

Reproducible and Automated Storage systems experimentation with Popper

Jayjeet Chakraborty

Mentored by Carlos Maltzahn, Ivo Jimenez, Jeff LeFevre

UC Santa Cruz



CENTER FOR RESEARCH IN
OPEN SOURCE SOFTWARE



UNIVERSITY OF CALIFORNIA
SANTA CRUZ

Problems in doing Systems experiments



Google Cloud Platform



Boot VMs or bare metal
Nodes

BUILD CEPH

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

INSTALLING CEPH

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

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 total disk space on the VM is at least 60 gigabytes.

Please also be aware that some distributions of Linux, like CentOS 7, LVM may reserve a large portion of disk space of a typical system. You may need to build Ceph source code, you need to install several

```
./install-deps.sh
```

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

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 CEPH

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

```
cd ceph
./do_configure.sh
cd build
make
```

Note: By default do_configure.sh will build a debug version of ceph. To build a release version, use the --release flag. To build a release version of ceph, use the --release flag. To build a release version of ceph, use the --release flag.

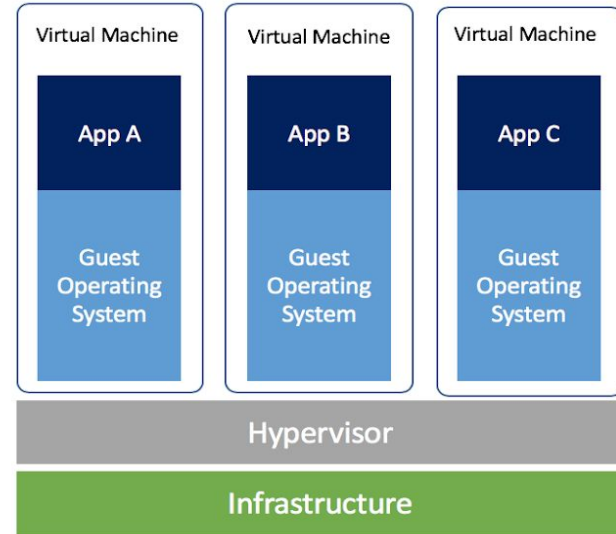
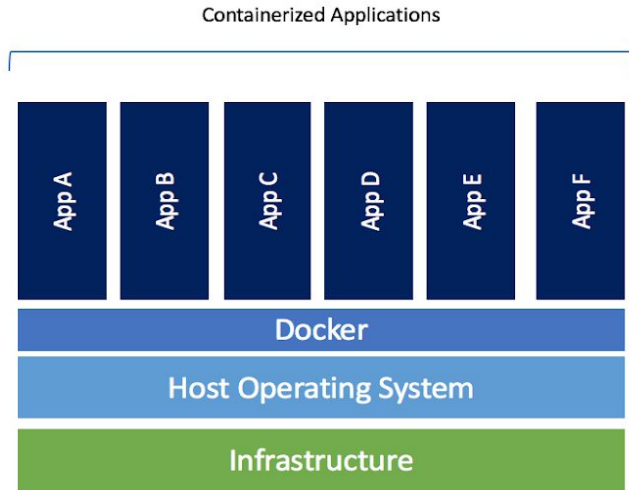
Build, Deploy and Run
experiments



Prepare plots and
notebooks

Should be platform independent and automated ! **Otherwise time consuming and error-prone !**

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 !

What is Popper ?

Popper

Containers

Operating System

Hardware



Popper



docker



circleci



podman



```
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]

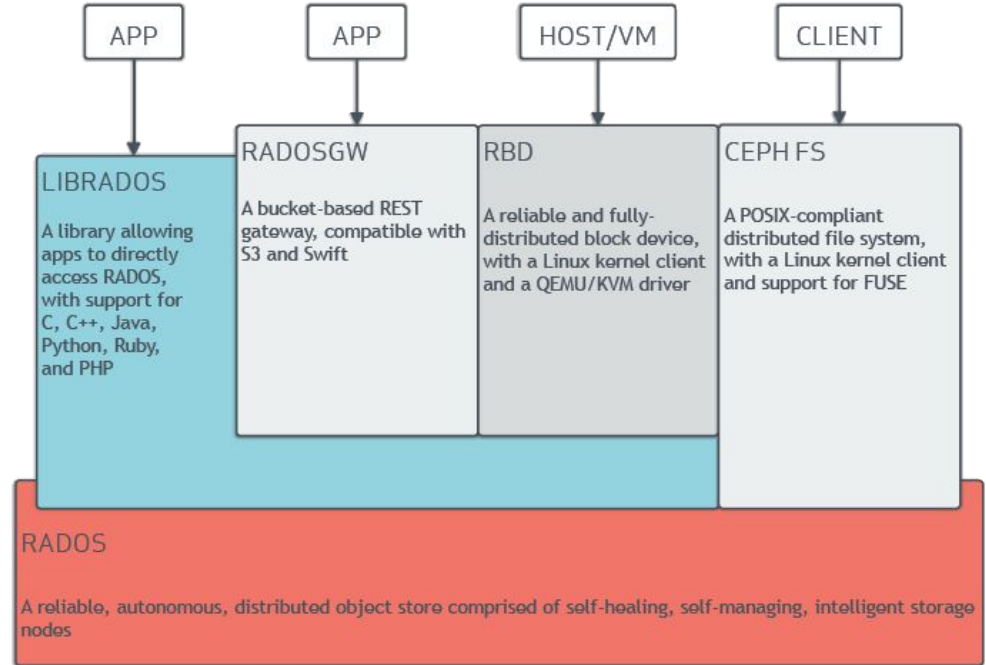
- id: generate sweep
  uses: jefftriplett/python-actions@master
  args: [
    "sweepj2",
    "--template", "./sweep/script.j2",
    "--space", "./sweep/space.yml",
    "--output", "./sweep/jobs/",
    "--make-executable"
  ]

- id: run sweep
  uses: popperized/spack@master
  args: [run-parts, ./sweep/jobs]
```

Reproducible and Scalable Ceph and SkyhookDM experimentation with Popper

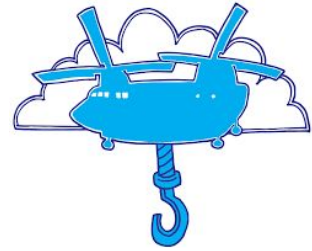
Ceph

1. Provides 3 types of storage interface: File, Object, Block
2. No central point of failure. Uses CRUSH maps that contains object - OSD mapping. A CRUSH map in each client. Client talks directly to OSD.
3. Highly extensible through plugins.

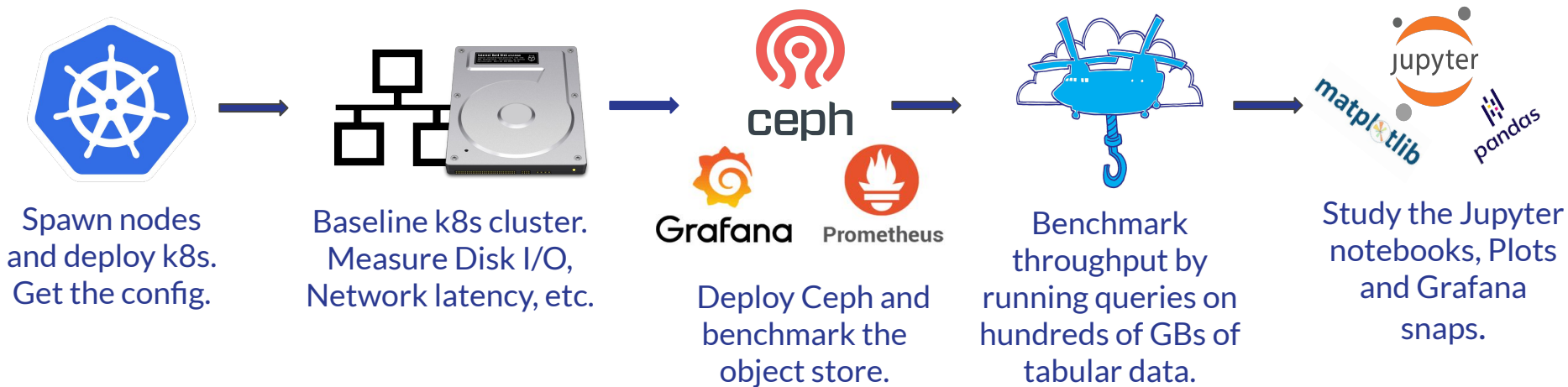


SkyhookDM

- Extends the Ceph object store using custom C++ Object classes for data management in the Storage layer.
- Allows push down of operations like SELECT, PROJECT, AGGREGATE to the storage layer.
- Supports querying both row oriented and column oriented data.



High-Level Workflow for SkyhookDM Experiments



Should be platform independent and automated !

“Popperizing” the SkyhookDM Experimentation Workflow

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-wpoot tpcdata; # here we will use postman=tpcdata for all of the below test queries.
# Alternatively, if that fails use this command:
bin/crunch-out post, create tpcdata 128 128 replicated; # +/- 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 $(seq 10); do
  wget https://users.soe.ucsc.edu/~jlefevr/skyhookdb/testdata/$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.$i
```

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!

```
!build;
# show available devices.
ddiskname; # choose a device with enough space.
subd mkdir -p $mnt/$disk; # USE WITH CAUTION
subd mount $mnt/$disk $mnt/$disk;
subd mount $mnt/$disk $mnt/$disk;
df -h;
```

Install dependencies and clone repo

```
PKG_MGR=apt-get; # for Ubuntu
PKG_MGR=dnf; # for Centos7
subd $PKG_MGR update;
subd $PKG_MGR install wget cmake git gnuip datat 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 checkout --recurse-submodules;
subd ./create_deps.sh;
./do_cmake.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.

Development and testing environrn

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

Skyhook development and testing only requires a Linux environrn 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.

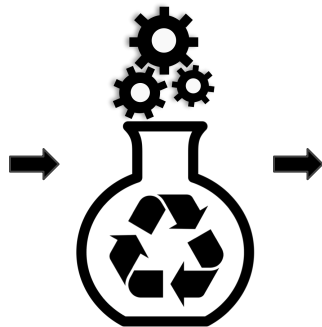
- On your machine's bios, enable Intel/AMD virtual execution
- Install VirtualBox or VMWare
- Create a new Ubuntu 18.04 LTS image, use settings as at le
- Start it. You can be any user you like, the user just needs s
 - mkdir -p /usr/local/repos/skyhook
 - cd /usr/local/repos/skyhook
- Go to the Build 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/cts:/usr/local/repos
 - docker run -ti -v /home/jp/repos/skyhook:/usr/local/repos

Note your chosen path in the container will be creat
- Now 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_id`
 - 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.



```
# 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
```

Building and Deploying Ceph on Bare metal

Building manually vs Building with Popper

```
$ git pull;  
$ git checkout luminous;  
$ git branch -a;  
$ git submodule update --init  
--recursive;   
$ sudo ./install-deps.sh;  
$ ./do_cmake.sh;
```

```
$ popper run dev-init  
$ popper run build
```

```
3 - id: dev-init  
4   uses: docker://alpine/git:v2.24.3  
5   runs: [sh, -ec]  
6   args:  
7   - |  
8     git clone \  
9       --recursive \  
10      --depth 1 \  
11      --shallow-submodules \  
12      --branch nautilus \  
13      https://github.com/ceph/ceph  
14      ln -s ../../../../src ceph/src/cls/tabular  
15      echo "add_subdirectory(tabular)" >> ceph/src/cls/CMakeLists.txt  
16  
17 - id: build  
18   uses: docker://uccross/skyhookdm-builder:nautilus  
19   runs: [bash]  
20   args: [scripts/build.sh]  
21   env:  
22     # travis config: 4 threads, release build  
23     CMAKE_FLAGS: "--DBOOST_J=4 -DCMAKE_BUILD_TYPE=Release -DWITH_MANPAGE=OFF  
24     BUILD_THREADS: "4"
```

Pack all these steps in a Popper workflow and let Popper handle the repetitive and error-prone work

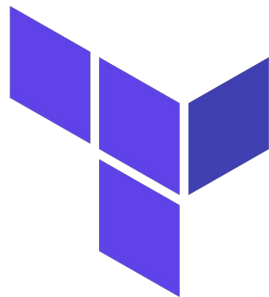
Deploying Ceph manually



ANSIBLE



puppet



ceph-deploy -- Deploy Ceph with minimal infrastructure

`ceph-deploy` is a way to deploy Ceph relying on just SSH access to the servers, `sudo`, and some Python. It runs fully on your workstation, requiring no servers, databases, or anything like that.

- *Learn writing and managing multiple playbooks, vars files, commands in READMEs, host files, etc. Overwhelming for new Ceph developers.*
- *Require knowledge about and using a series of different DevOps tools. Makes the entry barrier high.*

Deploying Ceph using Popper

Popper workflows abstract away tools like Ansible, Terraform, etc. Put all the scripts, playbooks, configuration in a repository and orchestrate with Popper.

```
popper run -f wf.yml
```

A single command to deploy Ceph while using your favourite DevOps tools !

```
.
├── README.md
├── ansible
│   ├── deploy-libcls_tabular.yml
│   ├── deploy-mons-and-osds.yml
│   ├── files
│   ├── group_vars
│   │   └── all.yml
│   └── purge_cluster.yml
├── geni
│   ├── cloudfab_cmd.py
│   └── config.py
└── wf.yml
```

Building and Deploying on Kubernetes via Rook

Even Easier !

Why Kubernetes ? Why Rook ?

Turns storage software into self-managing, self-scaling, and self-healing storage services.

Uses the facilities provided by Kubernetes. Provisioning, scaling, upgrading, disaster recovery.

Make experimentation container native.

Allow using tools from the Kubernetes ecosystem. Get an operator for everything.



Building Ceph and SkyhookDM inside Docker

```
noobjc@kubernetes04:~/skyhookdm-ceph-cls$
```

```
- id: build
  uses: docker://uccross/skyhookdm-builder:nautilus
  runs: [bash]
  args: [scripts/build.sh]
  env:
    # travis config: 4 threads, release build
    CMAKE_FLAGS: "-DBOOST_J=4 -DCMAKE_BUILD_TYPE=Release"
    BUILD_THREADS: "4"

- id: build-rook-img
  uses: docker://docker:19.03.10
  args:
    - build
    - --build-arg=CEPH_RELEASE=v14.2.9
    - --tag=uccross/skyhookdm-ceph:v14.2.9
    - --file=docker/Dockerfile.release
    - .
```

Deploying Ceph on Kubernetes with Rook



```
$ kubectl apply -f common.yml
$ kubectl apply -f operator.yml
$ kubectl apply -f cluster.yml
$ kubectl apply -f toolbox.yml
```

Select the correct config files, install **kubectl**, run the **kubectl apply**'s in proper order.

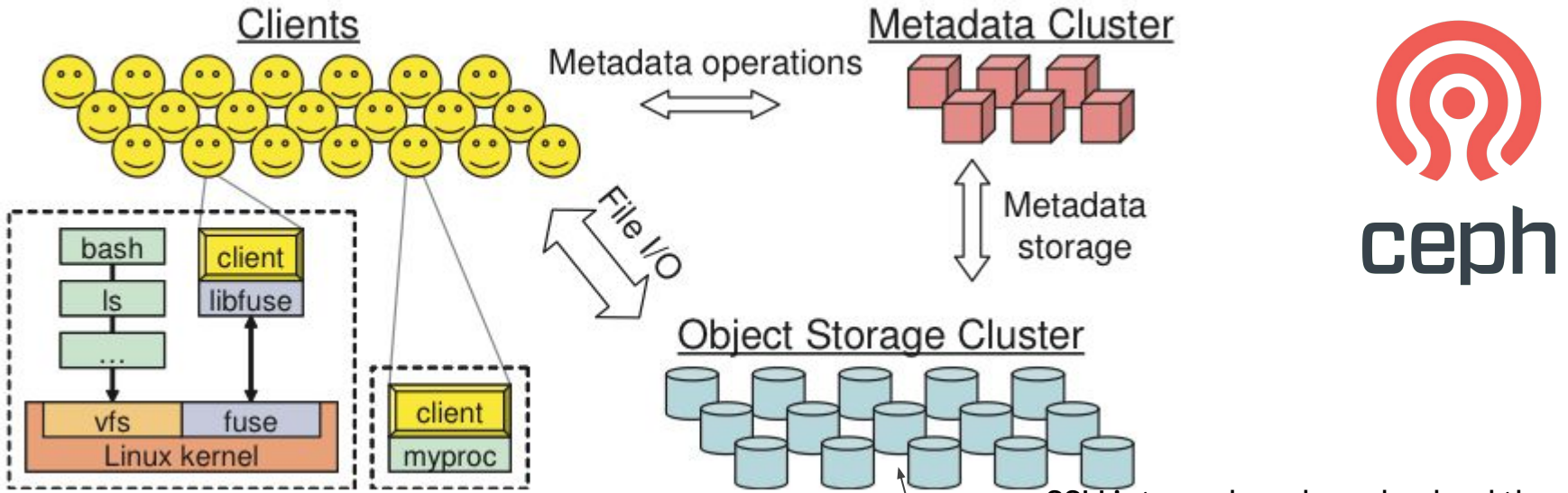
More time spent with Rook that is not really needed !

Making Rook Reproducible and Automated

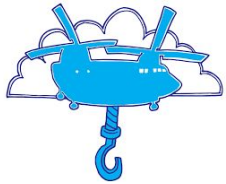
```
- id: setup-ceph
  uses: docker://bitnami/kubectl:1.17.4
  runs: [bash, -euc]
  args:
  - |
    kubectl apply -f ./rook/common.yaml
    kubectl apply -f ./rook/operator.yaml
    kubectl apply -f ./rook/cluster_ceph.yaml
    kubectl apply -f ./rook/toolbox.yaml
```

Clone the repository and `$ popper run -f rook.yml setup-ceph`

Easy upgrades with Rook than Bare metal



SSH into each node and upload the shared libraries at `/usr/lib64/rados-classes/`



`cls_tabular.cc`

`libcls_tabular.so`

Using Object class SDK

Upgrading Ceph to SkyhookDM manually

Use an existing Ceph cluster

Assumes Ceph luminous version. Just build and deploy our extensions library as below, there is no need to stop or restart the cluster.

1. Clone and build SkyhookDM as per our [Build](#) instructions, requires ~10 GB of disk space.
2. Copy the skyhookdm library file into the libcls directory on each of the OSDs, requires sudo
 - Centos7: `LIB_CLS_DIR=/usr/lib64/rados-classes/`
 - Ubuntu18: `LIB_CLS_DIR=/usr/lib/x86_64-linux-gnu/rados-classes/`

```
# from the BUILD dir:
for ((i = 0 ; i < $nosds ; i++)); do
  echo "copying shared lib to osd$i;"
  scp ./lib/libcls_tabular.so.1.0.0 osd$i:/tmp/;
  ssh osd$i "sudo cp /tmp/libcls_tabular.so.1.0.0 ${LIB_CLS_DIR}";
  ssh osd$i "cd ${LIB_CLS_DIR}; if test -f libcls_tabular.so.1; then sudo unlink libcls_tabular.so.1; fi";
  ssh osd$i "cd ${LIB_CLS_DIR}; if test -f libcls_tabular.so; then sudo unlink libcls_tabular.so; fi";
  ssh osd$i "cd ${LIB_CLS_DIR}; sudo ln -s libcls_tabular.so.1.0.0 libcls_tabular.so.1";
  ssh osd$i "cd ${LIB_CLS_DIR}; sudo ln -s libcls_tabular.so.1 libcls_tabular.so";
done;
```

```
osd_class_load_list = *
```

```
osd_class_default_list = *
```

Upgrade the configuration file and SSH into every OSD and copy the libraries at the `LIB_CLS_DIR` of Ceph. Restart OSDs.

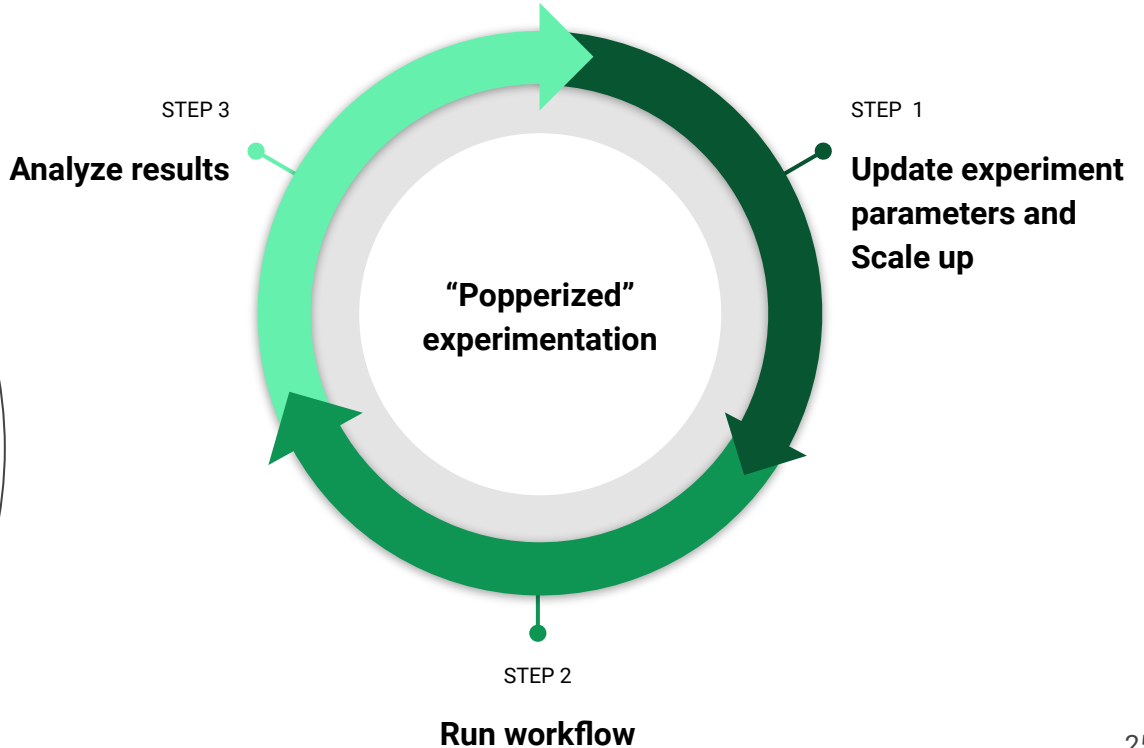
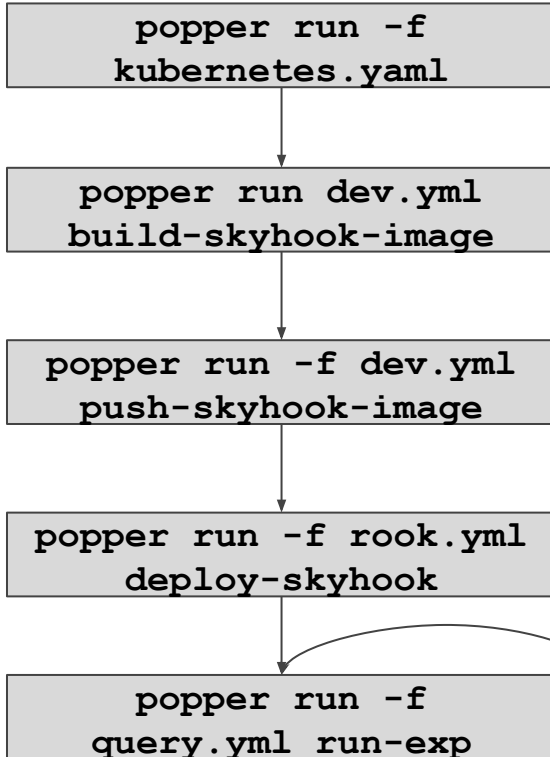
Heavy Manual Work !

Upgrade from Ceph to SkyhookDM on Rook using Popper

```
// deploy vanilla Ceph  
$ popper run deploy-ceph
```

```
// run a popper step to inject libcls_tabular.so and  
upgrade to SkyhookDM  
$ popper run deploy-skyhook-ceph
```

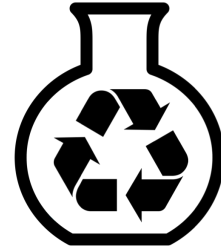

A Scalable experimentation Loop using Popper



Making Notebooks portable with Docker



```
docker run
  -p 8888:8888
  -v ~/notebooks:/home/jayjeet
  jupyter/tensorflow-notebook
  `nbconvert --execute
  --to=notebook
  ./run_query/notebook/plot.ipynb`
```



```
popper run plot-results
```

Let Popper take care of the
error-prone components !

Notebooks , Grafana dashboards and Plots

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

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

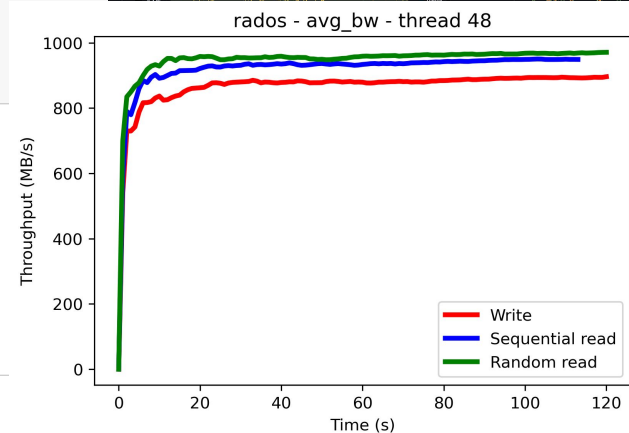
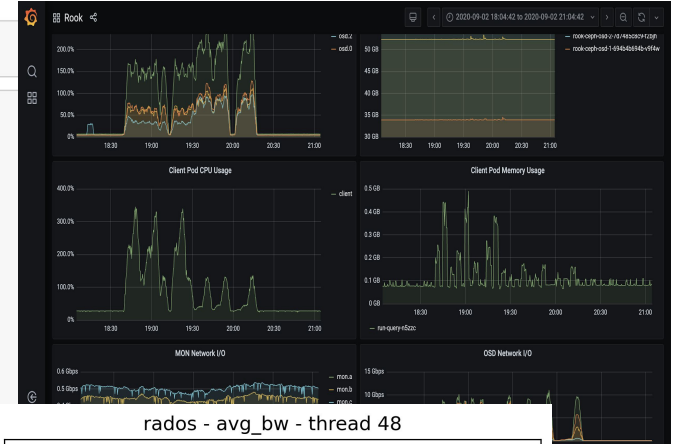
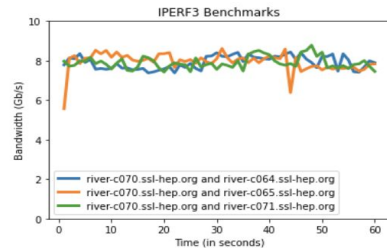
# iterate over each result file and plot the results
for file in files:
    if not file.endswith('.json'): continue
    with open(os.path.join(results_dir, file)) as f:
        results = json.load(f)

        seconds = []
        bandwidth = []

        for run in results["intervals"]:
            seconds.append(float(run["sum"]["end"]))
            bandwidth.append(float(run["sum"]["bits_per_second"])/(1000*1000*1000))

        plt.plot(seconds, bandwidth, markersize=10, linewidth=3.0, label=f"{os.environ['SERVER']} and {file[:-5]}")

# 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'), dpi=300, bbox_inches='tight')
plt.show()
```



Thank you !

Questions ?

