



# Article A Cloud-Based Data Collaborative to Combat the COVID-19 Pandemic and to Solve Major Technology Challenges

Max Cappellari <sup>1</sup>, John Belstner <sup>2</sup>, Bryan Rodriguez <sup>2</sup> and Jeff Sedayao <sup>2,\*</sup>

- <sup>1</sup> XPRIZE Foundation, Culver City, CA 90230, USA; Max.Cappellari@xprize.org
- <sup>2</sup> Intel Corporation, Santa Clara, CA 95054-1549, USA; john.belstner@intel.com (J.B.);
  - bryan.j.rodriguez@intel.com (B.R.)
- Correspondence: jeff.sedayao@intel.com; Tel.: +1-408-765-2935

**Abstract:** The XPRIZE Foundation designs and operates multi-million-dollar, global competitions to incentivize the development of technological breakthroughs that accelerate humanity toward a better future. To combat the COVID-19 pandemic, the foundation coordinated with several organizations to make datasets about different facets of the disease available and to provide the computational resources needed to analyze those datasets. This paper is a case study of the requirements, design, and implementation of the XPRIZE Data Collaborative, which is a Cloud-based infrastructure that enables the XPRIZE to meet its COVID-19 mission and host future data-centric competitions. We examine how a Cloud Native Application can use an unexpected variety of Cloud technologies, ranging from containers, serverless computing, to even older ones such as Virtual Machines. We also search and document the effects that the pandemic had on application development in the Cloud. We include our experiences of having users successfully exercise the Data Collaborative, detailing the challenges encountered and areas for improvement and future work.

**Keywords:** containers; virtual machines; cloud; COVID-19; serverless; analytics; software defined infrastructure

# 1. Introduction

The XPRIZE Foundation [1–3] designs and operates multi-million-dollar, global competitions to incentivize the development of technological breakthroughs that accelerate humanity toward a better future. The foundation's mission is to inspire and empower a global community of problem-solvers to positively impact our world. XPRIZE believes solutions to the world's problems can come from anyone, anywhere.

As part of that philosophy, XPRIZE looked to develop a Data Collaborative that will support the resolution of complex, global problems. It had three broad goals:

- 1. Comprising a collection of unique datasets and AI tools, this Data Collaborative would democratize data and the tools to analyze them, enabling virtually anyone, anywhere to use data to solve the world's grandest challenges.
- 2. It would provide access to the massive amounts of data that have been organized and cleaned so that it can be accessed in an accountable, transparent, and responsible manner.
- 3. Data owners across industries and topic areas will feel safe in sharing their data while maintaining ownership, privacy, and security. Data scientists, as well as Machine Learning and Artificial Intelligence systems, will all have access to the Data Collaborative, resulting in enhanced problem-solving models and innovative approaches to solving Grand Challenges posed by XPRIZE competitions.

Implementation of the Data Collaborative started in late 2019 but paused when the COVID-19 pandemic struck in early 2020. The XPRIZE Foundation decided to leverage existing Data Collaborative plans to make COVID-19 relevant data available from the XPRIZE



**Citation:** Cappellari, M.; Belstner, J.; Rodriguez, B.; Sedayao, J. A Cloud-Based Data Collaborative to Combat the COVID-19 Pandemic and to Solve Major Technology Challenges. *Future Internet* **2021**, *13*, 61. https://doi.org/10.3390/ fi13030061

Academic Editors: Nane Kratzke and Paolo Bellavista

Received: 11 January 2021 Accepted: 22 February 2021 Published: 27 February 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). Pandemic Alliance [4,5], which is a cooperative of different organizations sponsored by the foundation that would contribute data for analysis. The Data Collaborative would first be used for COVID-19 efforts and then be reused for other datacentric XPRIZE competitions and challenges.

A few existing environments and services have similarities to the Data Collaborative. Google's Kaggle [6] is a web-based environment that provides an analytics environment to users at large along with a large collection of publicly datasets. Similar to the Data Collaborative, Kaggle is used for hosting competitions [7], focusing on machine learning-oriented problems. Google also maintains Colab [8,9], which provides a Python-oriented environment. Available commercial services for hosting analytics environments include CoCalc [10,11] and Nextjournal [12]. The XPRIZE Foundation differs from these operations by not focusing on individual data scientists or consumers and instead focusing on fewer, higher impact technology challenges run by teams of people.

We present the design and implementation of the XPRIZE Data Collaborative as a descriptive case study [13,14] of a Cloud Native Application (CNA) [15–18], from requirements, design, and implementation to operating experiences. Many case studies and surveys of Cloud implementations focus on particular categories of Cloud Technology, such as microservices and container technologies [17,19,20] and serverless computing [21–23], but our examination asks the question of whether real CNAs are actually so narrowly focused and if they would use other Cloud technologies, such as Software Defined Infrastructure (SDI) [24] and even older, seemingly deprecated technologies such as Virtual Machines (VMs). As we describe our implementation choices, we match those choices with common, well-known design Cloud Design patterns [18,19,22,23,25]. Our case study also looks at whether the COVID-19 pandemic affected the Data Collaborative's design, development, and implementation.

In this introduction, we have described critical background information such as the XPRIZE Foundation's mission and goals for the Data Collaborative. We have listed similar and related services that currently exist and provided context regarding Cloud Native Applications with which we will frame our case study. Section 2 describes the methodology of our efforts, including the goals of our case study, detailed requirements of the Data Collaborative, and the infrastructure design chosen to meet those needs. Section 3 discusses the results of implementing our design, detailing the challenges we encountered and how we handled those challenges. Section 4 talks about insights about Cloud Computing implementations that we have learned from our experiences. Section 5 reviews our findings, describes, possible changes to the Data Collaborative for future competitions, and discusses implications to CNA development from what we learned.

#### 2. Materials and Methods

In this section, we describe the key research questions and methodology for our case study on the Data Collaborative. As part of the methodology, we examine the lifecycle of XPRIZE Contests and how these resulted in detailed requirements and individual components of the eventual solution. A holistic view of the Data Collaborative infrastructure is presented at the end of the section.

#### 2.1. Research Questions and Motivations

Our case study of the Data Collaborative Infrastructure attempts to answer the following questions:

- 1. Are applications developed in the Cloud moving to only use containers? How are the Cloud technologies ranging from VMs to containers SDIs used and why?
- 2. How has the COVID-19 pandemic affected application development and deployment?
- 3. Are there any noticeable trends in Cloud Computing deployment that became apparent during the development, and if so, what are they?

In addition, we are motivated to add to the literature of case studies for Cloud Computing application deployments. We mentioned in the Introduction that many case studies of Cloud implementations focus on specific technologies such as containers or serverless computing. We present a view of how these technologies could be used together and the tradeoffs and usage patterns appropriate to each, which are demonstrated in a real and potentially highly used CNA.

#### 2.2. Methodology

We conduct our case study by documenting the design, development, and operation of the XPRIZE Data Collaborative. From that documentation, we answer the research questions above. Understanding Data Collaborative requirements means understanding the XPRIZE Contest Lifecycle, which we describe in the next section. We use that information to progressively refine requirements to generate a list of the components that need to be implemented and what properties they need to have. For each component, we examine the possible design space that applies and select the best option.

In the Results section of this paper, we look at how well the design and implementation met requirements and then answer our research questions. In the Discussion section that follows, we make additional observations, while in the Conclusion, we describe future work that is contemplated for the Data Collaborative.

While the Data Collaborative was developed and deployed within one data center of a Cloud Service Provider, we do not use any Cloud Service Provider specific terminology or code. We discuss our implementation in a generic way. This allows others to replicate our work on any Cloud Service Provider that has the capabilities we mention (more than one does) and to validate or disprove our findings.

#### 2.3. Contest Lifecycle

We need to understand the lifecycle of an XPRIZE Foundation contest in order to design an optimal infrastructure for one. Figure 1, provided by the XPRIZE Foundation, shows the typical flow of one of their contests.



Figure 1. Typical lifecycle of an XPRIZE Competition.

Contests have a definitive start and finish, so apart from the obvious need to assemble the infrastructure necessary for a contest, that infrastructure needs to be maintained and then taken down. Once a competition is designed and funded through partnerships, the recruitment and team assemble phases require teams to register on an XPRIZE portal. During a competition, there are typically a few rounds where progress is judged, and the number of competitors is reduced. A contest could have as many as 5000 teams in the initial round, and there can be simultaneous contests of varying duration and in different stages of their lifecycle at any one time.

#### 2.4. Data Collaborative Detailed Requirements

Given the broad requirements in Section 1 and information about the contest/challenge lifecycle, we generate more detailed requirements of what needs to be implemented:

1. A highly scalable, accessible, and elastic infrastructure: A single competition could involve as many as thousands of teams at one time, and as each competition milestone is reached, the number could be reduced by one or two orders of magnitude.

Teams can be from almost anywhere in the world. The ability to rapidly scale up and down is required.

- 2. A stable and highly usable analytics software platform: Teams in XPRIZE competitions will need to be able to not just access data but analyze it. It should be familiar to many users and extendable.
- 3. **Isolated and secure analytics software platform**: Since awards in XPRIZE competitions can be multi-millions of dollars, teams and their work need to be effectively isolated from each. In addition, since we are providing the infrastructure for teams to run arbitrary code, we need to make sure that what we provide is not abused—not used as a base for attacks or noncompetition-related work such as crypto-mining.
- 4. A scalable analytics compute platform: While the common goal is to democratize data access and the ability to analyze that data, should a team wish to purchase more capacity, the platform should enable that.
- 5. **Control of data**: Contest sponsors and others want to be able to protect the data they provide and know who accesses it. The data should remain within the analytics platform, with it being difficult and time consuming to copy it outside of the analytics platform. Access to data needs to be recorded and attributable to a team.
- 6. **Manageability**: The common infrastructure needs to be maintainable by a relatively small staff.
- 7. **Reasonably fast implementation**: The infrastructure needs to be scaled up in a reasonably fast amount of time in order to make a difference in the pandemic.
- 8. **Costs**: Costs should be minimized when possible.

We will refer to these requirements as we make design decisions.

With more detailed requirements, we enumerate what decisions need to be made. The following components and processes need to be implemented in order to make the Data Collaborative viable:

- 1. User Analytics Software Platform
- 2. Team Isolation
- 3. Naming Design and Infrastructure
- 4. User Authentication
- 5. Protecting Data in Transit
- 6. Logging and Monitoring
- 7. Team Infrastructure Instantiation Process

We document the design options and final choices for each item above in the next section.

#### 2.5. Design Choices

In the previous sections, we generated a set of requirements and a set of components and processes that need to be designed in order to create the Data Collaborative. In this section, we go through each required component, discussing the possible design options, choosing an option, and then justifying that decision in terms of the requirements defined above. These decisions will be framed, wherever possible, in terms of Cloud Design patterns.

#### 2.5.1. User Analytics Software Platform

Requirement #2 (a stable and highly usable analytics software platform) drove our selection for an analytics software platform. Notebooks have become a crucial tool for data scientists [26] as a way of developing analytics code, visualizing analytics results, and sharing insights. NBView, a tool that estimates the number of publicly available notebooks, estimates that there are almost 10 million notebooks publicly available today on GitHub alone [27]. Platforms that offer similar functionality to our goals such Kaggle [6,7], CoCalc [10,11], and Nextjournal [12] all offer support for notebooks. Some sort of notebook seemed to be an obvious choice for the Data Collaborative.

We looked at two choices for the standard Data Collaborative notebook. The first choice was a Jupyter notebook [28], which is a web-based platform for analytics. The other

choice was Zeppelin [29], which is another web-based platform for analytics and an Apache project that runs on top of the Spark [30] analytics system. Jupyter notebooks are very widely used, stable, and have a larger community and eco-system around it. Kaggle, CoCalc, and Nextjournal offer some varying levels of support for Jupyter notebooks. While the community for the Zeppelin is growing, it is not yet as popular.

Another factor in analytics platform selection was requirement #4. Enabling teams to be able to purchase more resources for their project made it necessary to isolate notebooks into units that could use more resources. This was difficult to do quickly (requirement #7) in highly integrated multiuser implementations of Jupyter notebooks such as JupyterHub [31]. Zeppelin notebooks also are tightly integrated with Spark. Jupyter notebooks can be implemented on a standalone basis [28] (not part of JupyterHub). That property, along with the popularity of Jupyter notebooks, made the standalone Jupyter notebooks our platform of choice.

#### 2.5.2. Team Isolation

Teams competing for large monetary prizes need to be effectively isolated from each other (requirement #3). The two most viable placement choices for each team's notebook are within VMs or within a container. Note that putting each team on an individual server would provide maximum isolation but would meet neither requirement #8 (prohibitively expensive) nor requirement #1 (not scalable). Using VMs would provide better isolation than using containers but would be more resource intensive and therefore more expensive. In the end, we deemed that containers have sufficient isolation at a better price point.

To implement requirement #4, each container would be in its own resource group. That allows teams to add resources to the resource group by working directly with the Cloud Service Provider.

#### 2.5.3. Naming Design and Infrastructure

Our choice to use standalone Jupyter notebooks instantiated in individual containers, rather than an integrated environment such as JupyterHub, forced us to find an effective way to direct teams to their individual notebooks. If we had used JupyterHub, we could give all the teams a single domain name to connect to, and then they would be directed to their notebook after they authenticated. Instead, we had to direct users on an individual team to the individual container for that team. We chose to do this through DNS host names in the URL for each notebook.

Figure 2 shows our design of the infrastructure and process for connecting teams to their notebooks. We present each team that joins a competition with a unique DNS name for their notebook in the xprize.org domain. In the example below, "team1" gets a URL such as "http://team1.comp1.xprize.org/" as step 1. Rather than create a new DNS [32] record for each new team that was added or removed, we used a wild-card domain record for \*.comp1.xprize.org that points all teams in a competition to an application load balancer of our Cloud Service Provider.



Figure 2. Routing of web requests to team notebook container.

An HTTPS request for a notebook arrives at the application load balancer (labeled Application gateway), and in step 2, it is sent to a proxy server. In step 3, the proxy server routes the request to the container containing the team's notebook based on the name of the container that we gave to that team.

The container holding their notebook is in a private IP network space [33] that is not directly accessible to users on the Internet. For step 3 to work, we need a way to name the container and register that name into a private DNS space. Our Cloud Service Provider has a DNS service that can be manipulated programmatically, so we use a private DNS domain for our containers' names. In Figure 2, this host name for the container with team1's notebook is team1.privatexprize.org.

While this level of detail may seem excessive, explaining the design of the DNS namespace and architecture is important for our discussion. First, it shows how we applied the gatekeeper design pattern [25]. Second, it illustrates our use of SDI (Software Defined Infrastructure) [24], as the DNS infrastructure is allocated and loaded with domain information programmatically. Third, this infrastructure has some problematic features, and it should be understood in order to put some of the problems we experienced (to be described later) into context.

#### 2.5.4. User Authentication

Authentication of the users in a team became a design decision, since the official multiuser solution for Jupyter notebooks is JupyterHub, which we decided not to use. Since we chose to give teams a standalone Jupyter notebook isolated into a container, we needed to use one of the available authentication methods available for individual notebooks. Individual Jupyter notebooks have two options for user authentication [34]:

- 1. Password
- 2. Authentication token

As a team would have to share a password, this would be insecure, violating requirement #3. Password management (e.g., forgotten passwords, managing password change, and distribution software) would be additional work for the XPRIZE team to manage passwords, affecting requirement #6). XPRIZE already had a portal for teams to sign up that relied on third-party authentication services from OAuth identify providers such as Facebook and Google [35], and we decided use that portal to authenticate the users on a team, generate the token for a new notebook, and then distribute that token to the team members. Moreover, this matches a standard Cloud design pattern, federated identity [25].

# 2.5.5. Protecting Data in Transit

To enforce isolation between teams (requirement #3) and to help control access to data (requirement #5), we need to protect data in transit to and from the Data Collaborative. In addition, we needed to make sure that the token that we decide to use for authentication would not traverse the Internet or Intranets in the clear. In our case, this means encrypting data in transit using Transport Layer Security (TLS) [36]. As this is a common practice with websites, the chief design decision came in of where to put the certificate and do the encryption/decryption: the TLS endpoint. We looked at the following options:

- 1. On each notebook container
- 2. On a sidecar container
- 3. Proxy
- 4. Application gateway

Option 1 is viable using services such as Let's Encrypt [37] but is potentially harder to manage and consumes more resources, as adding the certificate process to each container increases the tasks and complexity to each container instantiation and rollout. The sidecar container pattern [38] could be used—this pattern would instantiate a new container next to each notebook container, which would hold the certificate and be a TLS endpoint, communicating the notebook container in the clear. This is easier to manage than the first

option, since all of those TLS sidecar containers could be identical, but it drives up the number of containers and costs (requirement #8) and results in twice as many containers to manage and monitor (affecting requirement #6). Putting the TLS endpoint on the proxies is a better option. This is a function that many proxies can perform, and there are far fewer proxies to manage for each competition. A better option is the last option—putting the endpoint on the Application Gateway. This leaves only one place to put the certificate—making container instantiation simpler, and it only needs to be done once. It lessens the amount of processing done on containers and the proxies, reducing cost (requirement #8). We selected this option, which also is a standard Cloud design pattern: gateway offloading [39].

#### 2.5.6. Logging and Monitoring

Putting each team in a notebook in a container enabled us to use Cloud Native monitoring functions available from our Cloud Service Provider. In particular, the provider has services that check for inappropriate uses of the compute facilities we provide, such as looking for crypto-mining. The provider also has centralized logging repositories to which we can send operational data, providing a single place to look for anomalies and investigate incidents. These are services not unique to our provider—other Cloud Service Providers have these two capabilities. We chose to utilize both.

#### 2.5.7. Instantiating Team Environment

Before a new team receives their notebook, there needs to some way of instantiating the environment for each team. Scripts need to be run that build the individual components for a team—the notebook container, the private DNS entries, and firewall rules. Teams would be notified and received the authorization token after their infrastructure is ready. Instantiating a new notebook for a new team happens sporadically. The main design question becomes where to run the scripts. The following options were available:

- 1. Run the instantiation scripts on XPRIZE portal.
- 2. Create a dedicated VM or container to run the scripts.
- 3. Define a function to run the scripts and launch it in the Cloud (serverless).

While option #1 is certainly possible and some processing must be done on the portal to provide the token, running the entire instantiation process ties up resources there and makes it a potential performance bottleneck and single point of failure. Option #2 of creating a whole VM or container is viable, but creating a new one every time a notebook must be instantiated will take time and cost money, while having one ready to process notebook instantiation requests leaves dedicated resources idle. We choose option #3, creating a function that instantiates that notebook and informs the team when it is ready. The requirements fit in well with the serverless design pattern for Asynchronous, Event-Driven processing [40]. We took advantage of serverless computing options offered by our Cloud Service Provider to implement this.

#### 2.6. Architecture and Component Layout

Figure 3, which is derived from our implementation, shows the overall architecture of a single Data Collaborative competition, showing the placement of the components we designed. Each team's notebook shares the same public IP address, which is an address on an application gateway. The gateway has multiple functions, such as load balancing traffic to the proxy servers, running a Web Application Firewall, and being a TLS endpoint.



Figure 3. Architecture and component layout for a single Competition.

We use reverse proxies on VMs to balance and filter traffic to the notebook containers, restrict the kind of operations permitted on data (requirement #5), and provide additional layers of isolation from the Internet. We made the choice of VMs rather than containers as we ran into issues containerizing the proxy application. In addition, as we were developing the application, we needed to log into the proxy periodically and make changes to configurations—this was much easier to do when they were in VMs rather than containers. Once we had the working configuration, we deployed using proxies in VMs.

The containers containing teams' notebooks are also on a separate segment. This segment has access lists that restrict where on the Internet the notebooks can connect. Teams can download software to their notebooks from a limited number of software repositories. Another firewall controls access from the containers to the datasets used as the basis for competition. These datasets are also stored in the Cloud on SDI.

The infrastructure in Figure 3 will be duplicated for each competition. This minimizes the opportunity for datasets from one competition leaking into another. It also reduces the effects of any single configuration mistake from affecting all competitions. This architectural decision matches the Cloud design pattern called the Bulkhead pattern [41].

#### 3. Results

In this section, we examine the results of our implementation efforts. We describe how long it took to create the Data Collaborative and what resources were needed to design and implement, and we check if the resulting infrastructure met requirements. After that, we revisit our research questions and answer them in the context of the results.

#### 3.1. Implementation Time and Resource Usage

Initial meetings with the members of the team that would build that on the Data Collaborative's infrastructure (listed above as authors) started in September 2019. A first workable prototype with deployment documentation was delivered in June of 2020. The four authors of this paper did not work on this project full time—all had other projects done during the time frame of this work.

#### 3.2. Evaluation of the Design and Implementation against Requirements

Our design and implementation of the Data Collaborative met the goal of creating a highly scalable analytics platform for competitions. It is in active use, for both COVID-19 work [4] and for general AI-related challenges [42]. We first completed a small competition with 25 teams and are currently in the testing phase of a competition that has hosted some 300 teams [4], so it definitely is an operational utility for the XPRIZE foundation. Table 1 shows a list of the requirements in Section 2.4 and whether the implementation met them.

In general, requirements were fulfilled. Some requirements are hard to prove, such as being "secure."

 Table 1. Data Collaborative high-level requirements and whether they were met.

Goal	Goal Met?	Title 3
1. Scalable Infrastructure	Yes	Scaled to 300
2. Usable Analytics Platform	Yes	Users for one complete competition and one in progress
3. Isolated Secure Platform	Possibly	Secure as far as we know—no incidents so far
4. Scalable Analytics Compute Platform	Probably	Put in place but not yet exercised
(self-service capacity adds)		
5. Control of data	Yes	Access to datasets is controlled
6. Manageability	Yes	Manageable by small XPRIZE team
7. Fast Implementation	Yes	Started in late 2019 and first competition started in mid-2020
8. Cost-controlled	Yes	No evidence of cost issues at this point

We conclude that the Data Collaborative infrastructure did meet its requirements.

### 3.3. Revisiting Our Research Questions

A successful implementation allows us to revisit our research questions.

#### 3.3.1. Are Applications Developed in the Cloud Moving to Only Use Containers?

As described in our design, the Data Collaborative uses a variety of technologies for very different circumstances. There are Cloud design patterns that dictate when Serverless computing, SDI, and containers are best used. Using a VM was a surprise, though.

# 3.3.2. How Has the COVID-19 Pandemic Affected Application Development and Deployment?

As mentioned in the introduction, the COVID-19 Pandemic pushed XPRIZE to focus on new challenges and to accelerate development. Another major effect of the pandemic was on the requirements. Another way to look at its influence is to look at the "nonrequirements" for the Data Collaborative, which were not priorities because of the need to implement relatively quickly:

- 1. **Portability between Clouds**: Extensive work has been done try to create technology that makes applications portable between Cloud Service Providers [43] and avoiding vendor lock-in [44,45]. This goal was simply not a priority in our efforts to get the Data Collaborative up and running. We chose to use the Cloud Service Provider that we were most familiar with to get the service to production as soon as possible.
- 2. **Ultra-low cost**: While a reasonable cost is a key requirement (#8 above), having an extremely low cost, especially at the expense of having a usable and working solution, was not.
- 3. **High Performance**: While our requirement for a usable and scalable platform (requirements #1 and #2) demands enough performance to be usable, it did not demand high performance. We need the Data Collaborative in a relatively short time frame that worked and was usable. In addition, the contest time frames should allow for time for analyses to complete.

Similar to many other IT development projects, developing the Data Collaborative had to be done remotely. There were face-to-face meetings and work sessions during late 2019, but the bulk of development was done remotely at home by the team in locations ranging from Northern California to Southern California and Arizona. We would occasionally encounter problems in working with our Cloud Service Provider, but it was impossible to know whether these were from Cloud Service Provider capacity problems or just network congestion from stay-at-home orders [46].

3.3.3. Are There Any Noticeable Trends in Cloud Computing Deployment That Became Apparent during the Development?

In addition to our observation about using multiple Cloud technologies, we have observed that even in the short time span of this project, Cloud Service Provider capabilities evolve extremely rapidly. We implemented the proxy servers because the application gateways lack certain functionality in rewriting URLs. By our second competition, the application gateway had already gained some of the functionality that was lacking, and we used it instead.

# 4. Discussion

In this section, we review possible limitations of this work. We also detail other issues we encountered, future work needed, and possible implications of our observations on Cloud Native Application development.

#### 4.1. Study Limitations

While we provided a detailed look into the requirements, design, and operation of a live, in-use Cloud Native Application, it should be noted that these are the results for one application. We saw that a variety of technologies ranging from VMs to containers to serverless computing were used, but this will not be applicable in all Cloud environments, as some Cloud Service Providers do not provide all of these services.

There is also the possibility that we may have missed better and easier ways to implement some of the functionality of the Data Collaborative. We matched our choices to known Cloud patterns, but there may be other options. At some point, the notebook functionality for our users may be a service offered by Cloud providers, which would really change the dynamic of how we use the Cloud.

#### 4.2. Other Issues Encountered

While the Data Collaborative proved to be scalable to a certain point, we did encounter issues. We had concentrated on setting up a competition environment under schedule pressure, but we did not focus much of what would take to deprovision it after a contest. After our first competition, some configured resources, such as containers and DNS records, had to be removed manually. This kind of problem is solvable through automation and scripting.

We also encountered problems regarding application state and containers. Since teams could add software packages to their notebook, any state changes such as this would be reflected in the current state of the container. However, if the container crashed for some reason and a replacement one was instantiated, any additions would be lost. This is because we did not separate out the stateful components of the notebook—additional software packages for instance. We fixed this by making sure that state changes such as this would be saved and restored if the container went down and had to be restarted.

A more subtle problem with container state happened because of the DNS architecture shown in Figure 2. A container's name in our private DNS space has a particular Time to Live (TTL). As an example, let us say that the record for the container named team1.privatexprize.org and IP address associated with 192.168.1.1 have a TTL of two hours. If the container goes down and comes back up again with IP address 192.168.1.2, an infrastructure component such as our proxy could attempt to connect to team1.privatexprize.org at 192.168.1.1 for up to two hours (until the TTL expires). The container state exists not only in the container but in infrastructure in things such as DNS. We reduced the impact of this problem by reducing the DNS TTL in our container private DNS.

The notebook per container architecture had many benefits, such as improving notebook isolation and using native Cloud container monitoring, but it had the drawback of being harder to maintain. If a software upgrade needed to be made to all notebooks in a competition, then we would have to go through each notebook, upgrade it in place if possible, and create new copies and then redeploy if not. We want to examine how much of requirement #4 is really needed, and whether we can use a more centralized approach using JupyterHub to simplify much of the architecture and operation.

### 5. Conclusions

The answers to our research question show that now, Cloud Native Applications can justify the use multiple Cloud technologies, ranging from virtual machines to containers to serverless computing. The COVID-19 pandemic prioritized getting a working version running above other concerns such as cost or performance. One observation relevant to our last research questions is that Cloud Services evolve very quickly, which has implications for future work described next.

#### 5.1. Future Work

We had put our proxies into VMs to expedite getting the Data Collaborative running but having to log into them to configure them, and to look at logs makes VMs unwieldy to work with. The proxy configuration and logging need to be automated, and the proxies need to be put into a container for easier deployment into contests and for maintenance. In the longer run, we are looking to see if the functionality of the proxies can be absorbed into Cloud infrastructure. We have observed that even in the short time span of this project, native Cloud infrastructure capabilities have started catching up to proxy functionality. In Section 3.3.2, we talked about the non-requirements of performance, low cost, and avoiding service provider lock-in. In the long run, those areas will need to be revisited.

# 5.2. Implications for Cloud Native Application Development

Much of the Cloud design pattern literature focuses on certain areas of Cloud technology, such as microservices/containers [17,19,20], serverless computing [21–23], or Software Defined Infrastructure [24]. Our experiences show that a Cloud Native Application can use any or all of these technologies, and even older Cloud technologies such as Virtual Machines. Design patterns need to be more inclusive and prescriptive about how and when to use them, as multiple technologies can apply in a single application implementation.

Our experiences with DNS and other issues with application and the file system state should warn developers deploying in the Cloud that maintaining state in a container is a bad idea. As shown with DNS, state problems can occur in subtle and unexpected ways. Tools to find and debug such problems would be very useful.

Dealing with service provider lock-in [44,45] was not a priority. Tools for reducing vendor lock-in need to be simple to use and standardized for them to be used.

**Author Contributions:** Conceptualization, M.C.; Data curation, M.C.; Funding acquisition, M.C.; Investigation, M.C., J.B., B.R. and J.S.; Methodology, J.B. and B.R.; Project administration, M.C.; Resources, M.C.; Software, M.C., J.B., B.R. and J.S.; Supervision, M.C.; Writing—original draft, J.S.; Writing—review & editing, M.C., J.B., B.R. and J.S. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: No Applicable, the study does not report any data.

**Acknowledgments:** We would like to thank Felix Labunsky of Microsoft for reviewing our architecture and providing some guidance. We would also like to thank Dan Gutwein of Intel for helping to drive this project.

Conflicts of Interest: The authors declare no conflict of interest.

# References

- 1. XPRIZE Foundation. Available online: https://www.xprize.org/ (accessed on 28 October 2020).
- 2. Hossain, M.; Kauranen, I. Competition-Based Innovation: The Case of the X Prize Foundation. J. Organ. Des. 2014, 3, 46–52. [CrossRef]
- 3. Haller, J.B.A.; Bullinger, A.C.; Möslein, K.M. Innovation Contests. Bus. Inf. Syst. Eng. 2011, 3, 103–106. [CrossRef]
- 4. Accelerating Radical Solutions to COVID-19 and Future Pandemics. Available online: https://www.xprize.org/fight-covid19 (accessed on 28 October 2020).
- Mackay, M.J.; Hooker, A.C.; Afshinnekoo, E.; Salit, M.; Kelly, J.; Feldstein, J.V.; Haft, N.; Schenkel, D.; Nambi, S.; Cai, Y.; et al. The COVID-19 XPRIZE and the need for scalable, fast, and widespread testing. *Nat. Biotechnol.* 2020, *38*, 1021–1024. [CrossRef] [PubMed]
- 6. Kaggle. Available online: https://kaggle.com/ (accessed on 5 November 2020).
- Yang, X.; Zeng, Z.; Teo, S.G.; Wang, L.; Chandrasekhar, V.; Hoi, S. Deep learning for practical image recognition: Case study on kaggle competitions. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, London, UK, August 2008*; ACM: New York, NY, USA, 2018; pp. 923–931.
- 8. CoLab. Available online: https://colab.research.google.com/notebooks/intro.ipynb (accessed on 5 November 2020).
- 9. Bisong, E. Google Colaboratory. In *Building Machine Learning and Deep Learning Models on Google Cloud Platform;* Apress: Berkeley, CA, USA, 2019; pp. 59–64.
- 10. CoCalc. Available online: https://cocalc.com/ (accessed on 5 November 2020).
- 11. Bouvin, N.O. From notecards to notebooks: There and back again. In *Proceedings of the 30th ACM Conference on Hypertext and Social Media, Hof, Germany, 2019; ACM: New York, NY, USA, 2019; pp. 19–28.*
- 12. Nextjournal. Available online: https://nextjournal.com/ (accessed on 5 November 2020).
- 13. Baxter, P.; Jack, S. *Qualitative Case Study Methodology: Study Design and Implementation For Novice Researchers;* The Qualitative Report: Fort Lauderdale, FL, USA, 2013; Volume 13, pp. 544–559.
- 14. Baskarada, S. Qualitative Case Study Guidelines; The Qualitative Report: Fort Lauderdale, FL, USA, 2014; Volume 19, pp. 1–25.
- Andrikopoulos, V.; Fehling, C.; Leymann, F. Designing for CAP—The Effect of Design Decisions on the CAP Properties of Cloud-native Applications. In Proceedings of the 2nd International Conference on Cloud Computing and Services Science (CLOSER 2012), Porto, Portugal, 18–21 April 2012; pp. 365–374.
- 16. Kratzke, N. A Brief History of Cloud Application Architectures. *Appl. Sci.* 2018, *8*, 1368. [CrossRef]
- 17. Kratzke, N.; Quint, P.C. Understanding cloud-native applications after 10 years of cloud computing-a systematic mapping study. *J. Syst. Softw.* **2017**, *126*, 1–6. [CrossRef]
- 18. Gannon, D.; Barga, R.; Sundaresan, N. Cloud-Native Applications. IEEE Cloud Comput. 2017, 4, 16–21. [CrossRef]
- 19. Balalaie, A.; Heydarnoori, A.; Jamshidi, P. Microservices architecture enables devops: Migration to a cloud-native architecture. *IEEE Softw.* **2016**, *33*, 42–52. [CrossRef]
- 20. Burns, B.; Oppenheimer, D. Design patterns for container-based distributed systems. In Proceedings of the 8th {USENIX} Workshop on Hot Topics in Cloud Computing (HotCloud 16), Denver, CO, USA, 20–21 June 2016.
- Baldini, I.; Castro, P.; Chang, K.; Cheng, P.; Fink, S.; Ishakian, V.; Mitchell, N.; Muthusamy, V.; Rabbah, R.; Slominski, A.; et al. Serverless Computing: Current Trends and Open Problems. In *Research Advances in Cloud Computing*; Springer: Singapore, 2017; pp. 1–20.
- 22. Eismann, S.; Scheuner, J.; Van Eyk, E.; Schwinger, M.; Grohmann, J.; Herbst, N.; Abad, C.L.; Iosup, A. Serverless Applications: Why, When, and How? *IEEE Softw.* 2021, *38*, 32–39. [CrossRef]
- 23. Hong, S.; Srivastava, A.; Shambrook, W.; Dumitraș, T. Go serverless: Securing cloud via serverless design patterns. In Proceedings of the 10th {USENIX} Workshop on Hot Topics in Cloud Computing (HotCloud 18), Boston, MA, USA, 15 June 2020.
- Kang, J.-M.; Lin, T.; Bannazadeh, H.; Leon-Garcia, A. Software-Defined Infrastructure and the SAVI Testbed. In Proceedings of the International Conference on Testbeds and Research Infrastructures, Guangzhou, China, 5–7 May 2014; Springer: Berlin/Heidelberg, Germany; pp. 3–13.
- 25. Homer, A.; Sharp, J.; Brader, L.; Narumoto, M.; Swanson, T. *Cloud Design Patterns: Prescriptive Architecture Guidance for Cloud Applications*; Microsoft Patterns & Practices: Redmond, WA, USA, 2014.
- 26. Meyers, A. Data Science Notebooks—A Primer. Available online: https://medium.com/memory-leak/data-science-notebooks-a-primer-4af256c8f5c6/ (accessed on 8 November 2020).
- 27. Estimate of Public Jupyter Notebooks on GitHub. Available online: https://nbviewer.jupyter.org/github/parente/nbestimate/ blob/master/estimate.ipynb (accessed on 10 January 2021).
- 28. Jupyter Notebook. Available online: https://jupyter.org/ (accessed on 6 November 2020).
- 29. Apache Zeppelin. Available online: https://zeppelin.apache.org/ (accessed on 7 November 2020).
- 30. Zaharia, M.; Chowdhury, M.; Franklin, M.J.; Shenker, S.; Stoica, I. Spark: Cluster computing with working sets. *HotCloud* **2010**, *10*, 95–101.
- 31. Jupyterhub. Available online: https://jupyter.org/hub (accessed on 8 November 2020).
- 32. Mockpetris, P. RFC 1035: Domain Names—Implementation and Specification; RFC Editor: Fremont, CA, USA, 1987.
- 33. Rekhter, Y.; Moskowitz, B.; Karrenberg, D.; Groot, G.D.; Lear, E. Rfc 1918: Address Allocation for Private Internets. Available online: https://www.rfc-editor.org/info/rfc1918 (accessed on 24 February 2021).

- 34. Security in the Jupyter Notebook Server. Available online: https://jupyter-notebook.readthedocs.io/en/stable/security.html (accessed on 10 January 2021).
- Sun, S.T.; Beznosov, K. The devil is in the (implementation) details: An empirical analysis of OAuth SSO systems. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security, Raleigh, NC, USA, October 2012*; ACM: New York, NY, USA, 2012; pp. 378–390.
- 36. Rescorla, E.; Dierks, T. RFC 8446: The Transport Layer Security (Tls) Protocol Version 1.3. Available online: https://tools.ietf.org/ html/rfc8446 (accessed on 11 November 2020).
- Aas, J.; Barnes, R.; Case, B.; Durumeric, Z.; Eckersley, P.; Flores-López, A.; Halderman, J.A.; Hoffman-Andrews, J.; Kasten, J.; Rescorla, E.; et al. Let's Encrypt: An Automated Certificate Authority to Encrypt the Entire Web. In *Proceedings of the 2019* ACM SIGSAC Conference on Computer and Communications Security, London, UK, November 2019; ACM: New York, NY, USA, 2019; pp. 2473–2487.
- 38. Sidecar Pattern. Available online: https://docs.microsoft.com/en-us/azure/architecture/patterns/sidecar (accessed on 9 November 2020).
- 39. Gateway Offloading Pattern. Available online: https://docs.microsoft.com/en-us/azure/architecture/patterns/gateway-offloading (accessed on 9 November 2020).
- 40. Asynchronous Background Processing and Message. Available online: https://docs.microsoft.com/en-us/dotnet/architecture/ serverless/serverless-design-examples#asynchronous-background-processing-and-messaging (accessed on 11 November 2020).
- 41. Bulkhead Pattern. Available online: https://docs.microsoft.com/en-us/azure/architecture/patterns/bulkhead (accessed on 9 November 2020).
- 42. AI and Data for the Benefit of Humanity. Available online: https://xprize.org/aialliance (accessed on 9 November 2020).
- 43. Bergmayr, A.; Breitenbücher, U.; Ferry, N.; Rossini, A.; Solberg, A.; Wimmer, M.; Kappel, G.; Leymann, F. A Systematic Review of Cloud Modeling Languages. *ACM Comput. Surv.* 2018, *51*, 1–38. [CrossRef]
- Opara-Martins, J.; Sahandi, R.; Tian, F. Critical review of vendor lock-in and its impact on adoption of cloud computing. In Proceedings of the International Conference on Information Society (i-Society 2014), London, UK, 10–12 November 2014; pp. 92–97.
- 45. Kratzke, N.; Quint, P.C.; Palme, D.; Reimers, D. Project Cloud TRANSIT—Or to Simply Cloudnative Applications Provisioning for SMEs by Integrating Already Available Container Technologies. In *Proceedings of the European Project Space on Smart Systems, Big Data, Future Internet—Towards Serving the Grand Societal Challenges, Rome, Italy, 2016;* Kantere, V., Koch, B., Eds.; SciTePress: Setubal, Portugal, 2017; pp. 2–26.
- 46. Liu, S.; Schmitt, P.; Bronzino, F.; Feamster, N. Characterizing Service Provider Response to the COVID-19 Pandemic in the United States. *arXiv* **2020**, arXiv:2011.00419.