
Vision for OpenStack Clouds

<https://review.openstack.org/592205> Patch Set 5

1 Purpose

This is a living document. Its purpose is to document the OpenStack community's vision for the output of the OpenStack project as a whole, as it evolves over time. It is also aspirational, not descriptive, in nature. That is, it describes the OpenStack that the community is committed to working toward, not the OpenStack that has existed at any particular point in time.

Project teams can refer to this document when evaluating proposed features and designing interfaces, to help ensure that their designs fit comfortably within a larger structure and contribute to the overall landscape of application deployment patterns.

Amongst other things, the Technical Committee will refer to this document when assessing new project applications to determine whether they fit with the overall technical direction of OpenStack. If you are working on a project that does not appear to fit with the vision as described here, that does not necessarily mean that your project can never become an official OpenStack project, but it may mean that you should submit a patch to this document to further flesh out the vision. If this is the case, you should do so as soon as possible - there is no need to wait until you are ready to submit a new project application.

The Technical Committee may also use this document to evaluate the suitability and priority of proposed project-wide goals.

2 Scope

The scope of this document is limited to the cloud services that an end-user interacts with. This corresponds to the main 'OpenStack' bucket and parts of the 'OpenStack Operations' bucket in the OpenStack project map. While OpenStack also has other kinds of official projects (e.g. deployment tools and client libraries), nothing at all can be inferred about our vision for them from this document.

3 The Pillars of Cloud

There are at least as many different opinions of what 'cloud' means as there are software developers. However, we can all agree that cloud does mean *something*. Cloud computing promotes more efficient utilization of resources by reducing the transaction costs involved in provisioning and deprovisioning infrastructure to near zero, and it is able to do so because it differs in qualitative ways from previous models of computing (including virtualization). We identify two in particular.

These concepts are not always applicable to all aspects of the system, but we expect all services in OpenStack to conform to them wherever they are applicable, either directly or by working in conjunction with other services.

3.1 Self-service

Clouds are self-service. They provide users with the ability to deploy applications on demand without having to wait for human action or review in the loop. The cloud has no ticket trackers. This requirement has a number of immediate consequences.

First, cloud services must provide robust multi-tenancy. In order to securely serve multiple users (or groups of users) without any human review, resources must be isolated between tenants of the system so that the resources controlled by one tenant have neither access to nor impact upon resources controlled by other tenants.

Cloud services must also have some mechanism to ensure that capacity is only utilized when the value to the user of doing so exceeds the opportunity cost to the operator of providing it. In public clouds this is typically accomplished by charging users for the resources consumed. Quotas are used to limit risk for both users and cloud operators. For

users, the risk is of an unexpectedly-high bill, particularly when the level of resource consumption is partially under automatic control. For operators it is reaching a level of utilization where opportunity costs increase non-linearly. Private clouds will often require the same sorts of monitoring and reporting capabilities used for billing, even if they do not make use of a chargeback mechanism and rely on quotas alone as the sole technical measure to control resource consumption.

3.2 Application Control

Clouds allow control of an application's infrastructure to be vested in the application itself. Just as clouds eliminate the need for a human approver in the loop, they also eliminate the need for a human user to be in the loop. While a cloud may have a user interface (graphical or otherwise), it *must* have an application programming interface. It should supply operationally relevant information in a form that is legible to applications, including event notifications where appropriate. It should also be designed to facilitate secure access to its APIs for applications that are running within the cloud itself, because no part of the application should need to reside outside of the cloud.

4 OpenStack-specific Considerations

Most proprietary clouds are operated by software designed by and for a single organisation. OpenStack is different - there are many OpenStack clouds, both public and private, each operated by a different organisation with different goals and making different decisions. These clouds may have overlapping sets of users. This leads to some requirements that are specific to OpenStack, and may not be shared by other clouds.

These concepts are not always applicable to all aspects of the system, but we expect all projects in OpenStack to conform to them wherever they are applicable.

4.1 Interoperability

The same application descriptions (in whatever form they may take) should be deployable, with only minimal, well-isolated modification, to a variety of public and private OpenStack clouds. Where requirements or common OpenStack deployment patterns differ (for example between public and private cloud), OpenStack should endeavour to design mechanisms that are usable in either environment so that applications can be ported from one OpenStack cloud to another.

4.2 Bidirectional Compatibility

When interacting with a single cloud, it is a reasonable assumption that changes in the cloud software are monotonic, proceeding from older to newer versions. When users interact with multiple clouds, however, this assumption no longer holds. A user interacting with a newer version of OpenStack may then move to another cloud running an older version of OpenStack. OpenStack services should therefore evolve in such a way that they either work correctly or fail gracefully with *both* older and newer clients.

4.3 Cross-Project Dependencies

Not every OpenStack cloud includes the same set of services, and deploying and managing new services requires additional work from cloud operators. Although we encourage services to reuse functionality from other services, we do not *require* them to maximise reuse by adding hard dependencies. Choosing to add a hard dependency always involves a trade-off between design simplicity versus operational flexibility. Projects should add a hard dependency when they judge it to be ultimately beneficial to users - for example, by reducing the surface area of security-sensitive code, reducing the possibility of duplicate bugs, enabling desirable properties such as scalability or resilience, or increasing the development speed of the team. Particularly high weight should be placed on security benefits to operators and users.

Soft dependencies, where a particular feature is only available in the presence of an optional service, represent a good solution to the trade-off in many cases.

4.4 Partitioning

A region in an OpenStack cloud is defined as a separate set of service endpoints in the Keystone service catalog, but a shared Keystone - allowing a registered user to access any region of the cloud starting from the same authentication URL. This meaning is controlled by the OpenStack software, and therefore tends to be consistent across clouds.

In contrast, groupings of resources that are defined by hardware or the physical topology of the data center are under the control of individual cloud operators. For example, many clouds include the concept of 'availability zones' - groupings within a region that share no common points of failure. The OpenStack software has no way to enforce this meaning across clouds, and there are numerous other reasons for cloud operators to want to group resources together. OpenStack projects are encouraged to move toward allowing operators to create arbitrary, hierarchical groupings of the resources they manage, and to avoid ascribing physical meanings to the groupings.

5 Design Goals

The following are design goals that we would like to see the OpenStack services as a whole fulfil. It is not expected that every service or feature would (or even could) conform to every category listed. Rather, any service that contributes to achieving one or more of the objectives below is likely to help further the mission of the OpenStack project.

5.1 Basic Physical Data Center Management

OpenStack does not assume the existence of an operating data center; it provides the tools to operate a data center and make its resources available to consumers. There is no required layer underneath OpenStack as a whole. It provides the abstractions needed to deal with external systems like compute, storage, and networking hardware, the Domain Name System, and identity management systems. The OpenStack APIs provide a consistent interface to these systems, which may each potentially be implemented by a variety of vendors and Open Source projects.

This broad base provides an abstraction that can host more specialised services - both those within OpenStack itself and those from third-parties.

5.2 Plays Well With Others

OpenStack supports and encourages additional layers of abstraction - including Platforms as a Service, Serverless compute platforms, and Container Orchestration Engines - between itself and end-user applications, running within the compute capacity that it provisions. OpenStack should provide tools to support the tight integration into an OpenStack cloud of popular third-party open source projects that provide these layers.

OpenStack projects that include an abstraction layer over multiple potential back-end services may also offer that abstraction layer as a standalone entity, to be consumed by external services independently of an actual OpenStack cloud.

5.3 Hardware Virtualisation

For any service that is typically provided by a specialised piece of hardware, OpenStack aims to provide a vendor-independent API that gives consumers software-defined control of allocating the resource in a multi-tenant environment. This is not limited to virtual servers, but may (for example) also include such things as storage, routers, load balancers, firewalls, HSMs, GPGPUs, FPGAs, ASICs (e.g. video codecs), and so on.

Some of these hardware categories may have pure-software equivalents that can be used behind the same API, allowing applications to be portable even to clouds that don't have specialised hardware in those cases.

5.4 Infinite, Continuous Scaling

OpenStack strives to provide application developers with interfaces that allow them, in principle, to scale efficiently from very small to very large workloads without rearchitecting their applications.

In part, this means allowing consumers to use capacity as needed and share the underlying resources with other applications and tenants, in preference to allocating discrete chunks to particular applications and wasting any excess capacity within the chunks that they do not utilize.

5.5 Built-in Reliability and Durability

In an environment full of unreliable (that is to say, real) hardware, making an application reliable is difficult and, for smaller applications in particular, expensive. (Typically components running across a minimum of three hypervisor nodes are required.) OpenStack aims to provide primitives (for example, reliable delivery of messages and durable storage) that allow developers to build reliable applications on top of it. The underlying resources can be shared between applications and tenants so that the cost is amortized across them, rather than requiring each application to pay the full cost.

The existence of these primitives allows some other services to be simpler and more scalable, for example by making use of eventual consistency.

5.6 Customisable Integration

OpenStack does not impose any particular deployment model or architecture on applications. Every application has unique requirements, and OpenStack accommodates them by allowing services to be wired together in 'userspace' - through public APIs - rather than hard-wired actions taken behind the scenes that support only pre-defined deployment models.

This allows the application developer to customise anything using client-side glue, but should not require it. OpenStack services should be sufficiently integrated that they can be connected together by the cloud consumer without requiring any client-side interaction beyond the initial wiring.

Security models must allow both kinds of interaction - between OpenStack services, and between applications and OpenStack services in both directions. They should also permit the cloud consumer to delegate only the minimal privileges necessary to allow the application to operate as designed, and allow for regular revocation and replacement of credentials to maintain as much security as possible in an environment where Internet-facing machines are likely to eventually be compromised.

5.7 Abstract Specialised Operations

Certain components of an application - for example, databases - often benefit from a specialist skillset to operate them. By abstracting the management of some of the most common of these components behind an API, OpenStack allows the relationship between them and the rest of the application to be formalised. For organisations that have access to the necessary specialists, this allows those specialists to cover more applications by working in a centralised manner. For other organisations, it allows them to access the specialised skills that they otherwise could not, via a public or managed OpenStack cloud.

Not every reusable component of an application warrants its own OpenStack service. Suitable candidates typically feature complex configuration, ongoing lifecycle management needs, and sophisticated OpenStack infrastructure requirements (such as managing clusters of virtual servers).

5.8 Graphical User Interface

A GUI is often the best way for new users to approach a cloud and for users in general to experiment with unfamiliar areas of it. Presenting options and workflows graphically affords discovery of capabilities in a way that reading through API or CLI documentation cannot. A GUI is also often the best way for even experienced users and cloud operators to get a broad overview of the state of their cloud resources, and to visualise relationships between them. For these reasons, in addition to the API and any other user interfaces, OpenStack should include a web-based graphical user interface.