

RESEARCH

Crane Cloud: a resilient multi-cloud service layer for resource constrained settings

Engineer Bainomugisha^{1*} and Alex Mwotil^{1,2}

*Correspondence:

baino@mak.ac.ug

¹Department of Computer Science, School of Computing and Informatics Technology, College of Computing and Information Sciences, Makerere University, Kampala, Uganda
Full list of author information is available at the end of the article

Abstract

Whereas the main cloud providers have set up cloud services on stable infrastructure, developers and users situated in low-resource settings access cloud services and platforms using low-end computing devices that often connect to the Internet via slow mobile connections. These settings require custom software abstraction layers that consider such bandwidth constraints and intermittent connections as a rule rather than the exception. In this paper, we identify key challenges for developing for and accessing cloud services in resource constrained settings, namely, (1) Frequent Internet partitions and bandwidth constraints, (2) Data jurisdiction restrictions, (3) Vendor lock-in, and (4) Poor quality of service. To address these challenges, we propose a set of design considerations and properties for a resilient multi-cloud service layer, that includes: (1) Containerisation and orchestration of applications, (2) Service scheduling and replication, (3) Portability and multi-cloud migration, (4) Resilience to network partitions and bandwidth constraints, (5) Automated service discovery and load balancing, (6) Localised image registry, and (7) Support for platform monitoring and management. We present a prototype validation case study, Crane Cloud, an open source multi-cloud service abstraction layer built on-top of Kubernetes that is designed with inherent support for resilience to network partitions, microservice orchestration (deployment, scaling and management of containerized applications)a localized image registry, support for migration of services between private and public clouds to avoid vendor lock-in issues and platform monitoring. We evaluate the performance and user experience of Crane Cloud by implementing and deploying a computational and bandwidth intensive machine learning system shows lower response time compared when hosted on other public clouds.

Keywords: Microservices; Kubernetes; Containers; Orchestration; Low-resource settings; Portable cloud apps; Cloud native platforms

1 Introduction

Cloud computing has become a popular model of delivering computing services over the Internet (“the cloud”) with flexible pricing models such as pay per use. It provides flexible on-demand access to an elastic computing resource base and represents the infrastructure, software, platforms, storage and application containers as a cloud where users can provision and access services over a network. Cloud Computing has a large number of deployed solutions in education [1] [2], big data computing [3], health [4, 5], private sector and government [6] [7] domains. The National Institute of Standards and Technology (NIST) describes the essential characteristics of cloud computing as an on-demand self service with a broad network access, computing

resource pooling with rapid elasticity and measurements done according to compute, network and storage usage attributes [8]. The computing resources can easily be scaled horizontally, vertically or both as and when the need arises. The compute, network and storage usage metrics provide a basis for billing by the cloud service providers. The NIST describes four deployment models (private, community, public and private clouds) and three primary cloud service models: (i) Infrastructure as a Service (IaaS): Provider offers hardware (compute, storage and network) and the associated software as a service [9, 10, 11], (ii) Platform as a Service (PaaS): A software platform which is an additional abstraction layer is provided [12, 13], (iii) Software as a Service (SaaS): Applications and software running in the cloud are offered as a service to clients for example Gmail and Facebook [12].

Recently cloud computing and application environments have evolved from monolithic to microservice architectures and platform support:

1.1 Monolithic Architecture

The monolithic architecture follows from an all-inclusive pattern where components of an application are packaged into one process to meet customer demands [14]. All the functional components of the system are deployed as a single application and the dependencies tightly coupled and multiple instances of the application may be run behind a load balancer to improve on service availability and for scalability [15]. The application components are packaged into a single-tier and deployed as a tightly coupled unit. Monolithic applications offer benefits such as ease in development, testing and deployment but also suffers from major drawbacks [16]:

- The resultant code base is difficult to maintain and understand greatly reducing the development and deployment speeds.
- Application changes may require a modification to the entire code base. A change in the programming language, platform, team or technology is extremely difficult to handle.
- Scaling of the application presents new challenges as this can only be done in one dimension. To scale the application, new instances may need to be dynamically provisioned to support the new load.

1.2 Microservice Architecture

Many organizations are now abandoning the monolithic model of developing and packaging applications partly driven by the advent of Service-Oriented Architectures (SOA), web services and Service-Oriented Computing (SOC). SOA introduced the concept of building and deploying individual service-based components linked together at runtime to provide the ‘whole’ system functionality [17]. In microservices architecture, a large application is built as a suite of loosely coupled specially-focused independently deployable services. The modules perform specific business goals over well-defined interfaces for communication with other modular parts of the system. This architecture offers several benefits:

- Developer teams can work independently and efficiently and use different programming languages of choice. This makes it easy to work with new and emerging technologies in delivering modular functionality where required.
- Continuous Integration/Continuous Deployment (CI/CD) is easy and flexible to work with and new members can easily be added/integrated to the teams.

- Microservices can easily be scaled with based on application demands

Cloud providers have set up cloud services with an assumption of constant availability of stable infrastructure for the developers and users. This assumption however is not true for the users and developers situated in low-resource settings that face major constraints when accessing their hosted applications and cloud platforms. In such settings, challenges such as frequent Internet partitions, unannounced power shutdowns, poor quality of services, among others, are the rule rather than the exception. We further elaborate these challenges and requirements using a concrete case study in Section 2. The paper presents:

- Description of the challenges and requirements for designing and operating a resilient multi-cloud model for low resource settings.
- The design considerations and properties for implementing a resilient multi-cloud and bare-metal application cluster such as application state, networking, loadbalancing, monitoring and service exposure for external user access.
- A prototype implementation of a resilient multi-cloud and bare-metal application cluster (Crane Cloud) which is a subset instantiation of the design and implementation options available. This will provide researchers and practitioners with a review point for further research and implementation cogitation respectively.
- Evaluation of performance and user experience of Crane Cloud platform by implementing and deploying a computational and bandwidth intensive machine learning system that shows lower response time compared when hosted on other public clouds.

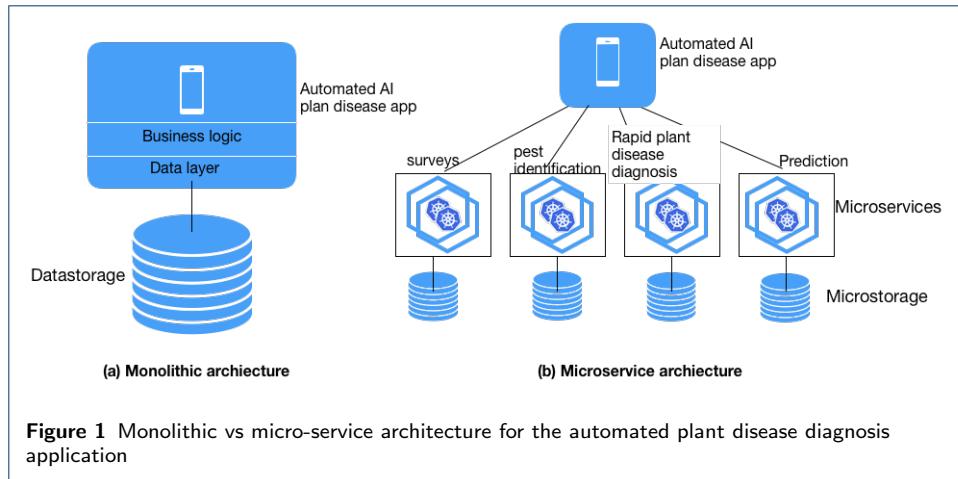
2 The Need for Resilient Multi-cloud Services

In this section, we present a motivating case study from which we derive and discuss the requirements of a resilient-multi cloud service abstraction layer.

2.1 Scenario: Containerised Machine Learning-based services for automated plant disease diagnosis

The Automated Plant Disease Diagnosis (APDD) system is a real-world case study of a machine learning system used by farmers and agricultural scientists in low-resource settings to diagnose plant diseases [18]. Specifically, it is designed to provide near real-time in-field diagnosis of plant diseases and identification of pests for cassava crops by farmers and agricultural experts in the East African region and national agricultural research organisations. Originally conceived as a monolithic system, the APDD is transitioning to a microservice architecture to optimise efficiency and work around the technology constraints (Figure 1). It consists of several microservices including, (1) pest surveys module, (2) in-field automated pest identification and count (such as white-flies) on plant leaves, (3) rapid plant disease diagnosis, and (4) prediction and spatial analysis of plant disease incidences and spread.

The data pipeline for APDD system involves huge training datasets of over 300GB including images of diseased and healthy plant leaves. The dataset is used for training and evaluation data for the machine learning models. The images are crowd sourced from local farmers using mobile phones and can be stored on available public cloud platforms as well as on local storage to allow online training and evaluation



of the machine learning models. The process flow typically involves uploading and downloading a huge datasets from the local storage and cloud systems. Such a dataset can take up to 24 hours or so to upload to a public cloud service provider over slow intermittent connections that characterise many low-resource settings.

2.2 Scenario Analysis and Requirements

The above case study signals one of many other similar systems faced with technical hurdles when delivering cloud-based solutions in a low resource setting. To an extent, it reveals key requirements and challenges that can be addressed by containerised cloud-based systems that are developed in and for use in such environments. If not addressed, they are likely to be significant barriers to development, adoption and use of cloud-based solutions such as machine learning services for many developers, researchers, startups, students, users and organisations based in such settings. We argue that it should be possible for users situation in such settings to take advantage of the benefits offered by cloud-based computing models through creation of appropriate abstraction service layers. We discuss the key requirements and challenges below:

2.2.1 Frequent Internet partitions and bandwidth constraints

Whereas the main cloud providers have set up cloud services on stable infrastructure, developers and users situated in low resource settings face major constraints when accessing their hosted applications and cloud platforms. In such settings, access to cloud services and platforms is often through low-end smartphones over slow and erratic 2G/3G/4G [19] connections that are characterised by frequent network partitions and bandwidth constraints. As can be observed from the above scenario, a developer uploading training dataset of 300GB images to a cloud platform could easily take several days over a slow connection. The in-field data collection by farmers to public cloud repository can be very slow because of sporadic Internet connectivity in rural areas. Access speed to cloud services would be faster if the cloud data centres were located near to the users, however, almost all the data centers of the popular cloud providers are located overseas [20]. For instance, on the African continent, almost 70% of the content and services consumed by the Internet users

is hosted and delivered from overseas data centers resulting in poor user experience resulting from high latencies. Furthermore, mobile data costs in Africa are significantly high in real and GDP-relative terms in the world with the prices averaging US\$ 7.04 per 1GB [21] [22] and speeds as low as 56 kbit/s. However, it is important to note that these issues are similar in any other locations with similar resource constraints. The design and setup of most public cloud platforms assume stable infrastructure across the users and leave the the issues of connectivity challenges to the application developers. This shifts the burden and unnecessary complexity to the application developers who must consider offering different function-trimmed variations of their app services for users situated in low-resource settings, for instance, Facebook Lite app [23], WhatsApp Lite [24], Uber Lite [25], Google Go [26] and Gmail Basic [27] for slow Internet connections and low-end devices.

The above issues form the motivation of the research work in this paper. In that cloud-platforms need to be designed and optimised to work for frequent Internet partitions and bandwidth constraints challenges since they are the rule rather than the exception in low-resource settings.

2.2.2 Data jurisdiction restrictions

Many countries are coming up with new laws and guidelines to regulate the location of data and access. Countries have recently enacted laws that enforce data sovereignty and prevent data from leaving the country's boundaries. This becomes difficult to implement in regions where public clouds do not have physical presence. For instance, on the African continent where public cloud data centers are still sparse, it becomes almost impossible to comply with such policies. The microservice architecture and cloud orchestration platforms such as Kubernetes [28] promises potential remedies to this challenge. In the above scenario, the APDD system can be broken into independent microservices each with different data jurisdiction policies. For instance, the data storage and management microservice needs to be enforced to remain within the boundaries of the country while the plant disease prediction service can run in a public cloud without restrictions and benefit from the rich machine learning libraries and tools. Such a setup would require a multi-cloud environment that spawns boundaries with support for data jurisdiction policies specific to a microservice and use case.

The popularity of cloud computing solutions has introduced gaps in key processes of the data management cycle (collection,storage, analysis and use/reuse). Most cloud solutions do not provide controls over where data should be stored and in cases where there is no infrastructure presence, users have to make exceptions at the expense of prescribed hosting recommendations. Cloud providers also distribute content over spatial infrastructure located in different regions to maintain the cloud Quality of Service (QoS) along dimensions of performance, availability and reliability. This leads to silos of data spanning different geographical regions that users may have no idea or control of.

Governments are now adopting data localization where a nation's data is collected, processed, and/or stored inside the country and data sovereignty where data is subject to the laws and governance structures within the nation it is collected. This has led countries and regions to enact Data Protection and Privacy laws such

as the European Union (EU)'s operational General Data Protection Regulation (GDPR) that impose stringent policy controls on the use of Personally identifiable information (PII).

2.2.3 Vendor lock-in

Vendor lock-in, which is a user difficulty of switching from one vendor to another, is regarded as one of the major deterrents in the adoption of cloud by developers as well as small and medium-sized enterprises (SMEs) [29]. These user categories may not have the local computing resources to run their workloads and most often resort to cloud providers but flexibility to switch/shift between providers is one of their desirable properties. Other than vendor lock-in, there are other variations including product lock-in, version lock-in, architecture lock-in, platform lock-in, skills lock-in and mental lock-in [30]. Public clouds offer provider-specific proprietary solutions to meet the market demands and this has resulted in an interoperability, integration and portability downside across the cloud divide. Consequently, the applications developed for a specific cloud provider such as Amazon Web Services (AWS) may not work out-of-the box with another cloud provider such as IBM cloud due to inherent dependencies of the underlying IT infrastructure (hardware and software), cloud semantics and non-standardized APIs [31] [32]. The migration of cloud services from one provider to another usually requires major reworks on the application that may be catastrophic for mission-critical systems. For instance, the above case study may use vendor-specific machine learning libraries and tools making it difficult to migrate to another cloud when there is need.

The vendor lock-in challenge emphasizes the need for new abstraction layers to alleviate the difficulty of migrating applications between clouds. New platforms and architectures such as Kubernetes [28] offer new possibilities to implement a vendor neutral layer on top of public and private clouds. However, the current offerings of managed Kubernetes layers assume migration of services in situations where there is stable connectivity and infrastructure and are not designed for data centres that may be characterised by frequent network partitions and bandwidth constraints.

2.2.4 Poor quality of service

In cases where cloud services are offered from data centres located overseas and far from the consumers, user experience can be poor compared to closely-located content. For a user located in Africa, the average round trip time to reach content hosted within Africa is **2.243 ms, 303.287 ms** for content hosted in the USA and **208.795 ms** for content hosted in Europe. Coupled with high costs of Internet access, limited bandwidth and high latencies mainly introduced by distance, the content load times are high and this negatively impacts the user experience. Consider for example in the above case study where a cloud-based service for farmers is situated on a public cloud with a data center situated in the USA while the target farmers are located on the African continent. The longer the distance, the higher the number of intermediary links which can act as failure points (bottlenecks) and potentially introduce network packet losses. Furthermore, there are applications that are delay-sensitive and these require optimal and stringent quality of service parameter values such as low latency, low jitter and minimal or no packet loss for

best performance. Currently public cloud providers attempt to solve this challenge by moving services closer to the user. This approach however assumes presence of data centers closer to the user. Unfortunately this is not always the case for users located in regions where public cloud data centers are sparse.

In the next section, we present the design options that need to be considered when developing a multi-cloud service abstraction layer to address the above challenges particularly in low resource settings. In the subsequent sections, we demonstrated the instantiation of the design considerations in a practical open source cloud project.

3 Design Considerations for a Resilient Multi-cloud Service

In this section, we present the design considerations and properties for a resilient multi-cloud service layer that is envisioned to meet the above requirements, namely, (1) Frequent Internet partitions and bandwidth constraints, (2) Data jurisdiction restrictions, (3) Vendor lock-in, and (4) Poor quality of service. In Section 4 we shall present the first prototype implementation that instantiates some of these design considerations (DC) and properties:

3.1 DC 1: Containerisation and orchestration of applications

As introduced in Section 1.2, many organisations are now adopting microservice architectures in place of traditional monolithic architectures so as to truly reap from the benefits of modern cloud services. Microservice architectures involve collaborations between different fine-grained and independently deployable modules usually without a centralized controller to achieve the desired overall functionality of the system [33] [34]. Driven by application features such as scalability, agility, performance and fault-tolerance, microservice architectures involve autonomous software development teams independently working to build loosely coupled application features and employing collaborative workflows and automation tools from version control systems to full scale production deployment [35]. A number of popular technology companies such as Uber, Spotify, Netflix, Amazon and Ebay are now using microservices at the core of their business processes and have achieved differing levels of reliability and scalability in their services [34]. As part of the inceptor team of the microservice terminology, Fowler and Lewis [36] identified the following key properties and benefits of microservice applications:

- 1 *Autonomous software components*: A complex system is decomposed into service-specific pluggable components along business service lines allowing for each service deployment and modification without impacting other functional facets of the application. The units are small, granular, manageable and loosely coupled and they communicate using well defined interfaces based on platform-independent data formats and technologies such as HTTP/REST or messaging solutions such as Kafka or RabbitMQ [34]. Microservice architectures thrive on the notion of 'small size' with reductions in the scope of the problem, task completion time, feedback response time and the size of the deployment unit. This in turn translates to an application's resilience to cascading failures, easier maintenance and seamless deployment. In addition, the decomposition yields small coherent components that are easy to understand and debug.

- 2 *Decentralization:* The services (components) are distributed, as there is no central controller and may store only data related to the supported business domain or different instances of the same database technology. Monolithic architectures usually have a single logical database for a range of applications. The teams are also decentralized and can adopt appropriate standards that allow them to deliver the business domain functionality without reinventing the wheel. The teams are responsible for the build it/run it cycle of the service and this improves on the quality of the code, fastens the deployment process, promotes component reusability and isolates the impact of changes on the schema.
- 3 *Technology independence:* The microservices can be built using different tools such as frameworks, programming languages and databases given the architecture supports heterogeneous technologies. The frontend service and reporting tools could, for example, be developed using a User Interface (UI) framework such as Node.js [1], the backend service could be written using Java [2], a real-time component could use C++ [3], the persistent data storage mechanisms could employ MongoDB [4]. The developers are at liberty and can independently choose a technology stack that best fits the work at hand hence placing the responsibility of development, maintenance and generally ownership to the teams.
- 4 *Automation:* One of the modern software development concepts is Continuous Integration (CI), Continuous Delivery (CD) where an application is manually deployed and Continuous Deployment (CD) where an application is automatically deployed. This has immensely changed how developers and testers ship software culminating into the famous CI/CD pipeline depicted in Figure 2. In CI, the development teams implement changes and write to version control systems with automated build and test scripts and in CD, the application is deployed/shipped to either a staging or production environment. This has been spurred by the introduction of Software Developers (Dev) and IT Operations (Ops) generally termed as DevOps model for software development: an end-to-end model for fast delivery of reliable applications and services involving cultures(by and for the people), automation (testing, feedback, deployment and performance benchmarks), quality measurement and sharing of ideas, processes and tools [37]. In an incremental migration and architectural refactoring of a commercial mobile backend (monolithic application) as a service, Balalaie et. al. [38] noted that the microservice architecture is an enabler for use of DevOps. Automation does not imply management overhead on introduction of new applications [39] but rather increased agility and reliability.

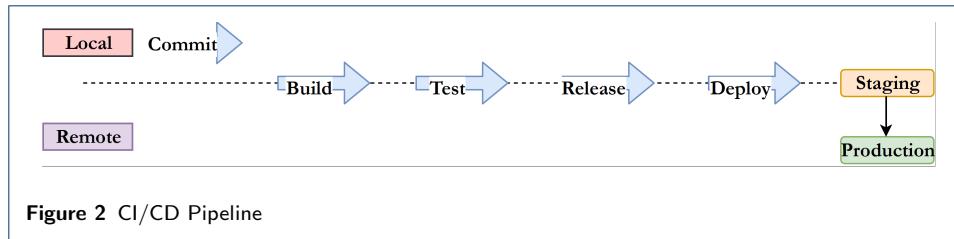
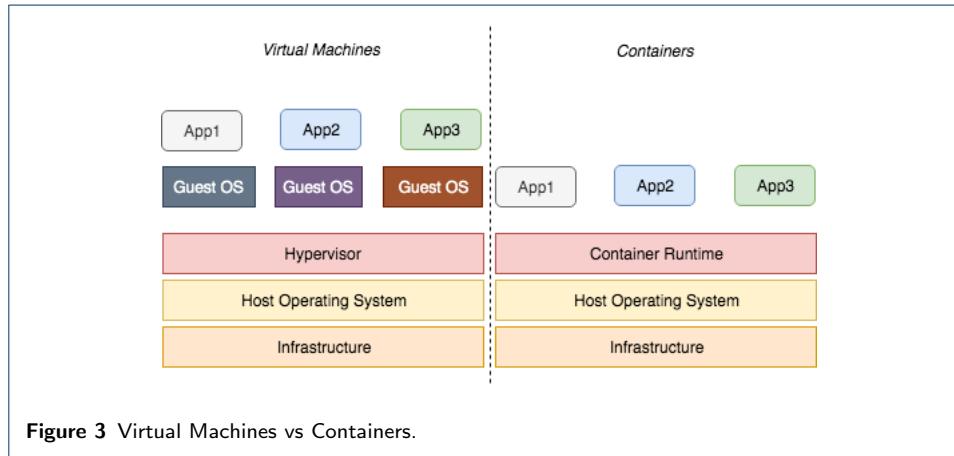
At the application development and deployment levels, developers require appropriate abstractions to efficiently package and deploy microservices on a cloud infrastructure. Recently, containers have emerged as a novel abstraction to deploy

[1]<https://nodejs.org>

[2]<https://www.java.com>

[3]<https://www.cplusplus.com>

[4]<https://www.mongodb.com>

**Figure 2** CI/CD Pipeline

microservices as opposed to traditional virtual machine approaches. Containers form the basic deployment units to release and ship microservices [33] given their modular design which closely maps with the functional decomposition of a complex business application. The lightweight nature of containerized applications coupled with functional and configuration encapsulation facilitates replication and portability, are cost-efficient and have a reduced overhead on the operation and maintenance line. It is for these reasons that containers emerged as the most suitable packaging toolset for microservices.

A container is a lightweight “virtual machine” based on the Linux Kernel Extension (LXC) [5] technology that allows running of many (up to hundreds) isolated and autonomous Linux environments on a single server or virtual machine. It is a collection of communicating components such as application code, runtime, system tools, system libraries and settings required to run an application. Container history dates back to the early 2000’s with the introduction of FreeBSD jails by Derrick T. Woolworth at R&D Associates for FreeBSD. The jails would provide a logical isolation environment through sandboxing for system features such as the filesystem, users and network mimicking a virtual machine but running on the same operating system [40]. In comparison with virtual machines, containers consume much less computing resources and take the least provisioning time because virtual machines require different instances of the operating system (guest OS) while containers use the same host operating system as shown in Figure 3. CoreOSrkt [6], Mesos Con-

[5]<https://linuxcontainers.org>

[6]<https://coreos.com/rkt/>

tainerizer [7], LXC Linux Containers, OpenVZ [8] and containerd [9] are examples of containerization technologies but Docker [10] is by far the most popular.

Fully leveraging the elasticity of using container-based virtualization requires *automated orchestration* and cluster management tools such as Kubernetes [11], Docker Swarm [12] and DC/OS [13] that provide an abstraction layer between computing resource pools and the applications. These tools provide a number of attributes important to implement service discovery, scalability, load balancing, service replication and provisioning of replicas across multiple compute nodes [39] [34]. Despite its complexity such as in installation, Kubernetes is the most widely adopted and powerful container orchestration tool owing to its immense scalability, performance and advanced automation features. It has inbuilt monitoring and logging libraries and processes which are lacking in the other tools [41]. In Kubernetes, the applications are packaged as containers and wrapped in a pod (a group of containers that form the basic deployable Kubernetes unit), which can then be deployed via a declarative manifest (YAML [14]) file. In this file, the user describes the desired state of the application such as name, replica count, labels, storage mounts and exposed ports and the deployment controller works to ensure that this state is maintained at all times for example by replacing pods that fail or are evicted from their nodes.

3.2 DC 2: Application placement and replication

In a multi-cluster/multi-node environment, placement involves determining what cluster/node should host an application based on factors such as application affinities, resource (storage, memory and network) availability, user preferences and data jurisdictions. The clusters/nodes could be located in zones/sites with varying availability and regulatory constraints. Application replication ensures that there is operational business continuity in the face of downtime incurred through computing equipment failures, natural disasters, planned maintenance operations [42], power outages, unreliable network connectivity, limited bandwidth and utilization surges. Replication can further be used to realize scalability, availability and fault-tolerance of an application under scheduled or unplanned downtime periods. The replica count, the number of clones of an application, depends on the reliability assurance as a requirement for an application and also the popularity of the service in a cluster/node region. Replication strategies fall into two broad categories:

- 1 Static replication strategy where the number of nodes and replicas is defined beforehand
- 2 Dynamic replication strategy where replicas are automatically created or destroyed based on changes such as user density, performance, storage utilization, loadbalancing features and bandwidth consumption.

[7]<http://mesos.apache.org/documentation/latest/containerizers/>

[8]<https://openvz.org>

[9]<https://containerd.io>

[10]<https://www.docker.com>

[11]<https://kubernetes.io>

[12]<https://docs.docker.com/engine/swarm/>

[13]<https://dcos.io>

[14]<https://yaml.org>

For the rest of this section, we shall consider the four (AI) microservices for the Automated Plant Disease Diagnosis (APDD) system: Surveys, pest identification, rapid plant disease diagnosis and prediction. A resilient multi-cloud service should provide an adaptive service replication approach that considers the following attributes:

3.2.1 User defined and cost-sensitive replication policies

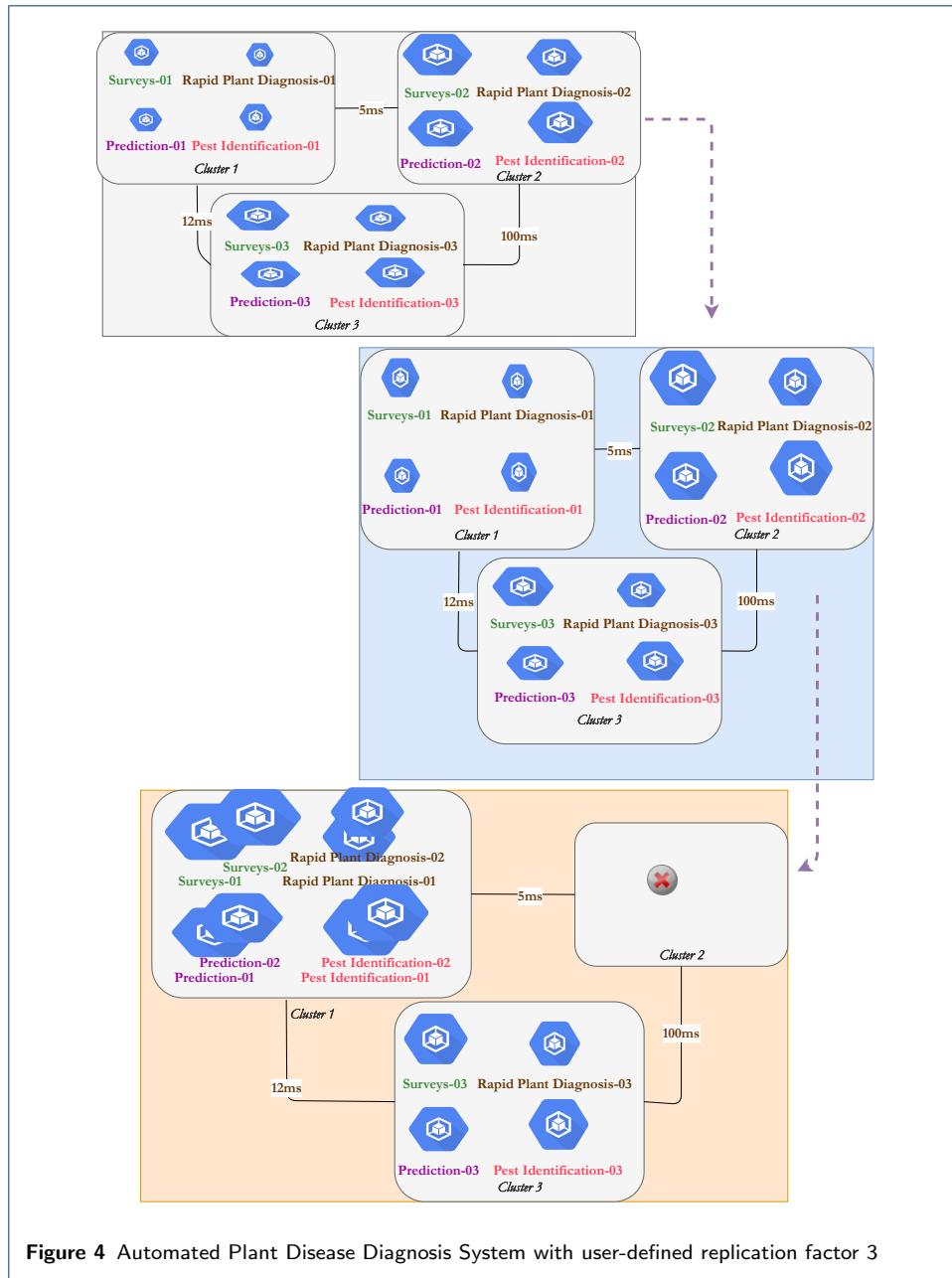
The cloud service should operate only within the user-defined replication limits but also ensuring minimal replication costs between the clusters and the target nodes. In the deployment of APDD, the user may specify a replication limit of 3: the cloud service should ensure that there are three instances of APDD microservices available at all times. Additionally, the replication approach in cases of downtime should consider the replication costs such as the impact on the network performance whenever provisioning is required. It should also be noted that the cloud service provider may impose restrictions based on, for example resource availability, which the user adheres to but regardless, the user will operate on a higher level of abstraction.

Figure 4 shows a 3-cluster cloud service located in different regions as shown by the link latencies. The user specified a replication count of 3 (classic 3-replica replication strategy [43]) and the services are initially deployed on each of the three clusters. Cluster 2 experiences a downtime and this requires scheduling of the four microservices into another cluster. The decision on which cluster the services will be provisioned on should consider the transmission cost and this will ultimately be Cluster 1. Considering the costs that the fixed-replica count strategy may impose, [44] presents a dynamic cost-effective replication algorithm for data in cloud data centers. This approach requires computation of a reliability requirement value which informs the replica count. The default replica count is 1 and this will be scaled upwards to meet the reliability requirement of a data intensive system. The initial costs of this approach are significantly low but may increase exponentially with provisioning of more replicas. This algorithm can further be enriched by introducing user-defined boundary limits while ensuring the reliability requirement is maintained.

3.2.2 Quality of service (QoS) and high availability

Some applications in heterogeneous clouds have higher QoS requirements in comparison with others for example a critical medical diagnosis system that should operate under stringent availability and consistency constraints. In distributed cloud environments that support service geo-replication, maintaining consistency and performance consecutively is desirable but not fully achievable according to the CAP's theorem [45]. Consistency may be achieved but at the expense of degraded system performance. A number of research works have been published in the line of QoS and high availability.

Esteves *et. al.* [46] developed a novel QoS consistency model for geo-replicated data in cloud computing. In this model, the consistency of an application module is either dynamically strengthened or weakened based on its criticality needs. This requires an in-memory cache system that has a monitor/control component, session manager, scheduler and a QoS engine. These components manage the replication



of a service based on the cache statistics specific to the application. Gao *et. al.* [47] proposed a lazy master update propagation model for transaction-based systems in cloud computing. This requires a master (registry host) replication site with replica management, an update receiver, propagator and executor components. The components manage replication through synchronization and message sends (transmission, receipt and update) with different sites. The update instruction is first committed at the master node and later replicated to other sites. This model also uses an immediate update propagation variant for data access to ensure that users only access a fully replicated service. The success of this approach requires a highly available master site and geographically nearby secondary sites as the updates di-

rectly impact network performance. Boru *et. al.* [48] asserted that most cloud applications interact with database systems that may be located locally or remotely and data queries are sent to locations nearest to the user for enhanced availability and improved user experience. The model presented is intended to minimize energy consumption, bandwidth utilization and communication delays in the network. The energy consumption model aggregates power usage from computing servers as well as core, aggregation, and access switches to build an optimal profile.

3.2.3 Location-Aware

Cloud-enabled applications necessitate different modes of delivering system functionality to the end users for improved experience and this includes location awareness: applications and data should follow the users. Application replication models should dynamically consider locations of the end-users based on their density, proximity and mobility and perform desirable provisioning *in situ*. Location awareness may also be in line with the QoS requirements expected (requires continuous monitoring for agreed QoS and replication should allow for reconfiguration of available resources so that minimum QoS requirements are achievable) of an application. If the main consumers of the APDD are in Uganda, then provisioning of the application should be done on the nearest clusters. At a certain point in time, the APDD consumers may be in South Africa and this requires resource reconfiguration and provisioning to serve the new user environment.

The design of a location aware system requires request tracking based on the Internet Protocol (IP) address, monitoring components, location sensing and prediction technologies and assorted geolocation APIs so that user requests take advantage of nearby computing servers to carry out demanding tasks. This allows users to have a more contextual and fulfilling experience while drastically reducing the costs of delivering compute, network and storage resources. It should also be noted that location-aware systems may impose serious privacy issues and should be handled appropriately.

3.3 DC 3: Portability and multi-cloud migration

Portability in cloud computing can be defined as the ability for movement of applications, workloads, processes and data from one cloud environment to another with the least disruption, whether manually or automatically. The least disruption should translate to lowest possible cost, effort and time. The movement of one service, such as the one instance of the prediction microservice for automated plant diagnosis system from Cluster 2 to Cluster 1 as shown in Figure 4, should cause minimal or no downtime and should not compromise the QoS attributes tagged to overall operation of the system. As noted earlier, cloud computing offers significant benefits such as scalability, disaster recovery, mobility and cost reduction in operation of an organization's IT infrastructure. This is evidenced in the introduction of different cloud computing technologies and deployments to make it easy for organizations to embrace and adopt this new wave of handling compute, storage and network workloads. One of the pertinent issues in the adoption of cloud computing is vendor lock-in (lack of portability and interoperability across cloud platforms) where providers work with specific technologies such as tools and programming interfaces.

Given the different deployment models and the cloud service models, organizations should be able to move cloud services from one provider to another without worries of complexities and infrastructure dependence.

Bozman *et. al.* [49] identified standardized programming interface, abstraction layers and management capabilities as some of the key enablers for portability and service migration between cloud providers. A standardized programming interface includes programming toolsets to support application movement, the abstraction layers insulates users from infrastructure complexities and dependencies and the management tools provide interfaces for operational activities such as application deployment, monitoring and troubleshooting. However, adopting a standardized approach to portability is a myth as it is extremely difficult for providers to agree and adopt a unified set of standards. This requires major rework of the proprietary APIs and file formats, and this also destroys the competition spirit which has been very effective in delivering high quality cloud services for the market [50].

Most of the research work geared to support portability across different cloud environments such as mOSAIC^[15] (Open-Source API and Platform for Multiple Clouds), Open Cloud Computing Interface (OCCI) have focused on abstraction layers and management tools and container-centric solutions. Containerization allows an application to be built once, placed inside a container image or series of images for a multi-service application and running it on any host operating system that supports the containerization technology in perspective such as Docker. It should however be noted that achieving full portability out-of-the-box and application storage persistence using containers has some limitations such as no support for cross operating system support - *a containerized Linux application requires a Linux host operating system, a windows one requires a Windows operating system.* Despite this limitation, containerization is a big step in ensuring applications can run uniformly and consistently across a plethora of computing platforms or cloud environments.

3.4 DC 4: Resilience to network partitions and bandwidth constraints

Low-resource environments usually experience unreliable and intermittent Internet connections due to power failures, few or no network redundancy points and the low internet penetration hindering access. In a multi-cluster setup, this can result in network partitions where some clusters are totally unreachable for prolonged periods of time. This setup also requires additional bandwidth to support, for example, synchronisation of services across these fault domains that could be distantly located. Network partitions can lead to ruinous system failures, some of which leaves the system in a continuous error state, that capitulates into data loss, corruption, unavailability and inconsistency, broken locks and system crashes [51].

A stateful application is a data-driven application that requires persistent storage across a set of multi-cloud clusters with strict data consistency demands. To allow for this, there is a need for consensus algorithms to ensure that cluster states are globally consistent thus providing for dynamic leader election approach (the cloud cluster to act as the leader and handle writes), replication for cluster consistency and safety in ensuring that client requests are served with the correct results [52] in

^[15]<https://occi-wg.org>

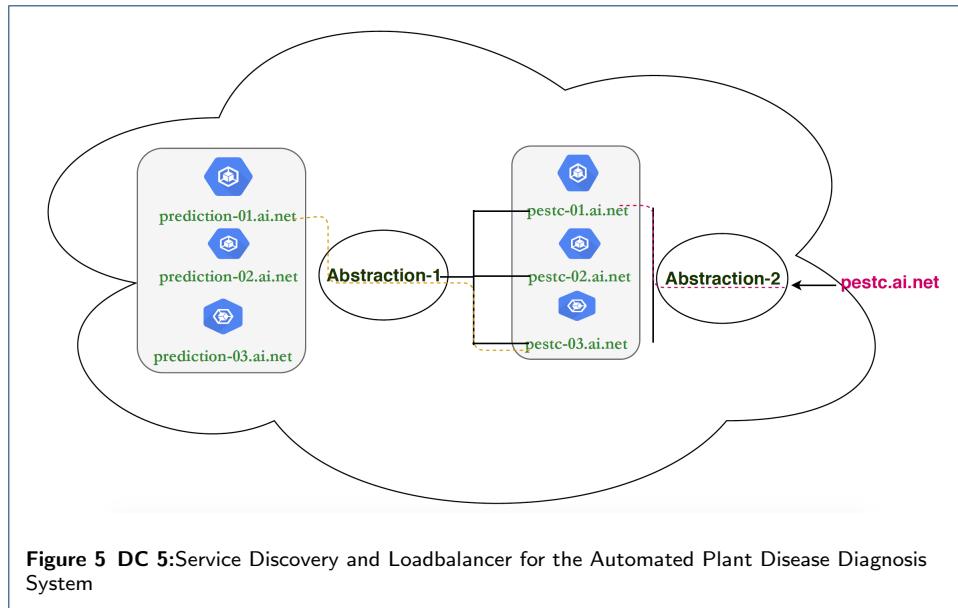
faces of complete, partial and simplex network partitions [51]. *Distributed consensus algorithms* are a well-studied research problem with a popular convergence in Raft [16], an implementation of Paxos, an easily understandable and practically implementable algorithm that guarantees a shared state among multiple servers for full operation of the system. Raft achieves this by decomposing the consensus problem into leader election, log replication and safety as independent tasks.

In resource-constrained environments, Internet traffic could be categorized into different priority classes such as sensitive(cluster replication), best-effort (service access) and undesired (other Internet traffic) to closely correspond to high, medium and low priority traffic which can then be dynamically allocated bandwidth. Assignment of optimal bandwidth to the different traffic classes will ensure that the defined QoS attributes such as availability and consistency of an application deployed on the clusters is achievable. This requires of course a fast and rigorous classification algorithm and considerable dynamic changes on the network. Another approach would be to route user requests to the closest cluster hence avoiding upstream bandwidth costs and limitations and also improved user experience. Multi-cluster support for application and data replication to achieve consistency, availability and tolerance to network partitions over Wide Area Networks (WANs) and especially geographically distant network points is still an open area of research. This requires a good and redundant connectivity between communication endpoints and a compromise in application properties such as data consistency, availability and user experience. In all this, a fair concession for near-efficient application demands should be achieved.

3.5 DC 5: Service discovery and load balancing

In a multi-cloud environment, applications may need to be scaled up by increasing application instances for improved user experience or scaled down by destroying excess instances to limit compute costs. In certain scenarios, an application may need to be moved from one cloud provider to another and rescheduled on a particular node in the new location. In the process, application settings such as its Internet Protocol (IP) addresses and Domain Name Service (DNS) attributes change and this necessitates an update to all reliant services in order to maintain application availability. Service discovery is the ability of a client component to discover healthy and available services (providers) that it can connect and communicate with. In addition, an application with multiple instances of different services spread across different clouds requires an internal and external load balancer solution to prevent network and node overload which in turn translates to optimal usage of computing resources. The internal load balancers consider microservice communication inside a cluster (cloud provider) while external load balancers consider routing of client requests to user facing endpoints of an application. In Figure 5, *Abstraction-1* and *Abstraction-2* represents the internal and external load balancers respectively. The prediction microservice and the pest identification microservice have to communicate and each has multiple instances and it is the role of the internal load balancer to route the traffic to the appropriate instances. A user visiting *pestc.ai.net* has no knowledge of what instance will respond to the request as this is abstracted by the external load balancer. *pestc-01.ai.net* at a certain point may be destroyed

[16]<https://raft.github.io>



or rescheduled in another node with possibly a new IP address and name and the abstraction layers have to update their registries so that the request/response cycle is always complete.

Kubernetes implements two options of service discovery: one based on environment variables available and the preferable DNS-based service discovery. Using environment variables, a service is identified by the IP address and port it is running on for example *PREDICTIONS SERVICE HOST = 10.233.11.2* and *PREDICTIONS SERVICE PORT = 5000*. In DNS, the services are identified by names that are specific to the cluster namespace that the application is deployed in for example, *prediction-01.dev.cluster1.local* and *prediction-02.prod.cluster2.local* to identify two instances of the prediction service deployed in the *dev* and *prod* namespaces of *Cluster 1* and *Cluster 2* respectively. The DNS option is more flexible as the environment variables are fixed for the lifetime of the service and changes require a redeployment of the application and/or service. Service discovery outside the Kubernetes cluster requires services to be exposed through NodePort, Ingress and LoadBalancer alternatives.

Most service discovery mechanisms in distributed computing and service-oriented architectures employ a centralized approach where a central server(s) maintains information about all the services that includes access credentials, protocols, versions numbers, service location and environment details. The service discovery process involves a distributed client querying the central registry for location and information of other services either using a client-side discovery (a client queries the service registry, selects an available instance and makes a request) or server-side discovery (a router acting on behalf of the client queries the service registry and forwards the request to an available instance) implementation. The popular centralized service registry/discovery solutions include Netflix's Eureka^[17], CoreOS's

^[17]<https://github.com/Netflix/eureka>

highly available etcd^[18] key-value distributed datastore, consul^[19] and Apache ZooKeeper^[20]. The major drawbacks of centralising the service registry/discovery are the introduction of points of failure, performance bottlenecks and possible network congestion. Distributing the nodes providing these services and ensuring there are multiple instances in a consistent way usually suffices. [53] [54] proposed an unstructured P2P(Peer to Peer)-enabled service discovery method for cloud environments based on Distributed Hash Tables (DHTs) with a decentralized index system. The peers maintain their own services and descriptions and a semantic-based matching rule is used to map the user requirement expressed in the query message to the desired service.

3.6 DC 6: Localised image registry

An image is an immutable file built according to instructions and can only be extended by building a layer on top of it. A container is an instance of an image with instructions on how to execute an application and operate as isolated environments with ability to interact with other containers and the host environment through well defined interfaces [55]. To fully effect containerization, an image of the application is created and pushed to a local or remote image registry through which a user, such as a DevOps engineer, can now pull and create an instance of it in a container host. A container image registry provides a storage location and distribution portal for images some with multiple versions identified by tags. Docker Hub^[21], Google's GCR^[22], Azure ACR^[23] and Amazon ECR^[24] are some examples of popular public/private image repositories.

According to [56], image pulling costs, workload network transition costs (such as bandwidth, latencies and connection limits) and energy conservation can significantly affect the scheduling and deployment of applications into container hosts. If the images are stored on the local container host, then the cost is zero otherwise extra costs shall be incurred depending on the sizes of the layers, image location and bandwidth restrictions involved in fetching the image from a remote repository. In resource-constrained settings, the image may be located thousands of kilometers from the local container hosts and this negatively impacts deployment, for example in cases of downtime where provisioning on another local cluster instance to ensure availability is required. Imagine a *10GB* image file located on <https://hub.docker.com> to be provisioned on a container host in Uganda with a dedicated upstream bandwidth limit of *2Mbps*. On average, there is a network latency of *357 ms* between <https://hub.docker.com> and Uganda. To fetch this image file, it will take close to *12 minutes* and this can have a huge negative impact on the availability QoS requirement. To reduce container schedule (download) times, images may need to be distributed across different cloud providers and located very close to the container hosts.

^[18]<https://github.com/etcd-io/etcd>

^[19]<https://www.consul.io>

^[20]<https://zookeeper.apache.org>

^[21]<https://hub.docker.com>

^[22]<https://cloud.google.com/container-registry>

^[23]<https://azure.microsoft.com/en-us/services/container-registry/>

^[24]<https://aws.amazon.com/ecr/>

3.7 DC 7: Platform monitoring and management

Monitoring is a critical and essential aspect of managing any IT infrastructure. Systems are susceptible to failure and without monitoring, it is difficult to ascertain the causes of failure and even anticipate future ones. Compared to traditional monolithic applications, monitoring of microservice applications requires intensive service reporting features especially given their distributed nature (services run as independent processes on possibly geographically different hosts) and dynamic behavior. Monitoring aids users in understanding the overall health of an application, gain insight into the performance of constituent services of an application and to ensure that APIs are available and performing as expected. The monitoring metrics divided into platform/host (CPU, RAM, threads and database connections) and application metrics (service availability, service and API endpoint latency, success of API endpoints, API endpoint response times, API request clients, errors and exceptions) should be collected at each stage of the deployment pipeline. [57] identifies four areas for microservice monitoring based on monitoring activities of information generation, processing, dissemination and presentation: *Generation and collection of monitoring data, storage, hosting and distribution of monitoring data, processing of the data to obtain platform and application metrics and presentation of need-to-know information via a dashboard to the relevant stakeholders.* In addition, a real-time monitoring component of a production-ready microservices application to detect current and imminent failures due to changes in key metrics.

A number of monitoring tools and frameworks exist but most are either native (Amazon Cloudwatch [25], Azure Monitor [26], Google Cloud's Operations Suite [27]) or virtualization type specific (such as cAdvisor [28]) or commercial (such as Datadog [29] and Dynatrace [30]). Given a plethora of monitoring options available and complexities of monitoring microservice applications, a monitoring framework should be designed to capture, report and alert stakeholders on performance and failures of an application based on critical metric data. [58] presents a framework for monitoring microservice-oriented cloud applications in heterogeneous virtualization environments. It is composed of mainly two components: a monitoring agent (*a cloud platform-independent software component that collects information from a microservice*) and a monitoring manager (*a software component that receives monitoring information from agents in heterogeneous cloud environments*).

4 Crane Cloud: an implementation of a resilient multi-cloud service layer

This section presents Crane Cloud, the first-cut prototype instantiation of design properties and considerations for a multi-cloud service layer presented in Section 3 and summarised in Table 1. Motivated by the unique requirements for low-resource settings in Subsection 2.2, the Crane Cloud is an open source project that attempts to encapsulate the intricacies of operating heterogeneous application clusters into a highly available unified platform for management and monitoring of

[25]<https://aws.amazon.com/cloudwatch/>

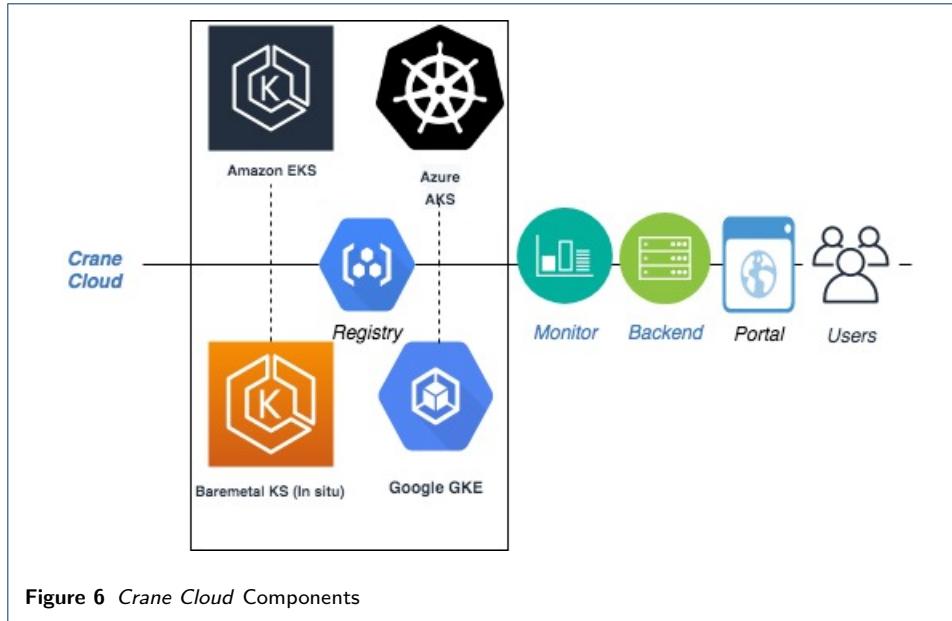
[26]<https://azure.microsoft.com/en-us/services/monitor/>

[27]<https://cloud.google.com/products/operations>

[28]<https://github.com/google/cadvisor>

[29]<https://www.datadoghq.com>

[30]<https://www.dynatrace.com>



a microservice application lifecycle. The target users of the platform include developers, researchers, students, and startups located in resource constrained environments. The public Github repository for the project is available on Github <https://github.com/crane-cloud/>.

Table 1 Design considerations for a resilient multi-cloud service model

| ID | Design Consideration |
|------|--|
| DC 1 | Containerisation and orchestration of applications |
| DC 2 | Service scheduling and replication |
| DC 3 | Portability and multi-cloud migration |
| DC 4 | Resilience to network partitions and bandwidth |
| DC 5 | Service discovery and load balancing |
| DC 6 | Localised image registry |
| DC 7 | Platform monitoring and management |

4.1 Architecture and Overview

Crane Cloud is an open source multi-cloud service layer designed to enable developers, organizations and researchers to set up reliable cloud-services in low resource setting. The Crane Cloud software layer was conceived to address the key hurdles of operating a cloud-service platform in resource constrained environments characterised by challenges identified in Section 2.2. Its main ingredients include resilience to network partitions, support for microservice orchestration, support for migration of services between private and public clouds to avoid vendor lock-in issues, seamless downtime and network traffic load distribution, monitoring metrics, and tools for transforming existing non-cloud compliant services into compliant cloud services. The multi-cloud service layer has five components (managed portal, authentication and authorization, monitoring and billing, local registry and the backend service) purposely designed taking into consideration features described in Section 3 and are shown in Figure 6.

4.1.1 Multi-cloud Cluster Support

Crane Cloud enables harmonization of clusters from different cloud providers (public or private) and bare metal environments. It provides for easy migration, replication and loadbalancing of services across different clusters and cloud providers to ensure high availability and improve the general user experience. The cloud providers are chosen based on location, API service offering and costs of running workloads in their data centers.

4.1.2 Managed Portal

The managed portal provides an interface for access to the rest of the abstracted Crane Cloud multi-cloud components. Developers can deploy and access their application services, monitor resources and running services, manage users and access the local private registry. Administrators can monitor the different clusters and nodes, view project details such as resource usage and its users, add or remove clusters, manage projects and accounts in the system and ensure that the clustered environment is performing optimally viz-a-viz the running services. The portal is a window to Crane Cloud features for management of resources in a clustered computing environment. As a developer, infrastructure setup complexities and application deployment intricacies are eliminated and focus shifts to software functionality. As a service consumer (user), the availability of a service and user experience regardless of location and underlying technologies is paramount.

4.1.3 Authentication and Authorization

Crane Cloud uses the concept of user projects to closely map with a cluster namespace. A namespace is a logical environment that supports resource management for users working in a team or across teams. In Crane Cloud, access to the cluster resources requires valid credentials and the right privileges mediated by an API. This involves creation of projects and accounts that correspond to specific privilege levels in the cluster. This ensures that users can only access what may be required to perform functions within the cluster without affecting other projects. Two types of accounts are supported: Project user accounts that are managed outside the cluster and service accounts inside the cluster. The service accounts are directly used to manage resources in the cluster while the project user accounts are mapped to service accounts but usually with defined privileges/roles over a namespace.

4.1.4 Platform and Service Monitoring Support

Monitoring is integral to the overall operation of Crane Cloud in terms of infrastructure and the distributed services hosted to ensure the QoS attributes are in check. The infrastructure includes the nodes while the services are the client applications and supporting tools. Monitoring coupled with an alert system also ensures that possible failures are averted early on before turning catastrophic. More specifically, the cluster monitoring involves the state of the cluster (collection of nodes) which is a constituent of node resource utilization parameters such as network bandwidth, disk utilization, CPU, and memory utilization while the service metrics include CPU, network, and memory usage irrespective of the nodes they are running on.

4.1.5 Local Registry

The localized and replicated registry provides a platform through which users can easily upload, store and deploy their applications fast. The registry also provides trust signing and vulnerability scanning of container images. This significantly reduces the costs of upload and download of container images from public registries. The registry is locally available on all the clusters and container images are replicated on all. A local registry also ensures that administrators have more control over it.

4.1.6 Crane Backend

The Crane Backend is the heartbeat of Crane Cloud providing abstractions and hooks for a number of services. The backend ensures that applications are appropriately scheduled on a cluster(s) or nodes, creation of service endpoints for communication between parts of an application and also ensure users can access the application, management of volumes for deployment of stateful applications, location-aware and user preference deployment of applications. It also provides endpoints for management of container images through the registry API and monitoring of deployed services.

4.2 Implementation

Crane Cloud is implemented using a combination of tools ranging from the design of the managed portal to the backend and from monitoring, registry and persistent volume management to the container orchestration. In most cases, open source solutions were experimented and used as much as possible.

4.2.1 Container Orchestration

Container orchestration tools are used to automate the deployment, management, scaling, and networking of containerized applications. These tools provide an abstraction layer between pools of resources and the application containers that run on those resources. Kubernetes, Docker swarm and Apache Mesos are the most popular tools with the former taking a fair share of the cloud-native market. With an impressively large community and functionality, backed by Cloud Native Computing Foundation (CNCF) and its open source nature, Crane Cloud uses Kubernetes for container orchestration. We considered five factors in selecting the most viable tool for our setup: Community (Support), Open Source, Scalability & Flexibility, Fault tolerance and Monitoring support. As shown in the comparison Table 2, Kubernetes is an open source project that impressively commands the cloud-native market with an adoption rate of 50% in the past 6 months and 87% market penetration supporting application scalability, fault tolerance and has inbuilt monitoring and logging tools and hence was the preferred choice for container orchestration implementation of Crane Cloud. Additionally, Kubernetes provides automated scheduling of applications, self healing capabilities, automated roll-out and rollback, service load-balancing and a higher density of resource utilization.

Table 2 A comparison of major container orchestration implementation tools

| Feature | Apache Mesos | Docker Swarm | Kubernetes |
|--------------------------------------|---|--|---|
| <i>Community</i> | Has a medium active community with 18,100 Github commits and 300 contributors [31] | Started in 2014, Docker Swarm has a relatively smaller community with 3,570 Github commits and 178 contributors [32] | Large community with over 100,000 Github commits and 3,500 contributors [33] making it one of the most active open source projects |
| <i>Open Source</i> | Yes | Yes but with an enterprise edition that detriments the open source version | Yes |
| <i>Scalability & Flexibility</i> | Automatic scaling but may require definitions in the application. | Manual scaling | Automatic scaling based on resource utilization |
| <i>Fault tolerance</i> | Yes | Yes | Yes |
| <i>Monitoring</i> | Has a diagnostic utility for health and other metrics but requires queries and aggregation through APIs | Uses basic out-of-the-box tools and supports other 3rd-party logging and monitoring tools | Uses inbuilt tools for logging and monitoring with support for third party integrations to keep track of logs and other performance metrics |

4.2.2 Developer Tools

The managed portal (frontend) was implemented using React.js [34], a fast, scalable, and simple JavaScript library for building user interfaces created and open-sourced by Facebook. It uses the component-based architecture and declarative approach hence simplifying the debugging process. It allows creation of simple reusable and stateful components which can be composed to build complex user interfaces. The Crane Cloud backend was implemented as a REpresentational State Transfer (REST) API, using Python Flask [35], for ease of use and integration while ensuring effective maintenance. Python Flask was used because it is simple, flexible and lightweight and now considered as one of the most popular Python web application frameworks by the programming community. PostgreSQL [36] is an open source object-relational database management system that Crane Cloud uses to maintain state for its internal workings. Applying multi-version concurrency control (MVCC) which allows several concurrent read/write operations, PostgreSQL can handle multiple tasks simultaneously and efficiently. In addition, PostgreSQL is SQL standards-compliant, highly programmable and extensible by many third-party tools and libraries.

4.2.3 Image Registry

In implementation of the image registry, Crane Cloud considered open source extensible tools that can secure, scan and sign container images and also support replication across clusters. Harbor [37] perfectly fitted into the picture, providing an extensible API that the backend service would easily consume. Harbor delivers a consistent experience across multiple clouds and works best for environments that

[34]<https://reactjs.org>

[35]<https://github.com/pallets/flask/>

[36]<https://www.postgresql.org>

[37]<https://goharbor.io>

may not want to rely on public registries but rather a private one packaged as an add-on. Harbor additionally provides features such as access control on registry images, image vulnerability scanners, image storage and replication using a clustering mechanism. Crane Cloud is a multi-cloud service layer that can work with cloud providers in different regions and availability zones and a zonal scalable registry with a replication service is cardinal.

4.2.4 Stateful Applications

Containerization technologies were originally designed to support stateless applications but considerable efforts have now been made to also support stateful applications owing to community adoption and contribution. This enables organizations to work with data-driven and legacy applications while leveraging the portability, scalability and highly available features of containers. Traditionally, Kubernetes used to provide support for manual attachment of cloud-backed storage to applications limiting usage outside the cloud provider but cloud native storage solutions have now been advanced. Crane Cloud uses OpenEBS [38], an open source Container Attached Storage (CAS) solution developed using the microservice architecture. Distributed, monolithic or streaming, OpenEBS allows deployment of storage technologies and optimizations appropriate to an application type using different storage engines. Additionally, OpenEBS is a multi-cloud storage solution that shares the same philosophy of Crane Cloud borderless computing.

4.2.5 Monitoring

Crane Cloud uses a combination of the inbuilt Kubernetes Metrics server, Prometheus [39] and Grafana [40]. The Metrics server is a cluster-wide aggregator of container resource metrics such as container CPU and memory usage exposed on each cluster node and available through the Metrics API. Prometheus is an open source time series database optimized to store monitoring metrics using a periodic data pull model and provides API endpoints. Providing basic visualizations and dashboards, Prometheus can be deployed alongside dedicated visualization and dashboard solutions such as Grafana. Using the Prometheus APIs, visualizations can be generated and customized for users.

4.2.6 Mapping of Crane Cloud, design considerations and challenges addressed

The Table 3 shows a mapping between components of Crane Cloud, the design considerations and low-resource computing environment challenges addressed. The registry, monitor and backend are continuously being refined for a seamless user experience. It should also be noted that resistance to network partitions and bandwidth constraints as a design consideration and how Crane Cloud can practically implement this is an ongoing research area.

[38]<https://openebs.io>

[39]<https://prometheus.io>

[40]<https://grafana.com>

Table 3 Mapping of Crane Cloud components and design considerations

| | Crane Cloud Component | Design Consideration | Requirements and Challenges |
|---|----------------------------------|--|---|
| 1 | Multi-cloud Cluster Support | Containerization and container orchestration engines. | Vendor lock-in, Poor quality of service |
| 2 | Managed Portal | Manual service deployment, scheduling and replication, monitoring and alerts portal | Data jurisdiction restrictions |
| 3 | Authentication and Authorization | Portability and multi-cloud migration | Cloud resource security |
| 4 | Crane Backend | Application containerization and orchestration, Ingress, Service discovery, scheduling, replication and load-balancing for QoS, location-aware and user defined policy management. | Frequent Internet partitions and bandwidth constraints, Poor quality of service |
| 5 | Crane Registry | Localized image registry | Poor quality of service |
| 6 | Crane Monitor | Platform monitoring, alerts and management | Poor quality of service |

5 Experiments and Results

To demonstrate the utility of the Crane Cloud platform, we deployed the rapid plant disease diagnosis (*mcrops*) microservice of APDD on the Crane Cloud platform. As introduced in Section 2, *mcrops* provides machine learning-based diagnostic tools to detect viral crop diseases in cassava plant using mobile and web-based technologies. In brief, the users upload images of suspected infected cassava root tubers through a mobile application or web browser and *mcrops* computes the Cassava Brown Streak Disease (CBSD) score to indicate the disease presence levels. The users can perform single uploads and/or multiple uploads of the image data. As shown in Section 2, *mcrops* is deployed as a monolith and for Crane Cloud deployment support, the application was containerized using Docker. The containerization process involved access to the *mcrops* source code, writing of the Dockerfile, iterative building and tagging of the Docker image and pushing it to the Crane Cloud image registry <https://registry.cranecloud.io/>. On the Crane Cloud portal, the deployment involved providing the application details such as name, container port and the docker image reference (as shown in Figure 7) from which *mcrops* is run and ingress resources for external access are subsequently created and availed for use.

5.1 Experiment setup

The purpose of the experiment was to evaluate the performance and user experience of the *mcrops* when deployed in a public cloud compared to the deployment on the Crane Cloud platform. We used Apache JMeter [41], a popular performance testing tool. Specifically we used JMeter tool to measure the response times of the two application setups against uploaded images over different mobile/wireless connections for a user situated in a bandwidth constrained setting. These testing settings and environments are a representative of the realities that developers and users work in in the low-resource settings. On JMeter, four thread groups representing the image upload settings and two samplers (the two application setups: *unet.mcrops.org* and *mcrops.cranecloud.io*) were used as shown in Tables 4 and 5. In **A1** and **B1**, we considered the simplest scenario of a user uploading a single image but varied the number of times to 10 and 20 respectively as represented by the loops. The

[41]<https://jmeter.apache.org/>

Deploy an app

| |
|---|
| mcrops |
| 1 |
| registry.cranecloud.io/mwotila/unet:0.5 |
| <input type="checkbox"/> Private Image |
| Entry Command (i) |
| 8000 |

Environment Variables (i)

| Name | Value | Remove |
|------|--------|---------------------------------------|
| name | mcrops | X |
| Name | Value | ADD |

CANCEL
DEPLOY

Figure 7 Deployment of mcrops on Crane Cloud

ramp-up time represents the seconds between successive user requests and we set this as same as the loops. In **A2** and **B2**, we randomly set the number of images to 10 to test a multi-upload use case as this feature is supported by the application. As the number of images and/or loops was increased, the response times were assumed to also rise hence the need to increase the ramp-up time. In all scenarios, 1 user was simulated because this is the typical field use case where a device has a non-shared connection to the next transmitting network device, for example an extension worker or rural smallholder farmer in a garden. The connection speeds for 2G, 3G, 4G and a local WiFi access point were determined by performing tests using the online Internet speed test tool *speedtest.net* [42] for upload and download speeds and the average computed and recorded in Table 5. The connections under use are motivated by real-world setups of resource constrained environments with no ideal and consistent network connectivity. Overall, we wanted to assess the utility of a platform like Crane Cloud for end-users situated in these environments. In this example, the students and researchers of the Computer Science Department at the University.

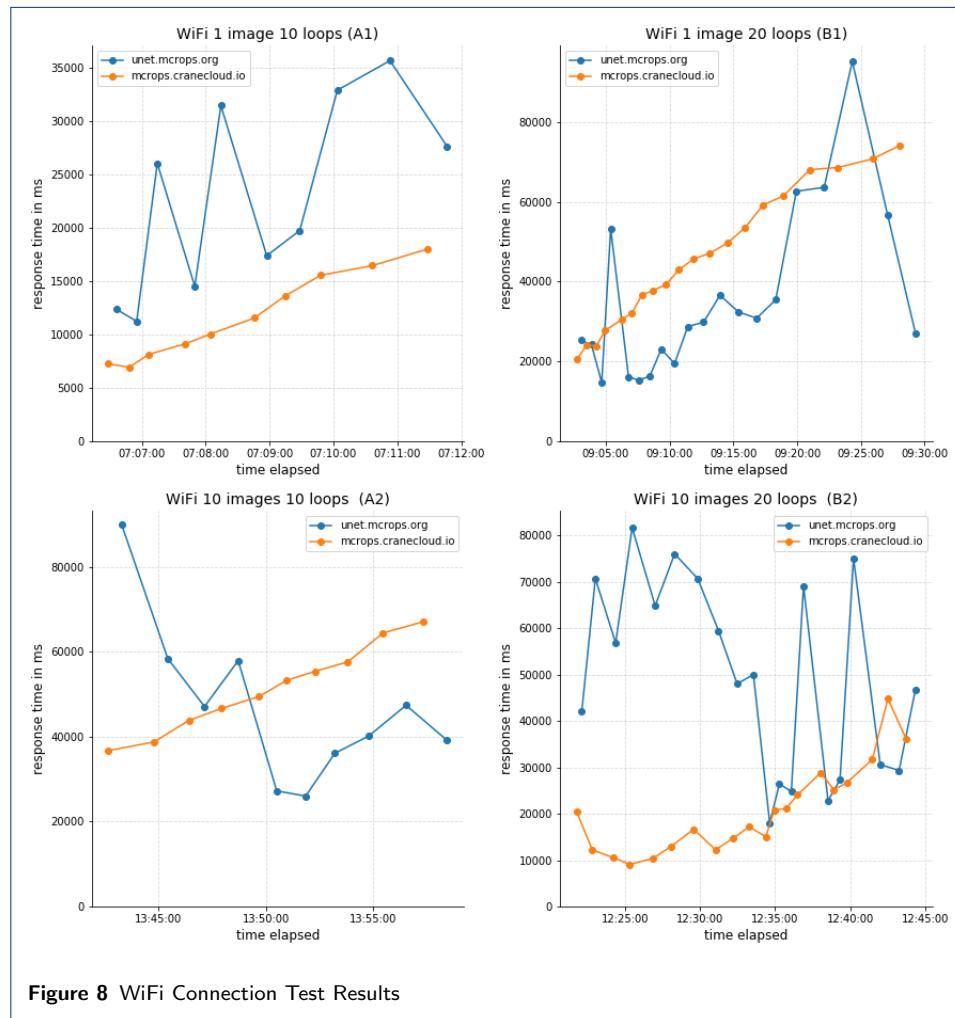
[42]<http://speedtest.net/>

Table 4 JMeter Test Settings

| Code/Setting | Images | Users | Loops | Ramp-Up time |
|--------------|--------|-------|-------|--------------|
| A1 | 1 | 1 | 10 | 10 |
| B1 | 1 | 1 | 20 | 20 |
| A2 | 10 | 1 | 10 | 10 |
| B2 | 10 | 1 | 20 | 20 |

Table 5 Experimental setup for mcrops tests

| Connection/Sampler | unet.mcrops.org | mcrops.cranecloud.io |
|--------------------|-----------------|----------------------|
| 2G (45Kbps) | A1,B1 | |
| 3G (4Mbps) | | A2,B2 |
| 4G (8Mbps) | | |
| WiFi (10Mbps) | | A1,B1 |
| | | A2,B2 |

**Figure 8** WiFi Connection Test Results

5.2 Results and Discussion

The results of the experiments are presented in Figures 8, 9, 10 and 11 and Table 6. In the analysis, the median and mean (average) times were used to conclude on which experiment had a shorter response time using the different connections for both *unet.mcrops.org* and *mcrops.cranecloud.io*. The median was given priority in cases where its difference compared to the mean is large since it is not inflated by the existence of outliers in the data collected.

WiFi connection test results. In the WiFi connection setup, the response times for experiment **B1** and **A2** were lower in *unet.mcrops.org* compared to *mcrops.cranecloud.io* using the median and mean times as shown in Table 6. For experiment **A1**, the results indicate that *unet.mcrops.org* has its lowest response time at *11.23 seconds* with an average of *22.9 seconds* compared to *mcrops.cranecloud.io* at *6.91 seconds* with an average of *11.66 seconds*. The mean times also indicate that responses are better (lower) in *mcrops.cranecloud.io* compared to *unet.mcrops.org*. The experiment **B2** response times are lower in *mcrops.cranecloud.io* compared to *unet.mcrops.org* using both median and mean. The erratic behaviour of the graphs for *unet.mcrops.org* in Figure 8 is partly attributed to the high number of network hops, packet losses (as shown by the completion rates in Table 7 and connection variations from the packet sources. It should also be noted that the completion rates for *mcrops.cranecloud.io* under WiFi was at 100% compared to *unet.mcrops.org* at 95%.

4G connection test results. Using the 4G connection, the response times for experiment **B1** and **B2** was lower in *unet.mcrops.org* compared to *mcrops.cranecloud.io* using both the median and mean as shown in Figure 9 and Table 6. In experiment **B1**, *unet.mcrops.org* had *28.01 seconds* and *29.69 seconds* while *mcrops.cranecloud.io* had *47.29 seconds* and *49.06 seconds* for the median and mean times respectively. For experiment **B2**, the results followed a similar pattern. However, experiments **A1** and **A2** performed significantly better under *mcrops.cranecloud.io* for example; in **A2**, *unet.mcrops.org* had *24.43 seconds* and *33.29 seconds* while *mcrops.cranecloud.io* had *10.81 seconds* and *10.999 seconds* for the median and mean times respectively. The completion rates for *mcrops.cranecloud.io* under WiFi was at 100% compared to *unet.mcrops.org* at 88.75%.

3G connection test results. Under the 3G connection, the average response time for experiments **A1**, **A2** and **B2** is lower in *mcrops.cranecloud.io* compared to *unet.mcrops.org* as shown in Figure 10 and Table 6. In experiment **A2**, for example, *unet.mcrops.org* had *64.37 seconds* and *60.05 seconds* while *mcrops.cranecloud.io* had *27.18 seconds* and *30.48 seconds* for the median and mean times respectively. This was quite significant as the response times for *mcrops.cranecloud.io* were half and very similar patterns for **A1** and **B2**. Despite the presence of an outlier in experiment **A1**, *mcrops.cranecloud.io* still performed better. However, for experiment **B1**, the response time was much lower at an average of *30.97 seconds* using *unet.mcrops.org* compared to *107.46 seconds* for *mcrops.cranecloud.io*. This could be attributed to a network congestion or server load at execution time especially given the positive results from **A1**, **A2** and **B2**. In general, the completion rate for *mcrops.cranecloud.io* under 3G was at 100% compared to *unet.mcrops.org* at 85%.

2G connection test results. Under the 2G connection, only two experiments (**A1** and **B1**) were completed successfully as shown in Figure 11 and Table 6. These experiments involved single image data with the loops varied. *mcrops.cranecloud.io* in both experiments had lower response times compared to *unet.mcrops.org* using

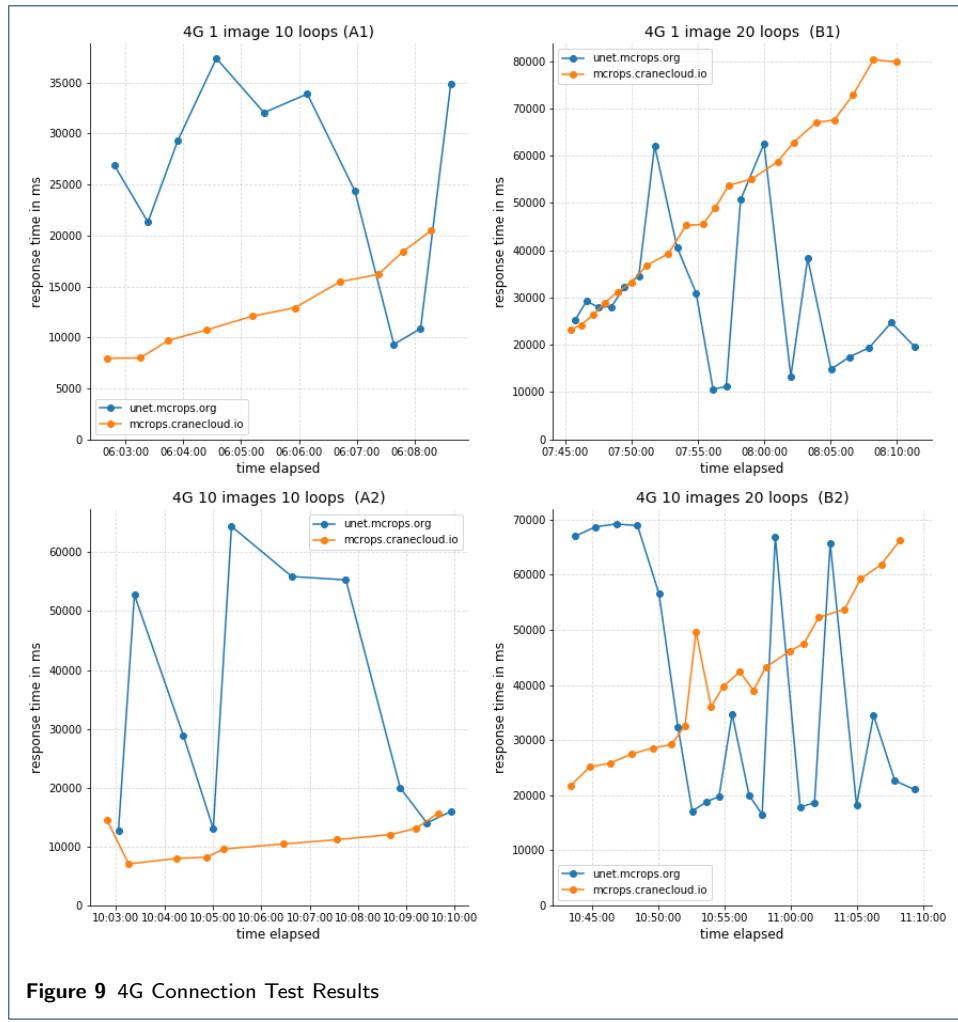
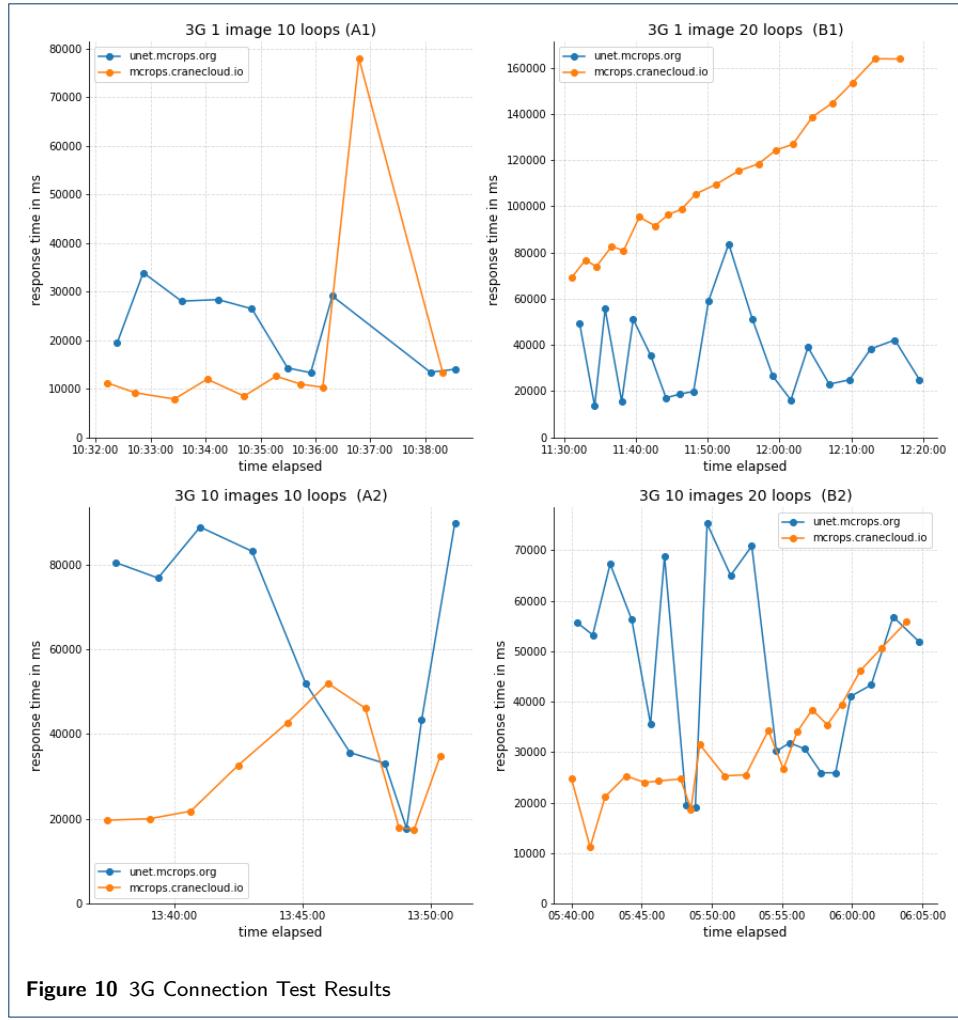


Figure 9 4G Connection Test Results

both the median and mean. In **A1**, *unet.mcrops.org* had 57.33 seconds and 56.61 seconds while *mcrops.cranecloud.io* had 38.81 seconds and 48.48 seconds for the median and mean times respectively. The spikes in the graphs are attributed to network unreliability especially under constrained capacities supported by 2G.

The computation of the CBSD score is a resource-intensive task as shown by the response times recorded in all the scenarios. The increasing response times are attributed to the ramp-up time where new requests are generated at specific intervals before some computations are concluded. From the results, it is also clear the *mcrops.cranecloud.io* is more consistent in the increasing response times and this is attributed to the completion rates of the execution as shown in Table 7 and Figure 12. In all instances, *mcrops.cranecloud.io* has a 100% completion rate compared to *unet.mcrops.org* at 89.64%. In cases where the public cloud hosted instance performs better, server errors such as 502 (Bad Gateway) and 504 (Gateway Timeout) were recorded. As expected and shown by the Internet speed results, the WiFi connection performs much better compared to the rest of the connectivity options.

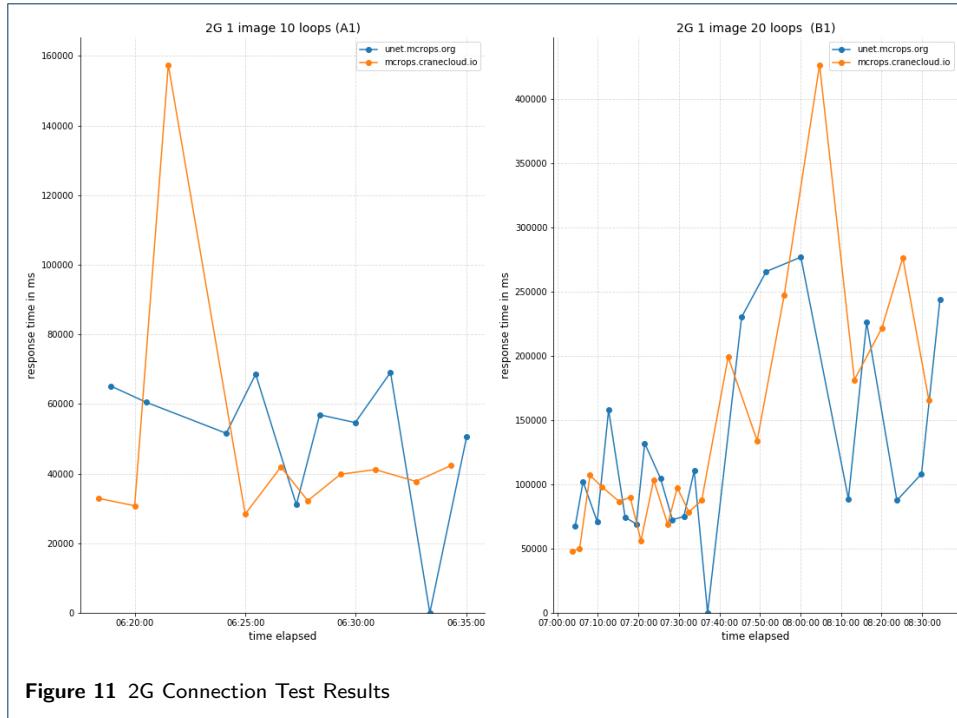
The behaviour of the trends as observed in the graphs is due to the Internet speed variations during the course of the experiments. For example, we noticed that experiments done in the morning provide better response times compared to

**Figure 10** 3G Connection Test Results

the later hours. This is explained by typical network usage patterns over a 24-hour period as shown in Figure 13.

6 Related Work

Van den Bossche *et. al.* addressed the challenge of cost-efficiently scheduling deadline constrained batch type applications on IaaS (virtual machines) hybrid clouds using custom heuristics. Properties such as high availability, scalability, fault-tolerance and monitoring are not discussed and the use of virtual machines may not be the most cost-effective approach to running application workloads [60]. Filip *et. al.* proposed a solution that considers a finite catalog of primitive microservices and designs a hybrid scheduling algorithm that matches tasks to resources based on task history and availability of resources [61]. In addition to benefits of using a microservice architecture, the paper asserted that costs can further be reduced by placing data closer to processing points based on user density. Mussig *et. al.* describes the concept of a high scalable microservice infrastructure using custom metrics in addition to commonly used ones such as CPU and RAM [62]. In this paper, custom metrics such as service utilization for scaling, load balancing and load prediction often results in better business-alignment of the scaling behavior as well as cost reduction. Guerrero



et. al. presented an optimization approach to reduce service cost, microservices repair time, and microservices network latency overhead in the orchestration process of containers in multi-cloud environments using the scale level of the microservices and their allocation in the virtual machines, the provider and virtual machine type selection and the number of virtual machines [63]. Sousa *et. al.* developed a framework for automated deployment of microservices applications in multi-cloud environments with containers. The application's multi-cloud requirements are defined and a systematic method for obtaining proper configurations that comply with the application's requirements and the cloud providers' constraints is adopted [64]. Rancher [43] and D2C [44] are two examples of container management platforms that simplify the process of operating container clusters on any cloud or infrastructure platform. The downside with these platforms is in addressing application requirements such as data processing and storage restrictions and the distinctive requirements for resource constrained settings. As more established cloud providers such as Microsoft, Google, Oracle and Amazon move towards hosted cloud-native platforms such as Kubernetes for easier configuration and management, the vendor-lockin issues are expected to exacerbate especially with no plans of integration tools or APIs. In summary, there is no standardized solution for implementation and operation of a multi-cloud service layer but rather blocks that independently address the design considerations in Section 3.

7 Conclusion and Future Work

In this paper, we presented Crane Cloud - a resilient multi-cloud service layer for resource constrained environments using Kubernetes and assorted management tools.

[43]<http://rancher.com/>

[44]<http://d2c.io>

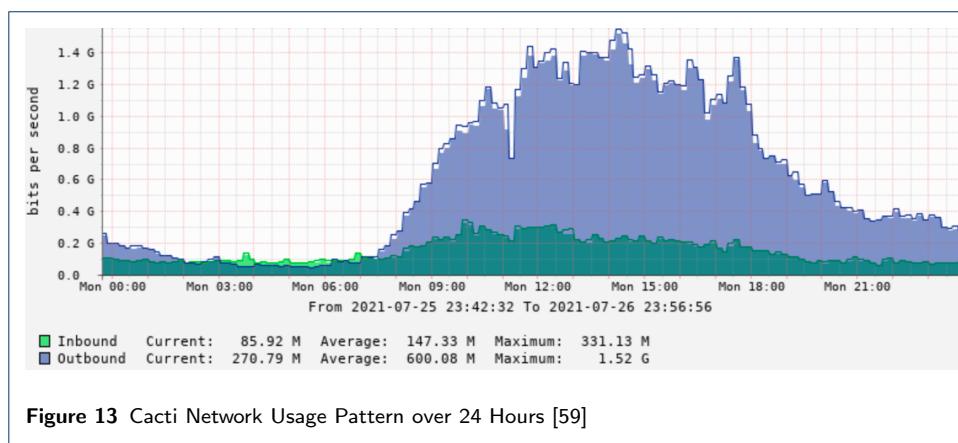
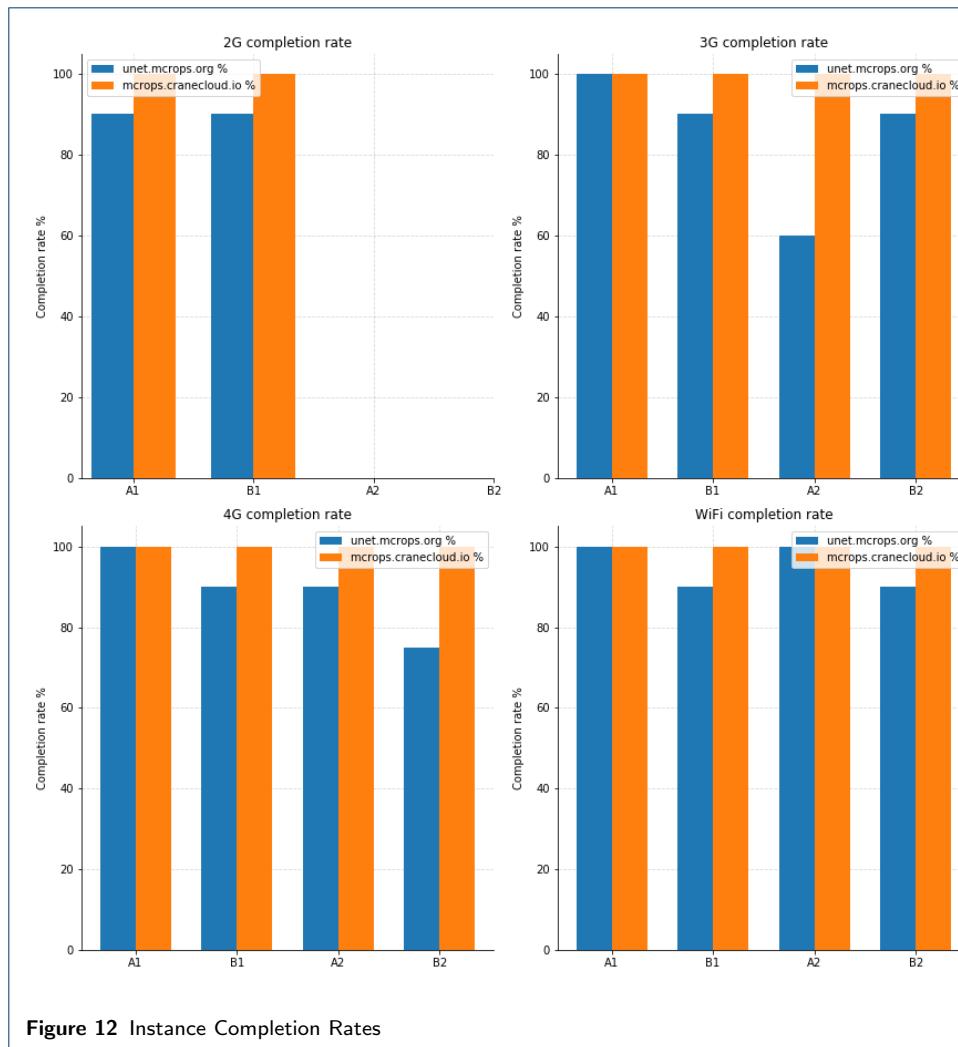
Table 6 Comparison of the Execution Response Times (*exp = Experiment, min = Minimum, med = Median, max = Maximum, avg = Average*)

| Sampler -> | | unet.mcrops.org (seconds) | | | | mcrops.cranecloud.io (seconds) | | | |
|------------|-----|---------------------------|--------|--------|--------|--------------------------------|--------|--------|--------|
| Connection | exp | min | med | max | avg | min | med | max | avg |
| 2G | A1 | 31.09 | 57.33 | 69.10 | 56.61 | 28.38 | 38.81 | 157.48 | 48.48 |
| | B1 | 67.34 | 106.29 | 303.91 | 143.39 | 47.70 | 100.76 | 426.60 | 141.19 |
| | A2 | - | - | - | - | - | - | - | - |
| | B2 | - | - | - | - | - | - | - | - |
| 3G | A1 | 13.31 | 22.94 | 33.91 | 22.03 | 7.88 | 11.11 | 78.03 | 17.41 |
| | B1 | 13.51 | 30.97 | 59.10 | 35.26 | 69.03 | 107.46 | 163.94 | 111.52 |
| | A2 | 17.69 | 64.37 | 89.80 | 60.05 | 17.52 | 27.18 | 51.97 | 30.48 |
| | B2 | 19.04 | 47.56 | 75.28 | 46.19 | 11.16 | 26.06 | 55.82 | 30.86 |
| 4G | A1 | 9.32 | 28.10 | 37.37 | 26.03 | 7.95 | 12.50 | 20.48 | 13.19 |
| | B1 | 10.65 | 28.01 | 62.56 | 29.69 | 23.16 | 47.29 | 80.31 | 49.06 |
| | A2 | 12.72 | 24.43 | 64.39 | 33.29 | 7.06 | 10.81 | 15.68 | 10.99 |
| | B2 | 16.5 | 27.53 | 69.18 | 37.75 | 21.69 | 41.05 | 66.18 | 41.37 |
| WiFi | A1 | 11.23 | 22.90 | 35.70 | 22.90 | 6.91 | 10.78 | 18.03 | 11.66 |
| | B1 | 14.62 | 29.24 | 95.23 | 35.31 | 20.42 | 44.30 | 74.13 | 45.67 |
| | A2 | 26.01 | 43.61 | 89.97 | 46.93 | 36.74 | 51.36 | 67.07 | 51.31 |
| | B2 | 18.04 | 49.04 | 81.70 | 49.55 | 9.09 | 18.80 | 44.81 | 20.55 |

Table 7 Execution Completion Rates (%)

| Connection/Sampler | Experiment | unet.mcrops.org (%) | mcrops.cranecloud.io (%) |
|--------------------|------------|---------------------|--------------------------|
| 2G | A1 | 90 | 100 |
| | B1 | 90 | 100 |
| | A2 | - | - |
| | B2 | - | - |
| 3G | A1 | 100 | 100 |
| | B1 | 90 | 100 |
| | A2 | 60 | 100 |
| | B2 | 90 | 100 |
| 4G | A1 | 100 | 100 |
| | B1 | 90 | 100 |
| | A2 | 90 | 100 |
| | B2 | 75 | 100 |
| WiFi | A1 | 100 | 100 |
| | B1 | 90 | 100 |
| | A2 | 100 | 100 |
| | B2 | 90 | 100 |

We highlighted the characteristics of a resource constrained environment that includes poor Internet connectivity, frequent Internet partitions and data center power cuts ultimately resulting in poor user experience or even service unavailability. Based on these challenges, we enumerated a number of design considerations and properties for a resilient multi-cloud service layer that would form the foundation for Crane Cloud. From easing terminal complexities of operating a cloud service, desirable scaling, availability, migration and loadbalancing to platform monitoring, Crane Cloud tries to provide an all-inclusive solution that best fits the resource constrained compute environment. As much as Crane Cloud directs implementations for the subject environment, it should be noted there are many moving parts some of which are under active development and research. The bandwidth constraints, for example, may require consensus algorithms for better handling of network splits but there are always trade-offs that should be considered. Management of persistent storage, replication and fault tolerance across geographically distant clusters accessible via Wide Area Networks (WANs) is still an open research area. There are always penalties introduced on the network especially with application data that has to traverse bottleneck links to maintain real-time application consistency. Further future work includes analysing and optimizing the scheduling processes for applications in production Crane Cloud clusters. In bare metal clusters, service load



balancing usually requires more investment in the network infrastructure which is not an option in a low resource setting and should be further explored.

Acknowledgements

The authors would like to acknowledge all persons who have contributed to the thought process of Crane Cloud and resources in funding or in-kind to bring the project to life. We thank Paul Maritz for the support and advice in the conceptualisation of the Crane Cloud project. We thank the Crane Cloud team for feedback and input. The authors would also like to acknowledge support from the Government of Uganda through the Makerere University Research Innovation Fund (RIF).

Availability of data and materials

The Crane Cloud source code is available on Github <https://github.com/crane-cloud/>.

Ethics approval and consent to participate

This work did not involve human subjects.

Competing interests

The authors declare that they have no competing interests.

Consent for publication

All authors consent to the publication of this work.

Authors' contributions

Engineer Bainomugisha contributed in conceptualisation, synthesis, supervision, writing, and review.
Alex Mwotil contributed in conceptualisation, implementation, and writing.

Author details

¹Department of Computer Science, School of Computing and Informatics Technology, College of Computing and Information Sciences, Makerere University, Kampala, Uganda. ²Department of Networks, School of Computing and Informatics Technology, College of Computing and Information Sciences, Makerere University, Kampala, Uganda.

References

- Alabadi, M.M.: Cloud computing for education and learning: Education and learning as a service (elaas). In: 2011 14th International Conference on Interactive Collaborative Learning, pp. 589–594 (2011). IEEE
- Sultan, N.: Cloud computing for education: A new dawn? International Journal of Information Management **30**(2), 109–116 (2010)
- Hashem, I.A.T., Yaqoob, I., Anuar, N.B., Mokhtar, S., Gani, A., Khan, S.U.: The rise of “big data” on cloud computing: Review and open research issues. Information systems **47**, 98–115 (2015)
- Nkosi, M., Mekuria, F.: Cloud computing for enhanced mobile health applications. In: 2010 IEEE Second International Conference on Cloud Computing Technology and Science, pp. 629–633 (2010). IEEE
- Rolim, C.O., Koch, F.L., Westphall, C.B., Werner, J., Fracalossi, A., Salvador, G.S.: A cloud computing solution for patient’s data collection in health care institutions. In: 2010 Second International Conference on eHealth, Telemedicine, and Social Medicine, pp. 95–99 (2010). IEEE
- Kshetri, N.: Cloud computing in developing economies. Computer **43**(10), 47–55 (2010)
- Zhang, W., Chen, Q.: From e-government to c-government via cloud computing. In: 2010 International Conference on E-Business and E-Government, pp. 679–682 (2010). IEEE
- Liu, F., Tong, J., Mao, J., Bohn, R., Messina, J., Badger, L., Leaf, D.: Nist cloud computing reference architecture. NIST special publication **500**(2011), 292 (2011)
- Bhardwaj, S., Jain, L., Jain, S.: Cloud computing: A study of infrastructure as a service (iaas). International Journal of engineering and information Technology **2**(1), 60–63 (2010)
- Dawoud, W., Takouna, I., Meinel, C.: Infrastructure as a service security: Challenges and solutions. In: 2010 the 7th International Conference on Informatics and Systems (INFOS), pp. 1–8 (2010). IEEE
- Manvi, S.S., Shyam, G.K.: Resource management for infrastructure as a service (iaas) in cloud computing: A survey. Journal of network and computer applications **41**, 424–440 (2014)
- Tao, F., Zhang, L., Venkatesh, V., Luo, Y., Cheng, Y.: Cloud manufacturing: a computing and service-oriented manufacturing model. Proceedings of the Institution of Mechanical Engineers, Part B: Journal of Engineering Manufacture **225**(10), 1969–1976 (2011)
- Vaquero, L.M., Rodero-Merino, L., Caceres, J., Lindner, M.: A break in the clouds: towards a cloud definition. ACM New York, NY, USA (2008)
- HAQ, S.U.: Introduction to Monolithic Architecture and MicroServices Architecture. 2018 (2019)
- Richardson, C.: Monolithic architecture. Microservices. io,[Online]. Available: <http://microservices.io/patterns/monolithic.html> [Accessed 3 October 2019] (2017)
- Anton, K.: Monolithic vs. microservices architecture. microservices.com,[Online]. Available: <https://articles.microservices.com/monolithic-vs-microservices-architecture-5c4848858f59> [Accessed 9 October 2019] (2015)
- Villamizar, M., Garcés, O., Castro, H., Verano, M., Salamanca, L., Casallas, R., Gil, S.: Evaluating the monolithic and the microservice architecture pattern to deploy web applications in the cloud. In: 2015 10th Computing Colombian Conference (10CCC), pp. 583–590 (2015). IEEE
- Mwebaze, E., Biehl, M.: Prototype-based classification for image analysis and its application to crop disease diagnosis. In: Advances in Self-Organizing Maps and Learning Vector Quantization, pp. 329–339. Springer, ??? (2016)
- von Wielligh, R.J., Grobler, M.J., Marais, H.-J.: Cellular iot capacity estimation for african smart cities. In: 2018 IEEE Global Conference on Internet of Things (GCIoT), pp. 1–6 (2018). IEEE
- Calandro, E., Chavula, J., Phokeer, A.: Internet development in africa: a content use, hosting and distribution perspective. In: International Conference on e-Infrastructure and e-Services for Developing Countries, pp. 131–141 (2018). Springer

21. Ecobank Research: The high cost of mobile data in Sub-Saharan Africa. High data costs are constraining Africa's digital revolution (2018). <https://www.ecobank.com/upload/publication/20180910054643018QJEBKEVZKD/20180910054635730h.pdf>
22. Gillwald, A., Mothobi, O.: After access 2018: A demand-side view of mobile internet from 10 african countries (2019)
23. Shankar, V.: Announcing facebook lite. Facebook Newsroom (2015)
24. DigitBin: WhatsApp Lite APK Download for Android. <https://www.digitbin.com/whatsapp-lite/>. Accessed: 2021-06-30
25. Uber: Uber Lite - Fast, Reliable, and Just 5MB. <https://www.uber.com/ug/en/u/uber-lite-app/>. Accessed: 2021-06-30
26. Google LLC: Lite but packs a punch: Google Go comes to Android everywhere. <https://www.blog.google/products/search/lite-packs-punch-google-go-comes-android-everywhere>. Accessed: 2021-06-30
27. Google LLC: See Gmail in standard or basic HTML version - Gmail Help. <https://support.google.com/mail/answer/15049?hl=en>. Accessed: 2021-06-30
28. Burns, B., Grant, B., Oppenheimer, D., Brewer, E., Wilkes, J.: Borg, omega, and kubernetes. Queue **14**(1), 70–93 (2016)
29. Sahandi, R., Alkhalil, A., Opara-Martins, J.: Cloud computing from smes perspective: a survey based investigation. Journal of Information Technology Management **24**(1), 1–12 (2013)
30. Hohpe, G.: Don't get locked up into avoiding lock-in (2019). <https://martinfowler.com/articles/oss-lockin.html>
31. Opara-Martins, J., Sahandi, R., Tian, F.: Critical analysis of vendor lock-in and its impact on cloud computing migration: a business perspective. Journal of Cloud Computing **5**(1), 4 (2016)
32. Kratzke, N., et al.: Lightweight virtualization cluster how to overcome cloud vendor lock-in. Journal of Computer and Communications **2**(12), 1 (2014)
33. Nadareishvili, I., Mitra, R., McLarty, M., Amundsen, M.: Microservice Architecture: Aligning Principles, Practices, and Culture. " O'Reilly Media, Inc.", ??? (2016)
34. Knoche, H., Hasselbring, W.: Drivers and barriers for microservice adoption-a survey among professionals in germany. Enterprise Modelling and Information Systems Architectures (EMISAJ)-International Journal of Conceptual Modeling **14**(1), 1–35 (2019)
35. Hasselbring, W., Steinacker, G.: Microservice architectures for scalability, agility and reliability in e-commerce. In: 2017 IEEE International Conference on Software Architecture Workshops (ICSAW), pp. 243–246 (2017). IEEE
36. Fowler, M., Lewis, J.: Microservices a definition of this new architectural term. URL: <http://martinfowler.com/articles/microservices.html>, 22 (2014)
37. Hüttermann, M.: DevOps for Developers. Apress, ??? (2012)
38. Balalaie, A., Heydarnoori, A., Jamshidi, P.: Microservices architecture enables devops
39. Dragoni, N., Dustdar, S., Larsen, S.T., Mazzara, M.: Microservices: Migration of a mission critical system. arXiv preprint arXiv:1704.04173 (2017)
40. Hope, P.: Using jails in freebsd for fun and profit. login: The Magazine of USENIX & SAGE **27**(3) (2002)
41. Modak, A., Chaudhary, S., Paygude, P., Ldate, S.: Techniques to secure data on cloud: Docker swarm or kubernetes? In: 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), pp. 7–12 (2018). IEEE
42. Levijarvi, E.S., Mitzev, O.S.: Private cloud replication and recovery. Google Patents. US Patent 8,930,747 (2015)
43. Li, W., Yang, Y., Chen, J., Yuan, D.: A cost-effective mechanism for cloud data reliability management based on proactive replica checking. In: 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012), pp. 564–571 (2012). IEEE
44. Li, W., Yang, Y., Yuan, D.: A novel cost-effective dynamic data replication strategy for reliability in cloud data centres. In: 2011 IEEE Ninth International Conference on Dependable, Autonomic and Secure Computing, pp. 496–502 (2011). IEEE
45. Brewer, E.A.: Towards robust distributed systems. In: PODC, vol. 7 (2000). Portland, OR
46. Esteves, S., Silva, J., Veiga, L.: Quality-of-service for consistency of data geo-replication in cloud computing. In: European Conference on Parallel Processing, pp. 285–297 (2012). Springer
47. Gao, A., Diao, L.: Lazy update propagation for data replication in cloud computing. In: 5th International Conference on Pervasive Computing and Applications, pp. 250–254 (2010). IEEE
48. Boru, D., Kliazovich, D., Granelli, F., Bouvry, P., Zomaya, A.Y.: Models for efficient data replication in cloud computing datacenters. In: 2015 IEEE International Conference on Communications (ICC), pp. 6056–6061 (2015). IEEE
49. Bozman, J., Chen, G.: Cloud computing: The need for portability and interoperability. IDC Executive Insights (2010)
50. Gonidis, F., Paraskakis, I., Kourtesis, D.: Addressing the challenge of application portability in cloud platforms. In: 7th South-East European Doctoral Student Conference, pp. 565–576 (2012)
51. Alquraan, A., Tahruri, H., Alfatafta, M., Al-Kiswany, S.: An analysis of network-partitioning failures in cloud systems. In: 13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18), pp. 51–68 (2018)
52. Ongaro, D., Ousterhout, J.: In search of an understandable consensus algorithm. In: 2014 {USENIX} Annual Technical Conference ({USENIX} {ATC} 14), pp. 305–319 (2014)
53. Zhou, J., Shi, Z.: Unstructured p2p-enabled service discovery in the cloud environment. In: International Conference on Intelligent Information Processing, pp. 173–182 (2010). Springer
54. Ranjan, R., Zhao, L., Wu, X., Liu, A., Quiroz, A., Parashar, M.: Peer-to-peer cloud provisioning: Service discovery and load-balancing. In: Cloud Computing, pp. 195–217. Springer, ??? (2010)
55. Jaramillo, D., Nguyen, D.V., Smart, R.: Leveraging microservices architecture by using docker technology. In:

- SoutheastCon 2016, pp. 1–5 (2016). IEEE
- 56. Zhang, D., Yan, B., Feng, Z., Zhang, C., Wang, Y.: Container oriented job scheduling using linear programming model. In: 2017 3rd International Conference on Information Management (ICIM), pp. 174–180 (2017)
 - 57. Haselböck, S., Weinreich, R.: Decision guidance models for microservice monitoring. In: 2017 IEEE International Conference on Software Architecture Workshops (ICSAW), pp. 54–61 (2017). IEEE
 - 58. Noor, A., Jha, D.N., Mitra, K., Jayaraman, P.P., Souza, A., Ranjan, R., Dustdar, S.: A framework for monitoring microservice-oriented cloud applications in heterogeneous virtualization environments. In: 2019 IEEE 12th International Conference on Cloud Computing (CLOUD), pp. 156–163 (2019). IEEE
 - 59. UbuntuNet: UbuntuNet Monitor. <https://monitor.ubuntunet.net/cacti/>. Accessed: 2021-07-21
 - 60. Van den Bossche, R., Vanmechelen, K., Broeckhove, J.: Cost-efficient scheduling heuristics for deadline constrained workloads on hybrid clouds. In: 2011 IEEE Third International Conference on Cloud Computing Technology and Science, pp. 320–327 (2011). IEEE
 - 61. Filip, I.-D., Pop, F., Serbanescu, C., Choi, C.: Microservices scheduling model over heterogeneous cloud-edge environments as support for iot applications. IEEE Internet of Things Journal 5(4), 2672–2681 (2018)
 - 62. Müssig, D., Stricker, R., Lässig, J., Heider, J.: Highly scalable microservice-based enterprise architecture for smart ecosystems in hybrid cloud environments. In: ICEIS (3), pp. 454–459 (2017)
 - 63. Guerrero, C., Lera, I., Juiz, C.: Resource optimization of container orchestration: a case study in multi-cloud microservices-based applications. The Journal of Supercomputing 74(7), 2956–2983 (2018)
 - 64. Sousa, G., Rudametkin, W., Duchien, L.: Automated setup of multi-cloud environments for microservices applications. In: 2016 IEEE 9th International Conference on Cloud Computing (CLOUD), pp. 327–334 (2016). IEEE