
Promenade Documentation

Release 0.1.0

Promenade Authors

Jul 11, 2019

Contents

1	User's Guide	3
1.1	Promenade Configuration Guide	3

Promenade is a tool for bootstrapping a resilient Kubernetes cluster and managing its life-cycle.

1.1 Promenade Configuration Guide

1.1.1 Developer On-boarding

Overview

Functionality:

- Airship Bootstrapping
- Core Kubernetes Management
- Misc.

Code structure:

1. Jinja templates (promenade/templates/**).
2. Helm charts for Kubernetes components, etcd, CoreDNS, Promenade charts/**.
3. Python support code.
 - API
 - Config access object
 - CLI
 - Certificate generation code

Since Promenade is largely templates + charts, unit testing is not enough to provide confidence for code changes. Significant functional testing is required to test small changes to, e.g. the `etcd` or `calico` charts, which can completely break deployment (or break reboot recovery). Developers can run the functional tests locally:

```
./tools/setup_gate.sh # Run once per machine you're on, DO NOT USE SUDO.  
./tools/gate.sh      # Runs a 4 node resiliency test.
```

This runs the test defined in `tools/g2/manifests/resiliency.json`. There are a few additional test scenarios defined in adjacent files.

There are helpful tools for troubleshooting these gates in `tools/g2/bin/*`, including `tools/g2/bin/ssh.sh`, which will let you ssh directly to a node to debug it, e.g.:

```
./tools/g2/bin/ssh.sh n0
```

Bootstrapping

Promenade is responsible for converting a vanilla Ubuntu 16.04 VM into a proper Airship.

How You Run It

Assuming you have a [valid set of configuration](#). Generate `genesis.sh`, which is a self-contained script for bootstrapping the genesis node.

```
promenade build-all -o output-dir config/*.yaml
```

What `genesis.sh` does:

1. Basic host validation (always room for more).
2. Drops pre-templated files in place:
 - Manifests to run initial Kubernetes components `/etc/kubernetes/manifests`
 - Basic components (apiserver, scheduler, controller-manager)
 - Etcd
 - Auxiliary Etcd
 - Docker configuration
 - Kubelet configuration
 - Apt configuration (proxy)
 - Bootstrapping Armada configuration
 - Dedicated Tiller
 - Dedicated Kubernetes API server
 - API server points at auxiliary etcd.
3. Installs some apt packages (docker + user-defined)
4. Starts Docker and Kubernetes.
5. Waits for bootstrapping services to be up (healthy Kubernetes API).
6. Applies configured labels to node.
7. Waits for Armada to finish bootstrapping deployment.
8. Final host validation.

When it's done, you should have a working Airship deployed as defined by your configuration (e.g. with or without LMA, keystone, etc) with no configuration loaded into Deckhand (via Shipyard).

How It Works

The templates that get dropped in place generally live in `promenade/templates/**`. The genesis node gets everything under `roles/genesis/**` and `roles/common/**` directly in place. Note that the templates under `roles/join/**` are used instead of the files under `genesis` for joining nodes to the existing cluster.

The “real” work happens inside kubelet managed “static” pods (defined by flat files in `/etc/kubernetes/manifests`), primarily via Armada.

Charts do a bunch of work to take control of essentially everything behind the scenes. Trickiest is `etcd`, for which we run multiple server processes to keep the cluster happy throughout bootstrapping + initial node join.

Note that we deploy two separate `etcd` clusters: one for Kubernetes itself, and one for Calico. The Calico one is a bit less sensitive.

Anchor Pattern

To provide increased resiliency, we do something a bit unusual with the core components. We run a `DaemonSet` for them which simply copy static `Pod` definitions into the `/etc/kubernetes/manifests` directory on the hosts (along with any supporting files/configuration). This ensures that these workloads are present even when the Kubernetes API server is unreachable. We call this pattern the `Anchor` pattern.

The following components follow this pattern:

- Kubernetes core components
 - API server
 - Scheduler
 - Controller Manager
- Kubernetes `etcd`
- Calico `etcd`
- HAProxy (used for API server discovery)

The HAProxy `DaemonSet` runs on every machine in the cluster, but the others only run on “master” nodes.

Kubernetes Cluster Management

Promenade is responsible for managing the Kubernetes lifecycle of nodes. That primarily consists of “joining” them to the cluster and adding labels, but also includes label updates and node removal.

Node Join

This is done via a self-contained script that is obtained by Drydock querying the Promenade API `GET /api/v1.0/join-scripts` (and providing a configuration link to Deckhand originally specified by Shipyard).

The join script is delivered to the node by Drydock and executed via a `systemd` unit. When it runs, it follows a similar pattern to `genesis.sh`, but naturally does not use any Kubernetes bootstrapping components or run Armada:

1. Basic host validation (always room for more).
2. Drops pre-templated files in place:
 - Docker configuration

- Kubelet configuration
 - Apt configuration (proxy)
3. Installs some apt packages (docker + user-defined)
 4. Starts Docker and Kubernetes.
 5. Waits for node to be recognized by Kubernetes.
 6. Applies configured labels to node.
 7. Final host validation.

After the node has successfully joined, the systemd unit disables itself so that it is not run again on reboot (though it would be safe to do so).

Other Management Features

Re-labeling and node removal API development has been delayed for other priorities, but is recently underway. While changing labels is generally easy, there are a few trickier bits around Kubelet and etcd management.

It is currently possible to fully de-label and remove a node from the cluster using a script that gets placed on each node (it requires `kubect1` so that must be in place), but that work is not exposed via API yet. The resiliency gate exercises this to reprovision the genesis node as a normal node.

Miscellaneous

Promenade does a few bits of additional work that's hard to classify, and probably don't belong in scope long term. Most notably is certificate generation.

Certificate generation is configured by the `PKICatalog` configuration document, which specifies the details for each certificate (CN, groups, hosts). Promenade then translates those requirements into calls to `cfssl`. The following will create a `certificates.yaml` file in `output-dir` containing all the generated certs:

```
promenade generate-certs -o output-dir config/*.yaml
```

If there are existing certs in `config/*.yaml`, then they will be used if applicable.

Troubleshooting

The context for this section is the functional gates described above. You can run them with:

```
./tools/gate.sh <gate_name>
```

When something goes wrong with this, you can ssh into individual nodes for testing (the nodes are named `n0` through `n3`):

```
./tools/g2/bin/ssh.sh <node_name>
```

When you get into a node and see various failures, or have an Armada error message saying a particular chart deployment failed, it is important to assess the overall cluster rather than just digging into the first thing you see. For example, if there is a problem with `etcd`, it could manifest as the Kubernetes API server pods failing.

Here is an approximate priority list of what to check for health (i.e. things higher up in the list break things lower down):

1. Kubernetes etcd

2. Kubernetes API Server
3. Other Kubernetes components (scheduler, controller-manager, kubelet).
4. Kubernetes proxy
5. Calico etcd
6. Calico node
7. DNS (CoreDNS)

For almost any other application, all of the above must be healthy before they will function properly.

1.1.2 Design

Promenade is a [Kubernetes](#) cluster deployment tool with the following goals:

- Resiliency in the face of node loss and full cluster reboot.
- Bare metal node support without external runtime dependencies.
- Providing a fully functional single-node cluster to allow cluster-hosted [tooling](#) to provision the remaining cluster nodes.
- [Helm](#) chart managed component life-cycle.
- API-managed cluster life-cycle.

Cluster Bootstrapping

The cluster is bootstrapped on a single node, called the genesis node. This node goes through a short-lived bootstrapping phase driven by static pod manifests consumed by `kubelet`, then quickly moves to chart-managed infrastructure, driven by [Armada](#).

During the bootstrapping phase, the following temporary components are run as static pods which are configured directly from Promenade's configuration documents:

- [Kubernetes](#) core components
 - `apiserver`
 - `controller-manager`
 - `scheduler`
- `Etcd` for use by the [Kubernetes](#) `apiserver`
- [Helm](#)'s server process `tiller`
- [CoreDNS](#) to be used for [Kubernetes](#) `apiserver` discovery

With these components up, it is possible to leverage [Armada](#) to deploy [Helm](#) charts to manage these components (and additional components) going forward.

Though completely configurable, a typical [Armada](#) manifest should specify charts for:

- [Kubernetes](#) components
 - `apiserver`
 - `controller-manager`
 - `proxy`

- scheduler

- Cluster DNS (e.g. [CoreDNS](#))
- [Etcd](#) for use by the [Kubernetes](#) `apiserver`
- A CNI provider for [Kubernetes](#) (e.g. [Calico](#))
- An initial under-cloud system to allow cluster expansion, including components like [Armada](#), [Deckhand](#), [Drydock](#) and [Shipyards](#).

Once these charts are deployed, the cluster is validated (currently, validation is limited to resolving DNS queries and verifying basic Kubernetes functionality including `Pod` scheduling log collection), and then the genesis process is complete. Additional nodes can be added to the cluster using day 2 procedures.

After additional master nodes are added to the cluster, it is possible to remove the genesis node from the cluster so that it can be fully re-provisioned using the same process as for all the other nodes.

Life-cycle Management

There are two sets of resources that require life-cycle management: cluster nodes and [Kubernetes](#) control plane components. These two sets of resources are managed differently.

Node Life-Cycle Management

Node life-cycle management tools are provided via an API to be consumed by other tools like [Drydock](#) and [Shipyards](#).

The life-cycle operations for nodes are:

1. Adding a node to the cluster
2. Removing a node from the cluster
3. Adding and removing node labels.

Adding a node to the cluster

Adding a node to the cluster is done by running a shell script on the node that installs the `kubelet` and configures it to find and join the cluster. This script can either be generated up front via the CLI, or it can be obtained via the `join-scripts` endpoint of the API (development of this API is in-progress).

Nodes can only be joined assuming all the proper configuration documents are available, including required certificates for `Kubelet`.

Removing a node from the cluster

This is currently possible by leveraging the `promenade-teardown` script placed on each host. API support for this function is planned, but not yet implemented.

Adding and removing node labels

Promenade provides `node-labels` API for updating node labels. For more information about updating node labels, please reference the [Promenade API Documentation](#).

It through relabeling nodes that key day 2 operations functionality like moving a master node are achieved.

Control-Plane Component Life-Cycle Management

With the exception of the `Docker` daemon and the `kubelet`, life-cycle management of control plane components is handled via `Helm` chart updates, which are orchestrated by `Armada`.

The `Docker` daemon is managed as an APT package, with configuration installed at the time the node is configured to join the cluster.

The `kubelet` is directly installed and configured at the time nodes join the cluster. Work is in progress to improve the upgradability of `kubelet` via either a system package or a chart.

Resiliency

The two primary failure scenarios Promenade is designed to be resilient against are node loss and full cluster restart.

`Kubernetes` has a well-defined `High Availability` pattern, which deals well with node loss.

However, this pattern requires an external load balancer for `apiserver` discovery. Since it is a goal of this project for the cluster to be able to operate without ongoing external dependencies, we must avoid that requirement.

Additionally, in the event of full cluster restart, we cannot rely on any response from the `apiserver` to give any `kubelet` direction on what processes to run. That means, each master node must be self-sufficient, so that once a quorum of `Etcd` members is achieved the cluster may resume normal operation.

The solution approach is two-pronged:

1. Deploy a local discovery mechanism for the `apiserver` processes on each node so that core components can always find the `apiservers` when their nodes reboot.
2. Apply the Anchor pattern described below to ensure that essential components on master nodes restart even when the `apiservers` are not available.

Currently, the discovery mechanism for the `apiserver` processes is provided by `CoreDNS` via a zone file written to disk on each node. This approach has some drawbacks, which might be addressed in future work by leveraging a `HAProxy` for discovery instead.

Anchor Pattern

The anchor pattern provides a way to manage process life-cycle using `Helm` charts in a way that allows them to be restarted immediately in the event of a node restart – even when the `Kubernetes` `apiserver` is unreachable.

In this pattern, a `DaemonSet` called the `anchor` that runs on selected nodes and is responsible for managing the life-cycle of assets deployed onto the node file system. In particular, these assets include a `Kubernetes` `Pod` manifest to be consumed by `kubelet` and it manages the processes specified by the `Pod`. That management continues even when the node reboots, since static pods like this are run by the `kubelet` even when the `apiserver` is not available.

Cleanup of these resources is managed by the `anchor` pods' `preStop` life-cycle hooks. This is usually simply removing the files originally placed on the nodes' file systems, but, e.g. in the case of `Etcd`, can actually be used to manage more complex cleanup like removal from cluster membership.

Pod Checkpointer

Before moving to the Anchor pattern above, the pod-checkpointer approach pioneered by the `Bootkube` project was implemented. While this is an appealing approach, it unfortunately suffers from race conditions during full cluster reboot.

During cluster reboot, the checkpointer copies essential static manifests into place for the `kubelet` to run, which allows those components to start and become available. Once the `apiserver` and `etcd` cluster are functional, `kubelet` is able to register the failure of its workloads, and delete those pods via the API. This is where the race begins.

Once those pods are deleted from the `apiserver`, the pod checkpointer notices that the flagged pods are no longer scheduled to run on its node and then deletes the static manifests for those pods. Concurrently, the `controller-manager` and `scheduler` notice that new pods need to be created and scheduled (sequentially) and begin that work.

If the new pods are created, scheduled and started on the node before pod checkpointers on other nodes delete their critical services, then the cluster may remain healthy after the reboot. If enough nodes running the critical services fail to start the newly created pods before too many are removed, then the cluster does not recover from hard reboot.

The severity of this race is exacerbated by:

1. The sequence of events required to successfully replace these pods is long (`controller-manager` must create pods, then `scheduler` can schedule pods, then `kubelet` can start pods).
2. The `controller-manager` and `scheduler` may need to perform leader election during the race, because the leader might have been killed early.
3. The failure to recover any one of the core sets of processes can cause the entire cluster to fail. This is somewhat trajectory-dependent, e.g. if at least one `controller-manager` is scheduled before the `controller-manager` processes are all killed, then assuming the other processes are correctly restarted, then the `controller-manager` will also recover.
4. `etcd` is somewhat more sensitive to this race, because it requires two successfully restarted pods (assuming a 3 node cluster) rather than just one as the other components.

This race condition was the motivation for the construction and use of the Anchor pattern. In future versions of [Kubernetes](#), it may be possible to use [built-in checkpointing](#) from the `kubelet`.

Alternatives

- [Kubeadm](#)
 - Does not yet support [HA](#)
 - Current approach to [HA Etcd](#) is to use the `etcd operator`, which recovers from cluster reboot by loading from an external backup snapshot
 - Does not support chart-based management of components
- [kops](#)
 - Does not support [bare metal](#)
- [Bootkube](#)
 - Does not support automatic recovery from a [full cluster reboot](#)
 - Does not yet support [full HA](#)
 - Adheres to different design goals (minimal direct server contact), which makes some of these changes challenging, e.g. [building a self-contained, multi-master cluster](#)
 - Does not support chart-based management of components

1.1.3 Getting Started

Note: This document is meant to give a general understanding of how Promenade could be exercised in a development environment or for general learning and understanding. For holistic Airship deployment procedures, refer to [Treasuremap](#)

Basic Deployment

This approach is quick to get started, but generates the scripts used for joining up-front rather than generating them in the API as needed.

Setup Build Machine

On the machine you wish to use to generate deployment files, install docker:

```
sudo apt -y install docker.io
```

This can be the same machine you intend to be the Genesis host, or it may be a separate build machine.

Generate Build files

To create the certificates and scripts needed to perform a basic deployment, you can use the following helper script on your build machine:

```
sudo ./tools/simple-deployment.sh examples/basic build
```

This will copy the configuration provided in the `examples/basic` directory into the `build` directory. Then, it will generate self-signed certificates for all the needed components in Deckhand-compatible format. Finally, it will render the provided configuration into directly-usable `genesis.sh` and `join-<NODE>.sh` scripts.

Genesis Host Provision

Install Ubuntu 16.04 on the machine intended to be the genesis host. Ensure the host has outbound internet access and DNS resolution. Ensure that the hostname matches the hostname specified in the `Genesis.yaml` file used to build the above configurations.

Execution

Perform the following steps to execute the deployment:

1. Copy the `genesis.sh` script to the genesis node and run it as `sudo`. In the event of runtime errors, refer to [Genesis Troubleshooting](#)
2. Validate the genesis node by running `validate-genesis.sh` on it.
3. Nodes for which `join-<NODE>.sh` scripts have been generated should be provisioned at this point, and need to have network connectivity to the genesis node. (This could be a manual Ubuntu provision, or a Drydock-initiated PXE boot in the case of a full fledged Airship deployment).
4. Join master nodes by copying their respective `join-<NODE>.sh` scripts to them and running them.

5. Validate the master nodes by copying and running their respective `validate-<NODE>.sh` scripts on each of them.
6. Re-provision the Genesis node
 - a) Run the `/usr/local/bin/promenade-teardown` script on the Genesis node:
 - b) Delete the node from the cluster via one of the other nodes `kubectl delete node <GENESIS>`.
 - c) Power off and re-image the Genesis node.
 - d) Join the genesis node as a normal node using its `join-<GENESIS>.sh` script.
 - e) Validate the node using `validate-<GENESIS>.sh`.
7. Join and validate all remaining nodes using the `join-<NODE>.sh` and `validate-<NODE>.sh` scripts described above.

API-Driven Deployment

This approach leverages the Promenade API to fetch join scripts as needed. This is the approach used in the functional testing discussed below.

Setup

Follow the setup instructions above for *Basic Deployment*. Then, start a webserver to serve configuration to Promenade.

```
cat build/*.yaml > promenade.yaml
mv promenade.yaml build/promenade.yaml
docker rm -fv promenade-nginx
docker run -d \
  -p 7777:80 \
  --restart=always \
  --name promenade-nginx \
  -v build:/usr/share/nginx/html:ro \
  nginx:stable
export DESIGN_REF=http://192.168.77.1:7777/promenade.yaml
```

Execution

Perform the following steps to execute the deployment:

1. Copy the `genesis.sh` script to the genesis node and run it.
2. Validate the genesis node by running `validate-genesis.sh` on it.
3. Generate join script for a node using:

```
URL=http://promenade-api.ucp.svc.cluster.local/api/v1.0/join-scripts?
URL="$URL&design_ref=${DESIGN_REF}"
URL="$URL&hostname=<HOSTNAME>&ip=<IP>"
URL="$URL&labels.dynamic=calico-etcd=enabled"
URL="$URL&labels.dynamic=kubernetes-apiserver=enabled"
URL="$URL&labels.dynamic=kubernetes-controller-manager=enabled"
URL="$URL&labels.dynamic=kubernetes-etcd=enabled"
URL="$URL&labels.dynamic=kubernetes-scheduler=enabled"
```

(continues on next page)

(continued from previous page)

```
URL="${URL}&labels.dynamic=ucp-control-plane=enabled"
curl -Lo join-<NODE>.sh "${URL}"
```

4. Copy the join script to the node, and run it via `bash join-<NODE>.sh`.
5. Repeat 3 and 4 until all nodes are joined.
6. Re provision the Genesis node by tearing it down as above in *Basic Deployment*, then generating and using a join script for it as done in 3 and 4.

Running Tests

Initial Setup of Virsh Environment

To setup a local functional testing environment on your Ubuntu 16.04 machine, run:

```
./tools/setup_gate.sh
```

Running Functional Tests

To run complete functional tests locally:

```
./tools/gate.sh
```

For more verbose output, try:

```
PROMENADE_DEBUG=1 ./tools/gate.sh
```

For extremely verbose output, try:

```
GATE_DEBUG=1 PROMENADE_DEBUG=1 ./tools/gate.sh
```

The gate leaves its test VMs running for convenience. To shut everything down:

```
./tools/stop_gate.sh
```

To run a particular set of functional tests, you can specify the set on the command line:

```
./tools/gate.sh <SUITE>
```

Valid functional test suites are defined by JSON files that live in `tools/g2/manifests`.

Utilities

There are a couple of helper utilities available for interacting with gate VMs. These can be found in `tools/g2/bin`. The most important is certainly `ssh.sh`, which allows you to connect easily to test VMs:

```
./tools/g2/bin/ssh.sh n0
```

Development

Using a Local Registry

Repeatedly downloading multiple copies images during development can be quite slow. To avoid this issue, you can run a docker registry on the development host:

```
./tools/registry/start.sh
./tools/registry/update_cache.sh
```

Then, the images used by the basic example can be updated using:

```
./tools/registry/update_example.sh
```

That change can be undone via:

```
./tools/registry/revert_example.sh
```

The registry can be stopped with:

```
./tools/registry/stop.sh
```

Building the image

To build the image directly, you can use the standard Docker build command:

```
docker build -t promenade:local .
```

To build the image from behind a proxy, you can:

```
export http_proxy=...
export no_proxy=...
docker build --build-arg http_proxy=$http_proxy --build-arg https_proxy=$http_proxy --
↳build-arg no_proxy=$no_proxy -t promenade:local .
```

For convenience, there is a script which builds an image from the current code, then uses it to generate certificates and construct scripts:

```
./tools/dev-build.sh examples/basic build
```

Using Promenade Behind a Proxy

To use Promenade from behind a proxy, use the proxy settings see [Kubernetes Network](#).

1.1.4 Configuration

Promenade is configured using a set of [Deckhand](#) compatible configuration documents and a bootstrapping [Armada](#) manifest that is responsible for deploying core components into the cluster.

Details about Promenade-specific documents can be found here:

Docker

Configuration for the docker daemon. This document contains a single *config* key that directly translates into the contents of the *daemon.json* file described in [Docker's configuration](#).

Sample Document

Here is a sample document:

```

schema: promenade/Docker/v1
metadata:
  schema: metadata/Document/v1
  name: docker
  layeringDefinition:
    abstract: false
    layer: site
data:
  config:
    live-restore: true
    storage-driver: overlay2

```

EncryptionPolicy

Encryption policy defines how encryption should be applied via Promenade. The primary use-case for this is to encrypt *genesis.sh* or *join.sh* scripts.

Sample Document

```

---
schema: promenade/EncryptionPolicy/v1
metadata:
  schema: metadata/Document/v1
  name: encryption-policy
  layeringDefinition:
    abstract: false
    layer: site
  storagePolicy: cleartext
data:
  scripts:
    genesis:
      gpg: {}
  ...

```

Scripts

The genesis and join scripts can be built with sensitive content encrypted. Currently the only encryption method available is gpg, which can be enabled by setting that key to an empty dictionary.

Genesis

Specific configuration for the genesis process. This document is a strict superset of the combination of *Kubernetes Node* and *HostSystem*, so only differences are discussed here.

Sample Document

Here is a complete sample document:

```
schema: promenade/Genesis/v1
metadata:
  schema: metadata/Document/v1
  name: genesis
  layeringDefinition:
    abstract: false
    layer: site
data:
  hostname: n0
  ip: 192.168.77.10
  armada:
    target_manifest: cluster-bootstrap
  labels:
    static:
      - calico-etcd=enabled
      - node-role.kubernetes.io/master=
    dynamic:
      - kubernetes-apiserver=enabled
      - kubernetes-controller-manager=enabled
      - kubernetes-etcd=enabled
      - kubernetes-scheduler=enabled
      - promenade-genesis=enabled
      - ucp-control-plane=enabled
  images:
    armada: quay.io/airshipit/armada:latest
    helm:
      tiller: gcr.io/kubernetes-helm/tiller:v2.14.0
    kubernetes:
      apiserver: gcr.io/google_containers/hyperkube-amd64:v1.11.6
      controller-manager: gcr.io/google_containers/hyperkube-amd64:v1.11.6
      etcd: quay.io/coreos/etcd:v3.0.17
      scheduler: gcr.io/google_containers/hyperkube-amd64:v1.11.6
  files:
    - path: /var/lib/anchor/calico-etcd-bootstrap
      content: ""
      mode: 0644
```

Armada

This section contains particular configuration options for bootstrapping with Armada. It currently only supports a single option: `target_manifest`, which specifies which `armada/Manifest/v1` to be used during Genesis.

Bootstrapping Images

Bootstrapping images are specified in the top level key `images`:

```

armada: <Armada image for bootstrapping>
helm:
  tiller: <Tiller image for bootstrapping>
kubernetes:
  apiserver: <API server image for bootstrapping>
  controller-manager: <Controller Manager image for bootstrapping>
  etcd: <etcd image for bootstrapping>
  scheduler: <Scheduler image for bootstrapping>

```

HostSystem

Sample Document

```

schema: promenade/HostSystem/v1
metadata:
  schema: metadata/Document/v1
  name: host-system
  layeringDefinition:
    abstract: false
    layer: site
data:
  files:
    - path: /opt/kubernetes/bin/kubelet
      tar_url: https://dl.k8s.io/v1.11.6/kubernetes-node-linux-amd64.tar.gz
      tar_path: kubernetes/node/bin/kubelet
      mode: 0555
  images:
    haproxy: haproxy:1.8.3
    helm:
      helm: lachlanevenson/k8s-helm:v2.14.0
    kubernetes:
      kubect1: gcr.io/google_containers/hyperkube-amd64:v1.11.6
packages:
  repositories:
    - deb http://apt.dockerproject.org/repo ubuntu-xenial main
  keys:
    - |-
      -----BEGIN PGP PUBLIC KEY BLOCK-----

      mQINBFWln24BEADrBl5p99uKh8+rpvqJ48u4eTt jeXAWbslJotmC/CakbNSqOb9o
      ddfzRvGveJVERT/Q/mlvEqgnyTQy+e6oEYN2Y2kqXceUhXagThnqCoxcEJ3+KM4R
      mYdoe/BJ/J/6rHOjq7Omk24z2qB3RU1uAv57iY5VGw5p45uZB4C4pNNsBJXoCvPn
      TGAs/7IrekFZDDgVraPx/hdiwopQ8NltSfZCyU/ jPpWFK28TR8yfVlzYFwibj5WK
      dHM7ZTqlA1tHIG+agyPf3Rae0jPMsHR6q+arXVwMccyOi+ULU0z8mHUJ3iEMlRrT
      X+80KaN/ZjibfsBOCjcfiJSB/acn4nxQQgNZigna32velafhQivsNREFeJpzENiG
      HOoyC6qVeOgKrRiKxzymj0FIMLru/iFF5pSwcQB7PYlt8J0G801AcPr6VCiN+4c
      NKv03SdvA69dCOj79Pu09I IvQsJXsSq96HB+TeEmmL+xSdpGtGdCJHHM1fDeCqkZ
      hT+RtBGQL2SEdWjxbF43oQopocT8cHvyX6Zaltn0svoGs+wX3Z/H6/8P5anog43U
      65c0A+64Jj00rNdr8j3lizhtQMRo892kGeQAaaxg4Pz6HnS7hRC+cOMHUU4HA7iM
      zHrouAdYeTZeZEQOA7SxtCME9ZnGwe2grxPXh/U/80WJGkzLFNcTKdv+rwARAQAB
      tDdEb2NrZXIgdUmVsZWZzZSBub29sICHyZWw1YXNlZG9ja2VyKSA8ZG9ja2VyQGRv
      Y2t1ci5jb20+IQI4BBMBAgAiBQJVpZ9uAhsVbG9sJCAcDAGYVCAIJCgsEFgIDAQIe
      AQIXgAAKCRD3YiFXLFJgnbRfEAC9Uai7Rv20QID1DogRzd+Vebg4ahyoUdj0CH+n
      Ak40RIoq6G26u1e+sdgjpCa8jF6vrX+smpgd1HeJdmpahUX0XN3X9f9qU9o+j9A4I
      1WDalRWJh+tP5Wnv2ySy6AwcP9QnjuBMRtTK27pk1sEMg9oJHK5p+ts8hlSC4S1

```

(continues on next page)

(continued from previous page)

```
uyMKH5NMVy9c+A9yqq9NF6M6d6/ehKfBFFLG9BX+XLBATvf1ZemGVHQusCQebTGv
0C0V9yqtdPdRWVIEhHxyNHATaVYOafTj/EF01DxLl6zDT6trRV5n9F1VCEh4Aa18
L5MxVPcIZVO7NHT2EkQgn8CvWjV3oKl2GopZF8V4XdJr190U/WDv/6cmfI08GkzD
YBHhS8ULWRFwGKobsSTyIvnbk4NtKdnTGyTJCQ8+6i52s+C54PiNgfj2ieNn6oOR
7d+bNCcG1CdOYY+ZXVOcsjl73UYvtJr00Rl/NpYERkZ5d/tzw4jZ6FCXgggA/Zxc
jk6Y1ZvIm8Mt8wLRFH9Nww+FVsCtaCXJLP8D1JLASMD9r15QS9Ku3u7ZNrr5HWXP
HXITX660jglyshch6CWeiUATqjIAzkeQom/kEnOrvJAtkypRJ59vYQOedZ1sFVEL
MXg2UCkD/FwojfnVtjzYaTCeGwFQeqzHmM241iuOmBYPeyTY5veF49aBJA1gEJOQ
TvBR8Q==
=Fm3p
-----END PGP PUBLIC KEY BLOCK-----
```

additional:

- curl
- jq

required:

- docker:** docker-engine=1.13.1-0~ubuntu-xenial
- socat:** socat=1.7.3.1-1

Files

A list of files to be written to the host. Files can be given as precise content or extracted from a tarball specified by url:

```
- path: /etc/direct-content
  content: |-
    This
    exact
    text
- path: /etc/from-tar
  tar_url: http://example.com/file
  tar_source: dir/file.txt
```

Images

Core Images

These images are used for essential functionality:

haproxy HAProxy is configured and used for Kubernetes API discovery during bootstrapping.

kubect1 Used for label application and validation tasks during bootstrapping.

Convenience Images

The `helm` image is available for convenience.

Packages

Repository Configuration

Additional APT repositories can be configured using the `repositories` and `keys` fields of the `SystemPackages` document:

repositories A list of APT source lines to be configured during genesis or join.

keys A list of public PGP keys that can be used to verify installed packages.

Package Configuration

The `required` key specifies packages that are required for all deployments, and the `additional` key allows arbitrary additional system packages to be installed. The `additional` key is particularly useful for installing packages such as *ceph-common*.

Kubelet

Configuration for the Kubernetes worker daemon (the Kubelet). This document contains two keys: `arguments` and `images`. The `arguments` are appended directly to the `kubelet` command line, along with arguments that are controlled by Promenade more directly.

The only image that is configurable is for the `pause` container.

Sample Document

Here is a sample document:

```
schema: promenade/Kubelet/v1
metadata:
  schema: metadata/Document/v1
  name: kubelet
  layeringDefinition:
    abstract: false
    layer: site
data:
  arguments:
    - --cni-bin-dir=/opt/cni/bin
    - --cni-conf-dir=/etc/cni/net.d
    - --eviction-max-pod-grace-period=-1
    - --network-plugin=cni
    - --node-status-update-frequency=5s
    - --v=5
  images:
    pause: gcr.io/google_containers/pause-amd64:3.0
```

Kubernetes Network

Configuration for Kubernetes networking during bootstrapping and for the kubelet.

Sample Document

```
schema: promenade/KubernetesNetwork/v1
metadata:
  schema: metadata/Document/v1
  name: kubernetes-network
  layeringDefinition:
```

(continues on next page)

```
abstract: false
layer: site
data:
  dns:
    cluster_domain: cluster.local
    service_ip: 10.96.0.10
    bootstrap_validation_checks:
      - calico-etcd.kube-system.svc.cluster.local
      - kubernetes-etcd.kube-system.svc.cluster.local
      - kubernetes.default.svc.cluster.local
    upstream_servers:
      - 8.8.8.8
      - 8.8.4.4

  kubernetes:
    apiserver_port: 6443
    haproxy_port: 6553
    pod_cidr: 10.97.0.0/16
    service_cidr: 10.96.0.0/16
    service_ip: 10.96.0.1

  etcd:
    container_port: 2379
    haproxy_port: 2378

  hosts_entries:
    - ip: 192.168.77.1
      names:
        - registry

  proxy:
    url: http://proxy.example.com:8080
    additional_no_proxy:
      - 192.168.77.1
```

DNS

The data in the `dns` key is used for bootstrapping and `kubelet` configuration of cluster and host-level DNS, which is provided by `coredns`.

bootstrap_validation_checks Domain names to resolve during the genesis and join processes for validation.

cluster_domain The Kubernetes cluster domain. Used by the `kubelet`.

service_ip The IP to use for cluster DNS. Used by the `kubelet`.

upstream_servers Upstream DNS servers to be configured in `/etc/resolv.conf`.

Kubernetes

The `kubernetes` key contains:

apiserver_port The port that the Kubernetes API server process will listen on hosts where it runs.

haproxy_port The port that HAProxy will listen on each host. This port will be used by the `kubelet` and `kube-proxy` to find API servers in the cluster.

pod_cidr The CIDR from which the Kubernetes Controller Manager assigns pod IPs.

service_cidr The CIDR from which the Kubernetes Controller Manager assigns service IPs.

service_ip The in-cluster Kubernetes service IP.

Kubernetes Node

Configuration for a basic node in the cluster.

Sample Document

Here is a sample document:

```
schema: promenade/KubernetesNode/v1
metadata:
  schema: metadata/Document/v1
  name: n1
  layeringDefinition:
    abstract: false
    layer: site
data:
  hostname: n1
  ip: 192.168.77.11
  join_ip: 192.168.77.10
  labels:
    static:
      - node-role.kubernetes.io/master=
    dynamic:
      - calico-etcd=enabled
      - kubernetes-apiserver=enabled
      - kubernetes-controller-manager=enabled
      - kubernetes-etcd=enabled
      - kubernetes-scheduler=enabled
      - ucp-control-plane=enabled
```

Host Information

Essential host-specific information is specified in this document, including the `hostname`, `ip`, and `join_ip`.

The `join_ip` is used to specify which host should be used when adding a node to the cluster.

Labels

Kubernetes labels can be specified under the `labels` key in two ways:

1. Via the `static` key, which is a list of labels to be applied immediately when the `kubelet` process starts.
2. Via the `dynamic` key, which is a list of labels to be applied after the node is marked as *Ready* by Kubernetes.

PKI Catalog

Configuration for certificate and keypair generation in the cluster. The `promenade generate-certs` command will read all `PKICatalog` documents and either find pre-existing certificates/keys, or generate new ones based on the given definition.

Sample Document

Here is a sample document:

```
schema: promenade/PKICatalog/v1
metadata:
  schema: metadata/Document/v1
  name: cluster-certificates
  layeringDefinition:
    abstract: false
    layer: site
data:
  certificate_authorities:
    kubernetes:
      description: CA for Kubernetes components
  certificates:
    - document_name: apiserver
      description: Service certificate for Kubernetes apiserver
      common_name: apiserver
      hosts:
        - localhost
        - 127.0.0.1
        - 10.96.0.1
      kubernetes_service_names:
        - kubernetes.default.svc.cluster.local
    - document_name: kubelet-genesis
      common_name: system:node:n0
      hosts:
        - n0
        - 192.168.77.10
      groups:
        - system:nodes
    - document_name: kubelet-n0
      common_name: system:node:n0
      hosts:
        - n0
        - 192.168.77.10
      groups:
        - system:nodes
    - document_name: kubelet-n1
      common_name: system:node:n1
      hosts:
        - n1
        - 192.168.77.11
      groups:
        - system:nodes
    - document_name: kubelet-n2
      common_name: system:node:n2
      hosts:
        - n2
```

(continues on next page)

(continued from previous page)

```

    - 192.168.77.12
    groups:
      - system:nodes
  - document_name: kubelet-n3
    common_name: system:node:n3
    hosts:
      - n3
      - 192.168.77.13
    groups:
      - system:nodes
  - document_name: scheduler
    description: Service certificate for Kubernetes scheduler
    common_name: system:kube-scheduler
  - document_name: controller-manager
    description: certificate for controller-manager
    common_name: system:kube-controller-manager
  - document_name: admin
    common_name: admin
    groups:
      - system:masters
  - document_name: armada
    common_name: armada
    groups:
      - system:masters
kubernetes-etcd:
  description: Certificates for Kubernetes's etcd servers
  certificates:
  - document_name: apiserver-etcd
    description: etcd client certificate for use by Kubernetes apiserver
    common_name: apiserver
  - document_name: kubernetes-etcd-anchor
    description: anchor
    common_name: anchor
  - document_name: kubernetes-etcd-genesis
    common_name: kubernetes-etcd-genesis
    hosts:
      - n0
      - 192.168.77.10
      - 127.0.0.1
      - localhost
      - kubernetes-etcd.kube-system.svc.cluster.local
  - document_name: kubernetes-etcd-n0
    common_name: kubernetes-etcd-n0
    hosts:
      - n0
      - 192.168.77.10
      - 127.0.0.1
      - localhost
      - kubernetes-etcd.kube-system.svc.cluster.local
  - document_name: kubernetes-etcd-n1
    common_name: kubernetes-etcd-n1
    hosts:
      - n1
      - 192.168.77.11
      - 127.0.0.1
      - localhost
      - kubernetes-etcd.kube-system.svc.cluster.local

```

(continues on next page)

(continued from previous page)

```
- document_name: kubernetes-etcd-n2
common_name: kubernetes-etcd-n2
hosts:
  - n2
  - 192.168.77.12
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
- document_name: kubernetes-etcd-n3
common_name: kubernetes-etcd-n3
hosts:
  - n3
  - 192.168.77.13
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
kubernetes-etcd-peer:
certificates:
- document_name: kubernetes-etcd-genesis-peer
common_name: kubernetes-etcd-genesis-peer
hosts:
  - n0
  - 192.168.77.10
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
- document_name: kubernetes-etcd-n0-peer
common_name: kubernetes-etcd-n0-peer
hosts:
  - n0
  - 192.168.77.10
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
- document_name: kubernetes-etcd-n1-peer
common_name: kubernetes-etcd-n1-peer
hosts:
  - n1
  - 192.168.77.11
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
- document_name: kubernetes-etcd-n2-peer
common_name: kubernetes-etcd-n2-peer
hosts:
  - n2
  - 192.168.77.12
  - 127.0.0.1
  - localhost
  - kubernetes-etcd.kube-system.svc.cluster.local
- document_name: kubernetes-etcd-n3-peer
common_name: kubernetes-etcd-n3-peer
hosts:
  - n3
  - 192.168.77.13
  - 127.0.0.1
  - localhost
```

(continues on next page)

(continued from previous page)

```

- kubernetes-etcd.kube-system.svc.cluster.local
calico-etcd:
  description: Certificates for Calico etcd client traffic
  certificates:
    - document_name: calico-etcd-anchor
      description: anchor
      common_name: anchor
    - document_name: calico-etcd-n0
      common_name: calico-etcd-n0
      hosts:
        - n0
        - 192.168.77.10
        - 127.0.0.1
        - localhost
        - 10.96.232.136
    - document_name: calico-etcd-n1
      common_name: calico-etcd-n1
      hosts:
        - n1
        - 192.168.77.11
        - 127.0.0.1
        - localhost
        - 10.96.232.136
    - document_name: calico-etcd-n2
      common_name: calico-etcd-n2
      hosts:
        - n2
        - 192.168.77.12
        - 127.0.0.1
        - localhost
        - 10.96.232.136
    - document_name: calico-etcd-n3
      common_name: calico-etcd-n3
      hosts:
        - n3
        - 192.168.77.13
        - 127.0.0.1
        - localhost
        - 10.96.232.136
    - document_name: calico-node
      common_name: calico-node
calico-etcd-peer:
  description: Certificates for Calico etcd clients
  certificates:
    - document_name: calico-etcd-n0-peer
      common_name: calico-etcd-n0-peer
      hosts:
        - n0
        - 192.168.77.10
        - 127.0.0.1
        - localhost
        - 10.96.232.136
    - document_name: calico-etcd-n1-peer
      common_name: calico-etcd-n1-peer
      hosts:
        - n1
        - 192.168.77.11

```

(continues on next page)

(continued from previous page)

```

- 127.0.0.1
- localhost
- 10.96.232.136
- document_name: calico-etcd-n2-peer
common_name: calico-etcd-n2-peer
hosts:
- n2
- 192.168.77.12
- 127.0.0.1
- localhost
- 10.96.232.136
- document_name: calico-etcd-n3-peer
common_name: calico-etcd-n3-peer
hosts:
- n3
- 192.168.77.13
- 127.0.0.1
- localhost
- 10.96.232.136
- document_name: calico-node-peer
common_name: calico-node-peer
keypairs:
- name: service-account
description: Service account signing key for use by Kubernetes controller-manager.

```

Certificate Authorities

The data in the `certificate-authorities` key is used to generate certificates for each authority and node.

Each certificate authority requires essential host-specific information for each node, including the hostname and ip as listed in each *Kubernetes Node* document.

The provided *Armada* manifest and will be applied on the genesis node as soon as it is healthy.

1.1.5 Troubleshooting

Genesis Troubleshooting

genesis.sh

Kubernetes services failures

Before the *Armada* manifests are applied, the `genesis.sh` script will bring basic kubernetes services online by starting docker containers for these services.

One of the first services to be brought up is the kubernetes API. If it fails to come up, you may see a repeated error as follows from the `genesis.sh` script:

```
.The connection to the server apiserver.kubernetes.promenade:6443 was
refused - did you specify the right host or port?
```

Check that the hostname in your Genesis.yaml matches the hostname of the machine you are trying to install onto. If they do not match, change one to match the other. If you change Genesis.yaml, then re-generate the Promenade payloads.

If the hostnames match, check the container logs under /var/log/pods to see the reason for the provisioning failure. (kubect1 logs function will not be available if the API container is not running).

Armada failures

When executing genesis.sh, you may encounter failures from Armada in the provisioning of other containers. For example:

```
CRITICAL armada [-] Unhandled error: armada.exceptions.tiller_exceptions.  
↳ReleaseException: Failed to Install release: barbican
```

Use kubect1 logs on the failed pod to determine the reason for the failure. E.g.:

```
sudo kubect1 logs barbican-api-5b8bccdf8f-x7sld --namespace=ucp
```

Other errors may point to configuration errors. For example:

```
CRITICAL armada [-] Unhandled error: armada.exceptions.source_exceptions.  
↳GitLocationException: master is not a valid git repository.
```

In this case, the git branch name was inadvertently substituted for the git URL in one of the chart definitions in bootstrap-armada.yaml.

Post-run failures

At its conclusion, the genesis script will output the list of containers provisioned and their status, as reported by kubernetes. It is possible that some containers may not be in a Running state. E.g.:

```
ucp   promenade-api-6696769cd-qwpzf   0/1   ImagePullBackOff   0   10h
```

For general failures, kubect1 logs may be used as in the previous section. In this case, it was necessary to run kubect1 describe on the pod to get the details of the image pull failure. E.g.:

```
kubect1 describe pod promenade-api-7dc54d47c-qw27m --namespace=ucp
```

In this particular incident report, the problem was a missing certificate on the bare metal node which caused the image download to fail. Installing the certificate, restarting the docker service, and then waiting for the container to retry resolved this particular issue.

1.1.6 Promenade API Documentation

/v1.0/health

Allows other components to validate Promenade's health status.

GET /v1.0/health

Returns the health status.

Responses:

- 204 No Content

/v1.0/join-scripts

Generates join scripts and for Drydock.

GET /v1.0/join-scripts

Generates script to be consumed by Drydock.

Query parameters

hostname Name of the node

ip IP address of the node

design_ref Endpoint containing configuration documents

dynamic.labels Used to set configuration options in the generated script

static.labels Used to set configuration options in the generated script

Responses:

- 200 OK: Script returned as response body
- 400 Bad Request: One or more query parameters is missing or misspelled

/v1.0/validatedesign

Performs validations against specified documents.

POST /v1.0/validatedesign

Performs validation against specified documents.

Message Body

href Location of the document to be validated

Responses:

- 200 OK: Documents were successfully validated
- 400 Bad Request: Documents were not successfully validated

/v1.0/node-labels/<node_name>

Update node labels

PUT /v1.0/node-labels/<node_name>

Updates node labels eg: adding new labels, overriding existing labels and deleting labels from a node.

Message Body:

dict of labels

```
{ "label-a": "value1", "label-b": "value2", "label-c": "value3" }
```

Responses:

- 200 OK: Labels successfully updated

- 400 Bad Request: Bad input format
- 401 Unauthorized: Unauthenticated access
- 403 Forbidden: Unauthorized access
- 404 Not Found: Bad URL or Node not found
- 500 Internal Server Error: Server error encountered
- 502 Bad Gateway: Kubernetes Config Error
- 503 Service Unavailable: Failed to interact with Kubernetes API

1.1.7 Promenade Exceptions