# Simplifying the creation of Slurm client environments

## A straw for your Slurm beverage

Pablo Llopis Sanmillán - FOSDEM 23
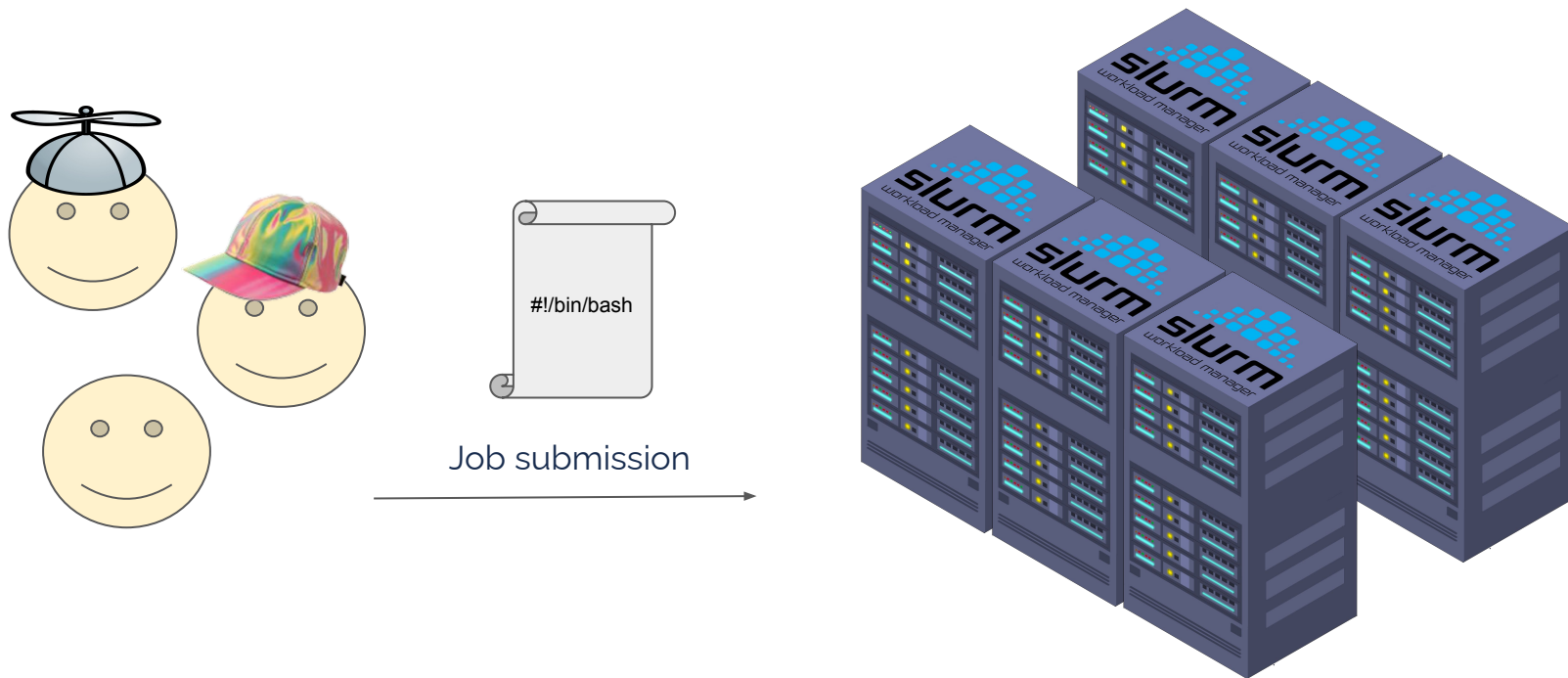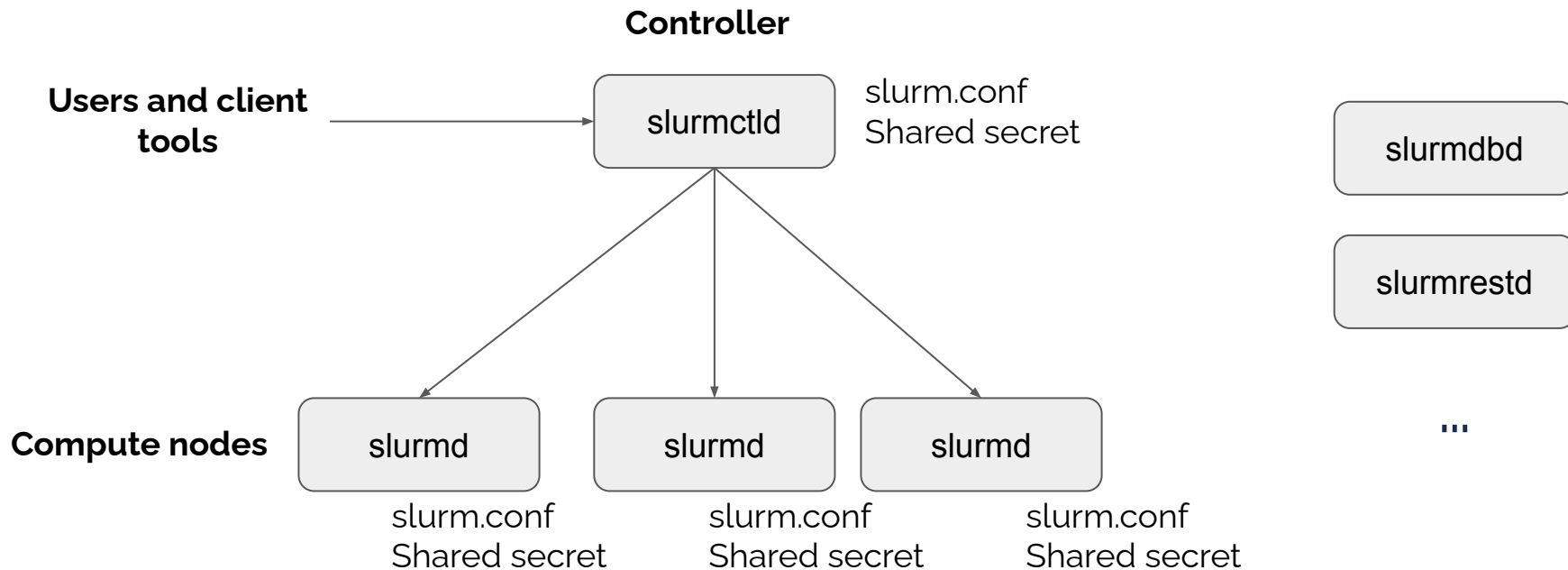
EPFL

# Slurm: a brief introduction

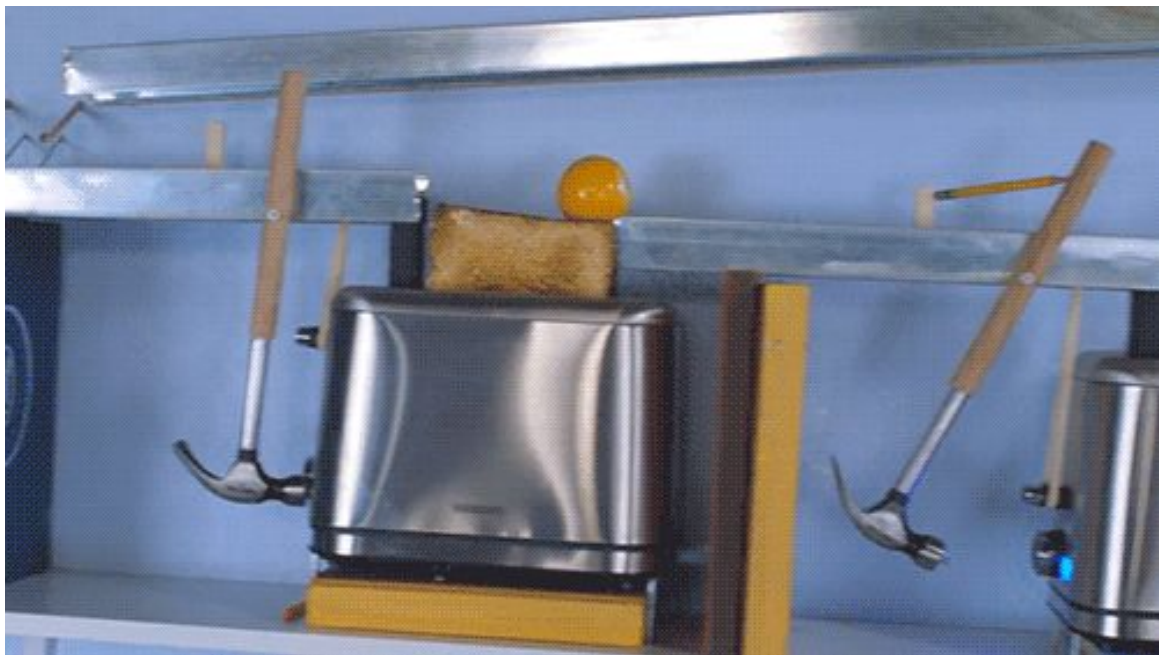Slurm is both a Resource Manager, and a Job Scheduler



#!/bin/bash

Job submission

# Slurm: a brief introduction

At its core, Slurm consists of a controller daemon, and client daemons

**Controller**

Users and client tools → slurmctld

slurm.conf
Shared secret

slurmdbd

slurmrestd

**Compute nodes**

slurmd          slurmd          slurmd

slurm.conf      slurm.conf      slurm.conf
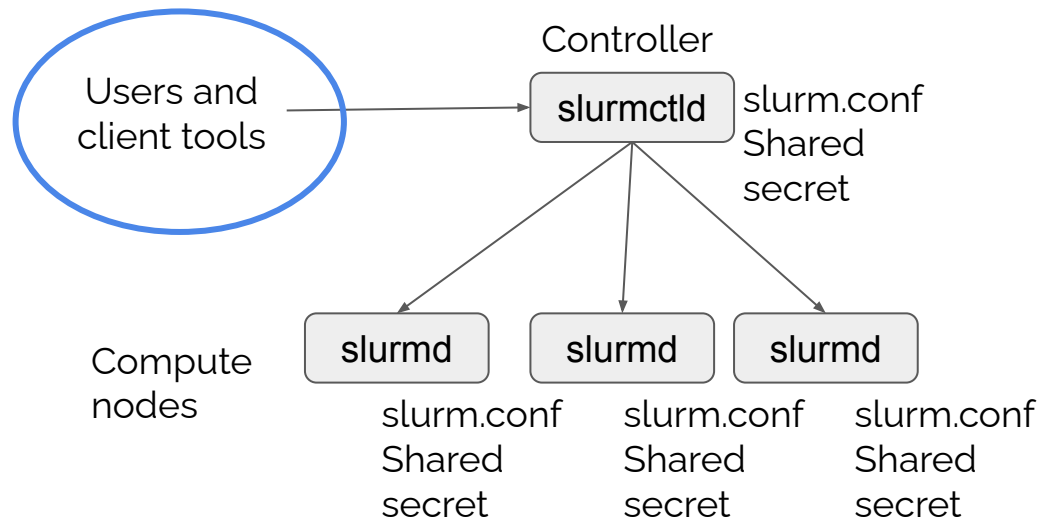Shared secret   Shared secret   Shared secret

...

# Containers

Containers are increasingly becoming a popular tool to run, automate deployments, and test modern infrastructure.
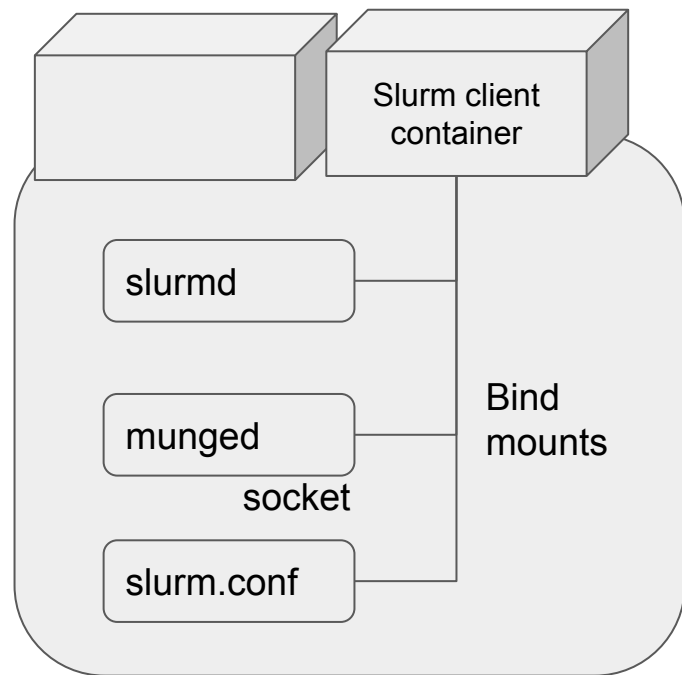
# Containers + Slurm: use cases

- Monitoring

- Health checks

- Accounting

- Integration with other
  services
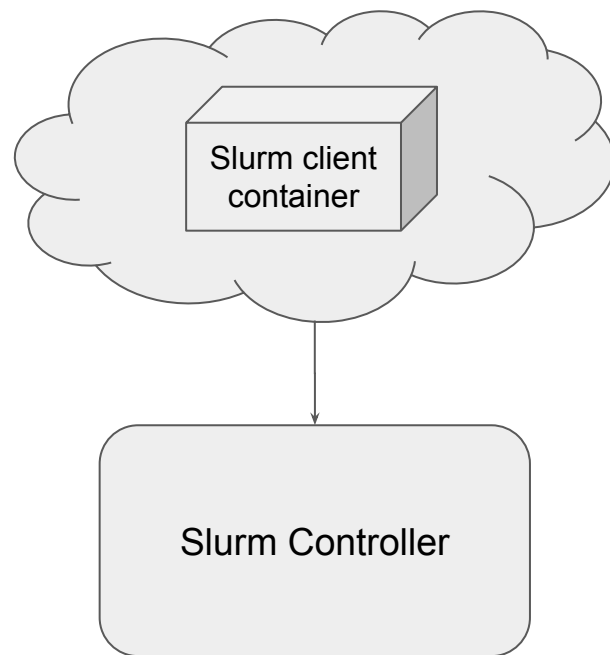
# Containers + Slurm

The local use-case

The distributed/remote use-case

Slurm client container

slurmd

munged

socket

slurm.conf

Bind mounts

Slurm client container

Slurm Controller
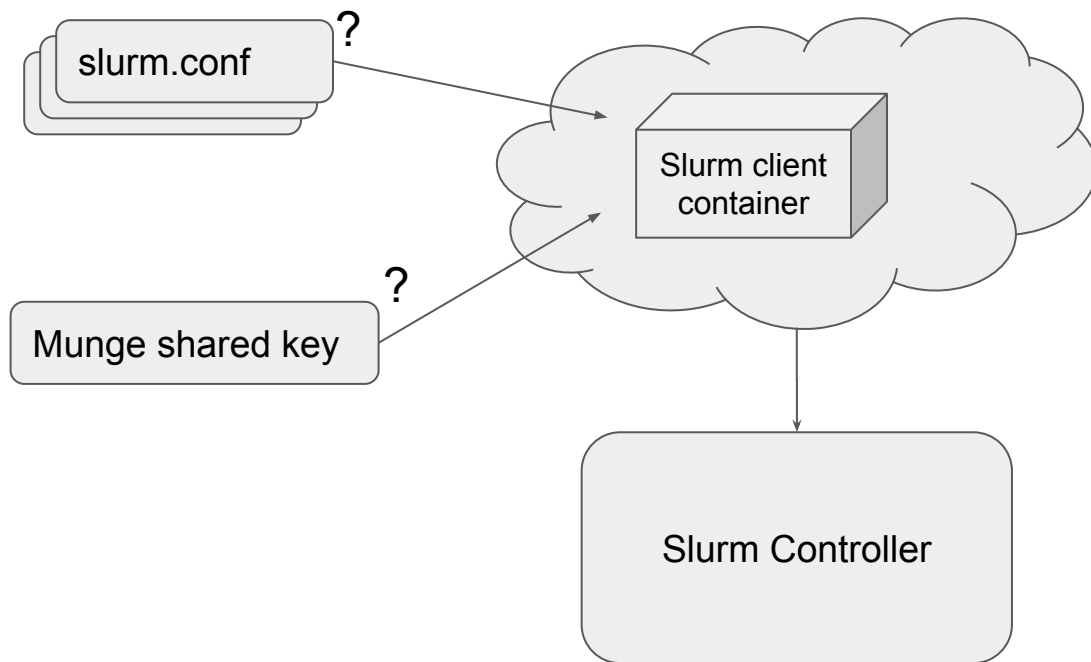
# Containers + Slurm

The distributed/remote use-case

# Containers + Slurm: the bad

```
87
88    COPY slurm.conf /etc/slurm/slurm.conf
```

This will absolutely work.  But it's often not necessarily good practice for maintainability.

# Containers + Slurm: the bad

```
87
88    COPY slurm.conf /etc/slurm/slurm.conf
```

This will absolutely work.  But it's often not necessarily
good practice for maintainability.

Solution: Use Slurm's configless* feature!    🙂 👍🏻

*Since Slurm 20.02

# Containers + Slurm: the ugly

```
25  # For some reason containers do not mount cgroups with file
26  # /sys/fs/cgroup/$subsystem/release_agent present.
27  # This file simply doesn't exist right after spawning the container (maybe due to container escape explo
28  # Slurm expects this file to exist. However, if we mount the cgroup subsystem manually
29  # after the container has already been spawned, the release_agent file will be there (??).
30  # We therefore umount them all, and then rely on Slurm's CgroupAutomount=yes to mount the cgroup subsyst
31  # ...
32  # Try containers, they said. It will be fun, they said.
33  umount /sys/fs/cgroup/freezer
34  umount /sys/fs/cgroup/cpuset
35  umount /sys/fs/cgroup/devices
36  umount /sys/fs/cgroup/cpuacct
37  umount /sys/fs/cgroup/memory
```

```
10  # More hacks needed for kubernetes:
11  # We will want to share the munge socket between containers.
12  # The main way to achieve this is to use the volume emptydir pattern in the pod,
13  # and let containers mount /var/run/munge via volumemounts.
14  # Munge checks and refuses to run if the directory has group write permissions enabled.
15  # But Kubernetes does not have a way to let us choose the directory mode and ownership.
16  # Therefore, make sure it is running in our deside mode and ownership.
17  # While we're at it, make sure we do not run into a "socket file already exists" in case
18  # the container is restarted independently from the pod.
19  # Isn't it nice when your infra is modern and declarative and doesn't need shell scripting everywher
20  rm -rf /var/run/munge/*
21  chmod 0755 /var/run/munge
22  chown munge:munge /var/run/munge
23  sudo -u munge /sbin/munged
```

```
3  # Another kubernetes-specific hack:
4  # Munge doesn't support secret as symlinks, and kubernetes forcibly
5  # presents secrets as symlinks. So we need to make a copy.
6  cp /secrets/munge.key /etc/munge/
```

# Containers + Slurm: the ugly

## Separate config files approach

- Manage a copy of slurm config files.

  (might be a challenge to keep a single, consistent source of truth)

- You will also need munged.

- And the munge key.

## Configless approach

- Add slurmd into your container to benefit from configless.

- You will also need munged.

- And the munge key.

# Containers + Slurm: the good?

A one-shot CLI tool that

- Authenticates to the controller (either munge or **JWT**)

- Fetches the Slurm config files

# Containers + Slurm: the good?

A one-shot CLI tool that

- Authenticates to the controller (either munge or **JWT**)

- Fetches the Slurm config files

Straw: A tool to fetch Slurm config files 🥤

https://github.com/pllopis/straw

# Straw in action

```
[pllopis@fedora straw]$ python straw.py -h
usage: straw.py [-h] [--auth {munge,jwt}] [-o OUTPUT_DIR] [-v] [-V]
[-l] server [server ...] version

positional arguments:
  server                  slurmctld server in server[:port] notation
  version                 Slurm major version that corresponds to that
of the slurmctld server (e.g. 22.05)

options:
  -h, --help              show this help message and exit
  --auth {munge,jwt}    Authentication method (default: jwt)
  -o OUTPUT_DIR, --output-dir OUTPUT_DIR
                          Existing output directory where config files
will be saved (default: ./)
  -v, --verbose           Increase output verbosity. Rrepetitions
allowed. (default: None)
  -V, --version           show program's version number and exit
  -l, --list              List available protocol versions (default:
False)
```

# Straw in action

```
[pllopis@fedora straw]$ python straw.py -l
22.05
21.08
20.11
[pllopis@fedora straw]$ echo $SLURM_JWT
eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9.eyJleHAiOjE3MDE2NDI5NjgsImlhdCI6
MTY3MDEwNjk2OCwic3ViIjoicm9vdCJ9.mhesNN8venwBDgXQNolzdi__QQbmV8jYm2BV
lTRi47c
[pllopis@fedora straw]$ python straw.py --auth jwt localhost 22.05
[pllopis@fedora straw]$ python straw.py -v --auth jwt localhost 22.05
Using authentication method: jwt
Trying localhost:6817...
SlurmdSpoolDir=/var/spool/slurm/d
[pllopis@fedora straw]$
```

# Conclusions

## Conclusions

- Straw can simplify the **cost** of Slurm client container creation.

- Straw can increase the **security** of Slurm integrations with other services.

## Caveats

- It would be even better if Straw didn't exist! Ideally this would be supported natively by Slurm.

- JWT tokens still need to belong to SlurmUser to be able to pull the config. (Slurm implementation limitation)

**Try it out!**   https://github.com/pllopis/straw