CERN**IT** Department

# Huawei Cloud Storage

Maitane Zotes Resines, CERN IT

Openlab Major Review Meeting
27. September 2012
CERN, Geneva

HUAWEI

CERN
openlab

- Huawei setup
- Remaining Issues after the first phase
  – see last major review
- Benchmark execution framework
- Initial results
- S3 access with ROOT and S3FS
- System upgrade
- New results and comparison
- Future workplan

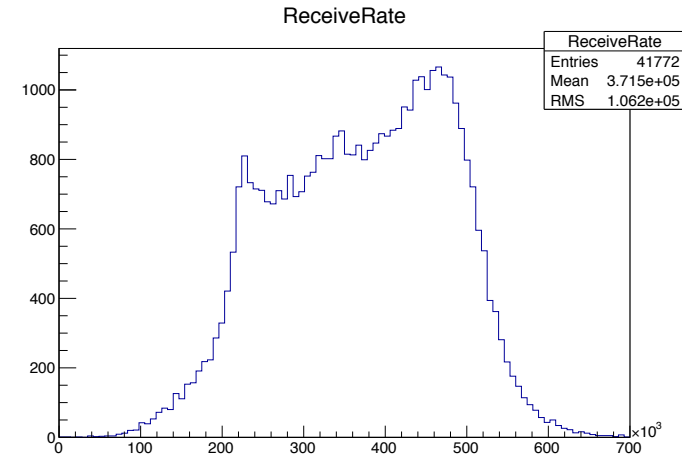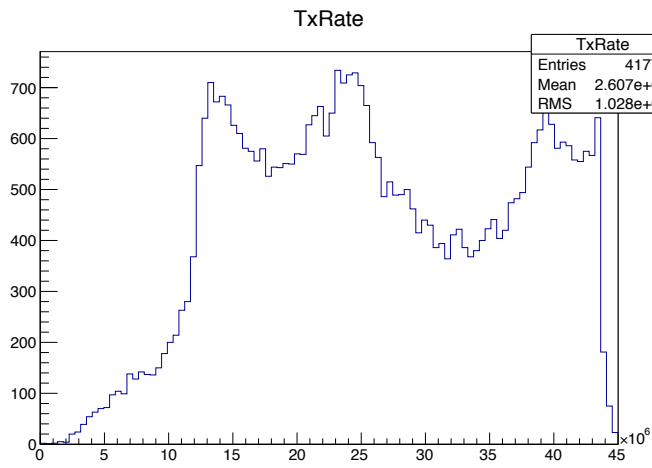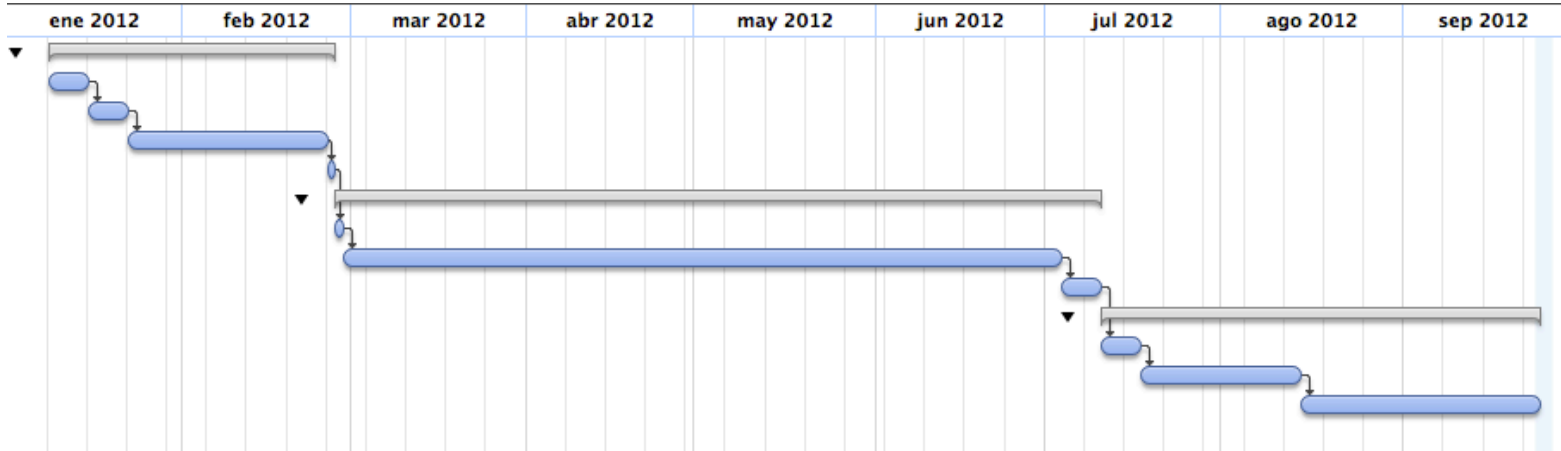# Huawei Setup

OSC

# Huawei Setup

SOD

OSC

SOD

**DSS**

- ## New framework integrated with ROOT
  - Histograms
  - Integrated with the Python benchmark
  - ssh connection to clients using lxplus

| | 1st Phase | 2nd Phase | 3rd Phase |
|---|---|---|---|
| | **1st Phase** | **2nd Phase** | **3rd Phase** |
| | Project start | Obtain 5 client boxes | Redo tests |
| | Mounting of system | Redo tests with new clients | Download issue found |
| | First tests | Upgrade of system | Fix issue |
| | First Review | | Redo download tests |
| | | | Obtain 20 client boxes |
| | | | Start 20 client tests |
| | | | Openlab Major Review |

# Limitations observed in first testing phase

- ## Upload limit
  - 50 files/second

- ## Download limit
  - 425 files/second

- ## Benchmark issues
  - Delayed closing of client sockets produces socket shortage
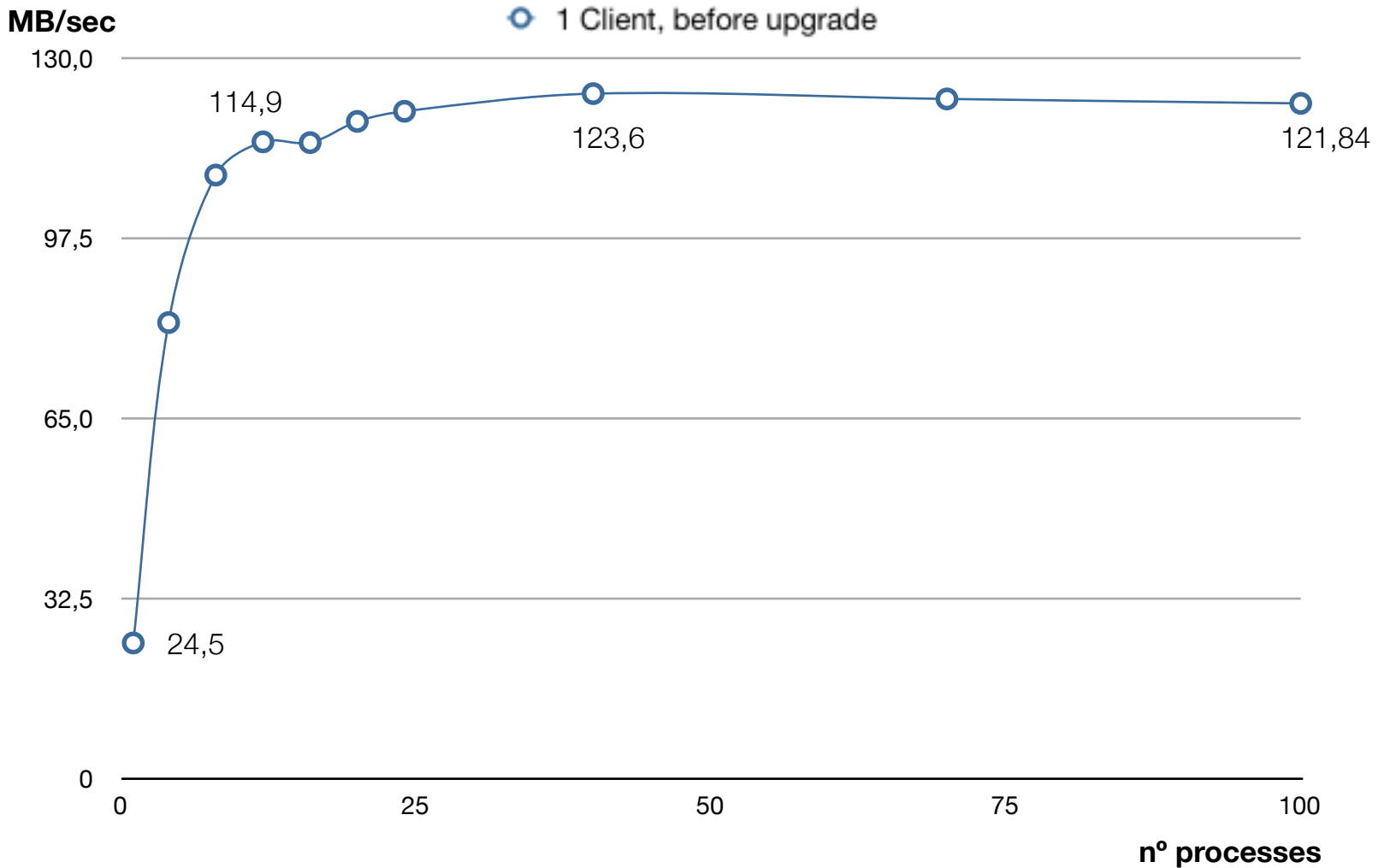  - Client bandwidth limit
  - Client memory limit
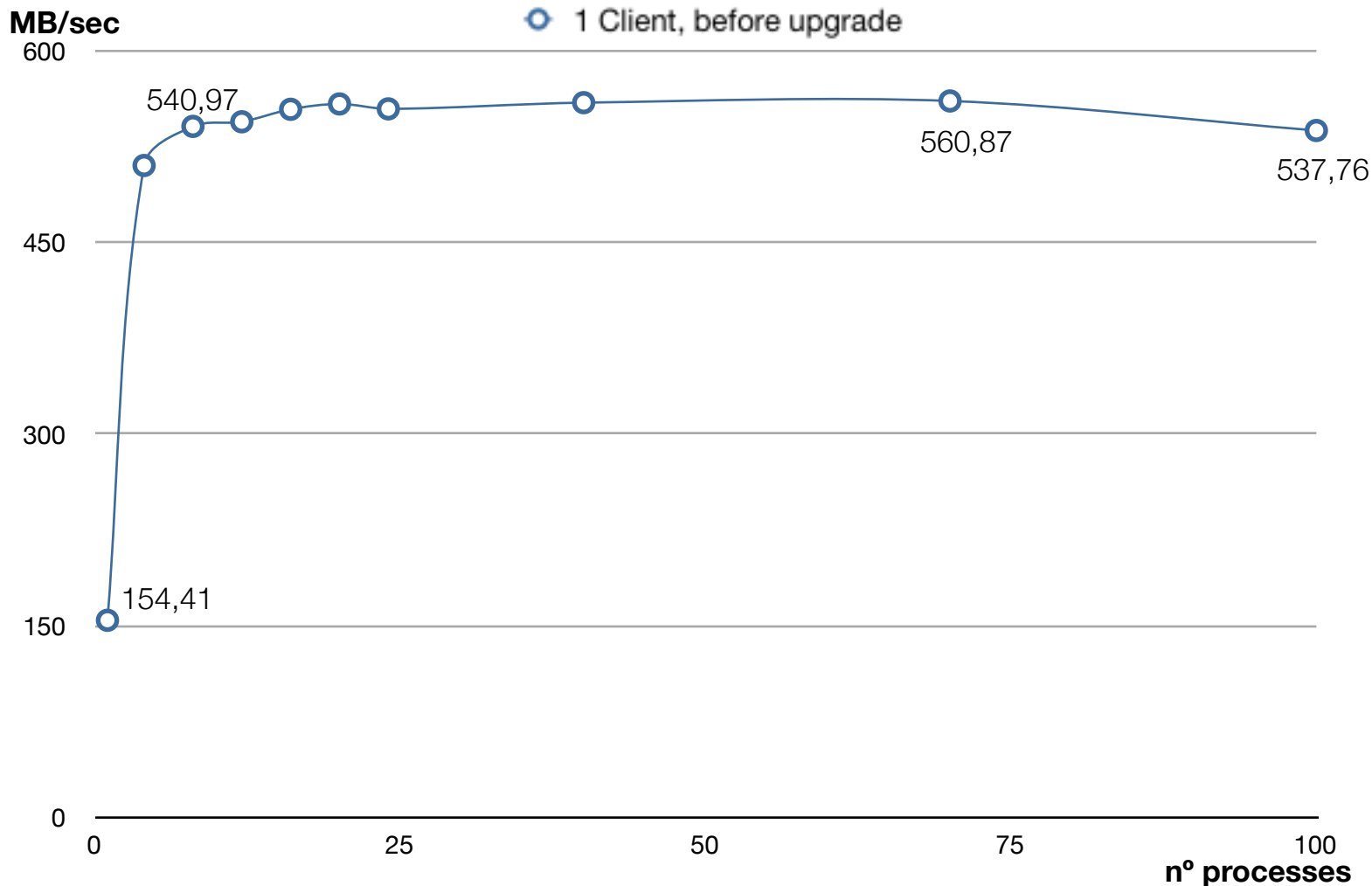
Succesfull scalability of metadata

# 4KB Download - Small Reads

files/sec

O 1 Client, before upgrade

- 52,3
- 438,4
- 654
- 689,6
- 682,8

n° processes

Successful scalability of metadata

Scales properly and bandwith limit of 5Gb reached

# 100MB Download - Sequential access of bulk data

**MB/sec**

○ 1 Client, before upgrade

- 540,97
- 560,87
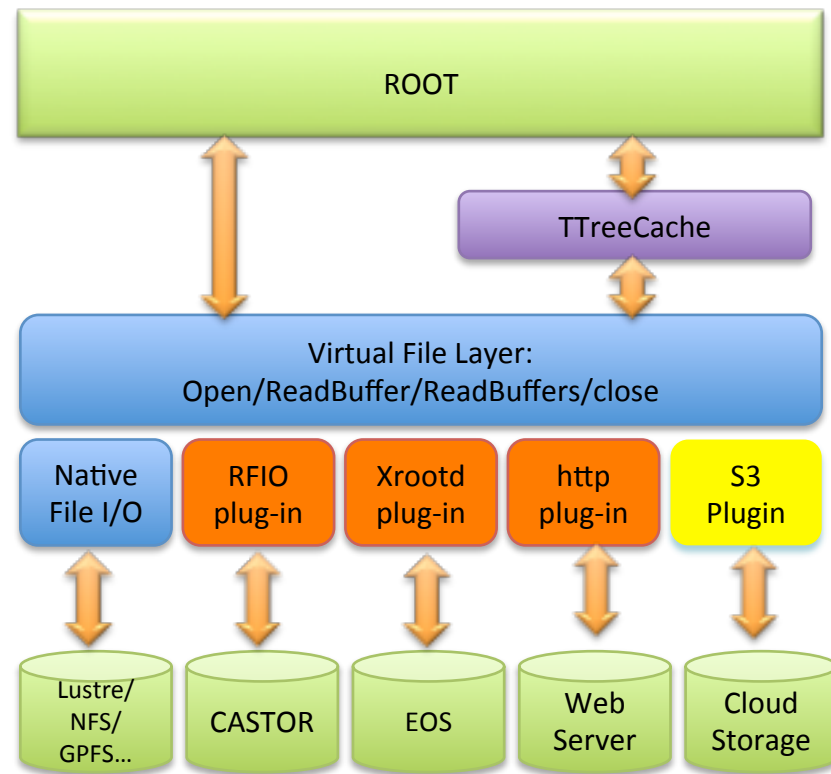- 537,76
- 154,41

**nº processes**

Scales properly and bandwith limit of 5Gb reached

**DSS**

- The system behaves as expected

- 4KB Tests
    - 110,7 files/sec upload limit
    - 682,8 files/sec download limit

- 100MB Tests
    - Reached the client bandwidth limit
    - Stable

- Huawei UDS is a typical cloud storage which provides Amazon S3 Interface

- ROOT is designed to support different I/O protocol
  - rfio, xrootd, posix ...

- Two solutions to run ROOT on cloud storage
  - develop a S3 Plug-in for ROOT (validated by tests)
  - develop a fuse-based POSIX file system to mount cloud storage
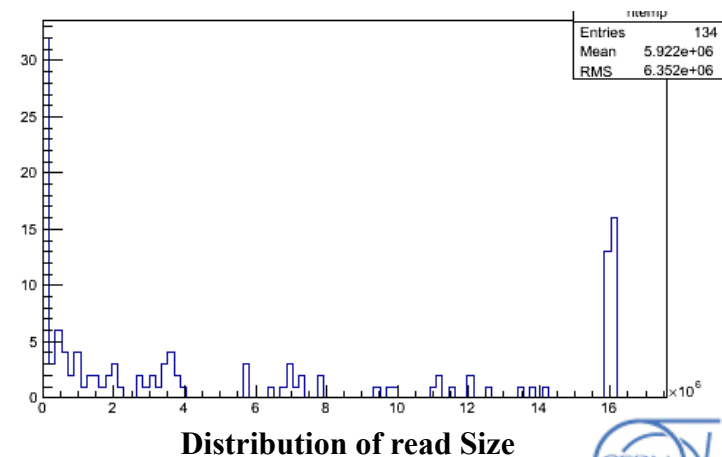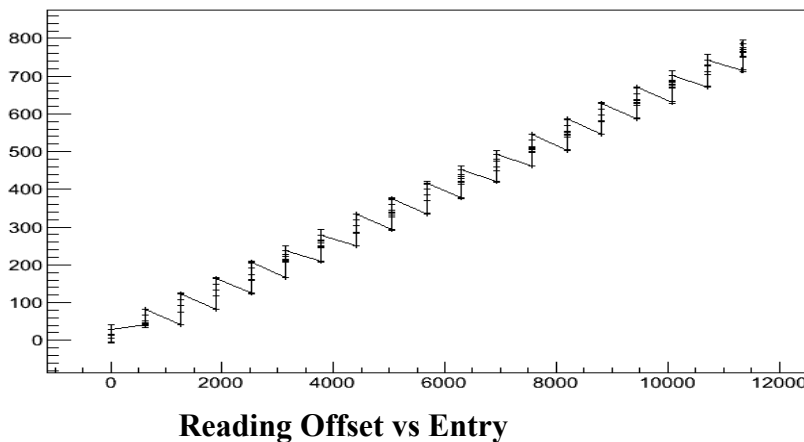
- Based on the http plug-in

- Adapts to S3 protocol

- Supports *vector read* requests issued by TTree cache
  - Huawei added multi-range read to S3 API

- Integrated with the distributed I/O test framework

- ## Open source prototype S3FS

  - http://code.google.com/p/s3fs/

- ## Current limitations:

  - Can only mount one single bucket instead of the whole system

  - Instead of remote I/O, file is downloaded to local cache during "open"
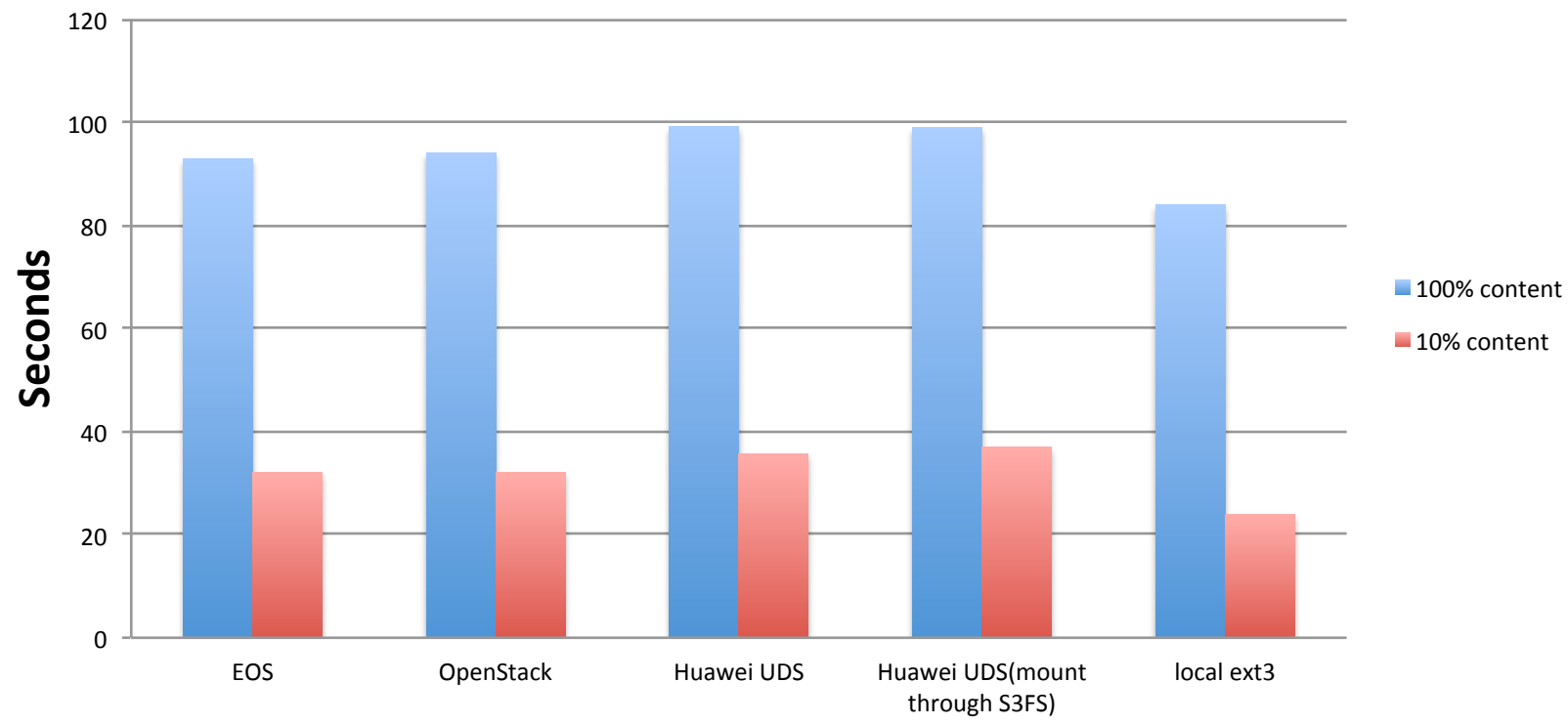
  - "df" returns not relevant information

- # A real ATLAS ROOT file

  - ## 793MB on disk, 2.11 GB after decompression

  - ## 11918 entries, 5860 branches ,cache size=30MB
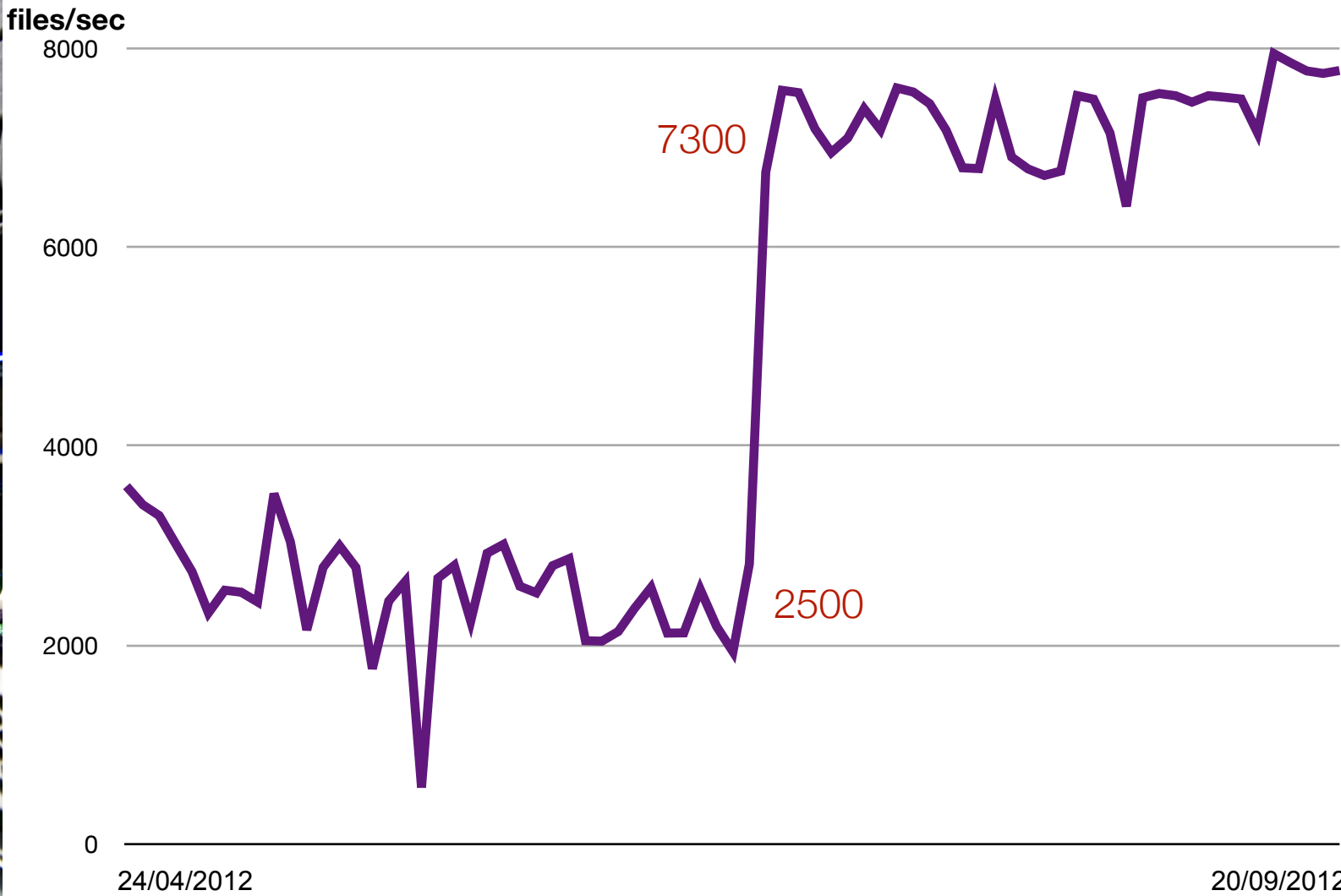
- # Reads in entries sequentially in "*physics tree*"



**Reading Offset vs Entry**



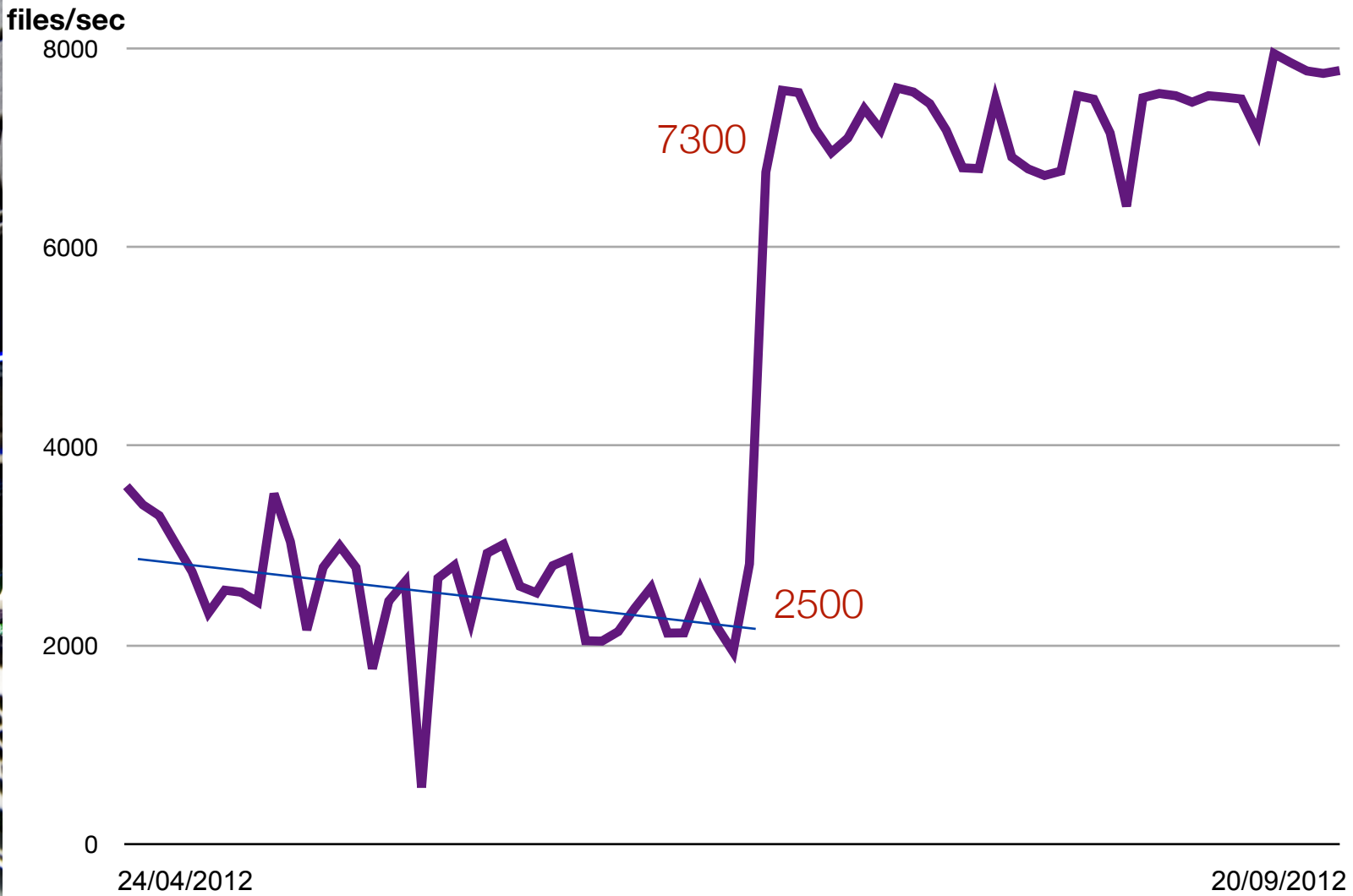**Distribution of read Size**

- Single client
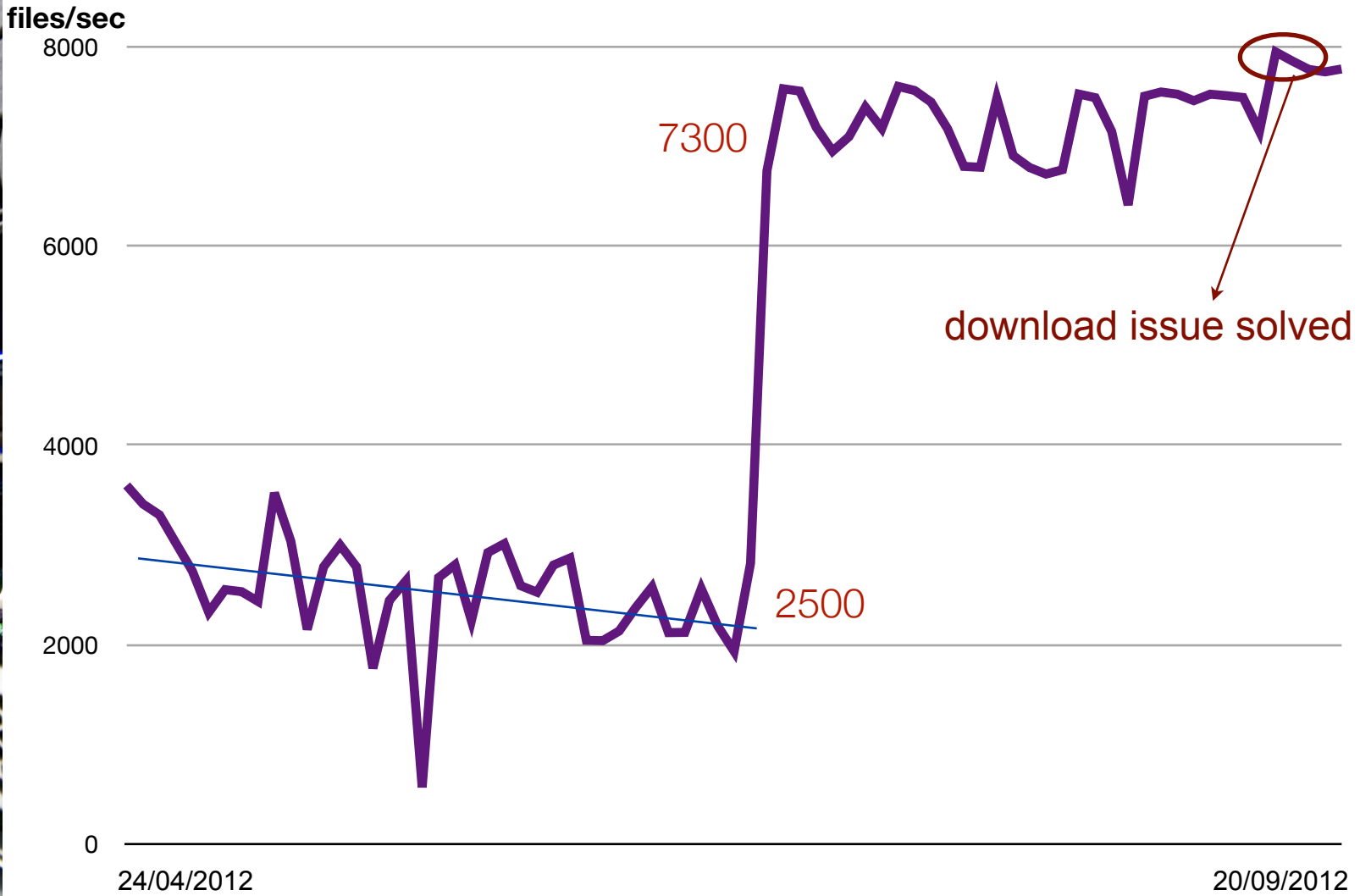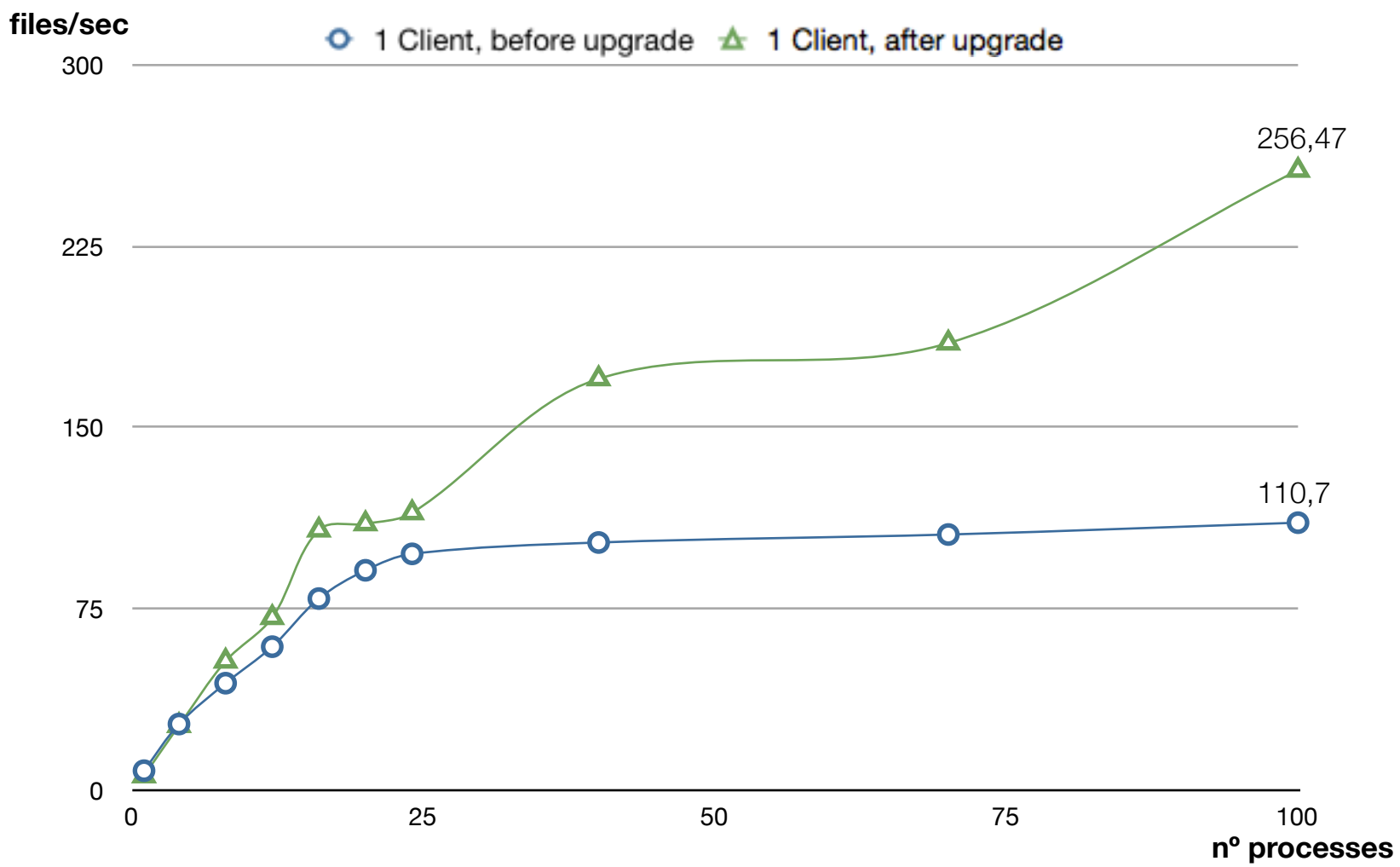
**Wall Time**

**DSS**

- Hardware improvements
- Upgrade of SOD node
- Performance optimization
- Billing function
- User management API
- More S3 features
- Cache turned on/off

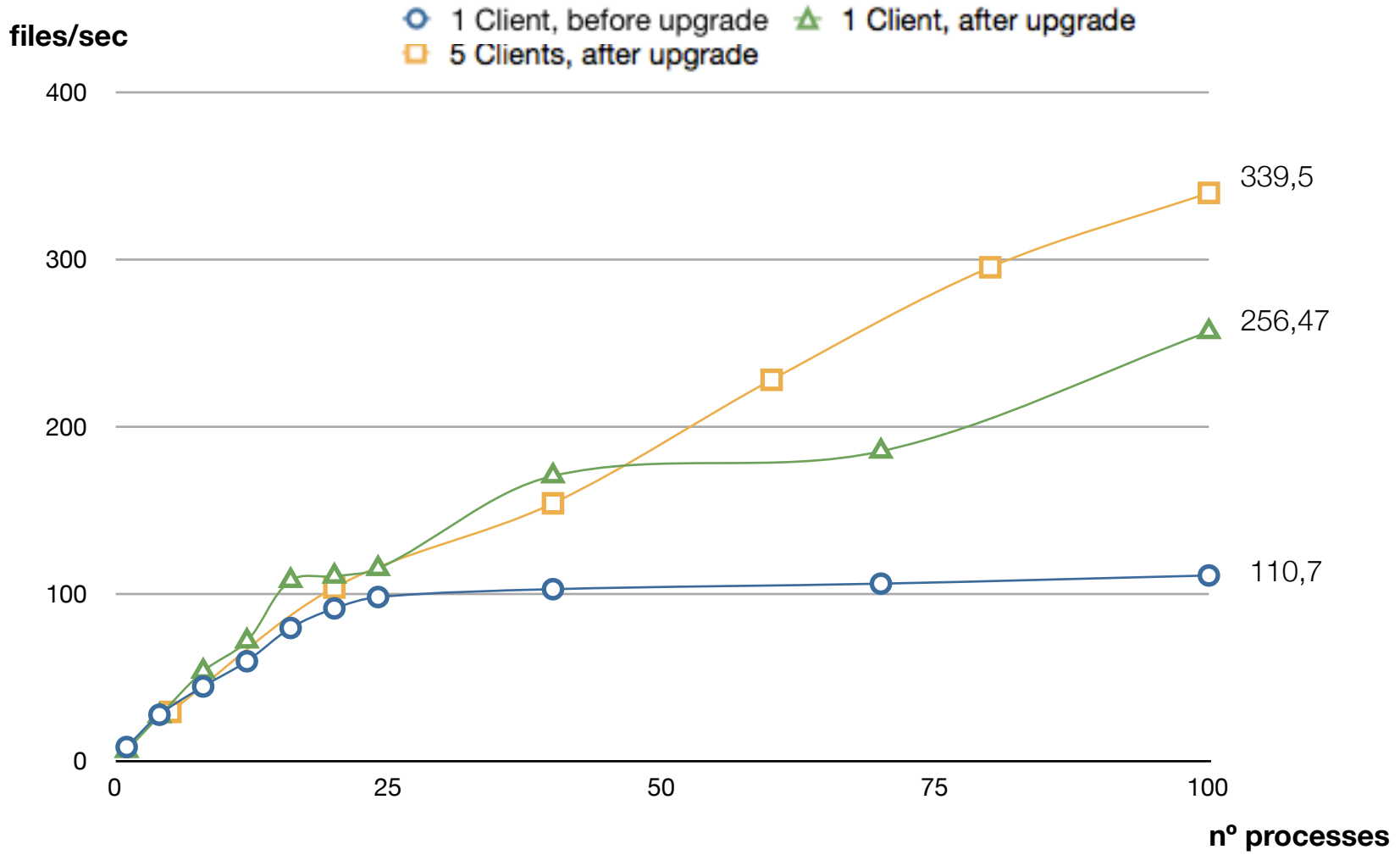DSS

CERN IT Department

**files/sec**

8000

7300

6000

4000

2500

2000

0

24/04/2012

20/09/2012

**date**

files/sec

7300

download issue solved

2500

24/04/2012

20/09/2012

**date**

# 4KB Upload

After the upgrade the performance doubles

382,45

**files/sec**

Legend:
- 1 Client, before upgrade
- 1 Client, after upgrade
- 5 Clients, after upgrade

339,5

256,47

110,7

**nº processes**

5 client tests improve the numbers

CERN**IT**
Department



files/sec

○ 1 Client, before upgrade △ 1 Client, after upgrade

3751

682

n° processes

Performance after the upgrade five times better

**files/sec**

Legend:
- ○ 1 Client, before upgrade
- △ 1 Client, after upgrade
- ☐ 5 Clients, after upgrade

7250

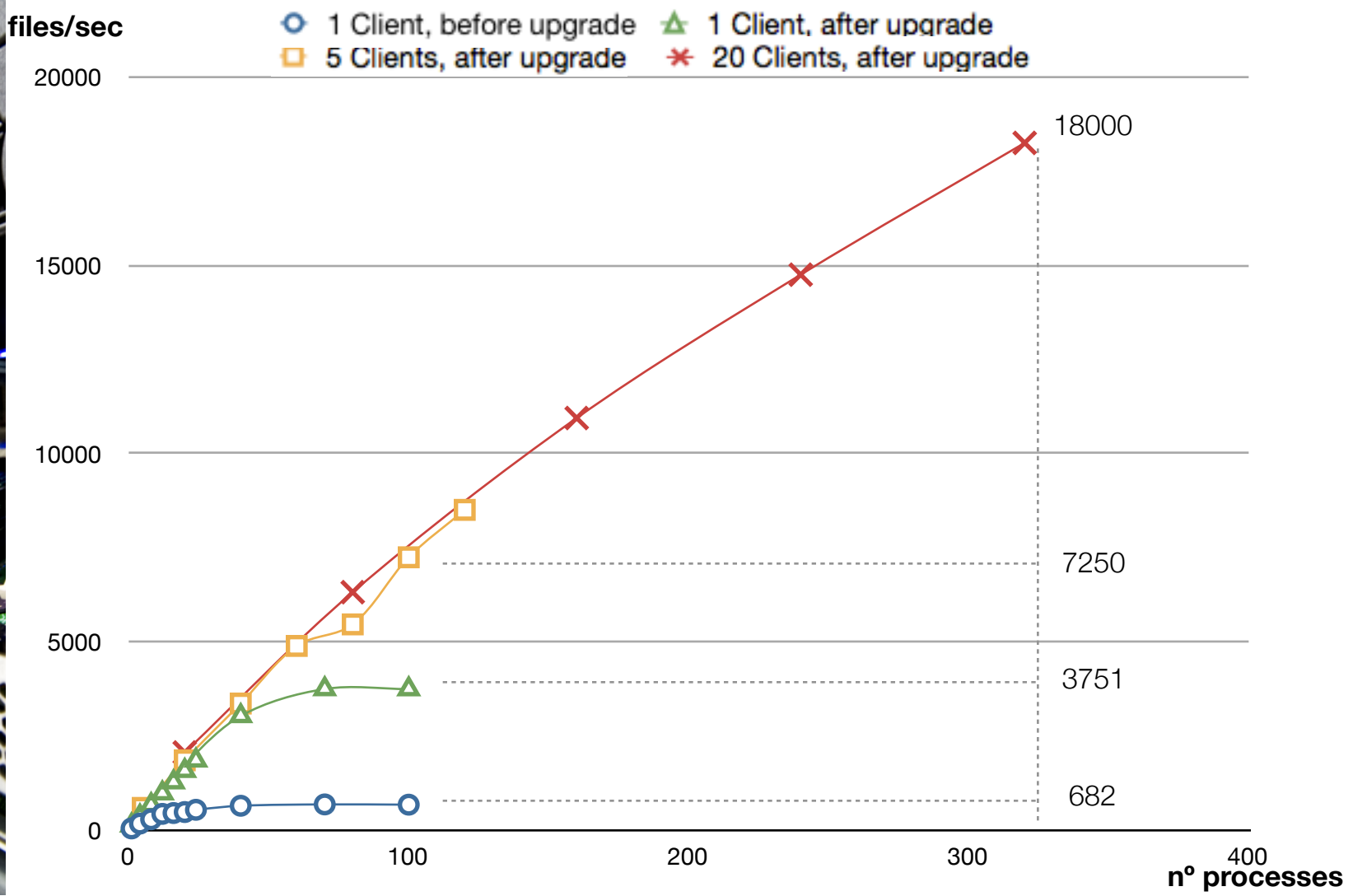3751

682

**nº processes**

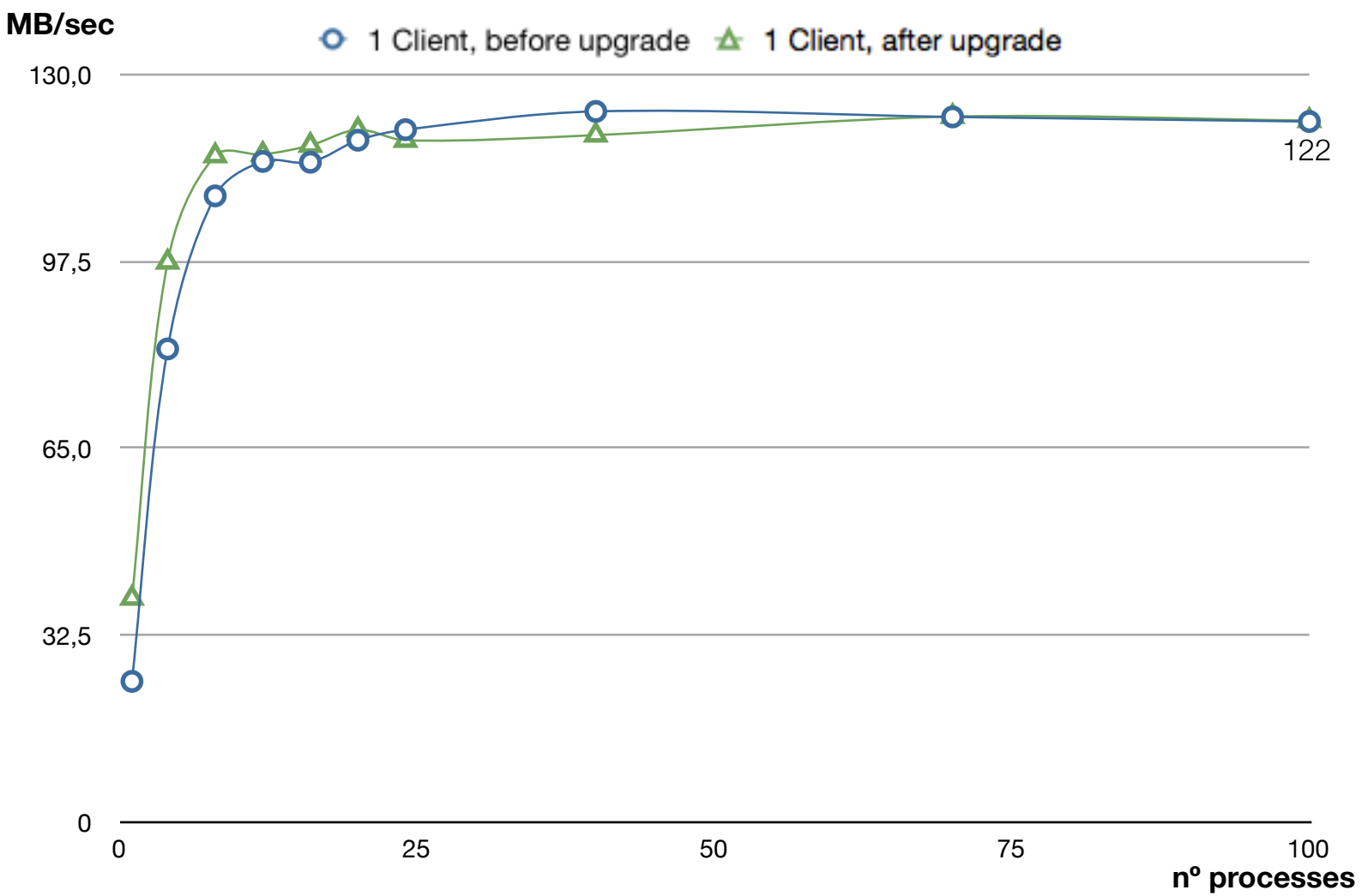5 client test doubles last better result

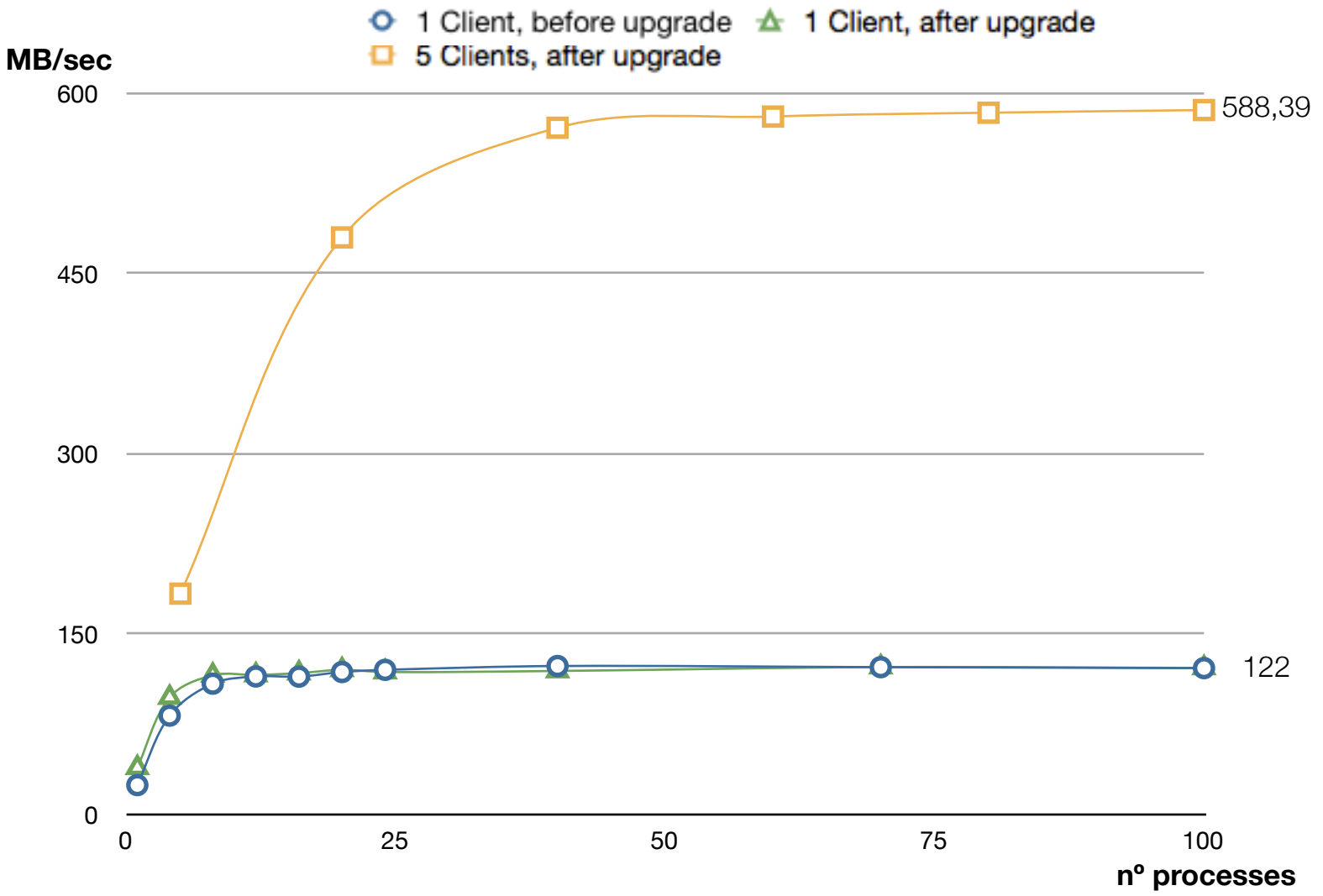Succesfull  scalability of metadata in 20 boxes

DSS

Succesfull scalability of metadata in 20 boxes

# 100MB Upload
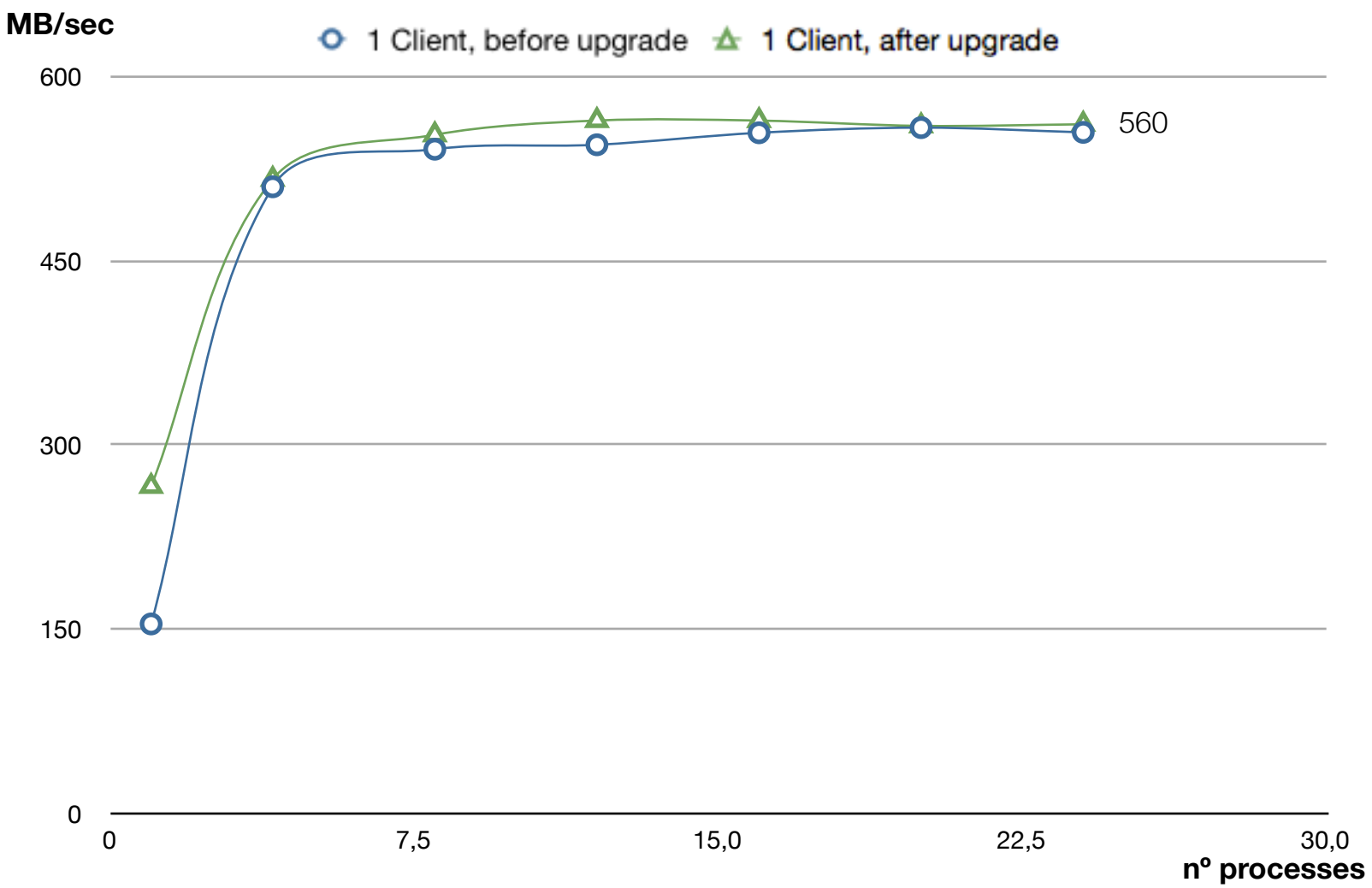
After upgrade slighly better, and both reach the bandwidth limit

5 clients also reach bandwidth limit of 5Gb

# 100MB Download

**MB/sec**

○ 1 Client, before upgrade    △ 1 Client, after upgrade

560

**nº processes**

After upgrade slighly better, and both reach the bandwidth limit

20 clients reach bandwidth limit of 18Gb

**MB/sec**

Legend:
- ○ 1 Client, before upgrade
- △ 1 Client, after upgrade
- □ 5 Clients, after upgrade
- ✳ 20 Clients, after upgrade

n° processes

2200
960
560

20 clients reach bandwidth limit of 18Gb

**DSS**

CERN IT Department

- Tested metadata performance up to 8k files/s
  - proved expected scalability of the system
- Tested total throughput up to 18 Gb/s
  - fully maxed out the 2 fibres available
  - balanced system with 350MB/s per OSC

- Fully achieved expected performance

- Minor technical problems found and resolved rapidly
  - very productive collaboration with clear benefits for CERN and HUAWEI

- **2012 - short term**
  - Further increase scalability range with additional network and client resources
  - Analyse performance impact of write cache and journal
  - Collect feedback ATLAS workload management system
  - Excercise transparent upgrade procedure

- **2013 - next year**
  - Multiple datacenter tests (eg with IHEP)
  - Erasure code impact on performance and space overhead
  - Prove transparent disk failure recovery with consumer drives

- **Ultimate goal 2013**
  - Prove TCO gains of the system as part of a production service

# Huawei Cloud Storage

Maitane Zotes Resines, CERN IT

Openlab Major Review Meeting
27. September 2012
CERN, Geneva