

Automatically configuring the network layer of data centers for cloud computing

C. Hu
M. Yang
K. Zheng
K. Chen
X. Zhang
B. Liu
X. Guan

With the requirement of very large data centers for cloud computing, the challenge lies in how to produce a scalable and inexpensive design to interconnect a large number of servers in data centers while providing fault tolerance and high network capacity. Internet Protocol (IP)-based data center networks (DCNs) can support a sufficiently large number of machines and dynamic multi-path routing. However, the complex and error-prone manual configuration of IP-based DCNs hinders its deployment. In this paper, we propose DCZeroconf, a fully automatic IP address configuration mechanism, to eliminate the burden of manual configurations of the IP addresses of servers and switches in data centers. We have implemented a prototype system to validate the effectiveness of DCZeroconf via extensive experiments and simulations. The evaluation results demonstrate that DCZeroconf supports different topologies, and the assigned IP addresses can be automatically adjusted upon dynamic topology changes. In addition, the entire automatic process can be completed in seconds or less.

Introduction

The design of a more efficient and scalable data center network (DCN) has attracted tremendous interest in both research and operation communities, where the physical topology, addressing, and routing mechanisms are accentuated as the three primary issues to be considered [1, 2]. In general, most of the recent research proposals follow four directions. First, the scaling and multi-path routing problem is solved by adding Layer 3 Internet Protocol (IP) switches or routers to inherit the scalability characteristics and equal-cost multi-path (ECMP) forwarding features from IP networks. Given the proven spectacular success of IP networks for interconnecting billions of hosts in the Internet, many practical data centers are already constructed from multiple IP subnets [3, 4]. The second direction involves the construction of DCNs based on Ethernet, because of its low cost, high bandwidth, preset addressing, and automatic configurations. However, traditional Ethernet and its spanning-tree-based forwarding cannot support two required features of data centers: large-scale characteristics (thousands of millions of servers) and multi-path forwarding [5]. Therefore, the Ethernet

proponents seek methods, such as EtherProxy [6], SEATTLE [7], and SPAIN [5], to scale the data centers and to utilize the redundant physical paths. The third direction is the new design of the physical DCN topology construction and addressing, as well as the routing algorithm, e.g., Fat-Tree [1], DCell [2], and BCube [8]. Each of these designs usually provides a unique interconnection topology, designs a specific addressing mechanism based on the topology, and proposes a specific routing mechanism according to the topology and the addressing. Most recently, a fourth direction of study focuses on a way to build hybrid electrical/optical data centers by making use of the very large bandwidth advantage of optical circuit switching technologies such as C-Through [9] and Helios [10].

Compared to other categories of approaches sketched above, IP-based DCNs can elegantly support a sufficiently large number of machines, along with dynamic, flexible load balancing and routing with available protocol stacks and devices. However, an underlying drawback of the IP-based large-scale DCNs is the complex configurations of its network layer. The IP addresses of the servers and intermediate Layer 3 switches or routers (in the remainder of this paper, we use the general term “switch” to refer to a Layer 3 switch) need to be configured before the routing

Digital Object Identifier: 10.1147/JRD.2011.2165751

© Copyright 2011 by International Business Machines Corporation. Copying in printed form for private use is permitted without payment of royalty provided that (1) each reproduction is done without alteration and (2) the Journal reference and IBM copyright notice are included on the first page. The title and abstract, but no other portions, of this paper may be copied by any means or distributed royalty free without further permission by computer-based and other information-service systems. Permission to republish any other portion of this paper must be obtained from the Editor.

0018-8646/11/\$5.00 © 2011 IBM

67 table can be automatically computed. The configuration
68 becomes more complex, difficult to manage, and error-prone
69 for the network administrators, especially when a large
70 number of virtual machines (VMs) dynamically migrate
71 among physical servers. In this paper, we investigate the
72 solution of automatically configuring the network layer of
73 IP-based DCNs.

74 Motivations and goals

75 Several benefits motivate us to seek a fully automatic
76 solution that can configure the IP addresses for VMs, servers,
77 and switches inside the IP-based DCNs. First, the cost of
78 maintaining a data center, especially the overhead on manual
79 IP address configuration, which currently accounts for
80 15%–30% of total capital expense,¹ can be significantly
81 reduced. Second, manually configuring today's data
82 networks is error-prone [11], and curbing manual
83 configuration activities in DCNs can substantially minimize
84 the risk of human errors. Third, the scale of DCN changes
85 according to the customer's demand on applications, which
86 requires frequent configuring and re-configuring of the IP
87 addresses of the machines inside the DCN; therefore, the cost
88 and overhead of evolving (i.e., changing the scale of) the
89 DCNs can be reduced by an automatic address configuration
90 mechanism.

91 The goal of this paper is to design a mechanism that
92 requires zero manual intervention during the configurations
93 on the IP addresses of DCNs with the following specific
94 features. First, each VM/server and each switch port can be
95 automatically and correctly configured with a proper IP
96 address regardless of what physical topology is employed to
97 interconnect the machines. Second, when the DCN topology
98 changes, e.g., due to plugging in or removing a VM/host/
99 switch, the IP addresses can be adaptively adjusted upon
100 the change. Third, the mechanism is scalable to a large
101 number of devices, e.g., hundreds of thousands of machines,
102 and the configuration time for such a large DCN is as fast
103 as seconds or less. Finally, the mechanism should be easy
104 to deploy in today's DCNs and future DCNs. The main
105 contribution of this paper is the design and implementation of
106 an automatic address configuration mechanism for IP-based
107 data centers, which achieves the above goals.

108 Related work and challenges

109 Dynamic Host Configuration Protocol (DHCP) [12] is the
110 most widely deployed scheme for automatic configuration of
111 IP addresses within the same subnet. One or more DHCP
112 servers are employed to record the available IP addresses and
113 to eliminate allocation conflicts of addresses to hosts. When
114 a new host joins the subnet, the host seeks a DHCP server
115 and then requests an unused IP address from the DHCP

¹This information comes from an internal IBM report on information technology (IT) for cloud computing in 2009.

server. The hosts and the DHCP servers should be in the
same subnet so that the broadcast of the DHCP protocol
messages can be received by a requesting host and the DHCP
server. To apply DHCP into data centers, the configuration of
the DHCP servers in each subnet, or the creation of virtual
network cards on a single DHCP server for each subnet, also
requires manual efforts in all the subnets, and the number
of the subnets could be in the thousands. In addition, if the
DHCP relay were employed in a large-scale data center with
hundreds of thousands of machines, the inefficient global
broadcast of DHCP messages would be a large traffic burden.

Zeroconf [13] can also be used for automatically assigning
IP addresses. A Zeroconf-enabled host first randomly
selects an address and validates its availability by
broadcasting queries to the network. The address will be
reserved for the host if no reply shows that the address has
already been occupied; otherwise, the host randomly selects
another address and repeats the validation process. Since
the number of servers in a DCN could be as large as
hundreds of thousands, naive modifications on switches to
allow Zeroconf broadcasting inside such a large DCN will be
costly and inefficient.

In the context of DCNs, Portland [14] and DAC [15]
can automatically assign addresses. Portland develops a
distributed location discovery protocol (LDP) to allocate the
physical media access control (PMAC) addresses of servers
and switches. PMAC is the private address dedicated for
multi-rooted tree topology, e.g., Fat-Tree [1]. Therefore, the
specific design of Portland on PMAC is not easily extended to
other network topologies. DAC has a centralized server that
learns the physical topology of the network and maps it to the
logical addresses with the help of a blueprint. DAC is a generic
method to configure all kinds of addressing on an arbitrary
topology; however, it focuses on the initial setup of a DCN
and does not seriously consider the cases in which the data
center topology changes. To solve the topology-change
problem with DAC, manual intervention is still required to
input a new blueprint, which depicts the changed topology
and corresponding addresses. Furthermore, the configurations
of VMs by DAC are also left unknown in its proposal.

By reviewing the related work in the literature, we
summarize the following challenges to achieve all the goals
mentioned above. First, for scalability reasons, large-scale
data centers are divided into a considerable number of
subnets (e.g., thousands of subnets), and this feature limits
the power of DHCP or Zeroconf in DCN, which can only
solve the address conflicts in a single subnet.

Second, assigning IP addresses to intermediate devices like
switches and routers also constitutes an indispensable
configuration procedure to enable the devices to forward
ordinary IP packets. However, the traditional DHCP or
Zeroconf only configures the IP addresses for end hosts
but not for the intermediate IP switches or routers used to
build up the underlying communication channels.

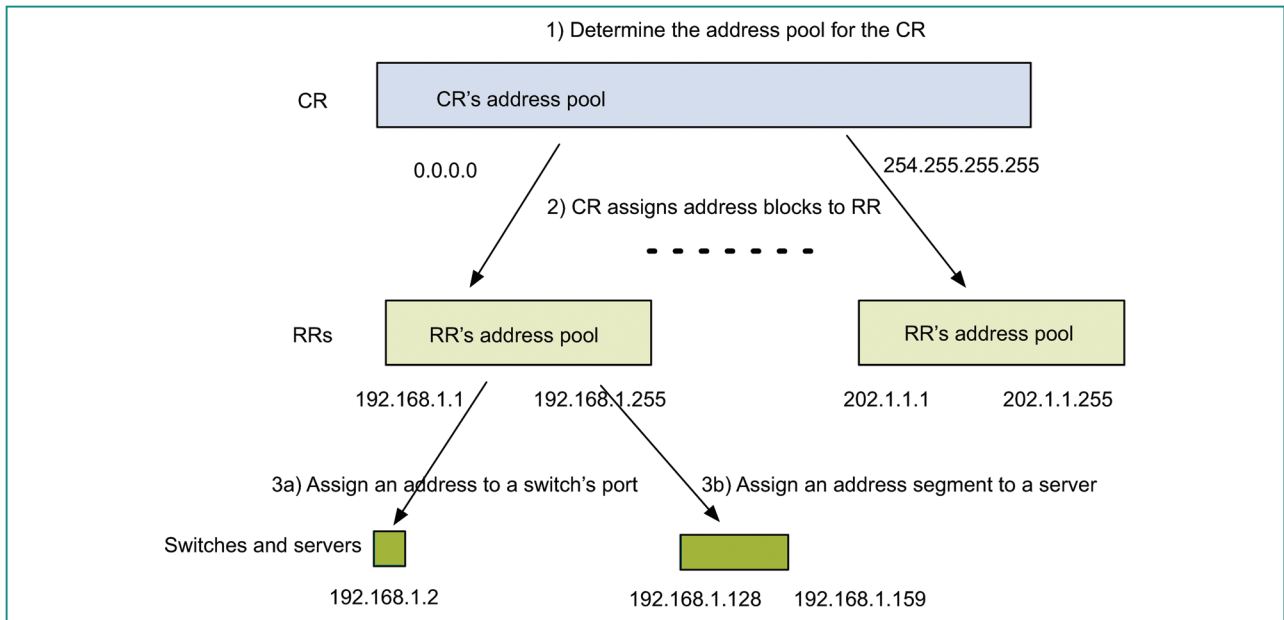


Figure 1

Addressing hierarchy.

170 Third, in the context of cloud computing, emerging
 171 applications or services will be constantly changed or
 172 reconfigured by the providers to adapt to new customer
 173 demands or marketing feedbacks [3]. Currently, the topology
 174 change is not uncommon in the data centers. New servers
 175 will be added to a data center if the computing demand on the
 176 data center is increased. As the demand continues to increase,
 177 the scale of the DCNs continues to evolve. Therefore, it is
 178 crucial, yet challenging, that the IP address configurations
 179 can automatically adapt to such topology changes.

System overview

181 In this paper, we propose DCZeroconf to automatically
 182 configure the network layer of DCN. Before delving into the
 183 design details, we first define the terminologies as follows.
 184 In modern data centers, the servers are placed in *racks* and
 185 the servers in the same rack are usually connected to the same
 186 switch. Such a switch is known as a *top-of-rack* (TOR)
 187 switch. The servers and TOR switches form the access layer
 188 of a data center, and to connect the TOR switches in an
 189 access layer, more Layer 3 switches are employed. For the
 190 topology used to connect the racks, DCZeroconf makes no
 191 specific assumptions about it.

192 DCN should support rapid VM migration, which in
 193 turn calls for separating names from locations. We follow
 194 the naming convenience discussed in [3], where an
 195 application-specific address (AA) and a location-specific
 196 address (LA) are maintained for each VM. DCZeroconf will
 197 assign a LA to a physical server, and the LA is routable

198 inside the data center. An AA will be assigned to a VM, with
 199 which the applications can identify a service running on a
 200 specific server/VM. When a VM is migrated from one server
 201 to another, the routable LA has to be changed, while the
 202 AA remains the same. Also, the mapping between the AA
 203 and LA will be changed, and the directory service will keep
 204 the latest mapping between AA and LA. With the mapping
 205 between an AA of a specific server/VM and the LA of
 206 the physical server it locates, one can find and reach the
 207 service running on that server/VM.

208 The DCZeroconf system is a two-tier mechanism that
 209 automatically configures the IP addresses inside the data
 210 center. We have a top configuration server called a *Central*
 211 *Regulator* (CR) and a configuration server called a *Rack*
 212 *Regulator* (RR) in each rack. Note that one RR is required for
 213 each broadcast domain if several broadcast domains are in a
 214 same rack. **Figure 1** provides an overview of how
 215 DCZeroconf works and can be understood as follows.

216 In the first step, the network administrator determines an
 217 available IP address pool that can be used in the data center
 218 and provides it as input to the CR, and this is the only manual
 219 effort required by DCZeroconf during the automatic
 220 configuration. For the example in Figure 1, the IP addresses
 221 for the servers inside the data center could be employed
 222 for internal use only, and a Network Address Translation
 223 (NAT) is functioning at the gateway; thus, the entire 32-bit
 224 address space (for IPv4) can be the pool for the CR.

225 Next, the CR partitions the available IP addresses into
 226 blocks and informs each RR regarding the block(s) that can

227 be used to configure the servers and switches within the
228 corresponding rack. For instance, the CR will send an
229 address block 192.168.1.1–192.168.1.255 to a RR when
230 receiving a request.

231 After the RR is configured, it then automatically assigns
232 the IP addresses to each VM, server, and switch. We may
233 continue with the example above for which the address block
234 192.168.1.1–192.168.1.255 is assigned to the RR. The RR
235 selects one address from the block to configure itself first,
236 e.g., 192.168.1.1. Then, the RR replies to the address
237 configuration request from the port of any switch with one of
238 the addresses in the block, e.g., 192.168.1.2. In addition, the
239 RR will send an IP address segment to any server in the same
240 rack after receiving an address configuration request from
241 the requesting server, e.g., 192.168.1.128–192.168.1.159.
242 Each VM in that server will be allocated with one IP address
243 from the address segment.

244 During the procedure associated with DCZeroconf, the
245 difficulty lies in how to coordinate the address assignment
246 between 1) CR and RR, 2) RR and servers/VM, and 3) RR
247 and switches.

248 **Address assignment between CR and RR**

249 To assist the negotiation between the CR and RR,
250 connections between each RR and the CR are required that
251 form a control plane path over the data path for traffic
252 forwarding. The introduction of such an extra control plane is
253 not associated with a large cost, and many literature studies
254 make use of a center regulator/scheduler to control the
255 switches/servers in a data center and indicate such a control
256 plane [3, 14, 16]. In fact, there is no critical bandwidth
257 requirement on a control plane, so that a commodity switch
258 fabric can be used to construct the communication channel.
259 Aside from connecting them with direct cables, a wireless
260 networking with encrypted communication channel for
261 configuration purposes could also be a viable option.

262 Constructing such a wireless control plane network does
263 increase capital investment, but it is relative small compared
264 to the cost of manual efforts that could be reduced. Only one
265 wireless network interface card (NIC) is needed for each
266 rack, and in each rack there could be 30–126 servers.
267 Suppose the cost of one server is \$1,000, and the cost of one
268 wireless NIC is \$10. Even without considering the cost on
269 the switches and routers, the increased cost compared to the
270 cost of servers is only 0.008%–0.03%. In addition, the
271 communication range of 802.11 could be 100 meters, i.e., it
272 can serve a 20,000-m² square room, which is sufficient to
273 accommodate a large DCN. Furthermore, there should be no
274 concern with the configuration on wireless networking for
275 each NIC. The same script could be copied to each RR and
276 CR to connect to the access point when the server boots.

277 A RR is a lightweight process that can be run in any of the
278 servers in a rack. A RR first requests an IP address block
279 (a number of continuous IP addresses) from the CR, which

280 can be assigned to VMs, servers, and switches. On receiving
281 a valid request, the CR assigns RR with the IP address block
282 as well as the corresponding lease (length of time when the
283 allocation is valid). The lease is employed to detect any
284 physical detachment of a certain rack, and each RR should
285 renew its contract with the CR before the lease expires.
286 Otherwise the CR will reclaim the assignment, and the
287 corresponding address block is available for reassignment.
288 The procedure is typically initiated immediately after RR
289 being booted, and a RR can start the assignment of the IP
290 addresses to VMs/servers in the same rack and switches upon
291 receiving the allocated address block from the CR. If any
292 message in each step is lost, retransmission will be activated
293 in order not to fail the procedure.

294 **Address assignment between the 295 RR and servers/VMs**

296 The servers/VMs and RR in a rack should be in a flat
297 Layer 2 network, i.e., the same Layer 2 broadcasting domain
298 such as Ethernet or a virtual local area network (VLAN).
299 When a RR is allocated with an IP address block from the
300 CR, it works in a similar way as a DHCP server to assign
301 IP addresses to the local servers/VMs. The only difference
302 with a DHCP server is that RR can assign multiple IP
303 addresses in a batch to a single server in the same rack,
304 while a DHCP server assigns only one IP address to a single
305 host. This difference is due to the existence of several VMs
306 in a single server.

307 As opposed to assigning one IP address for a VM each
308 time, we propose a “batch mode” that assigns multiple
309 addresses at once, which is more efficient. A server in a rack
310 first initiates a broadcasting *Discover* message seeking the
311 RR server. On receiving the *Discover* message, the RR sends
312 an *Offer* message to notify the requesting server. Then, the
313 server sends a *Request* message to request an address from
314 the RR, and the RR in turn replies with an Acknowledgement
315 (ACK) message to assign a number of addresses, i.e., an
316 *address blocks*. To identify an address block that can be
317 allocated to the VMs on the requesting server, the *ACK*
318 message contains the first IP address and the length of the
319 address block. If no address block is available, a *Negative*
320 *Acknowledgement* (NACK) message will be generated. Upon
321 receiving the address block on the server side, a hypervisor in
322 the server, which allows multiple VMs to run concurrently on
323 a server, performs the address assignment among the local
324 VMs in the server.

325 Original DHCP can be an alternative option for the
326 operation between the RR and servers. The servers and the
327 RR in a rack are in a same Layer 2 network. When the RR is
328 equipped with the IP address block from CR, in order to
329 assign IP addresses to servers in the same rack, it works like a
330 classic DHCP server. Servers in each rack act as DHCP-
331 enabled clients. If there are several VMs in one server, each
332 VM triggers an IP address assignment process via DHCP.

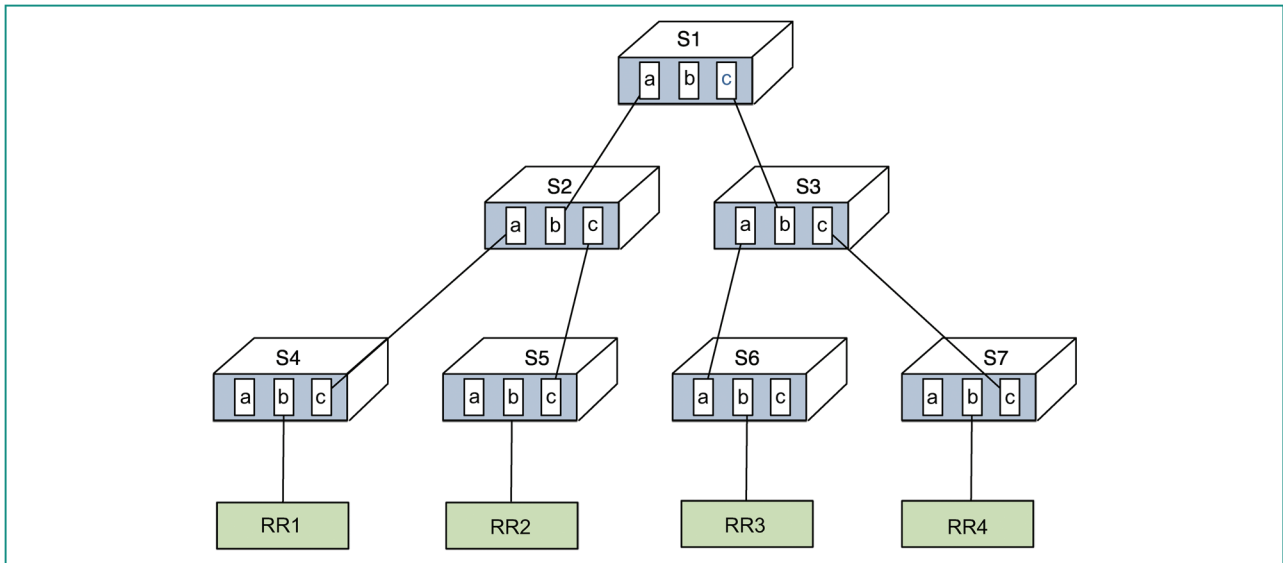


Figure 2

Example of a switch configuration. S1 through S7 in the figure are switches. In each switch, a, b, and c denote the different ports, and RR1 through RR4 are rack regulators.

333 Although the DHCP mechanism can be used between the RR
 334 and VMs, the “batch” mode mentioned provides a more
 335 efficient way that configures several machines
 336 simultaneously.

337 If one RR runs out of the available IP addresses that can be
 338 allocated to servers/VMs, it redirects the query message to
 339 another RR. The redirection follows a pre-determined order
 340 managed by a cyclic linked list and each RR is informed of
 341 the next RR to resort to when the CR configures the RR.

Address assignment between the RR and switches

342 In this section, we first describe how the RRs and switches
 343 exchange messages followed by the communication among
 344 them to assign addresses. For the ease of presentation yet
 345 without loss of generality, we depict a simple tree-based
 346 topology in **Figure 2** as an example topology. S1 through S7
 347 are switches, and a/b/c in each switch denotes the different
 348 ports. RR1 through RR4 are rack regulators.

349 The switches can query an arbitrary RR for obtaining
 350 available IP addresses. However, switches cannot form
 351 Layer 3 routing paths until they are assigned with IP
 352 addresses, and the switches also forbid the broadcast among
 353 themselves by default. A brute-force approach is to allow
 354 Layer 2 broadcasting of all the protocol messages in the
 355 entire DCN. Although the protocol messages in the control
 356 plane contribute relatively little traffic to the traffic in the data
 357 plane, brute-force broadcasting is inefficient and requires
 358 more time to perform the configuration. To enable a switch to
 359 communicate with a RR for requesting IP addresses,
 360 DCZeroconf follows a bottom-up approach and utilizes
 361

Layer 2 bridging to relay the DCZeroconf protocol messages
 among switches and RRs, akin to a path-vector protocol.

362 Each switch periodically sends a message to all its
 363 neighbors (i.e., servers and other switches), requesting “Can
 364 you reach a RR?” until the switch receives a reply indicating
 365 a valid next hop to a RR. During the bootstrapping phase,
 366 only the switches directly connected with a RR can receive a
 367 reply from the RR specifying, “I am a RR.” Then,
 368 recursively, such a switch S will answer the requests from its
 369 neighbors with the path from S to the RR that S has
 370 discovered (S prepends itself in the path). If there is more
 371 than one path received by a certain switch S, only the first
 372 received path is used by S to reach the RR and is further
 373 propagated to its neighbors. The first received path is likely
 374 to indicate the path with the least delay from the switch to a
 375 RR, since switches relay answers with path information,
 376 hop-by-hop from the RRs. For example, consider Figure 2.
 377 S1 through S7 periodically send path queries. S4, S5, S6, and
 378 S7 first receive replies from RR1, RR2, RR3, and RR4,
 379 respectively. S2 will then receive path information of “S4
 380 RR1” and “S5 RR2” from S4 and S5, respectively. S2 selects
 381 one of these two paths, and places itself in the path, i.e.,
 382 “S2 S5 RR2.”
 383
 384

385 RR only sends messages to switches after receiving
 386 messages from switches. When sending a message packet
 387 from a switch to the RR, the intermediate switches will add the
 388 path information to the packet. The path information can be
 389 treated as a stack. The source switch and all the intermediate
 390 nodes along the path to the RR will push their addresses to the
 391 stack. This path information is also kept in the packet back to

392 the switch from the RR. Thus, the reply messages from the RR
393 to the switch always can find their reverse way, making use of
394 the encapsulated path information in the packet. In the
395 example of Figure 2, when S1 sends a message to RR, the
396 packet carries the path information “S1 S2 S5 RR2.” RR (here,
397 in this example, is RR2) can compute the reverse path for
398 its reply message “S5 S2 S1.” The RR pops the first element
399 and send reply packet to S5. Now, when S5 receive the packet,
400 the reverse path information is “S2 S1.” S5 also pops an
401 element from the path stack and finds the next hop as S2. This
402 process continues, and eventually the packet can reach S1.
403 This part is similar to source-routing, and the switches do not
404 keep a forwarding table but just examine the packets for the
405 path information.

406 Address aggregation is usually used for the sake of
407 scalability and is able to keep the routing table small. To
408 utilize address aggregation for constructing a space-efficient
409 routing table in each switch, the addresses of the two ends of
410 a link between two switches are always within the same
411 /31 prefix. (Here, the number following the slash is the prefix
412 length, the number of shared initial bits, counting from the
413 most-significant bit of the address. Thus, the /31 prefix is
414 an address block with a 31-bit prefix.) Whenever a RR
415 receives an address request from a switch S for a certain
416 port p, RR returns two continuous IP addresses for this port
417 and the port in the neighboring switch of S that connected
418 with p. With the two returned IP addresses, the switch S
419 selects one of them and attempts to assign this address to the
420 connected port in the neighboring switch. The port of the
421 neighboring switch accepts the assignment and configures its
422 IP address as indicated in the message if this port has not
423 received any address assignment from any RR. In case a
424 certain port p of switch S receives two assignments from both
425 the RR and a neighboring switch N, switch S compares
426 the media access control (MAC) addresses of the two ends of
427 the link. The address assigned from the switch port with a
428 larger MAC address will be accepted. An ACK/NACK
429 message will be then triggered to accept/decline the
430 assignment from the RR or a neighboring switch.

431 Consider Figure 2 again as an illustration. S2.a requests
432 a pair of addresses from the RR and the RR replies with a
433 pair of addresses, A_1 and A_2 . S2.a first selects one address,
434 say A_2 , to allocate to S4.b. If S4.b returns an ACK, S2.a
435 configures itself with A_1 and sends an ACK to RR. In the
436 case where S4.b sends another IP address, e.g., A_3 to S2.a,
437 S2.a compares its MAC address with S4.b. If the MAC
438 address of S4.b is larger than S2.a, S2.a selects A_3 as its
439 address; otherwise, S2.a resends an address assignment
440 message to S4.b, and an ACK is expected to arrive later.

441 If the RR does not find any available IP addresses that
442 can be allocated, it redirects the query message to another
443 RR. The redirection follows a predetermined order managed
444 by the CR. The CR is aware of all the configured RRs,
445 and it keeps a cyclic linked list for the RRs.

Evaluations

446 The major performance metric of DCZeroconf is the time
447 it takes to configure a given DCN. In this section, we first
448 examine the configuration time on a small-scale testbed
449 and then check the larger-scale DCNs via simulations.
450 Moreover, experiments are also performed to investigate
451 the configuration time of DCZeroconf when topology
452 changes.
453

Experiments on testbed

454 We first examine the performance of DCZeroconf on an
455 eight-node prototyping testbed. In this testbed, we also
456 emulate two racks, each of which has one RR and a server
457 with two VMs. A CR is connected with two RRs, over a
458 wireless connection. The wireless connection among CR and
459 RRs is a 54-Mb/s 802.11g connection, and the connections
460 among CR, RR, servers, and switches are all with 100-Mb/s
461 Ethernet. In addition, three desktop computers are
462 deployed with the eXtensible Open Router Platform (XORP)
463 [17] and act as the Layer 3 switches. Switch no. 1 connects
464 with Switch no. 2 and Switch no. 3, and Switch no. 2 or 3 has
465 two links connecting with the two emulated racks.
466

467 The entire configuration of DCZeroconf contains four
468 phases: 1) *RR configuration phase*, in which the CR
469 configures the RRs; 2) *server configuration phase*, during
470 which servers are configured; 3) *communication channel*
471 *construction (CCC) phase*, which builds the communication
472 channel among switches and RRs for IP address allocation;
473 and 4) *switch configuration phase* that configures the IP
474 addresses for all the ports of each switch. We use T_{RR} , T_{server} ,
475 T_{CCC} , and T_{switch} to denote the configuration time of the
476 above four phases, respectively. Servers periodically send
477 discover messages to RR before RR is configured, and
478 switches will resend request messages until any reply is
479 received in an interval. These two processes in our
480 implementation have the same retry interval, which is
481 denoted as $T_{interval}$. By performing experiments on the
482 testbed, we are able to obtain the time of the four
483 configuration phases, as well as the total time to configure all
484 the switches/servers/VMs.
485

486 The experiments are divided into three groups. In each
487 group, we repeat the experiments 10 times, and $T_{interval}$ is set
488 to be 10 ms, 50 ms, or 100 ms in the three groups,
489 respectively. We average the measured results in each group
490 and illustrate them in **Table 1**. From the results, we make
491 several observations. First, the value of $T_{interval}$ does not
492 affect T_{RR} , T_{switch} , T_{CCC} , and T_{server} significantly, which is
493 in accord with the design of DCZeroconf.
494

495 Second, the time for the “CCC phase” and the “Switch
496 configuration phase” is much larger than the other two
497 phases, because the protocol messages in these two phases
498 are propagated hops away, while the “RR configuration
499 phase” and “Server configuration phase” only exchange
500 protocol messages in a same broadcast domain. In addition,
501

Table 1 Time (ms) consumed during the configuration on the testbed.

	T_{interval}	T_{CCC}	T_{RR}	T_{switch}	T_{server}	T_{total}
Experiment #1	10	180.6	2.5	12.2	2.5	200.1
Experiment #2	50	170.0	2.7	11.6	2.6	223.9
Experiment #3	100	171.9	2.5	12.4	2.4	266.5

CCC phase” needs to record the forwarding paths and port into its local memory. This is why the dominant time in the entire configuration appears to be the construction of the communication channel.

Third, we analyze the range of the configuration times and check whether our testing falls into the range. The “RR configuration phase” and the “CCC phase” can be started at the same time. Once RRs being configured, the “Server configuration phase” starts to assign IP addresses to servers/VMs, while the “Switch configuration phase” starts to work when “RR configuration phase” and “CCC phase” are completed. Further, with the desperate retry interval T_{interval} and the second observation, we derived the range of the entire configuration time as the following:

$$T_{\text{CCC}} + T_{\text{switch}} \leq T_{\text{total}} \leq T_{\text{CCC}} + T_{\text{switch}} + T_{\text{interval}}. \quad (1)$$

The measured results in Table 1 are obviously in this range. Moreover, as indicated in (1) and Table 1, increasing T_{interval} increases the total configuration time. A very short T_{interval} will increase the burden on servers, switches and the network, and we set it to be 100 ms in the remainder of the evaluation experiments.

Simulations

We analyze the parameters required in the simulation and then determine the values by measuring using the testbed. By substituting the values of these parameters, we estimate the configuration time of DCZeroconf using larger-scale DCNs.

In [15], the authors evaluate DAC under several data center topologies. To compare with DAC, we also examine the configuration time of DCZeroconf on the same topologies. The experimental topologies include BCube [8], Fat-Tree [1], VL2 [3], and DCell [2]. Although in the original proposals of BCube, Fat-Tree, and DCell, the authors have also introduced specific addressing methods instead of IP addressing, the physical topologies of BCube, Fat-Tree, and DCell could also be used to employ IP-based DCNs.

We use simulation to estimate the time for DCZeroconf to complete the configuration of an IP-based DCNs on different topologies, and the detailed results are illustrated in **Table 2**. The results of DAC are cited from [15],

Table 2 Configuration time (ms) in different topologies. The reader is referred to [15] for an explanation of the numbers in parentheses in the first column.

Topologies	No. of devices	DAC	DCZeroconf
BCube(4,4)	2,304	310	144.2
BCube(8,4)	53,248	2,000	166.9
Fat-tree(20)	2,500	262	192.9
Fat-tree(100)	262,500	6,223.6	247.3
VL2(20,100)	52,650	606.7	211.6
VL2(100,100)	252,650	2,032.3	227.6
DCell(2,3)	2,709	353	295.8
DCell(6,3)	3,807,349	45,578.1	3,276.3

which utilizes the specific addressing for the corresponding topology. The difference in addressing does differ with the protocol messages, but we believe that this effect is minor since the length of address does not vary very much. Furthermore, since DAC does not support server virtualization, the number of VMs is set to one in all the simulations using DCZeroconf.

Table 2 shows that the configuration time of DCZeroconf is shorter than the delay of DAC, and the gap is substantial for large topologies. The reason for the significant gaps in large topologies comes from the computation of a “mapping” between logical identification (ID) and physical ID for DAC. It is demonstrated in [15] that the time for mapping dominates the entire configuration time for large topologies. DCZeroconf does not consider such mapping during its configuration. However, as we mentioned in the introduction, it is just one of the useful characteristics of IP-based DC, which does not require us to embed the location information into the network ID.

Note that, when we set our simulation environment, all the parameters are measured from the test-bed with 100-Mb/s links. If the connections increase to 1 Gb/s as researchers in a DAC paper [3] used for their experiments, the configuration time of DCZeroconf could be shorter.

Configuration time when topology changes

It is quite common that a data center will enlarge its scale by gradually adding new servers and switches. In this section, we evaluate the configurations time in such cases.

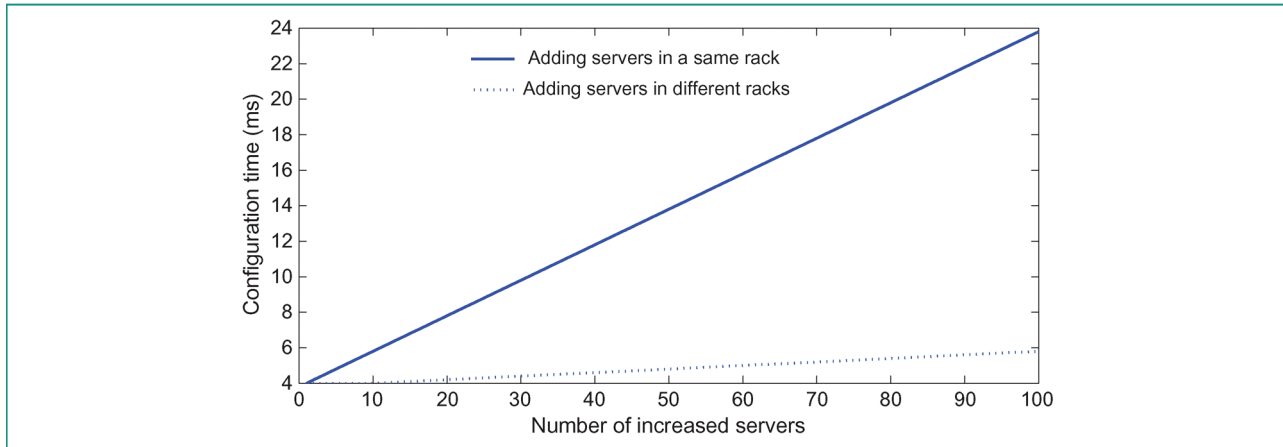


Figure 3
Configuration time when adding servers.

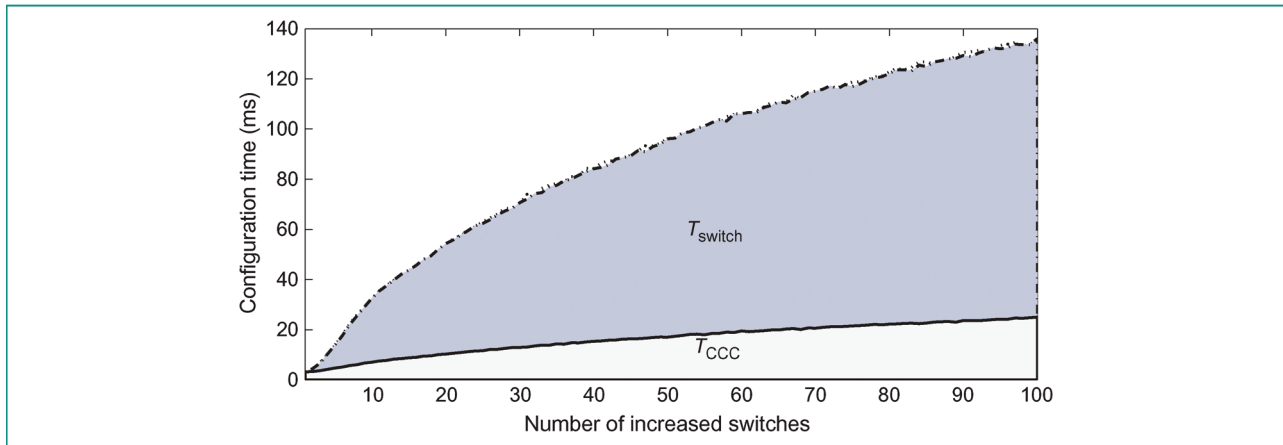


Figure 4
Configuration time when adding switches.

568 First, we study the case when adding servers. The required
 569 configuration time versus the number of the increased
 570 servers is depicted in **Figure 3**. This figure shows two
 571 extreme cases when adding servers. The solid line represents
 572 the case in which all the servers are added in a same rack.
 573 The dashed line indicates the results for the case when all the
 574 servers are evenly added to all the racks (in this experiment,
 575 the number of racks is 20). The increased configuration
 576 time is very small and even when we increase by 100 servers
 577 in one rack, the configuration time is less than 24 ms. The
 578 results demonstrate a linear bound of the scalability when
 579 adding servers: the configuration time linearly increased
 580 with the increase of the number of the servers. Since the
 581 configuration in different racks can be performed in parallel,
 582 the slope of the dash line is much smaller than the solid

583 line. Please note that the configuration time of any other
 584 tested cases is within the range of these two curves
 585 in Figure 3.

586 Second, we investigate the performance of DCZeroconf
 587 when adding switches. The required configuration time
 588 versus the number of the newly added switches is depicted in
 589 **Figure 4**. In this simulation, we first generate a partially
 590 connected DCN and then add switches to randomly selected
 591 locations. It is worth noting that in this experiment, T_{switch}
 592 and T_{CCC} do not apply, and the configuration time consists
 593 of T_{switch} and T_{CCC} . We run the experiment for 1,000 times,
 594 and the results in Figure 4 are the average value (the result
 595 is for a classical tree-based topology). The shaded area in
 596 Figure 4 indicates the difference between T_{switch} and T_{CCC} .
 597 The total configuration time, as well as T_{switch}/T_{CCC} ,

598 increases sublinearly with the number of the added switches,
599 and this fact indicates that DCZeroconf is scalable with
600 the enlargement of the data center. In addition, the entire
601 configuration time including T_{CC} and T_{switch} is less than
602 136 ms if we add 100 switches.

603 Conclusion

604 In this paper, we have designed and implemented the
605 DCZeroconf mechanism to automatically configure the
606 IP addresses inside the data center. DCZeroconf successfully
607 achieves the goals illustrated in the Introduction. First, it
608 eliminates the possible IP address conflicts in different
609 subnets without any assumption on the data center topology.
610 Second, as the topology changes, DCZeroconf does not
611 need to reconfigure the entire network; instead, RR
612 only assigns the new devices with IP addresses. Third,
613 the configuration of DCZeroconf is fast. It requires
614 200.1–266.5 ms to configure a real testbed with eight nodes
615 and about 3.3 seconds in the simulation with more than
616 3 million devices. The experiments on the topology
617 changes also indicate its scalability to large-scale DCNs
618 (e.g., clouds of clouds). Finally, DCZeroconf is
619 incrementally deployable in existing and future DCNs,
620 as it only requires minimal modifications on the software
621 stack of the current switches and servers.

622 Acknowledgments

623 This work was partly supported by National Nature Science
624 Foundation of China (NSFC) (60903182, 60921003,
625 60873250, 60736027), 973 plan (2007CB310702),
626 863 plan (2007AA01Z480), Tsinghua University Initiative
627 Scientific Research Program, open project of State Key
628 Laboratory of Networking and Switching Technology,
629 and 111 International Collaboration Program of China.

630 References

- 631 1. M. Al-Fares, A. Loukissas, and A. Vahdat, "A scalable,
632 commodity data center network architecture," in *Proc. SIGCOMM*,
633 Seattle, WA, 2008, pp. 63–74.
- 634 2. C. Guo, H. Wu, K. Tan, L. Shi, Y. Zhang, and S. Lu, "DCCell:
635 A scalable and fault-tolerant network structure for data centers," in
636 *Proc. SIGCOMM*, Seattle, WA, 2008, pp. 75–86.
- 637 3. A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim,
638 P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta, "VL2: A scalable
639 and flexible data center network," in *Proc. SIGCOMM*, Barcelona,
640 Spain, 2009, pp. 51–62.
- 641 4. D. Roisman, "Data center top-of-rack switch redundancy models,"
642 in *Proc. NANOG 46*, Philadelphia, PA, 2009.
- 643 5. J. Mudigonda, P. Yalagandula, M. Al-Fares, and J. Mogul,
644 "SPAIN: COTS data-center Ethernet for multipathing over
645 arbitrary topologies," in *Proc. NSDI*, San Jose, CA, 2010,
646 pp. 265–280.
- 647 6. K. Elmeleegy and A. Cox, "EtherProxy: Scaling Ethernet by
648 suppressing broadcast traffic," in *Proc. INFOCOM*, Rio de Janeiro,
649 Brazil, 2009, pp. 1584–1592.
- 650 7. C. Kim, M. Caesar, and J. Rexford, "Floodless in SEATTLE:
651 A scalable Ethernet architecture for large enterprises," in *Proc.*
652 *SIGCOMM*, Seattle, WA, 2008, pp. 3–14.
- 653 8. C. Guo, G. Lu, D. Li, H. Wu, X. Zhang, Y. Shi, C. Tian, Y. Zhang,
654 and S. Lu, "BCube: A high performance, server-centric network

- architecture for modular data centers," in *Proc. SIGCOMM*,
Barcelona, Spain, 2009, pp. 63–74. 655
9. G. Wang, D. G. Andersen, M. Kaminsky, K. Papagiannaki,
T. E. Ng, M. Kozuch, and M. Ryan, "C-Through: Part-time optics
in data centers," in *Proc. SIGCOMM*, New Delhi, India, 2010,
pp. 327–338. 657
10. N. Farrington, G. Porter, S. Radhakrishnan, H. H. Bazzaz,
V. Subramanya, Y. Fainman, G. Papen, and A. Vahdat, "Helios:
A hybrid electrical/optical switch architecture for modular data
centers," in *Proc. SIGCOMM*, New Delhi, India, 2010,
pp. 339–350. 661
11. C. Kim, "Scalable and efficient self-configuring networks,"
Ph.D. dissertation, Princeton Univ., Princeton, NJ, 2009. 662
12. R. Droms, "Dynamic host configuration protocol," RFC 2131,
1997. [Online]. Available: <http://www.ietf.org/rfc/rfc2131.txt> 663
13. B. A. S. Cheshire and E. Guttman, "Dynamic configuration of
IPv4 link-local addresses, Zeroconf," DFC 3927, 2002. [Online].
Available: <http://www.ietf.org/rfc/rfc3927.txt> 664
14. R. Niranjani Mysore, A. Pamboris, N. Farrington, N. Huang,
P. Miri, S. Radhakrishnan, V. Subramanya, and A. Vahdat,
"Portland: a scalable fault-tolerant layer-2 data center
network fabric," in *Proc. SIGCOMM*, Barcelona, Spain,
2009, pp. 39–50. 665
15. K. Chen, C. Guo, H. Wu, J. Yuan, Z. Feng, Y. Chen, S. Lu, and
W. Wu, "Generic and automatic address configuration for data
center networks," in *Proc. SIGCOMM*, New Delhi, India, 2010,
pp. 39–50. 666
16. M. Al-Fares, S. Radhakrishnan, B. Raghavan, N. Huang, and
A. Vahdat, "Hedera: Dynamic flow scheduling for data center
networks," in *Proc. NSDI*, San Jose, CA, 2010, p. 19. 667
17. XORPDownload code of Xorp. [Online]. Available: <http://www.xorp.org/downloads.html> 668

Received January 31, 2011; accepted for publication
March 19, 2011 688

Chengchen Hu MOE Key Lab for Intelligent Networks and
Network Security, Department of Computer Science and Technology,
School of Electronic and Information Engineering, Xi'an Jiaotong
University, Xi'an, 710049 China; also guest researcher with SKLNST
Lab, Beijing University of Posts and Telecommunications, Beijing
100876, China (huc@ieee.org). Dr. Hu received his Ph.D. degree from
the Department of Computer Science and Technology of Tsinghua
University in 2008. He worked as an assistant research professor in
Tsinghua University from June 2008 to December 2010 and is currently
an Associate Professor in the Department of Computer Science and
Technology of Xi'an Jiaotong University. His main research interests
include computer networking systems and network measurement and
monitoring. 690

Mu Yang Department of Computer Science and Technology,
Tsinghua University, Beijing 100084, China (yangmu266@gmail.com).
Mr. Yang is an undergraduate student in the Department of computer
science and technology, Tsinghua University. 703

Kai Zheng IBM Research Division, China Research Laboratory,
Beijing 100074, China (zhengkai@cn.ibm.com). Dr. Zheng received his
M.S. and Ph.D. degrees both in computer science from Tsinghua
University, Beijing, China, in 2003 and 2006, respectively. He is
currently working with IBM Research China. His research interests
include high-speed packet forwarding, pattern matching associated with
network-security issues, and new data center network architecture
design. 707

Kai Chen Department of Electrical Engineering and Computer
science, Northwestern University, Evanston, IL, 60201 USA
(kchen@northwestern.edu). Mr. Chen is currently a Ph.D. student in the
EECS Department at Northwestern University, Evanston, IL. Prior
to this, he received his B.S. and M.S. degrees in computer science in 715

720 2004 and 2007, respectively, both from University of Science and
721 Technology of China, Hefei, China. He is interested in finding simple
722 yet elegant solutions to real networking and system problems.

723 **Xin Zhang** *Computer Science Department, Carnegie Mellon*
724 *University, Pittsburgh, PA 15213 USA (xzhang1@cs.cmu.edu).*
725 Mr. Zhang is currently a Ph.D. candidate in the computer science
726 department at Carnegie Mellon University since 2006, working with
727 Dr. Adrian Perrig and Dr. Hui Zhang. He received his B.S. degree in
728 Department of Automation in 2006 from Tsinghua University. His
729 academic interests include network security, applied cryptography,
730 Internet routing, and clean-slate network architectural design. He is also
731 passionate about web technologies and web application development.

732 **Bin Liu** *Department of Computer Science and Technology,*
733 *Tsinghua University, Beijing 100084, China (liub@tsinghua.edu.cn).*
734 Professor Liu is currently a full Professor in the Department of
735 Computer Science and Technology, Tsinghua University. His current
736 research areas include high-performance switches/routers, network
737 processors, high-speed security, and making Internet use more energy
738 efficient. Professor Liu has received numerous awards from China
739 including the Distinguished Young Scholar of China.

740 **Xiaohong Guan** *MOE Key Lab for Intelligent Networks and*
741 *Network Security, School of Electronic and Information*
742 *Engineering, Xi'an Jiaotong University, Xi'an, 710049 China*
743 *(xhguan@mail.xjtu.edu.cn).* Dr. Guan received his B.S. and M.S.
744 degrees in control engineering from Tsinghua University, Beijing,
745 China, in 1982 and 1985, respectively, and his Ph.D. degree in
746 electrical engineering from the University of Connecticut in 1993.
747 He was a senior consulting engineer with PG&E (Pacific Gas and
748 Electric Company) from 1993 to 1995. He visited the Division of
749 Engineering and Applied Science, Harvard University, from January
750 1999 to February 2000. Since 1995, he has been with the Systems
751 Engineering Institute, Xi'an Jiaotong University, and was appointed
752 IEEE fellow in 2007, Cheung Kong Professor of Systems Engineering
753 in 1999, and dean of the School of Electronic and Information
754 Engineering in 2008. Since 2001, he has been the director of the Center
755 for Intelligent and Networked Systems, Tsinghua University, and
756 served as head of the Department of Automation, 2003–2008. He is
757 an Editor of *IEEE Transactions on Power Systems* and an Associate
758 Editor of *Automatica*. His research interests include optimization and
759 security of networked systems, computer network security, and sensor
760 networks.

AUTHOR QUERY

AUTHOR PLEASE ANSWER QUERY

Note that Ref. [17] has a web address that is not available. Please check.

END OF AUTHOR QUERY