# MODELS AND ALGORITHMS FOR BURST BUFFERS IN HPC SYSTEMS

Lionel Eyraud-Dubois, Olivier Beaumont, Guillaume Aupy

**IO congestion in HPC systems:**

- ▶ HPC applications are generating lots of data for PFS.
- ▶ Idea is to use a buffer when the I/O bandwidth is fully occupied
- ▶ The buffer can be emptied at a later time.

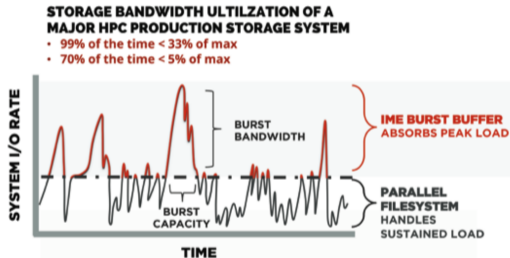Note: there are other uses of Burst-buffers



**STORAGE BANDWIDTH ULTILZATION OF A MAJOR HPC PRODUCTION STORAGE SYSTEM**
- 99% of the time < 33% of max
- 70% of the time < 5% of max

Figure: Burst-buffers to absorb IO peaks (DDN material)

Historically, Burst-Buffers were attached to IONodes (ION), used as buffers when the I/O Bandwidth was not enough (Gordon@SDSC).

But many other possible uses:

- For temporary data that may not be needed (e.g. fault-tolerance)
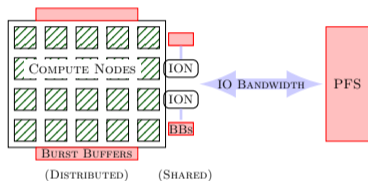- For intermediate data (e.g. BigData on HPC machine, In-situ/In-transit)
- For other uses?

Historically, Burst-Buffers were attached to IONodes (ION), used as buffers when the I/O Bandwidth was not enough (Gordon@SDSC).

But many other possible uses:

- ► For temporary data that may not be needed (e.g. fault-tolerance)
- ► For intermediate data (e.g. BigData on HPC machine, In-situ/In-transit)
- ► For other uses?

**How do we *design* and *dimension* our Burst-Buffer architecture depending on usage?**
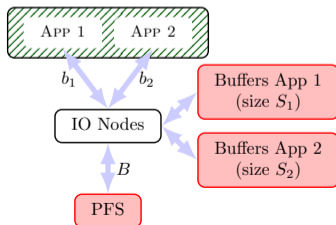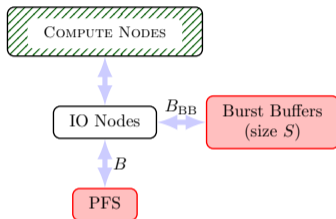
(Distributed)    (Shared)

Application Modeling:

► Compute and I/O behavior, buffer needs?

► Performance model of application?

► For both HPC and BigData applications

Algorithm design:

► I/O Scheduling, Data placement, Buffer sharing

► Dimensioning: what size / bandwidth / parameters?

► Explore different designs: Distributed vs Shared, Static vs Dynamic, ...

Shared buffers for I/O management:

- What bandwidth $B_{\mathrm{BB}}$?
- What size $S$?
- What filling/emptying policy?

Aupy, Beaumont, Eyraud-Dubois, What size should your Buffers to Disks be?, IPDPS'18

Static versus dynamic buffer sharing

- What buffer size to hide congestion?
- What overhead of static allocations?

We consider a unit time characteristic of the system.

**Machine** is characterized by:

- The Burst Buffer size $S$
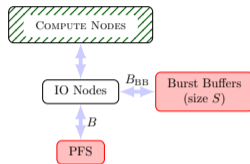- Its expected IO load: $\text{EXPECTEDLOAD} = \sum_i p_i b_i$;
- Its bandwidth to PFS: $B$

**Applications:** At any time unit, application $\mathcal{A}_i$ sends data:

- with probability $p_i$
- at bandwidth $b_i$.

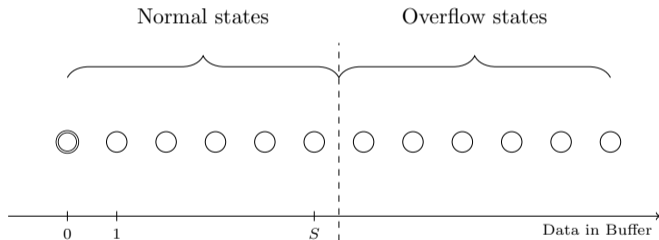$X_i$: random variable indicating whether $\mathcal{A}_i$ is sending I/Os.
$\rightarrow X_i = 1$ with proba $p_i$ and 0 with $1 - p_i$.
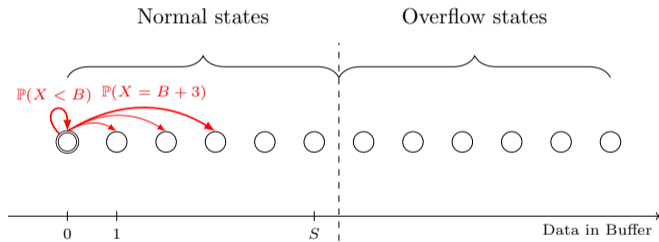
$$\text{Instant bandwidth } X = \sum_i b_i X_i$$

**Platform model:** when buffer full, stall all applications for one time unit

Normal states          Overflow states
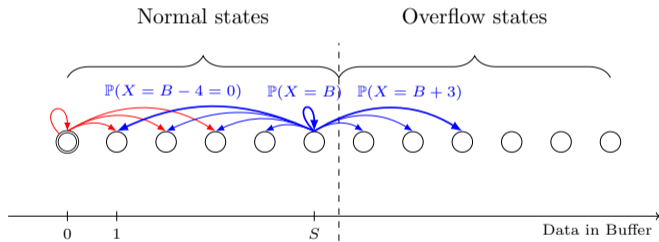


0   1                  S                          Data in Buffer

**Platform model:** when buffer full, stall all applications for one time unit

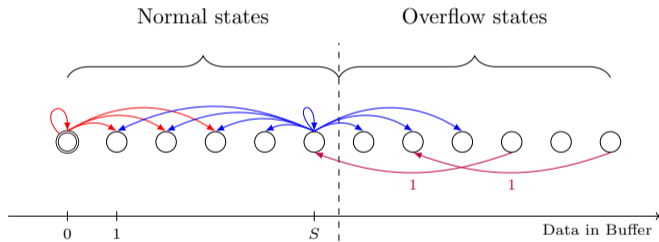**Platform model:** when buffer full, stall all applications for one time unit

**Platform model:** when buffer full, stall all applications for one time unit

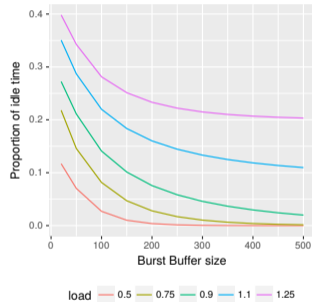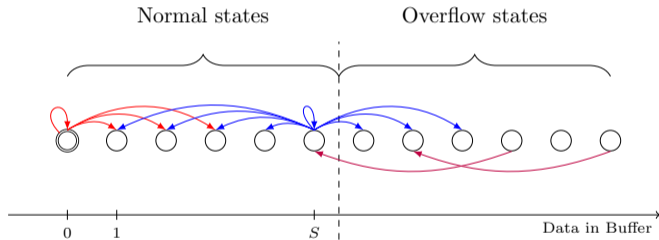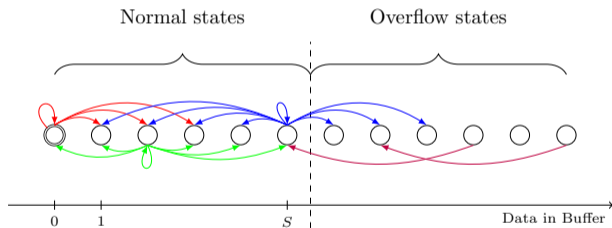**Platform model:** when buffer full, stall all applications for one time unit



Normal states      Overflow states

Data in Buffer

0   1       $S$



### Results

Can compute steady-state idle time for a given buffer size $S$

**Lazy Emptying [Cluster 2017]:**
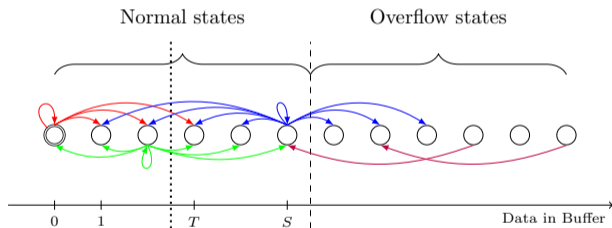Only empty the burst buffer when its load reaches a
threshold $T$ .

**Lazy Emptying [Cluster 2017]:**
Only empty the burst buffer when its load reaches a
threshold $T$ .



Normal states       Overflow states

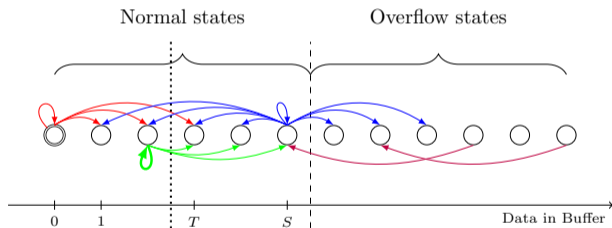0   1     $T$     $S$         Data in Buffer

**Lazy Emptying [Cluster 2017]:**
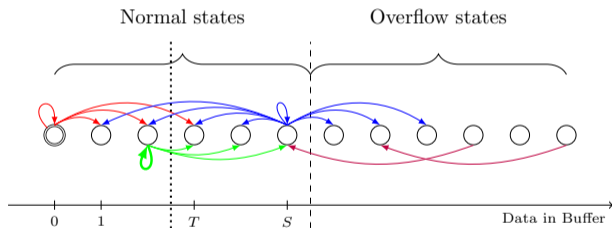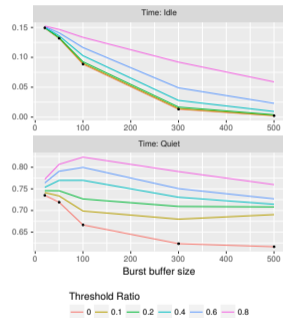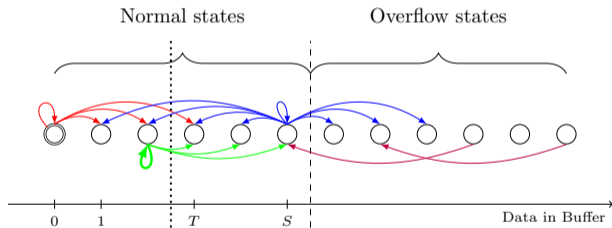Only empty the burst buffer when its load reaches a
threshold $T$ .

**Lazy Emptying [Cluster 2017]:**
Only empty the burst buffer when its load reaches a
threshold $T$ .



Normal states     Overflow states

0   1     $T$     $S$       Data in Buffer

**Lazy Emptying [Cluster 2017]:**
Only empty the burst buffer when its load reaches a threshold $T$.



Normal states          Overflow states

$0 \quad 1 \qquad T \qquad S$          Data in Buffer



Time: Idle

Time: Quiet

Burst buffer size

Threshold Ratio
— 0 — 0.1 — 0.2 — 0.4 — 0.6 — 0.8

### Results

Threshold ratio around 20-40% seems reasonable

**Application model:**

- divided into read-compute-write phases
- offline model: all data known in advance
- release dates

**Machine model:**

- Applications run independently, share the bandwidth $B$
- Each application communicates with bandwidth $b_i$
- Burst buffer is statically allocated

**Questions:** (solved with Linear Programming formulations)

- Buffer size to optimize an application by itself
- Additionnal buffer size to hide congestion

**Topics of interest**
- ▶ Burst buffer modeling and design
- ▶ Algorithms for dimensioning and/or scheduling

**Critics** and **suggestions** welcome!
- ▶ Interested in other people's view of Burst Buffers

## Questions?