

# Reproducible and Scalable Experiments with SkyhookDM Ceph

Jayjeet Chakraborty  
Carlos Maltzahn, Ivo Jimenez, Jeff LeFevre  
UC Santa Cruz



**CROSS**

CENTER FOR RESEARCH IN  
OPEN SOURCE SOFTWARE



# Generic Systems Experimentation Workflow



Google Cloud Platform



Boot VMs or bare metal  
Nodes

## BUILD CEPH

You can get Ceph software by retrieving Ceph source code and building environment, compile Ceph, and then either install it via

### BUILD PREREQUISITES

**Tip:** Check this section to see if there are specific prerequisites

A debug build of Ceph may take around 40 gigabytes. If you want 1 total disk space on the VM is at least 50 gigabytes.

Please also be aware that some distributions of Linux, like CentOS from LVM may reserve a large portion of disk space of a typical job

Before you can build Ceph source code, you need to install several

```
...lsaktal1-deps.sh
```

**Note:** Some distributions that support Google's memory profiler (perfstack).

## BUILD CEPH

Ceph is built using cmake. To build Ceph, navigate to your cloned 1

```
cd ceph
./configure.sh
make
```

**Note:** By default do, cmake sh will build a debug version of ceph loads. Pass `-DCMAKE_BUILD_TYPE=Release` to do an ceph executables instead.

## INSTALLING CEPH

There are several different ways to install Ceph. Choose the method that best suits your needs.

### RECOMMENDED METHODS

**Cephadm** installs and manages a Ceph cluster using containers and systemd, with tight integration with the CLI and dashboard GUI.

- cephadm only supports Octopus and newer releases.
- cephadm is fully integrated with the new orchestration API and fully supports the new CLI and dashboard features to manage cluster deployment.
- cephadm requires container support (podman or docker) and Python 3.

**Rook** deploys and manages Ceph clusters running in Kubernetes, while also enabling management of storage resources and provisioning via Kubernetes APIs. We recommend Rook as the way to run Ceph in Kubernetes or to connect an existing Ceph storage cluster to Kubernetes.

- Rook only supports Nautilus and newer releases of Ceph.
- Rook is the preferred method for running Ceph on Kubernetes, or for connecting a Kubernetes cluster to an existing (external) Ceph cluster.
- Rook supports the new orchestrator API. New management features in the CLI and dashboard are fully supported.

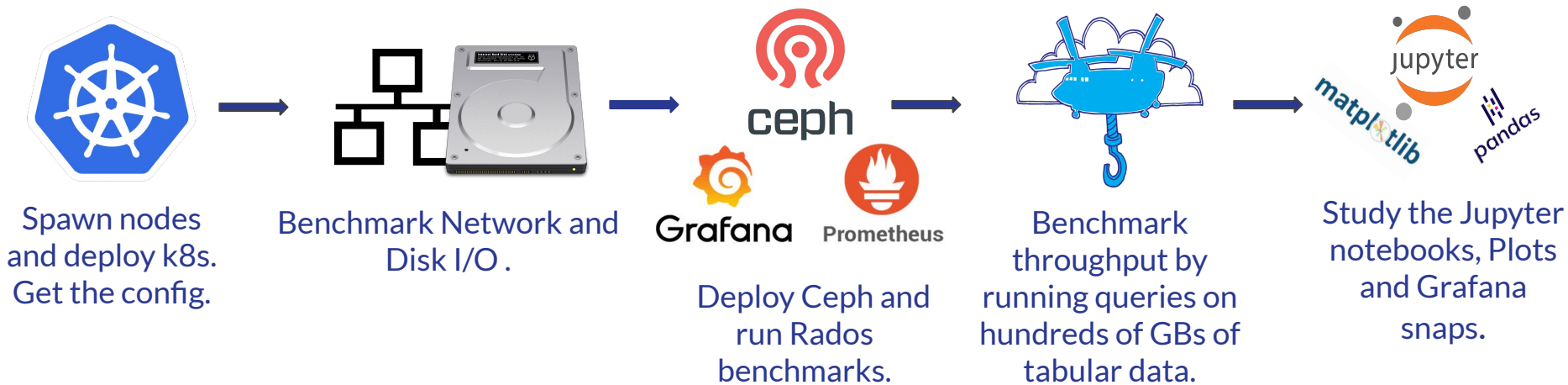


Build, Deploy and Run  
experiments

Prepare plots and  
notebooks

Doing manually is time consuming and error-prone !

# High-Level Workflow for SkyhookDM Experiments



Should be platform independent and automated !

# If done manually,

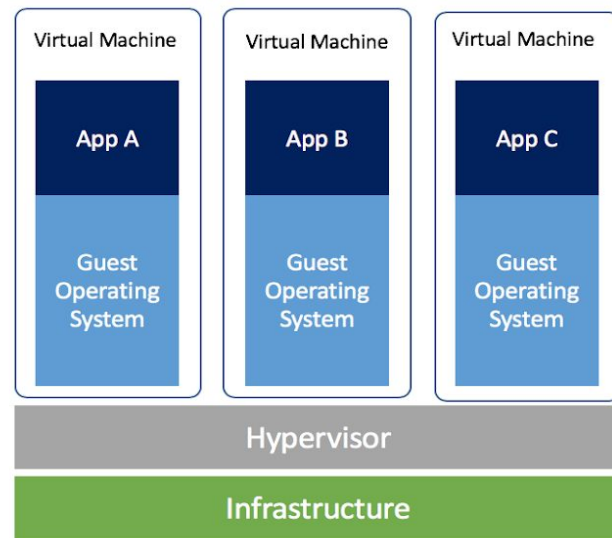
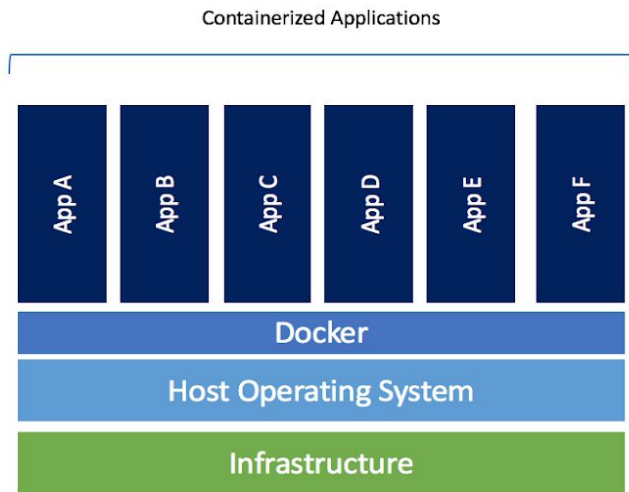
```
● ● ●
# manually build ceph everytime and update all osd, mon, etc.
$ git clone git://github.com/ceph/ceph
$ ./install-deps.sh
$ ./do_cmake.sh
$ cd build
$ make
$ sudo make install

# benchmarking the blockdevices in every node. take care of fio version and the parameters.
$ fio --filename=/dev/sdb --rw=randwrite --direct=1 --ioengine=libaio --bs=64k --numjobs=8 --runtime=120
$ fio --filename=/dev/sdb --rw=randwrite --direct=1 --ioengine=libaio --bs=128k --numjobs=8 --runtime=120

# running rados bench multiple times with different params.
$ ceph osd pool create testpool 32 32 replicated
$ rados bench --no-hints -t 1 -p testpool 120 seq
$ rados bench --no-hints -t 8 -p testpool 120 seq
$ rados bench --no-hints -t 32 -p testpool 120 seq

# plotting results
$ python plot_fio.py
$ pytnon plot_rados.py
```

# Overview of Containers



- Less resource usage than VMs
- Platform independent and portable software
- Consistent operation across environments
- Greater efficiency

# Containerizing Commands

```
$ docker run -e BLOCKDEVICE=sdb  
             -e IODEPTH=32  
             -v $PWD:/workspace  
             --rm  
             --entrypoint /bin/bash  
             -w /workspace  
             bitnami/kubectl:1.17.4  
             ./run_benchmarks.sh
```

Solves platform dependency.

But still lacks automation !

# Introducing Popper

Popper

Containers

Operating System

Hardware



Popper



docker



circleci



podman



```
~/ # cd github-actions-demo/  
~/github-actions-demo # popper run  
Enter the value for PHONY_SECRET:  
F00  
[popper] cloning action repositories  
[popper] - https://github.com/actions/bin@master  
[install] docker pull node:11.6.0  
[install] docker create node:11.6.0 npm install  
[install] docker start  
audited 195 packages in 4.857s  
found 0 vulnerabilities  
  
[test] docker pull node:11.6.0  
[test] docker create node:11.6.0 npm test  
[test] docker start  
  
> github-actions-demo@1.0.0 test /Users/ivo/github-actions-demo  
> mocha ./tests --recursive
```

```
steps:  
- id: install lulesh  
  uses: popperized/spack@master  
  args: [spack, install, -j8, lulesh+mpi]  
  
- id: delete existing jobs  
  uses: popperized/bin/sh@master  
  args: [rm, -fr, sweep/jobs]  
  
- id: install sweepj2  
  uses: popperized/python-actions@master  
  args: [pip, install, sweepj2]  
  
- id: generate sweep  
  uses: jefftriplett/python-actions@master  
  args: [  
    "sweepj2",  
    "--template", "./sweep/script.j2",  
    "--space", "./sweep/space.yml",  
    "--output", "./sweep/jobs/",  
    "--make-executable"  
  ]  
  
- id: run sweep  
  uses: popperized/spack@master  
  args: [run-parts, ./sweep/jobs]
```



# “Popperizing” the SkyhookDM Experimentation Workflow

## Development and testing environrn

Jeff LeFevre edited this page on 11 Aug 2019 - 39 revisions

Skyhook development and testing only requires a Linux environment Linux machine (Ubuntu/F preferred), a VM, or Docker container requires about 30GB of disk space.

Ubuntu 18.04 LTS is recommended, previous Ubuntu versions (1 Skyhook's additional library dependencies. Other Linux versions from major distribution of those supported by Ceph is likely to v

## Linux desktop instructions

- To use your own Linux machine, please go directly to the Bu

## VM Instructions (Virtual Box or V

This is not necessarily recommended due to resources required use Docker as below.

- On your machine's bios, enable Intel/AMD virtual execution
- Install VirtualBox or VMWare
- Create a new Ubuntu 18.04 LTS image, use settings as at le
- Start it. You can be any user you like, the user just needs su
  - mkdir -p /usr/local/repos/skyhook
  - cd /usr/local/repos/skyhook
- Go to the Build wiki page and continue from there, then you queries page.

## Docker instructions

- Install docker on your host machine. A Linux or Mac host mu see our Notes for installing Docker on Windows Home.
- On your host machine, create a dir path for your the skyhook about 30 GB of storage space to build skyhook-ceph. The s machine but visible for compile within the container by usin container will read/write access to this dir on your local ma docker on a local linux machine but if having trouble saving it may need some configuration:
  - mkdir -p /home/jp/repos/skyhook
- Start the container, note absolute paths are required
  - docker run -ti -v /path/to/local/machine/repos/roos /usr
  - docker run -ti -v /home/jp/repos/skyhook/usr/local/rep
  - Note the container should be running and you should be at a
    - cd /usr/local/repos/skyhook
- You can now detach from the running container ctrl-p ctrl-q container is up and running, note the container id
  - Reattach to the running container docker attach container
  - Do not type exit. In the container unless you really want t
- Now you can follow directions on the Build wiki page, then you queries page.
- Some useful docker commands
  - show all running container docker container ls
  - exit from inside a running container and terminate it exit
  - detach from running container ctrl-p ctrl-q
  - show all containers and their current status docker ps -a
  - attach to running container docker attach <container\_id>
  - stop a running container docker stop <container\_id>
  - show all images stored locally docker images -a

\*Thanks to Mark Seibel for help with testing these.

## Run test queries

Jeff LeFevre edited this page on 2 Jun - 76 revisions

Be sure you have started a real or virtual Ceph cluster as per the Build page.

CREATE STORAGE POOL for the test data The below commands assume you have built SkyhookDM and are in the build dir. For a non-virtual cluster, you can remove the bin/ prefix:

```
bin/rados-wgpool tpcdata; # here we will use poolname=tpcdata for all of the below test queries.
# Alternatively, if that fails use this command:
bin/crunch-out pool, create tpcdata 128 128 replicated; # -r- 256 placement groups works well with 1-8 OSDs in
```

GET TEST DATA. Each object contains 10 rows, and is formatted as per type indicated, where type is one of SkyFormatType. There are 2 test data objects for each supported data format.

```
# choose one object format type
OBJ_TYPE=SPF_TYPE;
OBJ_TYPE=SPF_ARROW;
OBJ_TYPE=SPF_FLATTEN_FILEX_ROW;
# get the sample data
OBJ_BASE_NAME=skyhook.$(OBJ_TYPE)-lineitem;
for i in $(seq 1 10); do
  wget https://users.soe.ucsc.edu/~jlefevr/skyhookdm/testdata/$OBJ_BASE_NAME.$i;
done;
```

STORE TEST DATA into Ceph objects. Setting the PATH variable is only needed when using a virtual dev cluster from the current build dir.

```
yes | PATH=$PATH:bin ./src/prog/rados-store-glob.sh tpcdata public lineitem skyhook.$OBJ_TYPE-1.lineit
```

## Build

Jeff LeFevre edited this page on 16 Jan - 61 revisions

Directions to clone and build. Tested on clean install of 64-bit Ubuntu 18.04 and Centos7. NOTE: requires at least 30GB disk space to build.

Be sure you have enough disk space to build.

On Cloudlab machines, the \$HOME dir is not large enough, so format and mount one of the larger disks. NOTE: do not wipe your primary disk!

```
lsblk; # show available devices.
ddiskname; # choose a device with enough space.
sudo mkdir -t ext4 /dev/$disk; # USE WITH CAUTION
sudo mkdir -p /mnt/$disk;
sudo mount /dev/$disk /mnt/$disk;
df -h;
```

## Install dependencies and clone repo

```
PKG_MGR=apt-get; # for Ubuntu
PKG_MGR=dnf; # for Centos7
sudo $PKG_MGR update;
sudo $PKG_MGR install wget cmake git g++ gcc python-otp >;
git clone https://github.com/ucscross/skyhookdm-ceph.git;
cd skyhookdm-ceph;
```

## Checkout latest and verify branch before running submodule update.

```
git pull;
git checkout skyhook-luminous;
git branch;
git submodule update --init --recursive;
sudo ./create1-deps.sh;
./do_make.sh;
```

## BUILD Ceph with Skyhook

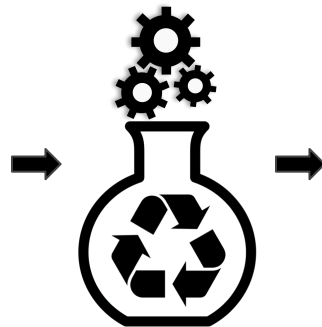
```
cd build;
make -jN cts_tabular_run_query sky_tabular_flatfile_writer ceph_test_skyhook_query_vstart;
# Add -jN to create n jobs i.e., if you have 12 cores use -j12 to compile with 12 cores
# Takes about 13 min with 3.2 GHz 12 core CPU
# All is not required, but make -j12 all, takes about 25 min with 3.2 GHz 12 core CPU
# To save time, just make cts_tabular_run_query for repeat builds, most Skyhook functionality is in there
```

## Start a virtual cluster for dev testing

After compiling vstart above, from the build dir, stop any previously running vstart and then start a new one.

```
./src/stop.sh; mkdir $HOME1 $HOME2 $HOME3 ./src/vstart.sh -d -x
```

IMPORTANT: anytime you recompile Skyhook you should also recompile vstart and stop/start the virtual cluster again.



```
# setup kubernetes
$ popper run -f kubernetes.yml

# baseline cluster
$ popper run -f iperf/fio.yml

# deploy ceph/skyhookdm-ceph
$ popper run -f rook.yml

# setup monitoring
$ popper run -f prometheus.yml

# run rados benchmarks
$ popper run -f radosbench.yml

# run experiment benchmarks
$ popper run -f run_query.yml
```

# Implementation Highlights

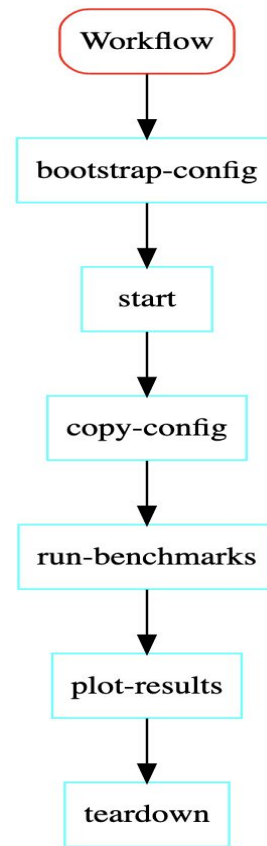
- Highly configurable and scalable workflows with parameter sweeps.
- Monitoring infrastructure with Prometheus + Grafana.
- Everything in Kubernetes. Rook.io, kube-prometheus, kubespray, kubestone, etc.



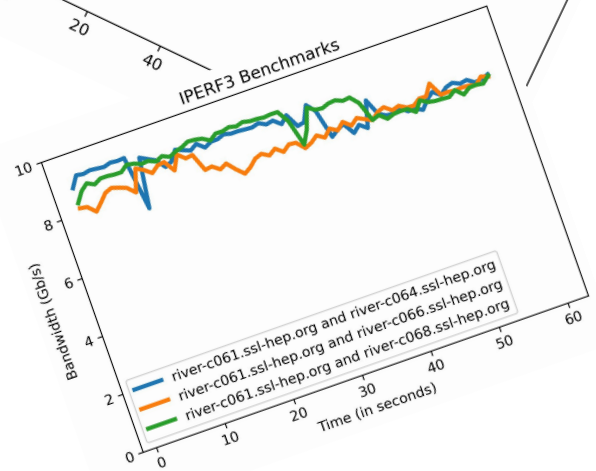
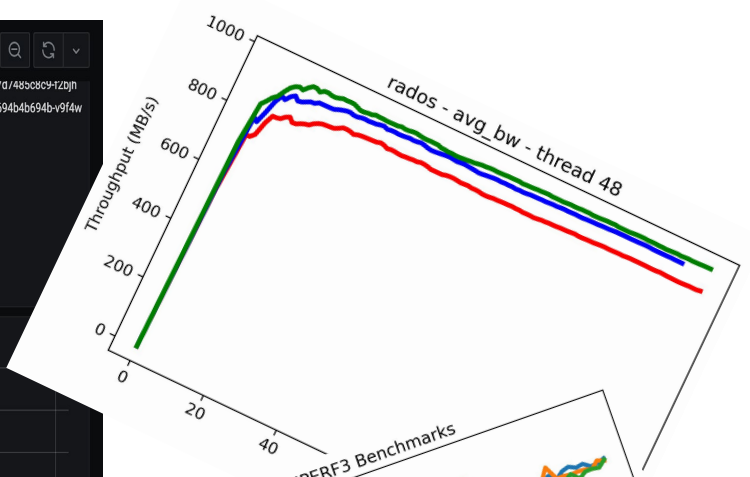
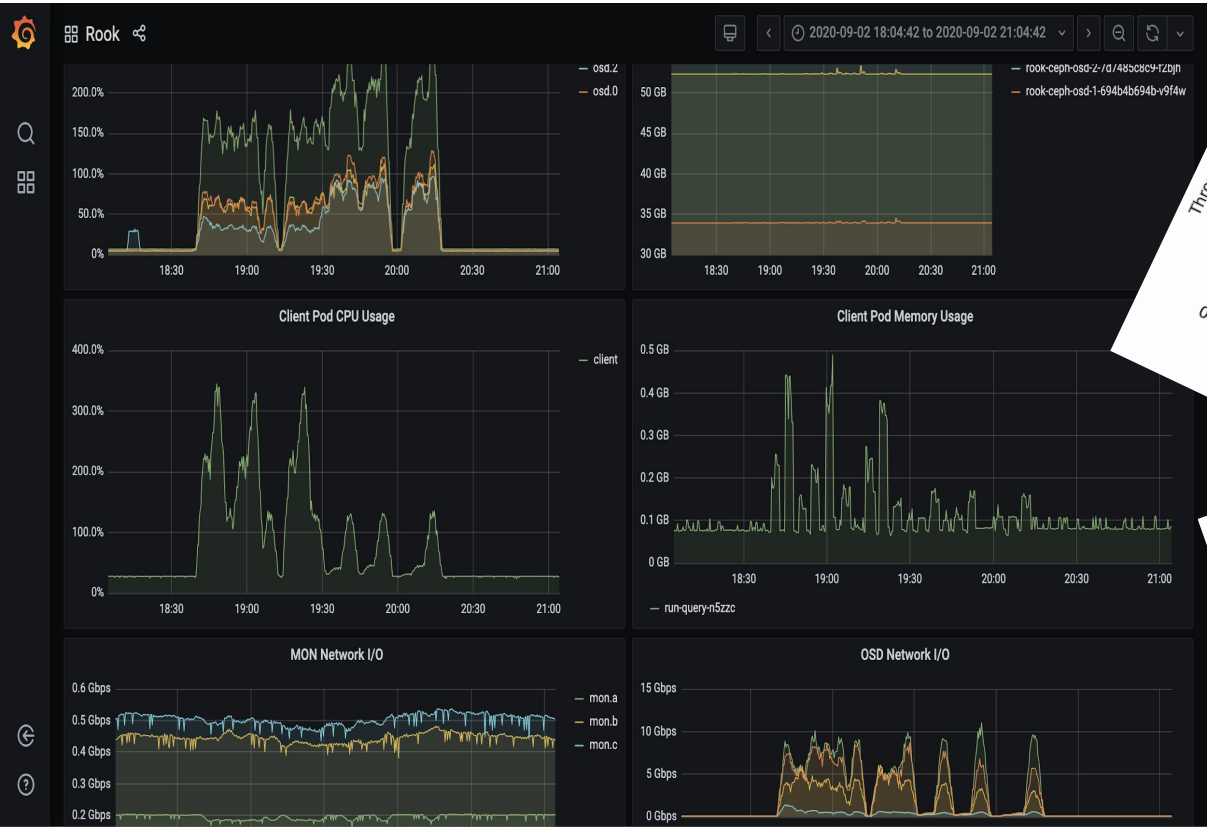
# Workflows for Every High-level Step

```
1 options:
2 env:
3   KUBECONFIG: ./kubeconfig/config
4   BLOCKDEVICES: $_BLOCKDEVICES
5   BLOCKSIZE: '4k 32k 128k 1m 4m'
6   PV_SIZE: 4Gi
7   IO_DEPTH: '32'
8   DURATION: '120'
9   IO_ENGINE: libaio
10  NUM_JOBS: '8'
11  MODES: 'read write randread randwrite'
12  HOSTNAME: $_HOSTNAME
13
14 steps:
15 - id: bootstrap-config
16   uses: docker://biowdl/pyyaml:3.13-py37-slim
17   runs: [python]
18   args: [./kubestone_fio/scripts/bootstrap.py]
19
20 - id: start
21   uses: docker://bitnami/kubectl:1.17.4
22   runs: [bash, -euc]
23   args:
24   - |
25     kubectl apply -n kubestone -f ./kubestone_fio/pv.yaml
26     kubectl apply -n kubestone -f ./kubestone_fio/pvc.yaml
27     kubectl apply -n kubestone -f ./kubestone_fio/job.yaml
28
29 - id: run-benchmarks
30   uses: docker://bitnami/kubectl:1.17.4
31   runs: [./kubestone_fio/scripts/run_benchmarks.sh]
32
33 - id: download-results
34   uses: docker://bitnami/kubectl:1.17.4
35   runs: [./kubestone_fio/scripts/download_results.sh]
36
37 - id: plot-results
38   uses: docker://jupyter/datascience-notebook:python-3.8.5
39   runs: [jupyter]
40   args: ["!nbconvert", "--execute", "--to=notebook", "./kubestone_fio/notebook/plot.ipynb"]
41   options:
42     ports:
43       8888/tcp: 8888
44
45 - id: teardown
46   uses: docker://bitnami/kubectl:1.17.4
47   runs: [bash, -euc]
48   args:
49   - |
50     kubectl delete -n kubestone --ignore-not-found -f ./kubestone_fio/job.yaml
51     kubectl delete -n kubestone --ignore-not-found -f ./kubestone_fio/pvc.yaml
52     kubectl delete -n kubestone --ignore-not-found -f ./kubestone_fio/pv.yaml
```

```
1 options:
2 env:
3   KUBECONFIG: ./kubeconfig/config
4   NAMESPACE: kubestone
5   WRITE_DURATION: '120'
6   READ_DURATION: '120'
7   THREADS: '1 8 32'
8   OBJECT_SIZE: '10M'
9   PG_SIZE: '128'
10  POOL_NAME: 'testbench'
11  POOL_TYPE: 'replicated'
12  REPLICATION_DISABLED: '1'
13  CLIENT: $_CLIENT
14
15 steps:
16 - id: bootstrap-config
17   uses: docker://biowdl/pyyaml:3.13-py37-slim
18   runs: [python]
19   args: [./radosbench/scripts/bootstrap.py]
20
21 - id: start
22   uses: docker://bitnami/kubectl:1.17.4
23   runs: [bash, -euc]
24   args:
25   - |
26     kubectl apply -n "$NAMESPACE" -f ./radosbench/deployment.yaml
27
28 - id: copy-config
29   uses: docker://bitnami/kubectl:1.17.4
30   runs: [./radosbench/scripts/copy_config.sh]
31
32 - id: run-benchmarks
33   uses: docker://bitnami/kubectl:1.17.4
34   runs: [./radosbench/scripts/run_benchmarks.sh]
35
36 - id: plot-results
37   uses: docker://jupyter/datascience-notebook:python-3.8.5
38   runs: [jupyter]
39   args: ["!nbconvert", "--execute", "--to=notebook", "./radosbench/notebook/plot.ipynb"]
40   options:
41     ports:
42       8888/tcp: 8888
43
44 - id: teardown
45   uses: docker://bitnami/kubectl:1.17.4
46   runs: [bash, -euc]
47   args:
48   - |
49     kubectl delete -n "$NAMESPACE" --ignore-not-found -f ./radosbench/deployment.yaml
```



# Monitoring with Prometheus and Grafana



# Jupyter Notebooks for further Exploration

```
In [1]: import os
import json
import matplotlib.pyplot as plt
```

```
In [2]: results_dir = '../results'
files = os.listdir(results_dir)

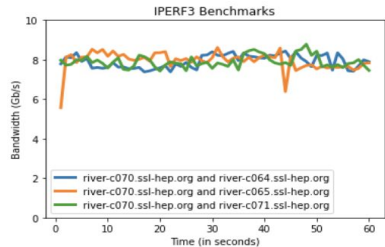
# iterate over each result file and plot the results
for file in files:
    if not file.endswith('.json'): continue
    with open(os.path.join(results_dir, file)) as f:
        results = json.load(f)

    seconds = []
    bandwidth = []

    for run in results["intervals"]:
        seconds.append(float(run["sum"]["end"]))
        bandwidth.append(float(run["sum"]["bits_per_second"])/10)

    plt.plot(seconds, bandwidth, markersize=10, linewidth=3.0, la

# pin the range between 0 and 10
plt.ylim(0.0, 10.0)
plt.xlabel('Time (in seconds)')
plt.ylabel('Bandwidth (Gb/s)')
plt.title('IPERF3 Benchmarks')
plt.legend()
plt.savefig(os.path.join(results_dir, './iperf-benchmarks.png'),
plt.show()
```



```
In [34]: import os
import json
import numpy as np
import matplotlib.pyplot as plt
from collections import defaultdict
from statistics import mean

In [35]: results_dir = "../results"
blockdevices = os.environ["BLOCKDEVICES"].split(" ")
nodes = os.environ["NODES"].split(" ")
iodelths = os.environ["IO_DEPTH"].split(" ")
jobnames = os.environ["JOB_NAMES"].split(" ")
blocksize = os.environ["BLOCKSIZE"].split(" ")
duration = os.environ["DURATION"]

In [36]: def get_mean(int):
    return "%2d" % mean(int)
def get_moving_avg(s, N):
    return np.convolve(s, np.ones(N,)/N)(N-1)
```

## Bandwidth Plots

```
In [37]: for blkdev in blockdevices:
    for iodelth in iodelths:
        for blksize in blksizes:
            for jobname in jobnames:
                base_dir = os.path.join(r"../results/(blkdev)/(blksize)")

                for mode in modes:
                    data = defaultdict(dict)
                    for jobnum in range(1, int(jobname) + 1):
                        time = []
                        value = []

                        with open(os.path.join(base_dir, f'{mode}-iodelth-{iodelth}-mode{job}-jobname_{jobnum}.log'), "r")
                            as f:
                                row = f.readlines()
                                for datapoint in row:
                                    time.append(int(datapoint.split(",")[0])/1000)
                                    value.append(int(datapoint.split(",")[1])/1000)

                                data[jobnum]['time'] = time
                                data[jobnum]['value'] = value

                                min_time_length = 100000000
                                min_val_length = 100000000

                                for jobnum in range(1, int(jobname) + 1):
                                    min_time_length = min(min_time_length, len(data[jobnum]['time']))
                                    min_val_length = min(min_val_length, len(data[jobnum]['value']))

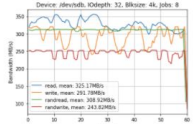
                                cur_time = []
                                cur_val = []

                                for idx in range(1, min_time_length):
                                    sum_time = 0
                                    for jobnum in range(1, int(jobname) + 1):
                                        sum_time += data[jobnum]['time'][idx]

                                    sum_val = 0
                                    for jobnum in range(1, int(jobname) + 1):
                                        sum_val += data[jobnum]['value'][idx]

                                    cur_time.append(sum_time)
                                    cur_val.append(sum_val)

                                moving_avg_val = list(get_moving_avg(cur_val, 3))
                                mean_val = get_mean(moving_avg_val)
                                plt.plot(x=idx, y=mean_val)
                                plt.xlabel('Time (seconds)')
                                plt.ylabel('Bandwidth (Gb/s)')
                                plt.title('IPERF3 Benchmarks')
                                plt.grid()
                                plt.legend()
                                plt.savefig(os.path.join(results_dir, f'./{io-bw-(blkdev)-(iodelth)-(blksize)-(jobname)}.png'), dpi=300, bbox_inches='tight')
                                plt.show()
```



```
1 import json
import os
import datetime

# add datascience libraries
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

[ ] sns.set(style="whitegrid")
results_dir = '../results'
results_arr = [ ]
```

```
[ ] with open(os.path.join(results_dir, 'result.json')) as f:
    # load the data
    data = json.loads(f.read())
    client_side_data = data["lineitem"]["fbx"]

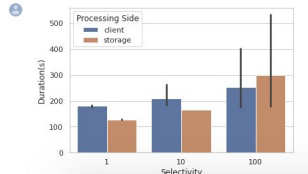
    # prepare the numpy array
    for key, value in client_side_data.items():
        points = value.split(",")
        for point in points:
            results_arr.append((key, point, 'client'))
```

```
1 with open(os.path.join(results_dir, 'result.json')) as f:
    # load the data
    data = json.loads(f.read())
    storage_side_data = data["lineitem"]["fbx_cls"]

    # prepare the numpy array
    for key, value in storage_side_data.items():
        points = value.split(",")
        for point in points:
            results_arr.append((key, point, 'storage'))
```

```
[ ] # convert to dataframes
df = pd.DataFrame(np.array(results_arr), columns=['Selectivity', 'Duration(s)', 'Processing Side'])
df[['Duration(s)']] = df[['Duration(s)']].apply(pd.to_numeric)
```

```
[ ] # plot
ax = sns.barplot(x="Selectivity", y="Duration(s)", hue="Processing Side", data=df)
ax.figure.savefig(os.path.join(results_dir, 'plot.png'), dpi=200)
ax.figure.show()
```



# Case Study: Benchmarking SkyhookDM on River SSL

- Baselined the River SSL Kubernetes cluster and performed performance benchmarks for SkyhookDM Ceph.
- Discovered bottleneck in Network I/O for 10GbE links.
- Discovered unbalanced CPU usage in OSDs due to unbalanced PGs.



**New avenues for further investigation !**

# Future Work

- Try capturing and creating workflows for other categories of Ceph benchmarks like CephFS and Ceph RBD benchmarks.
- Create experiment/benchmark workflows for other popular systems e.g. Key-value stores like RocksDB, databases like PostgreSQL, etc.

# Thank you !

Visit <https://github.com/uccross/skyhookdm-workflows/>

Questions ?

[jchakra1@ucsc.edu](mailto:jchakra1@ucsc.edu)

<https://twitter.com/heyjc25>