

Going Back and Forth: Efficient Multi-Deployment and Multi-Snapshotting on Clouds

Dr. Ashiquee Rasool Mohammad¹, Syeda Farhath Begum² and Dr. Kahalid Mohiuddin³

¹ King Khalid University

Received: 8 April 2012 Accepted: 2 May 2012 Published: 15 May 2012

Abstract

Cloud computing has changed the way people think of using resources. Especially, the IaaS (Infrastructure as a Service) allows users to make use of unlimited resources in pay per use fashion. Virtualization is the technology based on which the cloud service providers are able to provide or share computational resources and data centers to users. Though this approach is practical, it throws certain challenges in terms of designing and development of IaaS middleware. One such challenge is the need for deploying thousands of VM instances to meet the requirements of growing number of users. In the process another challenge is to snapshot multiple images and persisting them towards management tasks like stopping VMs temporarily and resuming them as and when required. The presence of data centers in different configurations enables the simultaneous deployment and snapshotting is important. This capability should be coupled with another feature that is the whole mechanism should be hypervisor independent. To achieve this, a new virtual file system is proposed in this paper. This is basing on lazy transfer scheme with VM optimization and object versioning that takes care of multi-snapshotting and multi-deployment simultaneously and effectively. The experiments have shown that the new filing system and related techniques have improved performance, and bandwidth utilization is reduced by 90

Index terms— Cloud Design, Cloud Storage Performance, Empirical Study, Multi-snapshotting, versioning, VM images, lazy propagation, cloning, multi-deployment.

1 Introduction

owadays, the emergence of Infrastructure as a Service (IaaS) cloud computing is a feasible substitute to the acquisition as well as physical resources management. With the help of IaaS, users can be able to lease storage and time of computation from datacenters that are very large. Leasing of computation time can be achieved by enabling users to deploy virtual machines (VMs) on the resources of the datacenter. As the user possess overall control on the configuration regarding Virtual Machines by making use of on-demand deployments, IaaS leasing is simply similar to purchase of hardware that is dedicated but with no long-term commitment as well as cost. The IaaS on-demand nature is complex to make such kind of leases more attractive, as it allows users for expanding or shrinking their resources with respect to their needs of computation, by making use of external resources for complementing their local resource base [15].

This emerging model results in new challenges with respect to the design as well as development of systems providing IaaS. One among frequently resulting patterns in the operation of IaaS is the necessity for deploying a huge number of VMs on most of the nodes relative to a datacenter at the same instant of time, starting from a collection of VM images that are stored previously in a fashion that is persistent. For instance, this pattern is occurred when the user needs the deployment of a virtual cluster that is used to execute a distributed application or a group of environments for supporting a workflow. This pattern is referred as multi deployment. Such kind

of large deployment of most of the VMs at a time can take a longer time. This problem is in particular acute for VM images that are used in scientific computing in which image are large in size (from small number of gigabytes up to greater than 10 GB). A conventional deployment contains hundreds or else thousands of such kind of images. Before starting the instances of VM, conventional techniques of deployment [23] broadcast the images to the nodes, a process which could take time ranging from tens of minutes to approximately hours, not taking into account the time for booting the operating system alone. This could make the time of response of the IaaS installation very longer than that is acceptable and remove the ondemand benefits obtained from cloud computing. Once the instances of the Virtual Machines are being run, a same kind of challenge is applied to snapshotting the deployment. Most of the VM images which were changed locally need to be transferred in a concurrent manner for making storage stable with the reason to N 1 (D D D D) B

capture the VM state for using later (for instance in check pointing or online migration to another cluster or cloud). This pattern is referred to as multisnapshotting. The technique of conventional snapshotting works definitely on custom VM image file formats [9] for storage of only incremental differences in a new file which rely on the original VM image similar to backing file. When taking regular snapshots for a huge number of VMs, such kind of approaches form a huge number of files as well as interdependencies among them, that are difficult for managing and that get in the way with the ease-of-use basis behind clouds. Moreover, with emerging datacenter trends as well as tendencies for Along with incurring delays that are significant and raising issues of manageability, these patterns can also form huge network traffic which comes in the way through the execution of applications on resources that are leased and results in greater costs of utilization for the user.

In this paper a virtual file system that is distributed specifically that is optimized for patterns of multideployment as well as multi-snapshotting. As the patterns are considered complementary, they are investigated in conjunction. Our proposal provides a proper balance between performance, storage space, and finally consumption of network traffic, while treating snapshotting in a transparent manner and revealing standalone and even raw image files (understood by many hypervisors) to the outside. The summary of our contributions are as follows: ? We present a flow of design principles which optimize patterns of multideployment as well as multisnapshotting and describe in which manner our design can be integrated with the resources of IaaS (Sections 2 and 3). ? We illustrate how to comprehend these principles of design by building a virtual file system which leverages distributed storage services that are versioning-based. To clear this point, we describe an implementation over BlobSeer, a service related to versioning storage particularly designed for maximum throughput under concurrency [17,24].

Our approach is evaluated in a sequence of experiments each of which is conducted over hundreds of nodes that are provisioned on the Grid’5000 testbed, by making use of synthetic traces as well as real-life applications.

2 II.

3 Related work

Multideployment which depends on complete broadcast-dependent pre-propagation is a commonly utilized technique [28,23,11]. While this technique prevents read contention to the repository, it can incur great overhead in network traffic as well as execution time, as mentioned in Section 5.2. Moreover, as the VM images are completely copied on the compute nodes locally, multisnapshotting will not be feasible: greater amounts of data have been duplicated unnecessarily and can cause transfer delays that are not acceptable, without mentioning huge space of storage and utilization of network traffic.

For alleviating this problem, most of the hypervisors offer support of native copy-on-write by giving definition of formats of custom VM image file [12,20] particularly designed for efficiently storing additional differences. Similar to our approach, this makes base images to be usable in the form of templates that are read-only for multiple logical instances that store modifications per instance. Moreover, deficiency of standardization and also the generation of more number of new files that are interdependent restrict the portability as well as manageability of the snapshots of VM image that result. Another approach that is different in nature for instantiating a huge number of VMs from the identical initial state has been proposed in [13]. The authors present a latest cloud abstraction: VM FORK. Basically this is considered as the equivalent of the fork call on operating systems like UNIX, cloning a VM at every instant into multiple replicas which are running on various hosts. While this is simply equal to CLONE followed by COMMIT in our method, the main concern is on reducing the time as well as traffic of the network for spawning and running, on the fly, new remote instances of VM that share the identical state of a VM that is already running. Local modifications have been assumed to be ephemeral, and no support is provided for storing the state persistently.

A similar one to our approach is Lithium [10], a replication system that is fork-consistent for virtual disks. Lithium supports instantaneous volume creation along with lazy space allocation and creation of writable snapshots instantaneously. Not similar to our approach is the one which is dependent on segment trees, Lithium is dependent on log structuring [22], that can potentially humiliate read performance when increasing the number of successive snapshots for the same image: the log of incremental differences is started growing, making it more costly for reconstructing the image.

Cluster volume managers for virtual disks like Parallel ax [16] allow compute nodes for sharing access to a block device that is single and globally visible, which in a collaborative manner managed for presenting individual

virtual disk images to the Virtual Machines. While this allows frequent snapshotting that is not efficient like our approach, image sharing is intentionally not encouraged so as to remove the requirement for a distributed lock manager that is claimed for dramatically simplifying the design. Most of the storage systems, like Amazon S3 [5] (backed by Dynamo [8]), are particularly designed as highly accessible key-value repositories for infrastructures of cloud. They may be building blocks that are valuable for block level storage volumes [1] which host images of virtual machine; moreover, they have not been optimized for snapshotting. The intention of our approach is to complement existing platforms of [12], configurations have become more and more varied. Custom image formats are not standardized and might be used with particular hypervisors alone that limits the ability for easily migrating VMs among various hypervisors. Hence, multisnapshotting should be handled in a transparent and portable style which hides the interdependencies of additional differences and exposes VM images that are standalone, by greater portability in various hypervisor configurations.

cloud computing, from industry (Amazon Elastic Compute Cloud: EC2 [4]) as well as from academia (Nimbus [2,12,24], Open Nebula [3]). While the particulars for EC2 are not available publicly, it has been widely accepted that all of these platforms depend on many of the techniques mentioned above. Claims for instantiating multiple VMs in -minutes, ? moreover, are not sufficient to meet our objectives of performance; So, our work is believed to be a welcome addition in this circumstance.

4 III.

Description of infrastructure and other components a) About Cloud Infrastructure Clusters are used in building IaaS cloud platforms. They are made up of hardware that makes use of less power and reduces cost per unit and provides high speed [4]. Many machines are interconnected and each machine is attached a disc storage. Virtualization technology is used in order to share physical resources well. The machines are able to run multiple VMs. Many nodes are dedicated for storage that is responsible for persistence. They might be having either distributed [5] or centralized [2] storage service. Such storage service is responsible to store VM instance images reliably. The manipulations of VMs include deleting, downloading, uploading and so on.

5 b) State of the Application

VM deployment state has two parts namely the state of all VM instances at any given point of time and the state of the channels between them meant for communications. They include sockets which have been open, network state and virtual topology. In order to make the state persistent for future reuse and maintenance, it is essential that the VM instances are to be persisted and at the same time hundreds of VM instances are to be created to meet increasing demands of cloud users. However capturing the global state of such channels is difficult [14]. To avoid this problem, the second model is to get sum of all VM instances. This model discards any in-transit traffic in the network and assumes that fault tolerant network is used.

Model 3, which is simplified version of model 2 is that the VM state is represented only by the virtual disk attached to it. It stores only minimal information pertaining to state and such information is reused later. It has the benefits like reduction in size and portability across systems. Model 3 is widely used mechanism in practice and the same is considered in this work.

6 c) Application Phases

Any VM may not access the whole image. Some utilities and applications are never used. To model this behavior the VM life cycle has been divided into three phases namely boot phase, application phase and shutdown phase. The boot phase reads configuration files, launches processes that represent initial state of VM. The application phase is in either negligible virtual disk access that need not be persisted or data-intensive which needs dedicated storage. The shutdown phase generates very negligible disk access and this phase is not there when VM instance terminates prematurely due to some hardware failure.

7 IV. Our methodology and architecture

In order to optimize the process of multisnapshotting and multideployment, a new filing system is proposed. The following sub sections describe it.

8 a) Overview of Design

The design of the proposed approach depends on the principles like optimizing multisnapshotting, reducing contention, optimizing VM disk access, and aggregating storage space.

i.

9 Aggregating local Storage

The existing approaches [5,2,3] are not capable of making use of storage space available in local hard discs of nodes. To overcome this shortcoming, the proposed approach aggregates storage space from local hard disks and forms a common pool which is used in a distributed fashion. Its advantages are high scalability and freeing memory for reducing overhead in managing VMs.

159 ii.

160 10 Optimizing VM access and Reducing Contention

161 On demand VM image mirroring facilitates to make use of locally available VM image for output. However, it
162 can get from global VM instance the required information in the form of mirroring. It improves performance.
163 Moreover our approach supports reduction of contention as the VM image is split into number of equal sized
164 pieces. While reading values if any piece is not available in the local disk, it is obtained from remote disk thus
165 reducing contention.

166 iii.

167 11 Optimizing Multi-Deployment and Snapshotting

168 When full VM image is saved every time, it consumes lot of resources even though small changes are made.
169 To avoid this certain file formats can be used to incrementally save to other virtual machine. Its drawbacks
170 include limitation of migration capabilities and also the risk of ending up with so many VM instances. By using
171 shadowing and cloning these problems are overcome by the proposed approach.

172 12 Proposed cloud architecture

173 Fig. 1 shows the architecture of the proposed system. It has cloud middleware, compute nodes or hypervisors,
174 clients mirroring modules. The cloud middleware facilitates communication to mirroring modules and also
175 hypervisor concurrently. COMMIT is used to save changes permanently while CLONE is used to make another
176 copy. Local disks are involved to form a distributed file system which improves the overall performance of
177 multisnapshotting.

178 13 VI.

179 14 Implimentation details

180 The proposed system implementation mainly has two modules namely distributed versioning storage service and
181 mirroring module. The former is meant for improving management of repository while the latter for trapping IO
182 access and runs in each compute node. a) Software Reused Some of the components are reused in the proposed
183 system. For instance BlobSheer [17,18,19] and FUSE are reused. The BlobSheer is meant for working with LOB
184 objects while the FUSE is meant for implementing mirroring module.

185 As can be seen in figure 2, the fuse module is made up of many components like hypervisor, cloud middleware,
186 BlobSheer etc.

187 15 Evaluations

188 Experiments and results on multi-deployment and multi-snapshotting are described in the following sub sections.
189 a) Emperical Setup Grid'5000 was used to perform experiments. iNancy with 120 clusters was used. Each one
190 is with x86 64 CPU with virtualization support, local HDD worth 250 GB and 8GB of RAM with Internet
191 connection. KVM 0.12.5 was the hypervisor and the OS is Red Hot Linux.

192 16 b) Multideployment Performance

193 The following sub sections throw light into the experimental results. The observations are done in a
194 multideployment pattern when a single VM is used to have -n? number of VM instances. Going Back and
195 Forth: Efficient Multi-Deployment and Multi-Snapshotting on Clouds i.

196 17 Propagation

197 As given in [21,23], it is part of cloud and has phases like broadcasting of VM image, and launching of VM
198 instances concurrently. The drawback in the propagation approach is the overhead incurred in the initialization
199 phase. Taktuk [7] has been used to overcome this downside. Taktuk is a broadcasting tool which is highly
200 scalable. NFS server is used to store VM images.

201 ii.

202 Comparing Qcow2 Over PVFS PVFS [6] is used to compare our work. This tool is meant for metadata
203 management with high performance. For comparison it was deployed in compute nodes. In order to initialize
204 VM instances qcow2 [9] images are created in the compute node in the local system while PVFS is used as backup
205 image. The performance is measured on average time take to boot each instance and total network traffic.

206 Figure ??, 7, 8 and 9 shows the results of comparison of other works and our approach.

207 18 Multi-Snapshotting Performance

208 The performance of our approach in case of multisnapshotting is described in this section. The comparison is
209 made between qcow2 over PVFS and our approach. Fig. 10 and 11 show the performance of multi-snapshotting

of our approach and qcow2 over PVFS. When overall performance is considered, our approach is taking relatively less time for instance creation and completion. Our approach is showing better results and it is intended to help cloud platforms such as EC2, Nimbus etc. We believe that our work can be used in any existing cloud platform in order to improve its performance in terms of managing virtual machines and improving performance by using our techniques pertaining to multi-snapshotting and multi-deployment. Figure12 (c) shows time taken to finish simulation using 100 VM instances VIII.

19 Conclusion and future work

Since cloud computing is becoming more popular and efficient management of VM images, like image propagation for computing nodes and image snapshotting for the purpose of check-pointing or migration is difficult. The performance of these kind of operations affects in a direct manner the usability of the benefits provided by systems of cloud computing. This paper presented various techniques which integrate with middleware of the cloud for handling two patterns efficiently. They are multideployment and multisnapshotting.

A lazy VM deployment scheme which fetches content of the VM image as required by the application that is executed in the VM, thereby minimizing the pressure on the storage service of VM for deployment requests that are heavily concurrent. Moreover, we leverage object versioning for saving local VM image differences alone back to persistent storage when a snap-shot is generated, yet offer the illusion that the snapshot is a different, completely independent image. This has two crucial benefits. First, it does the management of updates of the hypervisor in an independent manner, thus greatly enhancing the portability of VM images and providing compensation for the deficiency of standardization of the VM image format. Second, it manages snapshotting in a transparent manner at the level of the repository of the VM image, simplifying to a great extent the snapshots management. We have given the demonstration of the advantages of our approach via experiments on number of nodes by making use of benchmarks as well as applications of real-life. When compared with simpler approaches depending on pre-propagation, our approach gives a best improvement in execution time as well as resource usage: the total time for performing a multi-deployment got reduced approximately to a factor of 25, and the storage and bandwidth usage got reduced by approximately 90%. When compared with approaches which make use of copy-on-write images (i.e., qcow2) depending on raw backing images that are stored in a distributed file system (i.e., PVFS), a speedup of multideployment by a factor of 2 and multisnapshotting performance that is comparable are shown, each with the extra benefits of transparency as well as portability.

Depending on these results that are supported, we plan for exploring the multi-deployment as well as multi-snapshotting patterns in a more extensive manner. According to multideployment, one optimization that is possible is to build a scheme that is perfecting depending on last experience through the access pattern. According to multi-snapshotting, reductions that are interesting in time as well as storage space can be achieved by presenting deduplication schemes. We also intend for fully integrating the present work with Nimbus [2] and thereby explore its advantages for more critical applications of HPC in the real world.

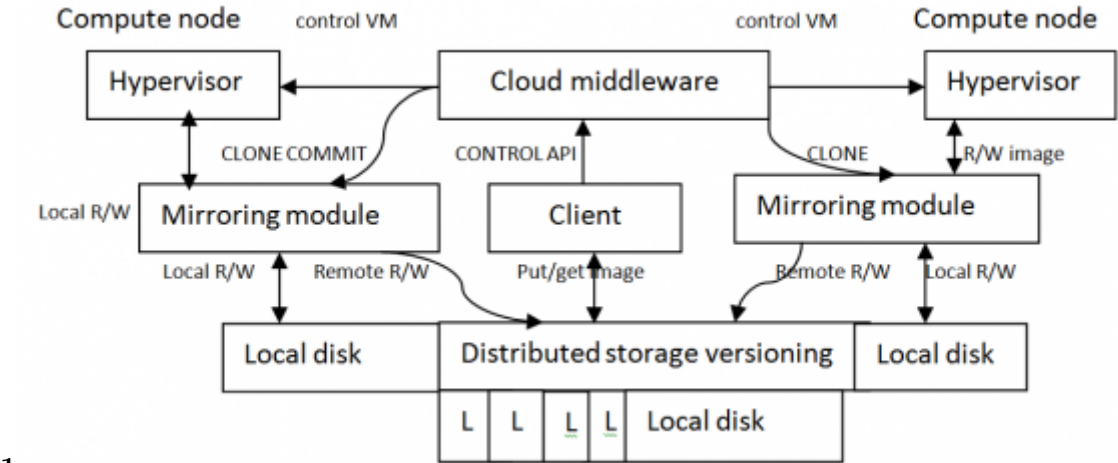
¹July

²© 2012 Global Journals Inc. (US)Global Journal of Computer Science and Technology

³July



Figure 1: ©



1

Figure 2: Fig. 1 :

2

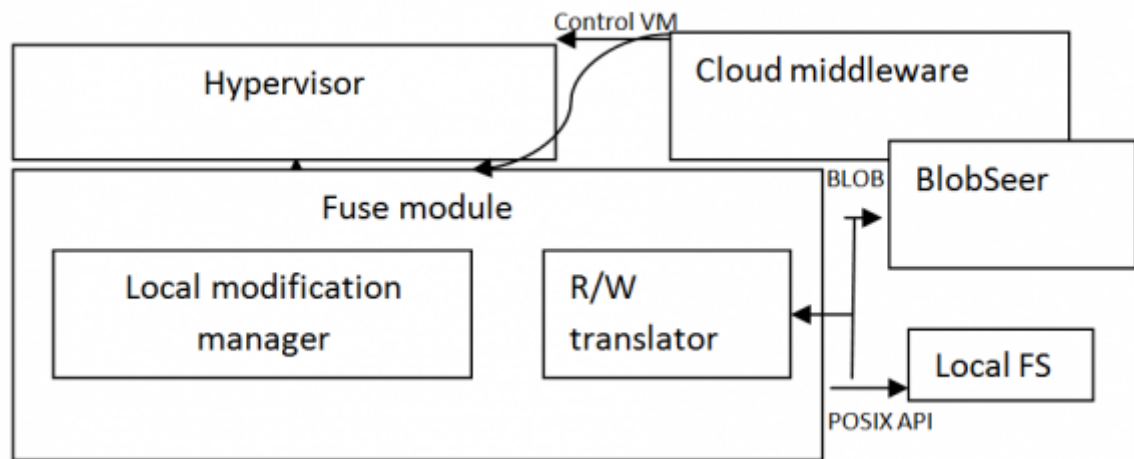


Figure 3: Fig. 2 :

35

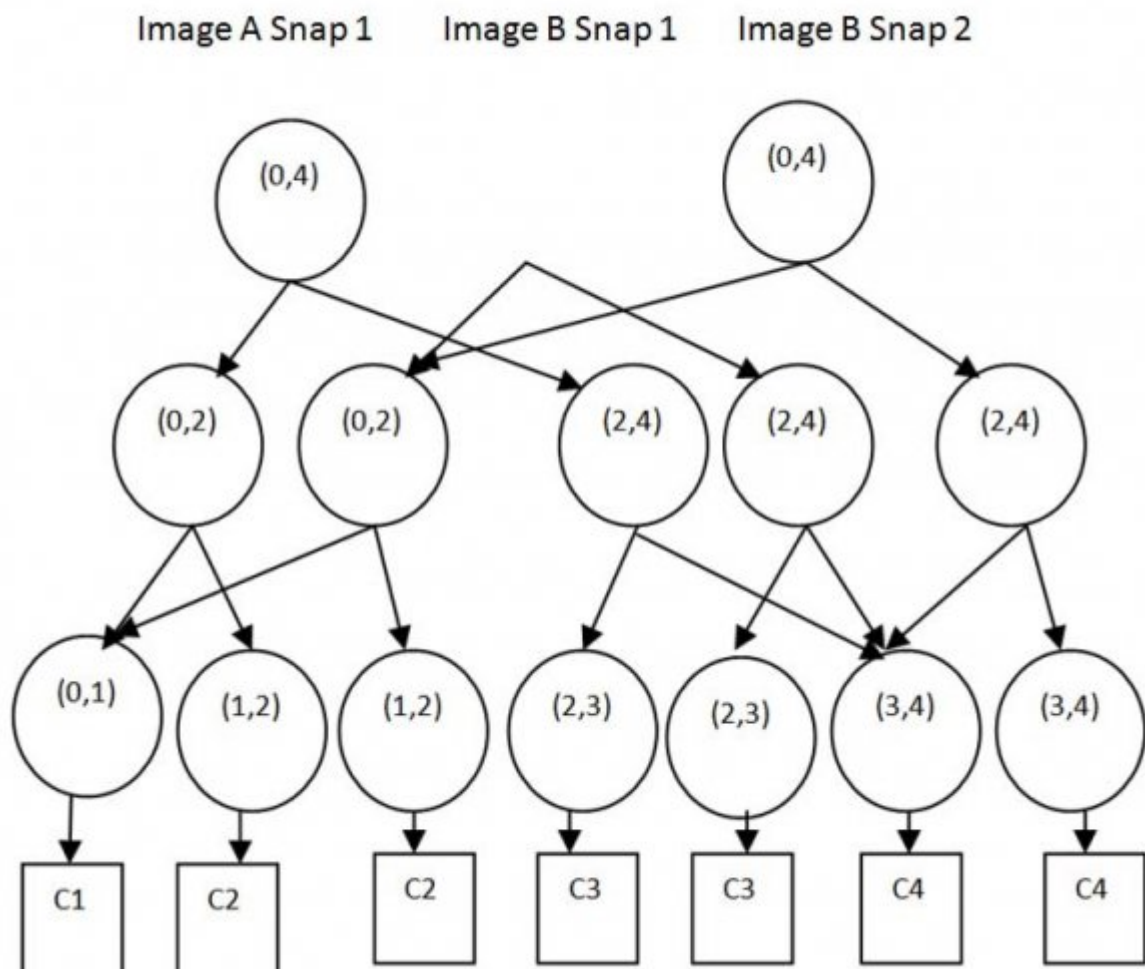


Figure 4: Fig. 3 :Fig. 5 :

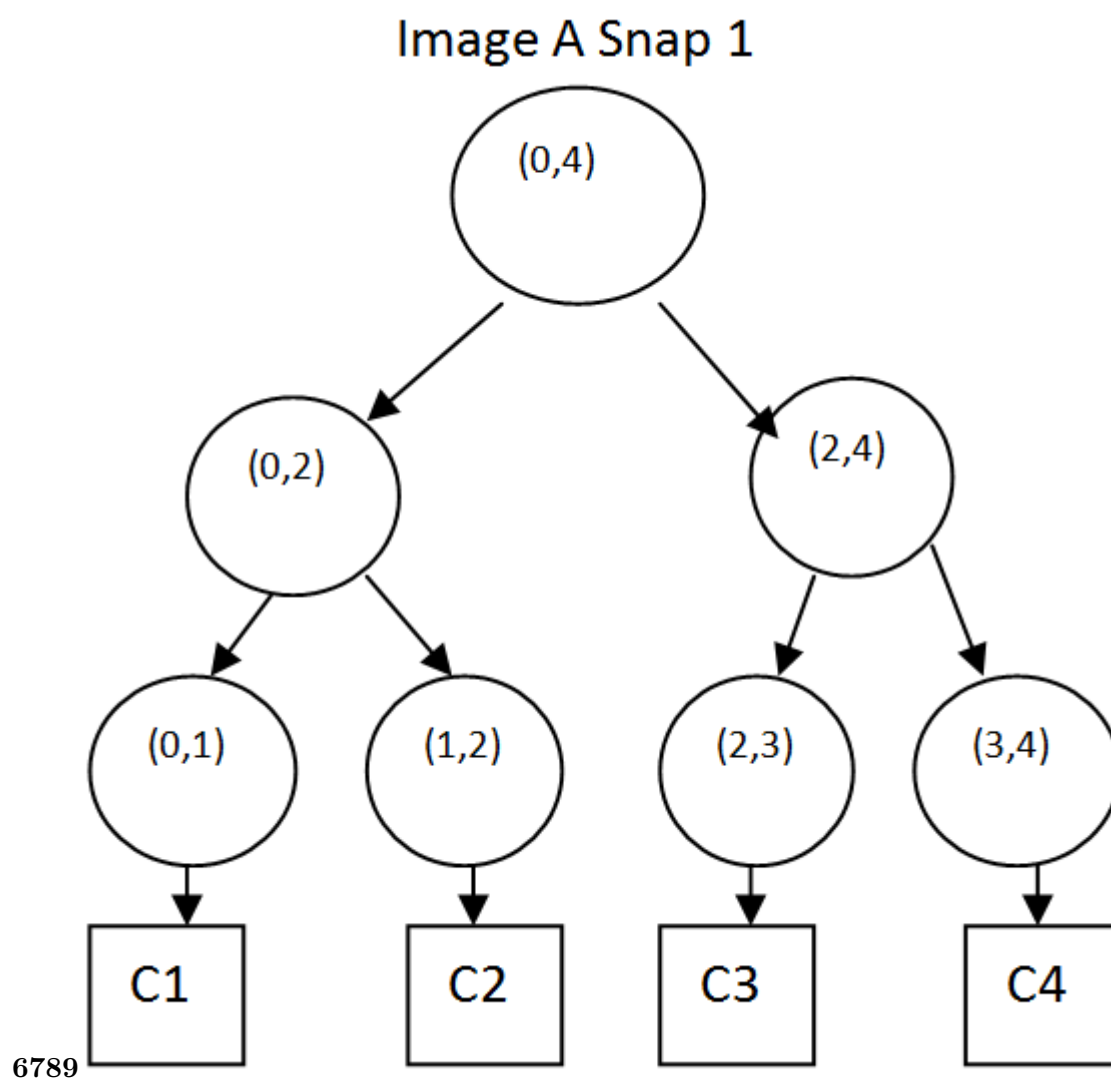


Figure 5: Fig. 6 :Fig. 7 :Fig. 8 :Fig. 9 :

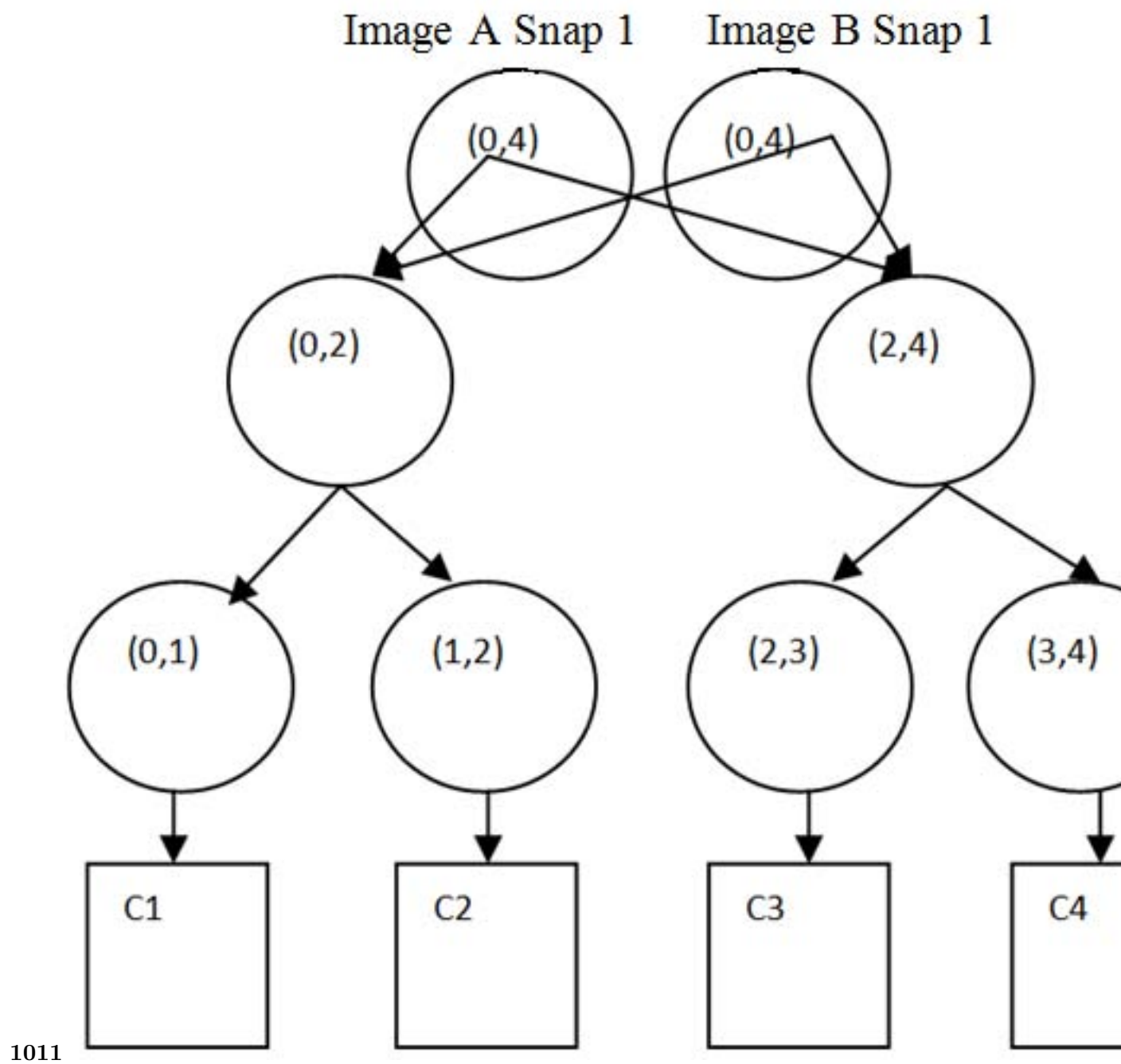


Figure 6: Fig. 10 :Fig. 11 :

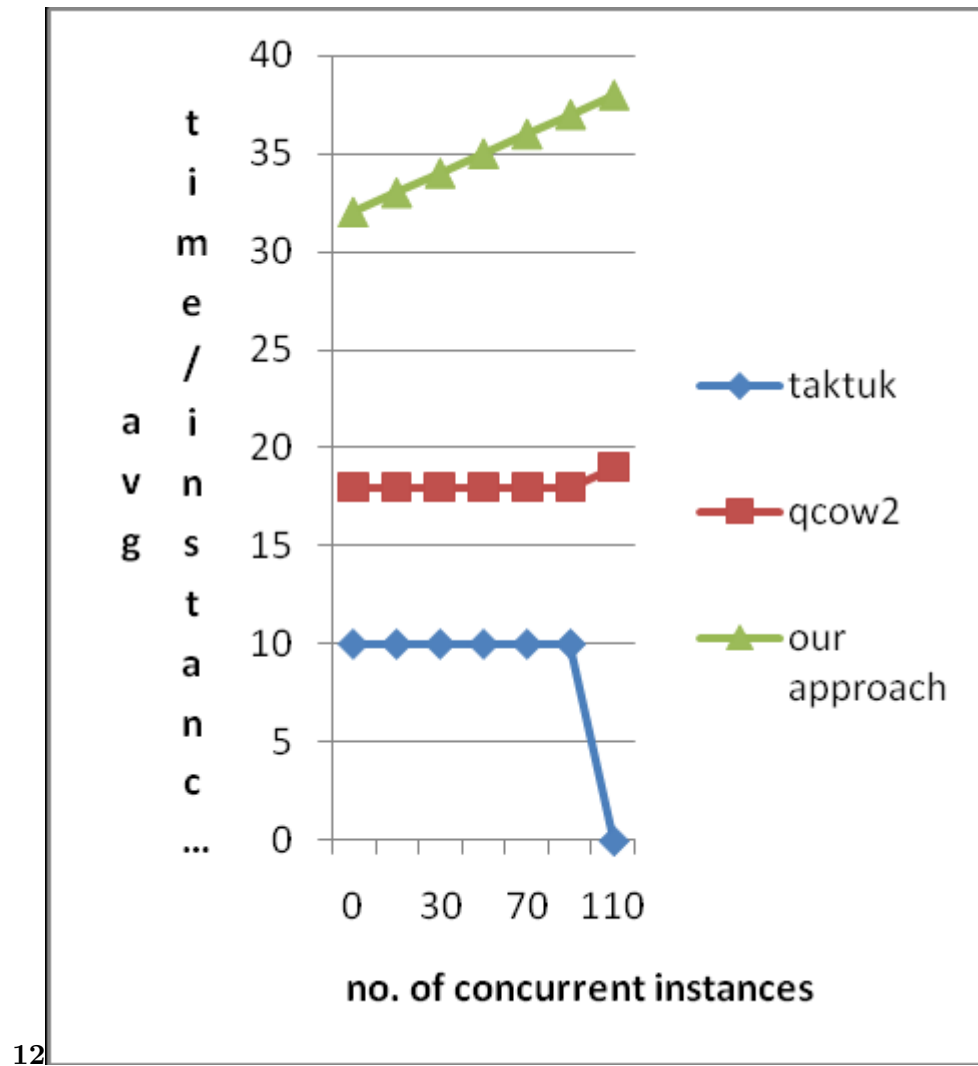


Figure 7: Fig. 12 :

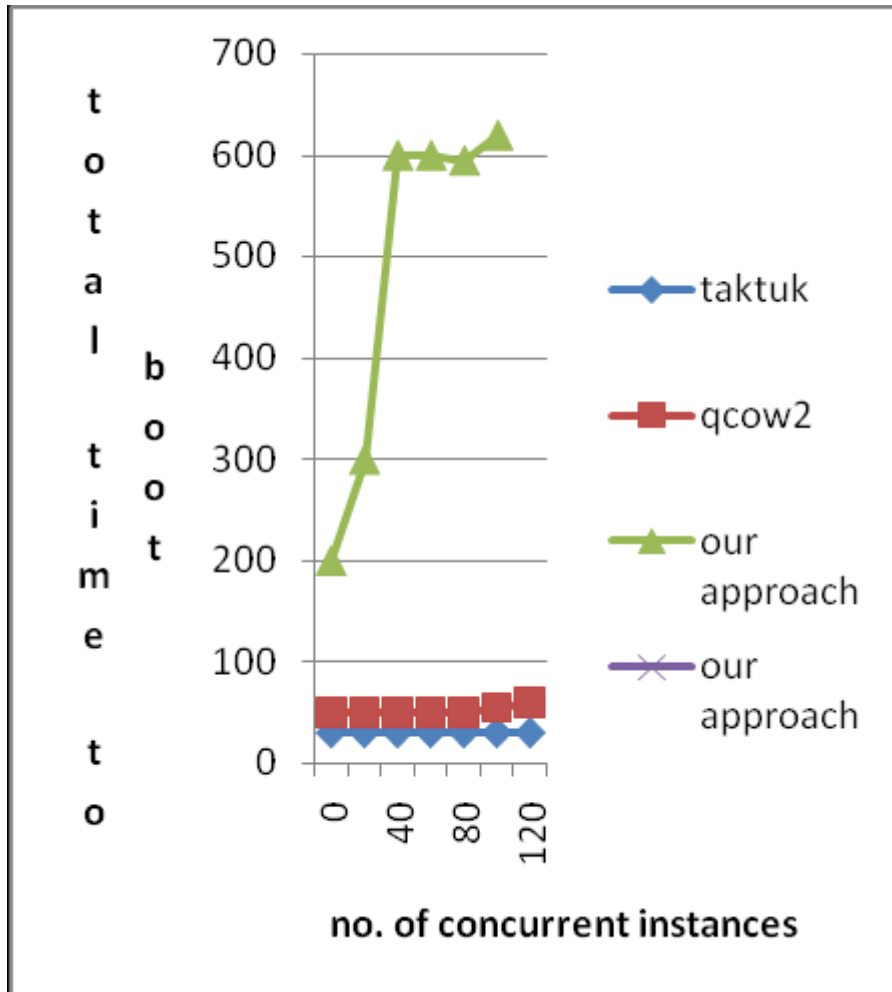


Figure 8:

244 [Nimbus] , Nimbus . <http://www.nimbusproject.org/>

245 [Amazon Simple Storage Service] , *Amazon Simple Storage Service (S3)* .

246 [Amazon Elastic Compute Cloud (EC2) Opennebula] ‘Amazon Elastic Compute Cloud (EC2)’. <http://www.opennebula.org/.4 Opennebula>,

247

248 [Nicolae et al. ()] ‘Blobseer: Bringing high throughput under heavy concurrency to Hadoop map/reduce

249 applications’. B Nicolae , D Moise , G Antoniu , L Boug’e , M Dorier . *IPDPS ’10: Proceedings of the*

250 *24th IEEE International Parallel and Distributed Processing Symposium*, (Atlanta, GA) 2010. p. .

251 [Nicolae et al. (2011)] ‘BlobSeer: Next-generation data management for large scale infrastructures’. B Nicolae ,

252 G Antoniu , L Boug’e , D Moise , A Carpen-Amarie . *J. Parallel Distrib. Comput* February 2011. 71 p. .

253 [Nicolae (2010)] *BlobSeer: Towards Efficient Data Storage Management for Large-Scale, Distributed Systems*, B

254 Nicolae . November 2010. University of Rennes 1 (PhD thesis)

255 [Gagn’e ()] ‘Cooking with Linux-still searching for the ultimate Linux distro?’. M Gagn’e . *Linux J* 2007. 2007.

256 (161) p. 9.

257 [Decandia et al. ()] ‘Dynamo: Amazon’s highly available key-value store’. G Decandia , D Hastorun , M Jampani

258 , G Kakulapati , A Lakshman , A Pilchin , S Sivasubramanian , P Voshall , W Vogels . *Proceedings of 21st*

259 *ACM SIGOPS Symposium on Operating Systems Principles*, (21st ACM SIGOPS Symposium on Operating

260 *Systems Principles* New York) 2007. ACM. p. .

261 [Marshall et al. ()] ‘Elastic site: Using clouds to elastically extend site resources’. P Marshall , K Keahey , T

262 Freeman . *CCGRID ’10: Proceedings of the 10th IEEE/ACM International Conference on Cluster, Cloud*

263 *and Grid Computing, CCGRID ’10*, (Washington,DC, USA) 2010. IEEE Computer Society. p. .

264 [Hibler et al. ()] ‘Fast, scalable disk imaging with Frisbee’. M Hibler , L Stoller , J Lepreau , R Ricci , C Barb .

265 *ATC ’03: Proceedings of the 2003 USENIX Annual Technical Conference*, (San Antonio, TX) 2003. p. .

266 [Keahey et al. ()] ‘Fortes. Sky computing’. K Keahey , M O Tsubawa , A M Matsunaga , JA B . *IEEE Internet*

267 *Computing* 2009. 13 (5) p. .

268 [Wartel et al. ()] ‘Image distribution mechanisms in large scale cloud providers’. R Wartel , T Cass , B Moreira

269 , E Roche , M Guijarro , S Goasguen , U Schwickerath . *CloudCom ’10: Proceedings 2nd International*

270 *Conference on Cloud Computing Technology and Science*, (Indianapolis, IN) 2010.

271 [Liu et al. ()] ‘Network state consistency of virtual machine in live migration’. X Liu , J Huai , Q Li , T Wo .

272 *SAC ’10: Proceedings of the 2010 ACM Symposium on Applied Computing*, (New York) 2010. ACM. p. .

273 [Reimer et al.] ‘Opening black boxes: Using semantic information to combat virtual machine image sprawl’. D

274 Reimer , A Thomas , G Ammons , T Mummert , B Alpern , V Bala . *VEE ’08: Proceedings of the 4th ACM*

275 *SIGPLAN/SIGOPS International Conference on*,

276 [Meyer et al. ()] ‘Parallax: Virtual disks for virtual machines’. D T Meyer , G Aggarwal , B Cully , G Lefebvre

277 , M J Feeley , N C Hutchinson , A Warfield . *SIGOPS Oper. Syst. Rev* 2008. 42 (4) p. .

278 [Carns et al. ()] ‘Pvfs: A parallel file system for Linux clusters’. P H Carns , W B Ligon , R B Ross , R Thakur

279 . *Proceedings of the 4th Annual Linux Showcase and Conference*, (the 4th Annual Linux Showcase and

280 *Conference* Atlanta, GA) 2000. USENIX Association. p. .

281 [Rodriguez et al. ()] ‘Resource selection for fast large-scale virtual appliances propagation’. A Rodriguez , J

282 Carretero , B Bergua , F Garcia . *ISCC ’09: Proceedings of 14th IEEE Symposium on Computers and*

283 *Communications*, 2009. p. .

284 [Hansen (2010)] ‘Scalable virtual machine storage using local disks’. J G Hansen , E . *SIGOPS Oper. Syst. Rev*

285 December 2010. 44 p. .

286 [Keahey and Freeman ()] ‘Science clouds: Early experiences in cloud computing for scientific applications’. K

287 Keahey , T Freeman . *CCA ’08: Proceedings of the 1st Conference on Cloud Computing and it’s Applications*,

288 2008.

289 [Lagar-Cavilla et al. ()] ‘SnowFlock: Rapid virtual machine cloning for cloud computing’. H A Lagar-Cavilla , J

290 A Whitney , A M Scannell , P Patchin , S M Rumble , E Lara , M Brudno , M Satyanarayanan . *EuroSys*

291 *’09: Proceedings of the 4th ACM European Conference on Computer Systems*, (New York) 2009. ACM. p. .

292 [Claudel et al. ()] ‘Taktuk, adaptive deployment of remote executions’. B Claudel , G Huard , O Richard . *HPDC*

293 *’09: Proceedings of the 18th ACM International Symposium on High Performance Distributed Computing*,

294 (New York) 2009. ACM. p. .

295 [Rosenblum and Ousterhout ()] ‘The design and implementation of a log-structured file system’. M Rosenblum ,

296 J K Ousterhout . *ACM Trans. Comput. Syst* 1992. 10 (1) p. .