# Huawei Cloud Storage

Seppo S. Heikkila
Maria Arsuaga Rios
CERN IT

Openlab Major Review Meeting

13th of February 2014

CERN, Geneva

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

HUAWEI

CERN openlab

CERN IT Department

## 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

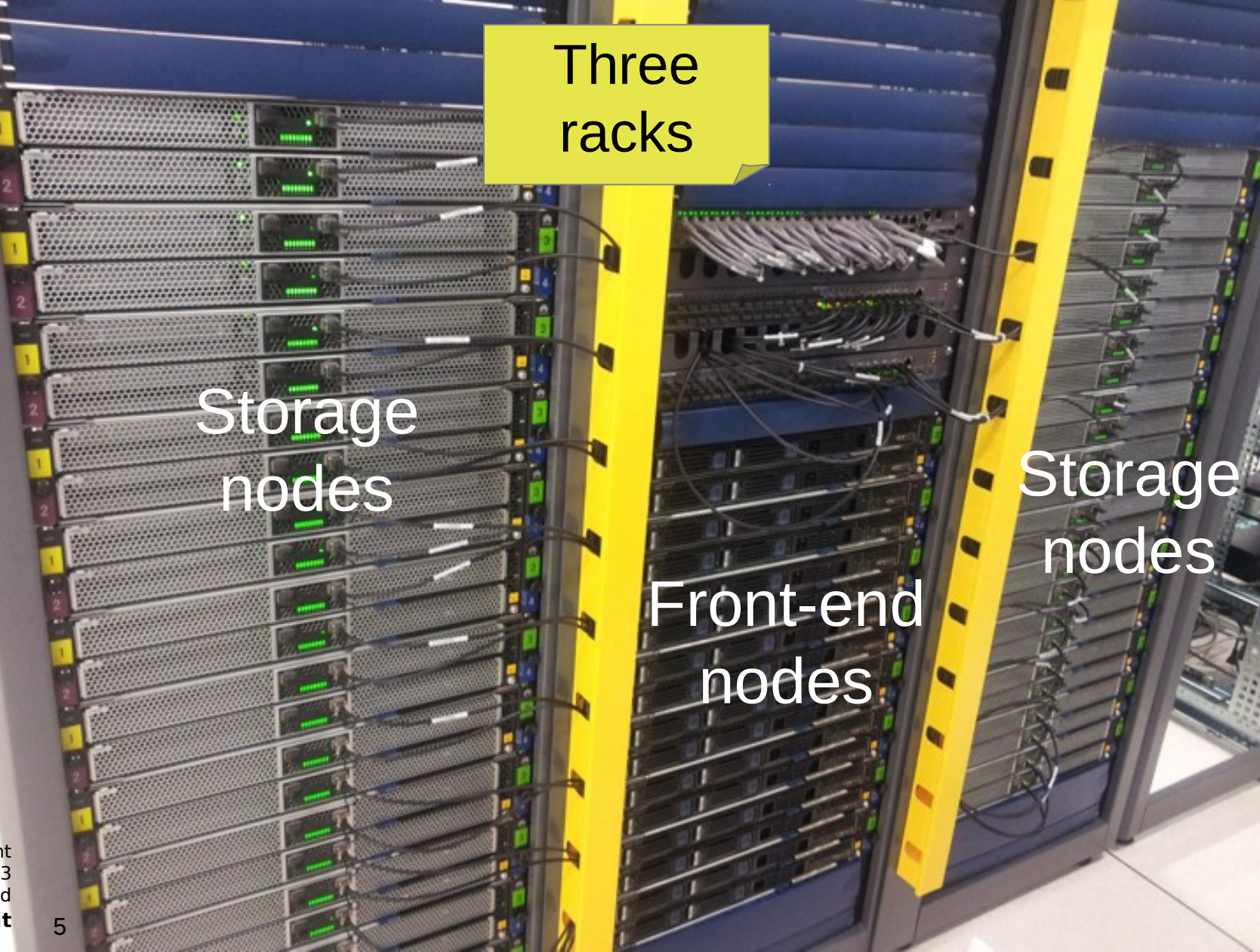– Evaluate scalability and fault-tolerance
– Test with real applications

CERN**IT** Department

# 2 years of Huawei...

Major Review    Major Review    Major Review    Major Review    Major Review

Minor Review    Board of Sponsors    Minor Review    Board of Sponsors    Minor Review

Maitane      Seppo      Maria

01/2012        01/2013        01/2014

Project starts    First tests    Upgrade of the system    Stress testing    File-system integration

Commissioning of the system    Failure recovery testing    Full-scale stress testing

Location: CERN Computer Center

"We are now here"

Cloud storage

# Huawei cloud storage setup

Three racks

Storage nodes

Front-end nodes

Storage nodes

# Huawei cloud storage setup

384 disks → 768 TB

Storage nodes

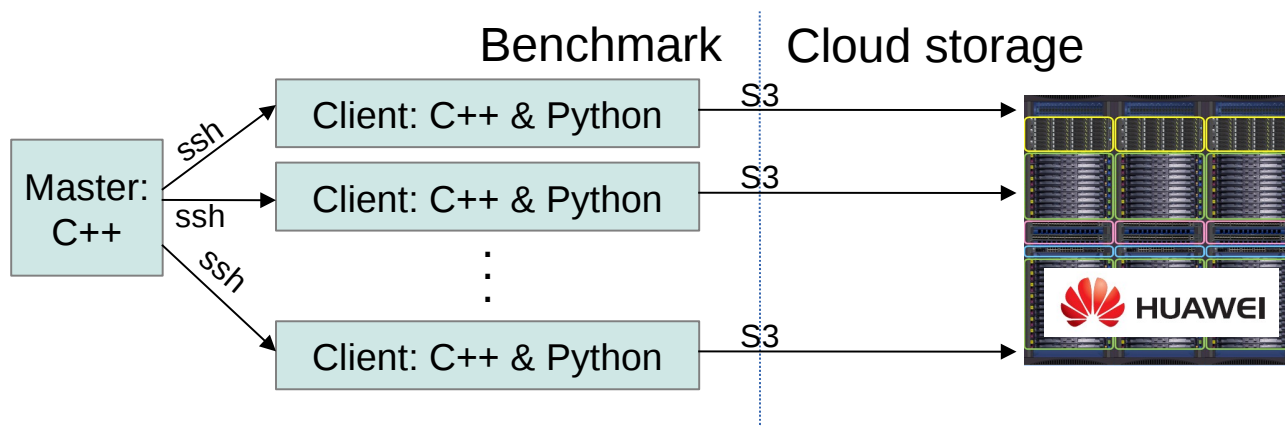S3 compatible → Buckets divide the namespace

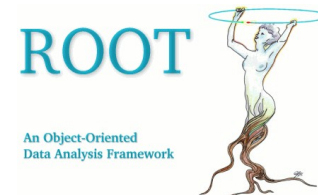Front-end nodes

Storage nodes

One chassis has two blades

Each blade has
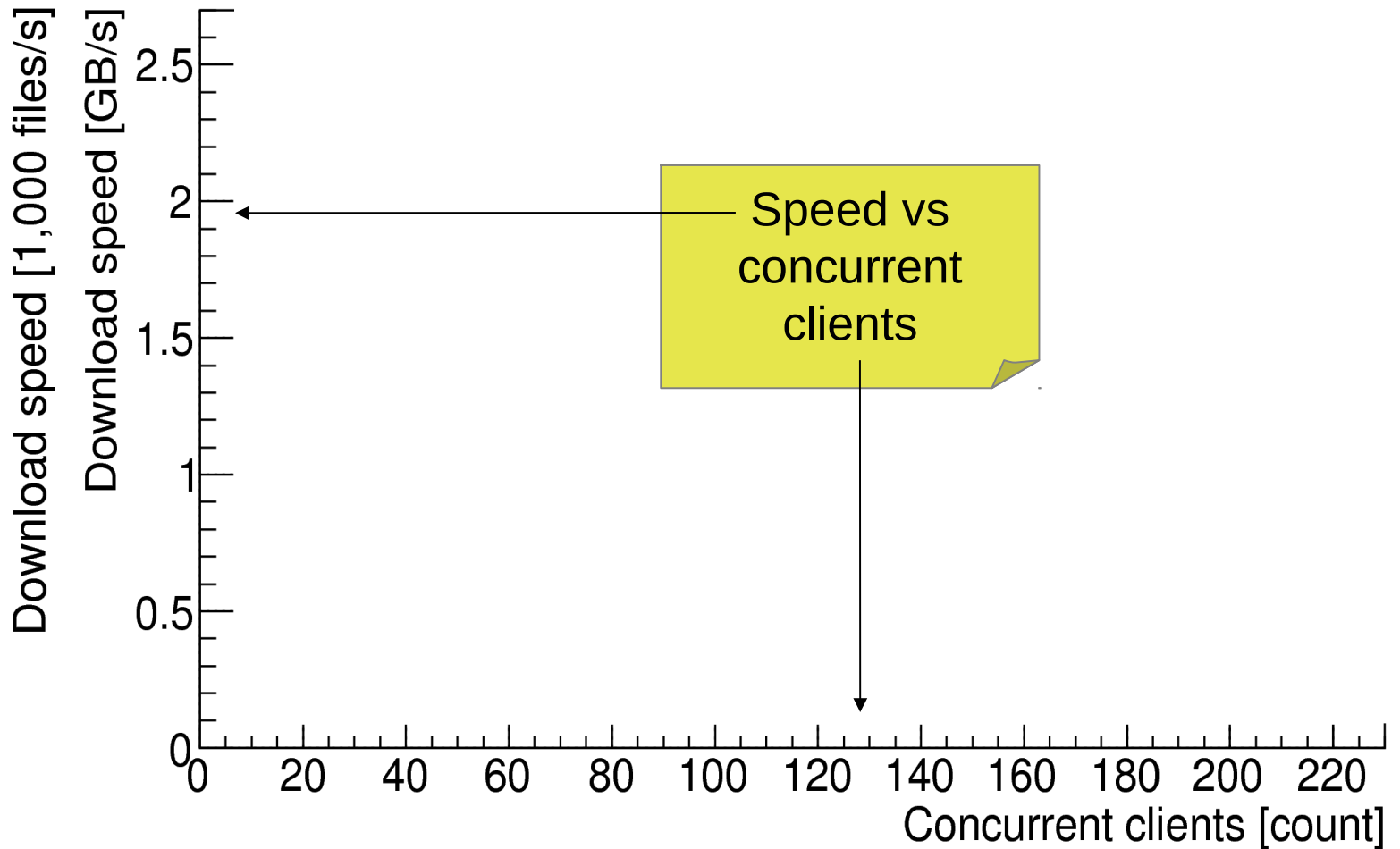eight storage nodes

# Distributed C++ benchmark

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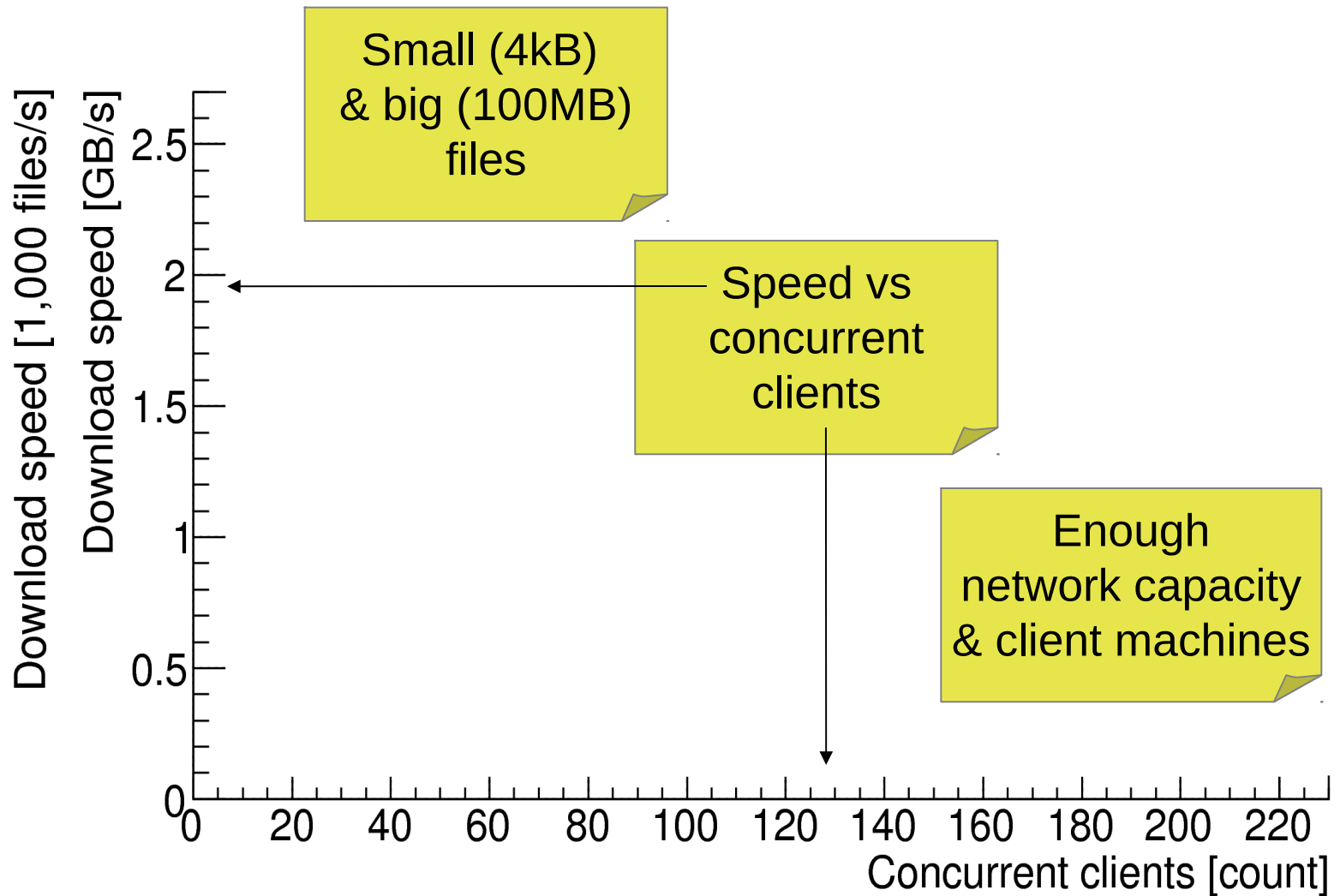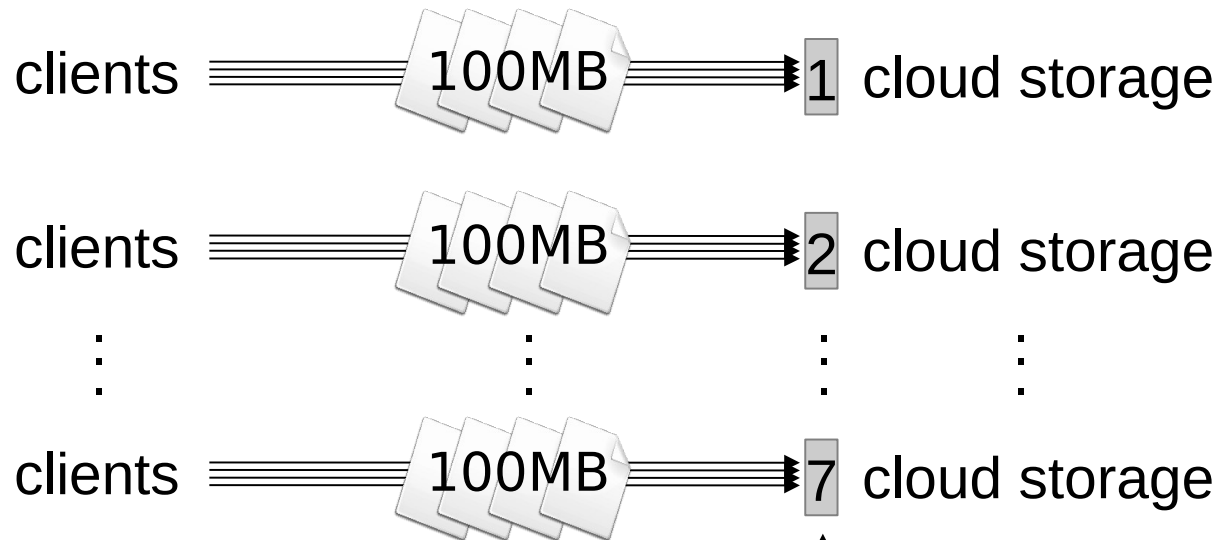  - Operation time, read/write speed, CPU/memory utilisation

"Zotes Resines M, Heikkila SS, Duelmann D, Adde G, Toebbicke R, Hughes J and Wang L. <u>Evaluation of the Huawei UDS cloud storage system for CERN specific data</u>, International Conference on Computing in High Energy and Nuclear Physics (CHEP) 2013, Amsterdam, The Netherlands, 14 October 2013"

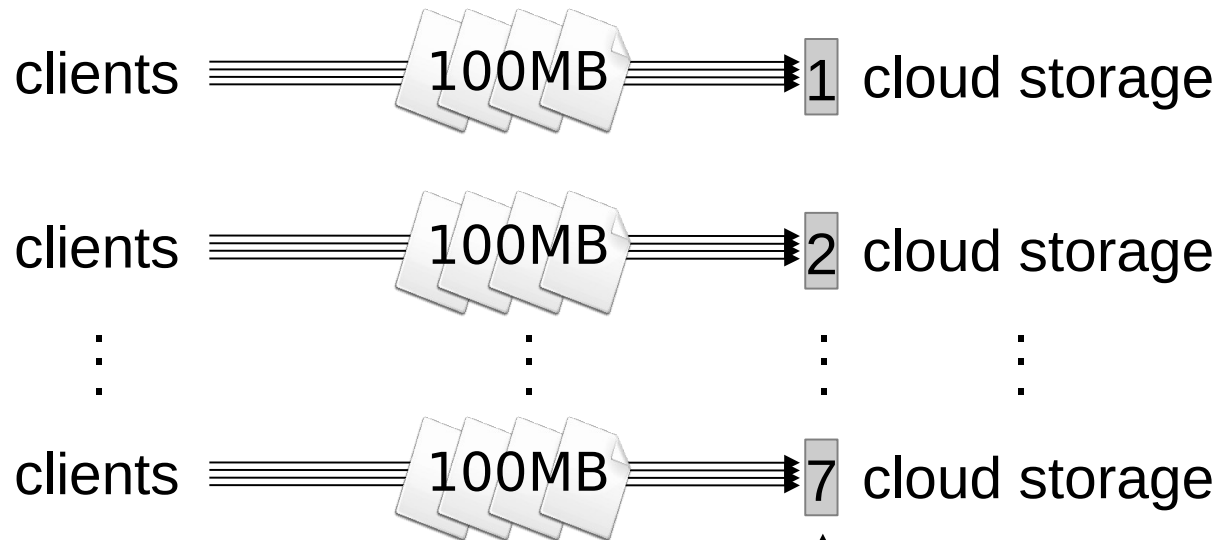Presented and submitted in October!
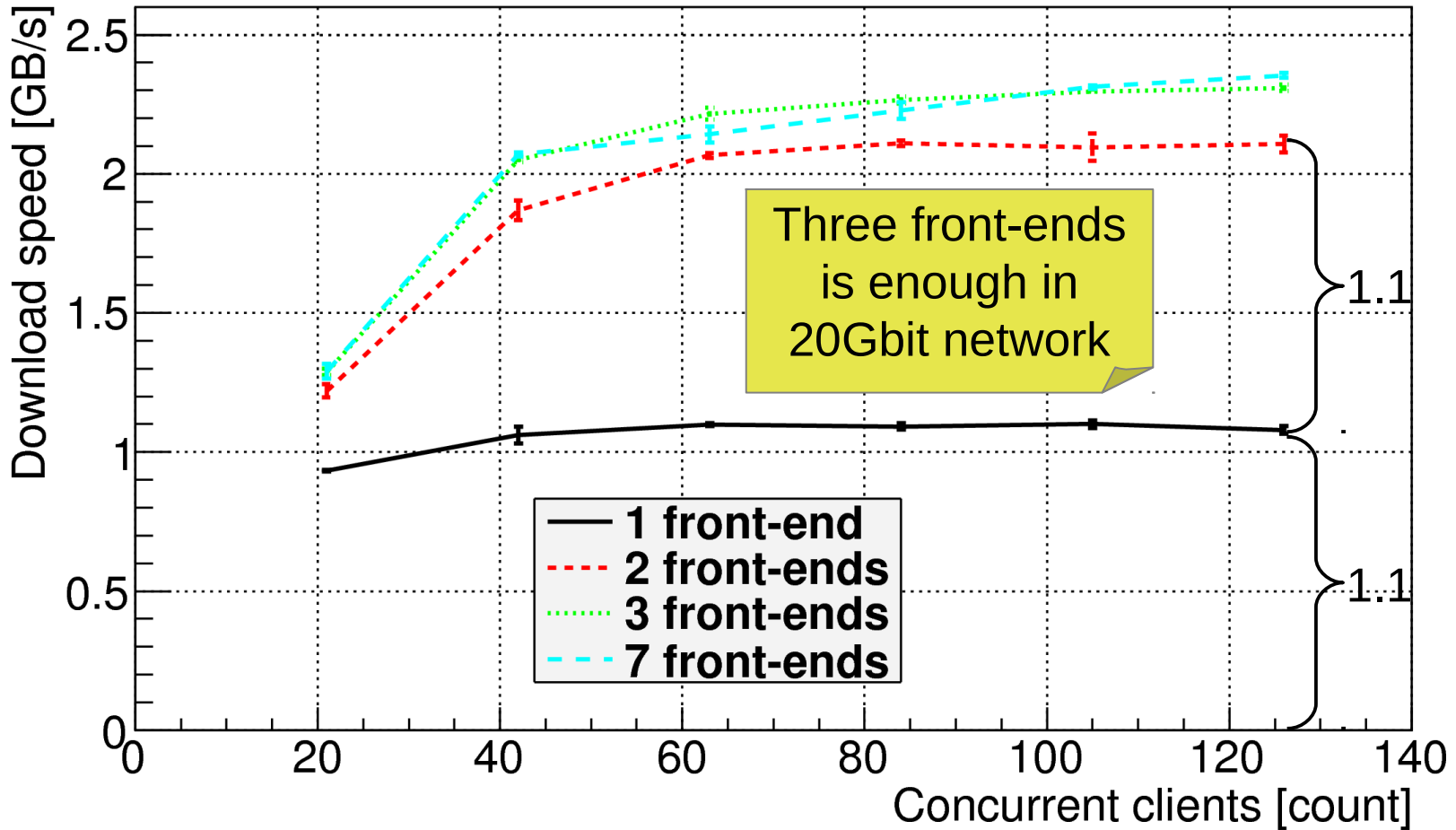
Revised version submitted in December!

# Benchmarking scalability

Speed vs concurrent clients

# Benchmarking scalability

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

Speed vs concurrent clients

Enough network capacity & client machines

y-axis: Download speed [1,000 files/s]    Download speed [GB/s]

2.5  2  1.5  1  0.5  0

x-axis: Concurrent clients [count]

0  20  40  60  80  100  120  140  160  180  200  220

# Front-end scalability

clients ===== 100MB =====▶ 1 cloud storage

clients ===== 100MB =====▶ 2 cloud storage

⋮              ⋮              ⋮        ⋮

clients ===== 100MB =====▶ 7 cloud storage

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

# Front-end scalability

clients ═══ 100MB ═══► 1 cloud storage

clients ═══ 100MB ═══► 2 cloud storage

⋮ ⋮ ⋮ ⋮

clients ═══ 100MB ═══► 7 cloud storage
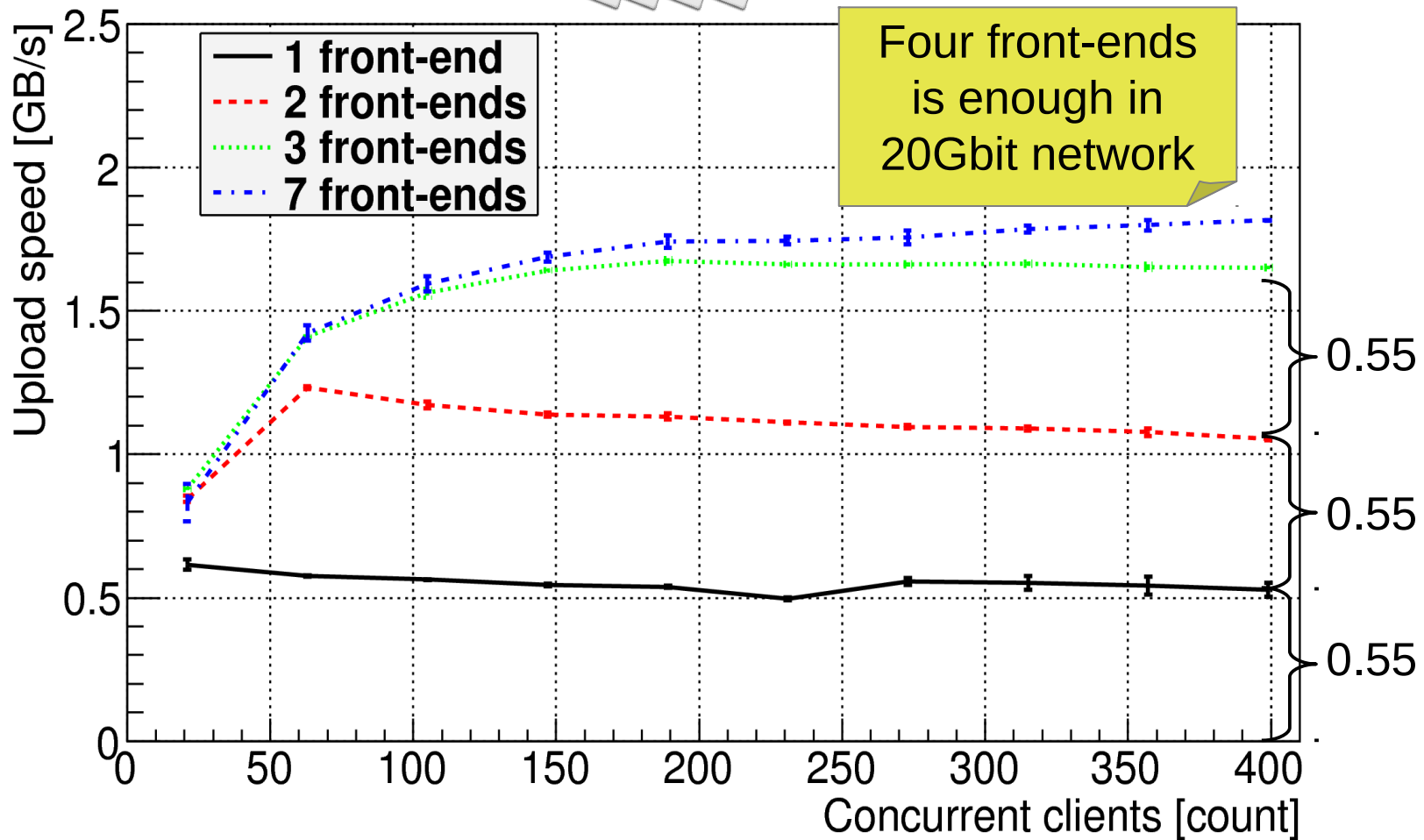
Small (4kB)
& big (100MB)
files

Uploads &
downloads

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

DSS

clients ⟵ 100MB ⟶ cloud storage



Three front-ends is enough in 20Gbit network

1.1

1.1

**Legend:**
- 1 front-end
- 2 front-ends
- 3 front-ends
- 7 front-ends

Y-axis: Download speed [GB/s]
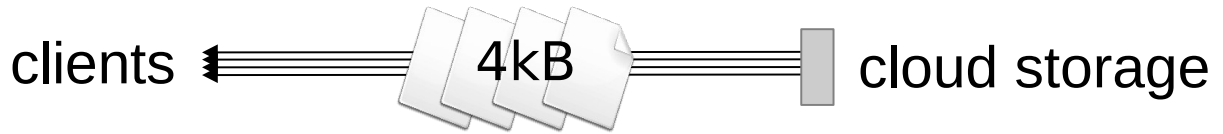X-axis: Concurrent clients [count]

# Data upload scalability

clients ⟹ 100MB ⟹ cloud storage

Four front-ends is enough in 20Gbit network

# Meta-data download scalability

clients ← 4kB → cloud storage

Linear scaling

# Meta-data upload scalability



clients ⟹ 4kB ⟹ cloud storage

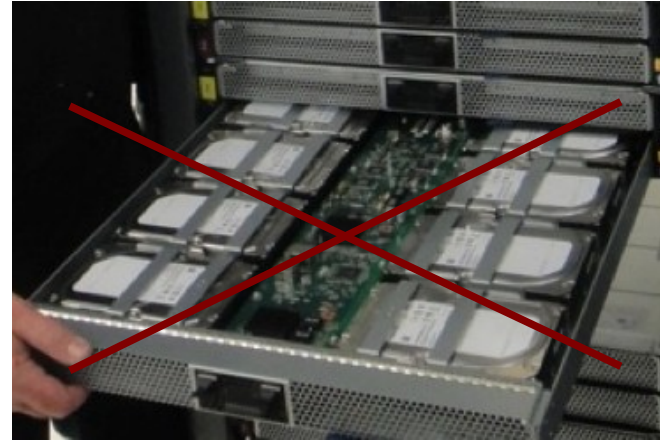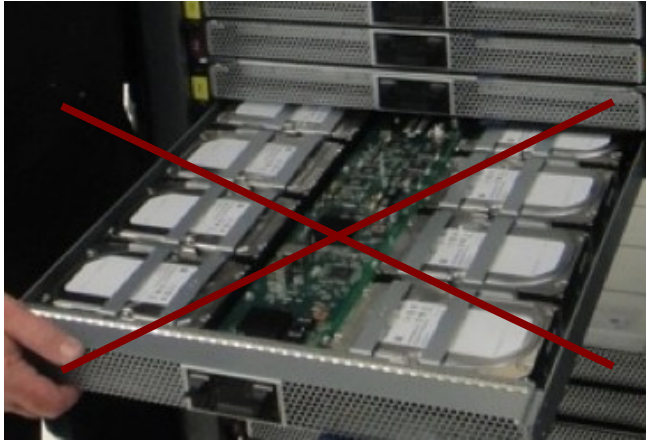Two front-ends would be enough

- Metadata (4kB) performance
  - 2,500 files/second upload
  - 25,000 files/second download

- Throughput (100MB) performance
  - 20Gbit network fully utilized

- Front-end scalability
  - One front-end can download 3500 files/s
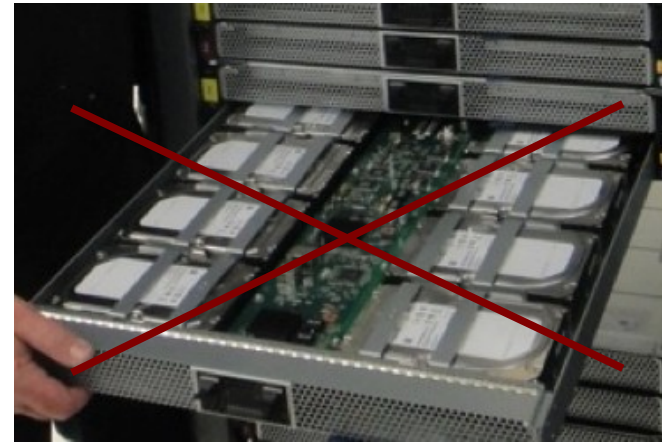  - Each front-end can upload 550 MB/s

# Chassis power off recovery

Two blades are powered off:




16 disks down

# Chassis power off recovery

## Two blades are powered off:



16 disks
down

Uploads and
downloads
continue
normally?

CERN IT Department



16 disks down

Chassis down

Chassis up

Write operation time [seconds]

Time (seconds)

# CERN IT Department

- ## 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 introduction

- **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

CernVM
File system

- 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
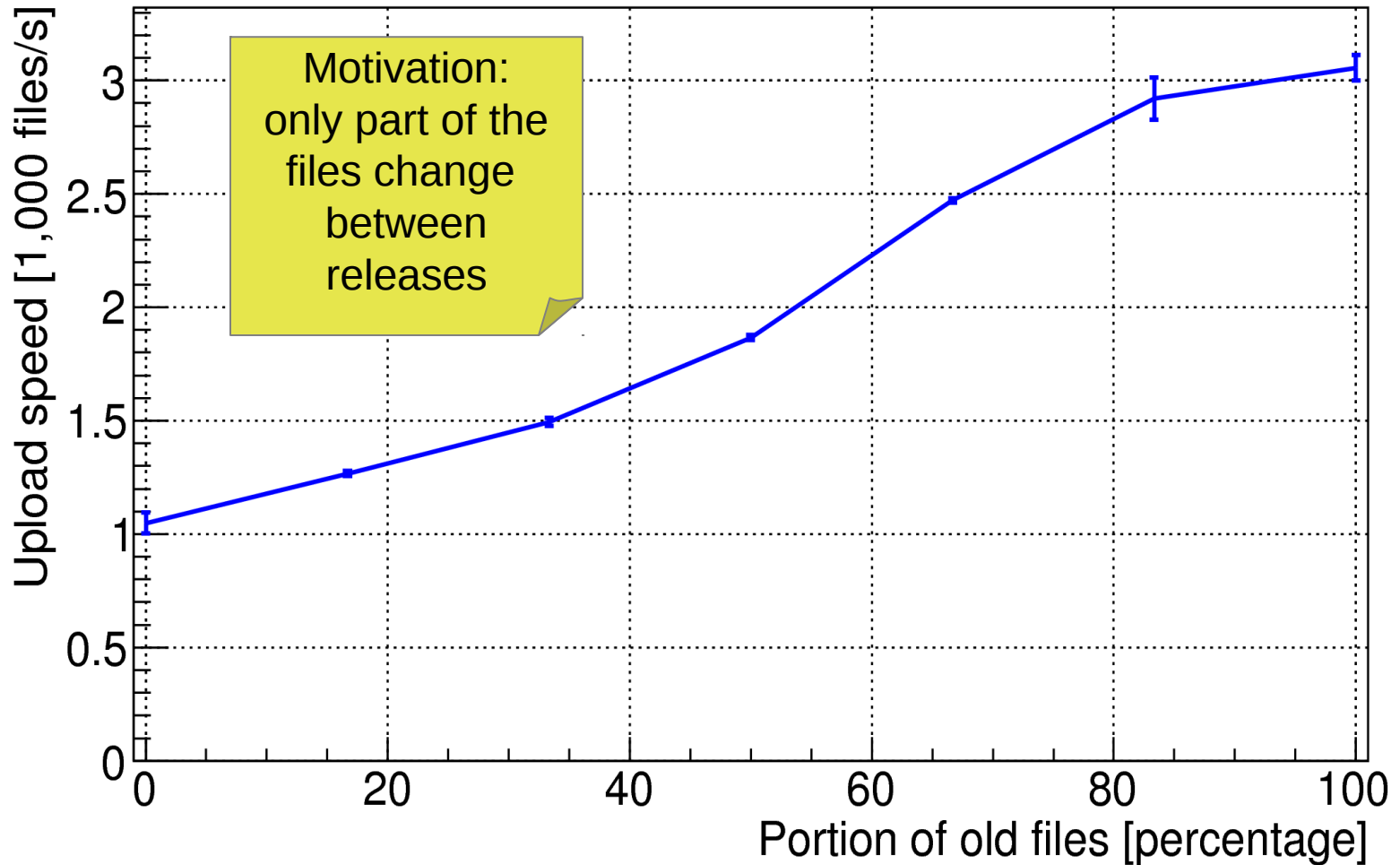    - ~~http://cloud.cern.ch/bucket-42/file001.bin~~
    - http://cloud.cern.ch/file001.bin
- # Result
  - Full publish procedure tested to work using 30,000 small files
  - Upload speed 1200 files/second (with 240 threads)

DSS

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

Motivation: only part of the files change between releases

# CHEP paper summary

- ## 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
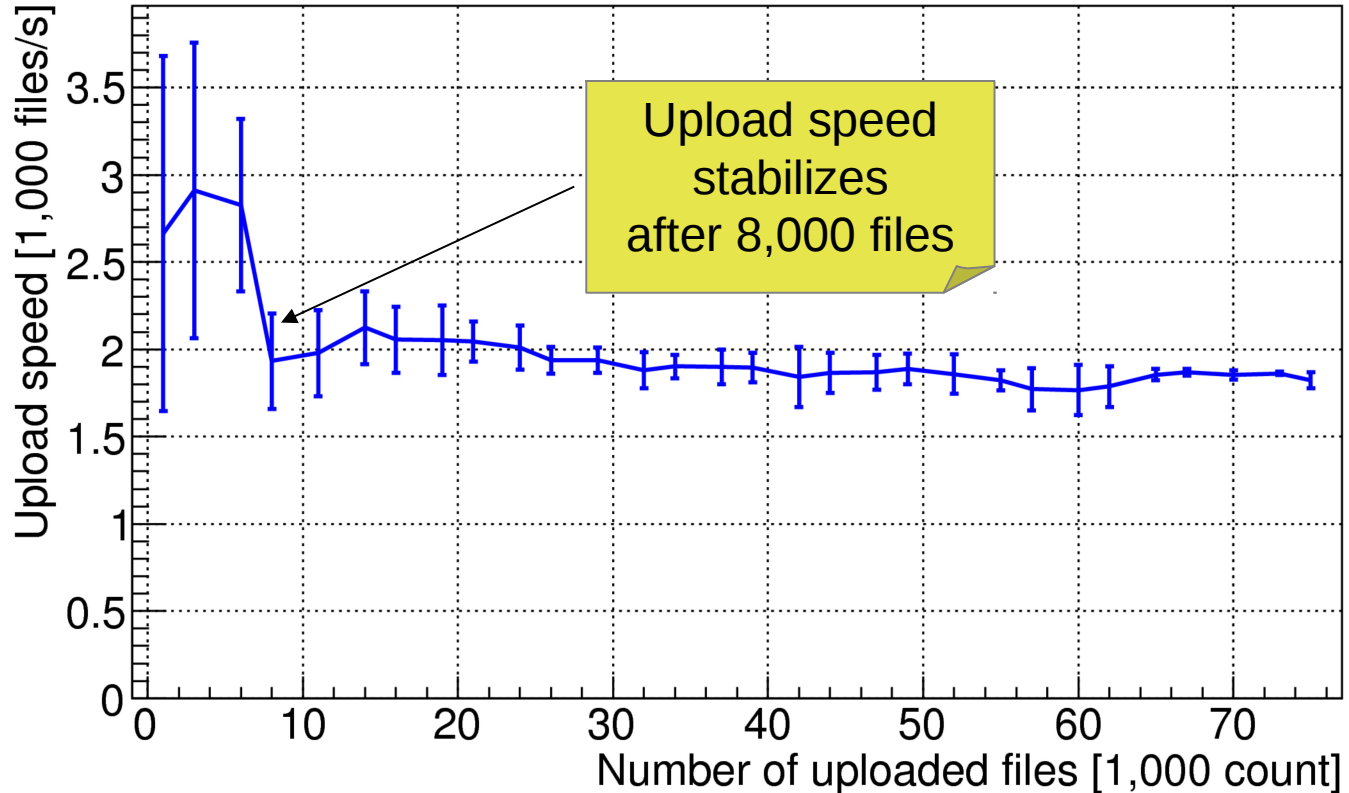  - Uploading of only new files feature tested

DSS

- Problem: no multi-part upload support
  - CVMFS software is foreseen to require multi-part uploads to S3 cloud storages in the near future

- Solution: supported in the new version
  - Current version of the Huawei cloud storage in CERN does not support multi-part uploads, but latest version does
  - New version will be tested when deployed in CERN
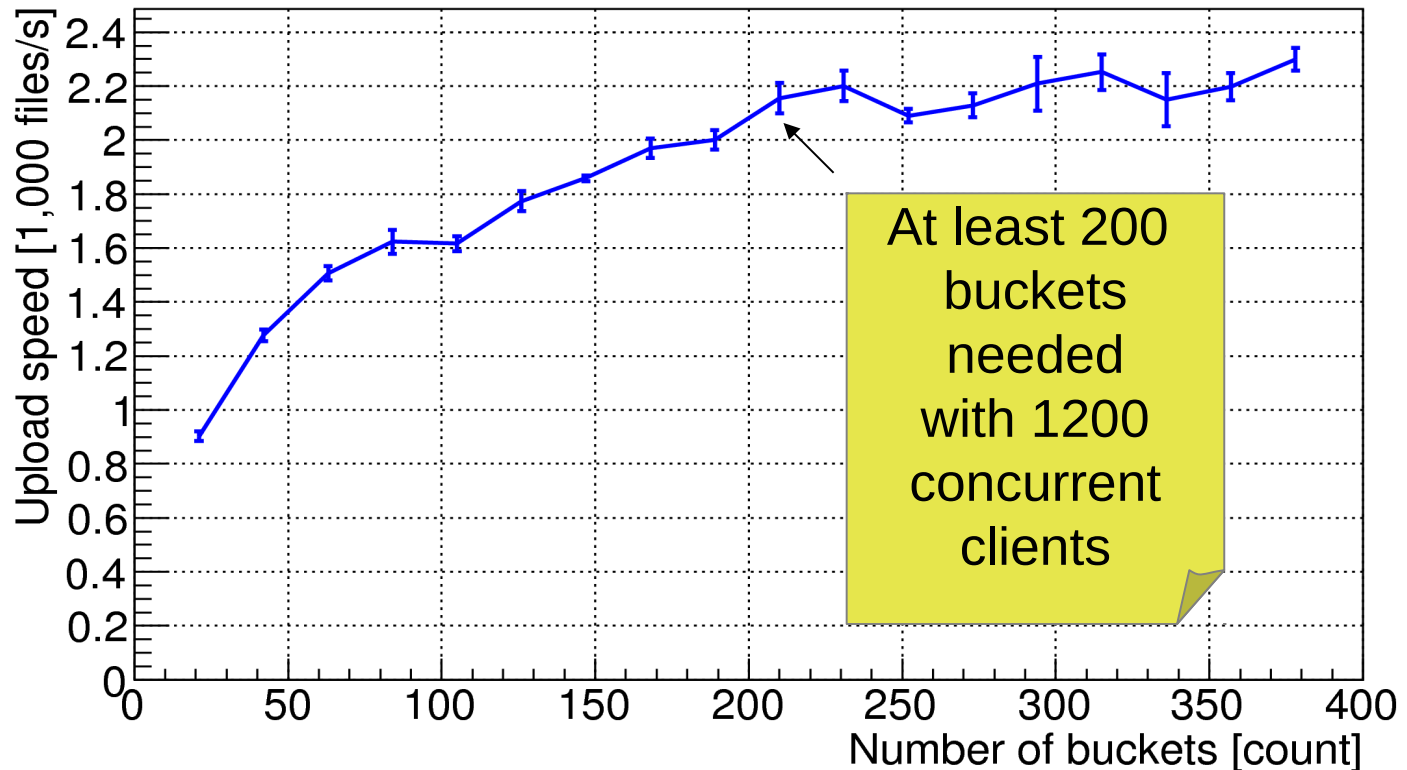
?

!

# DSS

- Problem: new ROOT S3 plugin
  - New ROOT S3 plugin was released, is it working properly with Huawei cloud storage?

- Solution: tested to work with one client
  - One client read-performance identical to the old ROOT S3 plugin
  - Multi-client stress tests are planned

?

!

CERN IT Department

- Problem: consecutive uploads
  - Does number of consecutive uploads affect the upload speed

?

!



Upload speed
stabilizes
after 8,000 files

- Problem: how many buckets needed **?**
  - How the number of used buckets affects the maximum achievable upload speed



At least 200 buckets needed with 1200 concurrent clients

**!**

CERN IT Department
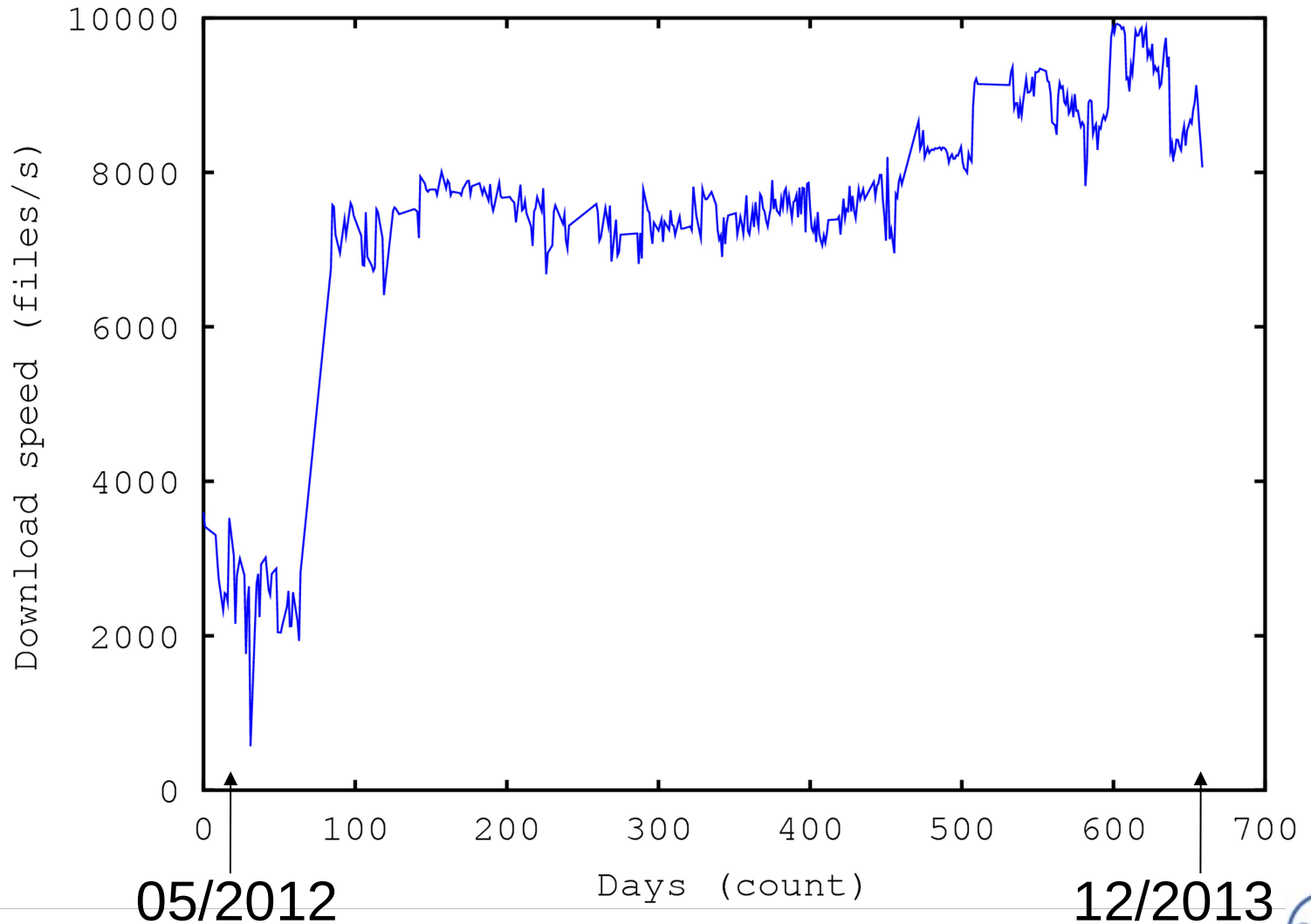
- Problem: bucket usage affects speed?
  - How the number of files in a bucket affects the maximum achievable upload speed



Not linear!
Fixed already
in the new
version of the
cloud storage

CERN IT Department

One thread downloads since 2012-04-24

DSS

CERN IT Department

- ## Identified new requirements

  - Multi-part file upload support
  - Bucket fullness should not affect the upload performance
  - Fast upload speeds should not require hundreds of buckets

  *Fixed already in the new version!*

- ## Test results

  - New ROOT S3 plugin worked without problems with the Huawei cloud storage
  - Long-term download stability good

- # Short term
  - – Benchmark CVMFS with real release data
  - – Test ROOT S3 plugin performance with multiple clients

- # Long term
  - – Second petabyte system with enterprise disks expected to arrive soon
  - – Replication tests between cloud storages
  - – Prove total cost of ownership (TCO) gains of the system as part of a production service

**DSS**

- # Short term

  - Benchmark CVMFS with real release data
  - Test ROOT S3 plugin performance with multiple clients

- # Long term

  - Second petabyte system with enterprise disks expected to arrive soon
  - Replication tests between cloud storages
  - Prove total cost of ownership (TCO) gains of the system as part of a production service

# Thank you!
## seppo.heikkila@cern.ch

# Huawei Cloud Storage

Seppo S. Heikkila
Maria Arsuaga Rios
CERN IT

Openlab Major Review Meeting

13th of February 2014

CERN, Geneva