# ATLAS disk usage patten at BNL
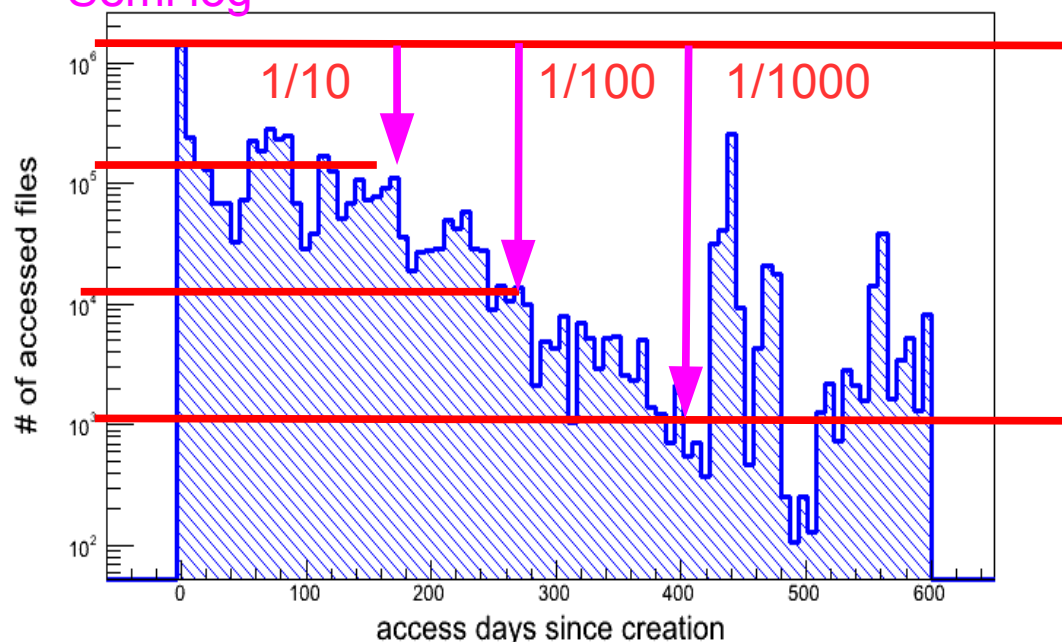
Hironori Ito
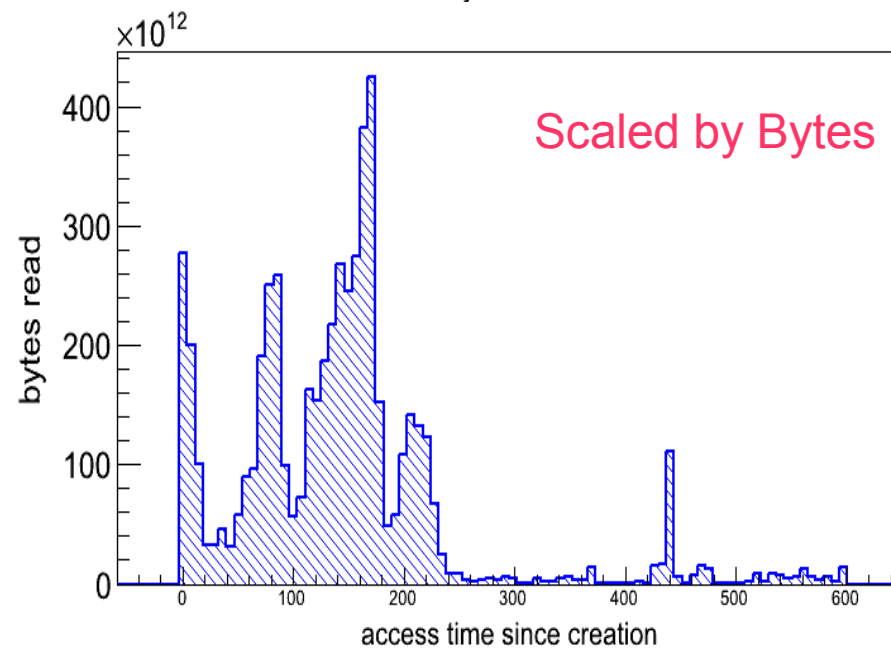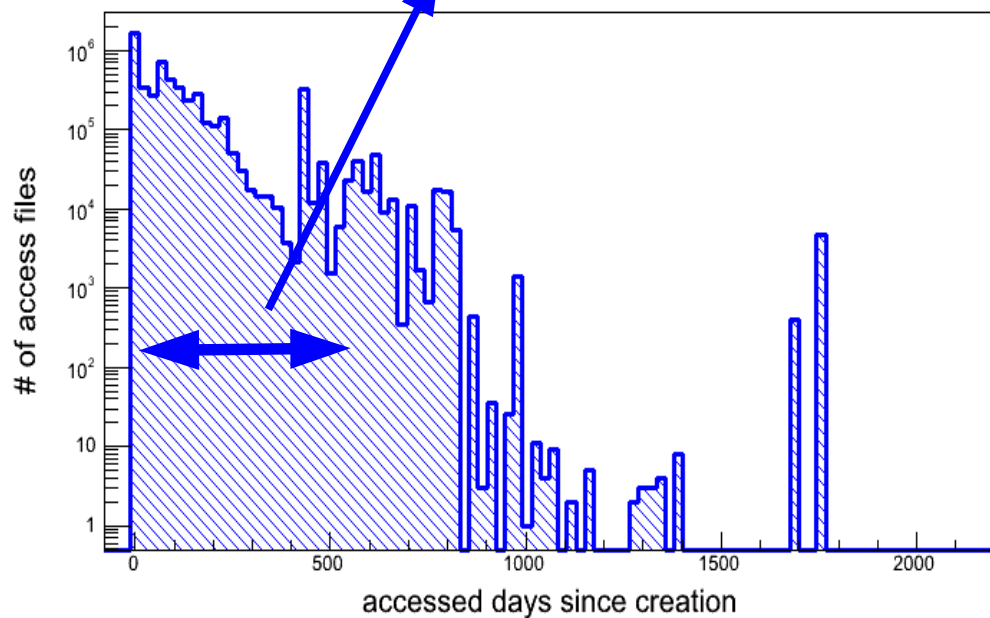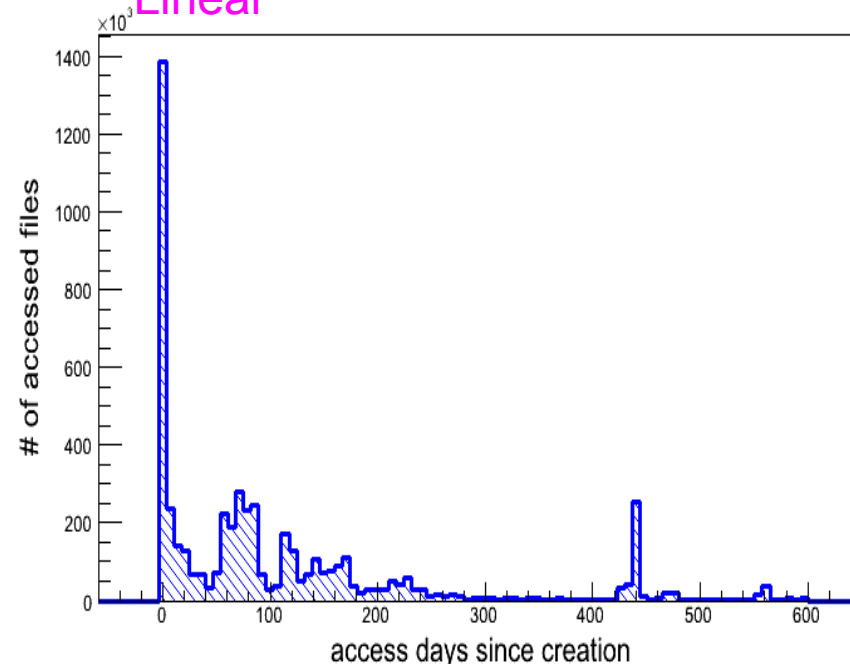Brookhaven National Laboratory

# Motivation and Methods

- Identify the current usage pattern of ATLAS data

  - How often data is read

  - How much data is read

  - Which kind of data is read

- Use dCache's internal Billing Database and Chimera Database as a source of information

  - dCache's billing db records all read/write activities in dCache in PostgreSQL

  - dCache's ChimeraDb records all existing file in dCache.

- Can we design the storage specifically for ATLAS data access pattern instead of the use of generic, all purpose storage.

  - Advantage/disadvantage in performance, cost, features

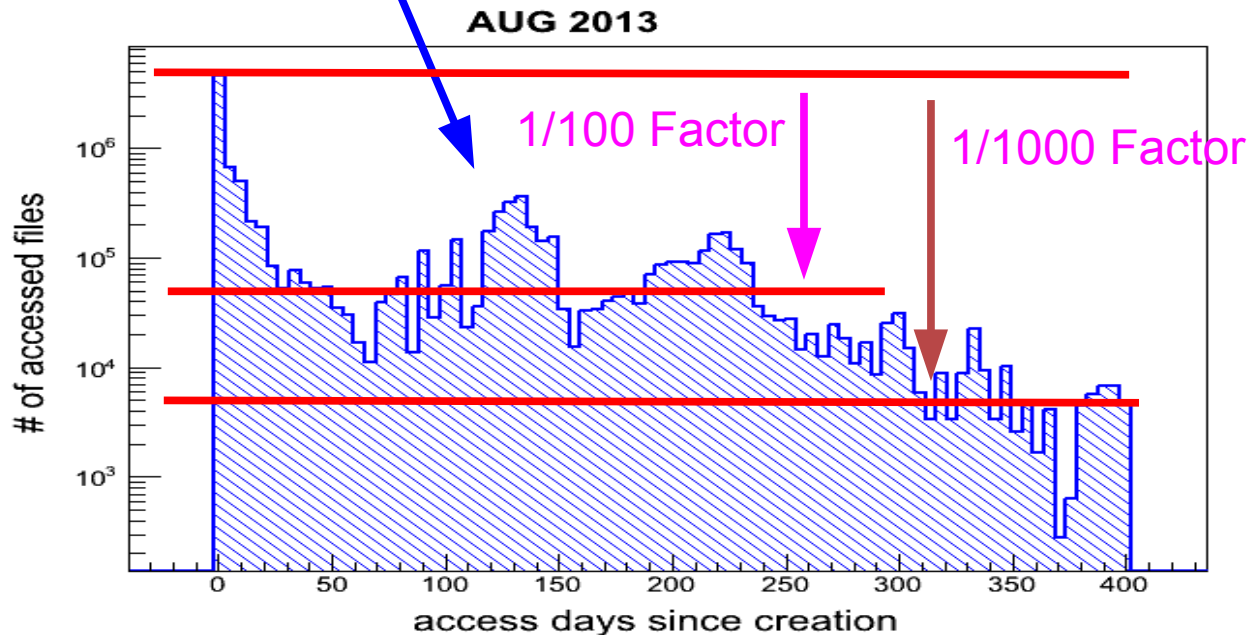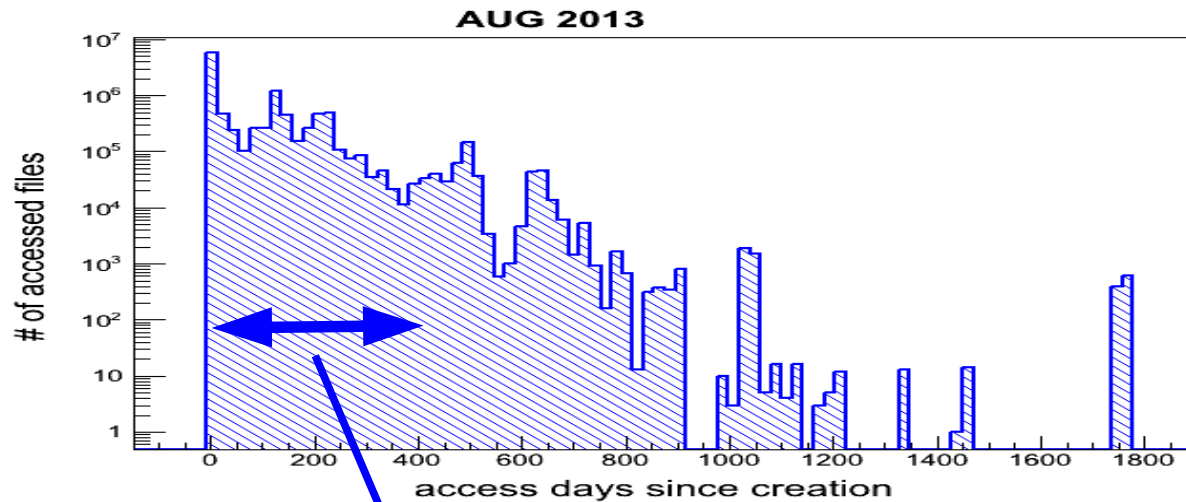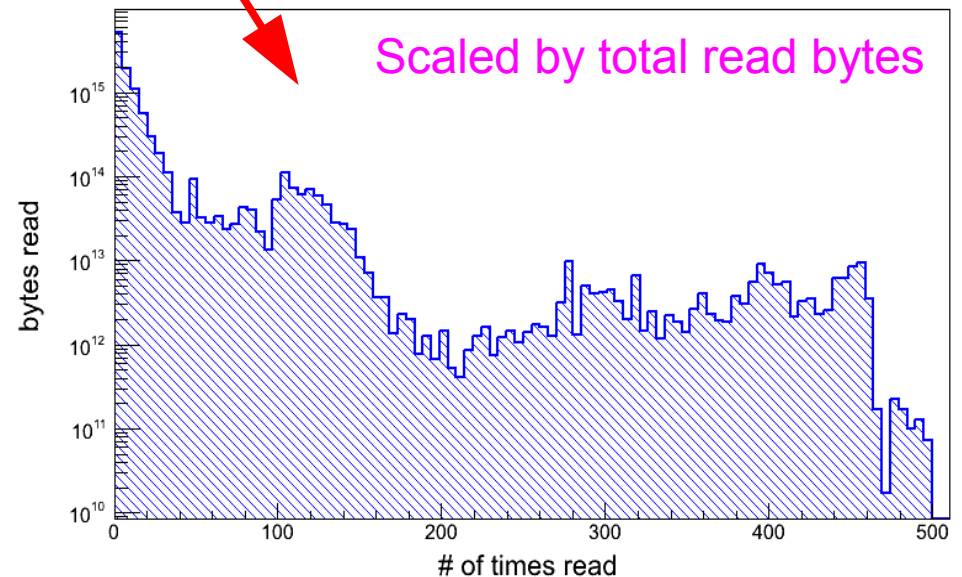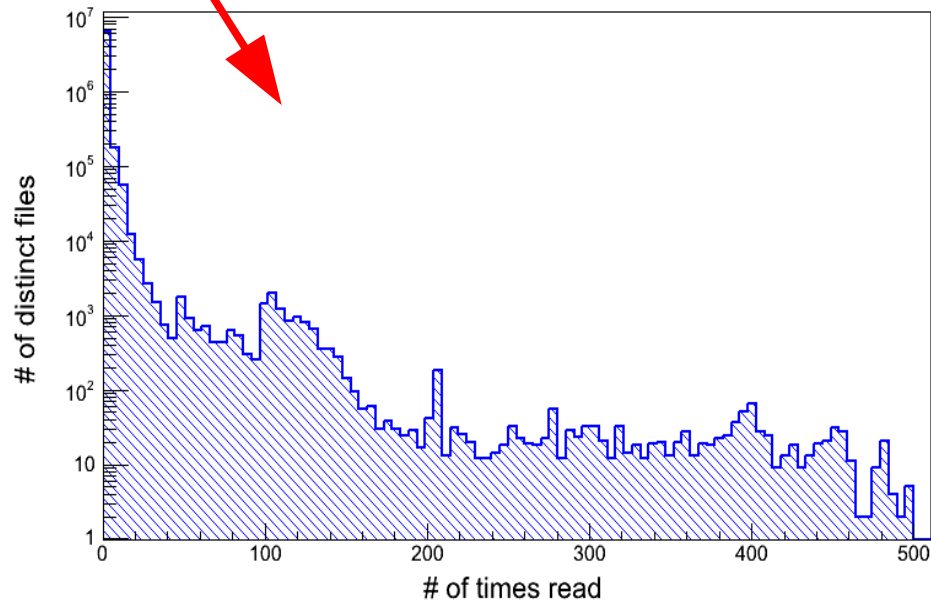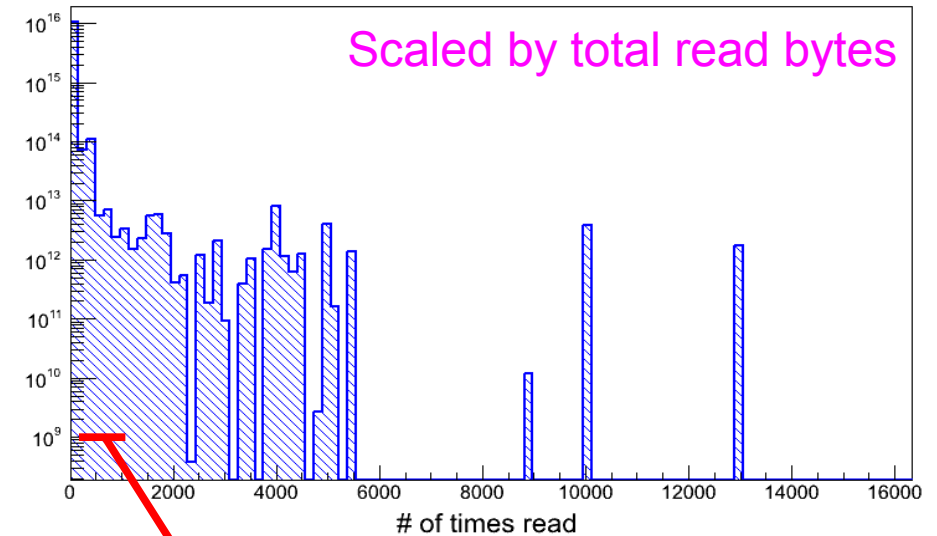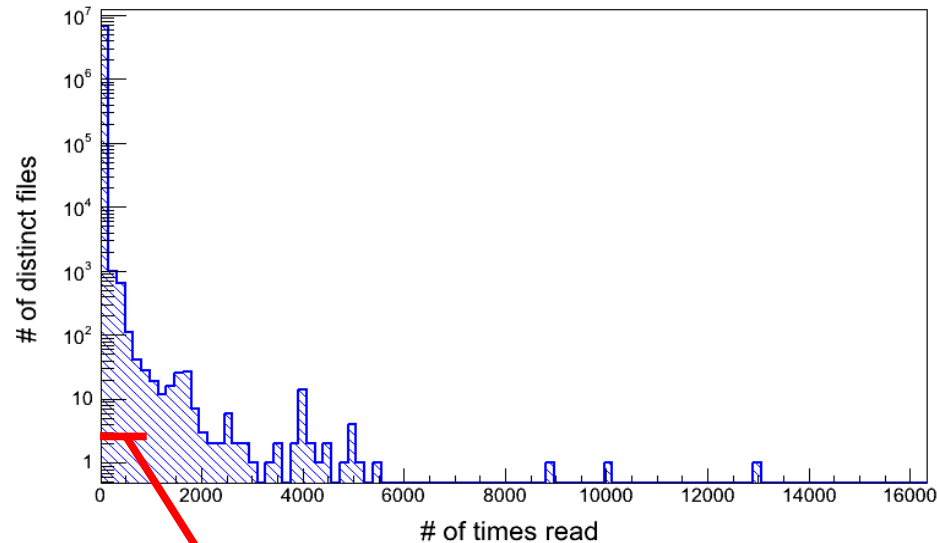# Age of files read during June 2013
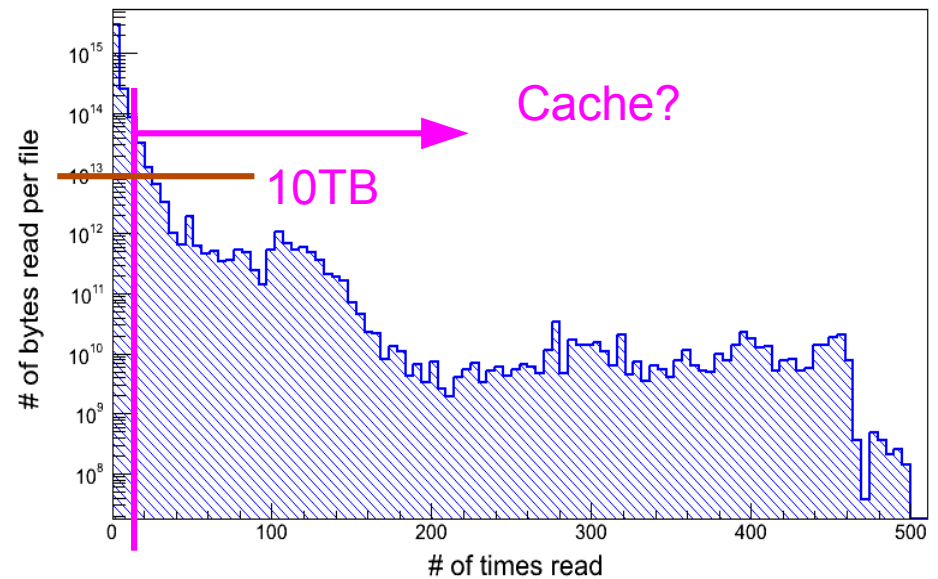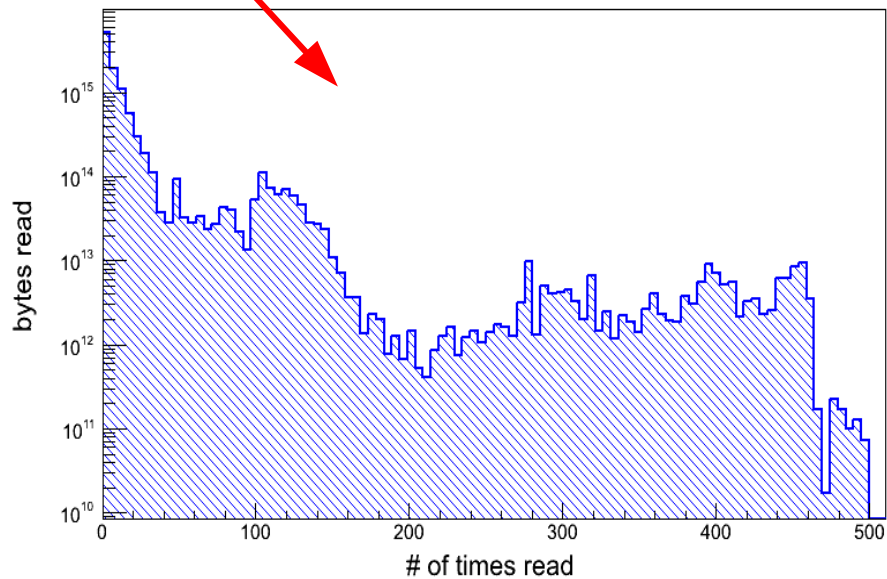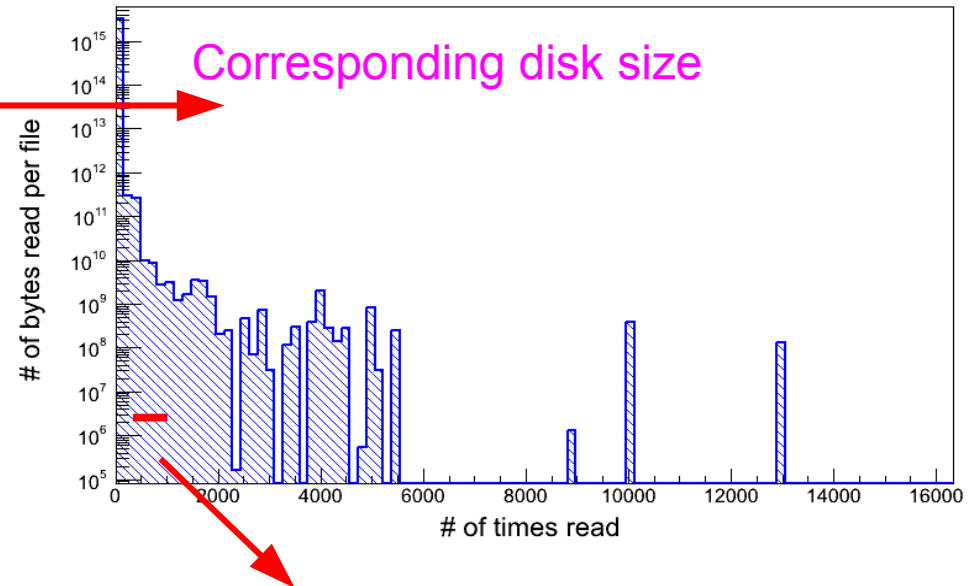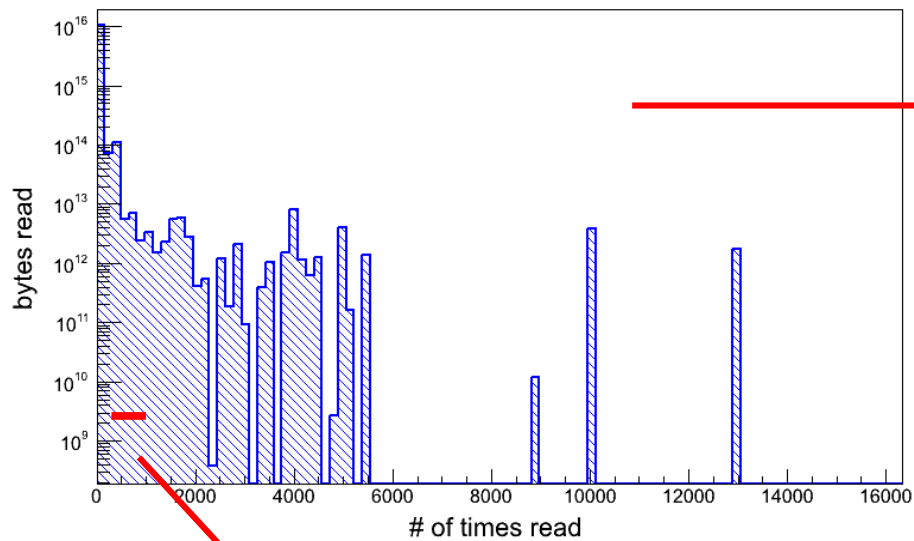
# Age of files read during Aug 2013

# Points for age of files

- There is a patten!  The age is not uniform.

- Most files are read immediately after they were written to storage.

- By the age of 1 year, the chance of being used is about $1/1000^{th}$

  - How much data do we have over 1 year?

    - Larger than a few PB.

    - Is it necessary to be located on the high performance disks without much use?

      - Can the high latency storage be acceptable?

# Frequency of file usage and total network data volume in June 2013
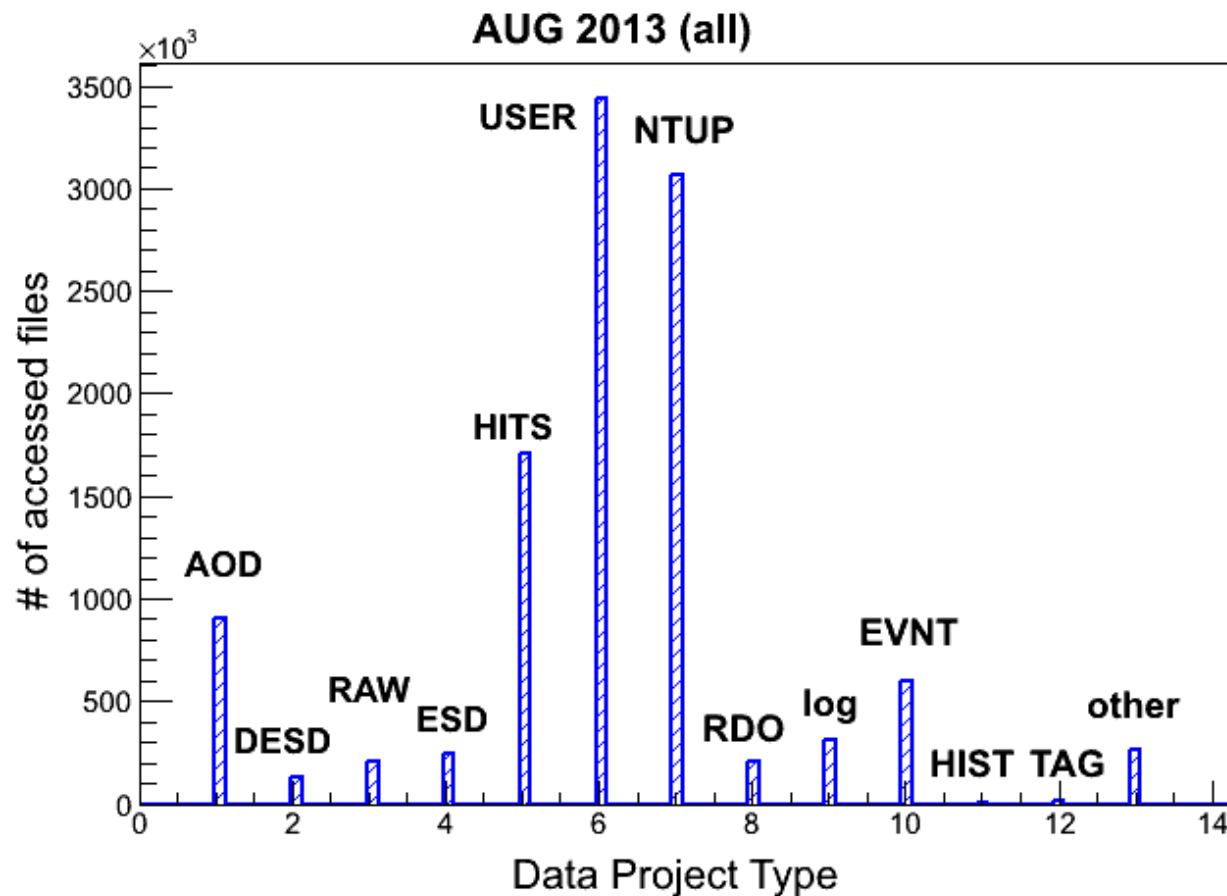
# Network vs Disk data volume in June 2013

# Points for frequency of file reuse

- Large fraction of files are reused only a few times at most.

- Small fraction of files are reused many times. Some of them are used over 1000 times.
    - Which one?

- The storage space required for highly reused data are not too large.
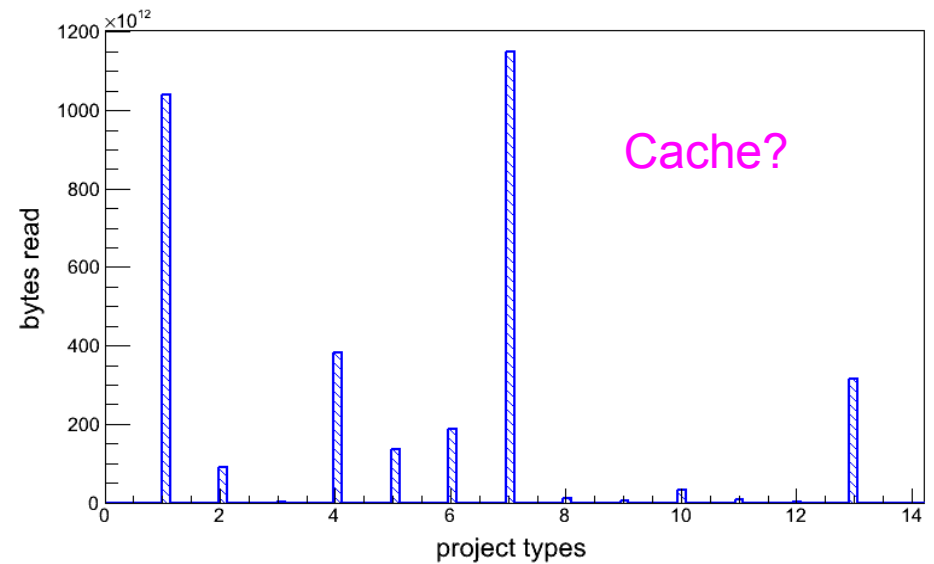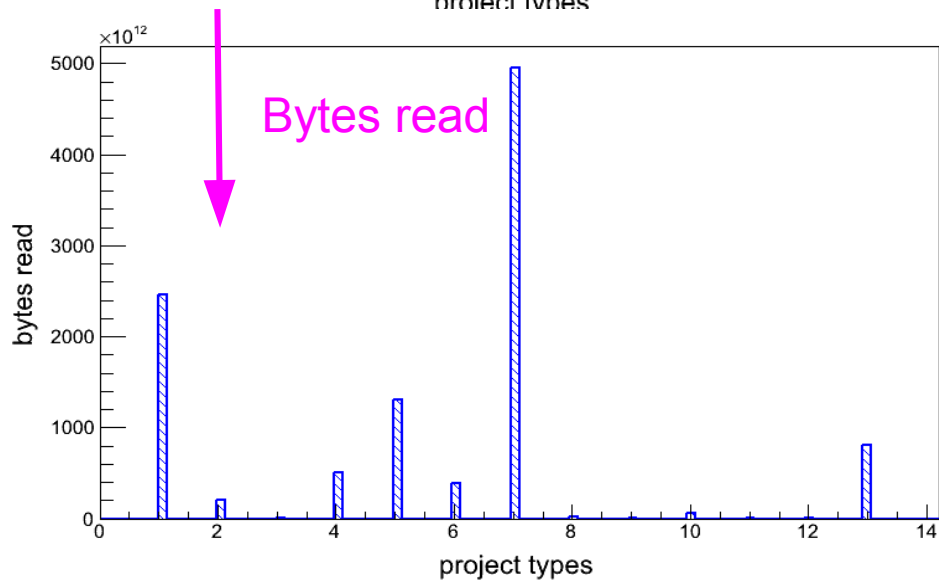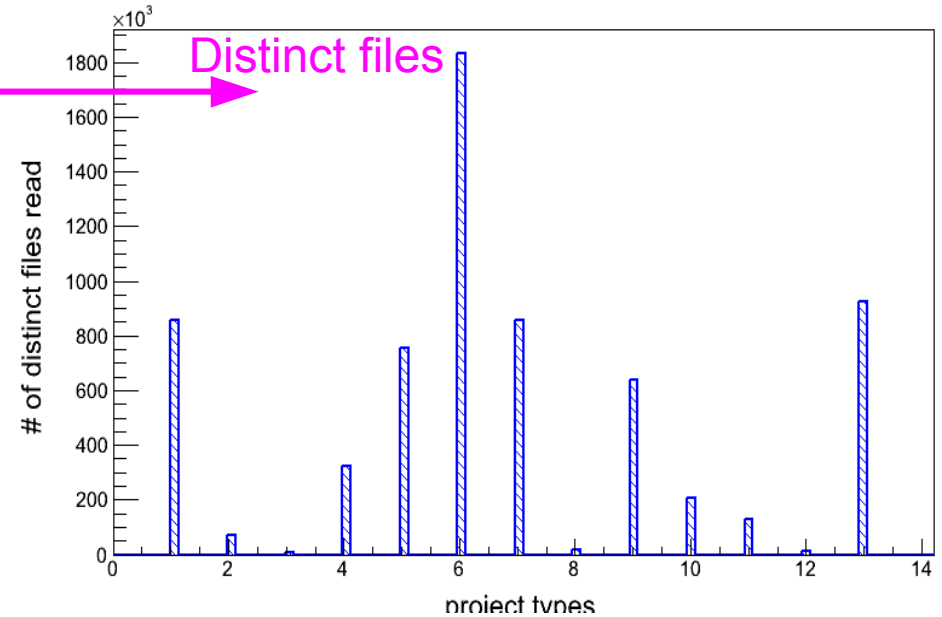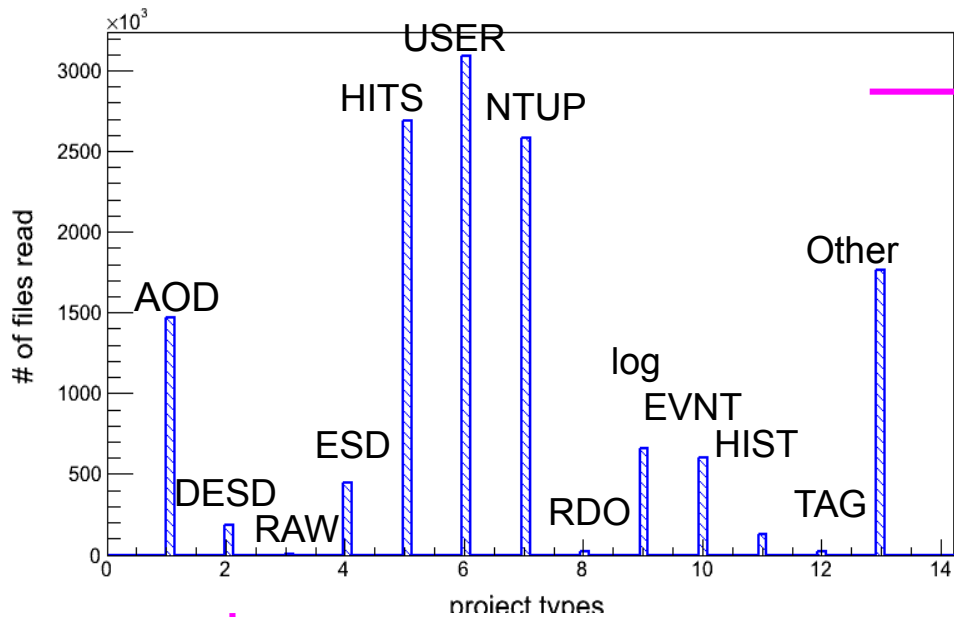    - How much?
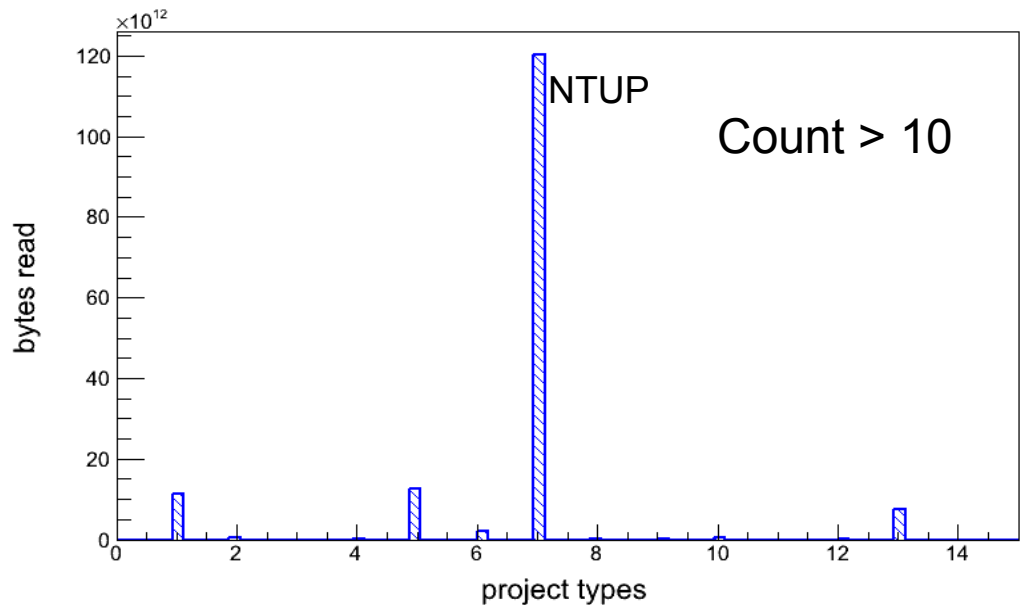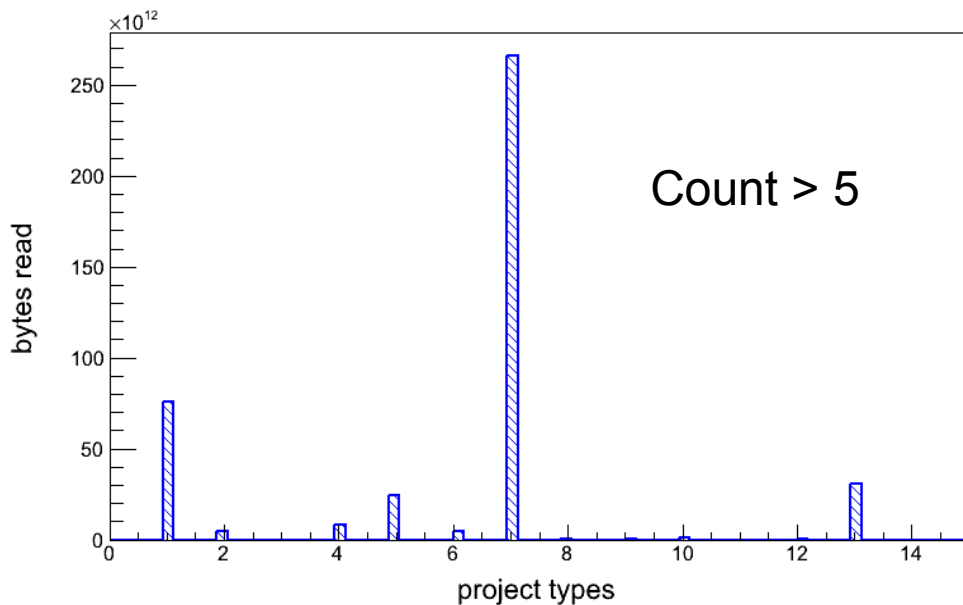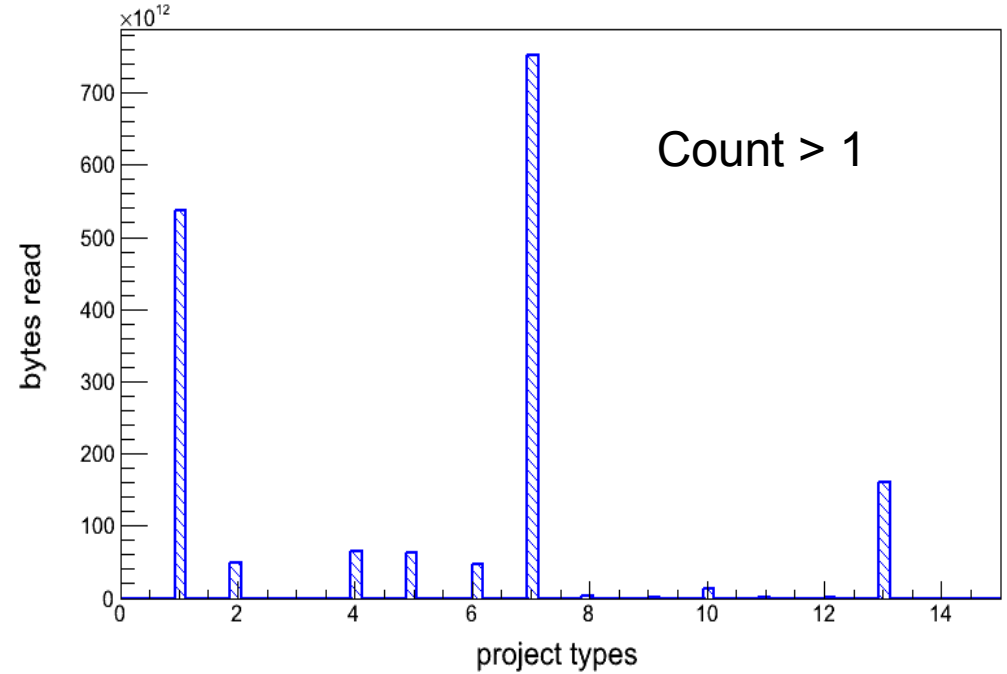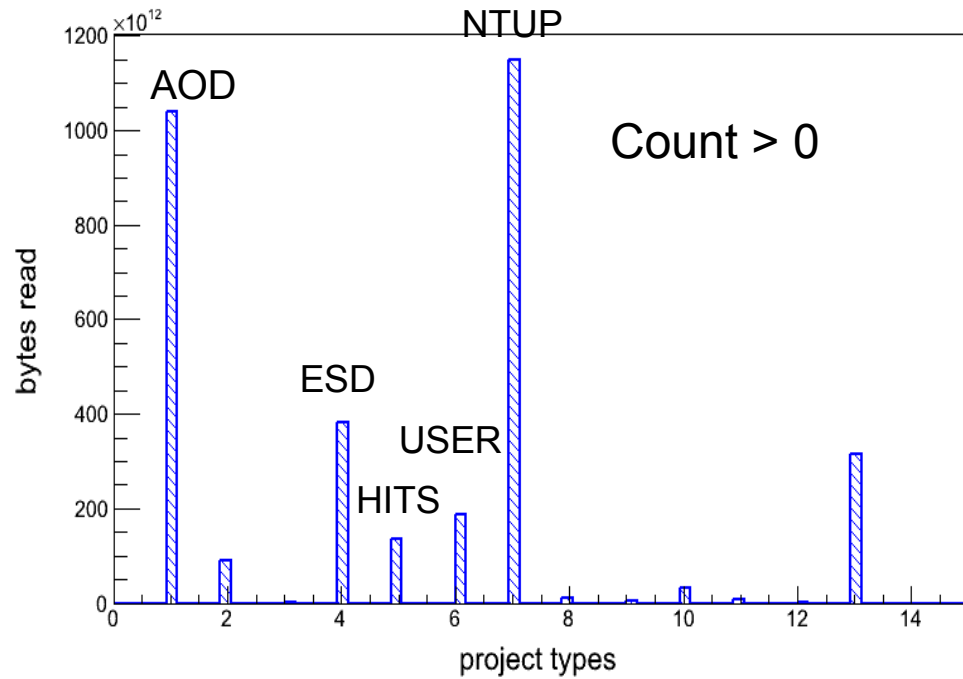
# Type of accessed data?



USER and NTUP are most accessed in terms of number of access.

# Type of accessed data
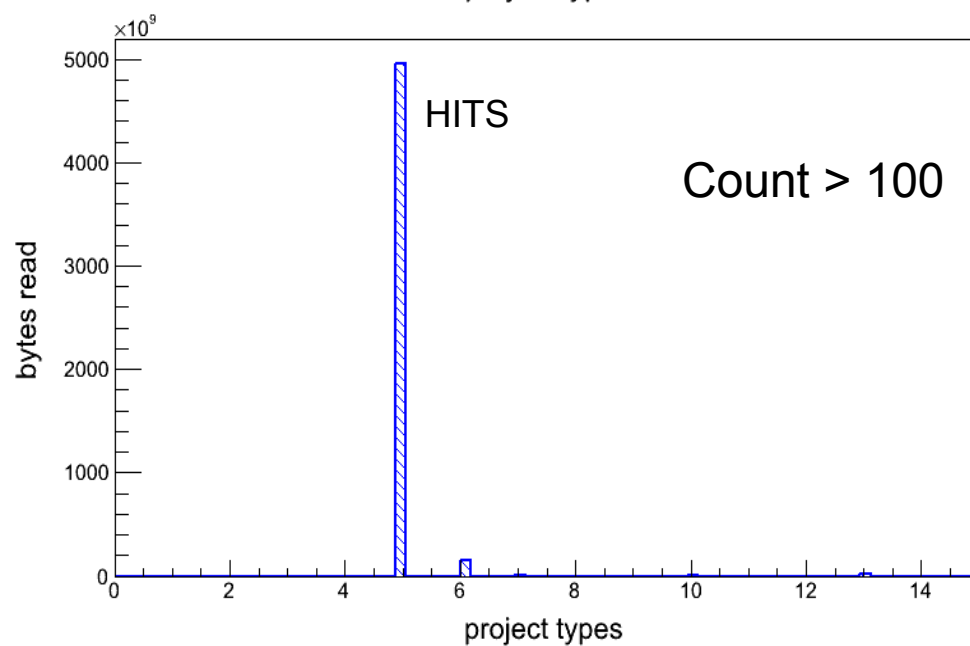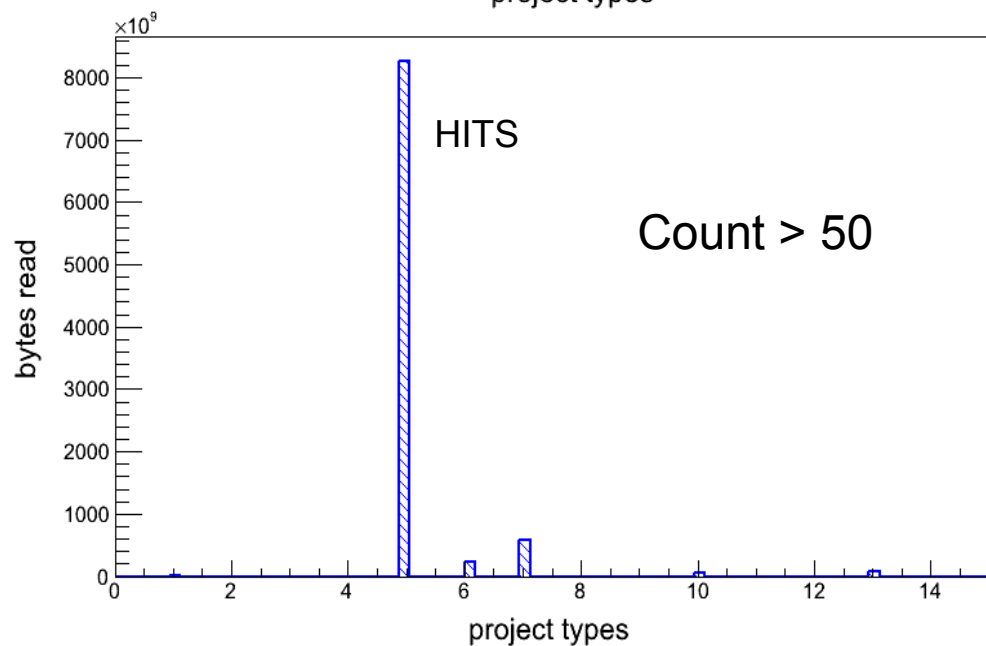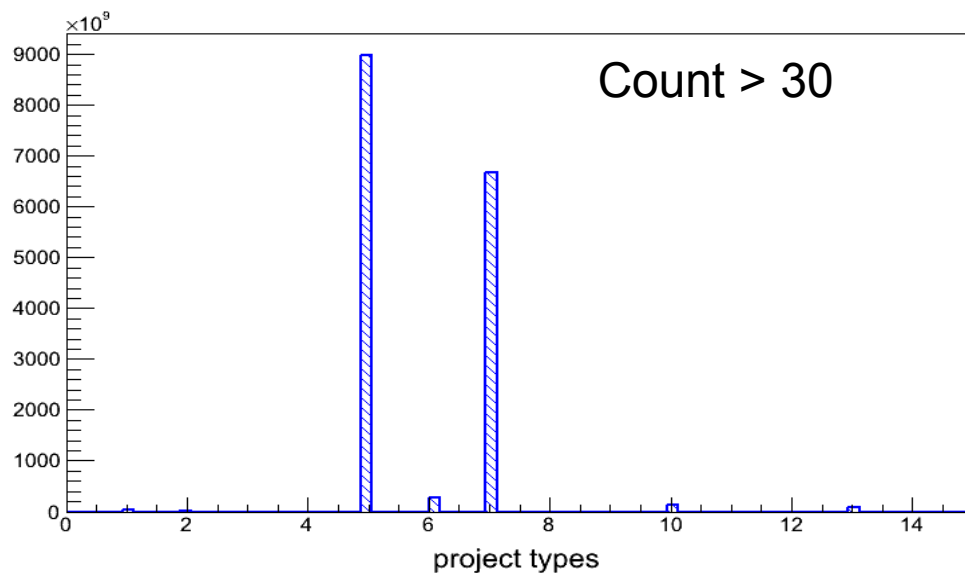# June 2013

# Variation of data types
# with different number of reuse counts
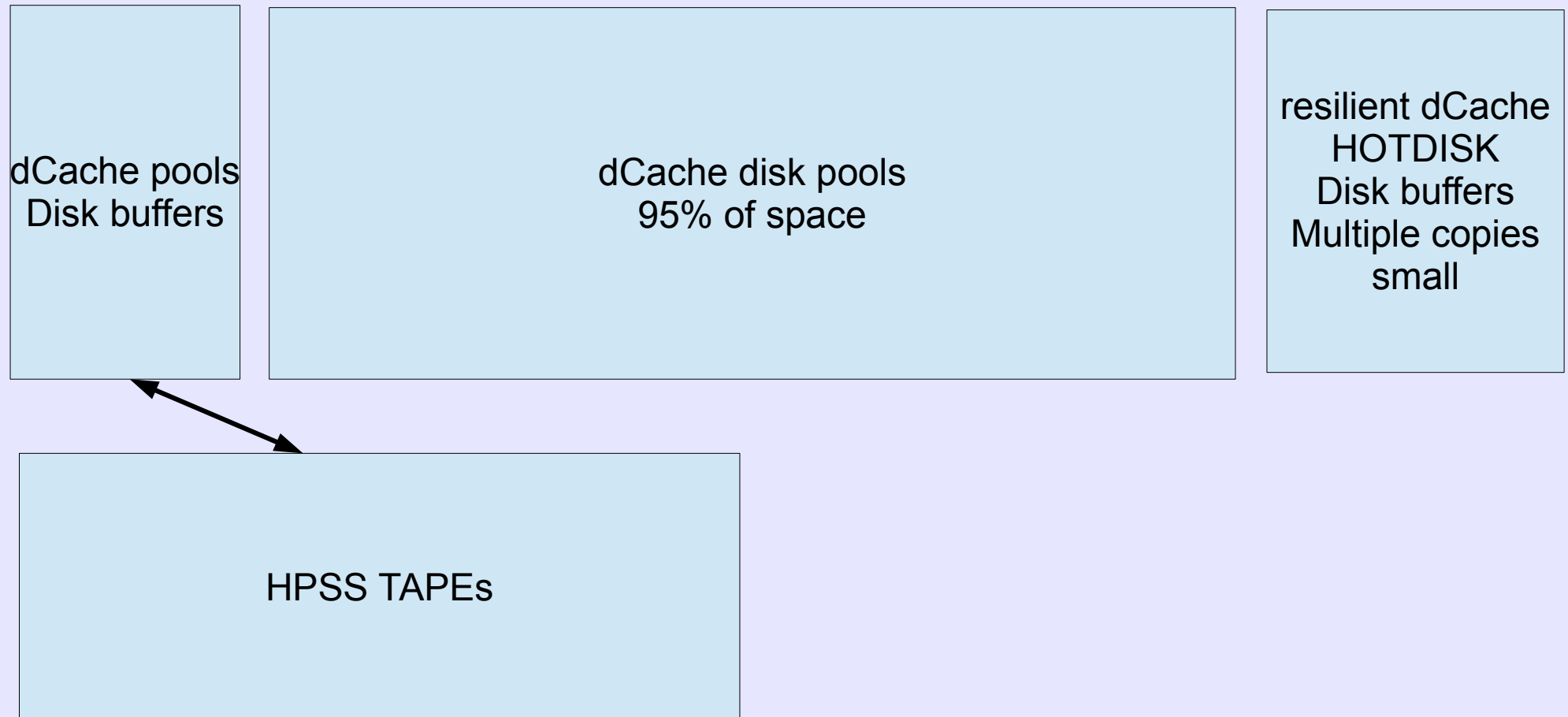
# Variation of data types with different number of reuse counts

- Only a few types of data requires storage for reuse.
  - NTUP has the largest fractions
  - For high reuse counts, HITS are most often used.

# Current BNL storage

## dCache

| dCache pools Disk buffers | dCache disk pools 95% of space | resilient dCache HOTDISK Disk buffers Multiple copies small |

HPSS TAPEs

# Possible modification to storage

## dCache

| dCache pools<br>Disk buffers | dCache disk pools<br>Resilent/distributed storage<br>Ceph/Mapr/Hadoop/etc... | High Performance disks<br>SSDs |
|---|---|---|

| HPSS TAPEs | Low performance disks<br>Resilient/distributed storage<br>Ceph/Hadoop |
|---|---|