PSI

# Efficient and fast container execution using image snapshotters

Max Fatouros, Derek Feichtinger, **Clemens Lange (PSI)**
Jakob Blomer, Amal Thundiyil, Valentin Völkl (CERN)
CHEP2024, 22nd October 2024

# Motivation

In its standard configuration,

*docker run <image name> <command>*

**downloads the entire container image** from the registry and unpacks it on disk before executing the actual command in the started container
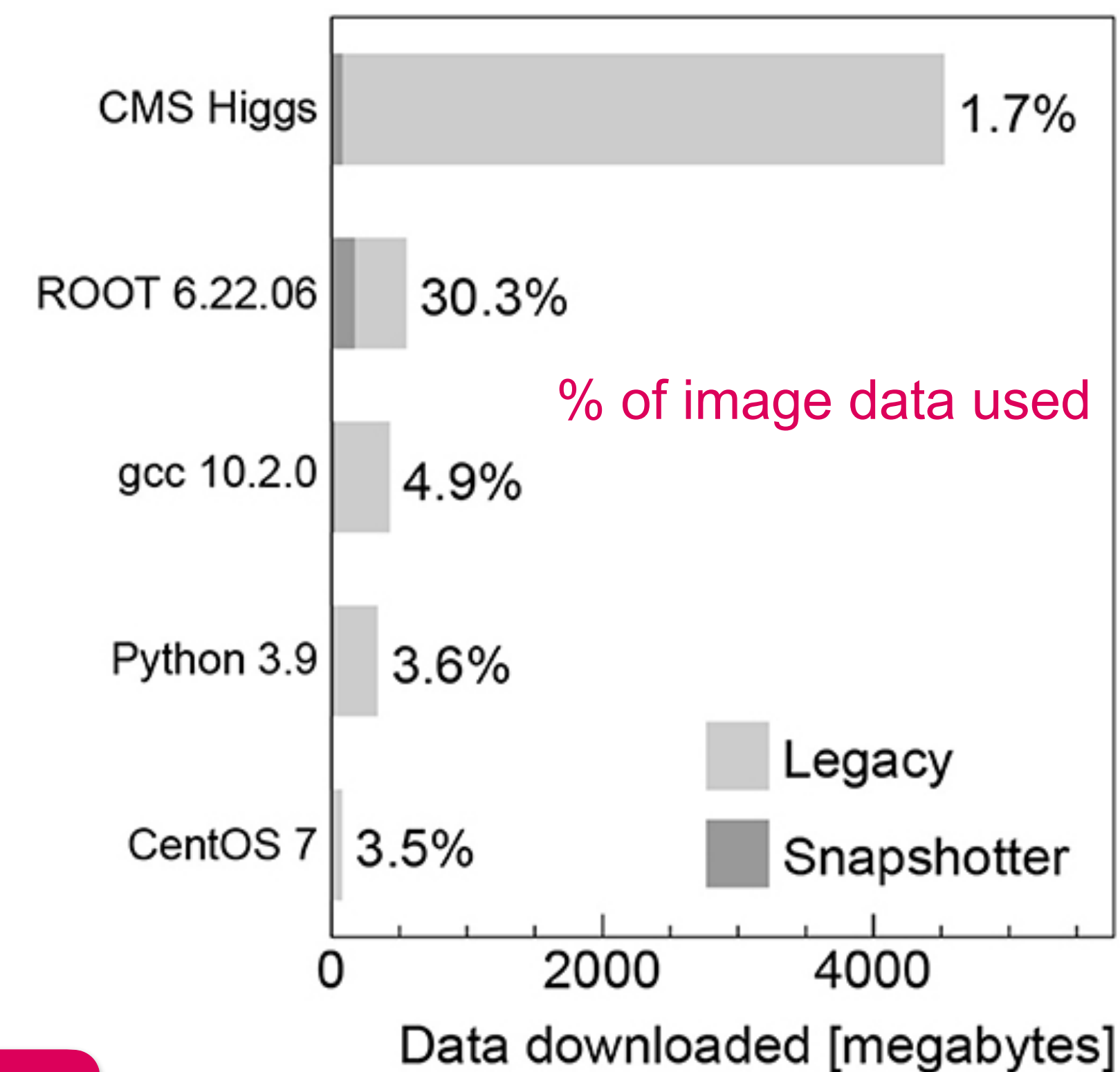
A typical container image used for physics analysis (important e.g. for analysis reusability) has a **size of ~gigabytes**

Executing containerised workloads on batch systems can therefore lead to **hundreds of parallel downloads** of several gigabytes of data

**However, only a fraction of the container image is actually needed**

**→ download only what is actually needed**

[Frontiers in Big Data Vol. 4 (2021) 673163]



% of image data used

Clemens Lange — Efficient and fast container execution using image snapshotters    22.10.2024
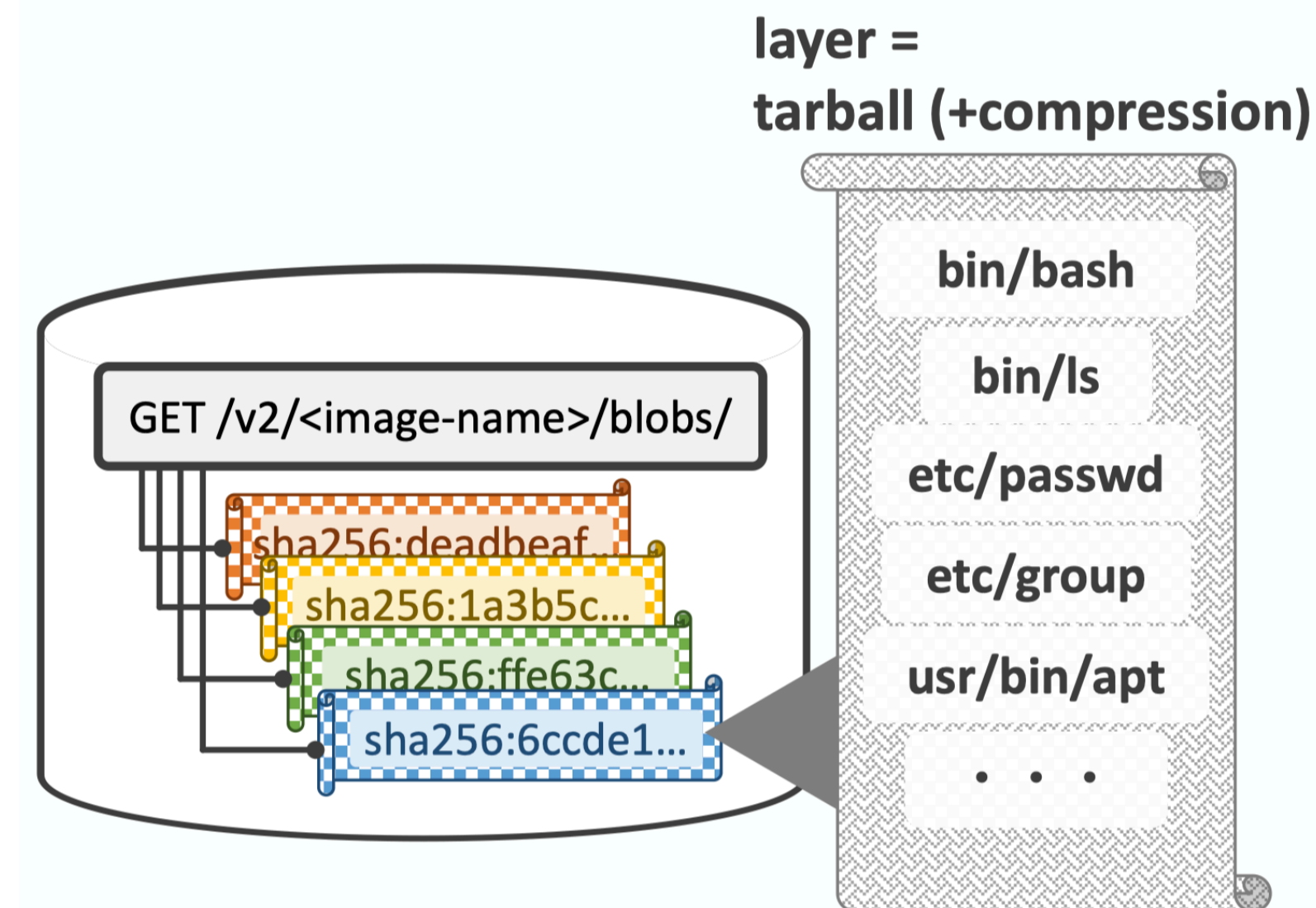
# Lazy-pulling of container images

Lazy-pulling = pull/download only **what is needed when it is needed**

Container reminder:
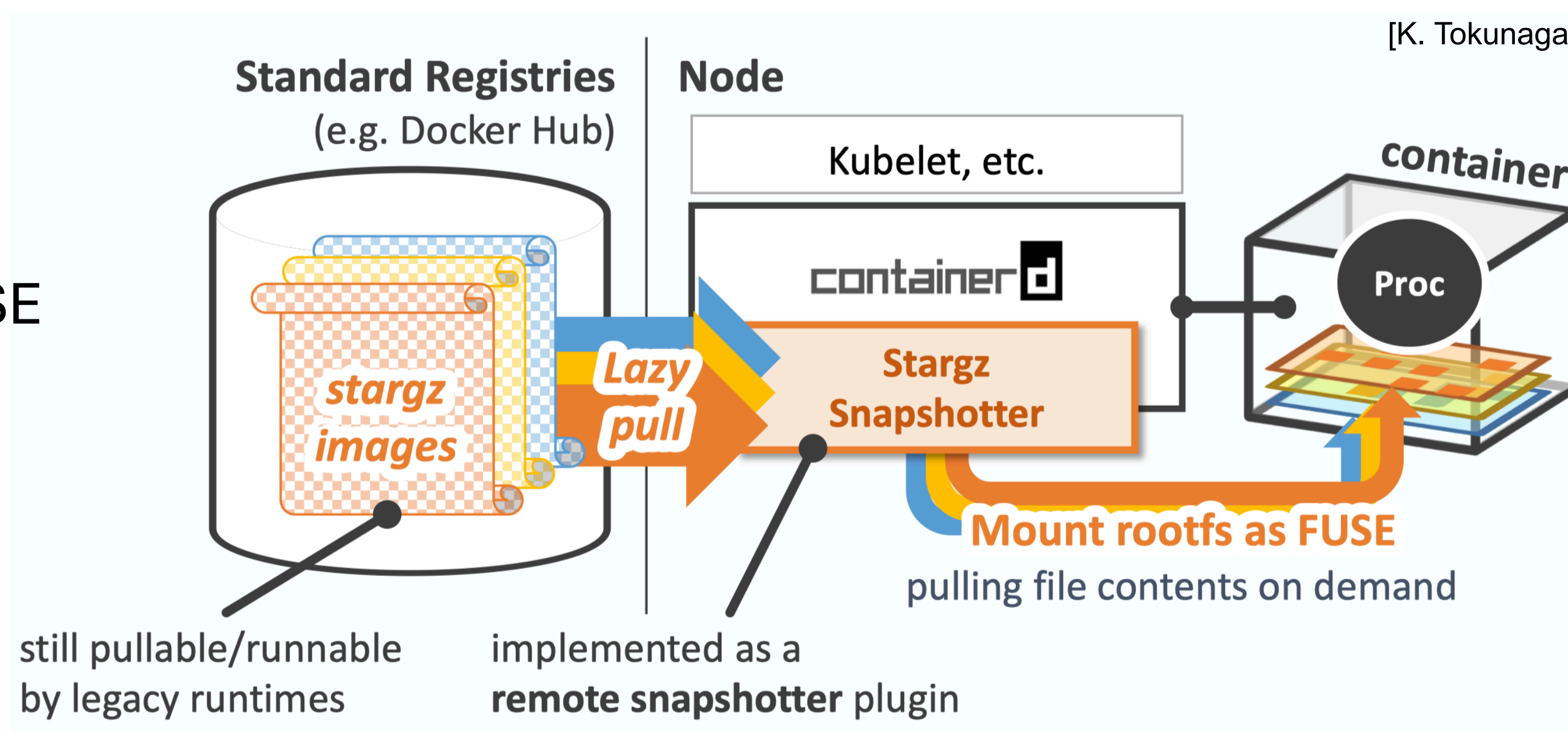
> A container is a set of tar-balls plus a manifest (list)

> Downloading and extracting the layers builds the container file system

Lazy pulling mounts (rootfs snapshots as FUSE and downloads) accessed file contents on-demand

> Can **start container** almost **immediately**

> Can be slower during execution



[K. Tokunaga]

Clemens Lange — Efficient and fast container execution using image snapshotters    22.10.2024

# Implementations of lazy pulling

**PSI**

Solutions are implemented as so-called image snapshotters for use with <u>containerd</u>

Evaluated tools:

> <u>Stargz snapshotter</u>: use images in **searchable tar.gz format**

> <u>SOCI snapshotter</u>: add **separate index artifact** to image (hosted in registry)

> <u>CVMFS snapshotter</u>: use **unpacked images on CVMFS**

> <u>Overlayfs snapshotter</u>: the default/legacy, **non-lazy-loading** snapshotter

All snapshotters will **fall back to legacy pulling** if image not available in required format

CVMFS Snapshotter can additionally do this at **layer level**:

> Enables use of **"protected" layers** based on public base images

> Mind: this is something Singularity/Apptainer cannot do

Side note: "lazy pulling" with Apptainer achieved through unpacked images on CVMFS

# Benchmarking approach

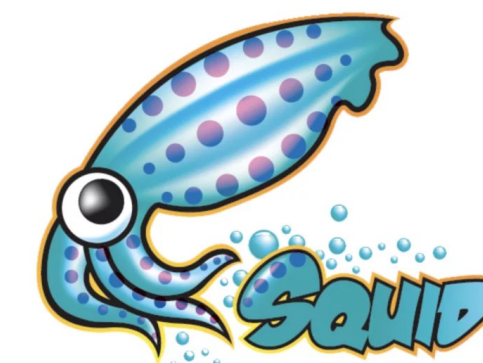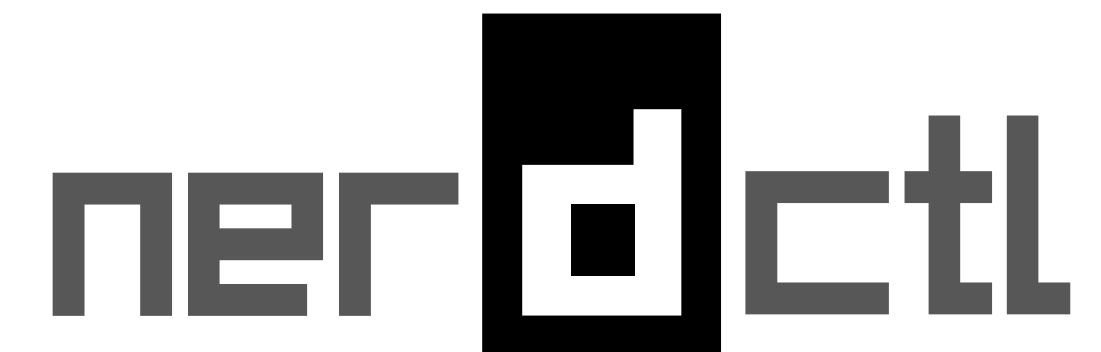Use typical particle physics tasks and container images, e.g.:

> ROOT

> Python



Using <u>nerdctl</u> to run workloads with the various snapshotters:

> Parse execution log files to extract timestamps

> Monitor traffic using network monitoring tools

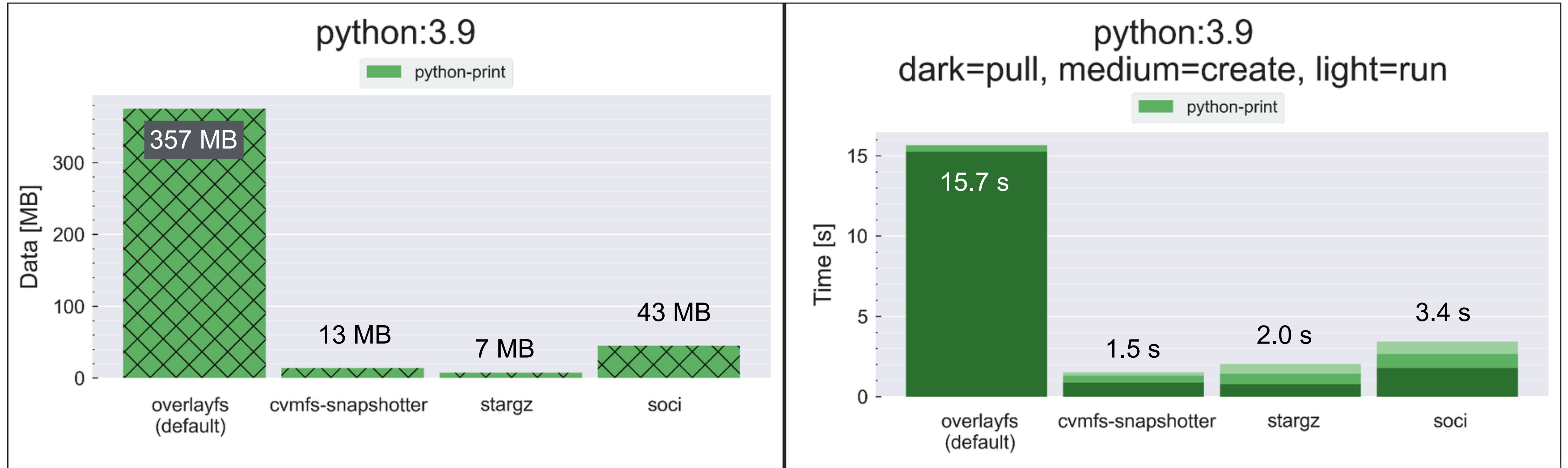> Repeat process several times, clear cache in between runs

Also compare to "legacy" approach pulling entire image before execution



Using local PC in connection with both a local (same machine) registry, and <u>Harbor registry</u> (in the same network).

Using <u>squid proxy</u> for CVMFS caching (with images pre-cached)

Clemens Lange — Efficient and fast container execution using image snapshotters          22.10.2024
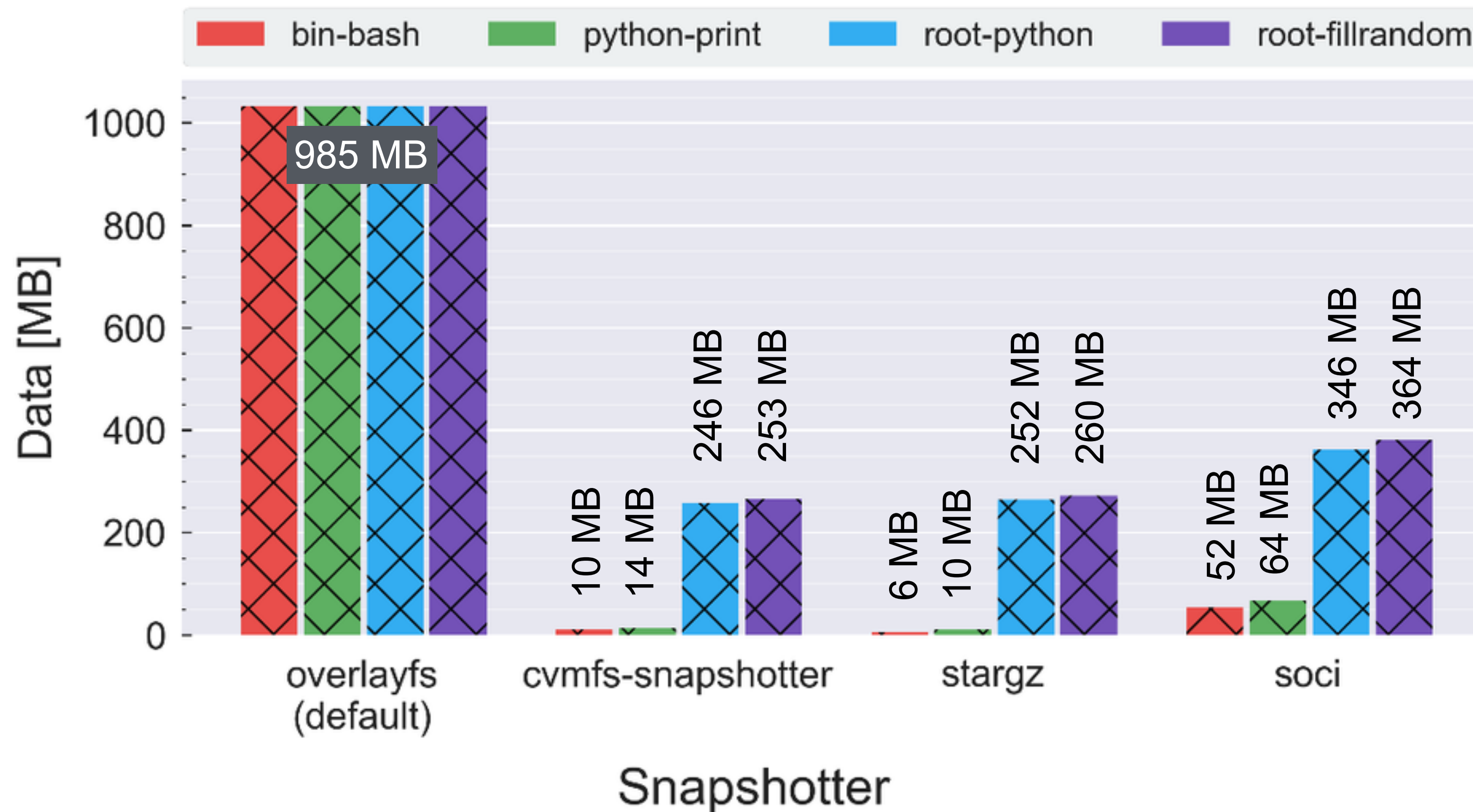
**PSI**



Observations:

> Time to start image drastically reduced for all lazy snapshotters

> Only a few megabytes downloaded

> SOCI loads more data because of layer minimum 10MB size requirement (configurable)

> Using a local registry is slightly faster (→ backup)

# Results: ROOT image (1)

Investigate performance with workloads of increasing complexity:

🟥 *bin/bash*

🟩 *print()* in python

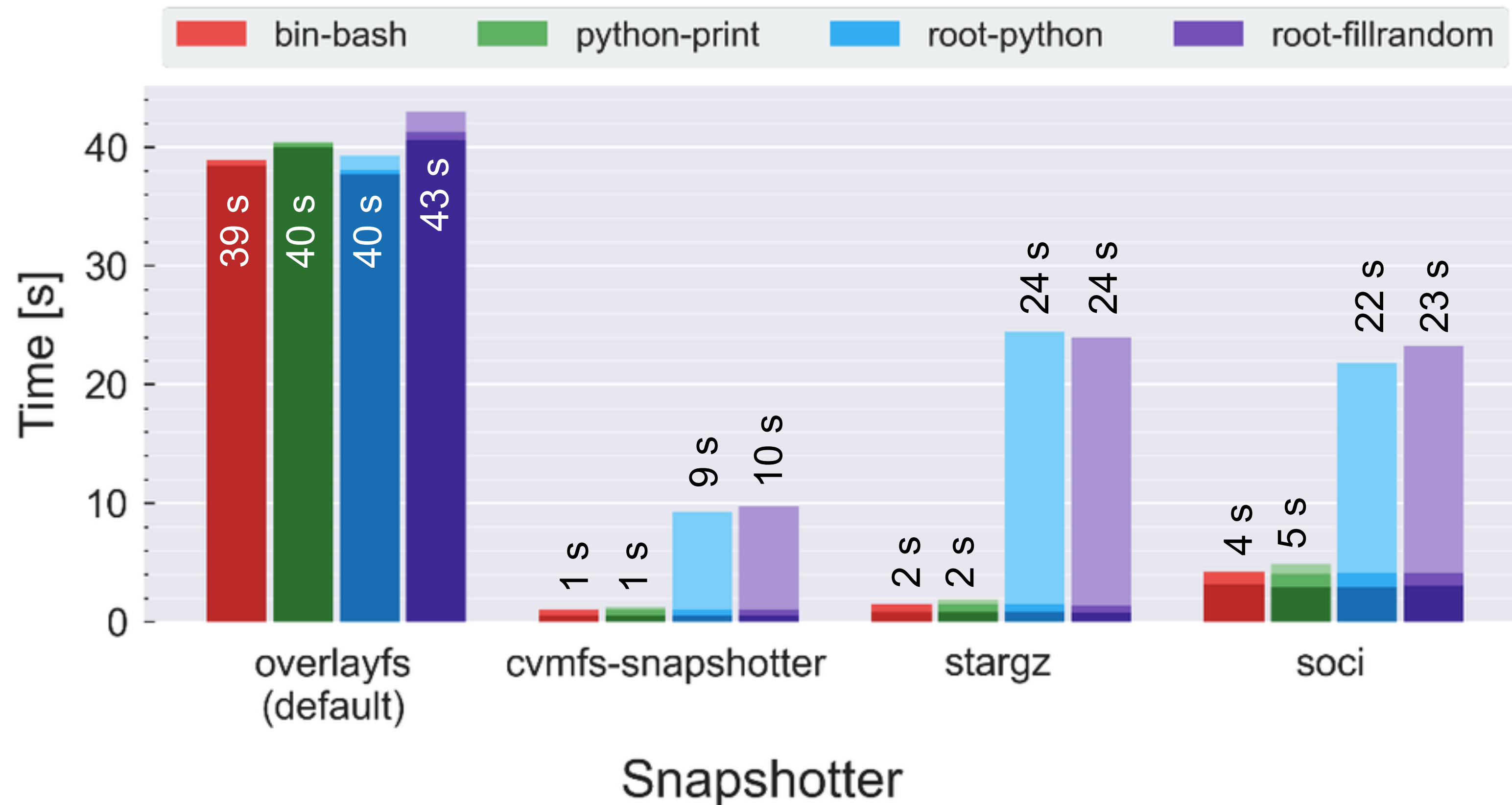🟦 *import ROOT* in python

🟪 *fillrandom.py* using pyROOT

Observations:

**>** Comparable performance

**>** *import ROOT* loads a lot of data

# Results: ROOT image (2)

root:6.32.04-ubuntu24.04
dark=pull, medium=create, light=run

bin-bash · python-print · root-python · root-fillrandom

Investigate performance with workloads of increasing complexity:

■ */bin/bash*

■ *print()* in python

■ *import ROOT* in python

■ *fillrandom.py* using pyROOT

Observations:

> CVMFS snapshotter faster than other two lazy snapshotters

> For complex workloads, pull time small compared to execution time (but mind significant data savings!)

Clemens Lange — Efficient and fast container execution using image snapshotters      22.10.2024

**PSI**

rootless docker available since RHEL 8 and kernel 4.18/5.11 — still requires (cluster) admin action

Use of Stargz snapshotter **requires images to be converted** (programatically) to a specific format
→ adoption might be slow/difficult

SOCI snapshotter only requires small addition to existing image—however, only certain registries support **additional artifacts** (Harbor ✅, GitLab ❌)

CVMFS snapshotter **requires images to be "unpacked"** → **delay** between building them and having them available (and they need to be added to the unpacker "sync" list)

# Example configurations for docker and kubernetes

**PSI**

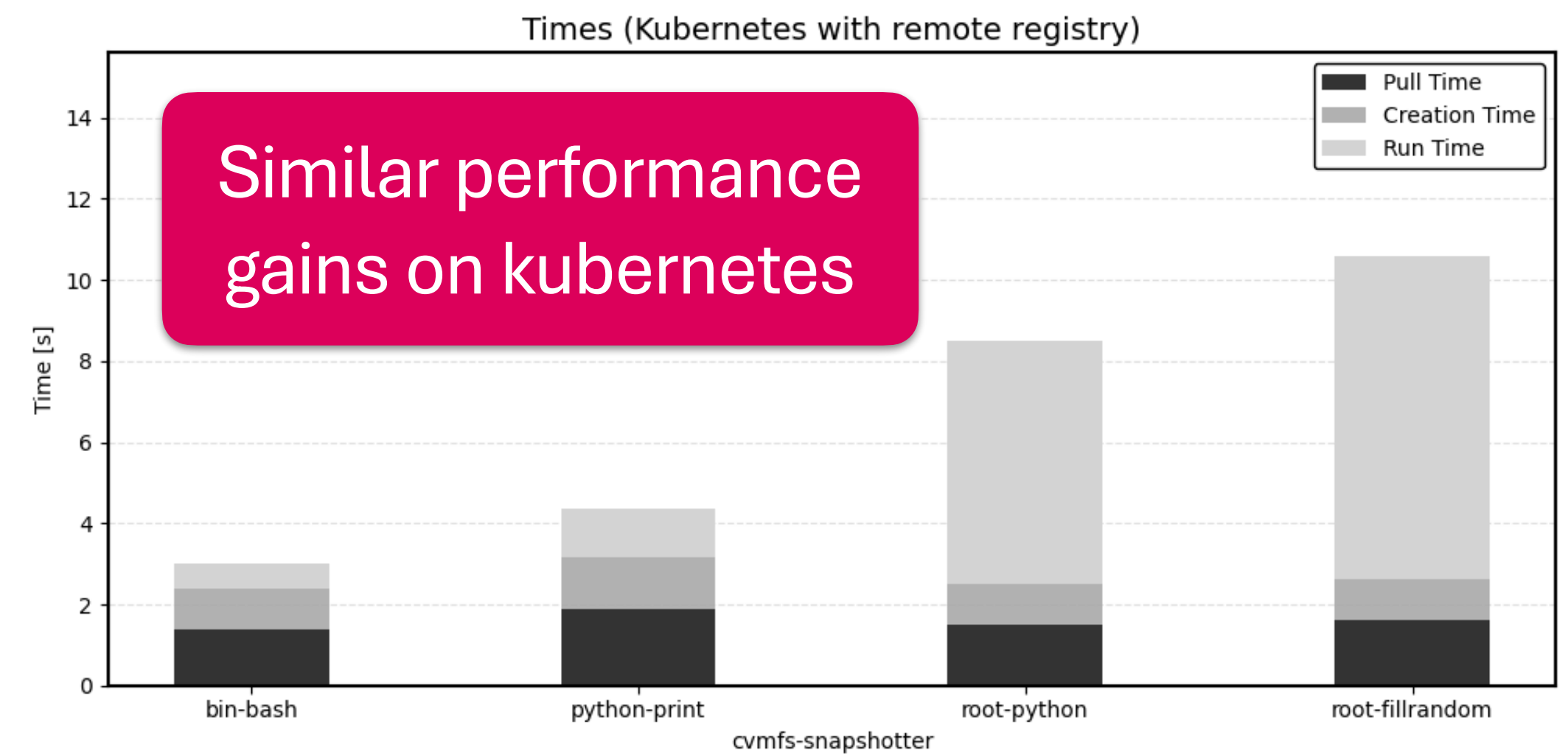containerd configs for CVMFS snapshotter:

```
1   # /etc/containerd/config.toml
2   version = 2
3
4   # Ask containerd to use this particular snapshotter
5   [plugins."io.containerd.grpc.v1.cri".containerd]
6       snapshotter = "cvmfs-snapshotter"
7       # important: the cvmfs snapshotter needs
8       # annotations to work.
9       disable_snapshot_annotations = false
10
11  # Set the communication endpoint between containerd
12  # and the snapshotter
13  [proxy_plugins]
14      [proxy_plugins.cvmfs-snapshotter]
15          type = "snapshot"
16          address =
17              "/run/containerd-cvmfs-grpc/containerd-cvmfs-grpc.sock"
```

kubernetes (k3s) CVMFS snapshotter config:

```
25  # /var/lib/rancher/k3s/agent/etc/containerd/config.toml.tmpl
26  version = 2
27  [plugins."io.containerd.grpc.v1.cri".containerd]
28      snapshotter = "cvmfs-snapshotter"
29      disable_snapshot_annotations = false
30  [proxy_plugins]
31      [proxy_plugins.cvmfs-snapshotter]
32      type = "snapshot"
33      address = "/run/containerd-cvmfs-grpc/containerd-cvmfs-grpc.sock"
34  [plugins."io.containerd.grpc.v1.cri".cni]
35      bin_dir = "/var/lib/rancher/k3s/data/current/bin"
36      conf_dir = "/var/lib/rancher/k3s/agent/etc/cni/net.d"
```

**Setup only requires a few steps/config changes**

```
19      # /etc/containerd-cvmfs-grpc/config.toml
20      # Source of image layers
21      repository = "unpacked.cern.ch"
22      absolute-mountpoint = "/cvmfs/unpacked.cern.ch"
```

**Similar performance gains on kubernetes**



Times (Kubernetes with remote registry)

Legend: Pull Time, Creation Time, Run Time

# Conclusions

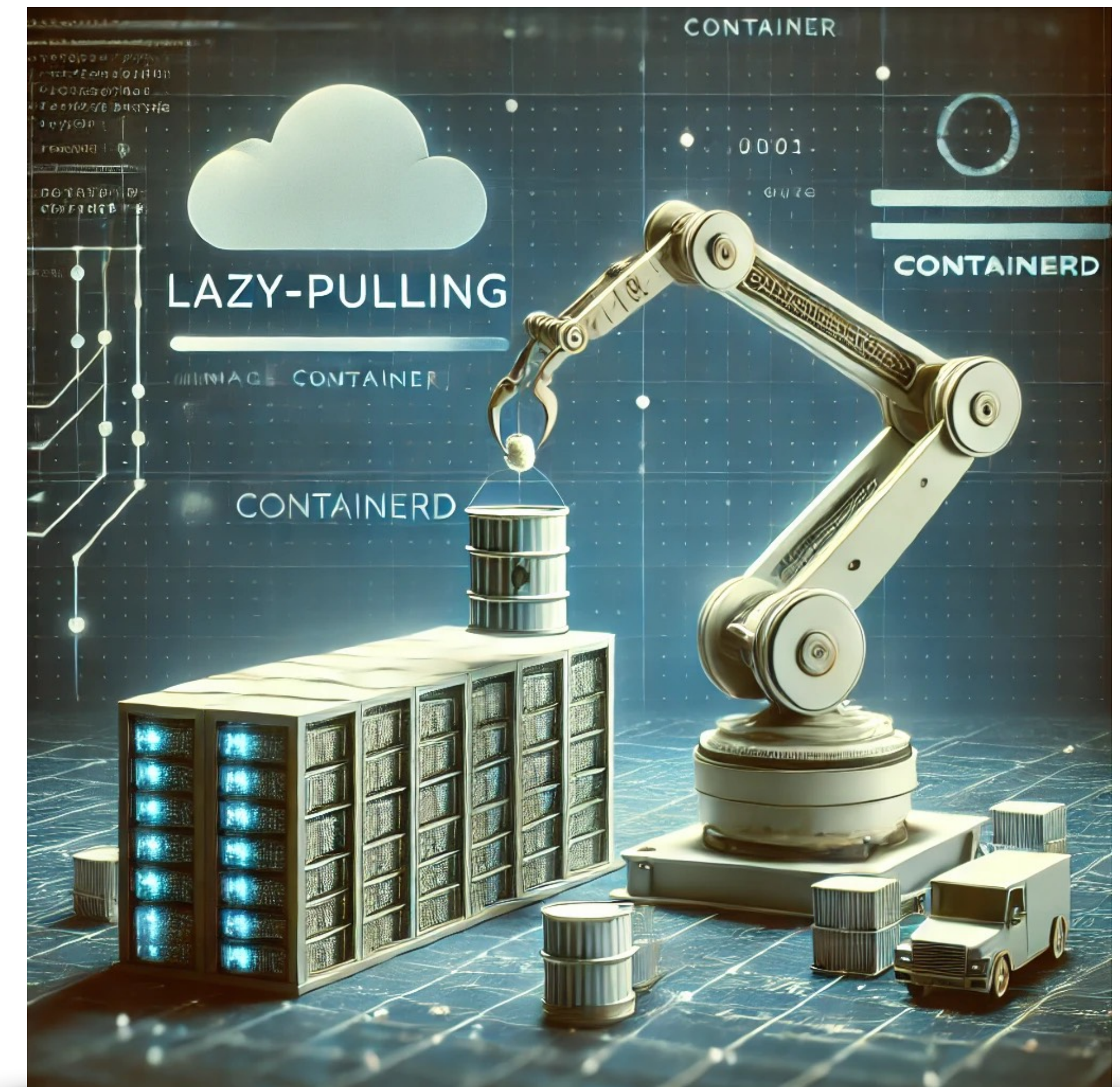Container image snapshotters open up new possibilities for image distribution and access

> Large bandwidth/data savings observed (here > 65%)

> Time saving depends on workload/image details

Overall, evaluated snapshotters all have advantages and disadvantages in usability/requirements
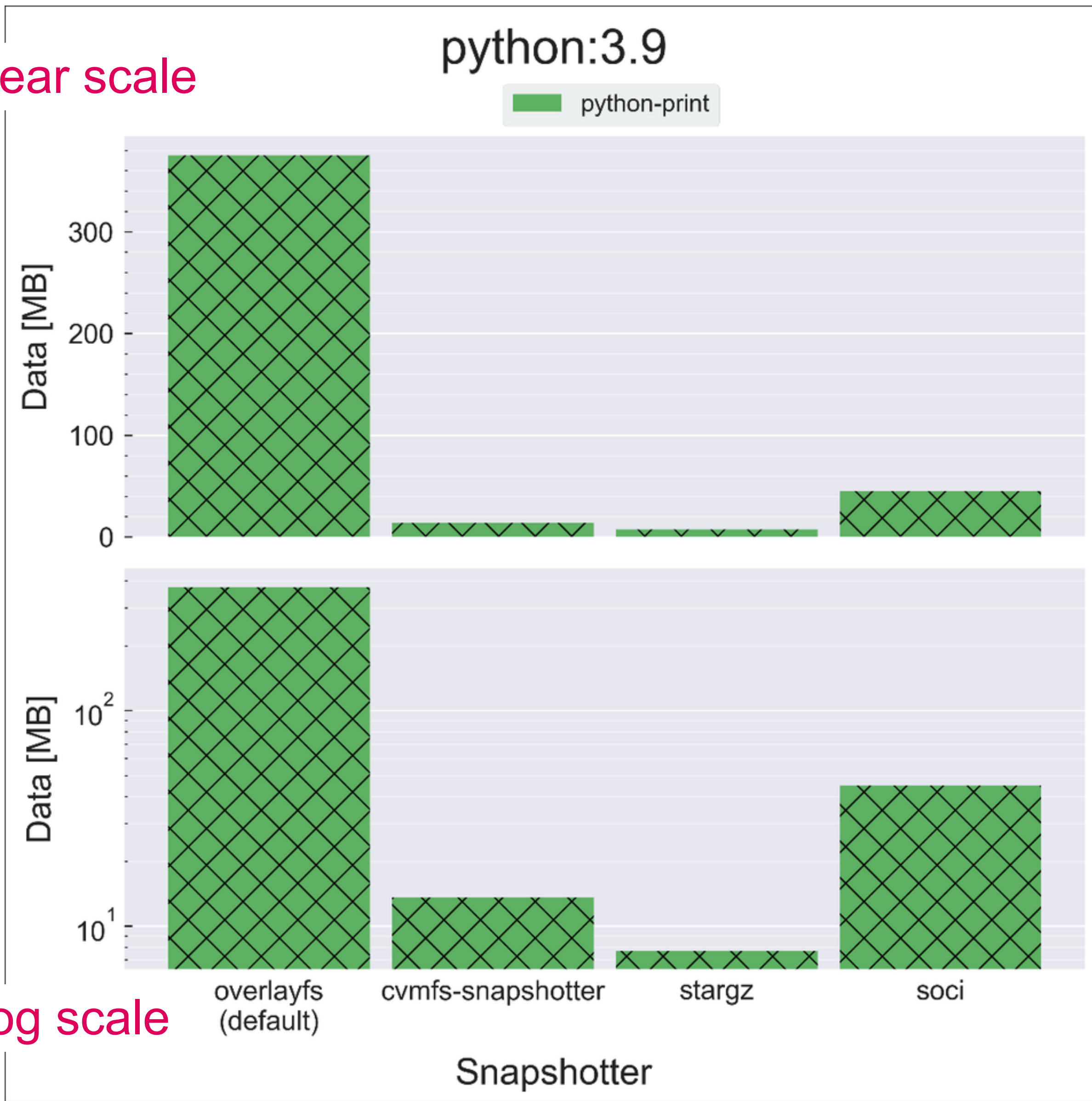
Performance similar

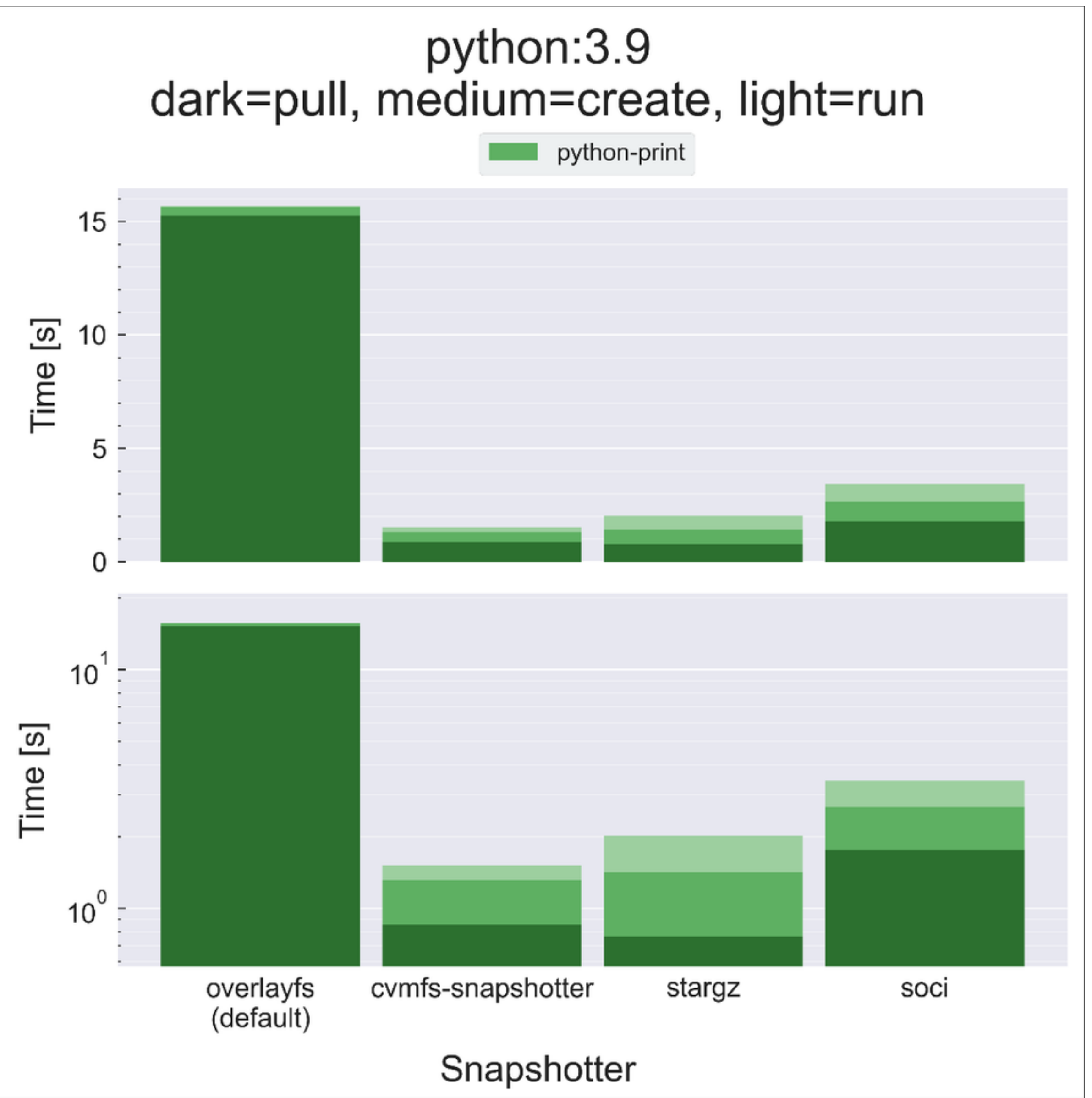> CVMFS snapshotter seems to be a bit faster than the other two snapshotters evaluated
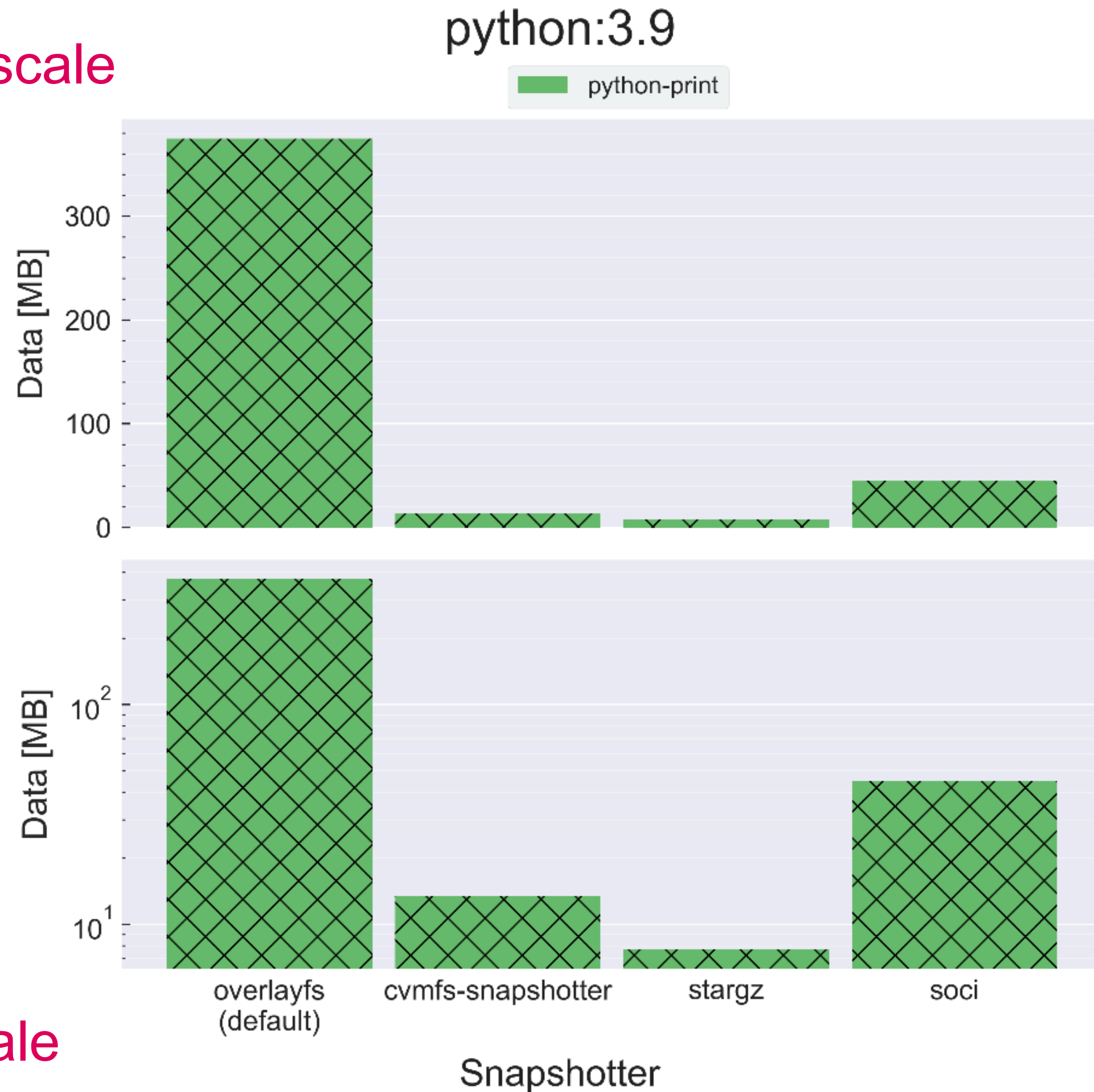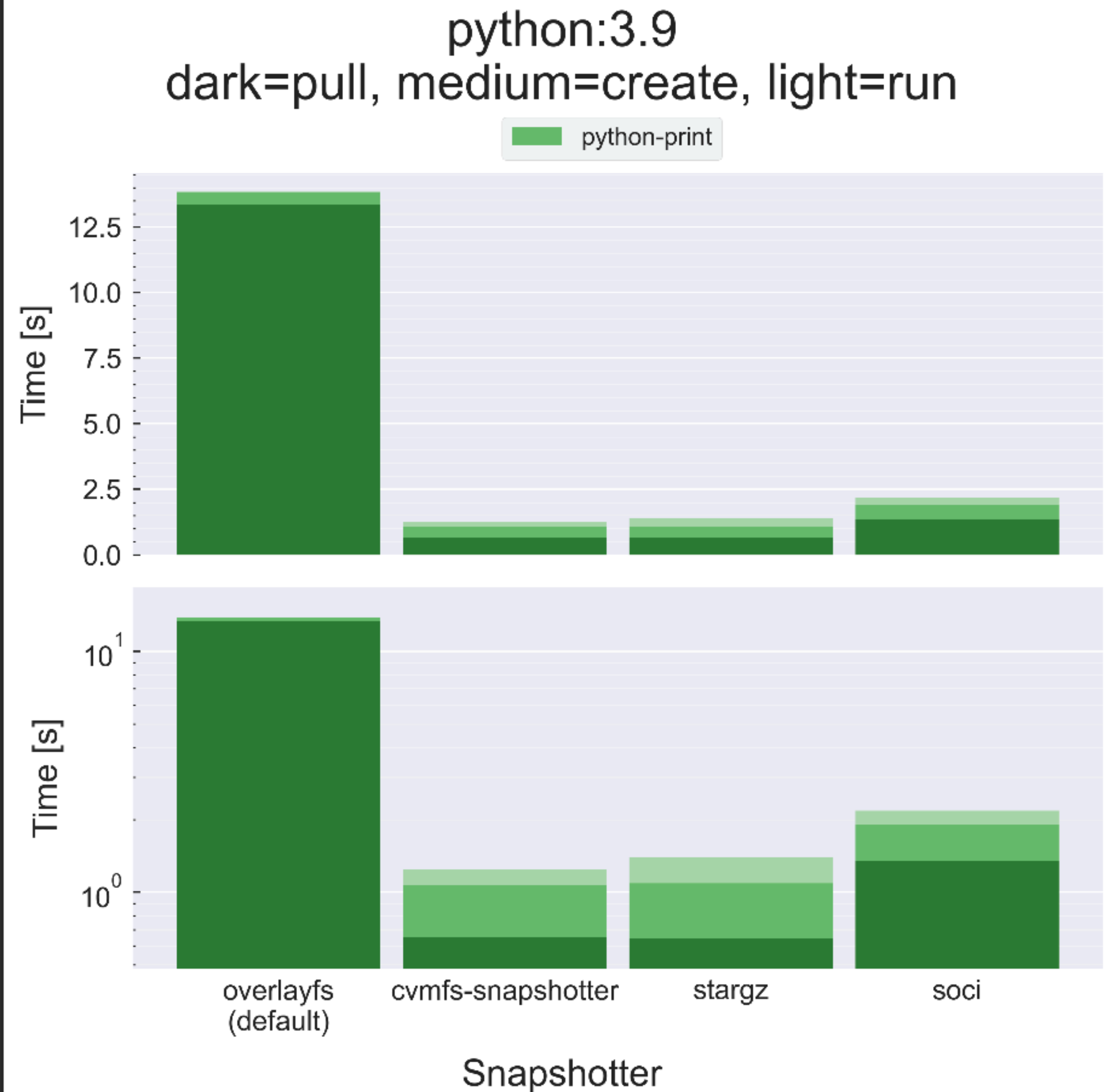
# Results: python image: print() — remote registry/cache

linear scale

log scale

# Results: python image: print() — local registry/cache
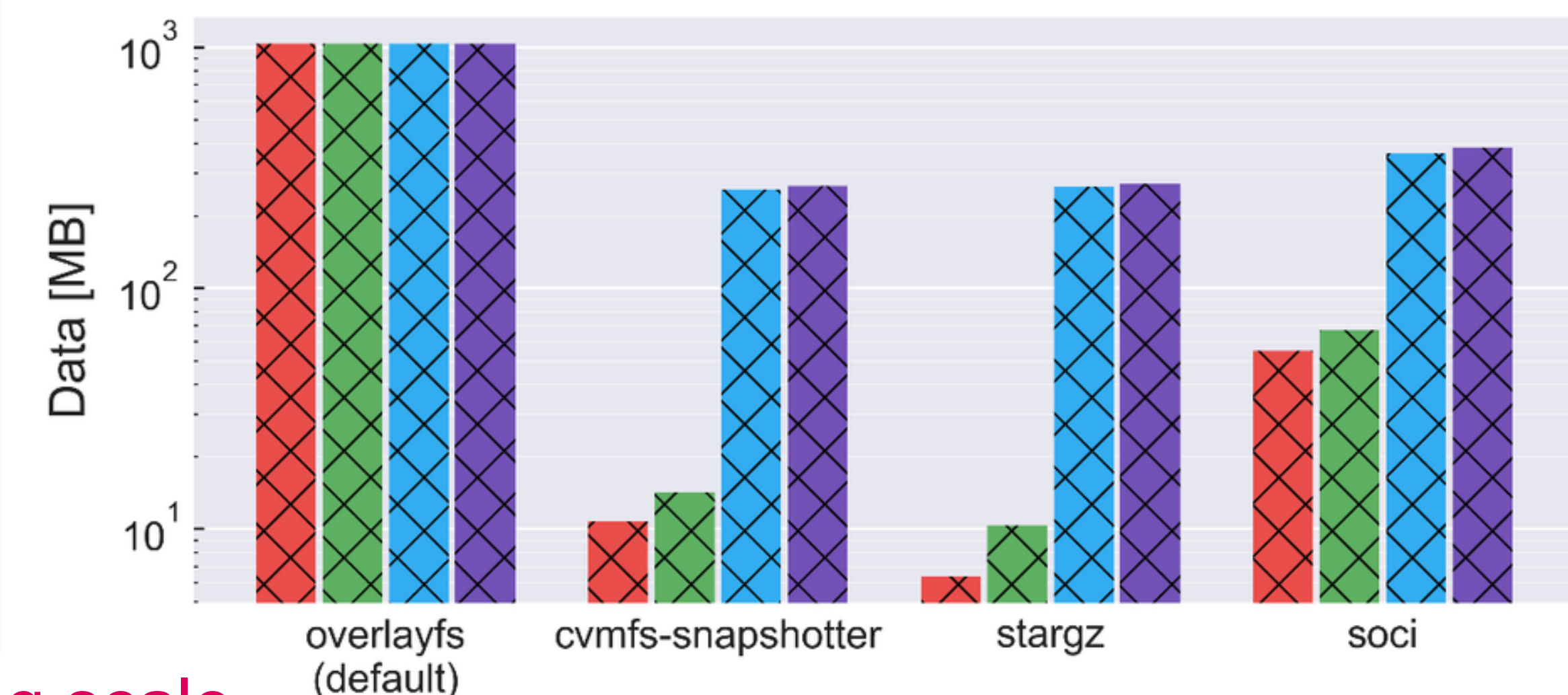
linear scale

log scale



Observations:
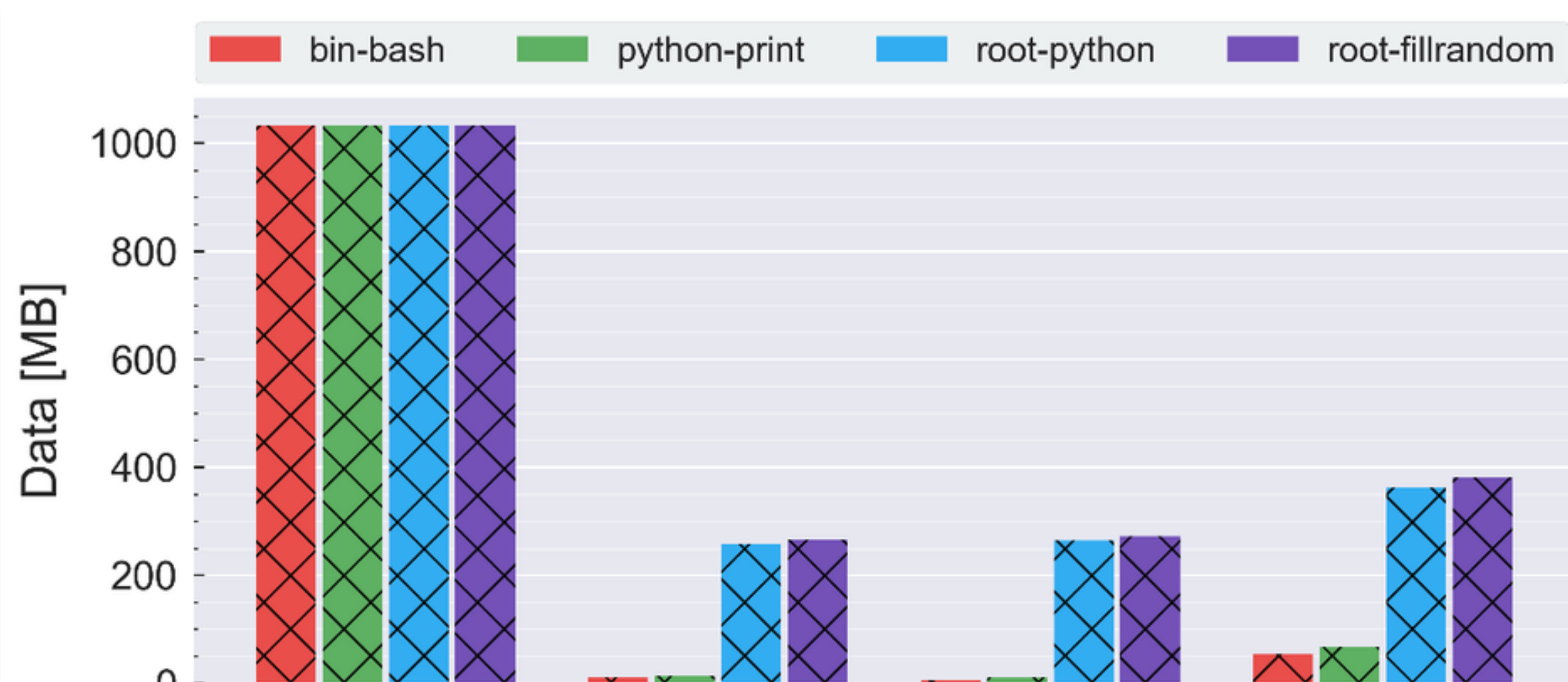
> Very similar behaviour w.r.t. remote registry/cache

> Slightly faster due to reduced network overhead

# Results: ROOT image

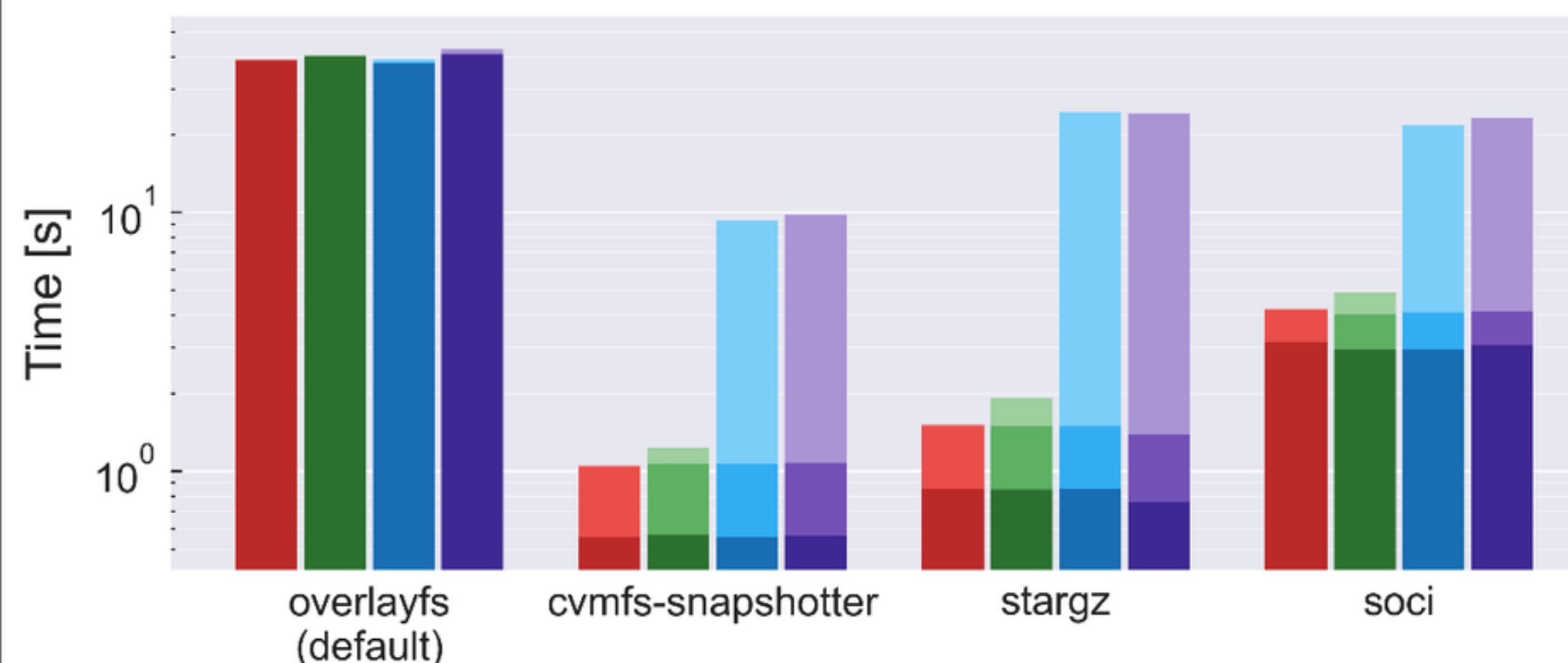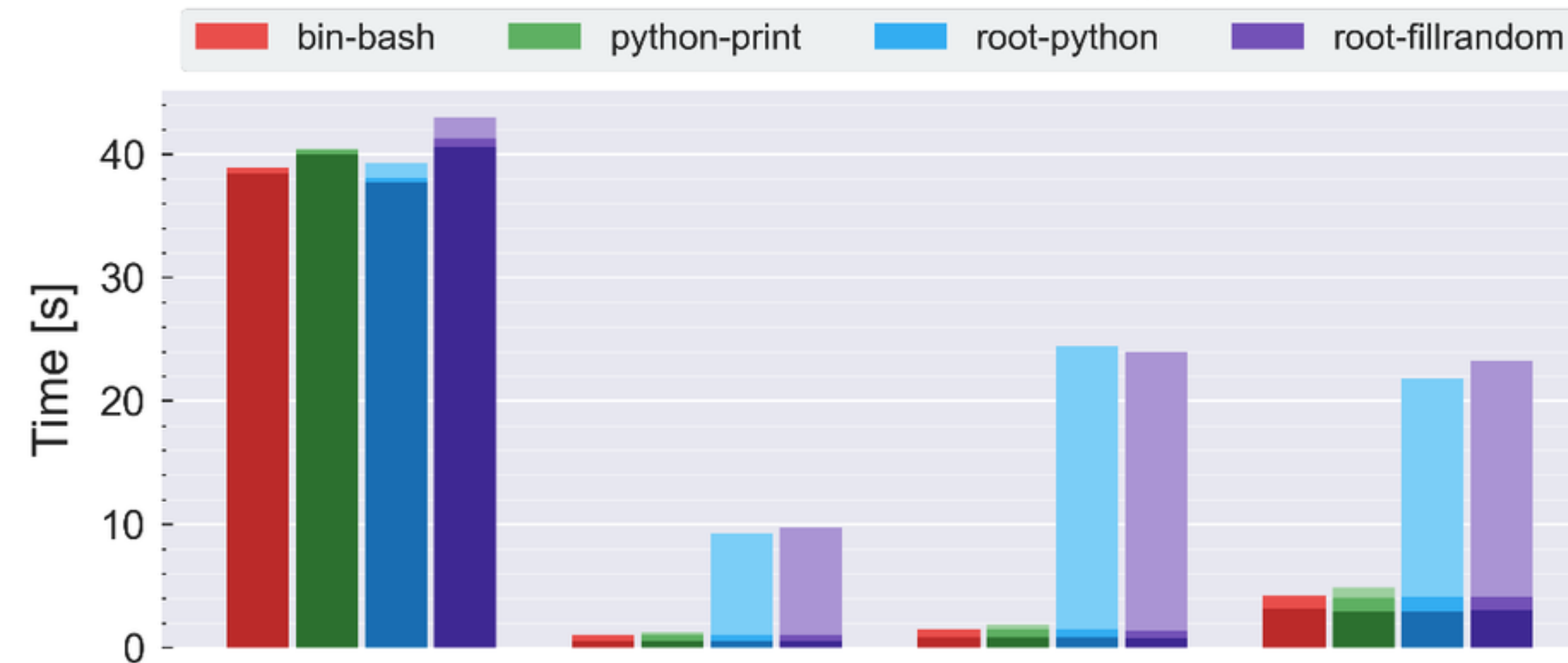Clemens Lange — Efficient and fast container execution using image snapshotters