



DynamicCloudSim: Simulating heterogeneity in computational clouds



Marc Bux*, Ulf Leser

Department of Computer Science, Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany

HIGHLIGHTS

- We present DynamicCloudSim, a simulator that models instability in cloud computing.
- We give an overview of studies on instability in compute clouds like Amazon EC2.
- We evaluate how well different scientific workflow schedulers handle instability.
- We compare the simulated execution of a workflow against actual execution on EC2.

ARTICLE INFO

Article history:

Received 15 February 2014

Received in revised form

8 August 2014

Accepted 19 September 2014

Available online 5 October 2014

Keywords:

Cloud computing

Simulation

Heterogeneity

Scientific workflows

Scheduling

ABSTRACT

Simulation has become a commonly employed first step in evaluating novel approaches towards resource allocation and task scheduling on distributed architectures. However, existing simulators fall short in their modeling of the instability common to shared computational infrastructure, such as public clouds. In this work, we present DynamicCloudSim which extends the popular simulation toolkit CloudSim with several factors of instability, including inhomogeneity and dynamic changes of performance at runtime as well as failures during task execution. As a validation of the introduced functionality, we simulate the impact of instability on scientific workflow scheduling by assessing and comparing the performance of four schedulers in the course of several experiments both in simulation and on real cloud infrastructure. Results indicate that our model seems to adequately capture the most important aspects of cloud performance instability. The source code of DynamicCloudSim and the examined schedulers is available at <https://code.google.com/p/dynamiccloudsim/>.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

Over the last decade, cloud computing emerged as a form of distributed computing, in which computational resources can be provisioned on-demand over the Internet [1]. In the Infrastructure-as-a-Service (IaaS) model of cloud computing, computational resources of any scale can be rented in the form of virtual machines (VMs) from commercial cloud providers like Amazon or Microsoft [2]. The convenience of its pay-as-you-go model along with the aggressive promotion by its providers has led to an exponential growth in the usage of cloud computing over the last years (see Fig. 1).

Tailoring highly scalable applications to make efficient use of cloud resources requires developers to be aware of both the performance of rented cloud infrastructure and the requirements of the to-be-deployed application. These characteristics are hard to

quantify and vary depending on the application and cloud provider. Since benchmarking a given application on cloud infrastructure of large scale repeatedly under various experimental conditions is both tedious and expensive, simulation constitutes a convenient and affordable way of evaluation prior to implementation and execution on real hardware [3–5].

Unfortunately, available cloud simulation toolkits like CloudSim [6] do not adequately capture inhomogeneity and dynamic performance changes inherent to non-uniform and shared infrastructures like computational clouds. The effect of these factors of uncertainty and instability is not negligible and has been repeatedly observed to strongly influence the runtime of a given application on commercial clouds such as Amazon's Elastic Compute Cloud (EC2) (e.g., [7–13]).

In this work, we present DynamicCloudSim, an extension to CloudSim which provides an array of capabilities to model the instability inherent to computational clouds and similar distributed infrastructures. We evaluate the applicability of DynamicCloudSim in a series of experiments involving the scheduling of two computationally intensive scientific workflows, one of which has been repeatedly used for evaluation purposes. We believe the field of

* Corresponding author. Tel.: +49 177 6263294.

E-mail addresses: buxmarcn@informatik.hu-berlin.de (M. Bux), leser@informatik.hu-berlin.de (U. Leser).

<http://dx.doi.org/10.1016/j.future.2014.09.007>

0167-739X/© 2014 Elsevier B.V. All rights reserved.

scientific workflow scheduling to be suitable for evaluating the capabilities of a cloud simulation framework for two reasons: (1) Scheduling computationally intensive scientific workflows on distributed and potentially shared architectures presents many opportunities for optimizing robustness to instability [14]; (2) Simulation has been repeatedly made use of for evaluating scientific workflow schedulers (e.g., [15,16]).

Scientific workflows are often represented as directed, acyclic graphs (DAGs), in which nodes correspond to data processing tasks and edges constitute data dependencies between these tasks. They have recently gained attention as a flexible programming paradigm for modeling, representing, and executing complex computations and analysis pipelines in many different areas of scientific research [17].

A variety of sophisticated algorithms for scheduling scientific workflows on distributed computational infrastructures have been developed (e.g., [18–20]). As an application example for DynamicCloudSim, we compare the performance of several established scientific workflow schedulers at different levels of instability. Since some of the investigated schedulers have been developed to handle heterogeneity, dynamic performance changes, and failure, we expect our experiments to replicate the advertised strengths of the different workflow schedulers. Results from an extensive number of simulation runs confirm these expectations, underlining the importance of elaborate scheduling mechanisms when executing workflows on shared computational infrastructure, in which resources are shared between multiple users and which are thus subject to resource contention and performance variation (e.g., shared clusters, grids, or clouds).

The remaining part of this paper is structured in the following way: The CloudSim framework is described in Section 2, whereas the features it has been extended with are described in Section 3. The setup of the scientific workflow scheduling experiments evaluating the extensions introduced by DynamicCloudSim is outlined in Section 4. The results of these experiments are presented and discussed in Section 5. Related work is summarized in Section 6. Finally, an outlook to future work and conclusions are given in Sections 7 and 8.

2. CloudSim

CloudSim is an extension of the GridSim [21] framework for simulation of resource provisioning and scheduling algorithms on cloud computing infrastructure developed by Calheiros et al. [6] at the University of Melbourne's CLOUDS Laboratory. It provides capabilities to perform simulations of assigning and executing a given workload on a cloud computing infrastructure under different experimental conditions. CloudSim for instance has been used to (1) measure the effects of a power-aware VM provisioning and migration algorithm on datacenter operating costs for real-time cloud applications [3], (2) evaluate a cost-minimizing algorithm of VM allocation for cloud service providers, which takes into account a fluctuating user base and heterogeneity of cloud VMs [4], (3) develop and showcase a scheduling mechanism for assigning tasks of different categories – yet without data dependencies – to the available VMs [5].

CloudSim operates event-based, i.e., all components of the simulation maintain a message queue and generate messages, which they pass along to other entities. A CloudSim simulation can instantiate several *datacenters*, each of which is comprised of *storage* servers and physical *host* machines, which in turn host multiple VMs executing several *tasks* (named *cloudlets* in CloudSim). For a detailed overview, refer to Fig. 2. A datacenter is characterized by its policy of assigning requested VMs to host machines (the default strategy being to always choose the host

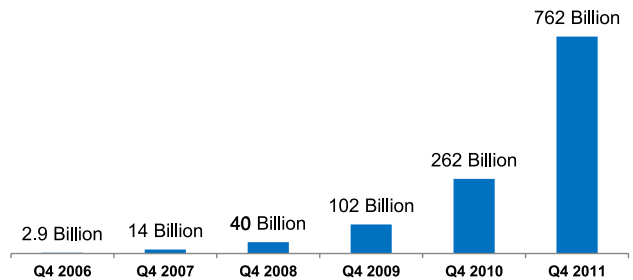


Fig. 1. Total number of objects stored in the Amazon Simple Storage Service (Amazon S3) since its introduction in 2006. The exponential growth of the largest cloud provider's data storage solution mirrors the trend of compute clouds increasing in size and popularity. Data has been published on the Amazon Web Services Blog in April 2012.¹

with the least cores in use). Each datacenter can be configured to charge different costs for storage, VM usage, and data transfer.

The computational requirements and capabilities of hosts, VMs, and tasks are captured in four performance measures: MIPS (million instructions per second per core), bandwidth, memory, and local file storage. Furthermore, each host has its own policy which defines how its computational resources are to be distributed among allocated VMs, i.e., whether VMs operate on shared or distinctly separated resources and whether over-subscription of resources is allowed. Similarly, each VM comes with a scheduling policy specifying how its resources are to be distributed between tasks. On top of this architecture, an application-specific datacenter broker supervises the simulation, requesting the (de-)allocation of VMs from the datacenter and assigning tasks to VMs.

One of the key aspects of CloudSim is that it is easily extensible. Several extensions have been presented, including (1) NetworkCloudSim [22], which introduces sophisticated network modeling and inter-task communication, (2) EMUSIM [23], which uses emulation to determine the performance requirements and runtime characteristics of an application and feeds this information to CloudSim for more accurate simulation, or (3) CloudMIG [24], which facilitates the migration of software systems to the cloud by contrasting different cloud deployment options based on the simulation of a code model in CloudSim.

3. DynamicCloudSim

CloudSim assumes provisioned virtual machines to be predictable and stable in their performance: Hosts and VMs are configured with a fixed amount of MIPS and bandwidth and VMs are assigned to the host with the most available MIPS. On actual cloud infrastructure like Amazon EC2, these assumptions do not hold. While most IaaS cloud vendors guarantee a certain processor clock speed, memory capacity, and local storage for each provisioned VM, the actual performance of a given VM is subject to the underlying physical hardware as well as the usage of shared resources by other VMs assigned to the same host machine. In this section, we outline the extensions we have made to the CloudSim core framework as well as the rationale behind them.

3.1. File I/O

In CloudSim, the amount of time required to execute a given task on a VM depends solely on the task's length (in MI) and the VM's processing power (in MIPS). Additionally, the external bandwidth (in KB/s) of VMs and their host machines can be specified, but neither have an impact on the runtime of a task. However, many data-intensive tasks are neither computational- nor communication-intensive, but primarily I/O-bound. Especially

¹ <http://tinyurl.com/6uh8n24>

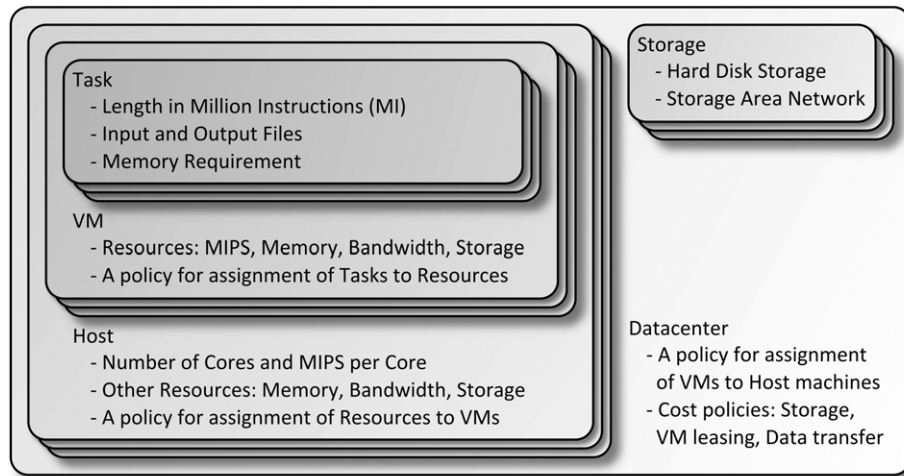


Fig. 2. The architecture of the CloudSim framework.

in database applications, a substantial amount of tasks involves reading or writing large amounts of data to local or network storage [7].

DynamicCloudSim introduces external bandwidth as a requirement of tasks and file I/O as an additional performance characteristic of tasks, VMs and hosts. It takes into account all performance requirements of a task when determining how long it takes to execute the task on a given VM. Hence, DynamicCloudSim allows the simulation of executing different kinds of tasks (CPU-, I/O-, bandwidth-bound) on VMs with different performance characteristics.

3.2. The need for introducing instability

In a performance analysis spanning multiple Amazon EC2 datacenters, Dejun et al. [7] observed occasional severe performance drops in virtual machines, which would cause the response time of running tasks to greatly increase. Furthermore, they reported the response time of CPU- and I/O-intensive application to vary by a factor of up to four on VMs of equal configuration. Notably, they detected no significant correlation between CPU and I/O performance of VMs. Zaharia et al. [8] found the I/O throughput of “small”-sized VM instances in EC2 to vary between roughly 25 and 60 MB/s, depending on the amount of co-located VMs running I/O-heavy tasks.

In a similar evaluation of Amazon EC2, Jackson et al. [9] detected different physical CPUs underlying similar VMs: Intel Xeon E5430 2.66 GHz, AMD Opteron 270 2 GHz, and AMD Opteron 2218 HE 2.6 GHz. They also observed network bandwidth and latency to depend on the physical hardware of the provisioned VMs. When executing a communication-intensive task on 50 VMs, the overall communication time varied between 3 and 5 h over seven runs, depending on the network architecture underlying the provisioned VMs. In the course of their experiments, they also had to restart about one out of ten runs due to the occurrence of failures. Similar observations have been made in other studies on the performance of cloud infrastructure [10,11].

Another comprehensive analysis of the performance variability in Amazon EC2 was conducted by Schad et al. [12]. Once per hour over a time period of two months they measured the CPU, I/O, and network performance of newly provisioned VMs in Amazon EC2 using established microbenchmarks. Performance was found to vary considerably and generally fall into two bands, depending on whether the VM would run on Intel Xeon or AMD Opteron infrastructure (see Fig. 3). The variance in performance of individual VMs was also shown to strongly influence the runtime of a

real-world MapReduce application on a virtual cluster consisting of 50 EC2 VMs. A further interesting observation of this study was that the performance of a VM depends on the hour of the day and day of the week. Iosup et al. [13] made similar observations when analyzing more than 250,000 real-world performance traces of commercial clouds.

Evidently, the performance of computational cloud infrastructure is subject to different factors of instability:

1. Heterogeneous physical hardware underlying the provisioned VMs (*Het*).
2. Dynamic changes of performance at runtime (*DCR*).
3. Straggler VMs and failed task executions (*SaF*).

In the remaining part of this section, we describe in detail how DynamicCloudSim attempts to capture these factors of instability.

3.3. Heterogeneity

Similar to Amazon EC2, the provisioning of resources to virtual machines in DynamicCloudSim is based on compute units instead of fixed performance measures. Different host machines provide a different amount of computing power per provisioned compute unit, effectuating in heterogeneity (*Het*) among VM performance. Furthermore, in contrast to CloudSim, DynamicCloudSim does not assign new VMs to the host with the most available resources, but to a random machine within the datacenter. Hence, VMs of equal configuration are likely to be assigned to different types of physical machines providing varying amounts of computational resources. Similar to a commercial cloud like Amazon EC2, the user is oblivious to the hardware underlying the provisioned VMs and has no control over the VM allocation policy. In Section 4, we describe how we set up a simulation environment resembling the inhomogeneous hardware configuration of Amazon EC2.

In addition to the heterogeneity achieved by assigning VMs to different types of hosts, DynamicCloudSim also provides the functionality to randomize the individual performance of a VM. If this feature is enabled, the performance characteristics (CPU, I/O, and bandwidth) of a newly allocated VM are sampled from a normal distribution instead of using the default values defined by the VM's host machine. The mean of this normal distribution is set to the host's default value of the performance characteristic and the relative standard deviation (RSD) can be defined by the user, depending on the desired level of heterogeneity. Schad et al. [12] found several of the performance measurements of VMs in Amazon EC2 – particularly random disk I/O and network bandwidth – to

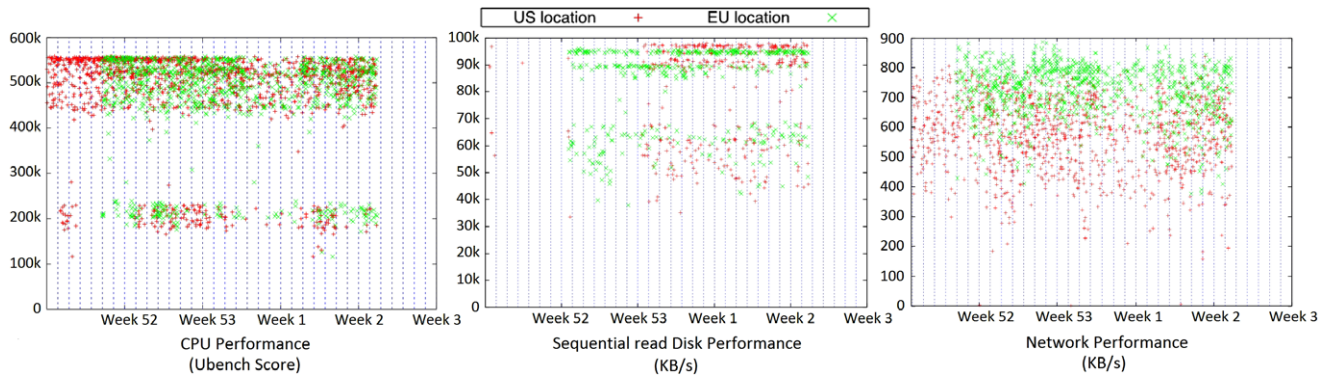


Fig. 3. CPU, sequential read disk, and network performance on Amazon EC2 large instances, measured over the course of several weeks. Each dot represents the performance of a single VM at a certain point in time.

Source: Images taken from [12] with friendly permission.

be normally distributed. Hence, while DynamicCloudSim supports sampling from nine different kinds of distributions, we have selected the normal distribution as the default distribution to sample from.

Dejun et al. [7] reported average runtimes between 200 and 900 ms – with a mean of roughly 500 and a standard deviation of about 200 – for executing a CPU-intensive task on 30 VMs across six Amazon EC2 datacenters. Based on these measurements, we set the default value for the RSD parameter responsible for CPU performance heterogeneity to 0.4. In the same way, we determined a default value for I/O heterogeneity of 0.15. Based on similar measurements taken by Jackson et al. [9] for communication-intensive tasks on Amazon EC2, we set the default value for network bandwidth heterogeneity to 0.2. These values are backed up by the performance measurements of Schad et al. [12], who observed an RSD of 0.35 between processor types in EC2 as well as RSD values of 0.2 and 0.19 for disk I/O and network performance.

3.4. Dynamic changes at runtime

So far we have only modeled heterogeneity, which represents permanent variance in performance of VMs due to differences in underlying hardware. Another important concept of instability inherent to cloud computing are dynamic changes of performance at runtime (DCR) due to external loads as a consequence of sharing common resources with other VMs and users, as for instance reported by Dejun et al. [7]. DynamicCloudSim attempts to capture two aspects of this concept: (1) Long-term changes in a VM's performance due to a certain event, e.g., the co-allocation of a different VM with high resource utilization on the same host. (2) Uncertainty or noise, which models short-term alterations in a VM's performance. While DynamicCloudSim does not model unforeseeable external factors affecting VM performance explicitly, it implicitly models their effects on performance.

To simulate the effects of long-term changes, DynamicCloudSim samples from an exponential distribution with a given rate parameter to determine the time of the next performance change. The exponential distribution is frequently used to model the time between state changes in continuous processes. In DynamicCloudSim, the rate parameter is defined by the user and corresponds to the average number of performance changes per hour. In light of the observations made by Schad et al. [12] and Iosup et al. [13], who found that the performance of a VM can change on an hourly basis (but does not necessarily do so), we assume the performance of a VM to change about once every other hour by default. Since this parameter is highly dependent on the particular computational infrastructure, we encourage users to adjust it to their respective environment.

Whenever a change of one of the performance characteristics has been induced on a VM, the new value for the given characteristic is by default sampled from a normal distribution, though DynamicCloudSim also supports the use of other distributions. The mean of this normal distribution is set to the baseline value of the given characteristic for this VM, i.e., the value that has been assigned to the VM at allocation time. The RSD of the distribution is once again set by the user. Higher values in both the rate parameter of the exponential distribution and the standard deviation of the normal distribution correspond to higher levels of dynamics.

Uncertainty (or noise) is the lowest tier of dynamic performance changes in DynamicCloudSim. It is modeled by introducing slight aberrations to a VM's performance whenever a task is assigned to it. As with heterogeneity and dynamics, this is by default achieved by sampling from a normal distribution with user-defined RSD parameter.

On Amazon EC2, Dejun et al. [7] observed relative standard deviations in performance between 0.019 and 0.068 for CPU-intensive tasks and between 0.001 and 0.711 for I/O-intensive tasks. We set the default values for the RSD parameter of long-term performance changes to the third quartile of these distributions, i.e., to 0.054 for CPU performance and 0.033 for I/O performance. Similarly, we set the default RSD value for the noise parameter to the first quartile, i.e., to 0.028 for CPU and 0.007 for I/O.

3.5. Stragglers and failures

In massively parallel applications on distributed computational infrastructure, fault-tolerant design becomes increasingly important [25]. For the purpose of simulating fault-tolerant approaches to scheduling, DynamicCloudSim introduces straggler VMs and failures (SaF) during task execution. Stragglers are virtual machines exhibiting constantly poor performance [8]. In DynamicCloudSim, the probability of a VM being a straggler can be specified by the user along with the coefficient that determines how much the performance of a straggler is diminished.

We propose default values of 0.015, respectively 0.5 for the straggler likelihood and performance coefficient parameters. These values are based on the findings of Zaharia et al. [8], who encountered three stragglers with performance diminished by 50% or more among 200 provisioned VMs in their experiments. The numbers are backed up by the observations of Zhang et al. [26], who encountered an amount of stragglers below 5% in an experiment in which VMs were allocated on a 160 node cluster using the large-scale VM provisioning toolkit VMThunder. The effect of these parameters is exemplarily shown in Fig. 4, which illustrates all of the introduced factors of instability (Het, DCR, SaF) in combination for the CPU performance of eight VMs, including one straggler, in an experiment of 12 h in DynamicCloudSim.

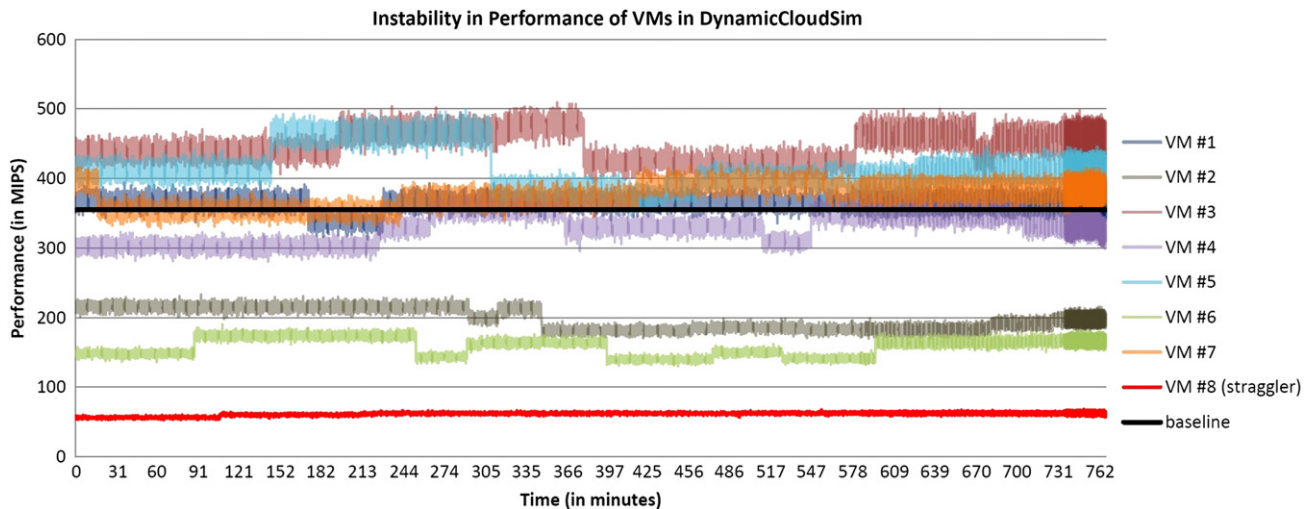


Fig. 4. CPU performance of a configuration of eight VMs (including one straggler) running for 12 h with default parameters in DynamicCloudSim. The black line represents how the VMs' performance would have looked like in basic CloudSim.

Failures during task execution are another factor of instability commonly encountered in distributed computing. DynamicCloudSim is currently confined to a basic method of failure generation: Whenever a task is assigned to a VM and its execution time is computed, DynamicCloudSim determines whether the task is bound to succeed or fail. This decision is based on the average rate of failure specified by the user. The default value for the rate of failed task executions is set to 0.002, based on the observations of Jackson et al. [9], who, in a series of experiments running on 50 VMs, had to restart every tenth run on average due to the occurrence of a failure on at least one VM. This parameter setting is reinforced by the SLA of Amazon,² which guarantees a machine uptime of 99.95%. We argue that a default task success rate of four times lower than the machine uptime guaranteed by Amazon is reasonable, since a failure during task execution can occur for different reasons, e.g., failure to retrieve data over the network, failure during machine startup, or failure due to intermittent machine hang.

There are various reasons for failed task execution, such as temporary performance breakdowns within a VM or the inability to access input data or write output data. Usually, such perturbations are not immediately recognized, hence resulting in severely increased runtimes. Consequently, in DynamicCloudSim the runtime of a failed task execution is determined by multiplying the task's execution time with a user-defined coefficient. The introduction of more sophisticated failure models on different levels (VM, storage, task) of workflow execution is left for future work (see Section 8).

4. Experimental validation

To evaluate DynamicCloudSim's adequacy in modeling instability encountered in computational clouds, we performed two experiments. In the first experiment, we simulate the execution of two workflows, the Montage workflow from astronomy and a genomic sequencing workflow from bioinformatics, using different mechanisms of scheduling and different levels of instability in the computational infrastructure. We expect the schedulers to differ in their robustness to instability, which should be reflected in diverging workflow execution times.

In the second experiment, we executed the Montage workflow on different collections of virtual machines in Amazon EC2 using

the same four workflow schedulers. We expect the observed behavior of the four workflow schedulers to be similar in simulation and on real cloud infrastructure. Furthermore, we expect the workflow execution times on EC2 to be comparable to the execution times of simulated runs in DynamicCloudSim with default parameters. Execution of the genomic sequencing workflow on EC2 was omitted since the workflow processes sensitive genomic data.

As the workflow execution engine for Amazon EC2, we used the Hi-WAY ApplicationMaster³ and Apache Hadoop 2.2.0 (YARN). Apache Hadoop⁴ has been developed as an open-source implementation of Google's MapReduce programming model [27] and distributed file system. Hadoop's newest advancement, Hadoop 2.2.0 (YARN), recently introduced the possibility to execute other programming models than MapReduce via custom ApplicationMasters. Hi-WAY has been developed as an ApplicationMaster for YARN that allows arbitrary scientific workflows of different languages, including Pegasus DAX, to be executed on top of any Hadoop installation (a detailed description of the Hi-WAY ApplicationMaster will appear elsewhere). Note that apart from being able to parse workflow execution traces generated by Hi-WAY, DynamicCloudSim is also able to parse synthetic workflows generated by the Pegasus Workflow Generator.⁵

In this section, we outline in detail the evaluation workflows, the schedulers which we used in our experiments, and the settings for both experiments.

4.1. Evaluation workflows

Over the course of the experiments outlined in this section, we made use of two distinct evaluation workflows: (1) the Montage workflow from the field of astronomy and (2) a genomic sequencing workflow from the field of bioinformatics. The two workflows were selected since they are from different scientific domains and have been repeatedly utilized for evaluating aspects of scientific workflow execution in the past. Also, they provide the desirable property of allowing to be scaled to nearly any desired size and degree of parallelism, by adjusting the amount of to-be-processed input data.

³ <https://github.com/marcbux/Hi-WAY>.

⁴ <http://hadoop.apache.org/>.

⁵ <https://confluence.pegasus.isi.edu/display/pegasus/WorkflowGenerator>.

² <http://aws.amazon.com/ec2/sla/>.

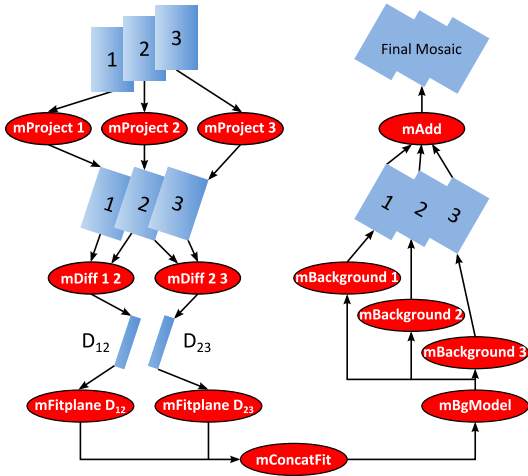


Fig. 5. A schematic instance of the Montage workflow, in which three input images are processed to generate a mosaic [28]. First, input images are projected to a common spatial scale (mProject). To rectify images to a common background level (mBackground), background radiation present in the images has to be captured and modeled. To this end, pairs of overlapping images are subtracted from one another (mDiff). The resulting difference images are then put on a plane (mFitplane) and concatenated (mConcatFit) before undergoing a modeling of common background radiation (mBgModel). Finally, the projected, background-corrected images are merged into the output mosaic (mAdd).

4.1.1. The Montage workflow

We constructed an evaluation workflow using the Montage toolkit [28]. Montage is able to generate workflows in Pegasus DAX format [29] for assembling high-resolution mosaics of regions of the sky from raw input data. It has been repeatedly utilized for evaluating scheduling mechanisms or computational infrastructures for scientific workflow execution (e.g., [30–33]). See Fig. 5 for a schematic visualization of the Montage workflow and Fig. 6 for an example of output generated by a Montage workflow.

In the first experiment, we used a Montage workflow which builds a large-scale (twelve square degree) mosaic of the m17 region of the sky. This workflow consists of 43,318 tasks reading and writing 534 GB of data in total, of which 10 GB are input and output files which have to be uploaded to and downloaded from the computational infrastructure. In the second experiment, we used a smaller Montage workflow that only builds a one square degree mosaic. This workflow exhibits similar runtime characteristics, yet consists only of 387 tasks reading and writing 7.3 GB of data of which 128 MB are input and output files.

4.1.2. The genomic sequencing workflow

As a second evaluation workflow, we used a genomic sequencing workflow outlined in detail by Pabinger et al. [34]. The workflow was implemented using the functional workflow language Cuneiform⁶ and is available online for reference⁷ (a detailed description of the Cuneiform workflow language will appear elsewhere). Similar implementations of this workflow or parts thereof have been provided as components of genome analysis toolkits such as GATK⁸ or ADAM.⁹ Genomic sequencing workflows have been previously proposed and made use of for the evaluation of scientific workflow scheduling [35].



Fig. 6. A one square degree mosaic of the m17 region of the sky. The image has been generated by executing the corresponding Montage workflow.

Genomic sequencing denotes the process that determines the ordered sequence of nucleotides within a given DNA molecule. So far, this can only be achieved by splitting the DNA into many short, overlapping fragments, which are called reads. The workflow in Fig. 7 compares such reads of two colorectal cancer cell lines, Caco-2 and GEO, to find genomic variants specific to one or the other. In the reference alignment step, the genomic reads of both cell lines are mapped onto a much larger reference genome to determine the reads' original position in the genome. The aligned reads are compared against the reference genome to detect variants, which are then compared against each other.

In our implementation of the workflow, we split the reads into distinct files of 5 MB, which we aligned against the human chromosome 22 using three different alignment tools—bowtie [36], SHRiMP [37], and PerM [38]. The resulting alignments were merged using samtools [39] and variants were detected using VarScan [40]. This resulted in an executable workflow comprising 4266 tasks reading and writing 436 GB of data in total.

4.2. Scientific workflow schedulers

Scheduling a scientific workflow denotes the process of mapping the workflow's tasks onto the available computational resources [19]. Most scheduling policies are developed with the aim to minimize overall workflow execution time. However, certain scenarios call for different scheduling objectives, e.g., the optimization of monetary cost or data security. In general, we differentiate between static and adaptive schedulers [14]. In static scheduling, a schedule is assembled prior to execution and then strictly abided by at runtime. Conversely, in adaptive scheduling, scheduling decisions are made on-the-fly at the time of execution. Here we consider (1) a static round robin scheduler, (2) the HEFT scheduling heuristic [18], (3) a greedy task queue, and (4) the LATE algorithm [8]. The schedulers outlined in this section were implemented in both DynamicCloudSim and Hi-WAY.

4.2.1. Static round robin scheduling

The round robin scheduler is often used as a baseline implementation of a static scheduler (e.g., [29]). It constructs a schedule by traversing the workflow from the beginning to the end, assigning tasks to computational resources in turn. This way, each resource

⁶ <https://github.com/joergen7/cuneiform>.

⁷ <https://raw.githubusercontent.com/marcbux/Hi-WAY/master/hiway-core/examples/variant-call.cf>.

⁸ <https://www.broadinstitute.org/gatk/>.

⁹ <https://github.com/bigdatagenomics/adam>.

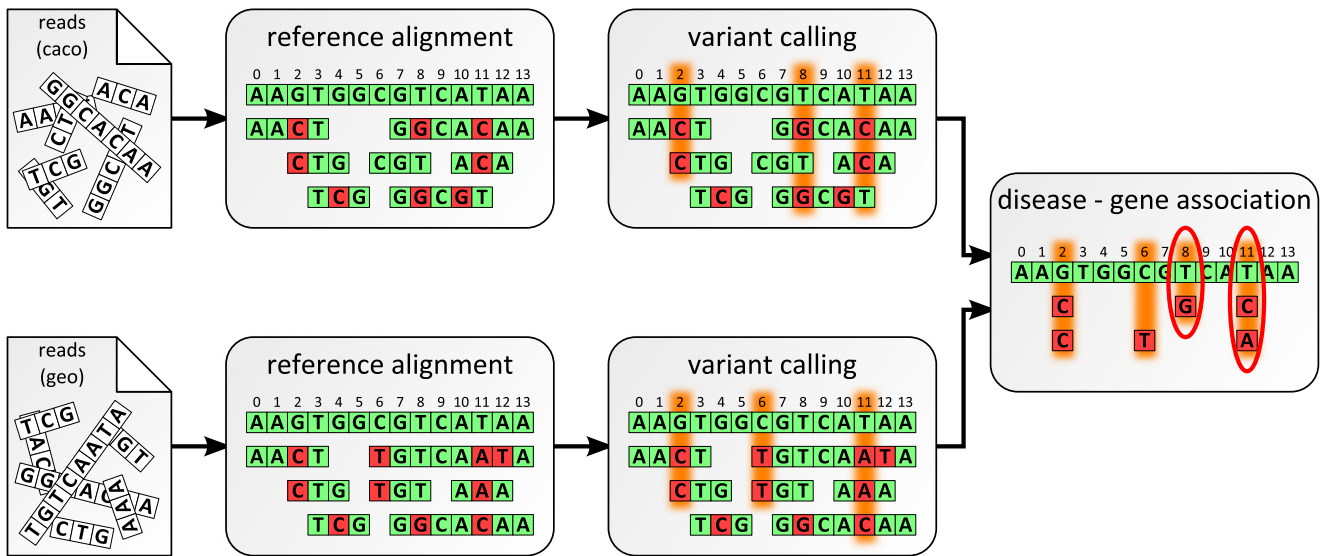


Fig. 7. An abstract representation of a genomic sequencing workflow that compares the genomic sequencing data of two colorectal cancer cell lines. The process of genomic sequencing produces data in the form of short substrings of the sequenced genome, so-called reads. A commonly employed first step in the analysis of these reads – and the first step in this workflow – is to align the reads against a reference genome to determine their original position in the genome. In the next step, the aligned reads are compared against the reference in order to detect variations from the reference. Finally, detected variants of both cancer cell lines are compared to one another to determine which of the detected variants are characteristic for either cell line.

will end up with roughly the same amount of tasks, independent of its computational capabilities or the tasks' workload. The static round robin scheduler is available in the SWfMS Pegasus [29]. We expect the static round robin scheduler to perform well in homogeneous and stable computational infrastructures. Adding heterogeneity (Het), dynamic changes at runtime (DCR) or stragglers and failures (SaF) to the experiment should heavily diminish its performance.

4.2.2. HEFT scheduling

A number of more sophisticated mechanisms than round robin for static workflow scheduling on heterogeneous computational architectures have been proposed (e.g., [20,19,16]). Among the more influential scheduling algorithms is the Heterogeneous Earliest Finishing Time (HEFT) heuristic, which has been developed by Topcuoglu et al. [18]. Similar to most sophisticated static schedulers, HEFT requires runtime estimates for the execution of each task on each computational resource, which are difficult to obtain.

The HEFT heuristic traverses the workflow from the end to the beginning, computing the upward rank of each task as the estimated time to overall workflow completion at the onset of this task. The computation of a given task's upward rank incorporates estimates for both the runtimes and data transfer times of the given task as well as the upward ranks of all successor tasks. The static schedule is then assembled by assigning each task in decreasing order of upward ranks a time slot on a computational resource, such that the task's scheduled finish time is minimized. Similar to static round robin scheduling, the HEFT scheduling heuristic has been integrated in the SWfMS Pegasus [29].

In our simulation experiments in DynamicCloudSim, HEFT is provided with accurate runtime estimates based on the execution time of each task on each CloudSim VM at the time of its allocation. Hence, we expect the HEFT scheduler to perform well in both homogeneous and heterogeneous infrastructures. However, we expect poor performance if dynamic changes (DCR) or failures in the computational infrastructure (SaF) are introduced and runtime estimates become inaccurate. In our experiments on Amazon EC2, we provide HEFT with accurate runtime estimates by executing and measuring the runtime of each task once on each machine prior to workflow execution.

4.2.3. Greedy task queue scheduling

The most intuitive approach to adaptive scheduling is a greedy task queue. Here, tasks are assigned to computational resources in first-come-first-serve manner at runtime. Whenever a resource has an available task slot, it fetches a task from a queue of tasks ready for execution. Task queues have been implemented in a variety of SWfMS, including Taverna [41] and Kepler [42]. The default scheduler of Hadoop [43] also employs a greedy queue. In our experiments, we expect a task queue scheduler to outperform static schedulers when dynamic changes (DCR, SaF) in the computational infrastructure are introduced.

4.2.4. LATE scheduling

The LATE (Longest Approximate Time to End) scheduler developed by Zaharia et al. [8] constitutes a well-established alteration of the default task queue. By speculatively replicating tasks progressing slower than expected, LATE exhibits increased robustness to the effects of straggler resources and failed task execution (SaF). LATE keeps track of the runtime and progress of all running tasks. By default, 10% of the task slots on resources performing above average are assigned speculative copies of tasks which are estimated to finish farthest into the future and have progressed at a rate below average. Intuitively, this approach maximizes the likeliness for a speculative copy of a task to overtake its original. LATE evidently follows a rationale similar to that of HEFT, since both scheduling heuristics prioritize the assignment of tasks with longest times to finish to well-performing computational resources.

LATE was implemented as an extension of Hadoop’s default scheduling algorithm. For Hadoop’s Sort benchmark executed on 800 virtual machines of an Amazon EC2 test cluster, the LATE scheduler has been shown to outperform the default scheduler of Hadoop by 27% on average [8]. In our experiments, we expect LATE to be robust even in settings with straggler resources and high rates of failures during task execution (SaF). However, due to 10% of the computational resources being reserved for speculative task execution, LATE should perform slightly inferior to a greedy queue on homogeneous and stable computational infrastructure.

As outlined above, LATE requires progress estimates of each currently running task. Unfortunately, reliable progress estimates are typically not available when executing a scientific workflow,

Table 1

CFP2006 benchmark results for processors found in Amazon EC2.

Machine	Cores	SPECfp® 2006 base score	Percentage of reference	URL
Intel Xeon E7-4870 2.4 GHz	10	51.0	100%	http://tinyurl.com/d3oghak
Intel Xeon E5430 2.66 GHz	8	18.1	35.5%	http://tinyurl.com/bckaqqw
AMD Opteron 2218 2.6 GHz	4	12.6	24.7%	http://tinyurl.com/ajqj3n3
AMD Opteron 270 2.0 GHz	4	8.89	17.4%	http://tinyurl.com/aug9xcq

in which tasks are black boxes. To account for a possible lack of progress estimates in a real-world scenario, DynamicCloudSim provides parameters to distort LATE's progress estimates. Since the tasks comprising the evaluation workflows do not provide progress estimates, we implemented a naive random cloning strategy in Hi-WAY, in which tasks are selected for speculative replication at random whenever there is an excess of task slots during workflow execution. We make use of these components in the second experiment when comparing simulated execution in DynamicCloudSim to actual execution on Amazon EC2.

4.3. Experimental settings—workflow scheduling at different levels of instability

In the first experiment, both the 43,318 task Montage workflow and the 4266 task genomic sequencing workflow were executed on a single core of a Dell PowerEdge R910 with four Intel Xeon E7-4870 processors (2.4 GHz, 10 cores) and 1 TB memory, which served as the reference machine of our experiments. Network file transfer, local disk I/O and the runtime of each task in user-mode were captured and written to a trace file.

We parsed the workflows and the trace they generated on the Xeon E7-4870 machine in DynamicCloudSim. The tasks were assigned performance requirements according to the trace file, i.e., a CPU workload corresponding to the execution time in milliseconds, an I/O workload equal to the file sizes of the task's input and output files, and a network workload according to the external data transfer caused by the task. When executing the workflows in CloudSim, all data dependencies were monitored. Thus, a task could not commence until all of its predecessor tasks had finished execution.

In our simulations, we attempt to mirror the computational environment of Amazon EC2. Hence, we obtained SPECfp® 2006 benchmark results for Intel Xeon E5430 2.66 GHz, AMD Opteron 270 2.0 GHz, and AMD Opteron 2218 HE 2.6 GHz, which Jackson et al. [9] detected in their evaluation of Amazon EC2 as underlying hardware. SPECfp® 2006 is the floating point component of the SPEC® CPU2006 benchmark suite. It provides a measure of how fast a single-threaded task with many floating point operations is completed on one CPU core. An overview of the benchmark results is displayed in Table 1.

A CloudSim datacenter was initialized with 500 host machines: 100 Xeon E5430, 200 Opteron 2218, and 200 Opteron 270. Since the Xeon E5430 has twice as many cores as the AMD machines, each type of machine contributes to the datacenter with an equal amount of cores and thus compute units. The CPU performance of each core of these machines was set to the ratio of the machine's SPECfp® 2006 score to the reference machine's score. For instance, the CPU performance of Xeon E5430 machines was set to 355, effectuating in a runtime of 28,169 ms for a task that took 10,000 ms on the Xeon E7-4870 reference machine.

We assume all of the data associated with the workflows – input, intermediate, and output – to be saved on shared storage such as Amazon S3. Different measurements of network throughput within Amazon EC2 and S3 ranging from 10 to 60 MB/s have been reported (e.g., [44,9,45]). We therefore set the default I/O throughput of virtual machines to 20 MB/s. The external bandwidth of virtual machines was set to 0.25 MB/s, based on the remote access

performance of S3 reported by Iamnitchi et al. [11] and Pelletineas et al. [45].

In the course of the experiments, we incrementally raised the level of instability in DynamicCloudSim. For each of the two evaluation workflows, we conducted four experiment runs, in which we measured the effect of heterogeneity (Het), dynamic performance changes at runtime (DCR), and straggler VMs and faulty task executions (SaF):

1. **Het:** We measure the effect of heterogeneous computational infrastructure on different approaches to workflow scheduling. To this end, the relative standard deviation (RSD) parameters responsible for inhomogeneity are incrementally set to 0, 0.125, 0.25, 0.375, and 0.5 (for CPU, I/O, and network performance). The simulation of dynamic performance changes at runtime (DCR) as well as straggler VMs and failed tasks (SaF) is omitted.
2. **DCR:** We examine how dynamic changes in the computational infrastructure affect workflow scheduling. Therefore, the RSD parameters responsible for long-term changes in the performance of a VM was varied between 0, 0.125, 0.25, 0.375, and 0.5. At the same time, the rate of performance changes is fixed at 0.5 and the RSD parameters for noise are set to 0.025 across all runs.
3. **SaF:** We determine the effect of straggler resources and failures during task execution. For this reason, the likelihoods of a VM being a straggler and of a task to fail are set to 0, 0.00625, 0.0125, 0.01875, and 0.025. The performance coefficient of straggler resources is set to 0.1 and the factor by which the runtime of a task increases in the case of a failure is set to 20.
4. **Extreme parameters:** In this setting, we utilize 1.5 times the maximum values for heterogeneity, dynamics and stragglers/failures from experiments 2, 3, and 4 (i.e., RSD parameters of 0.75 for Het and DCR; straggler and failure likelihoods of 0.0375 for SaF) to determine the effect of combining all introduced factors of instability at a very high level.

4.4. Experimental settings—validation on Amazon EC2

To verify that DynamicCloudSim adequately models the instability encountered on real cloud infrastructure, we performed a second run of experiments, in which we contrasted the execution of a Montage workflow in simulation to actual execution on Amazon EC2. Following the methods described in Section 4.3, we generated a workflow trace of the 387 task Montage workflow outlined in Section 4.1.1. Using this trace, we then simulated the execution of the workflow 20 times for each scheduler in DynamicCloudSim on eight virtual machines. In this experiment, we used the datacenter configuration described in Section 4.3 as well as DynamicCloudSim's default parameters as determined and presented in Sections 3.3–3.5.

We then executed the same workflow repeatedly on Amazon EC2 cloud infrastructure. The workflow runs were performed in different virtual compute clusters spread evenly across the EC2 datacenters of US East (Virginia), US West (California), and Europe (Ireland) and across different times of the day. Each cluster consisted of nine general purpose instances of type m1.small with 1 EC2 compute unit, 1.7 GB RAM and 8 GB elastic block storage each. The operating system on the cloud machines was a 64-bit Ubuntu Server 12.04.3 LTS.

Hadoop 2.2.0, Hi-WAY, and the Montage toolkit were installed on each of the rented EC2 instances. In each virtual cluster comprising nine instances, one instance served as the central master for the resource management and job tracking, while the other eight instances served as worker nodes. Across the provisioned virtual clusters, the workflow was executed 20 times for each scheduler and execution times were measured. To allow for a comparison to the measurements in DynamicCloudSim, the overhead introduced by the execution engine (initialization of a workflow run, communication between master and slaves, etc.) was deducted from runtime measurements.

5. Results and discussion

The experimental results are divided into two parts. As described in Section 4.3, we first present the results of scheduling two computationally intensive scientific workflows in DynamicCloudSim. In this experiment, we expect to replicate the known strengths and shortcomings of the scheduling mechanisms in question. As outlined in Section 4.4, we then compare the simulated execution of a Montage workflow in DynamicCloudSim against actual runs on Amazon EC2. By comparing the observations of simulated runs in DynamicCloudSim against workflow execution on actual hardware, we intend to showcase that DynamicCloudSim is able to adequately model the behavior of computational clouds.

5.1. Workflow scheduling at different levels of instability

For each configuration, both the 43,318 task Montage workflow described in Section 4.1.1 and the 4266 task genomic sequencing workflow outlined in Section 4.1.2 were executed 100 times on eight virtual machines and the average runtime was determined. The results of the experiments outlined in Section 4.3 are displayed in Figs. 8 and 9. Over the course of the entire experiments, average runtimes of between 296 and 13,195 min for Montage and between 143 and 1990 min for genomic sequencing have been observed. Evidently, the instability parameters provided by DynamicCloudSim, particularly the Het and SaF parameters, can have a considerable impact on execution time, especially for static schedulers.

In the experiment simulating the effect of heterogeneous resources (Het), all schedulers except the static round robin scheduler exhibit robustness to even the highest levels of variance (see Fig. 8(a)). The reasons for this observation are that HEFT has been designed specifically with inhomogeneous computational resources in mind and queue-based schedulers like the greedy scheduler and LATE are able to adapt to the computational infrastructure. All three schedulers effectively assign less tasks to slower resources. Conversely, the static round robin scheduler is oblivious to the computational infrastructure and simply assigns an equal amount of tasks to each resource, which results in faster resources having to idly wait for slower resources to finish.

Since LATE always reserves 10% of the available resources for speculative scheduling, we would expect runtimes slightly below a greedy queue, which we did not observe. The reason for this might be that HEFT and LATE have a slight edge over the greedy queue-based scheduler: If there is a computationally intensive task which blocks the execution of all successor tasks, of which there is one in Montage (mBgModel) and of which there are several in the sequencing workflow, HEFT and LATE are able to assign it to a well-suited computational resource, instead of simply assigning it to the first available resource. HEFT does this by consulting the accurate runtime estimates of all task–resource-assignments it has been provided with. LATE simply starts a speculative copy of the task on a compute node performing above average.

Finding only the static round robin scheduler to perform subpar in this experimental setting confirmed our expectations outlined in Section 4.2. Evidently, DynamicCloudSim is able to simulate the effect of inhomogeneous resources. Since heterogeneity is commonly encountered in distributed architectures like computational clouds, this is a very desirable property which will continue to be important going forward and has not been sufficiently supported by other cloud simulation toolkits.

In the second part of the experiment, we examined how dynamic changes in the performance of VMs (DCR) affect the runtime of the evaluation workflows achieved by the four scheduling mechanisms (see Fig. 8(b)). The results confirm our expectations of static schedulers like static round robin and HEFT not being able to handle dynamic changes. The major shortcoming of static schedulers lies in the fact that they assemble a schedule prior to workflow execution, which is then strictly abided by. Therefore, changes in the runtime environment make even elaborate static schedules suboptimal.

HEFT provided with accurate runtime estimates constitutes one of the most sophisticated static scheduling policies available, since it takes into account the suitability of a given task for each resource. The only scenarios, in which HEFT should perform substantially worse than a greedy task queue, should be ones in which all tasks have equal performance requirements (which is not the case in Montage) or the runtime estimates are inaccurate, e.g., due to alterations at runtime. Hence, the findings of the second experiment are a strong indicator of DynamicCloudSim being able to simulate dynamic changes in the performance of resources.

In the third part of the experiment, we measured how the appearance of straggler VMs and failed task executions (SaF) influence the performance of the four examined workflow schedulers. Fig. 8(c) confirms the robustness of the LATE scheduler even for high amounts of failures and stragglers. In contrast, the performance of all other schedulers diminished quickly in the face of failure. This is not surprising, since if critical tasks are assigned to straggler VMs or encounter a failure during execution, overall workflow execution time can increase substantially. Speculative replication of tasks with a low progress rate alleviates this problem.

As mentioned previously, the genomic sequencing workflow exhibits more blocking tasks than the Montage workflow. Such blocking tasks are particularly problematic when assigned to a straggler VM. For this reason, the performance degradation for high SaF values is even more apparent in the genomic sequencing workflow than in the Montage workflow. Notably, for very high DCR values, VM performance baseline changes can occasionally reach the extent of the VM temporarily behaving like a straggler. For workflows with frequent blocking tasks, this can lead to reduced performance of schedulers that can usually cope with DCR, as seen when executing the genomic sequencing workflow with the Greedy Queue scheduler at $DCR = 0.5$.

In the fourth part of the experiment, we examined how all three of the introduced factors of instability combined and taken to extremely high levels (RSD parameters of 0.75 for Het and DCR; likelihood parameters of 0.0375 for SaF) influence the workflow execution time. The results of this experiment for Montage are shown in Fig. 9. Once again, LATE is the only scheduler to exhibit robustness to even extreme parameter configurations. Furthermore and in contrast to the findings in the third experiment, the HEFT scheduler substantially outperforms the greedy job queue.

The combination of all factors, i.e., dynamic changes at runtime to inhomogeneous compute resources which can also be stragglers or subject to faulty task execution, can lead to cases, in which the execution of the Montage workflow task mBgModel (as described above) can take extremely long time. This is more problematic for a greedy task queue, which assigns a task to the first available computational resource, which might be a straggler. In contrast,

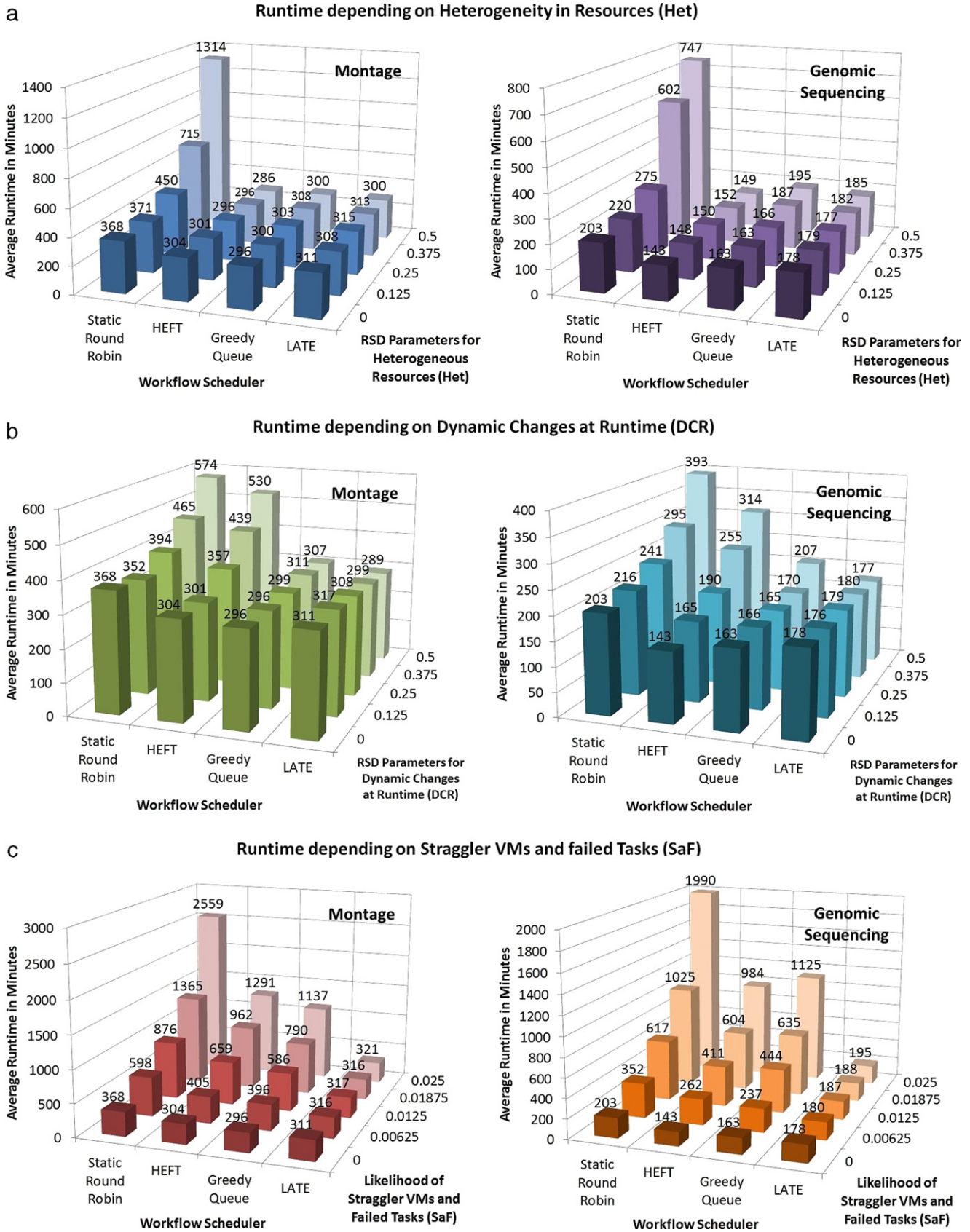


Fig. 8. Effects of (a) heterogeneity (Het), (b) dynamic changes of performance at runtime (DCR), and (c) straggler VMs and failed tasks (SaF) on execution time of the evaluation workflows using different schedulers. Runtimes of the Montage workflow are shown on the left, whereas runtimes of the genomic sequencing workflow can be found on the right.

Table 2

Mean values and standard deviations of Montage workflow execution runtimes in DynamicCloudSim as opposed to on Amazon EC2.

Configuration	DynamicCloudSim		Amazon EC2	
	Mean (min)	STD	Mean (min)	STD
Static round robin	16.166	16.674	8.194	1.166
HEFT	7.899	2.493	7.23	1.585
Greedy queue	8.754	2.423	6.905	0.489
LATE	7.485	1.080		
Distorted LATE/random cloning	7.721	0.898	6.592	0.517

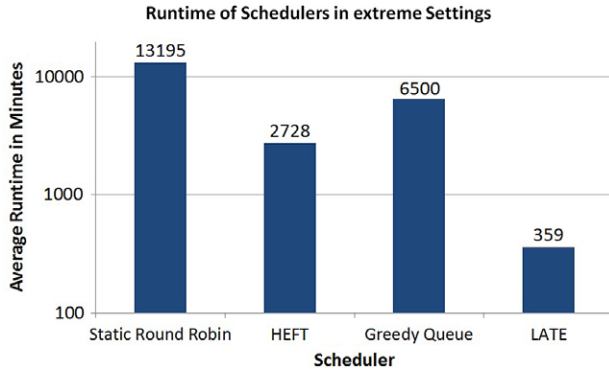


Fig. 9. Execution time (in log scale) of the Montage workflow in DynamicCloudSim in extreme cases of instability.

HEFT is at least able to handle heterogeneous and straggler resources by means of its accurate runtime estimates.

The last two experiments illustrated the severe effect of straggler VMs and failed tasks executions (SaF) on workflow runtime, confirming previous reports on the importance of fault-tolerant design in computationally intensive applications (e.g., [25,33]). While the simulation was able to replicate the advertised strengths of LATE, we acknowledge that more sophisticated failure models would be a desirable enhancement over the current implementation of DynamicCloudSim (see Section 8).

All in all, the experiments clearly confirmed the expectations described in Section 4.2. The simulations underline the importance of adaptive scheduling of scientific workflows in shared and distributed computational infrastructures like public clouds. Moreover, the results serve as an indicator of DynamicCloudSim being well-suited to simulate the dynamics and instability inherent to computational clouds. To further validate this claim, we provide experimental results of executing Montage on actual cloud infrastructure in the next section.

5.2. Validation on Amazon EC2

For each of the available schedulers, the execution of the 387 task Montage workflow was simulated 20 times in DynamicCloudSim and performed 20 times on Amazon EC2. See Table 2 for the mean and standard deviation values and Fig. 10 for box plots of the measured runtimes.

Slightly different mean workflow execution runtimes for the four schedulers were observed between DynamicCloudSim and EC2 (see Table 2). However, using a two-tailed, paired *t*-test, this difference was not found to be significant (*p*-value: 0.187). When compared to the workflow executions on Amazon EC2, a substantially higher variance in workflow runtime was observed in DynamicCloudSim for round robin and greedy queue scheduling. We attribute this finding to the appearance of stragglers and failures in DynamicCloudSim, which we did not encounter to a similar extent during our experiments on Amazon EC2. The lack of stragglers and failures on EC2 also contributes to the higher average runtime of the workflow in DynamicCloudSim when using

round robin scheduling, since this scheduler is particularly bad at handling stragglers and failures. When running the simulation without any straggler machines and failed tasks, the average runtime across simulations more closely resembled the average runtime of the actual executions.

In DynamicCloudSim, both the HEFT scheduling heuristic and greedy queue scheduling performed significantly better than the baseline round robin scheduler (*p*-values: 0.017 and 0.028). Furthermore, the LATE scheduler not only significantly outperformed greedy queue scheduling (*p*-value: 0.019), which served as the baseline for adaptive scheduling, but also exhibited less variance in workflow runtimes. Notably, heavy distortion of the progress estimates utilized by the LATE scheduler did not have a major impact on the scheduler's performance. Apparently, the selection strategy for speculative task replication is not essential for increasing robustness against instability.

Very similar observations were made when executing the workflow on Amazon EC2: Both HEFT and greedy queue scheduling provided significant improvements on workflow execution times when compared to static round robin scheduling (*p*-values: 0.017 and $2.6 \cdot 10^{-5}$ respectively). Furthermore, the random cloning strategy provided significant runtime improvements over the adaptive scheduling baseline, greedy queue scheduling (*p*-value: 0.045). These observations of HEFT and speculative task execution performing better than their respective baselines and adaptive scheduling performing better than static scheduling overall, were made in both the DynamicCloudSim simulation runs as well as in the workflow executions on EC2. We argue that this finding underlines the applicability of DynamicCloudSim in similar evaluation scenarios of scheduling and resource provisioning algorithms.

Fig. 11 shows the execution times of every task appearing in the Montage evaluation workflow on each of the provisioned EC2 cloud instances. These runtimes were captured to provide the HEFT scheduling heuristic with runtime estimates. Notably, we observed substantial heterogeneity across compute nodes, resulting in considerable variance in task runtimes (e.g., runtimes between 9.5 and 37.5 s for mAdd task instances). In contrast to our expectations, I/O-intensive tasks (mDiffFit, mlmgbl, mAdd) were subject to higher variance than CPU-intensive tasks (mProjectPP, mBgModel).

Over the course of the experiment, we observed that both the simulation runs on DynamicCloudSim and the actual workflow executions on EC2 provide similar answers regarding the strengths and weaknesses of the four investigated workflow schedulers as well as recommendations of which scheduler to utilize on a computational cloud. While the differences in observed variance might warrant a slight re-weighting of DynamicCloudSim's default parameters, we find that DynamicCloudSim provides an adequate model of the instability encountered in computational clouds like Amazon EC2.

6. Related work

Merdan et al. [46] and Hiraes-Carbajal et al. [47] developed simulation environments specifically for comparing different approaches to workflow scheduling on computational grids. They

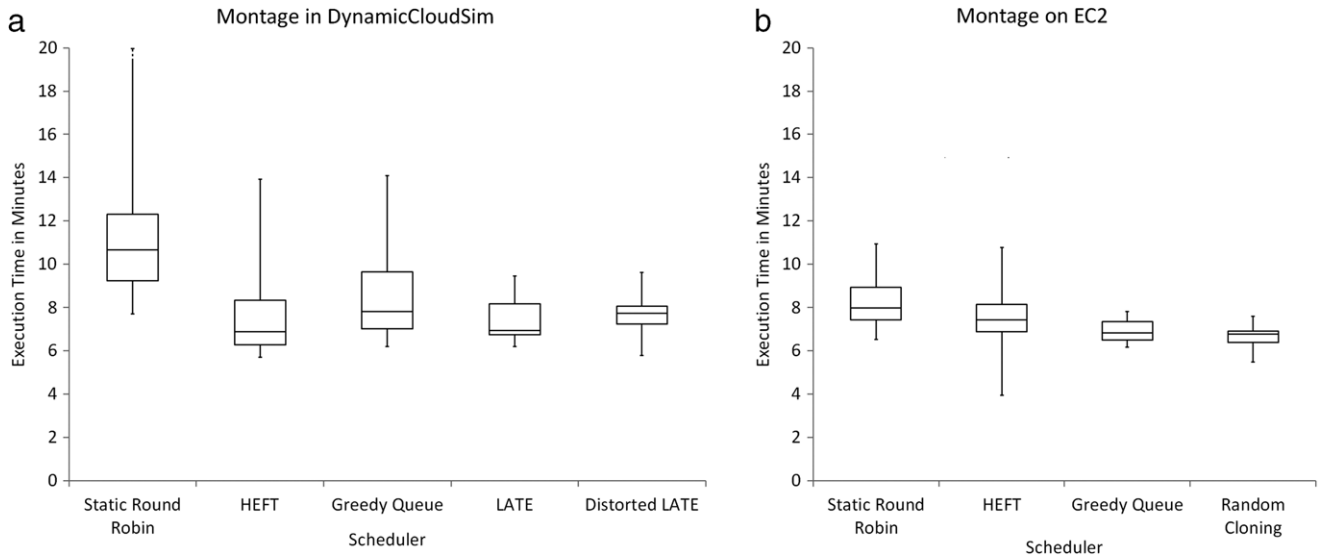


Fig. 10. The runtimes of (a) simulated Montage executions in DynamicCloudSim using default parameters contrasted to (b) runtimes of actual executions on Amazon EC2. Three of the workflow schedulers (round robin, HEFT, and greedy queue) were present both in simulation and on Amazon EC2 and can thus be compared against one another. We also compared LATE with heavily distorted progress estimates, in which tasks are selected for speculative execution nearly arbitrarily, against a random cloning strategy on EC2. Note that the measurements in DynamicCloudSim depend on its parameter configuration, which can be adjusted for different compute environments.

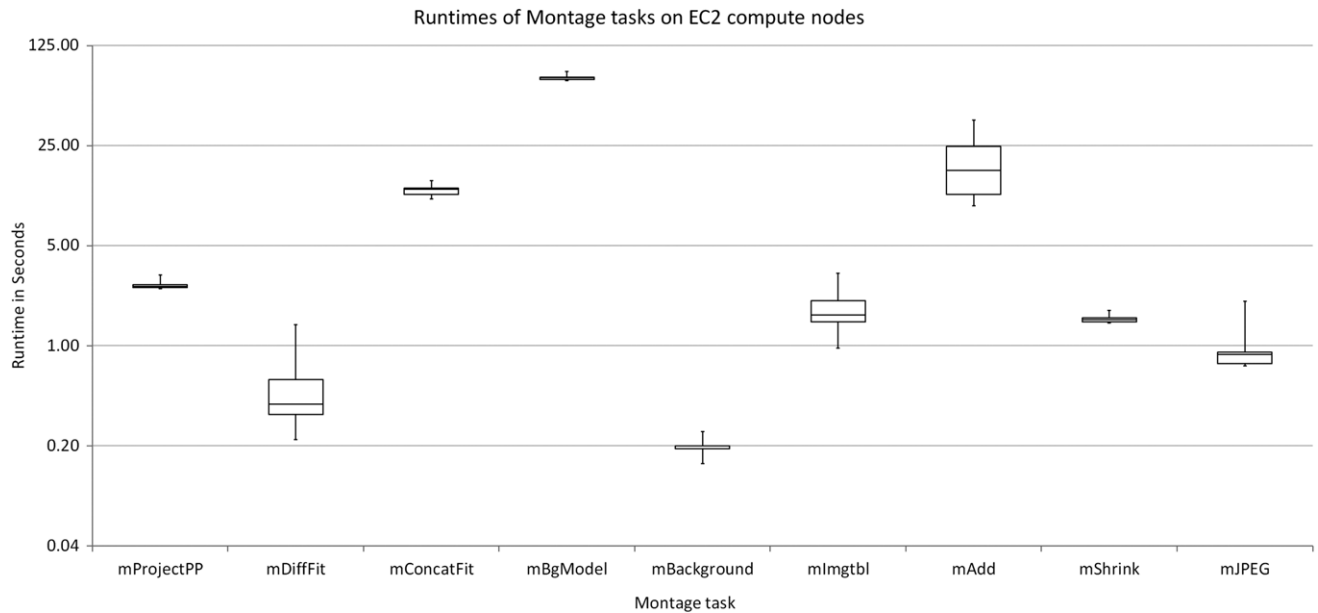


Fig. 11. Runtimes (in log scale) of every Montage task occurring in the evaluation workflow, executed once per provisioned EC2 compute node. Since the HEFT scheduling heuristic requires runtime estimates of each task on each machine, Montage tasks were launched once on each machine prior to HEFT scheduling. The measured runtimes that served as runtime estimates for HEFT are shown in this box plot.

also provide examples of possible experimental setups, yet omit the execution of these experiments. Our work differs from these publications in three ways: Firstly, by extending a universal simulation framework like CloudSim, DynamicCloudSim is not limited to the field of scientific workflows, but can be utilized for simulation of any cloud application. Secondly, our work puts a strong emphasis on instabilities in the computational infrastructure, which is important to achieve realistic results. Thirdly, we conduct an experimental validation of the changes added to the simulation toolkit.

Chen and Deelman [33] recently presented WorkflowSim as another extension to CloudSim. WorkflowSim is tightly bound to the SWfMS Pegasus [29] and adds to CloudSim (1) the workflow engine underlying Pegasus and DagMan [48], (2) an elaborate model

of node failures, (3) a model of delays occurring in the various levels of the Pegasus stack (e.g., queue delays, pre/post-processing delays, data transfer delays), and (4) the implementations of several workflow schedulers implemented in Pegasus (e.g., greedy task queue, HEFT [18], Min–Min, and Max–Min [19]). Parameters are directly learned from traces of real executions. WorkflowSim follows a quite different approach than DynamicCloudSim: WorkflowSim models delays in the Pegasus workflow stack and is thus tightly coupled to Pegasus. While it provides accurate means to simulate the execution of scientific workflows on Pegasus installations, it has no notion of heterogeneous hardware or variance in available resources. In contrast, DynamicCloudSim directly models instability and heterogeneity in the environment in which a workflow, or any other collection of computationally intensive tasks, is

Table 3
Features of CloudSim, WorkflowSim, and DynamicCloudSim.

Feature	CloudSim	WorkflowSim	DynamicCloudSim
Performance characteristics MIPS, bandwidth, memory	✓	✓	✓
Performance characteristic file I/O			✓
Runtime of a task depending on values other than MIPS			✓
Modeling of data dependencies	✓	✓	✓
Workflow parsing		✓	✓
Implementation of workflow schedulers		✓	✓
Modeling of delays at different layers of a SWfMS		✓	
Support for task clustering		✓	
Different VMs on different hosts	✓	✓	✓
Random assignment of new VM to a host			✓
Resource allocation based on compute units			✓
Dynamic changes of VM performance at runtime			✓
Modeling of failures during task execution		✓	(✓)
Introduction of straggler VMs			✓

executed. DynamicCloudSim is thus independent of the computational paradigm (e.g., scientific workflows) and the concrete system of execution (e.g., the SWfMS Pegasus). See Table 3 for a comparison of features available in CloudSim, WorkflowSim, and DynamicCloudSim.

Donassolo et al. [49] altered the SimGrid framework [50], another popular toolkit for the simulation of distributed systems such as grids, clouds, HPC or P2P systems, to increase the scalability of simulation runs. These improvements allow SimGrid to simulate the computation of workloads on tens or hundreds of thousands of heterogeneous and possibly volatile machines as encountered in volunteer computing. For the same reasons of increasing scalability, Ostermann et al. [51] also developed GroudSim, another toolkit for simulating the execution of scientific applications in a computational grid or cloud. Scalability has been previously reported to be an issue of GridSim [52] and thus also CloudSim. GridSim uses a multi-threaded core and can therefore reach the upper limit of threads supported by a typical Linux kernel when simulating very large infrastructures or concurrent user numbers beyond 10,000. In spite of this scalability limitation, we designed DynamicCloudSim as an extension of CloudSim, which brings with it the key advantage of being usable in conjunction with many of the valuable extensions of CloudSim, such as EMUSIM [23] or WorkflowSim [33].

7. Future work

In the experiments described in this paper, each virtual machine only processed one task at a time. Furthermore, issues of data locality were not incorporated yet, since we assumed all files (input, intermediate, and output) to be read from and written to shared network storage, such as Amazon S3. In future work, we would like to revisit the experiments, adding additional task slots per virtual machine and investigating the influence of storing files locally on each VM.

As mentioned in Section 5.2, when executing Montage on Amazon EC2, we did not encounter the amount of straggler machines and failed task executions suggested by DynamicCloudSim's default SaF parameters. While this observation suggests a downward correction of the parameters in question, the scale of the experiments on Amazon EC2 (20 runs on eight VMs each) was not large enough to capture reliable parameter values. In future work, we would like to run a similar experiment of substantially larger scale, in the context of which we plan to capture statistics on the rate of stragglers and failures, which we would then like to utilize for improving the default parameters of DynamicCloudSim. Another default parameter setting that we would like to re-evaluate in the context of this experiment is the frequency of VM performance baseline changes, as described in Section 3.4.

Another area of future research involves the integration of DynamicCloudSim with WorkflowSim [33] to harness the combined functionality of both CloudSim extensions. For instance, the elaborate and multi-layered failure models of WorkflowSim could further enhance the model of instability introduced by DynamicCloudSim.

8. Conclusion

We presented DynamicCloudSim as an extension to CloudSim, a popular simulator for evaluating resource allocation and scheduling strategies on distributed computational architectures. We enhanced CloudSim's model of cloud computing infrastructure by introducing models for (1) inhomogeneity in the performance of computational resources, (2) uncertainty in and dynamic changes to the performance of VMs, and (3) straggler VMs and failures during task execution.

We showed that applying these models to scientific workflow execution using four established scheduling algorithms and two evaluation workflows replicated the known strengths and shortcomings of these schedulers, which underlined the importance of adaptivity in scheduling of scientific workflows on shared and distributed computational infrastructures. Finally, we validated the models of instability introduced in DynamicCloudSim by comparing the simulated execution of a workflow in DynamicCloudSim against actual runs on Amazon EC2.

Acknowledgments

Marc Bux is funded by the Deutsche Forschungsgemeinschaft through graduate school SOAMED (GRK 1651). We further acknowledge support from the European Commission's Seventh Framework Programme (FP7) through the BiobankCloud project (project number 317871).

References

- [1] P. Mell, T. Grance, *The NIST Definition of Cloud Computing*, National Institute of Standards and Technology, 2009.
- [2] I. Foster, Y. Zhao, I. Raicu, S. Lu, Cloud computing and grid computing 360-degree compared, in: *Proceedings of the 1st Workshop on Grid Computing Environments*, Austin, Texas, 2008, pp. 1–10.
- [3] A. Beloglazov, R. Buyya, Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in cloud data centers, *Concurr. Comput.: Pract. Exper.* 24 (13) (2012) 1397–1420.
- [4] L. Wu, S.K. Garg, R. Buyya, SLA-based resource allocation for software as a service provider in cloud computing environments, in: *Proceedings of the 2011 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, IEEE, Newport Beach, California, USA, 2011, pp. 195–204.

- [5] S. Sadhasivam, N. Nagaveni, R. Jayarani, R.V. Ram, Design and implementation of an efficient two-level scheduler for cloud computing environment, in: Proceedings of the 2009 International Conference on Advances in Recent Technologies in Communication and Computing, Kottayam, India, 2009, pp. 884–886.
- [6] R.N. Calheiros, R. Ranjan, A. Beloglazov, C.A.F. De Rose, R. Buyya, CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms, *Softw.-Pract. Exp.* 41 (1) (2011) 23–50.
- [7] J. Dejun, G. Pierre, C.-H. Chi, EC2 performance analysis for resource provisioning of service-oriented applications, in: Proceedings of the 7th International Conference on Service Oriented Computing, Stockholm, Sweden, 2009, pp. 197–207.
- [8] M. Zaharia, A. Konwinski, A.D. Joseph, R.H. Katz, I. Stoica, Improving MapReduce performance in heterogeneous environments, in: Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation, San Diego, USA, 2008, pp. 29–42.
- [9] K.R. Jackson, L. Ramakrishnan, K. Muriki, S. Canon, S. Cholia, J. Shalf, H.J. Wasserman, N.J. Wright, Performance analysis of high performance computing applications on the Amazon Web services cloud, in: Proceedings of the 2nd International Conference on Cloud Computing Technology and Science, Indianapolis, USA, 2010, pp. 159–168.
- [10] S. Ostermann, A. Iosup, N. Yigitbasi, R. Prodan, T. Fahringer, D. Epema, An early performance analysis of cloud computing services for scientific computing, TU Delft, 2008.
- [11] M. Palankar, A. Iamnitchi, M. Ripeanu, S. Garfinkel, Amazon S3 for science grids: a viable solution? in: Proceedings of the 1st Workshop on Data-aware Distributed Computing, Boston, USA, 2008, pp. 55–64.
- [12] J. Schad, J. Dittrich, J.-A. Quiané-Ruiz, Runtime measurements in the cloud: observing, analyzing, and reducing variance, *Proc. VLDB Endow.* 3 (1) (2010) 460–471.
- [13] A. Iosup, N. Yigitbasi, D. Epema, On the performance variability of production cloud services, in: Proceedings of the 2011 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, Newport Beach, California, USA, 2011, pp. 104–113.
- [14] M. Bux, U. Leser, Parallelization in scientific workflow management systems, *CoRR abs/1303.7*, 2003.
- [15] T.D. Braun, H.J. Siegel, N. Beck, L.L. Boloni, M. Maheswarans, A.I. Reuther, J.P. Robertson, M.D. Theys, B. Yao, D. Hensgen, R.F. Freund, A comparison study of eleven static heuristics for mapping a class of independent tasks onto heterogeneous distributed computing systems, *J. Parallel Distrib. Comput.* 61 (2001) 810–837.
- [16] J. Blythe, S. Jain, E. Deelman, Y. Gil, K. Vahi, A. Mandal, K. Kennedy, Task scheduling strategies for workflow-based applications in grids, in: Proceedings of the 5th IEEE International Symposium on Cluster Computing and the Grid, vol. 2, Cardiff, UK, 2005, pp. 759–767.
- [17] Y. Gil, E. Deelman, M. Ellisman, T. Fahringer, G. Fox, D. Gannon, C. Goble, M. Livny, L. Moreau, J. Myers, Examining the challenges of scientific workflows, *IEEE Comput.* 40 (12) (2007) 24–32.
- [18] H. Topcuoglu, S. Hariri, M.-Y. Wu, Performance-effective and low-complexity task scheduling for heterogeneous computing, *IEEE Trans. Parallel Distrib. Syst.* 13 (3) (2002) 260–274.
- [19] A. Mandal, K. Kennedy, C. Koelbel, G. Marin, J. Mellor-Crummey, B. Liu, L. Johnson, Scheduling strategies for mapping application workflows onto the grid, in: Proceedings on the 14th IEEE International Symposium on High Performance Distributed Computing, Durham, USA, 2005, pp. 125–134.
- [20] J. Yu, R. Buyya, A budget constrained scheduling of workflow applications on utility grids using genetic algorithms, in: Proceedings of the 1st Workshop on Workflows in Support of Large-Scale Science, Paris, France, 2006, pp. 1–10.
- [21] R. Buyya, M. Murshed, GridSim: A toolkit for the modeling and simulation of distributed resource management and scheduling for grid computing, *Concurr. Comput.: Pract. Exper.* 14 (13–15) (2002) 1175–1220.
- [22] S.K. Garg, R. Buyya, NetworkCloudSim: Modelling parallel applications in cloud simulations, in: Proceedings of the 4th IEEE International Conference on Utility and Cloud Computing, IEEE, Melbourne, Australia, 2011, pp. 105–113.
- [23] R.N. Calheiros, M.A.S. Netto, C.A.F. De Rose, R. Buyya, EMUSIM: an integrated emulation and simulation environment for modeling, evaluation, and validation of performance of Cloud computing applications, *Softw.-Pract. Exp.* 43 (2013) 595–612.
- [24] S. Frey, W. Hasselbring, The CloudMIG approach: model-based migration of software systems to cloud-optimized applications, *Intern. Journal on Advances in SW 4* (3 and 4) (2011) 342–353.
- [25] B. Schroeder, G.A. Gibson, A large-scale study of failures in high-performance-computing systems, in: Proceedings of the 36th International Conference on Dependable Systems and Networks, Philadelphia, USA, 2006, pp. 249–258.
- [26] Z. Zhang, Z. Li, K. Wu, D. Li, H. Li, Y. Peng, X. Lu, VMThunder: fast provisioning of large-scale virtual machine clusters, *IEEE Trans. Parallel Distrib. Syst.* PP (99) (2014) 1–11.
- [27] J. Dean, S. Ghemawat, MapReduce: simplified data processing on large clusters, *Commun. ACM* 51 (1) (2008) 107–113.
- [28] G.B. Berriman, E. Deelman, J. Good, J. Jacob, D.S. Katz, C. Kesselman, A. Laity, T.A. Prince, G. Singh, M.-h. Su, Montage: a grid-enabled engine for delivering custom science-grade mosaics on demand, in: Proceedings of the SPIE Conference on Astronomical Telescopes and Instrumentation, vol. 5493, Glasgow, Scotland, 2004, pp. 221–232.
- [29] E. Deelman, G. Singh, M.-H. Su, J. Blythe, Y. Gil, C. Kesselman, G. Mehta, K. Vahi, G.B. Berriman, J. Good, A. Laity, J.C. Jacob, D.S. Katz, Pegasus: a framework for mapping complex scientific workflows onto distributed systems, *Sci. Program.* 13 (3) (2005) 219–237.
- [30] E. Deelman, G. Singh, M. Livny, B. Berriman, J. Good, The cost of doing science on the cloud: The montage example, in: Proceedings of the 2008 Conference on Supercomputing, IEEE, Austin, Texas, 2008, pp. 1–12.
- [31] C. Hoffa, G. Mehta, T. Freeman, E. Deelman, K. Keahey, B. Berriman, J. Good, On the use of cloud computing for scientific workflows, in: Proceedings of the 4th IEEE International Conference on eScience, Indianapolis, USA, 2008, pp. 640–645.
- [32] K. Lee, N.W. Paton, R. Sakellariou, E. Deelman, A.A.A. Fernandes, G. Mehta, Adaptive workflow processing and execution in Pegasus, *Concurr. Comput.: Pract. Exper.* 21 (16) (2009) 1965–1981.
- [33] W. Chen, E. Deelman, WorkflowSim: A toolkit for simulating scientific workflows in distributed environments, in: Proceedings of the 8th IEEE International Conference on eScience, Chicago, USA, 2012, pp. 1–8.
- [34] S. Pabinger, A. Dander, M. Fischer, R. Snajder, M. Sperk, M. Efreanova, B. Krabichler, M.R. Speicher, J. Zschocke, Z. Trajanoski, A survey of tools for variant analysis of next-generation genome sequencing data, *Brief. Bioinform.* 15 (2) (2014) 256–278.
- [35] G. Juve, A. Chervenak, E. Deelman, S. Bharathi, G. Mehta, K. Vahi, Characterizing and profiling scientific workflows, *Future Gener. Comput. Syst.* (2012).
- [36] B. Langmead, C. Trapnell, M. Pop, S.L. Salzberg, Ultrafast and memory-efficient alignment of short DNA sequences to the human genome, *Genome Biol.* 10 (3) (2009) R25.
- [37] M. David, M. Dzamba, D. Lister, L. Ilie, M. Brudno, SHRIMP2: Sensitive yet practical short read mapping, *Bioinformatics* 27 (7) (2011) 1011–1012.
- [38] Y. Chen, T. Souaiaia, T. Chen, PerM: efficient mapping of short sequencing reads with periodic full sensitive spaced seeds, *Bioinformatics* 25 (19) (2009) 2514–2521.
- [39] H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin, The Sequence Alignment/Map format and SAMtools, *Bioinformatics* 25 (16) (2009) 2078–2079.
- [40] D.C. Koboldt, K. Chen, T. Wylie, D.E. Larson, M.D. McLellan, E.R. Mardis, G.M. Weinstock, R.K. Wilson, L. Ding, VarScan: variant detection in massively parallel sequencing of individual and pooled samples, *Bioinformatics* 25 (17) (2009) 2283–2285.
- [41] T. Oinn, M. Addis, J. Ferris, D. Marvin, M. Senger, M. Greenwood, T. Carver, K. Glover, M.R. Pocock, A. Wipat, P. Li, Taverna: a tool for the composition and enactment of bioinformatics workflows, *Bioinformatics* 20 (17) (2004) 3045–3054.
- [42] B. Ludäscher, I. Altintas, C. Berkley, D. Higgins, E. Jaeger, M. Jones, E.A. Lee, J. Tao, Y. Zhao, Scientific workflow management and the Kepler system, *Concurr. Comput.: Pract. Exper.* 18 (10) (2006) 1039–1065.
- [43] T. White, Hadoop: The Definitive Guide, third ed., O'Reilly Media, Inc., Sebastopol, USA, 2012.
- [44] S.L. Garfinkel, An evaluation of Amazon's grid computing services: EC2, S3 and SQS, Technical Report TR-08-07, School for Engineering and Applied Sciences, Harvard University, MA, 2007.
- [45] C. Pellingeeas, Performance evaluation of virtualization with cloud computing (Master of Engineering thesis), Edinburgh Napier University, 2010.
- [46] M. Merdan, T. Moser, D. Wahyudin, S. Biffl, P. Vrba, Simulation of workflow scheduling strategies using the MAST test management system, in: Proceedings of the 10th International Conference on Control, Automation, Robotics and Vision, Hanoi, Vietnam, 2008, pp. 1172–1177.
- [47] A. Hiraes-Carbajal, A. Tchernykh, R. Röblitz, R. Yahyapour, A grid simulation framework to study advance scheduling strategies for complex workflow applications, in: Proceedings of the 24th IEEE International Symposium on Parallel & Distributed Processing, Workshops and PhD Forum, Atlanta, USA, 2010, pp. 1–8.
- [48] P. Couvares, T. Kosar, A. Roy, J. Weber, K. Wenger, Workflow management in condor, in: I.J. Taylor, E. Deelman, D. Gannon, M. Shields (Eds.), *Workflows for e-Science*, first ed., Springer, New York, USA, 2007, pp. 357–375.
- [49] B. Donassolo, H. Casanova, A. Legrand, P. Velho, Fast and scalable simulation of volunteer computing systems using SimGrid, in: Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing, Chicago, USA, 2010, pp. 605–612.
- [50] H. Casanova, A. Legrand, M. Quinson, SimGrid: A generic framework for large-scale distributed experiments, in: Proceedings of the Tenth International Conference on Computer Modeling and Simulation, Cambridge, UK, 2008, pp. 126–131.
- [51] S. Ostermann, K. Plankensteiner, R. Prodan, T. Fahringer, GroudSim: An event-based simulation framework for computational grids and clouds, in: CoreGRID/ERCIM Workshop on Grids, Clouds and P2P Computing in conjunction with EuroPAR 2010, Ischia, Italy, 2010, pp. 305–313.
- [52] W. Depoorter, N.D. Moor, K. Vanmechelen, J. Broeckhove, Scalability of grid simulators: An evaluation, in: Proceedings of the 14th international Euro-Par conference on Parallel Processing, Las Palmas de Gran Canaria, Spain, 2008, pp. 544–553.



Marc Bux is a Research Associate and Ph.D. candidate of the research training group SOAMED at Humboldt-Universität zu Berlin. He investigates adaptive scheduling techniques for scientific workflows. His research interests include cloud computing, scientific workflows, and high-throughput technologies in bioinformatics.



Ulf Leser is a Professor of Knowledge Management in Bioinformatics at Humboldt-Universität zu Berlin. His research interests include scalability and expressiveness of scientific workflows, integrated analysis of high-throughput data in bioinformatics, and biomedical data integration and text mining. He is a PI in more than half a dozen interdisciplinary research projects with colleagues from molecular biology, systems biology, and medicine. He regularly reviews for various journals, including *Oxford Bioinformatics*, *Briefings in Bioinformatics*, *Nucleic Acids Research*, *VLDB Journal*, *ACM Transactions on Database Systems*, *IEEE Transactions on Knowledge and Data Engineering*.