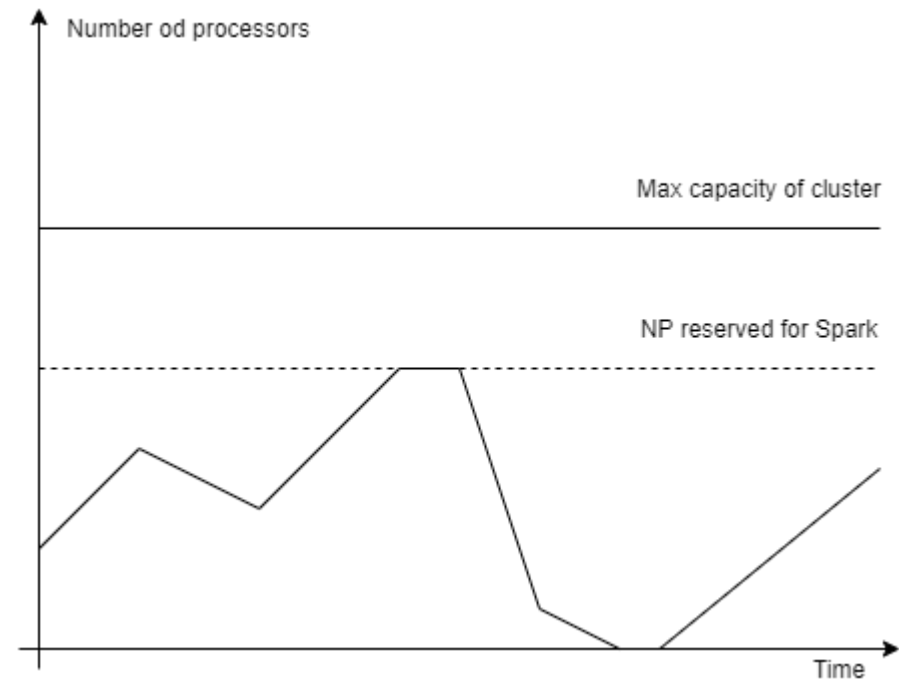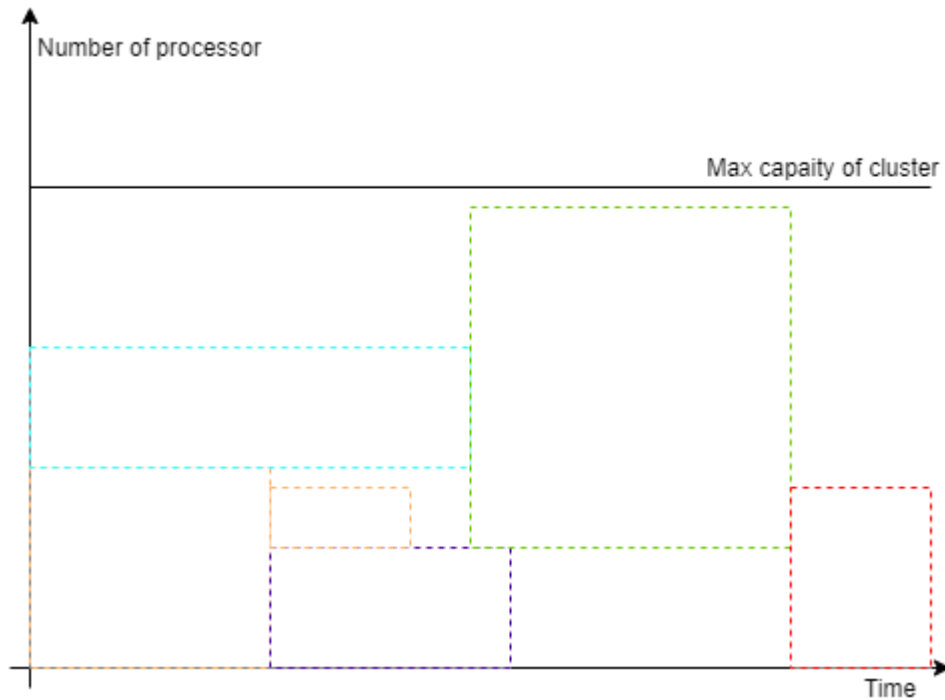# Self adjusted auto provision system at resource level

Weekly report 15th July
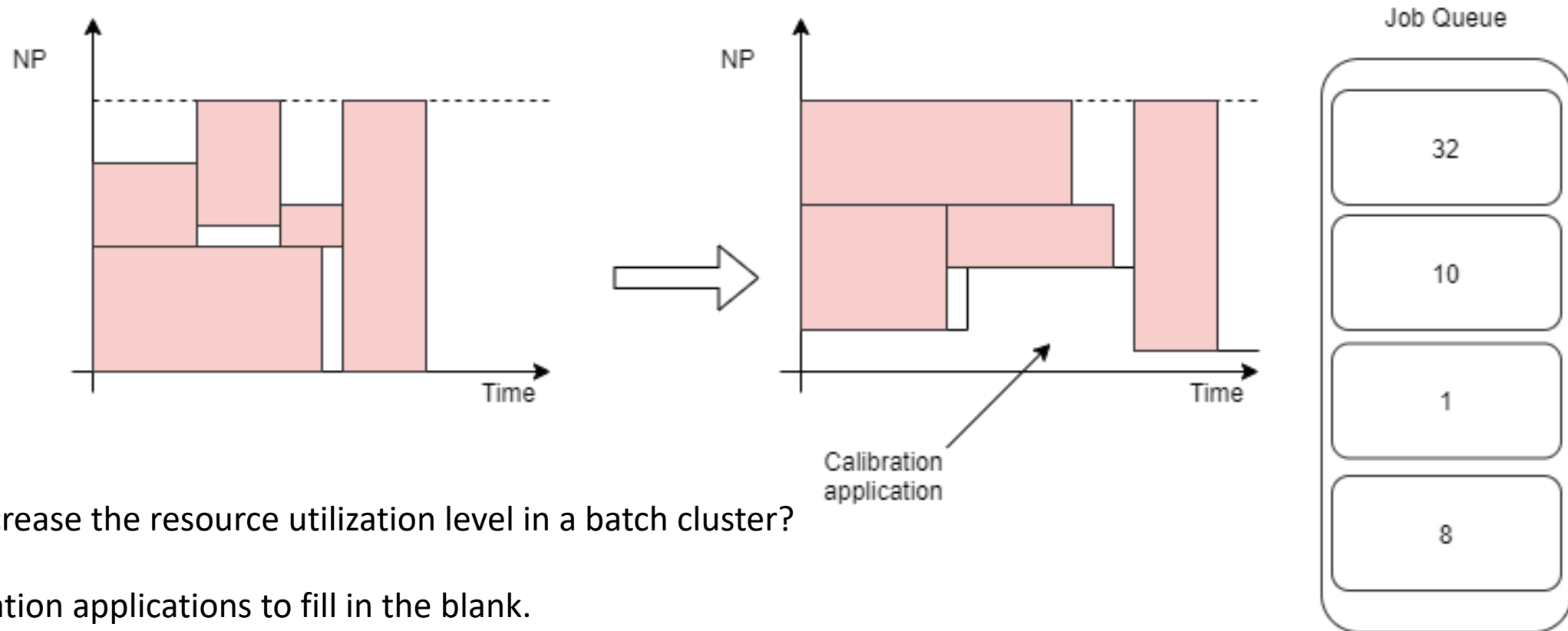
You Hu

# Recap – issue and research question



Existing solutions result in resource waste

# Recap – issue and research question
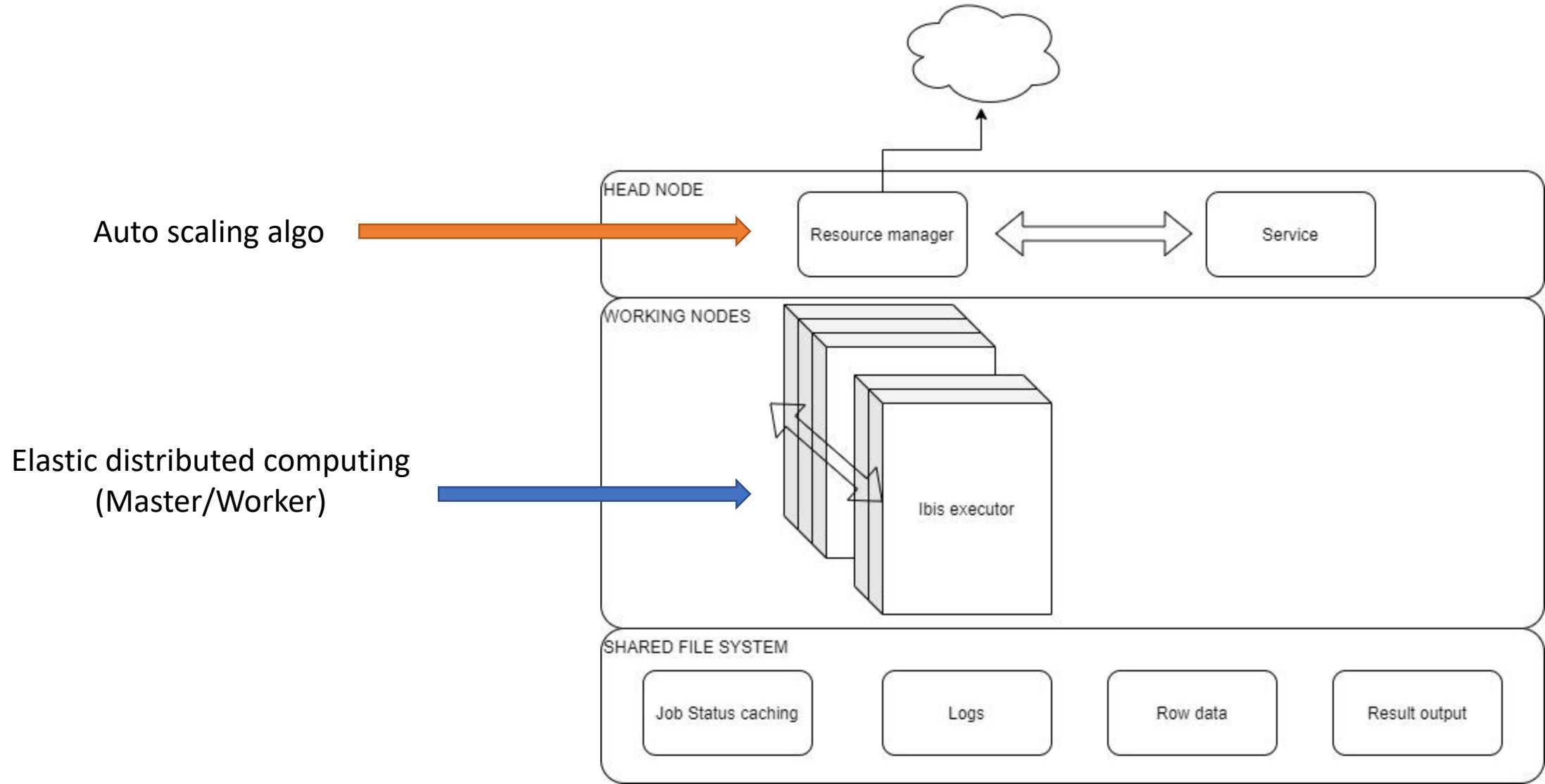


How to increase the resource utilization level in a batch cluster?

Use calibration applications to fill in the blank.

It requires: auto scaling/provisioning;
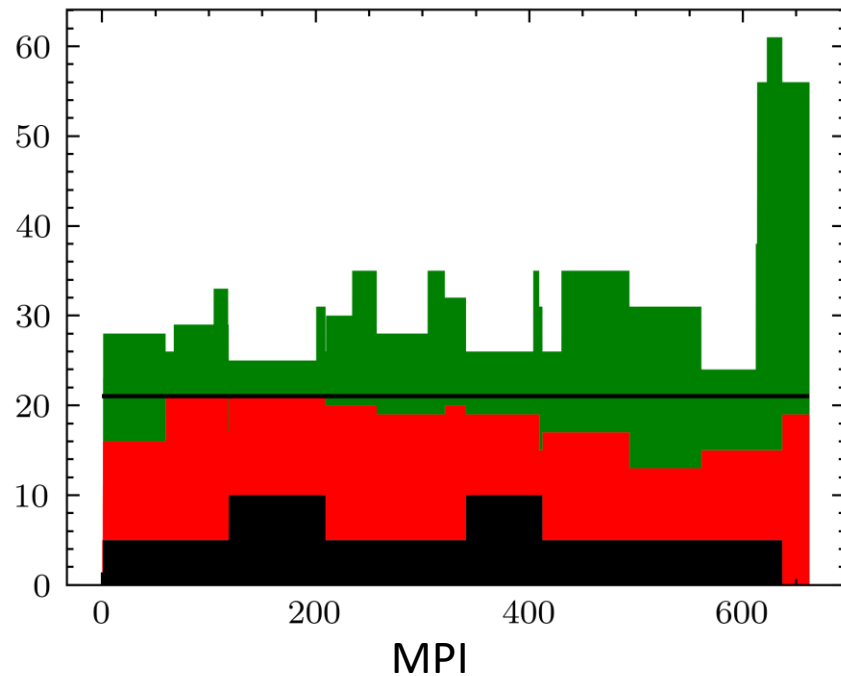
# Design-current layout

Auto scaling algo

Elastic distributed computing
(Master/Worker)

HEAD NODE

Resource manager ⟷ Service

WORKING NODES

Ibis executor

SHARED FILE SYSTEM

Job Status caching | Logs | Row data | Result output

# Auto scaling simulation

- Calibration resource: 47.8% increase
- Other resource: 1% increase
- Total resource: 16.8% increase

In total check point:39012
In total Calibration use:235280.0 AVG:6.030964831333948
In total Normal use:464868 AVG:11.91602583820363
In total use:700148.0 AVG:17.946990669537577
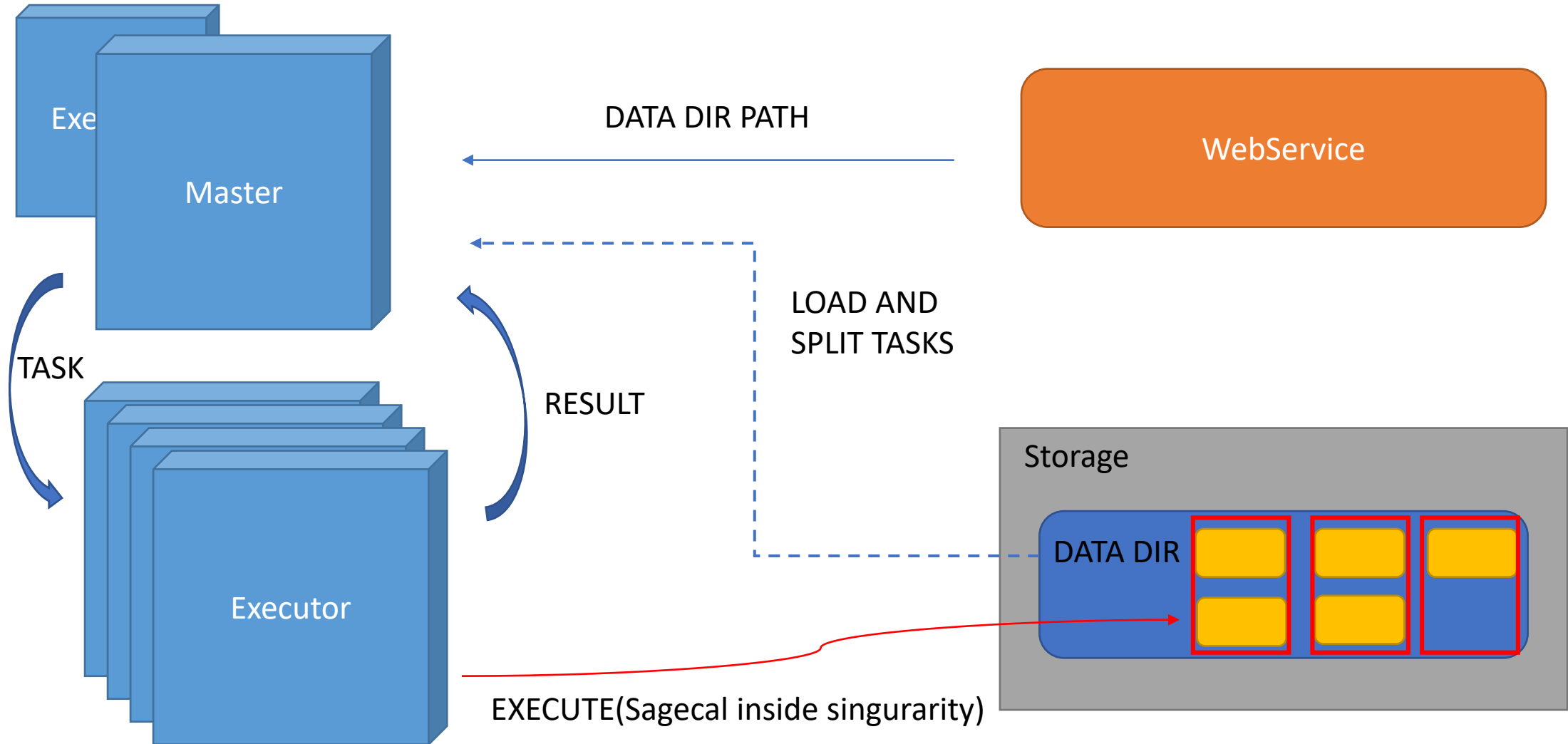
In total check point:38506
In total Calibration use:343185.0 AVG:8.912507141744143
In total Normal use:463826.0 AVG:12.045551342647899
In total use:807011.0 AVG:20.958058484392044

85.5% → 99.8%

Utilization rate

Pending jobs

Other Jobs

Calibration

Fixed mini node = 4

MPI

AUTO SCALING

# DynPrvDriver – compute layer

# How data processed

```
→ Calibration time singularity exec ~/DynPrvDriver/AppContainers/Sagecal/SagecalContainer.simg /opt/sagecal/bin/sagecal -d DATA0 -s 3c196.sky.txt -c 3c196.sky.txt.cluster -n 4 -t 10 -p sm.ms.solutions -e 4 -g 2 -l 10 -m 7 -x 30 -F 1 -j 5  -k -1 -B 1 -W
SAGECal 0.7.1 (C) 2011-2020 Sarod Yatawatta
 MS: DATA0
Selecting baselines > 30 and < 1e+08 wavelengths.
Using Robust noise model for solver with degrees of freedom [2,30].
Stations: 61 Baselines: 1830
Integration Time: 10.0139 s, Total timeslots: 125
Phase center (2.15374, 0.841552)
Only one MS
Got 2 clusters
Total effective clusters: 3
For 10 samples, solution time interval (s): 100.139
Freq: 152.733 MHz, Chan: 1 Bandwidth: 0.183105 MHz
nu=30
Timeslot: 10 Residual: initial=0.0870613,final=3.73346e-07, Time spent=0.0666667 minutes
nu=30
Timeslot: 20 Residual: initial=5.83557e-07,final=1.91819e-07, Time spent=0.0333333 minutes
nu=30
Timeslot: 30 Residual: initial=2.49905e-07,final=1.48017e-07, Time spent=0.0166667 minutes
nu=30
Timeslot: 40 Residual: initial=1.53844e-07,final=1.43131e-07, Time spent=0.0333333 minutes
nu=30
Timeslot: 50 Residual: initial=1.50543e-07,final=1.41432e-07, Time spent=0.0333333 minutes
nu=30
Timeslot: 60 Residual: initial=1.47163e-07,final=1.37958e-07, Time spent=0.0166667 minutes
nu=30
Timeslot: 70 Residual: initial=1.49591e-07,final=1.39044e-07, Time spent=0.0166667 minutes
nu=30
Timeslot: 80 Residual: initial=1.45653e-07,final=1.41048e-07, Time spent=0.0333333 minutes
nu=30
Timeslot: 90 Residual: initial=1.4735e-07,final=1.43614e-07, Time spent=0.0166667 minutes
nu=30
Timeslot: 100 Residual: initial=1.45031e-07,final=1.42376e-07, Time spent=0.0333333 minutes
nu=30
Timeslot: 110 Residual: initial=1.45285e-07,final=1.42938e-07, Time spent=0.0166667 minutes
nu=30
Timeslot: 120 Residual: initial=1.4827e-07,final=1.44877e-07, Time spent=0.0333333 minutes
Warning: Missing rows, got 9455 expect 18300 +- 610. (probably the last time interval, so not a big issue).
nu=30
Timeslot: 130 Residual: initial=4.23302e-08,final=2.58443e-08, Time spent=0.0166667 minutes
Done.
singularity exec ~/DynPrvDriver/AppContainers/Sagecal/SagecalContainer.simg    27.88s user 1.82s system 126% cpu 23.493 total
```

# Current data & speed

- Dataset sm.ms [https://github.com/nlesc-dirac/data](https://github.com/nlesc-dirac/data)
  - Size 69M
- Speed:
  - -n 2: 31s
  - -n 4: 26s
  - -n 6: 20.7s
- Work around:
  - Big/mid/small data set = duplicate sm.ms x 500/150/50 30G+/10G/3G+
    - estm. 150/45/15Min*Node
  - The assumption behind: an observation can be divided into sub data sets
- Limit:
  - 40 GB disk quotum on DAS5
  - No Sky Map for sagecal

# Current data

The data Souley is using
L232873_SB[0-19].dppp.ms, 4.9 GB
L232875_SB[0-19] ].dppp.ms, 228GB

The data Onno is using
L246909_SAP006_SB486_uv_001.MS_021
34b61.tar, 22GB
L250048_SB311_uv.dppp.MS_b8783fe7.tar
, 3.2GB
L250048_SB342_uv.dppp.MS_32c7efc3.tar,
3.2GB

What do these fields mean? SBXXX, dppp
How can I get the sky maps of them?

```
[yhu310@fs0 L232873]$ ls
L232873_SB000_uv.dppp.MS   L232873_SB005_uv.dppp.MS   L232873_SB010_uv.dppp.MS   L232873_SB015_uv.dppp.MS
L232873_SB001_uv.dppp.MS   L232873_SB006_uv.dppp.MS   L232873_SB011_uv.dppp.MS   L232873_SB016_uv.dppp.MS
L232873_SB002_uv.dppp.MS   L232873_SB007_uv.dppp.MS   L232873_SB012_uv.dppp.MS   L232873_SB017_uv.dppp.MS
L232873_SB003_uv.dppp.MS   L232873_SB008_uv.dppp.MS   L232873_SB013_uv.dppp.MS   L232873_SB018_uv.dppp.MS
L232873_SB004_uv.dppp.MS   L232873_SB009_uv.dppp.MS   L232873_SB014_uv.dppp.MS   L232873_SB019_uv.dppp.MS
[yhu310@fs0 L232873]$ du -h --max-depth=1 .
248M    ./L232873_SB000_uv.dppp.MS
246M    ./L232873_SB001_uv.dppp.MS
246M    ./L232873_SB002_uv.dppp.MS
246M    ./L232873_SB003_uv.dppp.MS
246M    ./L232873_SB004_uv.dppp.MS
246M    ./L232873_SB005_uv.dppp.MS
246M    ./L232873_SB006_uv.dppp.MS
246M    ./L232873_SB007_uv.dppp.MS
246M    ./L232873_SB008_uv.dppp.MS
246M    ./L232873_SB009_uv.dppp.MS
246M    ./L232873_SB010_uv.dppp.MS
246M    ./L232873_SB011_uv.dppp.MS
246M    ./L232873_SB012_uv.dppp.MS
246M    ./L232873_SB013_uv.dppp.MS
246M    ./L232873_SB014_uv.dppp.MS
246M    ./L232873_SB015_uv.dppp.MS
246M    ./L232873_SB016_uv.dppp.MS
246M    ./L232873_SB017_uv.dppp.MS
246M    ./L232873_SB018_uv.dppp.MS
246M    ./L232873_SB019_uv.dppp.MS
4.9G    .
```
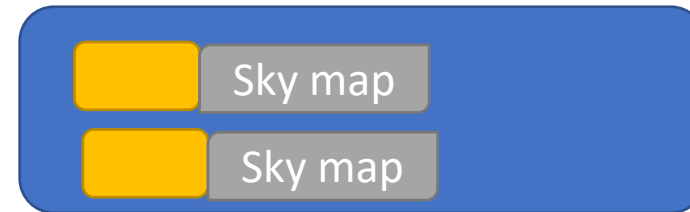
# How .MS data relates to sky map

L232873_SBXXX

One sky map vs multiple sub dataset

One sky map vs one sub dataset
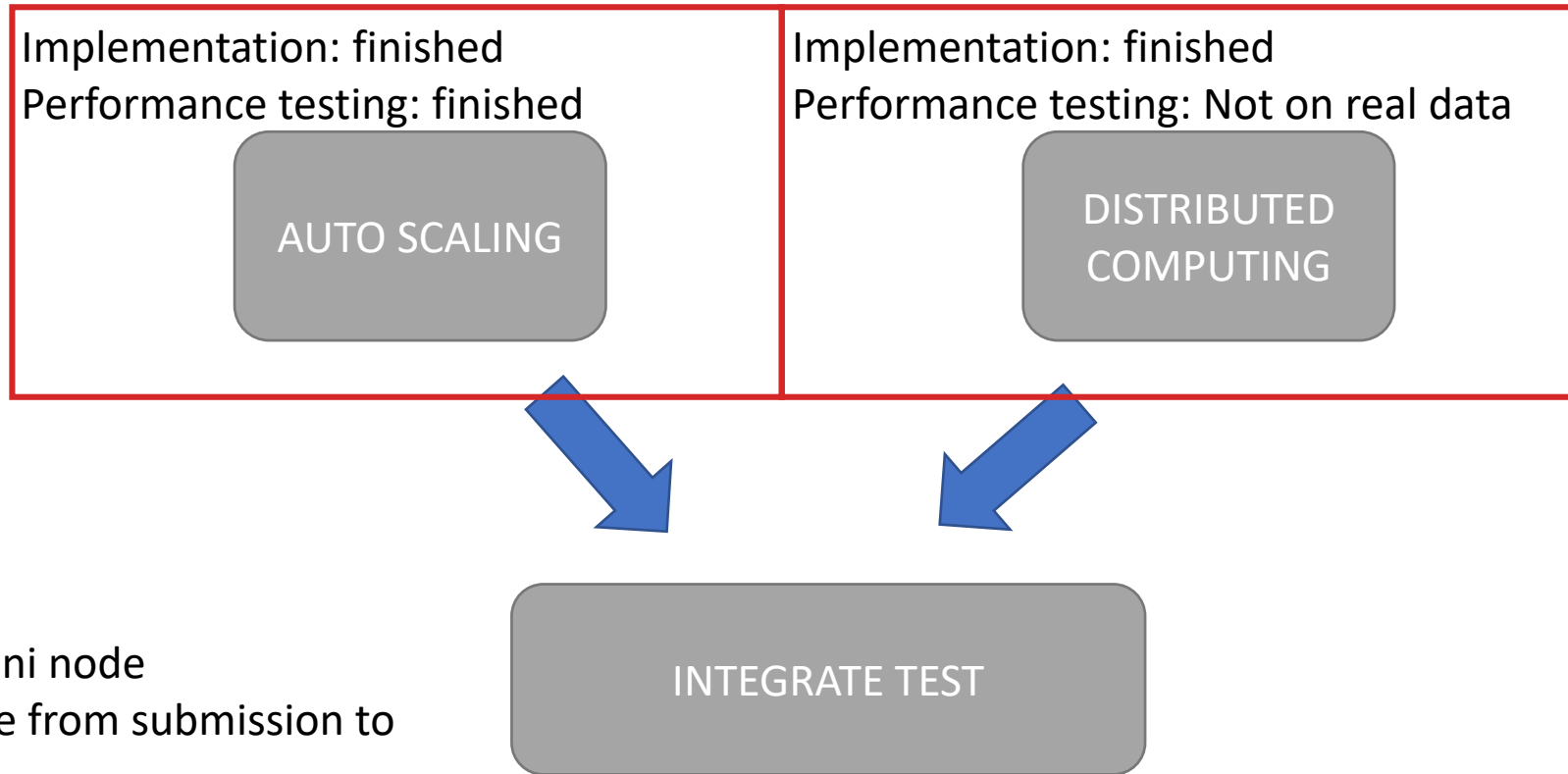
Sky Model Construction Using Shapelets

restore -f example.fits -i sky.txt -c sagecal_cluster.txt -l sagecal_sky.txt

# What is wrong?

How can we apply Sagecal on .MS data with skymap ?



```
[yhu310@fs0 data]$ singularity exec  -B /var/scratch/yhu310/:/var/scratch/yhu310/ ~/DynPrvDriver/AppContainers/Sagecal/S
agecalContainer.simg /opt/sagecal/bin/sagecal -d /var/scratch/yhu310/L232873/L232873_SB000_uv.dppp.MS -s ~/CalTest/Calib
ration/3c196.sky.txt -c ~/CalTest/Calibration/3c196.sky.txt.cluster  -n 4 -t 10 -p ~/sm.ms.solutions -e 4 -g 2 -l 10 -m
7 -x 30 -F 1 -j 5  -k -1 -B 1 -W 0
SAGECal 0.7.1 (C) 2011-2020 Sarod Yatawatta
 MS: /var/scratch/yhu310/L232873/L232873_SB000_uv.dppp.MS
Selecting baselines > 30 and < 1e+08 wavelengths.
Using Robust noise model for solver with degrees of freedom [2,30].
Stations: 62 Baselines: 1891
Integration Time: 2.01327 s, Total timeslots: 298
Phase center (2.15374, 0.841554)
Only one MS
Got 2 clusters
Total effective clusters: 3
For 10 samples, solution time interval (s): 20.1327
Freq: 120.311 MHz, Chan: 8 Bandwidth: 0.195312 MHz
terminate called after throwing an instance of 'casacore::TableError'
  what():  Table column CORRECTED_DATA is unknown
SIGABRT: abort
PC=0x473cdb m=0 sigcode=0
```

# Next phase: performance test

Implementation: finished
Performance testing: finished

**AUTO SCALING**

Implementation: finished
Performance testing: Not on real data

**DISTRIBUTED COMPUTING**

**INTEGRATE TEST**

Integrate test:
- Dynamic mini node
- Log the time from submission to the end