

Huawei Cloud Storage Evaluation and Testing with Prototype Services

Seppo S. Heikkila
CERN IT

Openlab Minor Review
29th of October 2013
CERN, Geneva



Motivation

- Cloud storage market is growing fast
- CERN uses custom made storage solutions

Question

“Are cloud storages able to meet the High Energy Physics (HEP) data storage requirements?”

Method

- Evaluate scalability and fault-tolerance
- Test with real applications

DSS

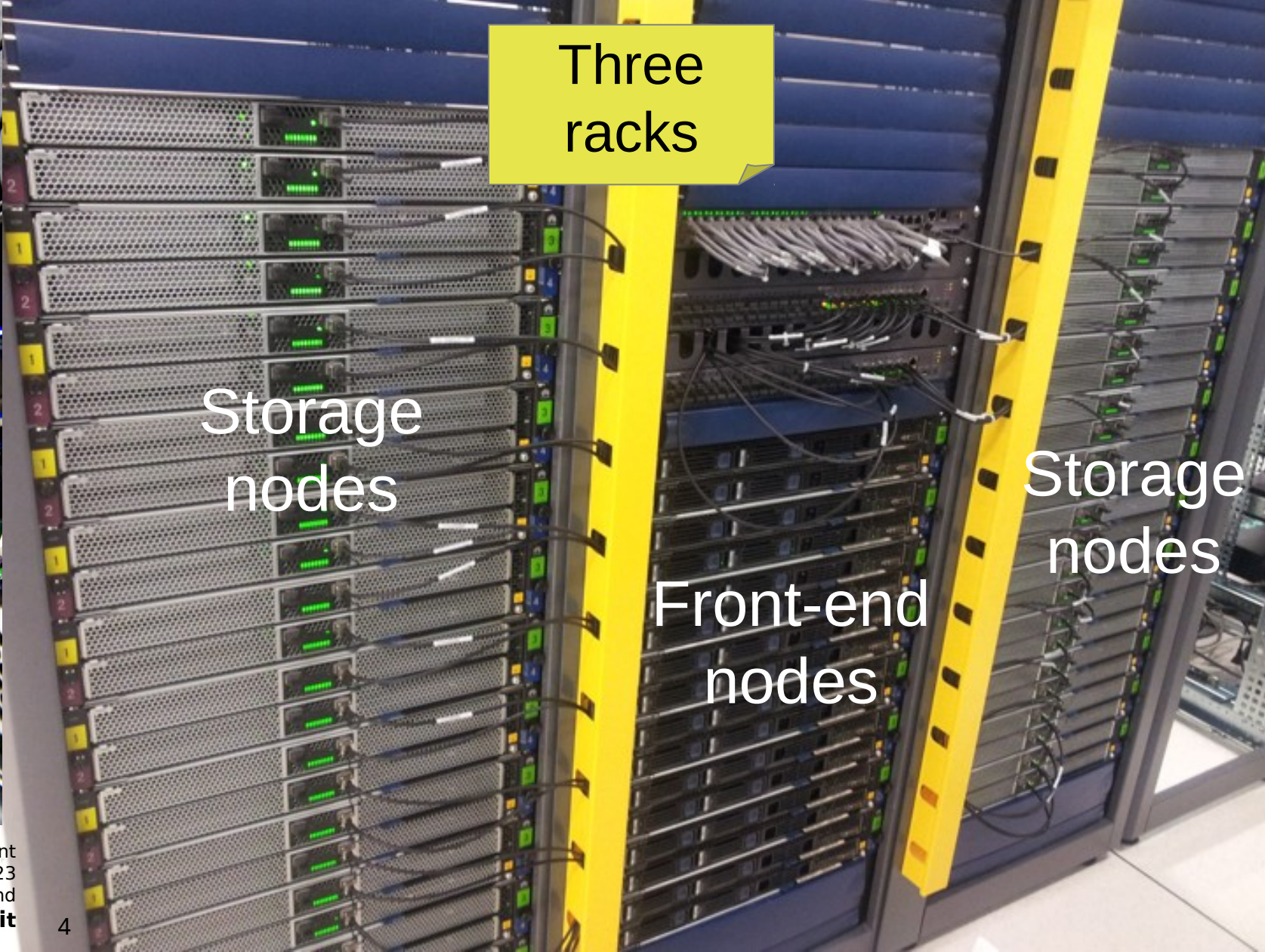
Huawei cloud storage

Location:
CERN
Computer
Center

“We are
now here”

“Cloud storage”



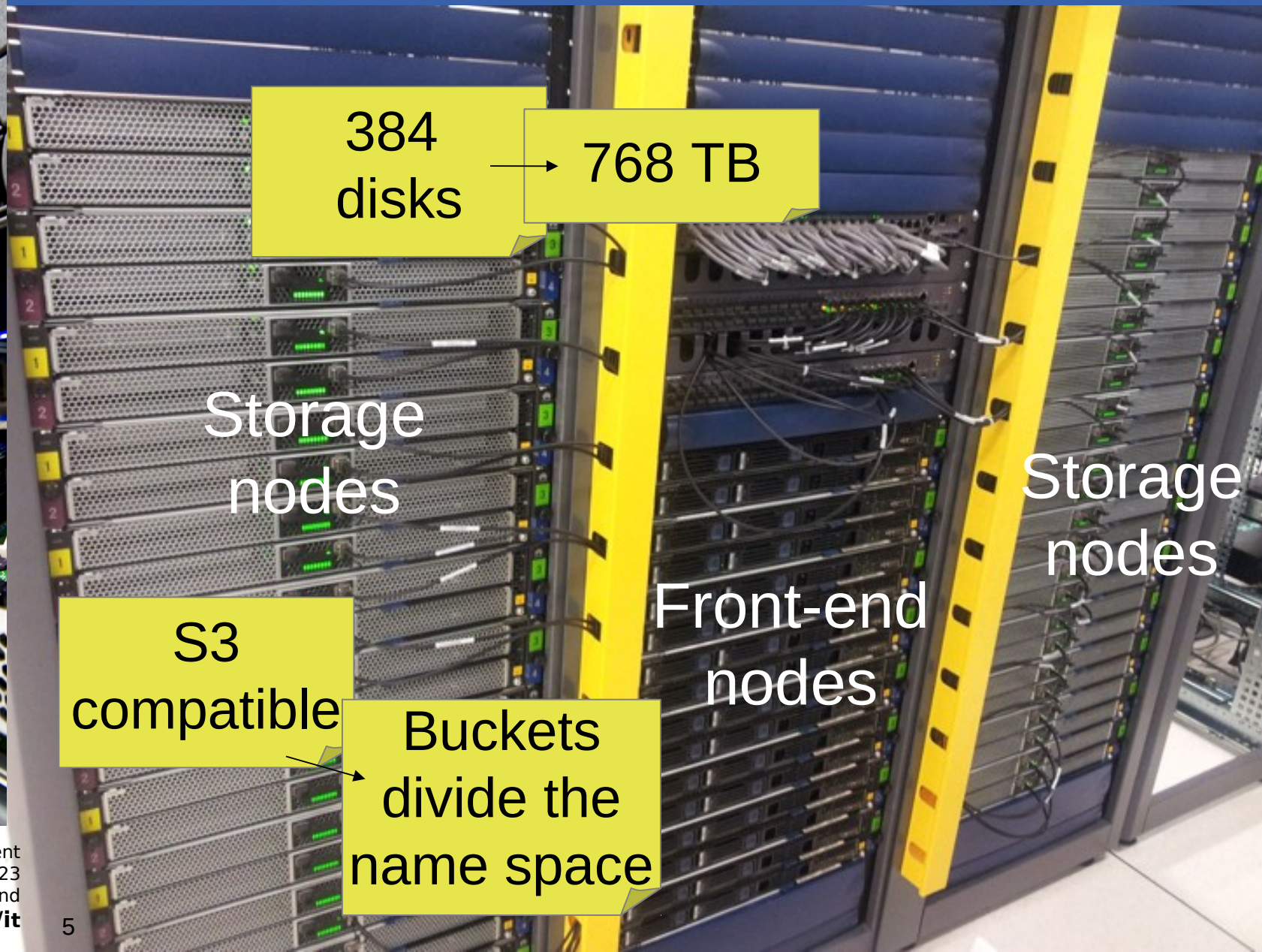


Three
racks

Storage
nodes

Front-end
nodes

Storage
nodes



384 disks → 768 TB

Storage nodes

Storage nodes

Front-end nodes

S3 compatible

Buckets divide the name space

Each blade has
eight storage nodes

One chassis has two
blades (16 disks)

1.5 years of Huawei...

Major
Review

Major
Review

Major
Review

Major
Review

Minor
Review

Board of
Sponsors

Minor
Review

Board of
Sponsors

Minor
Review

01/2012

01/2013

10/2013

Project
starts

First
tests

Upgrade
of the
system

Stress
testing

File-system
integration

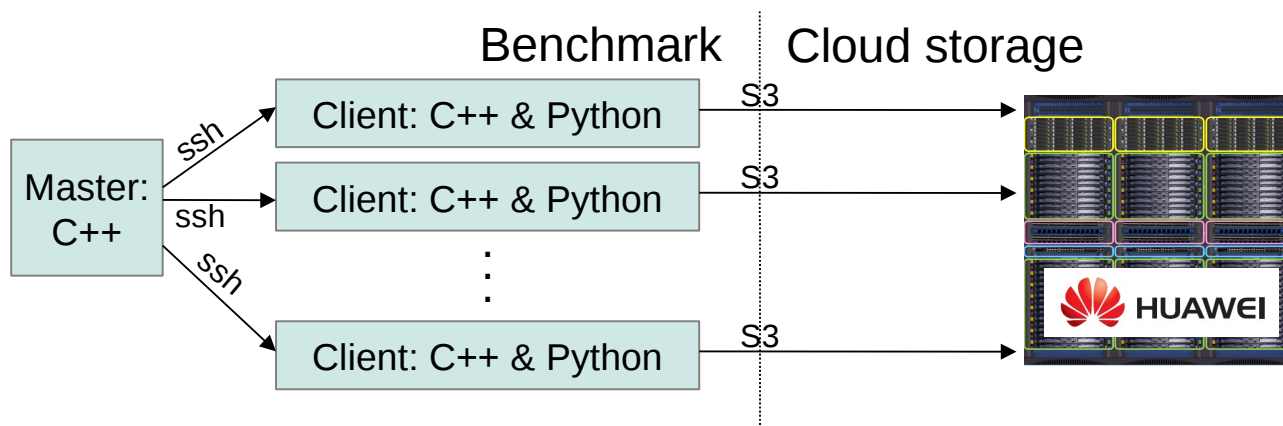
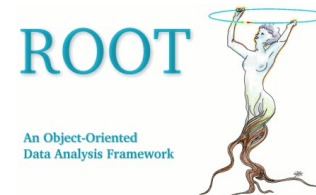
Commissioning
of the system

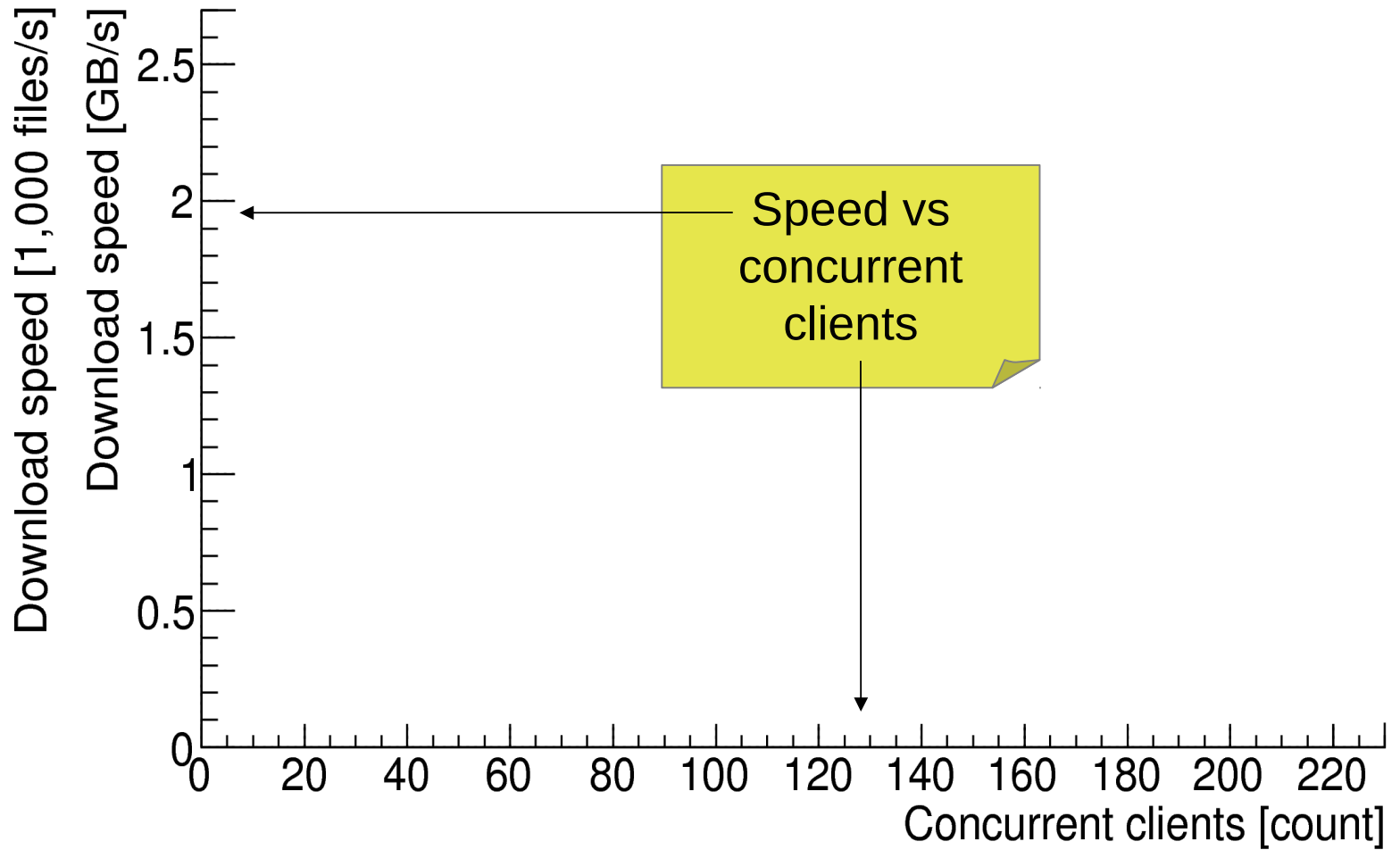
Failure
recovery
testing

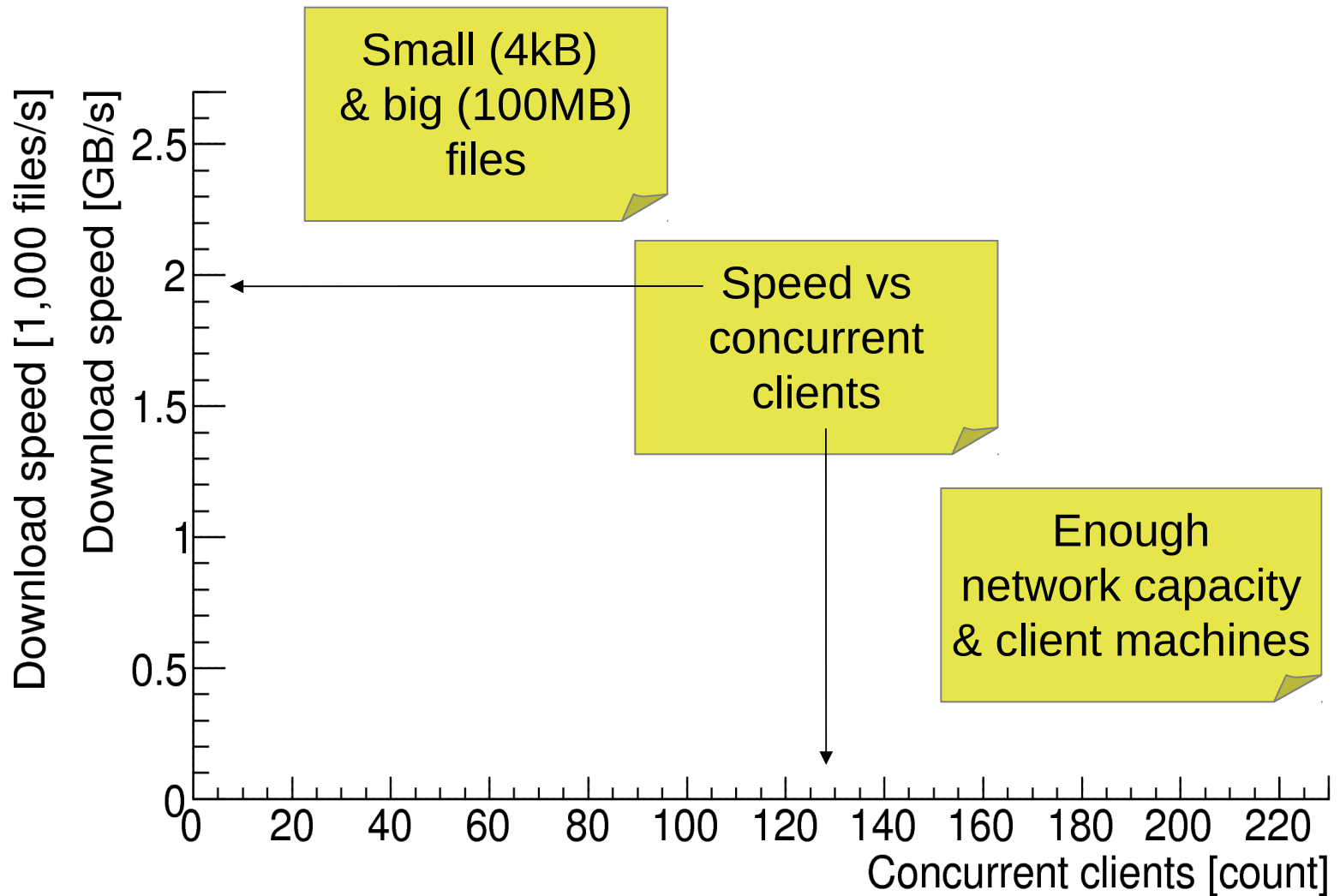
Full-scale
stress
testing

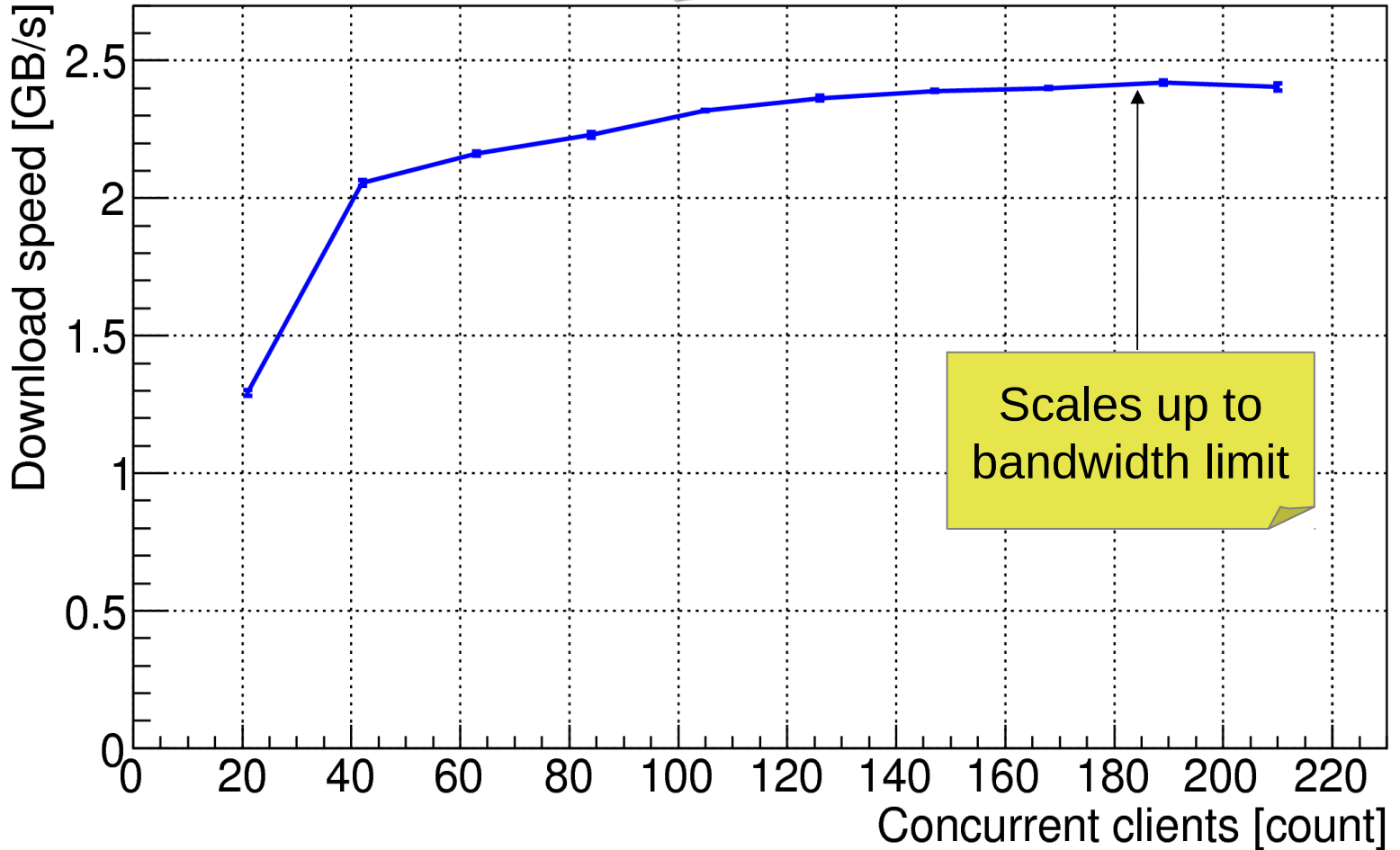
Distributed C++ benchmark

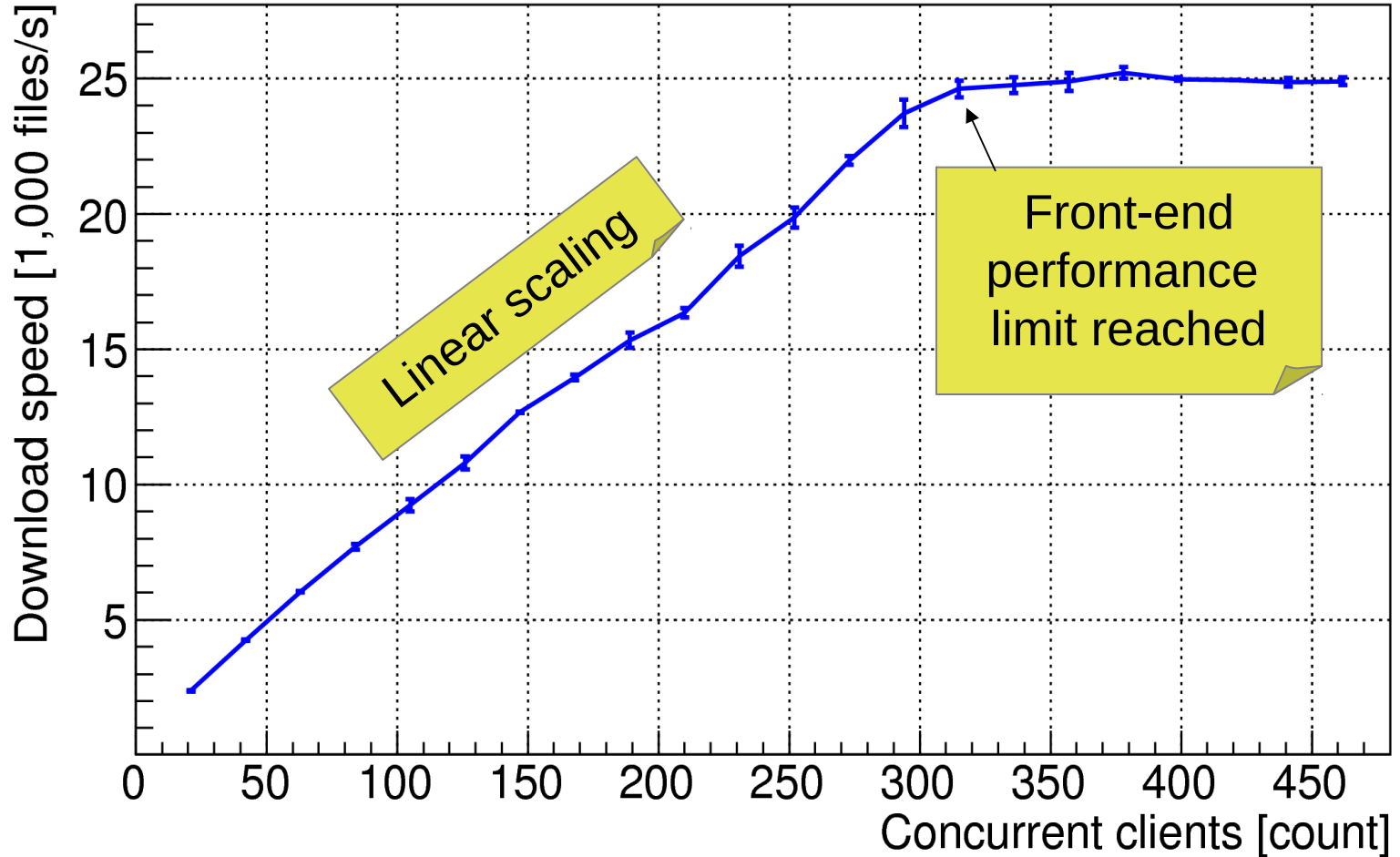
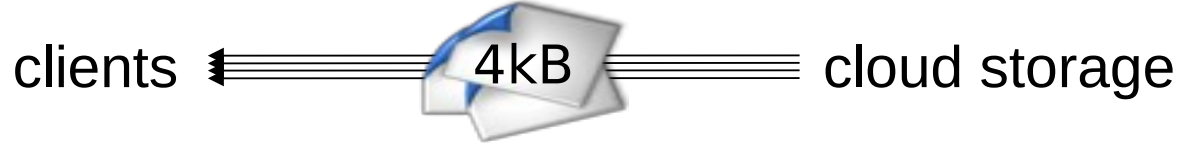
- Integrated with ROOT
- Client nodes connected with ssh
- S3 Python library to read and write files
- Histograms about specific metrics
 - Operation time, read/write speed, CPU/memory utilisation

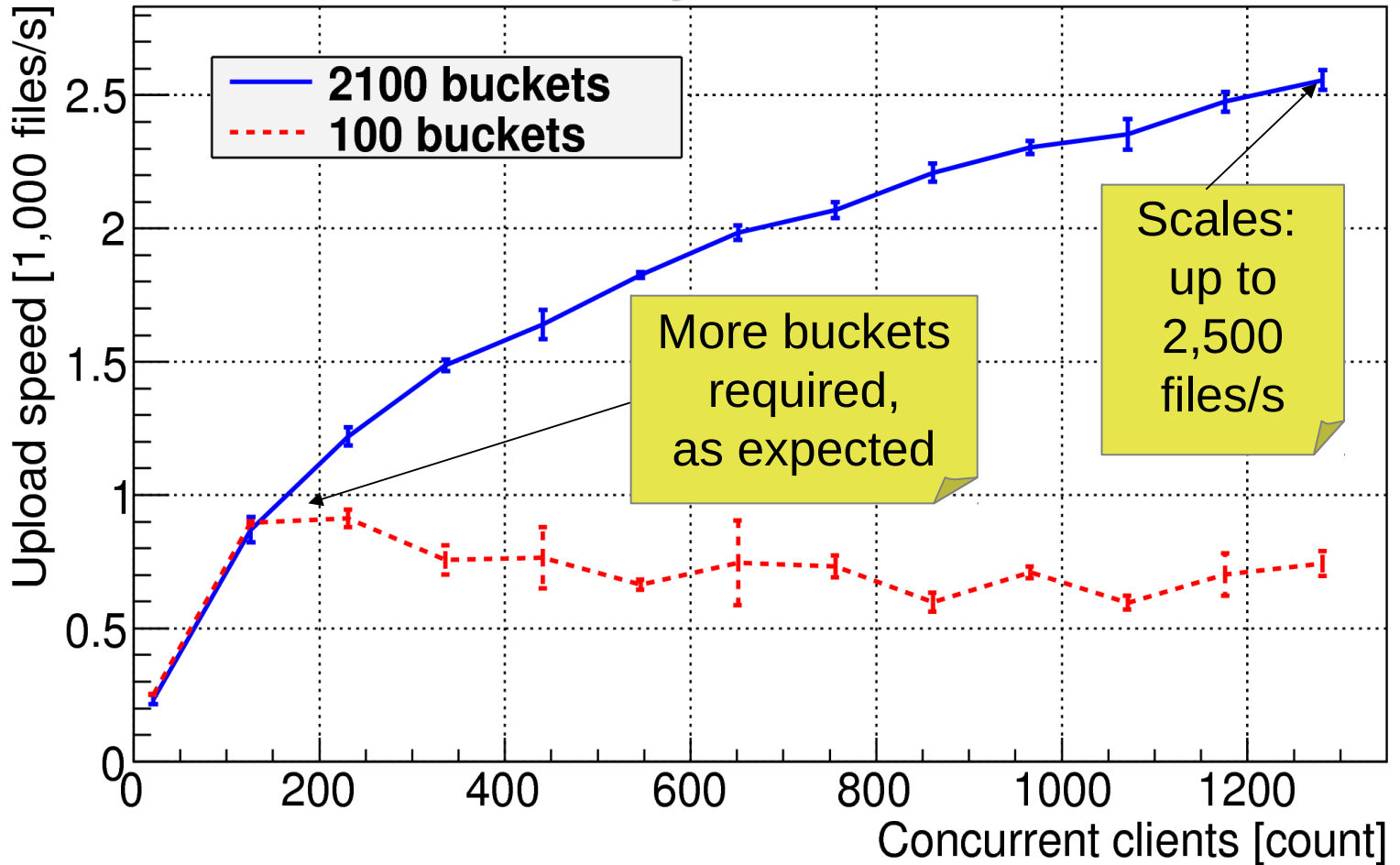
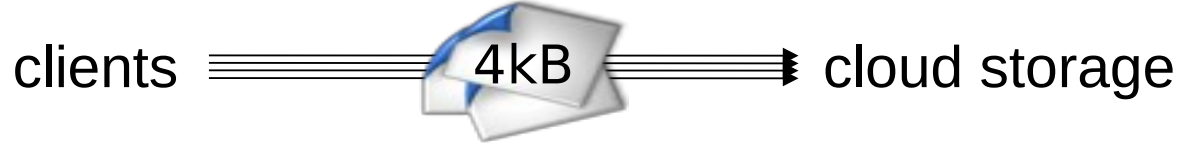


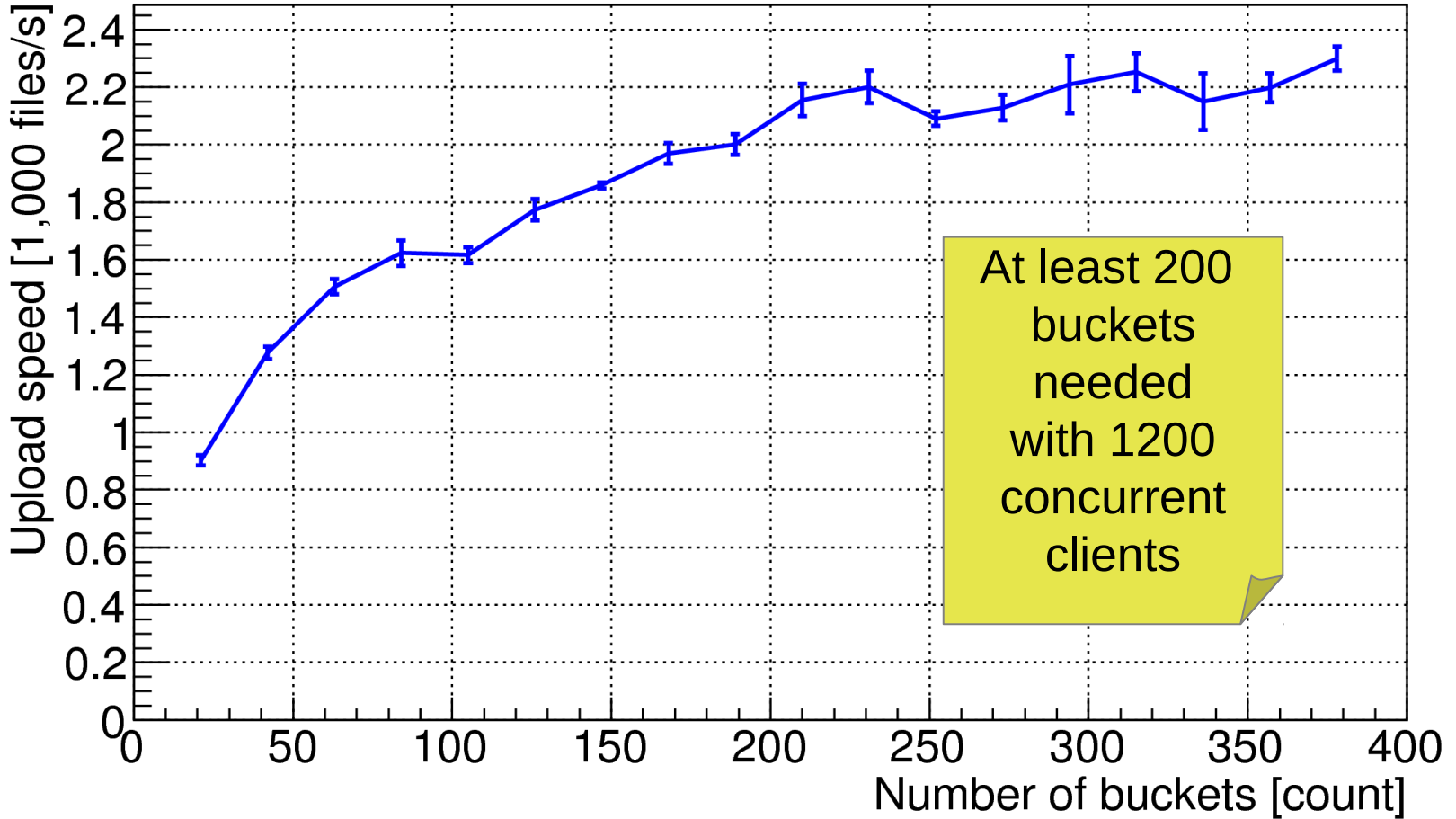
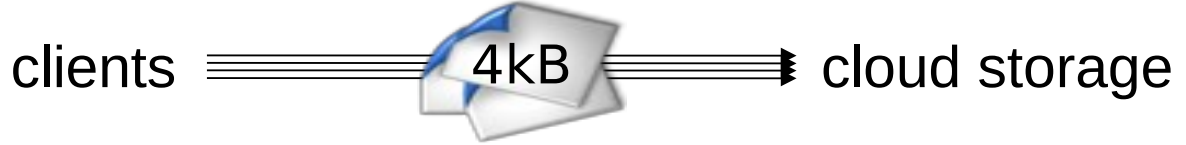




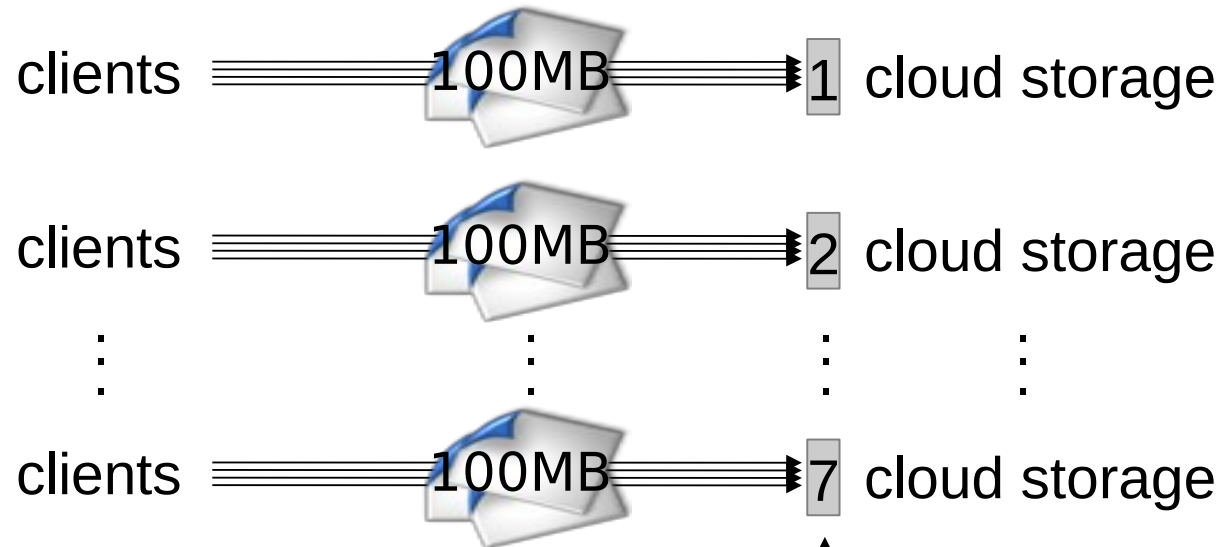




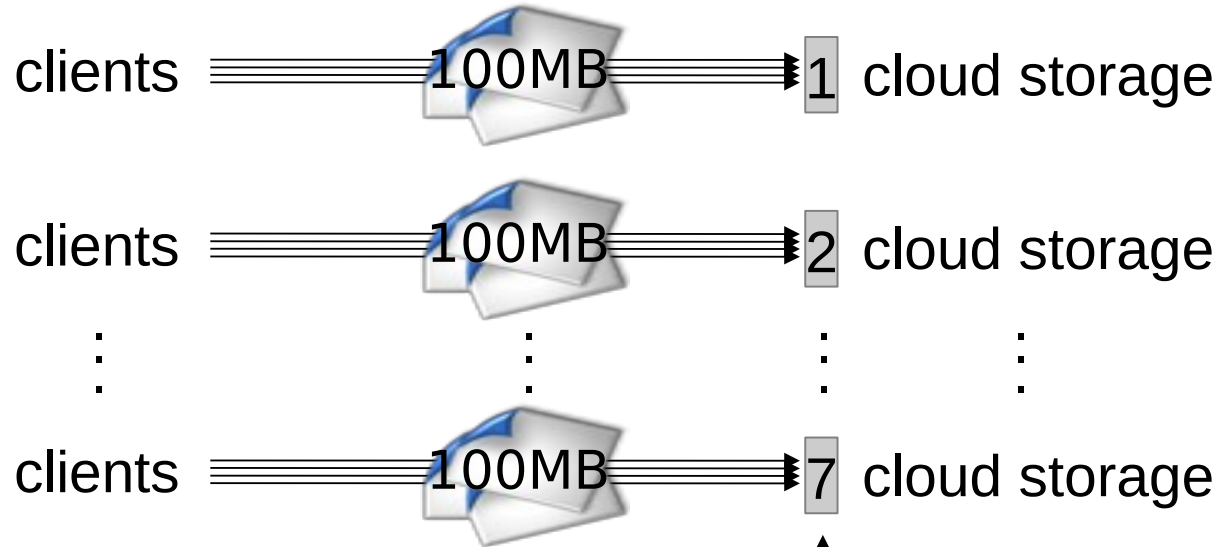




At least 200 buckets needed with 1200 concurrent clients



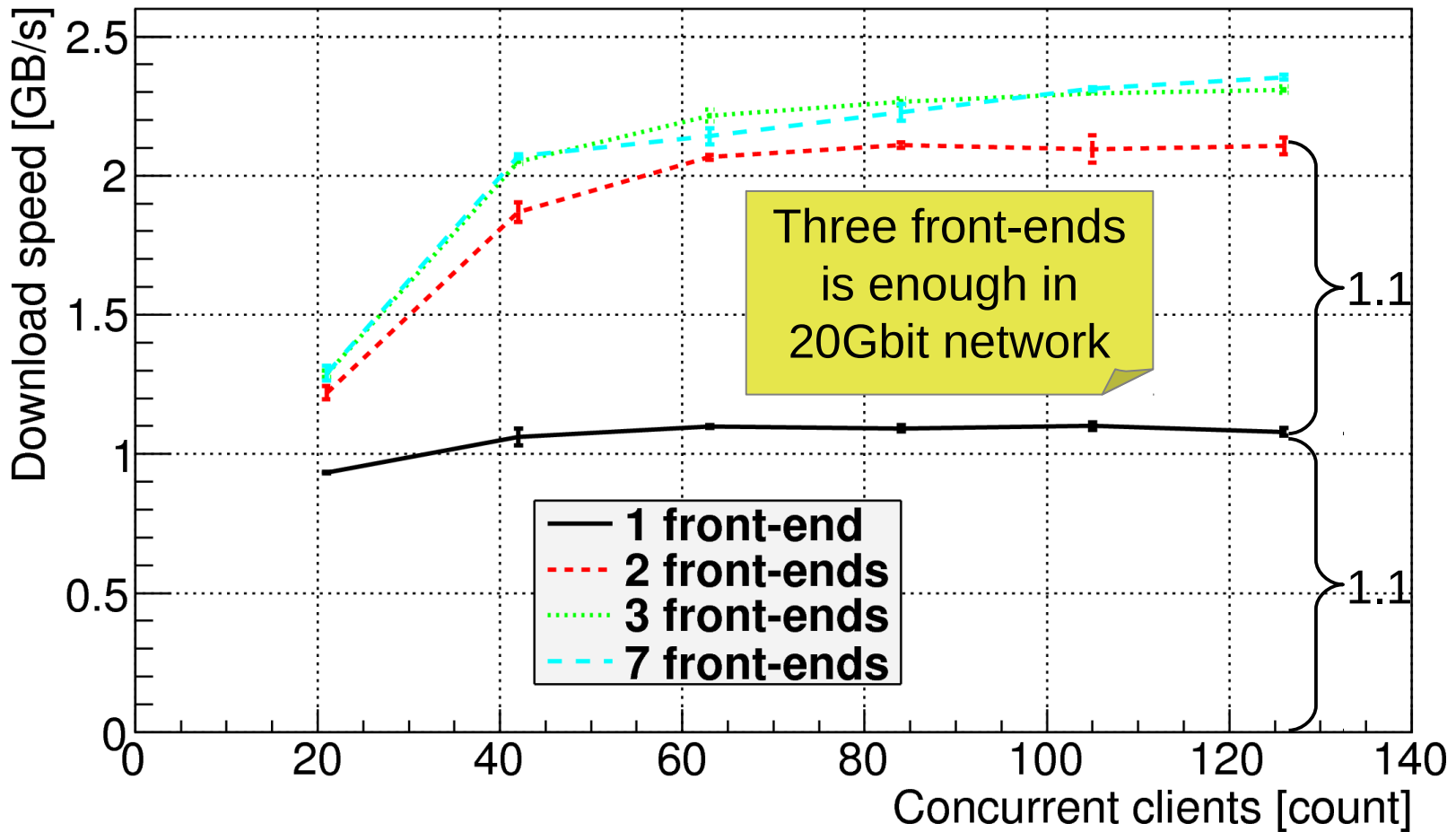
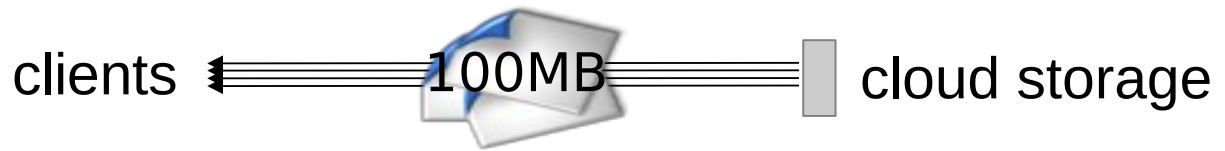
We use
different
numbers
of front-ends:
from 1 to 7

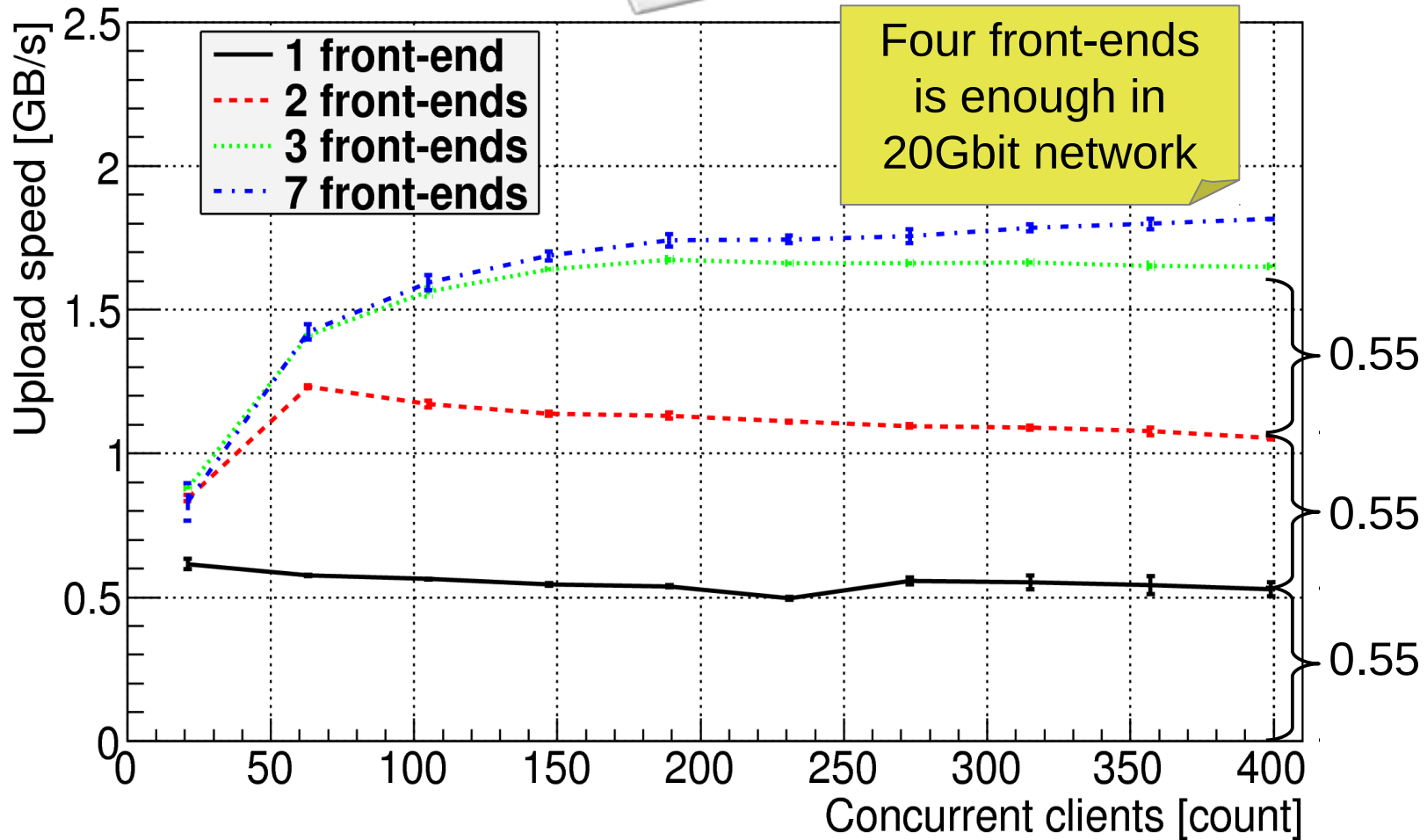


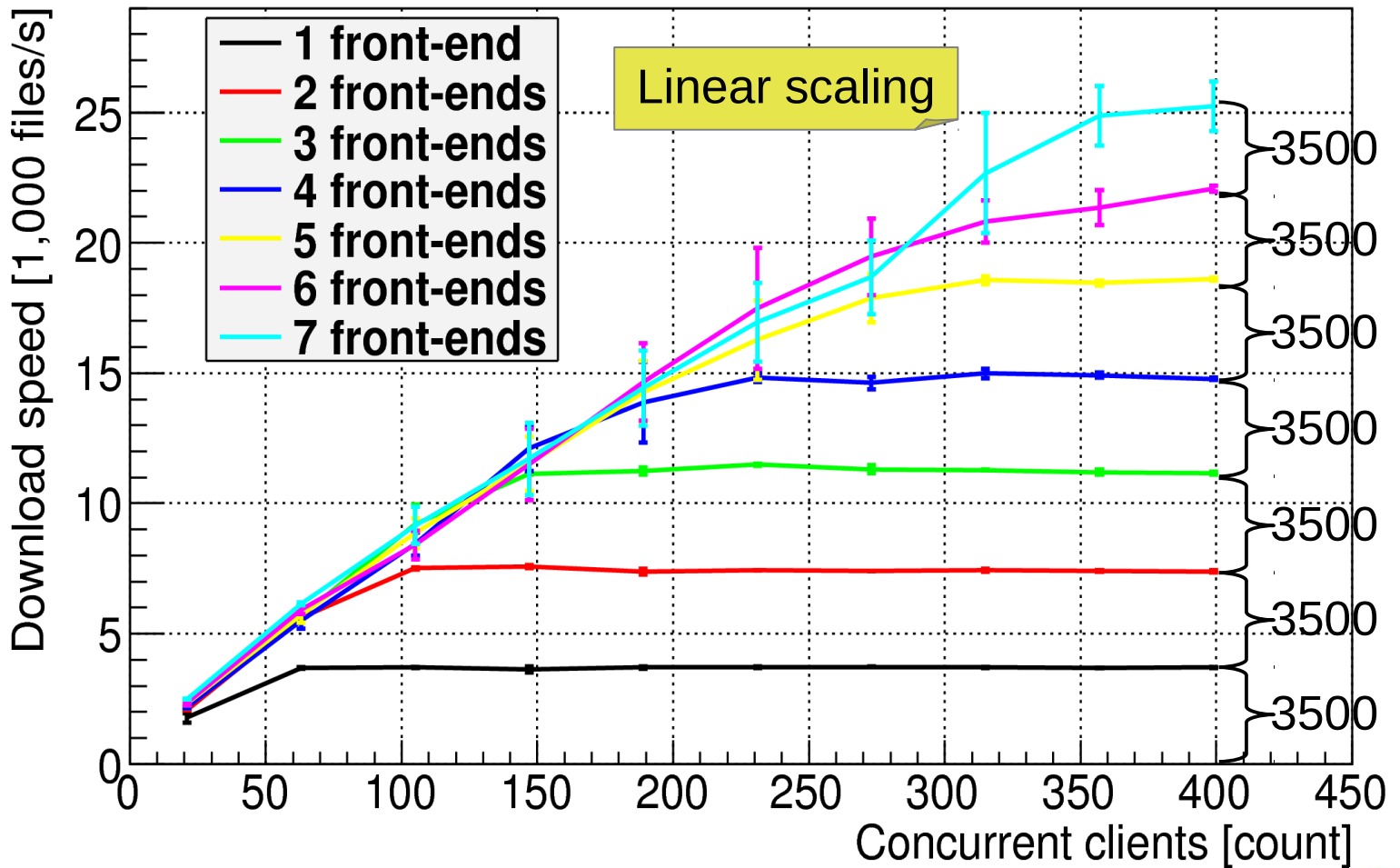
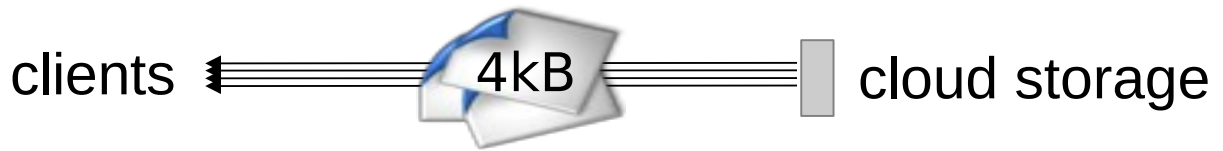
Small (4kB)
& big (100MB)
files

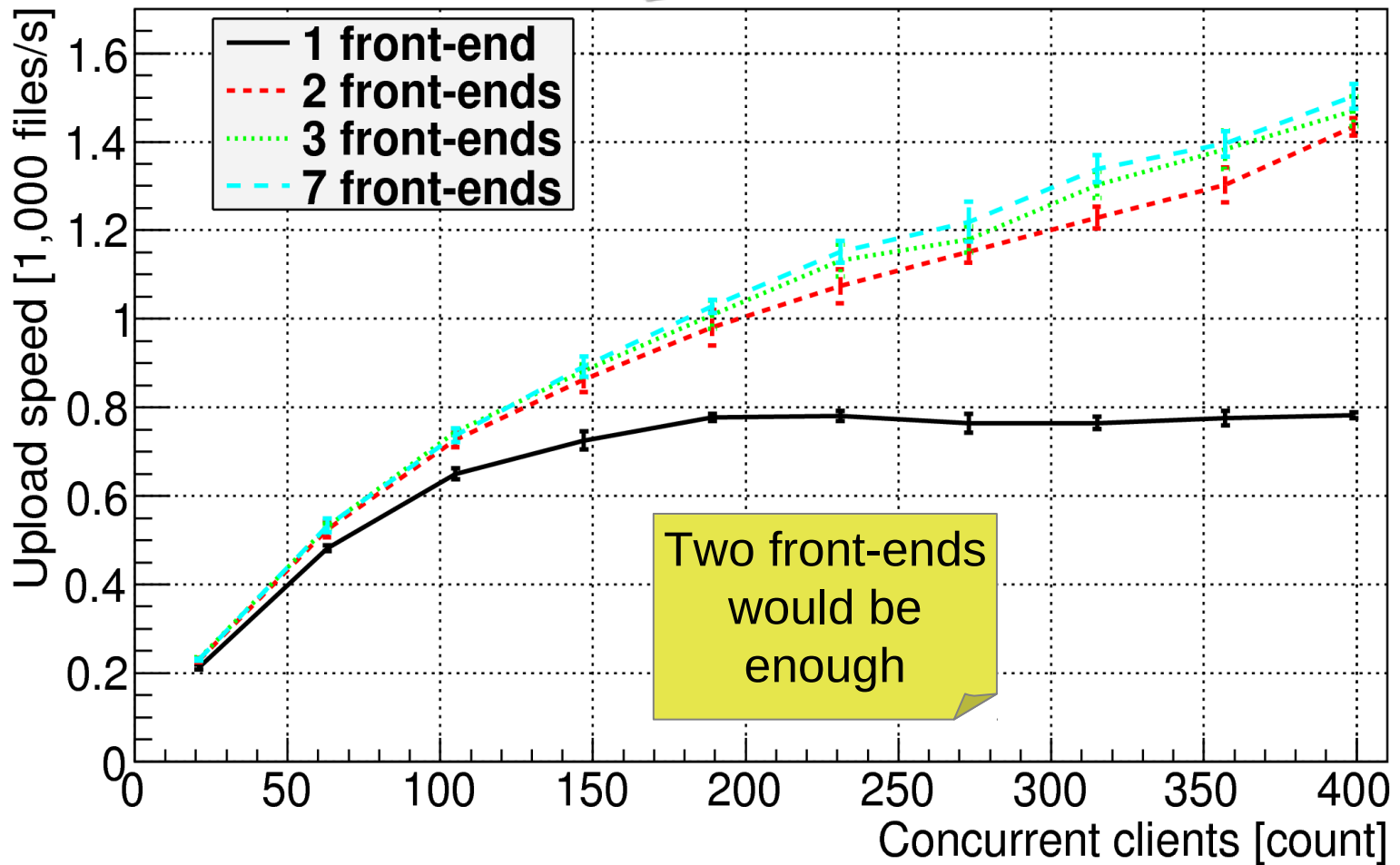
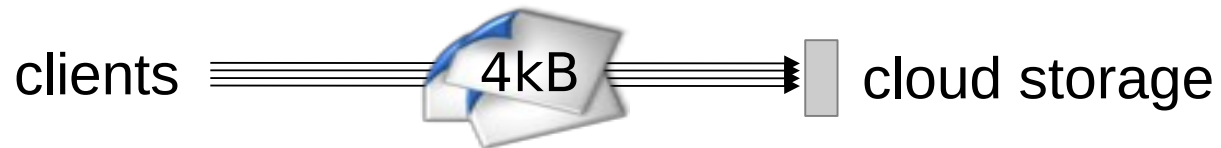
Uploads &
downloads

We use
different
numbers
of front-ends:
from 1 to 7





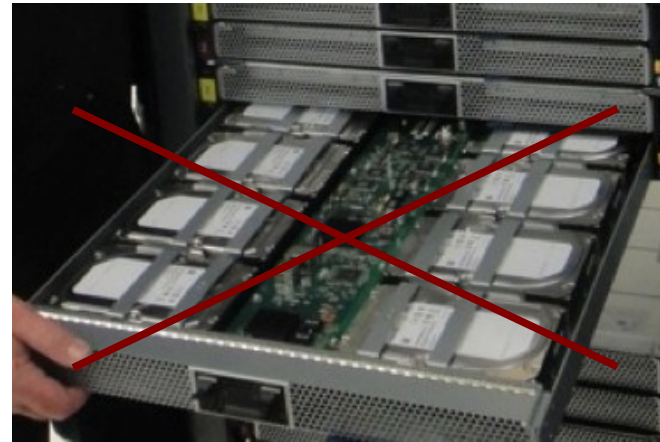
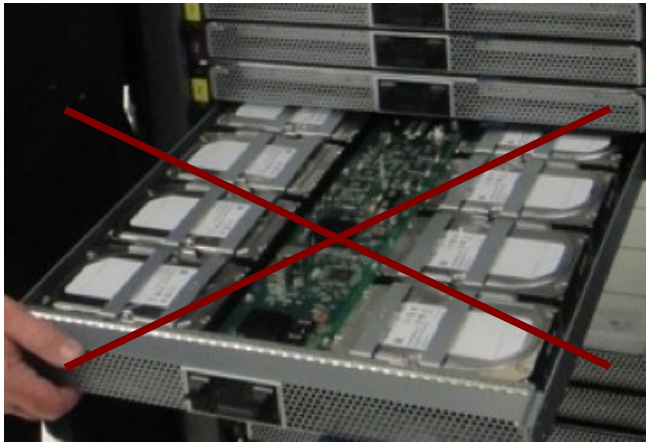




- Metadata (4kB) performance
 - 2,500 files/second upload
 - 25,000 files/second download
- Throughput (100MB) performance
 - 20Gbit network fully utilized
- Front-end scalability
 - Each front-end can download 3500 files/s
 - Each front-end can upload 550 MB/s

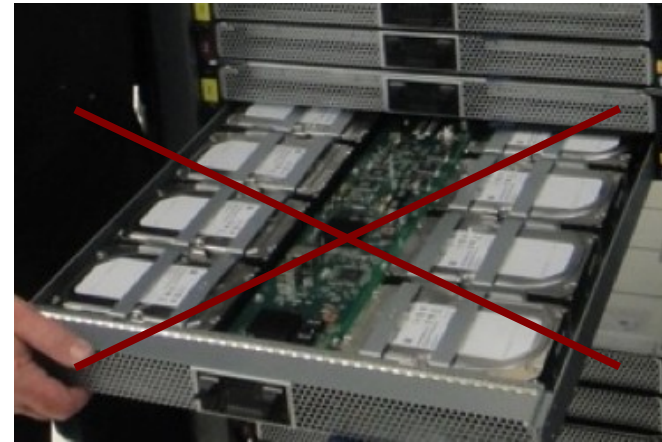
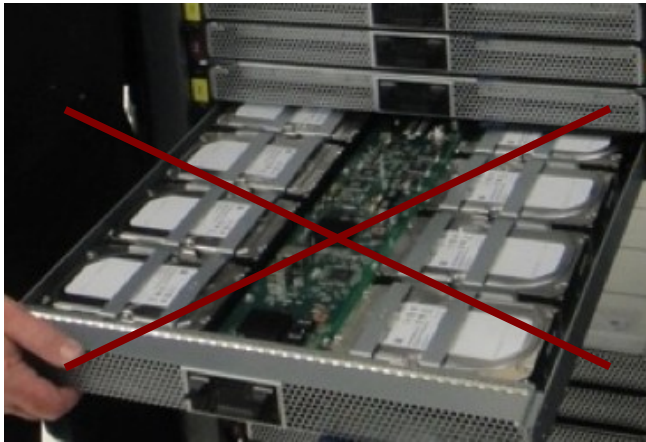


Two blades are powered off:



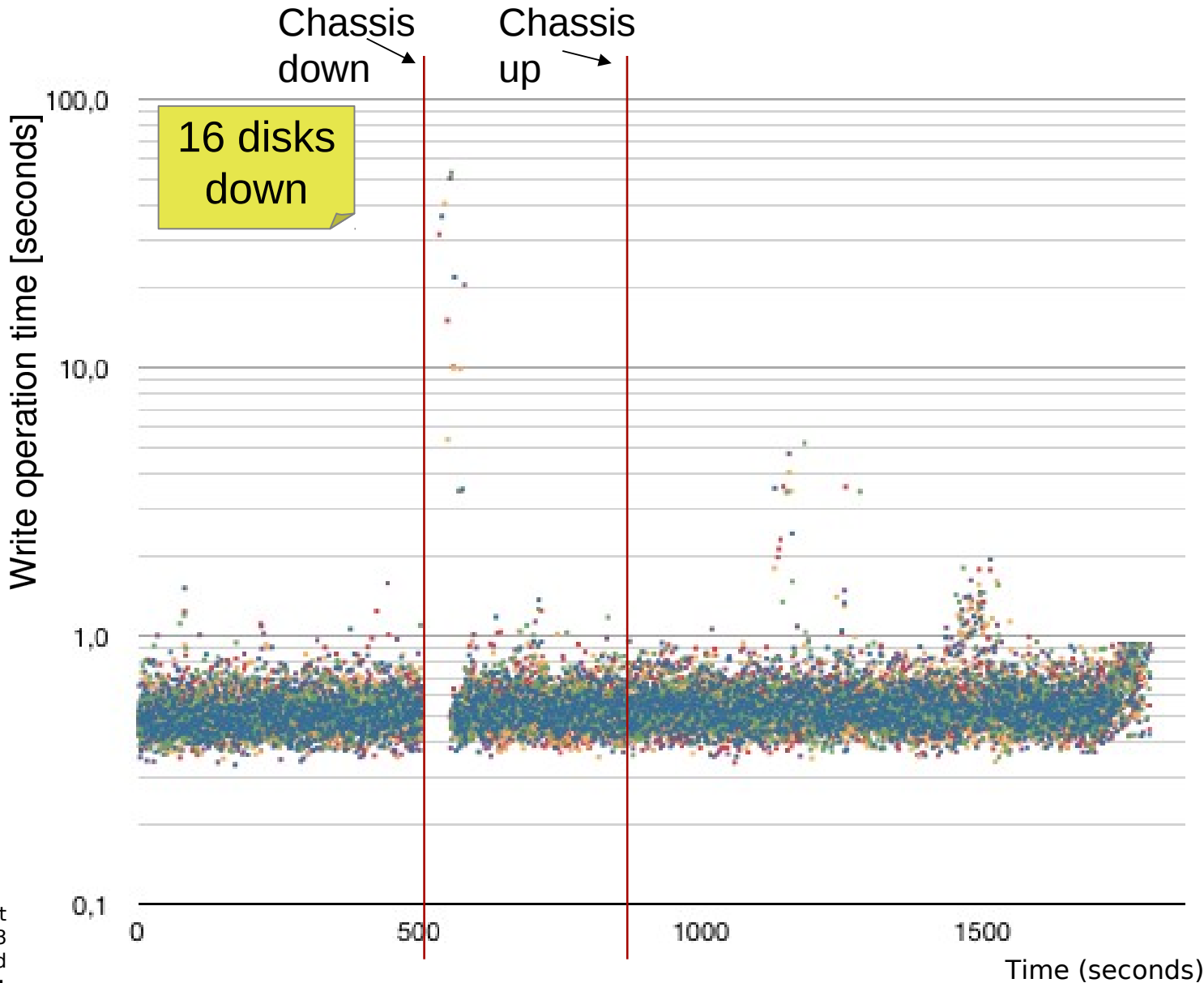
16 disks
down

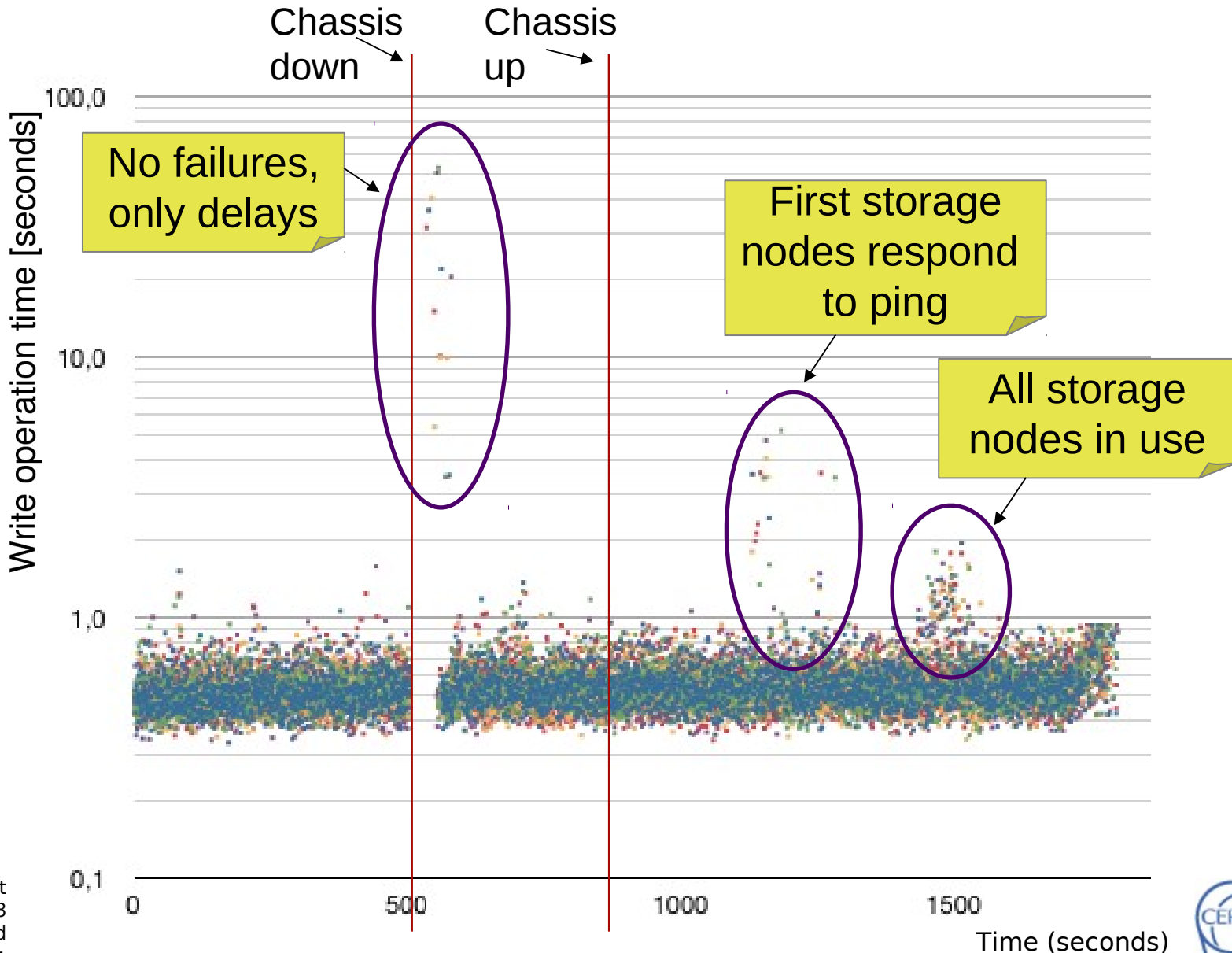
Two blades are powered off:



16 disks
down

Uploads and
downloads
continue
normally?





- What is CVMFS (CernVM File System)
 - Read only cached file system to deliver software
 - Widely used in WLCG (Worldwide LHC Computing Grid)
 - Mounted by users and files are downloaded on demand



- What is CVMFS (CernVM File System)
 - Read only cached file system to deliver software
 - Widely used in WLCG (Worldwide LHC Computing Grid)
 - Mounted by users and files are downloaded on demand



- CVMFS challenges
 - Publishing new software should be fast (upload tens of thousands of files)
 - Files should be accessed with HTTP protocol

- Implementation



- Files are uploaded to multiple buckets in the cloud storage
- Files are downloaded with unified name space
~~<http://cloud.cern.ch/bucket-42/file001.bin>~~
<http://cloud.cern.ch/file001.bin>



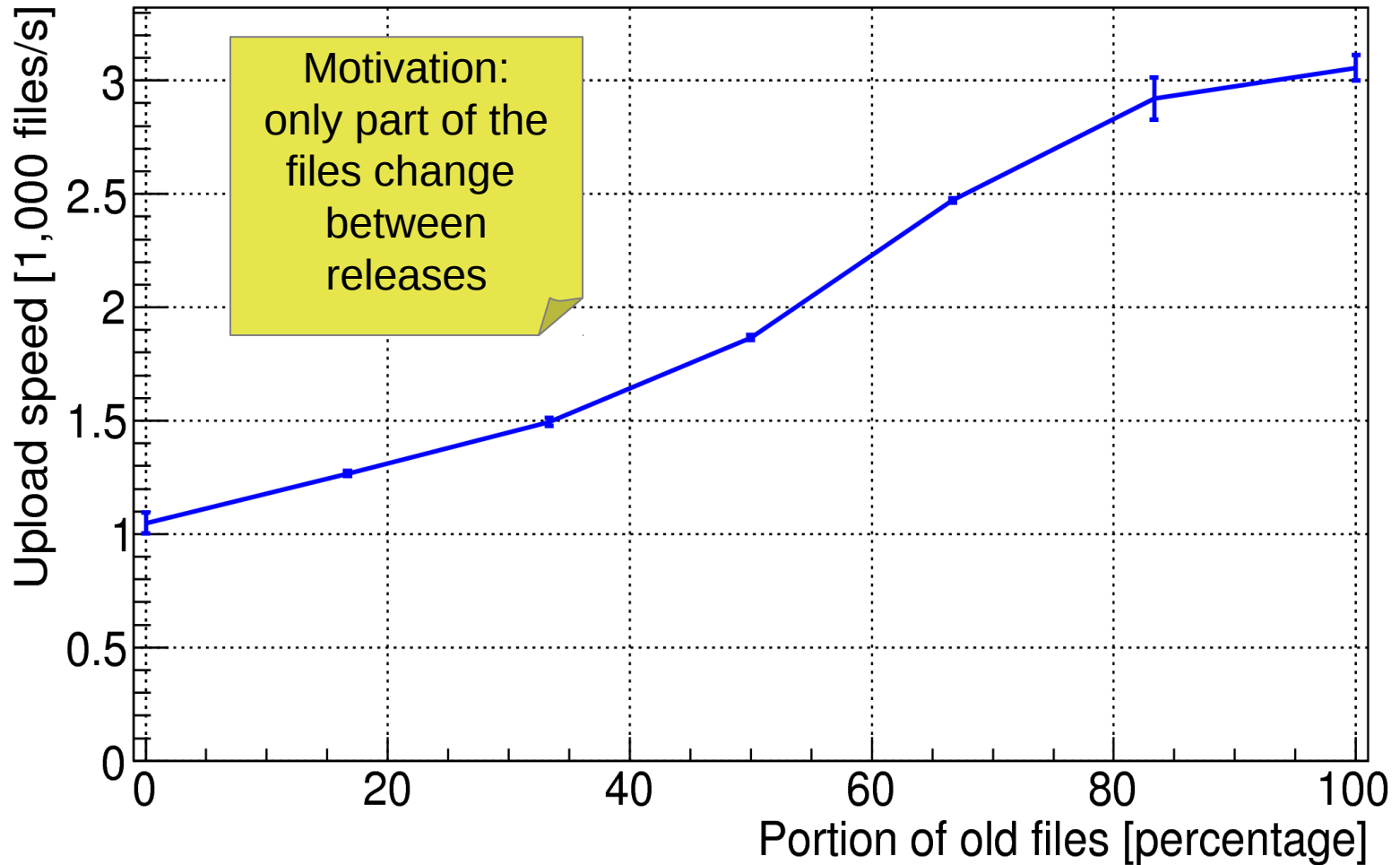
- Implementation

- Files are uploaded to multiple buckets in the cloud storage
- Files are downloaded with unified name space
~~<http://cloud.cern.ch/bucket-42/file001.bin>~~
<http://cloud.cern.ch/file001.bin>

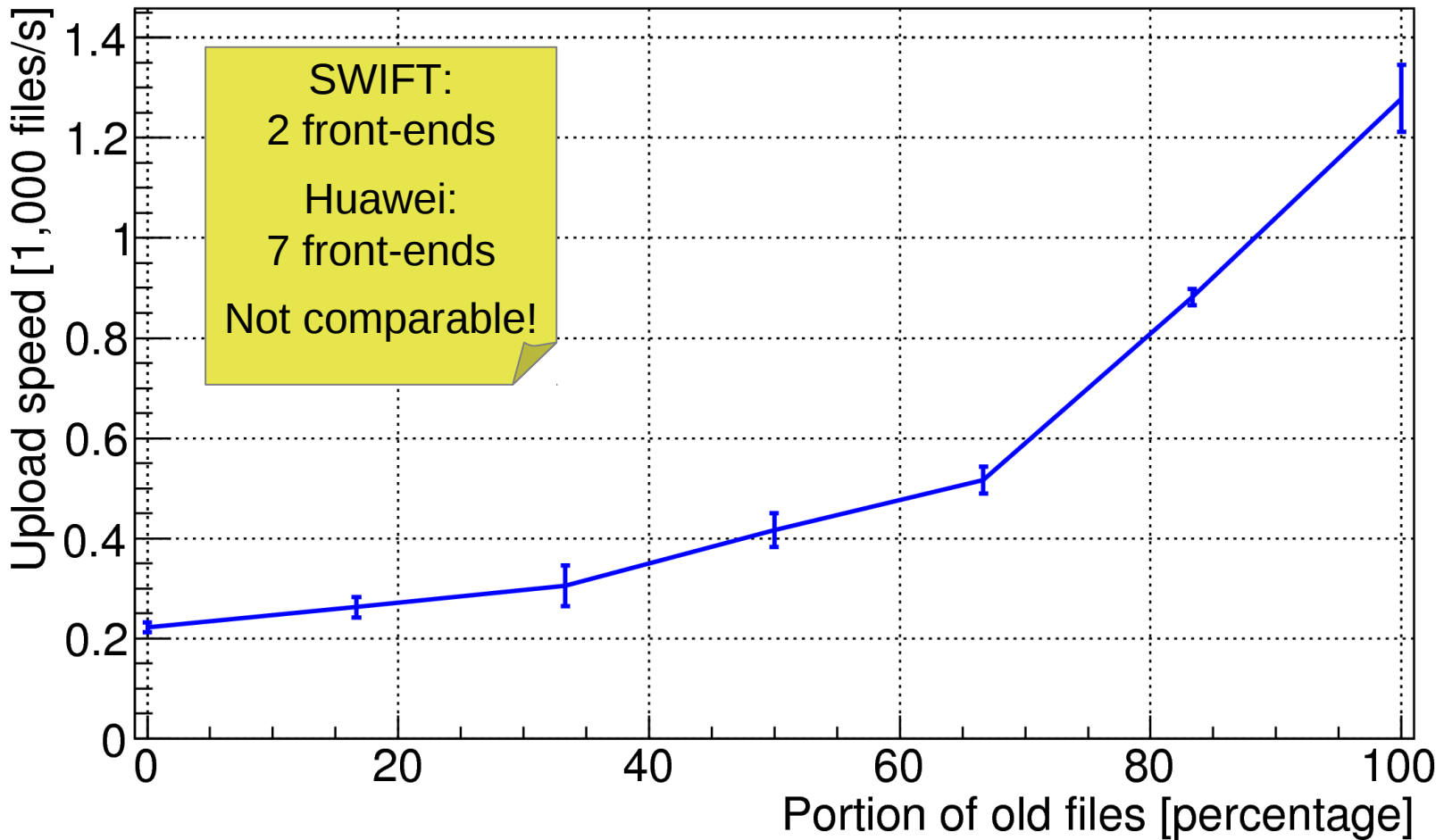
- Result

- Full publish procedure tested to work using 30,000 small files
- Upload speed 1200 files/second (with 240 threads)

Uploading 30,000 files (of average size 10kB)
to Huawei UDS



Uploading 30,000 files (of average size 10kB)
to OpenStack SWIFT



- Raw performance
 - Upload and download **scalability** demonstrated
 - Additional front-end nodes increased linearly the performance
- Fault tolerance: powering off a chassis
 - **Transparent** disk failure recovery demonstrated
- File system with cloud storage back-end
 - Full **publishing procedure** tested
 - Uploading of **only new** files feature tested

- Short term
 - Benchmark CVMFS with real release data
 - Test ROOT's new S3 plugin performance
- Long term
 - Second petabyte system with enterprise disks expected to arrive soon
 - Upgrade old Huawei cloud storage software version
 - Replication tests between two cloud storages
 - Prove total cost of ownership (TCO) gains of the system as part of a production service

- Short term
 - Benchmark CVMFS with real release data
 - Test ROOT's new S3 plugin performance
- Long term
 - Second petabyte system with enterprise disks expected to arrive soon
 - Upgrade old Huawei cloud storage software version
 - Replication tests between two cloud storages
 - Prove total cost of ownership (TCO) gains of the system as part of a production service

Thank you!

seppo.heikkila@cern.ch