed differently, considering locality as a primary concern. To this aim, we pro-  
9   pose to leverage any facilities available through the Internet in order to deliver widely

---

M. Bertier (✉) · F. Desprez · G. Fedak · A. Lebre · A.-C. Orgerie · J. Pastor · F. Quesnel ·

J. Rouzaud-Cornabas · C. Tedeschi

Inria, Campus universitaire de Beaulieu, 35042 Rennes, France

e-mail: Marin.Bertier@inria.fr

F. Desprez

e-mail: Frederic.Desprez@inria.fr

G. Fedak

e-mail: Gilles.Fedak@inria.fr

A. Lebre

e-mail: Adrien.Lebre@inria.fr

A.-C. Orgerie

e-mail: Anne-Cecile.Orgerie@inria.fr

J. Pastor

e-mail: Jonathan.Pastor@inria.fr

F. Quesnel

e-mail: Flavien.Quesnel@inria.fr

J. Rouzaud-Cornabas

e-mail: Jonathan.Rouzaud-Cornabas@inria.fr

C. Tedeschi

e-mail: Cedric.Tedeschi@inria.fr

© Springer International Publishing Switzerland 2014

Z. Mahmood, (ed.), *Cloud Computing*, Computer Communications and Networks,

DOI 10.1007/978-3-319-10530-7\_14



distributed UC platforms that can better match the geographical dispersal of users as well as the unending resource demand. Critical to the emergence of such locality-based UC (LUC) platforms is the availability of appropriate operating mechanisms. We advocate the implementation of a unified system driving the use of resources at an unprecedented scale by turning a complex and diverse infrastructure into a collection of abstracted computing facilities that is both easy to operate and reliable. By deploying and using such a LUC Operating System on backbones, our ultimate vision is to make possible to host/operate a large part of the Internet by its internal structure itself: a scalable and nearly infinite set of resources delivered by any computing facilities forming the Internet, starting from the larger hubs operated by ISPs, governments, and academic institutions to any idle resources that may be provided by end users.

**Keywords** Utility Computing • UC • Locality-based UC • Distributed Cloud Computing • IaaS • Efficiency • Sustainability

## 23      **14.1 Introduction**

**AQ1** The success of Cloud Computing has driven the advent of Utility Computing (UC). However, Cloud Computing is a victim of its own success. In order to answer the escalating demand for computing resources, Cloud Computing providers must build data centers (DCs) of ever-increasing size. As a consequence, besides facing the well-known issues of large-scale platform management, large-scale DCs have now to deal with energy considerations that limit the number of physical resources that one location can host.

Instead of investigating alternative solutions that could tackle the aforementioned concerns, the current trend consists in deploying larger and larger DCs in few strategic locations presenting energy advantages. For example, Western North Carolina, USA, an attractive area due to its abundant capacity of coal and nuclear power, brought about the departure of the textile and furniture industry [21]. More recently, several proposals suggested building next generation DCs close to the polar circle in order to leverage free cooling techniques, considering that cooling accounts for a big part of the electricity consumption [24].

### 39      **14.1.1 Inherent Limitations of Large-scale Data Centers**

40      Although, building large-scale DCs enables to cope with the actual demand, it is far  
41      from delivering sustainable and efficient UC infrastructures. In addition to requiring  
42      the construction and the deployment of a complete network infrastructure to reach  
43      each DC, it exacerbates the inherent limitations of the Cloud Computing model:

- 44      • The externalization of private applications/data often faces legal issues that re-  
45      strain companies from outsourcing them on external infrastructures, especially  
46      when located in other countries.

- 47 • The overhead implied by the unavoidable use of the Internet to reach distant  
48 platforms is wasteful and costly in several situations: Deploying a broadcasting  
49 service of local events or an online service to order pizzas at the edge of the polar  
50 circle, for instance, leads to important overheads since most of the users are *a*  
51 *priori* located in the neighborhood of the event/the pizzeria.
- 52 • The connectivity to the application/data cannot be ensured by centralized dedicated  
53 centers, especially if they are located in a similar geographical zone. The  
54 only way to ensure disaster recovery is to leverage distinct sites [23].

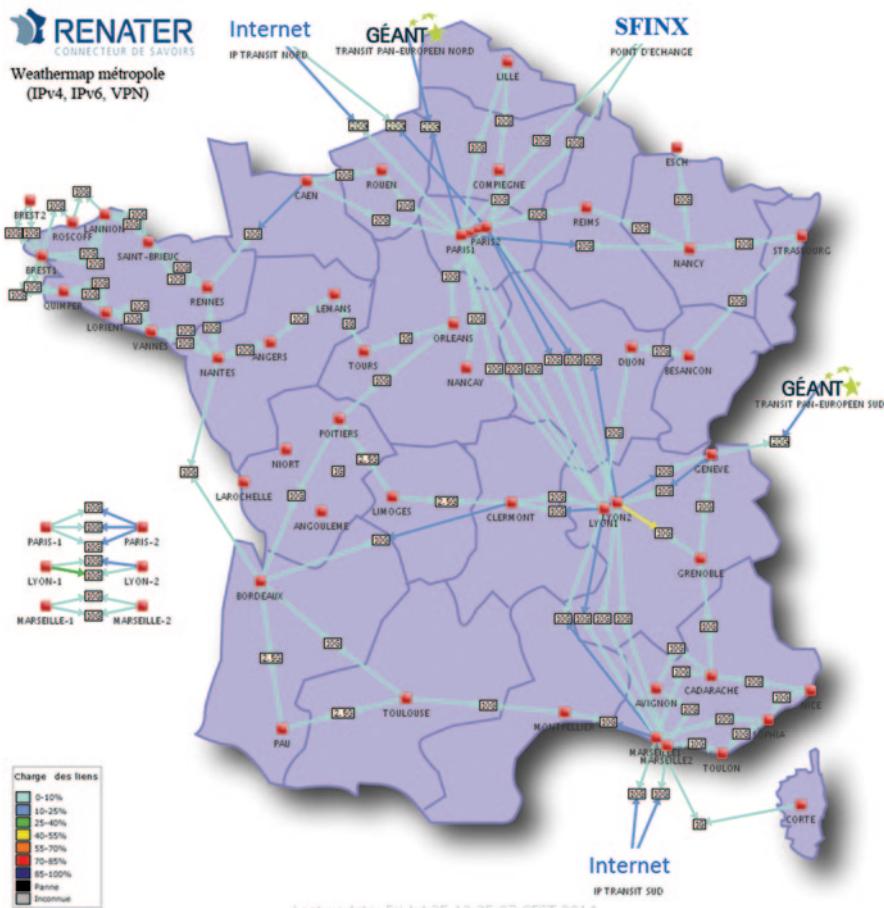
55 The two first points could be partially tackled by hybrid or federated Cloud solutions  
56 [4], that aim at extending the resources available on one Cloud with those of  
57 another one; however, the third one requires a disruptive change in the way UC  
58 resources are managed.

59 Another issue is that, according to some projections of a recent IEEE report [25],  
60 the network traffic has been doubling roughly every year. Consequently, bringing  
61 IT services closer to the end users is becoming crucial to limit the energy impact of  
62 these exchanges and to save the bandwidth of some links. Similarly, this notion of  
63 locality is critical for the adoption of the UC model by applications that need to deal  
64 with a large amount of data as getting them in and out using actual UC infrastruc-  
65 tures may significantly impact the global performance [18].

66 The concept of micro/nano DCs at the edge of the backbone [24] may be seen  
67 as a complementary solution to hybrid platforms in order to reduce the overhead of  
68 network exchanges. However, operating multiple small DCs breaks somehow the  
69 idea of mutualization in terms of physical resources and administration simplicity,  
70 making this approach questionable.

### 71 **14.1.2 Ubiquitous and Oversized Network Backbones**

72 One way to partially solve the mutualization concern enlightened by the defend-  
73 ers of large-scale DCs is to directly deploy the concept of micro/nano DCs upon  
74 the Internet backbone. People are (and will be) more and more surrounded by  
75 computing resources, especially those in charge of interconnecting the IT equip-  
76 ment. Even though these small- and medium-sized facilities include resources  
77 that are barely used [3, 8], they can hardly be removed (e.g., routers). Consider-  
78 ing this important aspect, we claim that a new generation of UC platforms can be  
79 delivered by leveraging existing network centers, starting from the core nodes of  
80 the backbone to the different network access points in charge of interconnecting  
81 public and private institutions. By such a mean, network and UC providers would  
82 be able to mutualize resources that are mandatory to operate network/data centers  
83 while delivering widely distributed UC platforms able to better match the geo-  
84 graphical dispersal of users. Figure 6.1 allows to better capture the advantages of  
85 such a proposal. It shows a snapshot of the network weather map of RENATER,  
86 the backbone dedicated to universities and research institutes in France. It reveals  
87 several important points:



**Fig. 6.1** The RENATER Weather Map on May 2013, the 27th, around 4 p.m. Each red square corresponds to a particular point of presence (PoP) of the network. The map is available in real-time at: <http://www.renater.fr/racourci>

- 88     • As mentioned before, most of the resources are underutilized (only two links are
- 89        used between 45 and 55%, a few between 25 and 40%, and the majority below
- 90        the threshold of 25%).
- 91     • The backbone was deployed and is renewed to match the demand: The density of
- 92        points of presence (PoPs, i.e., small- or medium-sized network centers), as well
- 93        as the bandwidth of each link, are more important on the edge of large cities such
- 94        as Paris, Lyon, or Marseille.
- 95     • The backbone was designed to avoid disconnections, since 95 % of the PoPs can
- 96        be reached by at least two distinct routes.

### 97 **14.1.3 Locality-Based Utility Computing**

98 This chapter aims at introducing locality-based UC (LUC) infrastructures, a new  
99 generation of UC platforms that solve inherent limitations of the Cloud Computing  
100 paradigm relying on large-scale DCs. Although, it involves radical changes in the  
101 way physical and virtual resources are managed, leveraging network centers is a  
102 promising way to deliver highly efficient and sustainable UC services.

103 From the physical point of view, network backbones provide appropriate infra-  
104 structures, i.e., reliable and efficient enough to operate UC resources spread across  
105 the different PoPs. Ideally, UC resources would be able to directly take advantage  
106 of computation cycles available on network active devices, i.e., those in charge of  
107 routing packets. However, leveraging network resources to make external computa-  
108 tions may lead to important security concerns. Hence, we propose to extend each  
109 PoP with a number of servers dedicated to hosting virtual machines (VMs). As it is  
110 natural to assume that the network traffic and UC demands are proportional, larger  
111 network centers will be completed with more UC resources than the smaller ones.  
112 Moreover, by deploying UC services on relevant PoPs, a LUC infrastructure will be  
113 able to natively confine network exchanges to a minimal scope, minimizing alto-  
114 gether the energy footprint of the network, the impact on latency and the congestion  
115 phenomena that may occur on critical paths (for instance Paris and Marseille on  
116 RENATER).

117 From the software point of view, the main challenge is to design a comprehen-  
118 sive distributed system in charge of turning a complex and diverse network of re-  
119 sources into a collection of abstracted computing facilities that are both reliable and  
120 easy to operate.

121 The design of the LUC Operating System (OS), an advanced system being able  
122 to unify many UC resources distributed on distinct sites, would enable Internet  
123 service providers (ISPs) and other institutions in charge of operating a network  
124 backbone to build an extreme scale LUC infrastructure with a limited additional  
125 cost. Instead of redeploying a complete installation, they will be able to leverage IT  
126 resources and specific devices such as computer room air conditioning units, invert-  
127 ers, or redundant power supplies already present in each center of their backbone.

128 In addition to considering *locality* as a primary concern, the novelty of the LUC  
129 OS proposal is to consider the VM as the basic object it manipulates. Unlike ex-  
130 isting research on distributed operating systems designed around the process con-  
131 cept, a LUC OS will manipulate VMs throughout a federation of widely distributed  
132 physical machines. Virtualization technologies abstract out hardware heterogeneity,  
133 and allow transparent deployment, preemption, and migration of virtual environ-  
134 ments (VEs), i.e., a set of interconnected VMs. By dramatically increasing the flex-  
135 ibility of resource management, virtualization allows to leverage state-of-the-art  
136 results from other distributed systems areas such as autonomous and decentralized  
137 systems. Our goal is to build a system that allows end users to launch VEs over a  
138 distributed infrastructure as simply as they launch processes on a local machine,  
139 i.e., without the burden of dealing with resources availability or location.

#### 140 14.1.4 Chapter Outline

141 Section 14.2 describes the key objectives of a LUC OS and the associated chal-  
142 lenges. Section 14.3 explains why our vision differs from current and previous UC  
143 solutions. In Section 14.4, we present how such a unified system may be designed  
144 by delivering the premises of the DISCOVERY (DIStributed and COoperative  
145 framework to manage Virtual EnviRonments autonomously) system, an agent-  
146 based system enabling the distributed and cooperative management of virtual envi-  
147 ronments over a large-scale distributed infrastructure. Future work and opportunities  
148 are addressed in Sect. 14.5. Finally, Sect. 14.6 concludes this chapter.

## 149 14.2 Overall Vision and Major Challenges

150 Similar to traditional operating systems (OSes), a LUC OS is composed of many  
151 mechanisms. Trying to identify all of them and establishing how they interact is an  
152 on-going work (see Sect. 14.4). However, we have pointed out the following key  
153 objectives to be considered when designing a LUC OS:

- 154 • *Scalability*: A LUC OS must be able to manage hundreds of thousands of virtual  
155 machines (VMs) running on thousands of geographically distributed computing  
156 resources. These resources are small- or medium-sized computing facilities and  
157 may become highly volatile according to the network disconnections.
- 158 • *Reactivity*: To deal with the dynamicity of the infrastructure, a LUC OS should  
159 swiftly handle events that require to perform particular operations, either  
160 on virtual or on physical resources. This has to be done with the objective of  
161 maximizing the system utilization while meeting the Quality of Service (QoS)  
162 expectations of VEs. Some examples of operations that should be performed as  
163 fast as possible include: (i) the reconfiguration of VEs over distributed resources,  
164 sometimes spread across wide area networks (WANs), or (ii) the migration of  
165 VMs, while preserving their active connections.
- 166 • *Resiliency*: In addition to the inherent dynamicity of the infrastructure, failures,  
167 and faults should be considered as the norm rather than the exception at such a  
168 scale. The goal is therefore to transparently leverage the underlying infrastruc-  
169 ture redundancy to: (i) allow the LUC OS to keep working despite node failures  
170 and network disconnections (LUC OS robustness), and to (ii) provide snapshot-  
171 ting as well as high availability mechanisms for VEs (VM robustness).
- 172 • *Sustainability*: Although the LUC approach would reduce the energy footprint of  
173 UC services by minimizing the cost of the network, it is important to go one step  
174 further by considering energy aspects at each level of a LUC OS and propose  
175 advanced mechanisms in charge of making an optimal usage of each source of  
176 energy. To achieve such an objective, the LUC OS should take account of data  
177 related to the energy consumption of the VEs and the computing resources, as  
178 well as the environmental conditions (computer room air conditioning unit, loca-  
179 tion of the site, etc.).

- 180 • *Security and Privacy:* Similar to resiliency, the security and privacy issues affect  
181 the LUC OS itself and the VEs running on it. Regarding the LUC OS, the goals  
182 are: (i) to create trust relationships between different locations, (ii) to secure the  
183 peer-to-peer layers, (iii) to include security and privacy decisions and enforce-  
184 ment points in the LUC OS, and (iv) to make them collaborate through the se-  
185 cured peer-to-peer layers to provide end-to-end security and privacy. Regarding  
186 the VEs, users should be able to express their requirements in terms of security  
187 and privacy; the LUC OS would then enforce these requirements.

188 In addition to the aforementioned objectives, working on a virtual infrastructure re-  
189 quires to deal with the management of VM images. Managing VM images in a dis-  
190 tributed way across a wide area network (WAN) is a real challenge that will require  
191 adapting state-of-the-art techniques related to replication and deduplication. Also,  
192 the LUC OS must take into account VM images location, for instance: (i) to allocate  
193 the right resources to a VE, or (ii) to prefetch VM images, to improve deployment  
194 performance or VM relocations.

195 Finally, one last scientific and technical challenge is the lack of a global view  
196 of the infrastructure. Maintaining a global view would indeed limit the scalability  
197 of the LUC OS, which is inconsistent with our objective to manage large-scale  
198 geographically distributed systems. Therefore, we claim that the LUC OS should  
199 rely on decentralized and autonomous mechanisms, which can match and adapt to  
200 the volatile topology of the infrastructure. Several decentralized mechanisms are  
201 already used in production on large-scale systems; for instance, Amazon relies on  
202 the Dynamo service [15] to create distributed indexes and recover from data in-con-  
203 sistencies and Facebook uses Cassandra [29], a massive scale structured store that  
204 leverages peer-to-peer techniques. In a LUC OS, decentralized and self-organizing  
205 overlays will enable to maintain the information about the current state of both vir-  
206 tual and physical resources, their characteristics, and availabilities. Such informa-  
207 tion is mandatory to build higher level mechanisms ensuring the correct execution  
208 of VEs throughout the whole infrastructure.

### 209 14.3 Background

210 Several generations of UC infrastructures have been proposed and still coexist  
211 [19]. However, neither Desktop, Grid, nor Cloud Computing platforms provide a  
212 satisfying UC model. Contrary to the current trend that promotes large offshore-  
213 centralized DCs as the UC platform of choice, we claim that the only way to achieve  
214 sustainable and highly efficient UC services is to target a new infrastructure that  
215 better matches the Internet structure. As it aims at gathering an unprecedented  
216 amount of widely distributed computing resources into a single platform providing  
217 UC services close to the end users, a LUC infrastructure is fundamentally different  
218 from existing ones. Keeping in mind the aforementioned objectives, recycling UC  
219 resource management solutions developed in the past is doomed to failure.

220 As previously mentioned, our vision significantly differs from hybrid Cloud  
221 Computing solutions. Although these research activities address important concerns  
222 related to the use of federated Cloud platforms, such as interface standardization  
223 for supporting cooperation and resource sharing, their propositions are incremental  
224 improvements of existing UC models. Recent investigations on hybrid Clouds and  
225 Cloud federation are comparable in some ways to previous works done on Grids,  
226 since the purpose of a Grid middleware is to interact with each resource manage-  
227 ment system composing the Grid [11, 49, 55].

228 By taking into account the network issues, in addition to traditional computing and  
229 storage concerns in Cloud Computing systems, the European SAIL project [50] is  
230 probably the one which targets the biggest advances with regard to previous works on  
231 Grid systems. More concretely, this project investigates new network technologies to  
232 provide end users of hybrid/federated Clouds with the possibility to configure and vir-  
233 tually operate the network backbone interconnecting the different sites they use [37].

234 More recently, the *Fog Computing* concept has been proposed as a promising  
235 solution to applications and services that cannot be put into the Cloud due to local-  
236 ity issues (mainly the latency and mobility concerns) [10]. Although it might look  
237 similar to our vision as they propose to extend the Cloud Computing paradigm to the  
238 edge of the network, *Fog Computing* does not target a unified system but rather pro-  
239 poses to add a third party layer (i.e., the *Fog*) between Cloud vendors and end users.

240 In our vision, UC resources (i.e., Cloud Computing ones) should be repacked in  
241 the different points of presence of backbones and operated through a unified system,  
242 the LUC OS. As far as we know, the only system that investigated whether a widely  
243 distributed infrastructure can be operated by a single system was the XtreemOS  
244 Project [36]. Although this project shared some of the goals of the LUC OS, it did  
245 not investigate how the geographical distribution of resources can be leveraged to  
246 deliver more efficient and sustainable UC infrastructures.

247 To sum up, we argue for the design and the implementation of a kind of dis-  
248 tributed OS, manipulating VEs instead of processes, and considering locality as  
249 a primary concern. Referred to as a LUC OS, such a system will include most  
250 of the mechanisms that are common to current UC management systems [17, 32,  
251 35, 39–41]. However, each of them will have to be rethought in order to leverage  
252 peer-to-peer algorithms. While largely unexplored for building operating systems,  
253 peer-to-peer/decentralized mechanisms have the potential to achieve the scalability  
254 required to manage LUC infrastructures. Using this technology for establishing the  
255 base mechanisms of a massive-scale LUC OS will be a major breakthrough from  
256 current static, centralized, or hierarchical management solutions.

## 257 14.4 Premise of a LUC OS: The DISCOVERY Proposal

258 In this section, we propose to go one step further by discussing preliminary in-  
259 vestigations around the design and implementation of a first LUC OS proposal:  
260 the DISCOVERY system. We draw the premises of the DISCOVERY system by  
261 emphasizing some of the challenges as well as some research directions to solve

262 them. Finally, we give some details regarding the prototype that is under development  
263 and how we are going to evaluate it.

#### 264 **14.4.1 Overview**

265 The DISCOVERY system relies on a multi-agent peer-to-peer system deployed on  
266 each physical resource composing the LUC infrastructure. Agents are autonomous  
267 entities that collaborate with one another to efficiently use the LUC resources. In  
268 our context, efficiency means that a good trade-off is found between users' expecta-  
269 tions, reliability, reactivity, and availability, while limiting the energy consumption  
270 of the system and providing scalability.

271 In DISCOVERY, each agent has two purposes: (i) maintaining a knowledge base  
272 on the composition of the LUC platform, and (ii) ensuring the correct execution of  
273 VEs. This includes the configuration, deployment, and monitoring of VEs as well  
274 as the dynamic allocation or relocation of VMs to adapt to changes in VEs require-  
275 ments and physical resources availability. To this end, agents will rely on dedicated  
276 mechanisms related to:

- 277 • The localization and monitoring of physical resources
- 278 • The management of VEs
- 279 • The management of VM images
- 280 • Reliability
- 281 • Security and privacy

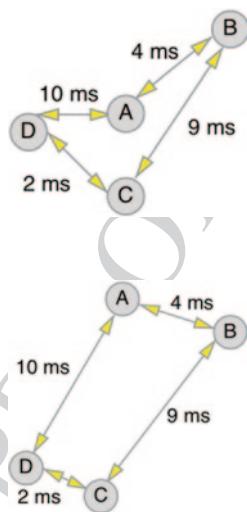
#### 282 **14.4.2 Resource Localization and Monitoring Mechanisms**

283 Keeping in mind that DISCOVERY should be designed in a fully decentralized fash-  
284 ion, its mechanisms should be built on top of an overlay network able to abstract out  
285 changes that occur at the physical level. The specific requirements of this platform  
286 will lead to the development of a novel kind of overlay networks based on locality  
287 and a minimalistic design. More concretely, the first step is to design, at the lowest  
288 level, an overlay layer intended to hide the details of the physical routes and com-  
289 puting utilities, while satisfying some basic requirements such as locality and avail-  
290 ability. This overlay needs to enable the communications between any two nodes in  
291 the platform. While overlay computing has been extensively studied over the last de-  
292 cade, we emphasize here on minimalism, especially on one key feature to implement  
293 a LUC OS: retrieving nodes that are geographically close to a given departure node.

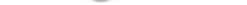
##### 294 **14.4.2.1 Giving Nodes a Position**

295 The initial configuration of the physical network can take an arbitrary shape.  
296 We choose to rely on the Vivaldi protocol [14]. *Vivaldi* is a distributed algorithm

**Fig. 6.2** Vivaldi plot before updating positions. Each node pings other nodes. Each node maintains a map of distance



**Fig. 6.3** Vivaldi plot after updating positions. The computed positions of other nodes have been updated



297 as-signing coordinates in the plane to nodes of a distributed system. Each node is  
 298 equipped with a *view* of the network, i.e., a set of nodes it knows. This view is ini-  
 299 tially assumed as random. Coordinates obtained by a node reflect its *position* in the  
 300 network, i.e., close nodes in the network are given close coordinates in the plane. To  
 301 achieve this, each node periodically checks the round trip time between itself and  
 302 another node (randomly chosen among nodes in its view) and adapts its distance (by  
 303 changing its coordinates) with this node in the plane accordingly. Refer to Figs. 6.2  
 304 and 6.3 for an illustration of four nodes (A, B, C, and D) moving according to the  
 305 Vivaldi protocol. A globally accurate positioning of nodes can be obtained if nodes  
 306 have a few long-distance nodes in their view [14]. These long-distance links can be  
 307 easily maintained by means of a simple gossip protocol.

#### 308 14.4.2.2 Searching for Close Nodes

309 Once the map is achieved (each node knows its coordinates), we are able to decide  
 310 whether two nodes are *close* by calculating their distance. However, the view of  
 311 each node does not a priori contain its closest nodes. Therefore, we need additional  
 312 mechanisms to locate a set of nodes that are close to a given initial node—Vivaldi  
 313 gives a *location* to each node, but not to the neighborhood. To achieve this, we use  
 314 a modified distributed version of the classic Dijkstra's algorithm used to find the  
 315 shortest path between two nodes in a graph. The goal is to build a *spiral* intercon-  
 316 necting the nodes in the plane that are the closest ones to a given initial node. Note  
 317 that the term *spiral* is here a misuse of language, since the graph actually drawn in  
 318 the plane might contain crossing edges. The only guarantee is that when following  
 319 the path constructed, the nodes are always further from the initial node.

320 Let us consider that our initial point is a node called *I*. The first step is to find  
 321 a node to build a two-node spiral with *I*. Such a node is sought in the view of *I* by

322 selecting the node, say  $S$ , having the smallest distance with  $I$ .  $I$  then sends its view  
323 to  $S$ ,  $I$  stores  $S$  as its successor in the spiral, and  $S$  adds  $I$  as its predecessor in the  
324 spiral. Then  $I$  forwards its view to  $S$ .  $S$  creates a new view by keeping the  $n$  nodes  
325 which are the closest to  $I$  in the views of  $I$  and  $S$ . This last view is then referred to  
326 as the *spiral view* and is intended to contain a set of nodes among which to find the  
327 next step of the spiral. Then  $S$  restarts the same process: Among the spiral view, it  
328 chooses the node with the smallest distance to  $I$ , say  $S'$ , and adds it in the spiral— $S$   
329 becomes the predecessor of  $S'$  and  $S'$  becomes the successor of  $S$ . Then, the spiral  
330 view is sent to  $S'$  which updates it with the nodes it has in its own view. The process  
331 is repeated until we consider that enough nodes have been gathered (a parameter  
332 sent by the application).

333 Note that one risk is to be blocked by having a spiral view containing only nodes  
334 that are already in the spiral, leading to the impossibility to build the spiral further.  
335 However, this problem can be easily addressed by forcing the presence of a few  
336 long-distance nodes whenever it is updated.

#### 337 14.4.2.3 Learning

338 Applying the protocol described above, the quality of the spiral is questionable in  
339 the sense that the nodes that are actually close to the starting node  $s$  may not be in-  
340 cluded. The only property ensured is that one step forward on the built path always  
341 takes us further from the initial node.

342 To improve the *quality* of the spiral, i.e., reduce the average distance between the  
343 nodes it comprises and the initial node, we add a learning mechanism coming with  
344 no extra communication cost: when a node is contacted to become the next node  
345 in one spiral, and receives the associated spiral view, it can also keep the nodes  
346 that are the closest to itself, thus potentially increasing the quality of a future spiral  
347 construction.

#### 348 14.4.2.4 Routing

349 In the context of a LUC infrastructure, one crucial feature is to be able to locate an  
350 existing VM. Having the same strategy consisting in improving the performance of  
351 the overlay based on the activity of the application, we envision a routing mecha-  
352 nism which will be improved by past routing requests. By means of the spiral mech-  
353 anism, a node is able to contact its neighboring nodes to start routing a message.

354 This initial routing mechanism can be very expensive as the number of hops  
355 can be linear in the size of the network. However, from previous communications,  
356 a node is able to memorize long links to different locations of the network. Conse-  
357 quently, from each routing request, the source of the request and each node on the  
358 path to the destination are able to learn long links, which will significantly reduce  
359 the number of hops of future requests. We are currently studying the amount of  
360 requests needed to get close to a logarithmic routing complexity. More generally,

361 we are working on the estimation if the activity of the application is required to: (i)  
362 guarantee the constant efficiency of the overlay, and (ii) converge, starting from a  
363 random configuration, to a fully efficient overlay network.

### 364 ***14.4.3 VEs Management Mechanisms***

365 In the DISCOVERY system, we define a VE as a set of VMs that may have specific  
366 requirements in terms of hardware, software, and also in terms of placement: For in-  
367 stance, some VMs must be on the same node/site to cope with performance objectives  
368 while others should not be collocated to ensure high-availability criteria [26]. As op-  
369 erations on a VE may occur in any place from any location, each agent should provide  
370 the capability to configure and start a VE, to suspend/resume/stop it, to relocate some  
371 of its VMs if the need arises, or simply to retrieve the location of a particular VE. Most  
372 of these mechanisms are provided by current UC platforms. However, as mentioned  
373 before, they should be revisited to leverage peer-to-peer mechanisms to correctly run  
374 on the infrastructure we target (i.e., in terms of scalability, resiliency, and reliability).

375 As a first example, placing the VMs of a VE requires the ability to find the avail-  
376 able nodes that fulfill the VM's needs (in terms of resource requirements as well as  
377 placement constraints). Such a placement can start locally, close to the client appli-  
378 cation requesting it, i.e., in its local group. If no such node is found, a simple navi-  
379 gation ensures that the request will encounter a bridge, leading to the exploration  
380 of further nodes. This navigation goes on until an adequate node is found. A similar  
381 process is performed by the mechanism in charge of dynamically controlling and  
382 adapting the placement of VEs during their lifetime. For instance, to ensure the  
383 particular needs of a VM, it can be necessary to relocate other VMs. According to  
384 the predefined constraints of VEs, some VMs might be relocated on far nodes while  
385 others would prefer to be suspended. Such a mechanism has been deeply studied in  
386 the DVMS framework [16, 46]. DVMS (Distributed Virtual Machine Scheduler) is  
387 able to dynamically schedule a significant number of VMs throughout a large-scale  
388 distributed infrastructure while guaranteeing VM resource expectations.

389 A second example regards the networking configuration of VEs. Although it might  
390 look simple, assigning the right IP to each VM as well as maintaining the intracon-  
391nectivity of a VE becomes a bit more complex than in the case of a single network  
392 domain, i.e., a single site deployment. Keeping in mind that a LUC infrastructure  
393 is, by definition, spread WANwide, a VE can be hosted between distinct network  
394 domains during its lifetime. No solution has been chosen yet. Our first investiga-  
395 tions led us to leverage techniques such as the IP over P2P project [20]. However,  
396 software-defined networking becomes more and more important; investigating pro-  
397 posals such as the Open vSwitch project [44] looks promising to solve such an issue.

### 398 ***14.4.4 VM Images Management***

399 In a LUC infrastructure, VM images could be deployed in any place from any other  
400 location. However, being in a decentralized, large-scale, heterogeneous, and widely

401 spread environment makes the management of VM images more difficult than with  
402 conventional centralized repositories. At coarse grain: (i) the management of the  
403 VM images should be consistent with regard to the location of each VM in the  
404 DISCOVERY infrastructure, and (ii) each VM image should remain reachable or at  
405 least recoverable in case of failures. The envisioned mechanisms to manage VM im-  
406 ages have been classified into two categories. First, some mechanisms are required  
407 to efficiently upload VM images and replicate them across many nodes, to ensure  
408 efficiency as well as reliability. Second, other mechanisms are needed to schedule  
409 VM image transfers. Advanced policies are important to improve the efficiency of  
410 each transfer that may occur either during the first deployment of a VM or during  
411 its relocations.

412 Regarding the storage and replication mechanisms, an analysis of an IBM Cloud  
413 concludes that a fully distributed approach using peer-to-peer technology is not the best  
414 choice to manage VM images, since the number of instances of the same VM image is  
415 rather small [42]. However, central or hierarchical solutions are not suited for the in-  
416 frastructure we target. Consequently, an improved peer-to-peer solution working with  
417 replicas and deduplication has to be investigated to provide more reliability, speed, and  
418 scalability to the system. For example, analyzing different VM images shows that at  
419 least 30% of the image is shared between different VMs [28]. This 30% can become a  
420 30% reduction in space, or a 30% increase in reliability, or in transfer speed. Depend-  
421 ing on the situation, we should decide to go from one scenario to another.

422 Regarding the scheduling mechanisms, a study showed that VM boot time could  
423 be increased from 10 to 240 s when multiple VMs running I/O intensive tasks use  
424 the same storage system [53]. Some actions can provide a performance boost and  
425 limit the overhead that is still observed in commercial Clouds [33], like providing  
426 the image chunks needed to boot first [54], defining a new image format, and paus-  
427 ing the rest of the I/O operations.

428 More generally, the amount of data linked with VM images is significant. Ac-  
429 tions involving data should be aware of consequences on metrics like (but not  
430 limited to): energy efficiency, reliability, proximity, bandwidth, and hardware us-  
431 age. The scheduler could also anticipate actions, for instance, moving images when  
432 the load is low or the energy is cheap.

#### 433 **14.4.5 Reliability Mechanisms**

434 Although, we can expect that the frequency of failures on LUC resources should be  
435 similar to that in current UC platforms, it is noteworthy to mention that the expected  
436 mean time to repair failed equipment might be much higher since resources will be  
437 highly distributed. For these reasons, specific mechanisms should be designed to  
438 manage failures transparently with a minimum downtime.

439 Ensuring the high availability of the DISCOVERY system requires the ability to  
440 autonomously relocate and restart any service on a healthy node in case of failure.  
441 Moreover, a Cassandra-like framework [29] is required to avoid losing or corrupt-  
442 ing information belonging to stateful services, since it provides a reliable and highly  
443 available back-end.

444 Regarding the VEs reliability, leveraging periodical VM snapshotting capabilities can provide a first level of fault tolerance. In case of failure, a VE can be  
445 restarted from its latest snapshot. Performing VM snapshotting in a large-scale,  
446 heterogeneous, and widely spread environment, is a challenging task. However, we  
447 believe that adapting ideas that were recently proposed in this field [38] would allow  
448 us to provide such a feature.  
449

450 Snapshotting is not enough for services that should be made highly available, but  
451 a promising solution is to use VM replication [43]. To implement VM replication in  
452 a WAN, solutions to optimize synchronizations between replicas [22, 47] should be  
453 investigated. Also, we think that a LUC infrastructure has a major advantage over  
454 other UC platforms, since it is tightly coupled with the network infrastructure. As  
455 such, we can expect *low* latencies between nodes which would enable us to provide  
456 a strong consistency between replicas while achieving acceptable response time for  
457 the replicated services.

458 Reliability techniques will of course make use of the overlays for resource lo-  
459 calization and monitoring. Replicated VMs should be hosted on nodes that have a  
460 low probability to fail simultaneously. Following the previously defined overlay  
461 structure, this can be done through a navigation scheme where at least one bridge  
462 is encountered. A replica can then be monitored by a *watcher*, which is in the same  
463 local group as the replica.

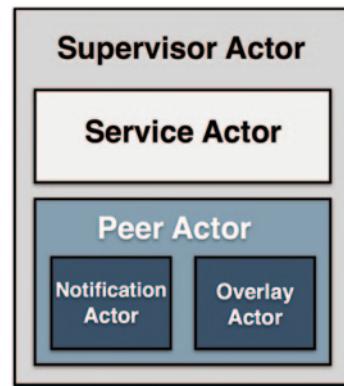
#### 464 **14.4.6 Security and Privacy Mechanisms**

465 To be successful, DISCOVERY needs to provide mechanisms and methods to con-  
466 struct trust relationships between resource providers. Trust relationships are known  
467 to be complex to build [34]. Providing strong authentication, assurance, and certifi-  
468 cation mechanisms to providers and users is required, but is definitely not enough.  
469 Trust covers socioeconomic aspects that must be addressed but are out of the scope  
470 of this chapter. The challenge is to provide a trusted DISCOVERY base.

471 As overlays are fundamentals to all DISCOVERY mechanisms, another chal-  
472 lenge is to ensure that they are not compromised. Recent advances [12] might en-  
473 able to tackle such concerns.

474 The third challenge will consist in: (i) providing end users with a way to define  
475 their own security and privacy policies, and (ii) ensuring that these policies are  
476 enforced. The expression of these policies itself is a complex task, as it requires  
477 to improve the current trade-off between security (and privacy) and usability. To  
478 ease the expression of these policies, we are currently designing a domain-specific  
479 language to define high-level security and privacy requirements [9, 30]. These poli-  
480 cies will be enforced in a decentralized manner, by distributed security and privacy  
481 decision and enforcement points (SPDEPs) during the lifetime of the VEs. Imple-  
482 menting such SPDEP mechanisms in a distributed fashion will require conducting  
483 specific research, as currently there are only prospective proposals for classic UC  
484 infrastructures [5, 51]. Therefore, we need to investigate whether such proposals  
485 can be adapted to the LUC infrastructure by leveraging appropriate overlays.

**Fig. 6.4** The *Peer Actor* Model. The *Supervisor Actor* monitors all the actors it encapsulates while the *Peer Actor* acts as an interface between the services and the overlay



#### 486 14.4.7 Towards a First Proof of Concept

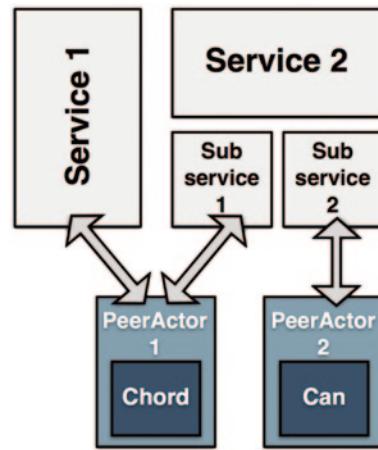
487 The first prototype is under heavy development. It aims at delivering a simple  
 488 mock-up for integration/collaboration purposes. Following the coarse-grained ar-  
 489 chitecture described in the previous sections, we have started to identify all the  
 490 components participating in the system, their relationships, as well as the resulting  
 491 interfaces. Conducting such a work now is mandatory to move towards a more  
 492 complete as well as more complex system.

493 To ensure a scalable and reliable design, we chose to rely on the use of high-level  
 494 programming abstractions. More precisely, we are using distributed complex event  
 495 programming [27] in association with the actor model [1]. This enables us to easily  
 496 switch between a push- and a pull-oriented approach depending on our needs.

497 Our preliminary studies showed that a common building block is mandatory  
 498 to handle resiliency concerns in all components. Concretely, it corresponds to a  
 499 mechanism in charge of throwing notifications that are triggered by the low level  
 500 network overlay each time a node joins or leaves it. Such a mechanism makes the  
 501 design and the development of higher building blocks easier as they do not have to  
 502 provide specific portions of code to monitor infrastructure changes.

503 This building block has been designed around the *Peer Actor* concept (see  
 504 Figs. 6.4 and 6.5). The *Peer Actor* serves as an interface between higher services  
 505 and the communication layer. It provides methods that enable to define the behav-  
 506 iors of a service when a resource joins or leaves a particular peer-to-peer overlay  
 507 as well as when neighbors change. Considering that several overlays may coexist  
 508 in the DISCOVERY system, the association between a *Peer Actor* and its *Overlay*  
 509 *Actor* is done at runtime and can be changed on the fly if need be. However, it is  
 510 noteworthy that each *Peer Actor* takes part to one and only one overlay at the same  
 511 time. In addition to the *Overlay Actor*, a *Peer Actor* is composed of a *Notification*  
 512 *Actor* that processes events and notifies registered actors. As illustrated in Fig. 6.5,  
 513 a service can use more than one *Peer Actor* (and reciprocally). Mutualizing a *Peer*  
 514 *Actor* enables for instance to reduce the network overhead implied by the main-  
 515 tenance of the overlays. In the example, the first service relies on a *Peer Actor*

**Fig. 6.5** A Peer Actor instantiation. The first service relies on a *Peer Actor* implementing a Chord overlay while the second service uses an additional *Peer Actor* implementing a CAN structure



516 implementing a Chord overlay [52], while the second service uses an additional  
 517 *Peer Actor* implementing a CAN structure [48].

518 By such a mean, higher-level services can take the advantage of the advanced  
 519 communication layers without dealing with the burden of managing the different  
 520 overlays. As an example, when a node disappears, all services that have been reg-  
 521 istered as dependent on such an event are notified. Service actors can thus react  
 522 accordingly to the behavior that has been specified.

523 Regarding the design and the implementation of the DISCOVERY system, each  
 524 service is executed inside its own actor and communicates by exchanging messages  
 525 with the other ones. This ensures that each service is isolated from the others: When  
 526 a service crashes and needs to be restarted, the execution of other services is not af-  
 527 fected. As previously mentioned, we consider that at the LUC infrastructure scale,  
 528 failures are the norm rather than the exception; hence, we decided that a *Supervisor*  
 529 *Actor* would monitor each actor (see Fig. 6.4). DISCOVERY services are under  
 530 the supervision of the DISCOVERY agent: This design allows to precisely define  
 531 a strategy that will be executed in case of service failures. This will be the way to  
 532 introduce self-healing and self-organizing properties to the DISCOVERY system.

533 This building block has been fully implemented by leveraging the Akka/Scala  
 534 framework [2], and is available online at <https://github.com/BeyondTheClouds>.

535 As a proof of concept (POC), we are implementing a first high-level service in  
 536 charge of dynamically scheduling VMs across a LUC infrastructure by leveraging  
 537 the DVMS [46] proposal (see Sect. 14.4.3). The low-level overlay that is being cur-  
 538 rently implemented is a robust ring based on the Chord algorithm combined with  
 539 the Vivaldi positioning system: It enables services to select nodes that have low  
 540 latency, so that collaboration will be more efficient.

541 To validate the behavior, the performance as well as the reliability of our proof of  
 542 concept, we are performing several experiments on the Grid'5000 test bed [6] that  
 543 comprises hundreds of nodes distributed on 10 computing sites that are geographi-  
 544 cally spread across France. To make experiments with DISCOVERY easier, we  
 545 developed a set of scripts that can deploy thousands of VMs throughout the whole

46 infrastructure in a *one-click* fashion [64]. By deploying our POC on each node and by  
547 leveraging the VM deployment scripts, we can evaluate real scenario by injecting  
548 specific workloads in the different VMs. The validation of this first POC is almost  
549 completed. The resulting system will be the first to provide reactive, reliable, and  
550 scalable reconfiguration mechanisms of virtual machines in a fully distributed and  
551 autonomous way. This new result will pave the way for a complete proposal of the  
552 DISCOVERY system.

## 553 14.5 Future Work/Opportunities

### 554 14.5.1 *Geo-diversification as a Key Element*

555 The Cloud Computing paradigm is changing the way applications are designed. In  
556 order to benefit from elasticity capabilities of Cloud systems, applications integrate  
557 or leverage mechanisms to provision resources, i.e., starting or stopping VMs, ac-  
558 cording to their fluctuating needs. The ConPaaS system [45] is one of the promising  
559 systems for elastic Cloud applications. At the same time, a few projects have started  
560 investigating distributed/collaborative ways of hosting famous applications such as  
561 Wikipedia or Facebook-like systems by leveraging volunteer computing techniques.  
562 However, considering that resources provided by end users were not reliable enough,  
563 only few contributions have been done yet. By providing a system that will enable to  
564 operate widely spread but more reliable resources closer to the end users, the LUC  
565 OS proposal may strongly benefit to this research area. Investigating the benefit of  
566 locality provisioning (i.e., combining elasticity and distributed/collaborative host-  
567 ing) is a promising direction for all Web services that are embarrassingly distributed  
568 [13]. Image sharing systems, such as Google Picasa or Flickr, are examples of ap-  
569 plications where leveraging locality will enable to limit network exchanges: Users  
570 could upload their images on a peer that is close to them, and images would be trans-  
571 ferred to other locations only when required (pulling versus pushing model).

572 LUC infrastructures will allow envisioning a wider range of services that may  
573 answer specific SMEs requests such as data archiving or backup solutions, while sig-  
574 nificantly reducing the network overhead as well as legal concerns. Moreover, it will  
575 make the deployment of UC services easier by relieving developers of the burden of  
576 dealing with multi-Cloud vendors. Of course, this will require software engineering  
577 and middleware advances to easily take advantage of locality. But proposing LUC OS  
578 solutions, such as the DISCOVERY project, is the mandatory step before investigat-  
579 ing new APIs enabling applications to directly interact with the LUC OS internals.

### 580 14.5.2 *Energy, a Primary Concern for Modern Societies*

581 The energy footprint of current UC infrastructures, and more generally of the  
582 Internet, is a major concern for the society. Although we need to conduct deeper

583 investigations, we clearly expect that by its design and the way to operate it, a  
584 LUC infrastructure will have a smaller impact with a better integration in the whole  
585 Internet ecosystem.

586 Moreover, the LUC proposal is an interesting way to deploy the data furnaces pro-  
587 posal [31]. Concretely, following the Smart City recommendations (i.e., delivering  
588 efficient and sustainable ICT services), the construction of new districts in metrop-  
589 olises may take advantage of each LUC/Network PoP in order to heat buildings  
590 while operating UC resources remotely by means of a LUC OS. Finally, taking into  
591 account recent results about passive data centers, such as solar-powered micro-data  
592 centers, might extend this idea. The idea behind passive computing facilities is to  
593 limit as much as possible the energy footprint of major hubs and DSLAMS by  
594 taking advantage of renewable energies to power them, and by using the heat they  
595 produce as a source of energy. Combining such ideas with the LUC approach would  
596 allow reaching an unprecedented level of energy efficiency for UC platforms.

## 597 14.6 Conclusion

598 Cloud Computing has entered into our daily life with a great speed. From classic  
599 high performance computing simulations to the management of huge amounts of  
600 data coming from mobile devices and sensors, its impact can no longer be dis-  
601 regarded. While a lot of progress has already been made in Cloud technologies,  
602 there are several concerns that limit the complete adoption of the Cloud Computing  
603 paradigm.

604 In this chapter, we have outlined that, in addition to these concerns, intrinsic  
605 issues limit the current model of UC. Instead of following the current trend by try-  
606 ing to cope with existing platforms and network interfaces, we proposed to take a  
607 different direction by promoting the design of a system that will be efficient and  
608 sustainable at the same time, putting knowledge and intelligence directly into the  
609 network backbone itself.

610 The innovative approach, we introduced, will definitely tackle and go beyond  
611 Cloud Computing limitations. Our objective is to pave the way for a new genera-  
612 tion of Utility Computing infrastructures that better match the Internet structure by  
613 means of advanced operating mechanisms. By offering the possibility to tightly  
614 couple UC servers and network backbones throughout distinct sites and operate  
615 them remotely, the LUC OS technology may lead to major changes in the design  
616 of UC infrastructures as well as in their environmental impact. The internal mecha-  
617 nisms of the LUC OS should be topology-dependent and resources-efficient. The  
618 natural distribution of the nodes through the different points of presence should  
619 be an advantage, which allows to process a request according to its scale: Local  
620 requests should be computed locally, while large computations should benefit from  
621 a large number of nodes.

622 Finally, we believe that LUC investigations may contribute to fill the gap between  
623 the distributed computing community and the networked ones. This connection

624 between these two communities has already started with the different activities  
625 around Software-Defined Networking and Network as a Service. In the long term,  
626 this may result in a new community dealing with UC challenges where network and  
627 computational concerns are fully integrated. Such a new community may leverage  
628 the background of both areas to propose new systems that are more suitable to ac-  
629 commodate the needs of our modern societies.

630 We are well aware that the design of a complete LUC OS and its adoption by  
631 companies and network providers require several big changes in the way UC infra-  
632 structures are managed and WANs are operated. However, we are convinced that  
633 such an approach will pave the way towards highly efficient and sustainable UC  
634 infrastructures, coping with heterogeneity, scale, and faults.

## 635 References

- 636 1. Agha G (1986) Actors: a model of concurrent computation in distributed systems. MIT Press,  
637 Cambridge
- 638 2. Akka (2013) Build powerful concurrent & distributed applications more easily. <http://www.akka.io>. Accessed: March 2013
- 639 3. Andrew O (2003) Data networks are lightly utilized, and will stay that way. review of net-  
640 work economics. Rev Netw Econ 2(3):210–237
- 641 4. Armbrust M, Fox A, Griffith R, Joseph AD, Katz R, Konwinski A, Lee G, Patterson D, Rabkin  
642 A, Stoica I, Zaharia M (2010) A view of cloud computing. Commun ACM 53(4):50–58
- 643 5. Bacon J, Evans D, Evers DM, Migliavacca M, Pietzuch P, Shand B (2010) Enforcing end-  
644 to-end application security in the cloud (big ideas paper). In: Proceedings of the ACM/IFIP/  
645 USENIX 11th International Conference on Middleware, Springer-Verlag, Berlin, Middle-  
646 ware'10, pp 293–312
- 647 6. Balouek D, Carpen Amarie A, Charrier G, Desprez F, Jeannot E, Jeanvoine E, Lébre A,  
648 Margery D, Niclausse N, Nussbaum L et al (2013) Adding virtualization capabilities to the  
649 Grid'5000 testbed. In: Ivanov I, Sinderen M, Leymann F, Shan T (eds) Cloud computing and  
650 services science, Springer, Berlin
- 651 7. Balouek D, Lebre A, Quesnel F (2013) Flauncher and DVMS: deploying and scheduling  
652 thousands of virtual machines on hundreds of nodes distributed geographically. In: The sixth  
653 IEEE international scalable computing challenge (collocated with CCGRID), Delft, The  
654 Netherlands
- 655 8. Benson T, Akella A, Maltz DA (2010) Network traffic characteristics of data centers in the  
656 wild. In: Proceedings of the 10th ACM SIGCOMM Conference on Internet Measurement,  
657 ACM, New York, IMC'10, pp 267–280
- 658 9. Blanc M, Briffaut J, Clevy L, Gros D, Rouzaud-Cornabas J, Toinard C, Venelle B (2013)  
659 Mandatory protection within clouds. In: Nepal S, Pathan M (eds) Security, privacy and trust  
660 in cloud systems, Springer, Berlin
- 661 10. Bonomi F, Milito R, Zhu J, Addepalli S (2012) Fog computing and its role in the internet of  
662 things. In: Proceedings of the first edition of the MCC workshop on mobile cloud computing,  
663 ACM, New York, USA, MCC'12, pp 13–16
- 664 11. Buyya R, Ranjan R, Calheiros RN (2010) InterCloud: utility-oriented federation of cloud  
665 computing environments for scaling of application services. In: Proceedings of the 10th inter-  
666 national conference on algorithms and architectures for parallel processing, Springer-Verlag,  
667 Berlin, ICA3PP'10, pp 13–31
- 668 12. Castro M, Druschel P, Ganesh A, Rowstron A, Wallach DS (2002) Secure routing for struc-  
669 tured peer-to-peer overlay networks. SIGOPS Oper Syst Rev 36(SI):299–314
- 670

- 671 13. Church K, Greenberg A, Hamilton J (2008) On delivering embarrassingly distributed cloud  
672 services. In: HotNets
- 673 14. Dabek F, Cox R, Kaashoek MF, Morris R (2004) Vivaldi: a decentralized network coordinate  
674 system. In: Proceedings of the 2004 conference on applications, technologies, architectures,  
675 and protocols for computer communications, SIGCOMM'04, pp 15–26
- 676 15. DeCandia G, Hastorun D, Jampani M, Kakulapati G, Lakshman A, Pilchin A, Sivasubramanian  
677 S, Vosshall P, Vogels W (2007) Dynamo: Amazon's highly available key-value store. In: Pro-  
678 ceedings of twenty-first ACM SIGOPS symposium on operating systems principles, ACM,  
679 SOSP'07, pp 205–220
- 680 16. Discovery (2013) Distributed VM scheduler. <http://beyondtheclouds.github.io/DVMS/>.  
681 Accessed: March 2013
- 682 17. CloudStack (2013) CloudStack, open source cloud computing. <http://cloudstack.apache.org>.  
683 Accessed: March 2013
- 684 18. Foster I (2011) Globus online: accelerating and democratizing science through cloud-based  
685 services. IEEE Internet Comput 15(3):70–73
- 686 19. Foster I, Kesselman C (2011) The history of the grid. In: Foster I, Gentzsch W, Grandinetti L,  
687 Joubert GR (eds) Advances in parallel computing—volume 20: HPC: from grids and clouds  
688 to exascale. IOS Press, Amsterdam
- 689 20. Ganguly A, Agrawal A, Boykin PO, Figueiredo R (2006) IP over P2P: enabling self-configur-  
690 ing virtual IP networks for grid computing. In: Proceedings of the 20th international confer-  
691 ence on Parallel and distributed processing, IEEE Computer Society, Washington, DC, USA,  
692 IPDPS'06
- 693 21. Gary Cook JVH (2013) How dirty is your data? Greenpeace International report
- 694 22. Gerofi B, Ishikawa Y (2012) Enhancing TCP throughput of highly available virtual machines  
695 via speculative communication. In: Proceedings of the 8th ACM SIGPLAN/SIGOPS confer-  
696 ence on Virtual Execution Environments, NY, USA, VEE'12
- 697 23. Gigaom Consortium (2012) Amazon outages—lessons learned. <http://gigaom.com/cloud/amazon-outages-lessons-learned/>. Accessed: 23 Feb. 2014
- 698 24. Greenberg A, Hamilton J, Maltz DA, Patel P (2008) The cost of a cloud: research problems  
699 in data center networks. SIGCOMM Comput Commun Rev 39(1):68–73
- 700 25. Group IEW (2012) IEEE 802.3TM industry connections ethernet bandwidth assessment
- 701 26. Hermenier F, Lawall J, Muller G (2013) BtrPlace: a flexible consolidation manager for highly  
702 available applications. IEEE Transactions on Dependable and Secure Computing
- 703 27. Janiesch C, Matzner M, Muller O (2011) A blueprint for event-driven business activity man-  
704 agement. In: Proceedings of the 9th international conference on business process manage-  
705 ment, Springer-Verlag, BPM'11, pp 17–28
- 706 28. Jin K, Miller EL (2009) The effectiveness of deduplication on virtual machine disk images.  
707 In: Proceedings of SYSTOR 2009: the Israeli Experimental Systems Conference, ACM, New  
708 York, USA, SYSTOR'09, pp 7:1–7:12
- 709 29. Lakshman A, Malik P (2010) Cassandra: a decentralized structured storage system. SIGOPS  
710 Oper Syst Rev 44(2):35–40
- 711 30. Lefray A, Caron E, Rouzaud-Cornabas J, Zhang HY, Bousquet A, Briffaut J, Toinard C  
712 (2013) Security-aware models for clouds. In: Poster Session of IEEE Symposium on High  
713 Performance Distributed Computing (HPDC)
- 714 31. Liu J, Goraczko M, James S, Belady C, Lu J, Whitehouse K (2011) The data furnace: heating  
715 up with cloud computing. In: Proceedings of the 3rd USENIX conference on hot topics in  
716 cloud computing, HotCloud'11
- 717 32. Lowe S (2011) Mastering VMware vSphere. Wiley: Indianapolis
- 718 33. Mao M, Humphrey M (2012) A performance study on the VM startup time in the cloud. In:  
719 Proceedings of the 2012 IEEE Fifth International Conference on Cloud Computing, IEEE  
720 Computer Society, CLOUD'12, pp 423–430
- 721 34. Miller KW, Voas J, Laplante P (2010) In trust we trust. Computer 43:85–87
- 722 35. Moreno-Vozmediano R, Montero R, Llorente I (2012) IaaS cloud architecture: from virtual-  
723 ized datacenters to federated cloud infrastructures. Computer 45(12):65–72

- 725 36. Morin C (2007) XtreemOS: a grid operating system making your computer ready for participating  
726 in virtual organizations. In: Proceedings of the 10th IEEE International Symposium  
727 on Object and Component-Oriented Real-Time Distributed Computing, IEEE Computer So-  
728 ciety, ISORC'07, pp 393–402
- 729 37. Murray P, Sefidcon A, Steinert R, Fusenig V, Carapinha J (2012) Cloud networking: an infra-  
730 structure service architecture for the wide area. HP Labs Tech Report-HPL-2012-111R1
- 731 38. Nicolae B, Bresnahan J, Keahey K, Antoniu G (2011) Going back and forth: efficient mul-  
732 tiple deployment and multisnapshotting on clouds. In: Proceedings of the 20th international  
733 symposium on High performance distributed computing, ACM, New York, USA, HPDC'11,  
734 pp 147–158
- 735 39. Nimbus (2013) Nimbus is cloud computing for science. <http://www.nimbusproject.org>. Accessed: March 2013
- 736 40. OpenNebula (2013) Open source data center virtualization. <http://www.opennebula.org>. Accessed: March 2013
- 737 41. OpenStack (2013) The open source, open standards cloud. <http://www.openstack.org>. Accessed:  
738 March 2013
- 739 42. Peng C, Kim M, Zhang Z, Lei H (2012) VDN: virtual machine image distribution network  
740 for cloud data centers. In: INFOCOM, 2012, pp 181–189
- 741 43. Petrovic D, Schiper A (2012) Implementing virtual machine replication: a case study using  
742 Xen and Kvm. In: Proceedings of the 2012 IEEE 26th International Conference on Advanced  
743 Information Networking and Applications, pp 73–80
- 744 44. Pfaff B, Pettit J, Koponen T, Amidon K, Casado M, Shenker S (2009) Extending networking  
745 into the virtualization layer. In: ACM HotNets
- 746 45. Pierre G, Stratan C (2012) ConPaaS: a platform for hosting elastic cloud applications. IEEE  
747 Internet Comput 16(5):88–92
- 748 46. Quesnel F, Lebre A, Sudholt M (2012) Cooperative and reactive scheduling in large-scale  
749 virtualized platforms with DVMS. Concurr Comput Pract Exp 25(12):1643–1655
- 750 47. Rajagopalan S, Cully B, O'Connor R, Warfield A (2012) SecondSite: disaster tolerance as a  
751 service. In: Proceedings of the 8th ACM SIGPLAN/SIGOPS conference on Virtual Execu-  
752 tion Environments, ACM, VEE'12, pp 97–108
- 753 48. Ratnasamy S, Francis P, Handley M, Karp R, Shenker S (2001) A scalable content-address-  
754 able network. In: SIGCOMM'01: Proceedings of the conference on applications, technolo-  
755 gies, architectures, and protocols for computer communications, ACM, New York, USA,  
756 SIGCOMM'01, pp 161–172
- 757 49. Rochwerger B, Breitgand D, Levy E, Galis A, Nagin K, Llorente IM, Montero R, Wolfsthal  
758 Y, Elmroth E, Caceres J, Ben-Yehuda M, Emmerich W, Galan F (2009) The reservoir model  
759 and architecture for open federated cloud computing. IBM J Res Dev 53(4):4:1–4:11
- 760 50. Sail Consortium (2012) Scalable and adaptive internet solutions—European Project FP7 pro-  
761 gram. <http://www.sail-project.eu>. Accessed: March 2014
- 762 51. Sandhu R, Boppana R, Krishnan R, Reich J, Wolff T, Zachry J (2010) Towards a discipline  
763 of mission-aware cloud computing. In: Proceedings of the 2010 ACM workshop on cloud  
764 computing security workshop, pp 13–18
- 765 52. Stoica I, Morris R, Karger D, Kaashoek MF, Balakrishnan H (2001) Chord: a scalable peer-  
766 to-peer lookup service for internet applications. In: Proceedings of the 2001 conference  
767 on Applications, technologies, architectures, and protocols for computer communications,  
768 ACM, New York, USA, SIGCOMM'01, pp 149–160
- 769 53. Tan T, Simmonds R, Arlt B, Arlitt M, Walker B (2008) Image management in a virtualized  
770 data center. SIGMETRICS Perform Eval Rev 36(2):4–9
- 771 54. Tang C (2011) FVD: a high-performance virtual machine image format for cloud. In:  
772 Proceedings of the 2011 USENIX conference on USENIX annual technical conference,  
773 USENIX Association, USENIXATC'11
- 774 55. Zhao H, Yu Z, Tiwari S, Mao X, Lee K, Wolinsky D, Li X, Figueiredo R (2012) CloudBay:  
775 enabling an online resource market place for open clouds. In: Proceedings of the 2012 IEEE/  
776 ACM Fifth International Conference on Utility and Cloud Computing, UCC'12

## Chapter 14: Author Query

---

AQ1. The author "Marin Bertier" has been marked as the corresponding author. Please confirm.



AQ2. "[7]" is not cited in the text. Please provide the citation or delete the entry from the reference list.

