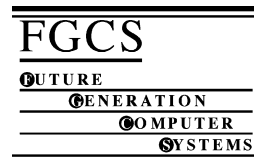




ELSEVIER

Available online at www.sciencedirect.com

Future Generation Computer Systems ■■■■■ ■■■■■ ■■■■■

www.elsevier.com/locate/fgcs

Powerful resource discovery for Arigatoni overlay network[☆]

Raphael Chand^a, Michel Cosnard^b, Luigi Liquori^{b,*}

^a University of Geneva, Switzerland

^b Inria Sophia Antipolis, France

Received 26 January 2007; received in revised form 16 February 2007; accepted 27 February 2007

Abstract

Arigatoni is a structured multi-layer overlay network providing various services with variable guarantees, and promoting an intermittent participation in the overlay since peers can appear, disappear and organize themselves dynamically. Arigatoni provides fully decentralized, asynchronous and scalable resource discovery; it also provides mechanisms for dealing with an overlay with a dynamic topology. This paper introduces a nontrivial improvement of the resource discovery protocol by allowing the registration and request of *multiple instances of the same service*, *service conjunctions*, and *multiple services*. Adding multiple instances is a nontrivial task since the discovery protocol must keep track (when routing requests) of peers that accept to serve and peers that deny the service. Adding service conjunctions allows a single peer to offer different services *at the same time*. Simulations show that it is efficient and scalable.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Overlay networks; Resource discovery; Virtual organizations; Dynamic graphs; Peer-to-peer; Global computing; Grid computing

1. Introduction

The explosive growth of the Internet gives rise to the possibility of designing large *overlay networks* and *virtual organizations* consisting of Internet-connected *global computers*, able to provide a rich functionality of services that makes use of aggregated computational power, storage, information resources, etc. Arigatoni [1] is a structured multi-layer overlay network which provides resource discovery with variable guarantees in a virtual organization where peers can appear, disappear and organize themselves dynamically. In a nutshell, the main units in Arigatoni are:

- A *Global Computer Unit*, GC, i.e. the basic peer of the global computing paradigm; it is typically a small device, like a PDA, a laptop or a PC, connected through IP in various ways (wired, wireless, etc.).

- A *Global Broker Unit*, GB, i.e. the basic unit devoted to subscribe and unsubscribe GCs, to receive service queries from client GCs, to contact potential server GCs, to negotiate with them services, to authenticate clients and servers, and to send all the information necessary to allow the client GC and the servers GCs to communicate. Every GB controls a *colony* of collaborating global computers. Hence, communication intra-colony is initiated via only one GB, while communication inter-colonies is initiated through a chain of GB-2-GB message exchanges whose security is guaranteed via PKI mechanisms. In both cases, when a client GC receives an acknowledgment of a service request from the direct leader GB, then the GC is served directly by the server(s) GC, i.e. without a further mediation of the GB, in a pure peer-to-peer fashion. Registrations and requests are performed via a simple query language *à la* SQL and a simple *orchestration language à la* LINDA, or BPEL.
- A *Global Router Unit*, GR i.e. the basic unit close to GCs and GBs that is devoted to send and receive packets, using the resource discovery protocol [2,3], and to forward the “payload” to the units which are connected with this router. The connection GB-GR-GC is ensured via a suitable API.
- A *Colony* is a simple virtual organization composed of exactly one leader GB and a set (possibly empty) of

[☆] This work is supported by the AEOLUS FET Global Computing Proactive IST-015964, *Algorithmic Principles for Building Efficient Overlay Computers*.

* Corresponding author. Tel.: +33 4 92 38 71 93; fax: +33 4 92 38 79 71.

E-mail addresses: Raphael.Chand@cu.unige.ch (R. Chand), Michel.Cosnard@inria.fr (M. Cosnard), Luigi.Liquori@inria.fr, Luigi.Liquori@gmail.com (L. Liquori).

individuals. Individuals are global computers (think it as an *Amoeba*) or subcolonies (think it as a *Protozoa*). The two main characteristics of a colony are:

- (1) A colony has *exactly* one leader GB and at least one individual (the GB itself);
- (2) A colony contains individuals (GCs, or other subcolonies).

The main challenges in Arigatoni lie in the management of an overlay network with a dynamic topology, the routing of queries, and the discovery of resources in the overlay. In particular, resource discovery is a nontrivial problem for large distributed systems featuring a discontinuous amount of resources offered by global computers and an intermittent participation in the overlay. Thus, Arigatoni features two protocols: The *virtual intermittent protocols*, VIP, and the *resource discovery protocol* RDP. The VIP protocol deals with the *dynamic topology* of the overlay, by allowing individuals to login/logout to/from a colony. This implies that the routing process may lead to failures, because some individuals have logged out, or are temporarily unavailable, or because they have been *manu militari* logged out by the broker because of their poor performance or greediness [4].

The total decoupling between GCs in *space* (GCs do not know each other), *time* (GCs do not participate in the interaction at the same time), and *synchronization* (GCs can issue service requests and do something else, or may be doing something else when being asked for services) is a major feature of Arigatoni overlay network. Another important property is the encapsulation of resources in colonies. All those properties play a major role in the scalability of Arigatoni's RDP.

The version V1 of the RDP protocol [2] enabled one service at the time to be requested, e.g. a CPU or a specific file. In [3], the protocol was enhanced (V2) to take into account *multiple instances of the same service*. Adding multiple instances is a nontrivial task because the broker must keep track (when routing requests) of how many resource instances were found in its own colony before delegating the rest of the instances to the surrounding colonies.

The version V3, presented in this paper, adds *multiple services* and *service conjunctions*. Adding service conjunctions allows a global computer to offer several services *at the same time*. Multiple services requests can be also asked to a GB; each service is processed sequentially and independently of others. As an example of multiple instances, a GC may ask for three CPUs, *or* four chunks of 1GB of RAM, *or* one chunk of 10 GB of HD, *or* one gcc compiler; as an example of a service conjunction, a GC may ask for another GC offering *at the same time* one CPUs, *and* one chunk of 1GB of RAM, *and* one chunk of 10 GB of HD *and* one gcc compiler.

If a request succeeds, then via the orchestration language of Arigatoni (not described in this paper), the GC client can synchronize all resources offered by the servers GCs. To sum up, the contributions of this paper are:

- A complete description of the resource discovery protocol RDP V3, which allows multiple instances, multiple services and service conjunctions.

- A new version of the simulator taking into account the nontrivial improvements in the resource discovery protocol.
- Simulation results that show that our enhanced protocol is scalable.

The rest of the paper is structured as follows: after Section 2 describing the main machinery underneath the protocol features, Section 3 introduces the pseudocode of the protocol; then Section 4 shows our simulation results and finally Section 5 provides related work analysis and concluding remarks. This paper is an extended and improved version of [3].

2. Resource discovery protocol RDP V3

Suppose a GC X registers to its GB and declares its availability to offer a service S, while another GC Y issues a request for a service S'. Then, the GB looks in its *routing table* and *filters* S' against S. If there exists a solution to this filter equation, X can provide a resource to Y. For example, $S \triangleq [\text{CPU} = \text{Intel}, \text{Time} < 10 \text{ sec}]$ filters against $S' \triangleq [\text{CPU} = \text{Intel}, \text{Time} > 5 \text{ sec}]$, with attribute values Intel and Time between 5 and 10 seconds. In RDP V2, a global computer asks not only for a service S, but also for a certain number of instances of S; this is denoted by SREQ : [(S, n)]. In RDP V3:

- Every GC registers in the colony with a *tuple* of (*services, instances*) like SREG : $[(S_i, n_i)]^{i=1\dots h}$, and may ask for a tuple like SREQ : $[(S_j, n_j)]^{j=1\dots k}$. Each service is processed sequentially and independently of others. This is achieved by wrapping the RDP V2 code inside a **for each** $j = 1 \dots k$ **do** ... V2code ... **end foreach**.
- A service request may also have the shape SREQ : $[(\bigwedge_{i=1\dots n} S_i), n]$, i.e. the system is no longer asked to find n occurrences of a single service, but rather n occurrences of a conjunction of services. That is, the system has to look for n distinct GCs, each GC being able to provide all the services in $\bigwedge_{i=1\dots n} S_i$.

Each GB maintains a *routing table* \mathcal{T} representing the *services* that are registered in its colony. The table is updated according to the *dynamic registration and unregistration* of GC in the overlay. For a given S, the table has the form $\mathcal{T}[S] = [(P_j, m_j)]^{j=1\dots k}$, where $(P_j)^{j=1\dots k}$ are the address of the *direct children in the GB's colony*, and $(m_j)^{j=1\dots k}$ are the instances of S available at P_j . For a single atomic service request SREQ : [(S, n)], the steps are:

- Look for q *distinct* GCs able to provide S in the local GB's colony;
- If $q < n$, then search $r \leq (n - q)$ remaining instances in local subcolonies;
- If $r < (n - q)$, then delegate $(n - q - r)$ remaining instances to the leader of the colony.

A GC receiving a service request chooses the services that it *accepts/rejects* to serve; then, it generates a SRESP message containing the lists of accepted/rejected services, and sends it

to its GB. The response messages are then propagated back in the overlay, following the reverse path.

A service request SREQ : [(S, n)] may arrive bottom-up to the GB directly from its colony, or top-down from its own leader. In both cases, the GB tries to locate n distinct GC that can provide S. More precisely, the list $[(P_j, m_j)]^{j=1\dots k}$ contains all the direct children in GB's colony that can provide S (child P_j with m_j instances of S).

The discovery protocol features two search modes, *selective* and *exhaustive*. Let SREQ : [(S, n)], and $\mathcal{T}[S] = [(P_j, m_j)]^{j=1\dots k}$.

- The selective search mode is resource conservative at the price of important delays in case of low acceptance rates. The selective mode consist in:
 - If $\sum_{i=1}^k m_i \geq n$, then there are enough resources in the GB's colony to provide S. Let $y \leq k$ be the smallest index such that $\sum_{i=1}^y m_i \geq n$, and $\sum_{i=1}^{y-1} m_i < n$. Then, SREQ : [(S, m_i)] is sent to all $P_i (i \leq y - 1)$, and SREQ : [(S, $n - \sum_{i=1}^{y-1} m_i$)] is sent to P_y .
 - If $\sum_{i=1}^k m_i < n$, then there are not enough GCs in the GB's colony to provide S. Then, SREQ : [(S, m_i)] is sent to all $P_i (i \leq k)$, and SREQ : [(S, $n - \sum_{i=1}^k m_i$)] is delegated to the GB's leader. The rationale is that one first try to ask for *as many resources* in GB's colony, and then ask GB's leader for the *remaining resources*.
- The exhaustive mode is resource eager, but is independent of the acceptance rate. The exhaustive search mode consists in sending SREQ : [(S, $\min(m_i, n)$)] to all $P_i (1 \leq i \leq k)$, and to delegate SREQ : [(S, $n - \sum_{i=1}^k \min(m_i, n)$)] to the GB's leader. The rationale is to first ask for *all* resources in the GB's colony, and then ask the GB's leader for the remaining resources.

A Service Response SRESP : ACC : [(S, a)], or SRESP : REJ : [(S, d)], may follow service requests for services S. That is, “a” GCs accepted to provide S, and “d” denied. Due to the asynchrony of Arigatoni, more replies can arrive to the colony's leader (i.e. $a + d \geq n$). As for requests, there exists two modes that tell the way the acceptances are propagated back to the leader of the colony. In the *selective reply* mode, at most the number of instances of S that were asked by the leader are returned, whereas in the *exhaustive reply* mode, *all* acceptances are returned.

As for acceptances, there exist two modalities that determine the way those acceptances are propagated back to the colony's leader.

- In the *selective search* mode, the *whole colony* is asked for n instances of S, at most. This implies that exactly d instances of S must now be looked for to fulfil the original request. Hence, one first try to find d instances of S in other subcolonies. One then delegate the instances that could not be found to the colony's leader. Finally, the remaining instances are reported back as rejected.
- In the *exhaustive search* mode, each *sub-colony* is asked for n instances of S, at most. Hence, there may be other subcolonies that have not replied yet, and which may reply

with enough acceptances to fulfil the request. The remaining instances must be delegated to the colony's leader.

3. RDP pseudo-code

In this section, we detail the pseudo-code of the RDP V3. Five global variables are used for each Arigatoni's interaction “ask-route-reply-route-back”: *Path*, *asked*, *downstream*, *upstream*, and *SendList*. Each message (SREQ or SRESP) contains a unique identifier *id*, which is initially set to the address of the GC that sends the initial SREQ message. The variable *Path* is a simple hash “keyed” by the identifier of the message. The other variables are double hashes which first key is the identifier of the message, and second key is a given service S. The intuitive meaning of those variables is listed below.

- *Path*{*id*}: Peer address: identifies the peer from which the original SREQ message came from.
- *asked*{*id*}{S}: Integer: instances of S asked and not replied, i.e. the remaining number of instances of S to find to fulfil the request.
- *downstream*{*id*}{S}: Integer: instances of S asked in colony and not replied.
- *upstream*{*id*}{S}: Integer: instances of S delegated but not replied.
- *SendList*{*id*}{S}: (Peer address, Integer)*: the list of direct children that are potentially able to provide S.

The pseudo-code of RDP V3 is showed in Algorithms [1–8].

Algorithm 1 Receiving SREQ_{id}:[(S_i, n_i)]^{i=1...k} from P_{from} (executed by P)

```

1: Path{id} ← Pfrom // To trace back the reverse route
2: for each (S, n) ∈ SREQ do
3:   if SendList{id}{S} = ∅ then
4:     SendList{id}{S} ← Filter(S, Pfrom) // Filter S in P's routing table
5:   end if
6:   (RoutingList, remaining) ← Route(Pfrom, S, n, search_mode) // Build a routing // list
7:   asked{id}{S} ← asked{id}{S} + n
8:   if remaining ≠ 0 then // Remaining instances to find
9:     if L ≠ ∅ and L ≠ Pfrom then // L exists and is different from Pfrom
10:      Insert L:(S, remaining) in RoutingList
11:      upstream{id}{S} ← upstream{id}{S} + remaining
12:     else // P's colony is isolated
13:      Send SRESPid:REJ:[(S, remaining)] to Pfrom
14:      asked{id}{S} ← asked{id}{S} - remaining
15:     end if
16:   end if
17: end for
18: for each Q:(S, m) ∈ RoutingList do
19:   Send SREQid:[(S, m)] to Q // Send SREQid to every element in RoutingList
20: end for

```

Case of service request (Algorithm 1). Consider a global broker P receiving a service request SREQ_{id} from a neighbour P_{from}, and let L be P's leader. The same steps are performed for each tuple (S, n) ∈ SREQ.

- In line 1, the originator of the request is first recorded in *Path*{*id*}, so as to allow reply messages to follow the reverse path.
- In line 4, the *Filter* function (Algorithm 6) determines the *SendList*{*id*}{S} corresponding to service S, i.e. the list of P's direct children that are potentially able to provide S.

- In line 6 the *Route* function (Algorithm 8) builds (*RoutingList*, *remaining*), i.e., the list of children that will receive a particular service request, according to the selected search mode, and the positive number of the remaining instances for which no server has been found. The *RoutingList* contains a list of mappings of the form $Q : [(S, m)]$ which means that we send a service request $SREQ : [(S, m)]$ to a neighbour Q .
- In line 9, if L exists and is not the originator of the request (to avoid routing loops), then the entry $L : (S, remaining)$ is appended to *RoutingList* (line 10), and the *upstream* counter is incremented, accordingly (line 11); else (line 12, L exists and it is the originator of the request), since servers can be found for *remaining* instances of service S , a rejection reply is sent back to the originator of the request (line 13), and the *asked* counter is decremented, accordingly (line 14).
- In line 19, a service request is sent to each neighbour Q having an entry in the *RoutingList*.

Algorithm 2 Receiving $SRESP_{id}:ACC:[(S_i, a_i)]^{i=1\dots k}$ from P_{from} (exec. by P)

```

1: case search_mode is
   "selective" :
2:   Send  $SRESP_{id}:ACC:(S, a)$  to  $Path\{id\}$  // Forward the SRESP
3: "exhaustive" :
4:   for each  $(S, n) \in SRESP$  do
5:     if  $P_{from} = L$  then // Top-down request
6:        $upstream\{id\}\{S\} \leftarrow \max(upstream\{id\}\{S\} - a; 0)$ 
7:     else // Bottom-up request
8:        $downstream\{id\}\{S\} \leftarrow \max(downstream\{id\}\{S\} - a; 0)$ 
9:     end if
10:    if  $asked\{id\}\{S\} \geq a$  then // More instances asked than accepted
11:       $asked\{id\}\{S\} \leftarrow asked\{id\}\{S\} - a$ 
12:       $acc\_return \leftarrow a$ 
13:    else // More instances accepted than asked
14:       $acc\_return \leftarrow asked\{id\}\{S\} - a$ 
15:       $asked\{id\}\{S\} \leftarrow 0$ 
16:    end if
17:   case reply_mode is
18:     "selective" :
19:       Send  $SRESP_{id}:ACC:(S, a)$  to  $Path\{id\}$  // Accepted "a" instances
20:     "exhaustive" :
21:       Send  $SRESP_{id}:ACC:(S, acc\_return)$  to  $Path\{id\}$  // Accepted
22:       // "acc_return" instances
23:   end case
24: end for
25: end case

```

Case of service response (Algorithms 2, 3). Consider a global broker P receiving a reply message $SRESP_{id}$ from a neighbour P_{from} . The operation of the resource discovery algorithm is explained hereafter. The same steps are performed for each tuple in $SRESP$.

- *Acceptance* (Algorithm 2). For each $(S, a) \in SREQ$, let $SRESP_{id} : ACC : [(S, a)]$ arrive from P_{from} at P , i.e. " a " global computers in P 's colony accepted to provide S .

If the *selective search* mode is used to route the original service request $SREQ_{id} : (S, n)$, issued by $Path\{id\}$, then the *whole colony* is asked for at most n instances of S . Hence, no more than n acceptances may arrive from P 's colony. Thus, the reply message is simply forwarded back to $Path\{id\}$ (line 2).

If the *exhaustive search* mode is used, then *each child* is asked for at most n instances of S . Hence, it is possible that a number of acceptances higher than n arrives from P 's

Algorithm 3 Receiving $SRESP_{id}:REJ:[(S_i, d_i)]^{i=1\dots k}$ from P_{from} (exec. by P)

```

1: if  $P_{from} = L$  then // Return rejections
2:   Send  $SRESP_{id}:REJ:(S, d)$  to  $Path\{id\}$ 
3:    $asked\{id\}\{S\} \leftarrow asked\{id\}\{S\} - d$ 
4:    $upstream\{id\}\{S\} \leftarrow upstream\{id\}\{S\} - d$ 
5: else // Retry at other children or delegate
6:   case search_mode is
7:     "exhaustive" : // Try to delegate or reject
8:       for each  $(S, n) \in SRESP$  do
9:          $downstream\{id\}\{S\} \leftarrow \max(downstream\{id\}\{S\} - d; 0)$ 
10:        if  $asked\{id\}\{S\} \leq downstream\{id\}\{S\} + upstream\{id\}\{S\}$  then
11:          // Fewer instances asked than down/upstream'd
12:          Wait for more replies from other children
13:        else // More instances asked than down/upstream'd
14:           $remaining \leftarrow asked\{id\}\{S\} - downstream\{id\}\{S\} - upstream\{id\}\{S\}$ 
15:          if  $L \neq \emptyset$  and  $L \neq Path\{id\}$  then
16:             $upstream\{id\}\{S\} \leftarrow upstream\{id\}\{S\} + remaining$ 
17:            Send  $SREQ_{id}:(S, remaining)$  to  $L$ 
18:          else
19:             $asked\{id\}\{S\} \leftarrow asked\{id\}\{S\} - remaining$ 
20:            Send  $SRESP_{id}:REJ:(S, remaining)$  to  $Path\{id\}$ 
21:          end if
22:        end if
23:      end for
24:      Remove  $P_{from}$  from  $SendList\{id\}\{S\}$ 
25:    end for
26:    "selective" : // Try other children, delete or reject
27:      for each  $(S, n) \in SRESP$  do
28:        Remove  $P_{from}$  from  $SendList\{id\}\{S\}$  // Don't send requests to  $P_{from}$ 
29:        // anymore
30:         $(RoutingList, remaining) \leftarrow Route(P_{from}, S, d, search\_mode)$ 
31:        if  $remaining \neq 0$  then // Still remaining instances to treat
32:          if  $L \neq \emptyset$  and  $L \neq P_{from}$  then // L exists and is different from  $P_{from}$ 
33:            Insert  $L:(S, remaining)$  in  $RoutingList$ 
34:             $upstream\{id\}\{S\} \leftarrow upstream\{id\}\{S\} + remaining$ 
35:          else // P's colony is isolated
36:            Send  $SRESP_{id}:REJ:(S, remaining)$  to  $Path\{id\}$ 
37:             $asked\{id\}\{S\} \leftarrow asked\{id\}\{S\} - remaining$ 
38:          end if
39:        end if
40:      end for
41:    end case
42:  end if

```

colony. To do this, counters *asked*, *upstream*, *downstream* and *acc_return* are updated, accordingly (lines 6–15).

The *selective reply* mode simply replies back to $Path\{id\}$ with a acceptance instances (line 18), while the *exhaustive reply* mode replies with *acc_return* instances (line 20).

- *Rejections* (Algorithm 3). For each $(S, d) \in SREQ$, let $SRESP_{id} : REJ : [(S, d)]$ arrive from P_{from} at P , i.e., " d " global computers in P 's colony refused to provide S . This implies that *all* global computers in P 's colony have received a request for a service S .

If the sender of the message is the leader L , then no other potential servers for the d instances of S can be found. Consequently, the rejection message is simply forwarded back (line 2), and counters *asked* and *upstream* are updated, accordingly (lines 3 and 4).

If L is not the sender of the rejected message, then there may be other potential servers in the colony or in other surrounding colonies. The operation of the protocol depends on the search mode that is used.

– (*exhaustive search mode*) Then there are no other potential servers in P 's colony but there may be in other surrounding colonies. Hence, the number of instances of S that need to be found to fulfill the request is first determined.

If $asked \leq downstream + upstream$ (line 9), then there are enough potential servers in the colony or in surrounding colonies that have not replied yet, to fulfil the request. Consequently, we simply wait for more replies (line 11).

In contrast, if $asked \geq downstream + upstream$, then one looks for more potential servers in order to fulfil the request. Then, there are $(asked - downstream - upstream)$ of them to be found (line 13). As said before, servers may be found by delegating to the leader L. Hence, the latter receives a request for the remaining instances of S, if possible, (line 16), or a rejection is sent back to the original sender of the request (line 19). The *upstream* or *asked* counters are updated, accordingly (lines 15 and 18).

– (*selective search mode*) Then there may be other potential servers in P's colony. The process is the same as in Algorithm 1, except that one do not consider children that have already received a request (line 22,24). For that purpose, one use the *SendList* that is originally created by the *Filter* function (during the processing of the original service request message), and produce another *RoutingList* with the *Route* function (line 27).

Finally, one proceeds as in Algorithm 1 (lines 28–41).

Algorithm 4 Receiving SREQ:[(S_i, 1)] from L (executed by a GC)

```

1: for each  $i = 1 \dots k$  do
2:   if accept then
3:      $Acc \xleftarrow{append} S_i$ 
4:   end if
5: end for
6: Send SRESP:ACC:[(Si, 1)] $i \in Acc$  to L
7: Send SRESP:REJ:[(Si, 1)] $i \notin Acc$  to L

```

Algorithm 5 Receiving SRESP:ACC:[(S, a)] from L (executed by a GC)

```

1: Initiate P2P negotiation with GCs (embedded in message)

```

RDP embedded in GCUs (Algorithms 4, 5). We show the cases of receiving a service request and a positive service response. The case of negative service response is trivial since the GC do simply nothing. Note that each reply message is formally of the form SRESP : ACC : [(S, P_i)] ^{$i=1\dots k$} where the P_i are the GCs that accepted to provide S (the same for rejections). Those algorithms are quite intuitive and need not be commented.

Algorithm 6 The *Filter*(S, P_{from}) function for RDP V2

```

1: for each entry  $T[S'] = [(P_j, n_j)]^{j=1\dots k}$  in T do
2:   if S filters S' then
3:     for each  $j = 1 \dots k$  such that  $P_j \neq P_{from}$  do
4:        $SendList\{id\}\{S\}\{P_j\} \leftarrow SendList\{id\}\{S\}\{P_j\} + n_j$  // Add/update
//  $SendList\{id\}\{S\}\{P_j\}$ 
5:     end for
6:   end if
7: end for
8: return  $SendList\{id\}\{S\}$ 

```

The *filter* function for V2 builds the *SendList*{id}{S} corresponding to the request id for a service S, i.e. the direct list of P's children that are potentially able to serve

the request for S coming from P_{from}. The function parses all the services in the routing table, accordingly. The *Filter* function for V3 enables service conjunctions and for this it has to be modified. For a service request of the form SREQ : [($\bigwedge_{i=1\dots n} S_i$), n], the system is no longer asked to find n occurrences of a single service, but rather n occurrences of a conjunction of services. That is, the system has to look for n distinct GCs, each GC being able to provide all the services in $\bigwedge_{i=1\dots n} S_i$. A conjunction of services is treated atomically, i.e. as a single service S. Both algorithms are quite intuitive and they are described in Algorithms 6 and 7.

The *Route* function of Algorithm 8 builds *RoutingList*, i.e., the list of neighbors that ask for a particular service, according to the selected search mode; it has the form {(P_i : (S, n_i)) ^{$i=1\dots h$} }, that is neighbors P_i will receive a request for n_i instances of S. The function also returns the remaining instances for which no server has been found.

4. Protocol evaluation

The actual Arigatoni's topology is tree-based with a routing complexity of $O(\log N)$ (N being the number of nodes). However, in each GB, an extra complexity is required in order to solve the filter equation between the service request and the routing table T containing the mapping between peers and resources; this complexity is usually linear in the size of S.

To assess the effectiveness and the scalability of the protocol, we have conducted simulations using large numbers of units and service requests. For lack of space, we only present the results that correspond to the new features of the protocol, namely, the ability to specify multiple instances of a service, service conjunctions, and multiple services.

We have generated a network topology of 103 GBs, using the transit-stub model of the Georgia Tech Internetwork Topology Models package [5], on top of which we added the Arigatoni overlay network. We considered a finite set of services S₁ . . . S_r of size $r = 128$, with an exact filtering policy (i.e., S_i filters S_j and no other services), and we defined the *overlap interval* $1 \leq L \leq 128$, as the interval of indices inside which services filter each other, that is, for all $(i, j) \in L^2$, S_i filters against S_j. If $L = 128$, then all services filter each other;

Algorithm 7 The *Filter*(S \triangleq ($\bigwedge_{i=1\dots n} S_i$), P_{from}) function for RDP V3

```

1: for each  $i = 1 \dots n$  do
2:    $tmp \leftarrow 0$  // Auxiliary vector
3:   for each entry  $T[S'] = [(P_j, n_j)]^{j=1\dots k}$  in T do
4:     if Si filters S' then // Handle all conjunctions
5:       for each  $j = 1 \dots k$  such that  $P_j \neq P_{from}$  do
6:          $tmp[j] \leftarrow tmp[j] + n_j$ 
7:       end for
8:     end if
9:   end for
10:  if  $SendList\{id\}\{S\} = \emptyset$  then
11:     $SendList\{id\}\{S\} \leftarrow tmp$ 
12:  else
13:    for each  $j = 1 \dots k$  do
14:       $SendList\{id\}\{S\}\{P_j\} = \min(SendList\{id\}\{S\}\{P_j\}, tmp[j])$ 
15:    end for
16:  end if
17: end for
18: return  $SendList\{id\}\{S\}$ 

```


Algorithm 8 $Route(P_{\text{from}}, S, n, \text{search_mode})$

```

1: remaining ← n
2: RoutingList ← ∅
3: for each (Q, f) ∈ SendList{id}{S} do
4:   if Q = Pfrom or Q = Path{id} then
5:     continue // Go to next iteration in loop
6:   end if
7:   case search_mode is
8:     "exhaustive" :
9:       if n ≥ f then // More instances asked than offered
10:        Insert Q:(S, f) in RoutingList
11:        remaining ← remaining - f
12:        downstream{id}{S} ← downstream{id}{S} + f
13:        Remove (Q, f) from SendList{id}{S}
14:       else // More instances offered than asked
15:        Insert Q:(S, n) in RoutingList
16:        remaining ← 0
17:        downstream{id}{S} ← downstream{id}{S} + n
18:        f ← f - n
19:       end if
20:     "selective" :
21:       if remaining ≥ f then // More instances asked than offered
22:        Insert Q:(S, f) in RoutingList
23:        remaining ← remaining - f
24:        Remove (P, f) from SendList{id}{S}
25:       else // More instances to offer than asked
26:        Insert Q:(S, remaining) in RoutingList
27:        f ← f - remaining
28:        remaining ← 0
29:       end if
30:       if remaining = 0 then // No more instances to treat
31:         break // Break loop
32:       end if
33:     end case
34:   end for
35: return (RoutingList, remaining)

```

if $L = 1$, then each service only filters with itself. At each GB, we added a number of GCs chosen randomly between zero and 100.

At each GB, we added a random number of GCs chosen uniformly at random between zero and 100. To simulate subscription load, we then randomly registered at each GC each service with a probability ρ denoting the *global availability of services*, or as the density of population of GCs (since the more the number of GCs, the more likely it is that a given service is provided). The routing tables were updated, accordingly.

We then issued 50,000 service requests at GCs chosen uniformly at random. Each request contained either a certain number of instances l of a service, or one instance of a conjunction of services, also chosen uniformly at random. Each service request is then handled by the RDP V3. We used a service acceptance probability of $\alpha = 75\%$, which corresponds to the probability that a GC, receiving a request for a S (and offering S), accepts to provide it.

Upon completion of all the requests, we measured for each GB its load as the number of requests (messages) it received. We then computed the average load as the average value over the population of GBs in the system. We also computed the maximum load as the maximum value of the load over all the GBs in the system.

We computed the average and maximum load fractions as the average and maximum loads divided by the number of requests. The average load represents the average load of a GB due to the completion of the n requests. The average load fraction represents the fraction of requests that a GB served, on average. The maximum fraction represents the maximum fraction of the requests that a GB served. Since a GB receives

at most one request message corresponding to a given service request, the average load fraction can be seen as the fraction of GBs in the system involved in a service request, in average.

We computed the average service acceptance ratio as follows. For each GC, we computed the local acceptance ratio as the number of service requests that yielded a positive response (i.e. the system found at least one GC), over the number of service requests issued at that GC. A service request that contained multiple instances of a service counts as a positive response only if the system found as many GCs as the number of instances specified in the request.

We then computed the average acceptance ratio as the average value over the number of GC (that issued at least one service request). Fig. 1(a) shows the influence of the number of instances l in service requests on the average and maximum load fraction and on the average success rate. It is obtained with a value of ρ of 0.12%. Unsurprisingly, we observe that asking for more instances of a service requires more resources from the system. Indeed, for each instance, the system tries to find a different GC able to provide the service. We observe that low level GBs participate more, since there are more delegations. For values l higher than seven, the average and maximum load fractions stabilize, as the average success rate keeps decreasing; this means that there are not enough resources in the system to completely fulfil the request (i.e. not enough GCs able to provide the requested service).

Fig. 1(b) shows the influence of the number of services in a conjunction. It is obtained with a value of ρ of 3%. The phenomenon and its explanation is mostly similar to that of Fig. 1(a), except that it happens at a much greater scale. Indeed, the system must find a GC that can provide (and accepts) all the services in a conjunctive service request, which requires to probe a much greater portion of the network than if a single service is requested.

5. Related work and conclusions

Many technologies, algorithms and protocols have been proposed recently for resource discovery. Some of them focus on Grid or P2P oriented applications [6], but none targets the full generality as Arigatoni does. Indeed, Arigatoni deals with generic resource discovery for building an overlay network of global computers, structured in a virtual organization of variable topology, with clear distinct roles between leader GBs and individuals (GCs or subcolonies).

5.1. Discussion on closest overlay architectures (from [7])

The main challenges of “pervasive computing” are *how to build* an overlay network with dynamic topology, and *how to route queries* and *discover resources* efficiently.

In an overlay network, any message is routed through the full overlay; as such, the topology adopted in the overlay strongly affects routing algorithms and their complexity. The overlay is built on top of the physical one, and, thus, two neighbour nodes in the overlay network may be many links apart in the physical network. The Arigatoni topology is a dynamic

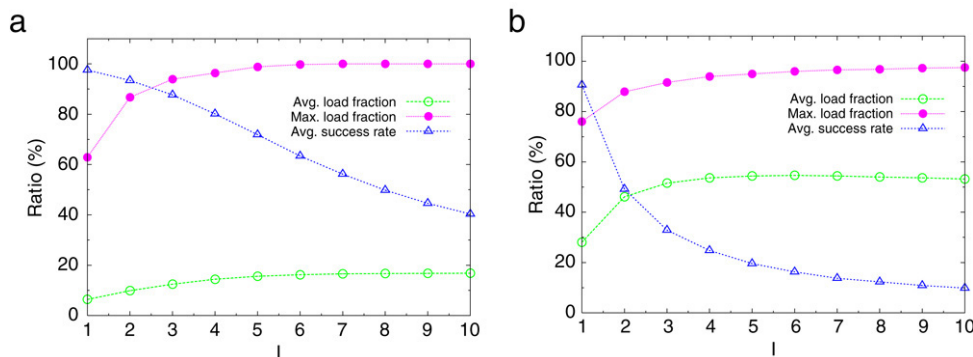


Fig. 1. Average and maximum load fraction, and average success rate w.r.t. (a) number of instances (b) number of services in conjunction.

hierarchical n-layer tree. To assist lookup, structured overlays map (key of) data item to nodes (our GBs). Hence, the mapping is usually done through hashing the key space of the data item to the id space of nodes. In Arigatoni, routing tables denoting the set of resources are stored in GBs; thus, each GB maintains a partition of the data space. When a GC asks for a resource, the query is *filtered* against the first direct GB's routing table; in case of *filter-failure*, the query is recursively forwarded to the direct super-GB. Any answer of the query must follow the reverse path. Thus, lookup overhead reduces when a query is satisfied in the current colony. Most structured overlays guarantee lookup operations that are logarithmic in the number of nodes. To improve performance of lookup, caching and replication of either data, search paths or both is possible. Besides improving routing, replication assists in providing load balancing, improves fault tolerance and the durability of data items.

In the literature, there are essentially the following types of overlays: structured (tree, ring or grid), unstructured, hybrid overlays (a combination of the two above), and multi-layer (or n-layer) overlays. Arigatoni falls in the latter category that is widely used in many P2P systems.

In a nutshell, in a *n-layer* overlay network, the responsibility assigned to individuals differs (think of the different roles between GBs and GCs), since super-peers (GBs) serving as a server for a subset of all peers. Ordinary peers (GCs) submit queries to their super-peers and receive results from it. Super-peers are also connected to each others; they route messages over the overlay network, submit, delegate and answer queries on behalf of their peers. This structure is replicated *recursively*, creating a *n-layer topology*, where peers become super-peers with decreasing responsibilities.

Typical issues in n-layer overlays are the size of each colony, and the internal coherence of the resources offered and requested by each colony. Typical bottlenecks of n-layers are reliability, service availability (related to few points of failure) and load balancing. Classical solutions to cope with these problems are adding redundancy at the broker-layer.

Historically, the most related tree topologies are BATON [8] and P-GRID [9], whereas the closest n-layer topologies are the one of Canon [10] and Coral [11]. We summarize the closest topologies.

- (BATON) is a balanced binary tree that features a left and a right routing table, both contained in each node (denoted by a single logical id). Nodes may join or leave the network at any time, provided the tree remains balanced. The node receiving a join can forward the join towards a node which has less children or which is a leaf node. This implies that a GC can become a GB. Leaving the network is constrained to not breaking the balanced tree unless finding a substitute. As such, load balancing can be costly.
- (P-GRID) is a distributed dynamic binary search tree, such that the search space is partitioned between peers. The salient feature of P-GRID is the separation of concerns between id and position in the network. All peers maintain a partial routing table of the search space, that *negotiated* with the closest peers. Multiple peers can be responsible for the same path, resulting in a non uniqueness of routing and a robustness under peer failure.
- (Canon) is a multi-layer overlay where routing is based on a hierarchical DHT. As in Arigatoni, the search space is partitioned into *domains*; in contrast, routing inside a domain is DHT-based, and topology is static.
- (Coral) is another hierarchical DHT. The search space is partitioned into three *clusters*, based on latency; a regional cluster, a continental cluster and a planet-wide cluster. It also comes with algorithm for self-organizing, merging and splitting clusters, to ensure acceptable diameters.

5.2. Conclusions

In this paper, we describe the version V3 of Arigatoni's generic resource discovery protocol. The new improved protocol RDP presented in this paper allows for *multiple instances*, *multiple services*, and *service conjunctions*. Other main achievements are the complete decoupling between the different units in the system, and the encapsulation of resources in local colonies, which enable Arigatoni to be potentially scalable to very large and heterogeneous populations.

The reliability of the RDP V3 itself, although desirable, is of lesser importance, given the fact that service provision is not guaranteed at all in Arigatoni (indeed it is not a requirement). In other words, when a GC issues a service request, it is possible that no individuals are found for some of the services included in the request. This happens, for example, if those services have

not been declared by any GCs in the system, or if all the GCs that have declared themselves as potential servers refuse those services.

However, at the cost of memory and bandwidth requirements, it is still possible (future work) to implement *reliable* resource discovery by using a reliable transmission protocol (e.g. TCP), an applicative *acknowledgment scheme* in combination with a retransmission buffer, and persistent data storage, and leader's replication.

As part of our ongoing research, we are also working on a more complete mathematical study of our system, based on more elaborate statistical models and realistic assumptions, as well as the possibility to include hierarchical DHT in addition to the routing tables. The possibility to change the Arigatoni topology from a hierarchical tree to a graph is also intriguing. We are currently working on the implementation of a actual prototype and the subsequent deployment on the PlanetLab experimental platform [12], and/or on GRID5000, the experimental platform available at the INRIA [13].

Acknowledgments

We warmly thank Pierre Lescanne and the anonymous referees for the useful comments and multiple constructive suggestions.

Work partly done while the first author was at INRIA Sophia Antipolis, France.

References

- [1] D. Benza, M. Cosnard, L. Liquori, M. Vesin, *Arigatoni*: A simple programmable overlay network, in: Proc. John Vincent Atanasoff International Symposium on Modern Computing, IEEE, 2006, pp. 82–91.
- [2] R. Chand, M. Cosnard, L. Liquori, Resource discovery in the Arigatoni overlay network, in: I2CS: International Workshop on Innovative Internet Community Systems, in: LNCS, Springer, 2006 (in press). Also available as RR INRIA 5928.
- [3] R. Chand, M. Cosnard, L. Liquori, Improving resource discovery in the Arigatoni overlay network, in: ARCS: International Conference on Architecture of Computing Systems, in: LNCS, vol. 4415, Springer, 2007, pp. 98–111.
- [4] M. Cosnard, L. Liquori, R. Chand, Virtual organizations in Arigatoni, in: DCM: International Workshop on Development in Computational Models, ENTCS (in press).
- [5] E. Zegura, K. Calvert, S. Bhattacharjee, How to model an internet, in: Proc. of INFOCOM, IEEE, 1996, pp. 594–602.
- [6] P. Trunfio, D. Talia, H. Papadakis, P. Fragopoulou, M. Mordacchini, M. Pennanen, K. Popov, V. Vlassov, S. Haridi, Peer-to-Peer resource discovery in grids: Models and systems, Future Generation Computer Systems (2007) (in press). Available online 21 December 2006.
- [7] E. Pitoura, AEOLUS, Deliverable D2.1.1: Resource discovery: State of the art survey and algorithmic solutions, Tech. Rep., University of Ioannina, <http://aeolus.ceid.upatras.gr>, 2006.
- [8] H. Jagadish, B.Q. Vu, BATON: A balanced tree structure for Peer-to-Peer networks, in: Proc. of VLDB, ACM, 2005, pp. 661–672.
- [9] K. Aberer, P-Grid: A self-organizing access structure for P2P information systems, in: Proc. of CoopIS, in: LNCS, vol. 2172, Springer, 2001, pp. 179–194.
- [10] P. Ganesan, P. Krishna, H. Garcia-Molina, Canon in G-major: Designing DHTS with hierarchical structure, in: Proc. of ICDCS, IEEE, 2004, pp. 263–272.
- [11] M.J. Freedman, D. Mazières, Sloppy hashing and self-organizing clusters, in: Proc. of IPTPS, in: LNCS, vol. 2735, Springer, 2003, pp. 45–55.
- [12] Planet Lab Consortium. <http://www.planet-lab.org>.
- [13] The Grid 5000 Consortium, <http://www.grid5000.org>.



Raphael Chand is currently a postdoctoral researcher in the Department of Informatics at the University of Geneva, Switzerland. From 2005 to 2006, he was a postdoctoral researcher in the MASCOTTE project team at INRIA Sophia Antipolis, France. In 2001, he received an engineering diploma from Telecom INT, France. In 2002, he received a M.S. in Networks and Distributed Systems from the University of Nice Sophia Antipolis, France. He received his Ph.D. in computer science from Eurecom Institute and the University of Nice Sophia Antipolis, France, in 2005. His research interests include scalable dissemination of information, publish/subscribe systems and sensor networks.



Michel Cosnard, MS 1975 Cornell University, Doctorat d'Etat 1983 Université de Grenoble, served as Professor at ENS Lyon and from 1997, as director of the INRIA Research Unit in Lorraine. In 2001, he was nominated director of the INRIA Research Unit in Sophia Antipolis and served as Professor at the University of Nice-Sophia Antipolis. In 2006, he was appointed Chairman and CEO of INRIA. Michel Cosnard served as Editor in Chief of PPL, and editor of IEEE TPDS, Parallel Computing, Mathematical Systems Theory. He received a prize from the French Academy of Science, the IFIP Silver Core and IPDPS Babbage award.



Luigi Liquori, MS 1990 Udine University, Ph.D. 1996 University of Turin, served as Lecturer at the Ecole des Mines of Nancy from 1999. Since 2001, he has been a senior researcher of INRIA. Luigi Liquori research fields range from logics and foundations of mechanical proof assistants, to semantics of object orientated programming languages, until foundations of Overlay Networks and Pervasive Computing.