

Position Paper: Edge Clouds Multiple Control Planes Data Replication Challenges

Bohdan Dobrelia
OpenStack platform
Red Hat
Poznan, Poland
bdobreli@redhat.com

Abstract—Fog computing is an emerging paradigm aiming at bringing cloud functions closer to the end users and data sources. Its control plane and deployments data synchronization is a major challenge. Autonomy requirements expect even the most distant edge sites always manageable, available for monitoring and alerting, scaling up/down, upgrading and applying security fixes. Whenever temporary disconnected sites are managed locally or centrally, some changes and data need to be eventually synchronized back to the central site(s) with having its merge-conflicts resolved for the central data hub(s). While some data needs to be pushed from the central site(s) to the Edge, which might require resolving data collisions at the remote sites as well. In this paper, we position the outstanding data synchronization problems for OpenStack platform becoming a cloud solution number one for fog computing. We define the inter-cloud operational invariants based on that Always Available autonomy requirement. We show that a causal consistent data replication is the best match for the outlined operational invariants and there is a great opportunity for designing such a solution for Edge clouds. Finally, the paper brings vision of unified tooling to solve outstanding state synchronization problems the same way for infrastructure owners, cloud operators and tenants running stateful workloads hosted on OpenStack IaaS or OpenShift/Kubernetes PaaS deployed in Edge clouds as multi-cloud workloads abstraction and unification layer, to make it truly cloud-vendors agnostic and portable.

Index Terms—Open source software, Edge computing, Distributed computing, System availability, Design

I. INTRODUCTION

OpenStack is an Infrastructure-as-a-Service platform number one for private cloud computing, and it becomes being so for fog computing as well. Hybridization and multi-cloud trends for private clouds interconnected with public clouds and Platform-as-a-Service (PaaS) solutions, like OpenShift/Kubernetes, allow the containerization of micro-services oriented workloads to emerge in a highly portable, self-contained and the hosting cloud-agnostic way. Giving it massively distributed scale of fog computing and bringing the data it operates closer to end users, opens great opportunities for Internet of Things (IoT) and 5G telecommunication technologies, which first of all requires low-latency and highly responsive interfaces always available for end users.

Speaking of an always available, back to the system administration realities, the Edge clouds control and management plane capabilities in such a massively distributed world shall not fall behind as well. In Fog environments, expectations for

those capabilities are very different to the traditional cloud environments. There you have a highly available control plane, and a management layer for your cloud. From the geographical standpoint, it is also a mostly still layout. Control/management nodes are hosted in well suited datacenters and communicate over low-latency reliable connections, where network partitions are considered as disasters. While in Fog clouds, being offline is rather a normal state of things, where state can be synchronized eventually. Think of autonomous unmanned drones controlled via compact mobile datacenters that cannot always maintain its upstream links up. In 5G networks, smart mobile applications, like speech or video recognition, go far beyond the computing powers of tiny devices running it. So data have to travel with its consumers to be always at hand and accessed via a low-latency connection from a nearest compact datacenter. It has to be compact as maintaining thousands of classic datacenters won't fit neither the most generous business model nor if limited space/weight et al requirements. And such autonomy imposes multiple control and configuration management planes to maintain local operations and manageability of systems while being offline. In turn, that brings outstanding problems of data replication and consistency.

II. BACKGROUND CONCEPTS

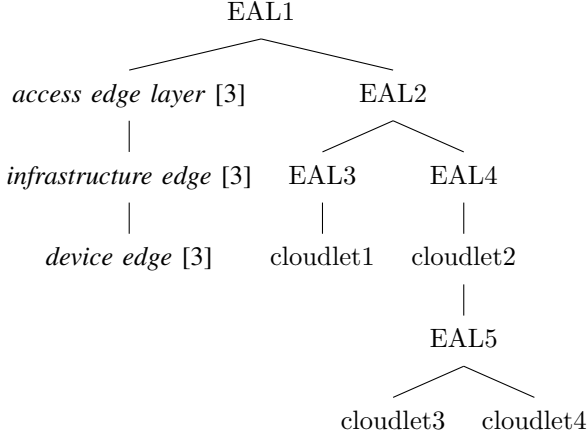
Aside of the established terms [3], we define a few more for the data processing and operational aspects:

Deployment Data: data that represents the configuration of *cloudlets* [3], like API endpoints URI, or numbers of deployed *edge nodes* [3] in *edge clouds* [3]. That data may represent either the most recent state of a deployment, or arbitrary data chunks/operations required to alter that state. When there is unresolved data merging conflicts or queued operations pending for future execution, the most recent state becomes the *best known* one.

Cloud Data: similarly to deployment data, represents the most recent or the best known internal and publicly visible state of cloudlets, like cloud users or virtual routers. Cloud data also includes logs, performance and usage statistics, state of message queues and the contents of databases. It may as well represent either data chunks or operations targeted for some state S transitioning into a new state S' .

Control Plane: corresponds to any operations performed via cloudlets API endpoints or command line tooling. E.g. starting

Fig. 1. Example Replication Topology



a virtual machine instance or a Kubernetes pod, or creating an OpenStack Keystone user. Such operations are typically initiated by cloud applications, tenants or operators.

Replication Topology: represents a hierarchy for allowed state replication flows and targets for operations on interconnected cloudlets, including *edge aggregation layers* [3]. So effectively there is only a limitation for horizontally interconnected cloudlets cannot be replicating its state nor targeting operations to each other. This corresponds to the acyclic graph (tree) topology.

For the example graph (Fig. 1), the infrastructure edge can replicate data, like contributing its local stats into the global view of hardware and performance metrics, to the the edge aggregation layer EAL1 omitting the access edge layer. Meanwhile the latter can be pushing something, like deployment data changes, to the device edge and infrastructure edge. On the right side of the graph, cloudlet3 and cloudlet4 cannot replicate to each other, but can do to the edge aggregation layer EAL1 either directly or consequently via EAL5, cloudlet2, EAL4 and EAL2, which is totally the replication implementation specific. We also say that EAL2 and EAL1 are hierarchically upper associated with EAL4 and EAL5, or just EAL1 is an *upper layer* for EAL2, when we mean that the latter is a potential subordinate of the former. Finally, we can say that EAL4 controls/manages cloudlet2 and anything sitting down of it.

That is, a control/manage-subordinate hierarchical association only defines influence domains and does not imply state replication as a mandatory thing nor imposes bidirectional data synchronization, but rather provides an opportunity of one-way or bidirectional state replications and regulates a possibility of issuing control/management operations. State replication may be performed from subordinates to its upper associates, think of sending reports. Or vice versa, think of sending directives to subordinates. While control/management operations may only flow up-down, think of planning future work for subordinates.

Management Plane: corresponds to administrative actions performed via configuration and lifecycle management sys-

tems. Such operations are typically targeted for cloudlets, like edge nodes, *edge datacenters* [3], or edge clouds. E.g. upgrading or reconfiguring cloudlets in a *virtual datacenter* [3], or scaling up edge nodes. And typically initiated by cloud infrastructure owners. For some cases, like Baremetal-as-a-Service, tenants may as well initiate actions executed via the management plane. Collecting logs, performance and usage statistics for monitoring and alerting systems also represents the management plane operations, although it operates with the cloud data.

When we refer to an edge aggregation layer and cloudlets under its control/management, we mean exactly any operations executed via the control/management planes of that edge aggregation layer and targeted for cloudlets sitting down the nested connections graph. And a replication topology regulates allowed targets for operations and state replication.

Data Replication Conflict: according to [1], two operations on the same target are in conflict if they are not related by causality.

Always Available: the operational mode of the control and management planes that corresponds to *sticky available causal consistency* [4] data replication models, i.e. RTC (*Real-Time Causal* [2]), or *causal+* [1]. Depending on the consistency model choices, there may be additional constraints:

- the stickiness property is a mandatory constraint for sticky available causal consistent systems. That is: “on every non-faulty node, so long as clients only talk to the same servers, instead of switching to new ones” [4].
- the real-time constraint is keeping the system time synchronized for all cloudlets. That is a mandatory constraint for RTC.
- *one-way convergence* [2] is a mandatory for RTC.

Causal+ and RTC consistency ensure ordering of relative operations, i.e. all causal related writes can be seen in the same order by all processes (connected to the same server). It is also known that “the causal consistency supports non-blocking operations; i.e. processes may complete read or write operations without waiting for global computation. Therefore, the causal consistency overcomes the primary limit of stronger criteria: communication latency” [6]. All that provides the best causal consistency guarantees we can get for today for always available systems.

III. ANALYSIS AND DISCUSSION

A. Autonomy Requirements

We define always available autonomy for cloudlets as the following strict requirements:

1) *Maintain multiple control/management planes:* Autonomous cloudlets may only operate as always available when having multiple control/management planes. For example, when an offline cloudlet cannot start virtual machines or containers, that violates the autonomy requirements.

2) *Provide no read-only or blocking access for inter-cloudlet operations:* Any operations performed on cloudlets state, despite its aliveness/failure conditions as it’s shown in

the global view of upper edge aggregation layers, fit consistency models that allow the involved control/management planes operating as always available, and there is no limitations, like read-only or blocking access. Operations may be queued for future processing in order to meet this autonomy requirement though. Internal state of cloudlets is allowed to keep its failure modes unchanged. E.g. for a local or distributed Galera/MariaDB database or a Etcd cluster, its standard quorum requirements apply for internal data transactions. That is a transitioning requirement until the internal cloudlets state can be migrated to causal consistent data storages as well.

3) *Provide a local control/management plane, whenever possible:* Operations on cloudlets can be scheduled at any given moment of time. For offline/partitioned autonomous cloudlets, that can be done via local control/management plane, if it exists and not failed. The same transitioning considerations for internal cloudlets state apply.

4) *Support fully autonomous (offline) cloudlets:* Aggregation edge layer cloudlets should allow for arbitrary or all of its managed/controlled cloudlets running fully autonomous long-time or indefinitely long, queuing any operations targeted for such cloudlets. That is, to have failure domains size of a 1. For a permanently disconnected cloudlet it may make more sense though to detach it from its adjacent aggregation layer and/or reorganize its place taken in the replication topology.

5) *Queue operations to keep it always available at best effort:* Queued operations have to be eventually replied if/after the control/management plane capabilities restored, or dropped (e.g. expired) otherwise. That poses a lazy replication principle. If there is intermediate aggregation edge layers down the way to the target of the queued operations, the queue may be shared across each of the involved aggregation edge layers or optionally, queued operations may be distributed across not shared queues. That should reduce the associated memory and disc pressure for aggregation layers.

6) *Provide a global view for cloudlets aliveness state:* Global view of cloudlets needs to be periodically presented for at least one of upper aggregation edge layers. For example, with the state marks, like “unknown”, or “autonomous”; “synchronizing”, “connected”; “failed”, or “disconnected”, or “fenced”, if and only if it is confirmed as failed, or manually disconnected, or fenced automatically. That poses the aliveness of the control/management planes principle.

B. Operational Invariants

To be always available as we defined it, control and management planes of cloudlets should provide the following operational capabilities (*invariants* hereafter):

1) *Keep control planes always available at best effort:* CRUD (Create, Read, Update and Delete) operations on cloud data can always be requested via API/CLI of local cloudlets or upper layers. The same queuing requirements apply as it is defined for the autonomy requirements.

2) *Do not wait for edge aggregation layers control planes for local operations:* Local CRUD operations for offline cloudlets, if its control plane exists and not failed, can be

processed without waiting for the upper aggregation layers to recover its control over the cloudlets. When a cloudlet is running only compute/storage resources, it cannot meet this requirement.

3) *Allow local scaling of infrastructure edge nodes without waiting for management planes of upper aggregation layers:* Similarly the to control plane operations, deployed infrastructure edge nodes can always be scheduled for scaling up/down by the cloudlets local management planes, if it is possible (a cloudlet may be relying on the remote configuration management only), or via its upper layers. Same queuing requirements apply for operations.

4) *Allow hotfixes and kernel/software updates applied locally for cloudlets:* Security patches and minor system updates, including kernel upgrades, can always be scheduled for installation by the same meanings (via operations issued locally or by the associated aggregation layers, including the same queuing requirements).

5) *Allow major software versions upgrades applied locally for cloudlets:* Similarly, major versions of system components, like OpenStack or OpenShift/Kubernetes platforms, can be always scheduled for upgrades, using the same meanings as above.

6) *Provide an extended global view for cloudlets:* Additionally to the aforementioned global view for cloudlets control/management plane aliveness state marks, there needs to be a periodically updated global view for each of the edge aggregation layers into its controlled/managed cloudlets, at least the adjacent ones, for the key system administration aspects, like hardware status, power management, systems state logging, monitoring and alerting events, performance and metering statistics.

C. State Replication Consistency Requirements

As it follows from the defined always available autonomy requirements and operational invariants, we define the following data replication requirements¹:

1) *Incorporate convergent conflict handling [1]:* Data replication conflicts can be resolved automatically, or by hand and maintained as causal related. The conflicts resolving strategies and rules should be customizable, like “last writer wins” or “return them all”. After the conflicts resolved, the data may be considered causal related, that is by definition [1] of the data conflicts in eventually consistent systems.

2) *Prefer one-way convergence in replication topologies:* As far as the replication topology and queuing capabilities allow that, causal related data can be replicated across cloudlets. Prefer one-way convergence and avoid bidirectional replication whenever possible.

3) *Bidirectional replication is only a nice to have requirement:* Bidirectional (two-way convergence) data replication is not a strict requirement but is nice to have. Indeed, some

¹when we refer to just *data* or *state*, we intentionally do not differentiate either that is deployment or cloud data, or queued API/CLI operations, to be replicated/replayed for management or control planes. That poses the **unified approach principle**

state needs to be replicated one-way from aggregation edge layer to cloudlets under its control/management, like virtual machine or hardware provisioning images data. While logs, performance and metering statistics may be collected only from cloudlets to its upper layers.

OpenStack/Kubernetes, have yet support for neither causal consistent storage backends for its cloud/deployment data, nor client libraries that could drive replication of casual related state. That poses an open opportunity for developers and system architects to design and implement such state replication tooling for multiple control/management planes.

OpenStack cloud data is normally stored in databases via transactions based on stronger than causal *unavailable* [4] data consistency models, e.g. *serializable* [4], or *repeatable read* [4]. OpenShift/Kubernetes clusters state, some of SDN (Software Defined Network) solutions are backed with Etcd, which also only supports the stronger than causal consistency models. Those cannot tolerate high network latency, serve two-way convergence only, therefore do not scale for potentially dozens of data replicas.

From the other side, the weaker consistency models, which is *total available* [4], provide a poor alternative that brings greatly increased implementation complexity. E.g. *monotonic atomic view* [4] would require handling of *fuzzy reads* [4], *phantoms* [4], discarded write-only transactions, empty state returned for any reads.

D. Data Replication Challenges

All that brings us to challenges that need to be addressed for multiple control/management planes:

- categorizing control/management operations and data flows associated with it, then grouping those into particular replication topologies. Such groups may be identified by multiple metrics, like communication latency, network partition duration tolerance for offline cloudlets, one- or two-way convergence based, a shared causal data storage or client libraries implementation specific. For the latter case, either it should allow replicating data low-level, or replaying queued operations at API level. Finally, that grouping should also be done for built-in state replication capabilities of components. Depending on the constraints for latency or maximum numbers of replicas, a distributed control/management plane may be a fit for some types of virtual data centers [3].
- in the end, the final design should not bring excessive operational overhead, like if maintaining all of the identified replication topologies simultaneously for a deployment, and require not too much of human care, but still meet the unified approach principle as we defined it earlier.
- picking/combining identified replication topologies to use with each of the involved system components, like identity provider or images serving services. Ideally, (almost) stateless components will not require data replication aside of built-in capabilities of a single distributed control plane.

- designing strategies and rules for conflicts resolving based on picked replication topologies identified the previous steps. A “last writer wins” may be a good fit for database/KVS conflicting data synchronization, while “return them all” seems the best choice for manual or artificial intelligence driven/assisted “smart” conflict resolvers.
- keeping state replication topologies efficient, e.g. distributing locally queued operations targeted for offline cloudlets by nesting more of the edge aggregation layers underneath. Also, stateless or almost stateless components may be happy with only built-in capabilities of a single distributed control plane.
- abiding the unified approach/architecture principles for Edge clouds IaaS, optional PaaS, and edge-native workloads as well.

E. Vision of a Unified Deployment/Cloud State Replication Design

The definition we made for always available distributed systems self-explains why the causal consistent state replication is the best match for the massively distributed cloudlets autonomy requirements and operational invariants as we defined those.

The vision of the unified architecture for future state replication tooling imposes it should be solving the multiple control/management planes data synchronization problems for IaaS, PaaS and end users consuming it as Replication-as-a-Service. Although generic version control systems, like Git, might fit all cases for deployment data replicating and conflicts resolving, that would break the unified design approach for cloud data/state replication.

Client libraries implementing causal consistent data replication and customizable conflicts resolving rules may provide a unification layer for different underlying databases/KVS (Key Value Storage). The replication will be effectively acting as database/KVS-to-database/KVS data synchronization tooling syncing data at a database level. The main benefit for such an approach is no a shared data storage needed. Instead, the underlying local to cloudlets data storages may keep operating as is, share nothing and provide unavailable consistency models stronger than causal consistency. And cloudlets may keep using different solutions for its local data storages as far as the state replication tooling may support such backends.

Additionally, client libraries may replicate not only data but operations at an API level² as well, i.e. resolve conflicting operations on-fly, then apply the resulting causal related operations for its original targets, effectively replicating changes at an API level. Operations queued by the control/management planes, including those targeted for offline cloudlets, may be also processed that way.

²For OpenStack Nova, there had been an example for such an API level replication, that is a Cells V1 protocol. But it had been deprecated as real cells v1 deployments required constant human care and feeding operationally

IV. RELATED WORK

COPS [1] formally proves implementation of a client library and highly scalable tooling for causal+ data operations. By design, it does not impose any real-time constraints and supports a single edge datacenter failure. The real tooling made off that base, may be operating on top of the not shared local cloudlets databases, or KVS, that provide the stronger consistency guarantees by the costs of reduced availability for local services. That would work as weaker consistency guarantees work well, when built on top of the stronger ones, and provide an always available global view of cloudlets for upper layers. Replicating the state changes via causal related operations and conflicts resolving via custom handlers is that COPS covers as well.

Global causal consistent databases [6] describes alternative solutions, if one can serve on the large scale and provide two-way convergence over high latency networks. Or otherwise, multiple shared instances of it, each serving to its dedicated virtual datacenter. A causal consistent data storage is the most unified and also the most simple way of exposing it for edge-native workloads as Replication-as-a-Service. The downside is such a thing has to be supported as a control/management planes data backend for IaaS and PaaS itself. So locally, cloudlets would have to be switched to use that data backend as well.

RainbowFS [7] expands on the consistency models and building tools topics.

- **Open questions:** does COPS work for multiple datacenters failure events, up to failure domains of a size of a 1? Is COPS applicable for two-way convergent systems, in terms of [2], for bidirectional causal+ replications? Any open source implementations exist for RTC/COPS?
- **Open questions:** which generic replication cases can be performed one-way? That simplifies implementation a lot: “Although most implementations use bidirectional communication, the communication from the update-receiver to the update-sender is just a (significant) performance optimization used to avoid redundant transfers of updates already known to the receiver. One-way convergence is also important in protocols that transmit updates via delay tolerant networks” [2].
- **Open questions:** anything to propose for bidirectional state synchronization, without breaking the autonomy requirements? Alternatively, instead of bidirectionally replicating data, may an API-to-API based synchronization mechanisms become a replacement for such cases, like predictable cloud users/projects ID require no databases synchronized?
- **Open questions:** which of the existing causal consistent databases [6] or file systems [7] can be integrated as a global solution that fits all control/management cases for cloudlets and may benefit workloads consuming these as a service?

V. CONCLUSION

We defined autonomy requirements for multiple control/management planes of massively distributed Fog environments and imposed operational invariants off it. That brought us to consistency requirements for cloudlets state replication and associated challenges. We introduced a replication topology building concepts. Finally, we posed vision of key design principles, like queuing and lazy replication, aliveness of the control and management planes and a unified approach for the subject tooling. Possible solutions may be based on either shared causal consistent databases, file systems, or client libraries that replicate not shared local data low-level. Or as well client libraries that replicate operations at an API high-level.

We want to position the unification principle as the most important thing and the greatest opportunity for developers to do it “the right way”, which is to bring the best of two IaaS and PaaS worlds for end users whom such data replication tooling might benefit as a service, i.e. cloud tenants, infrastructure owners and operators. That also applies to any fog-based system and is not limited to OpenStack or Kubernetes ecosystems.

Finally, we have to admit that because of outstanding future research work and potentially huge amount of changes for IaaS/PaaS cloud components required to implement or adopt either of the aforementioned options for multiple control/management planes state replication, it perhaps should be a next iteration after a minimum viable product (MVP) done as a single distributed control plane separated from the management plane. That enables early implementations for a very restricted subsets of operational invariants without autonomy supported, but “retain workloads operational as the best effort”.

ACKNOWLEDGMENT

Thanks to all who reviewed this paper, including but not limited to: Reedip Banerjee, from Red Hat; Flavia Delicato and Paulo Pires, from Federal University of Rio de Janeiro. Special thanks to all of the participants of OpenStack Summit at Berlin, 2018.

REFERENCES

- [1] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen, “Don’t settle for eventual: Scalable causal consistency for wide-area storage with COPS,” Proc. 23rd ACM Symposium on Operating Systems Principles (SOSP 11), Cascais, Portugal, October 2011.
- [2] P. Mahajan, L. Alvisi, and M. Dahlin, “Consistency, availability, and convergence,” Technical Report TR-11-22, Univ. Texas at Austin, Dept. Comp. Sci., 2011.
- [3] The Linux Foundation, “Open Glossary of Edge Computing,” [Online]. Available: <https://github.com/State-of-the-Edge/glossary>
- [4] K. Kingsbury, “Consistency Models,” [Online]. Available: <https://jepsen.io/consistency>
- [5] M. Bayer, “Global Galera Database,” [Online]. Available: <https://review.openstack.org/600555>
- [6] M. M. Elbushra, J. Lindstrom, “Causal Consistent Databases,” Open Journal of Databases (OJDB), Volume 2, Issue 1, 2015.
- [7] PRCE (Projet de recherche collaborative Entreprises), “RainbowFS: Modular Consistency and Co-designed Massive File” [Online]. Available: <http://rainbowfs.lip6.fr/data/RainbowFS-2016-04-12.pdf>