

# jump trading

Slurm User Group – September 2024  
Matthieu Hautreux, Larry Pezzaglia



# Jump Trading

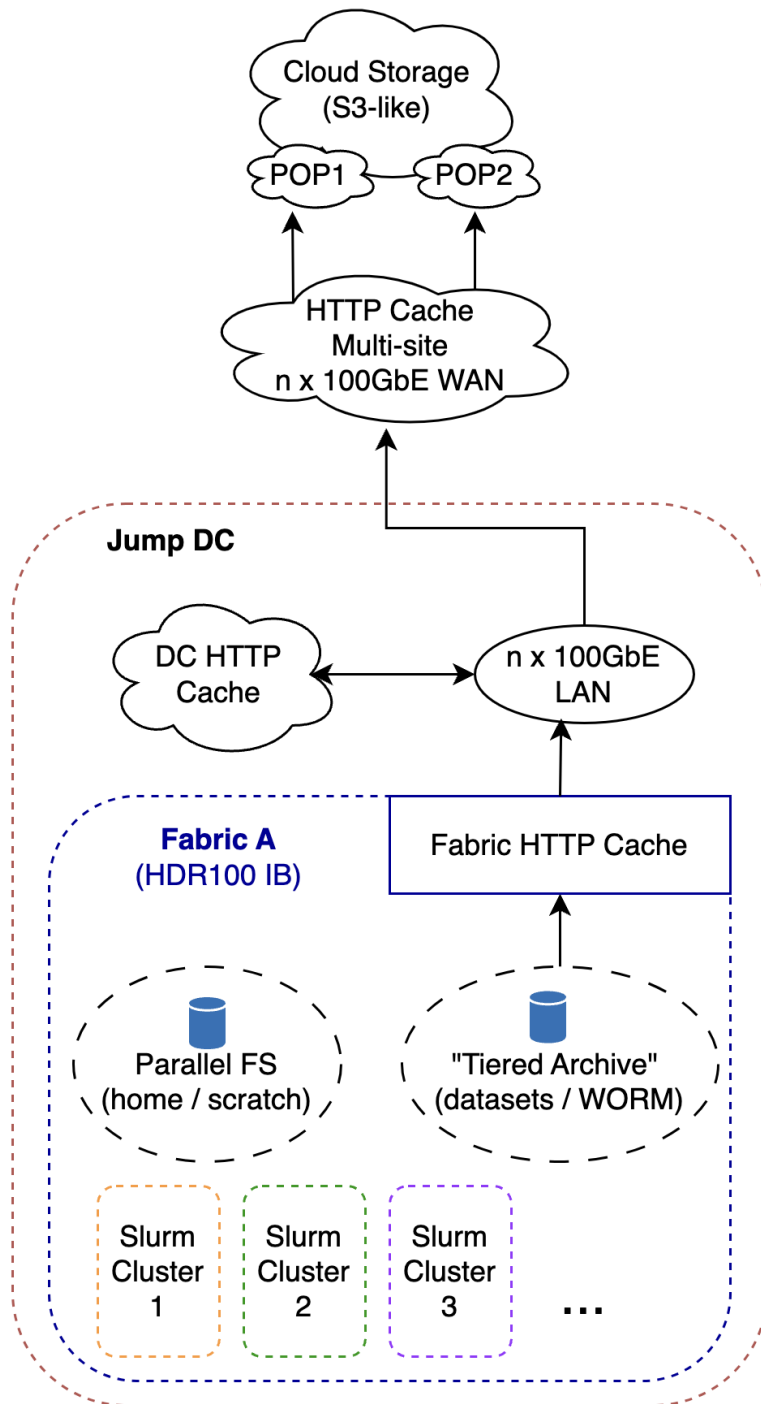
- Jump Trading is a proprietary trading firm committed to world-class research
- Full migration to Slurm in 2023
- Enabling research at scale requires automation, flexibility, self-service systems, and a mindset of constant improvement

# Slurm at **Jump Trading**

## Agenda

- **Jump Trading HPC overview**
- Workload Characteristics and Slurm migration
- Slurm features sponsored by Jump
- Slurm enhancements for high job throughput
- What's next ?

# Jump HPC Fabric

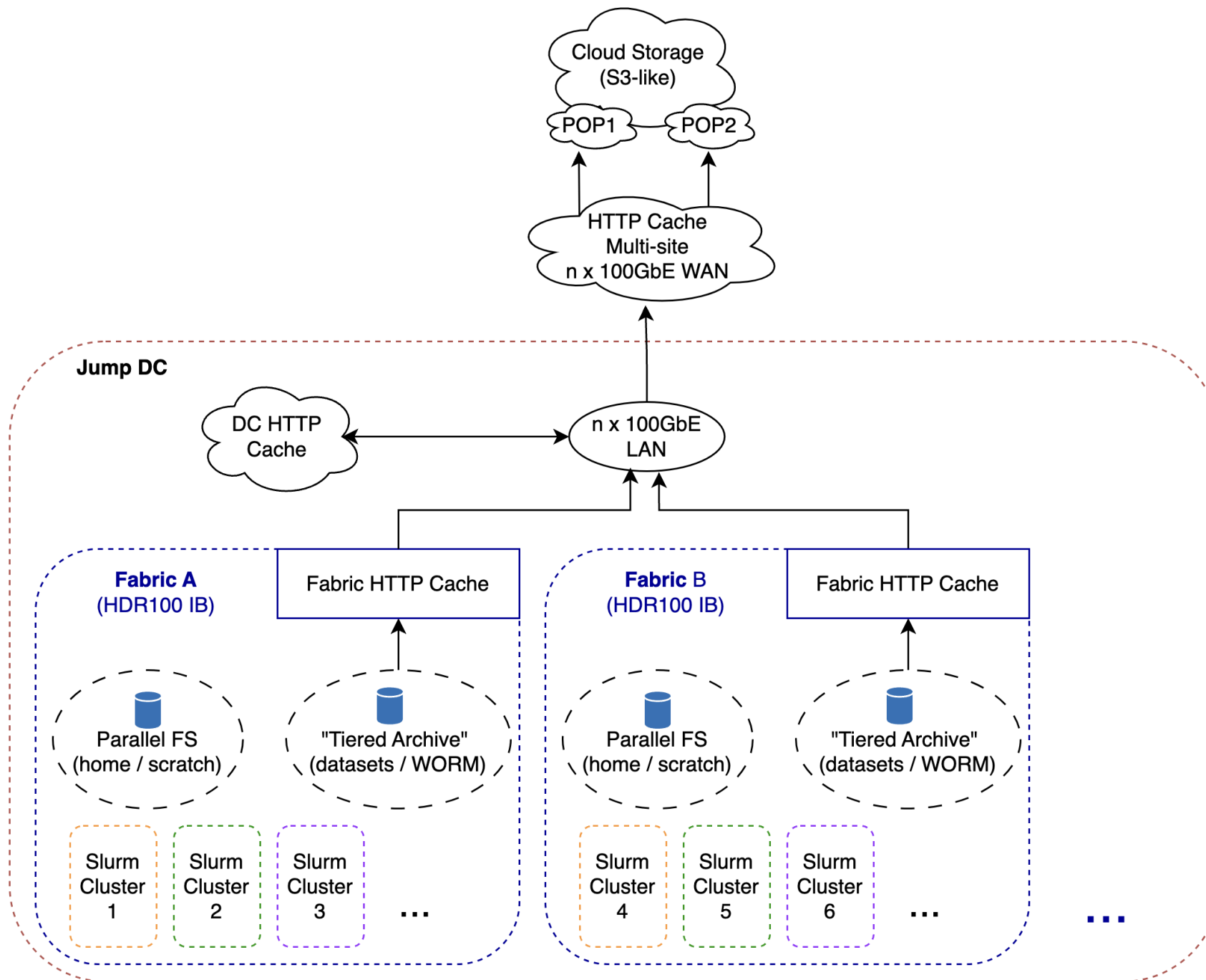


- Textbook HPC components
  - RDMA-capable fabric
  - Parallel filesystem
  - Workload manager (Slurm)
- Add: Many Slurm clusters (1 or 30+ per fabric)
  - Separate clusters per internal team or project
  - Dynamically resizable by users via self-service automation
- Add: “Tiered Archive” storage system<sup>1</sup>
  - Read-only filesystem presentation (CVMFS)
  - Backed by HTTP caches and cloud storage
  - rsync-like write interface for users

[1] <https://indico.cern.ch/event/1377701/contributions/5863778/attachments/2837970/4959865/HEPIX-CVMFS-Presentation-JUMP.pdf>

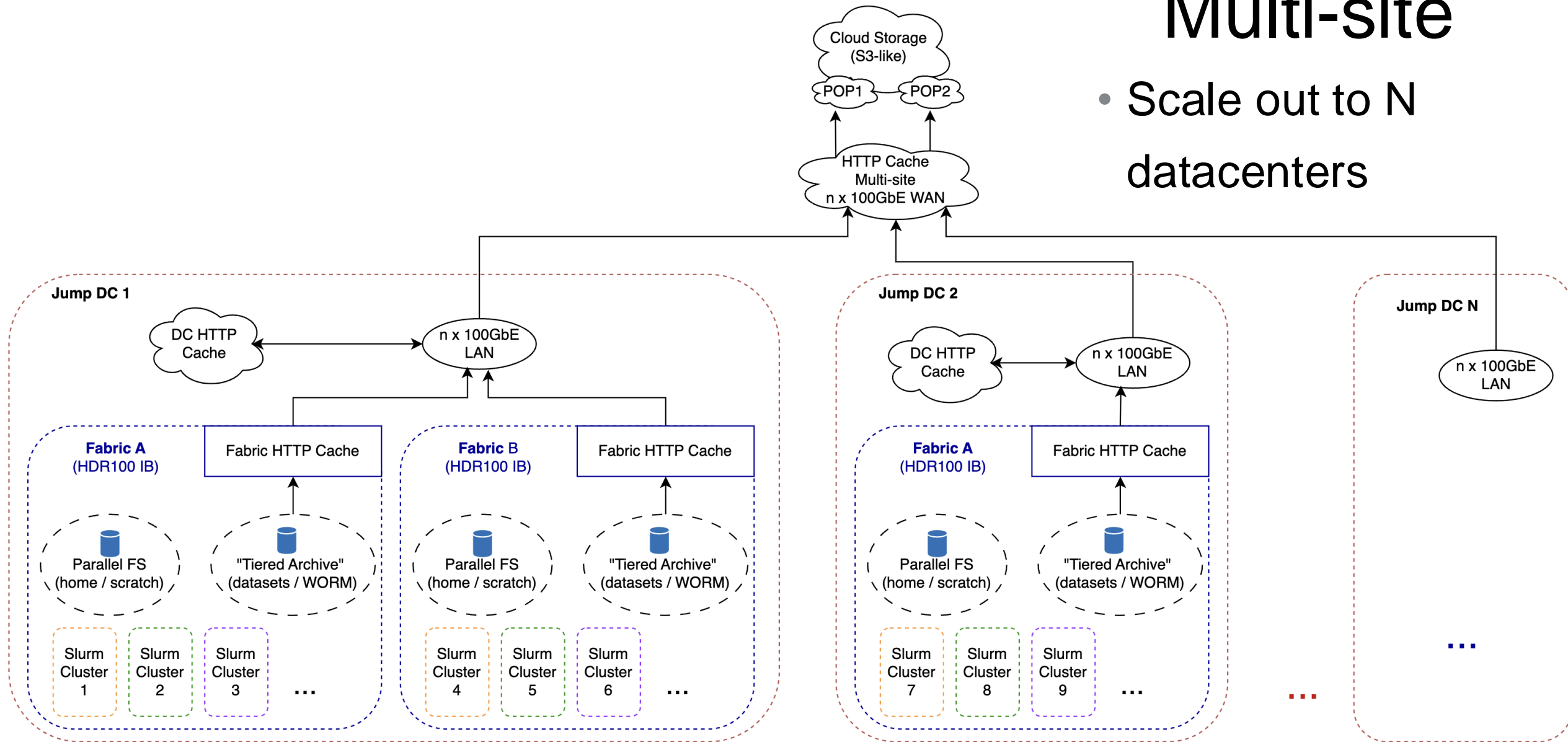
# Multi-fabric

- Scale out to N fabrics in a DC
- Contains blast radius of fabric and parallel FS issues
- Data sharing among fabrics via Tiered Archive only



# Multi-site

- Scale out to N datacenters



# Slurm at Jump Trading

## Agenda

- Jump Trading HPC overview
- **Workload Characteristics and Slurm migration**
- Slurm features sponsored by Jump
- Slurm enhancements for high job throughput
- What's next ?

# Workload Characteristics

- **Most** jobs are short, **single-core**, and data-intensive
  - Nodes run many tasks from many different jobs and users
  - Expectation of “immediate” task starts when resources are available
  - High Slurm throughput required to keep clusters full
- **Most** pipelines orchestrated by in-house workflow generator
  - Complex workflows, thousands of work units, submitted in job arrays
  - Abstracts away Slurm details
  - Performance sensitive to delta between submit time and start time



# Workload Characteristics

## **Many** exceptions and complications

- Complicated resource requests (both node-local and cluster-wide), including GPUs (both partial and exclusive access)
- Large jobs (including multi-node) versus small, short jobs
  - Some large jobs use MPI
  - Want to avoid starvation of large jobs by small jobs
- Deadline-sensitive pipelines (e.g., must finish by midnight)
- Some “power users” prefer direct Slurm access

# Slurm Migration Playbook

- Provide teams with a Slurm cluster alongside the incumbent batch system's cluster
- Teams grow (via self-service automation) their Slurm cluster and shrink the incumbent cluster, eventually to zero
  - Bump in the road? Reverse this process
  - Many issues triggered by high throughput and only surfaced at scale
- Migration completed by end of 2023

# Slurm at Jump Trading

## Agenda

- Jump Trading HPC Overview
- Workload Characteristics and Slurm Migration
- **Slurm features sponsored by Jump**
- Slurm enhancements for high job throughput
- What's next ?

# Sponsored Features

Challenge	Slurm Solution
Many dynamic logical clusters, user-sized	<b>Dynamic nodes (22.05) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• No node definitions in slurm.conf</li><li>• Pairs well with configless mode</li></ul>
Existing user practice: “soft” feature requests	<b>Preferred Features (22.05) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• --prefer: Like –constraint, but best effort</li><li>• Users tag nodes with features (Jump automation)</li><li>• Enables “spillover” (reserve N nodes for a pipeline, but allow use of other free nodes)</li></ul>
Per-task resource requirements <ul style="list-style-type: none"><li>• Many GRES (including GPU) requests previously were per-node in Slurm</li><li>• But our nodes are heterogeneous and can run tasks from many jobs</li></ul>	<b>--tres-per-task (23.02) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• Per-task requests like resources (GPUs, etc.)</li></ul>

# Sponsored Features

Challenge	Slurm Solution
Existing user practice: time-slicing GPUs	<b>GPUs <u>sharding</u> (22.05) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• Users may request full GPUs or shards of GPUs</li><li>• Some teams use MIG (Multi-Instance GPU) instead</li></ul>
Granular control over GPU shards	<b>Per-task shards (23.11) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• MULTIPLE_SHARING_GRES_PJ SelectType</li><li>• Jobs need multiple “sharing” (e.g., shards) GRES per node</li></ul> <b>Shard {anti-,}affinity (23.11) – Sponsored by Jump</b> <ul style="list-style-type: none"><li>• Require shards on same/different physical GPUs</li><li>• Toggleable per job</li></ul>

# Sponsored Features

Challenge	Slurm Solution
<p>Runtimes hard to estimate</p> <ul style="list-style-type: none"><li>• Research problem maps well to one job array task per data slice<ul style="list-style-type: none"><li>• Data slices are different sizes</li><li>• External factors (e.g., FS performance)</li></ul></li><li>• Users react by overestimating <code>--time</code> by 100x</li></ul>	<p><b>Soft Time Limits (23.11) – Sponsored by Jump</b></p> <ul style="list-style-type: none"><li>• <code>--time-min</code> used for backfill</li><li>• <code>--time</code> still used for hard time limit</li></ul>
<p>High job throughput</p> <ul style="list-style-type: none"><li>• Core counts per node keep increasing</li><li>• Need other “expensive” features (preemption, backfilling) enabled</li><li>• Avoid starvation of large jobs</li></ul>	<p>Development in progress – <b>Sponsored by Jump</b></p>

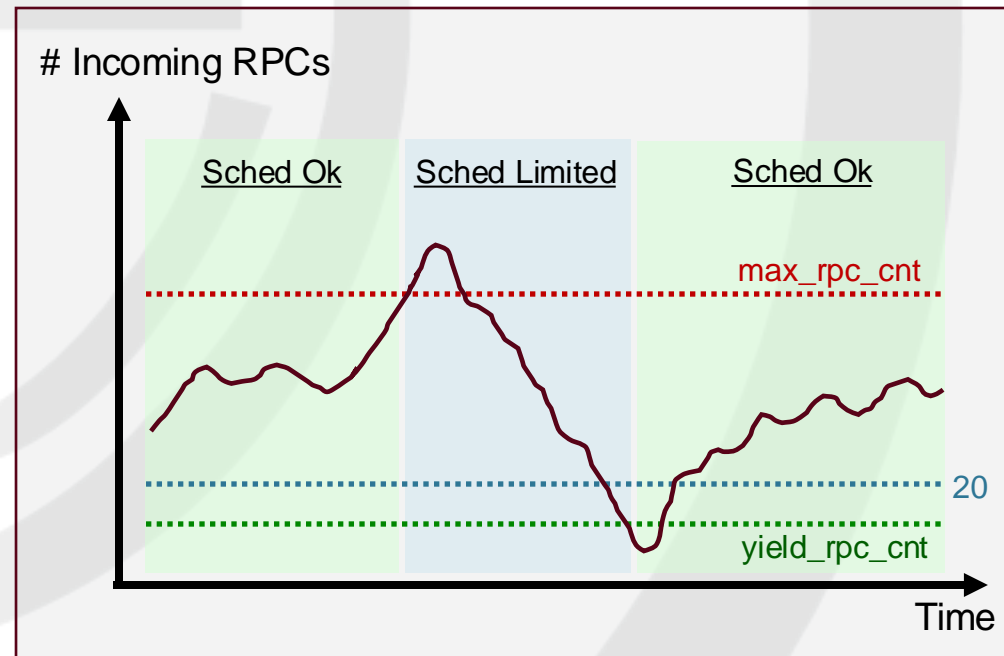
# Slurm at Jump Trading

## Agenda

- Jump Trading HPC Overview
- Workload Characteristics and Slurm Migration
- Slurm features sponsored by Jump
- **Slurm enhancements for high job throughput**
- What's next ?

# Slurm enhancements for High-throughput

- Slurm controller locks contention decrease High-throughput perf
  - Slurm Controller is highly/multi-threaded
  - Slurm internal data structures are protected by read/write locks
    - Internal Slurm components acquire the read / write locks they need
    - Internal Slurm components timing may vary depending on the scale / complexity
- Locks arbitration logic mostly based on the # of RPC
  - Schedulers (Main or backfill) stop / yield under high incoming RPC load until reaching a low level again
    - $\text{yield\_rpc\_cnt} = \text{MAX}((\text{max\_rpc\_cnt} / 10), 20)$
  - Heavy RPC processing serialized+delayed under outgoing RPC load
- This may lead to suboptimal scheduling perf under HT
  - Free/available resources not being used fast enough





# Slurm enhancements for High-throughput

- Hidden existing `enable_rpc_queue` parameter
  - Create a dedicated processing queue+thread for certain types of RPCs
  - Experimental / incomplete native logic in Slurm
  - Queued RPCs not counted in `max_rpc_count`
- Local enhancement of the `rpc_queue` logic
  - Additional RPC types managed via `rcp_queue`
  - **per-RPC-type credits|time-based concurrency** model
  - Associated tunables per RPC type
    - `max_per_cycle`, `max_usec_per_cycle`, `yield_sleep`, `interval`, `max_queued`, ...
  - per-RCP-type disabling is possible
    - (read-lock only RPCs may be better served with the native logic)

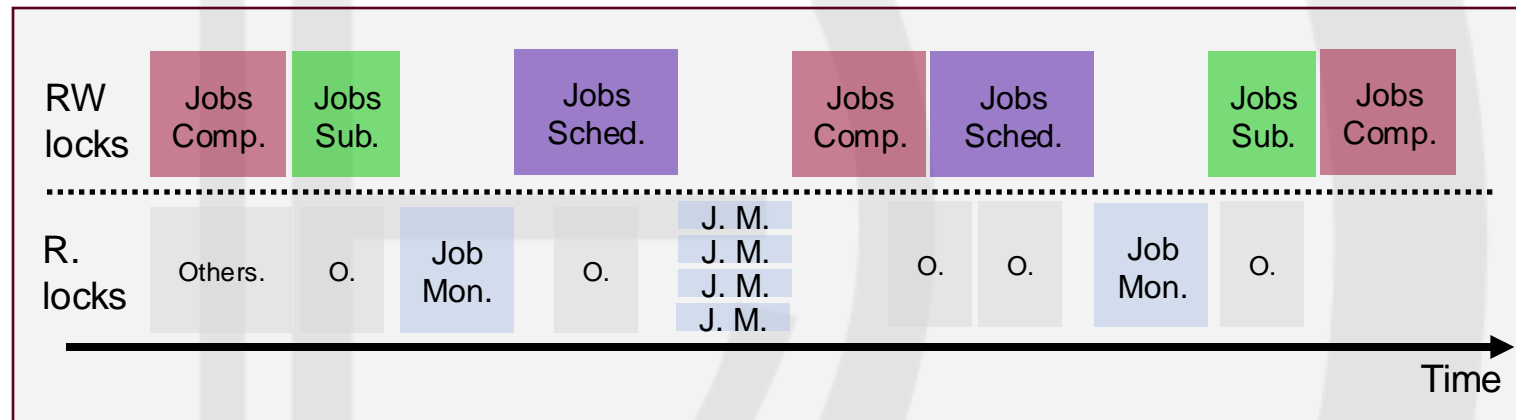
## Queue-enabled RPCs :

REQUEST\_JOB\_INFO  
REQUEST\_JOB\_USER\_INFO  
REQUEST\_JOB\_INFO\_SINGLE  
REQUEST\_FED\_INFO  
REQUEST\_NODE\_INFO  
REQUEST\_PARTITION\_INFO  
REQUEST\_COMPLETE\_PROLOG  
REQUEST\_COMPLETE\_BATCH\_SCRIPT  
REQUEST\_JOB\_STEP\_CREATE  
MESSAGE\_NODE\_REGISTRATION\_STATUS  
REQUEST\_SUBMIT\_BATCH\_JOB  
REQUEST\_STEP\_COMPLETE



# Slurm enhancements for High-throughput

- Enhanced rpc\_queue logic is not perfect but greatly helps
  - More deterministically spread locks time across the main areas of activities
    - Jobs submissions
    - Jobs states monitoring
    - Jobs scheduling
    - Jobs completion & epilog



- Queued RPCs still keep their TCP socket opened

- Need to properly configure queues tunable, especially max queue size (to trigger client backoff retries and/or abort)

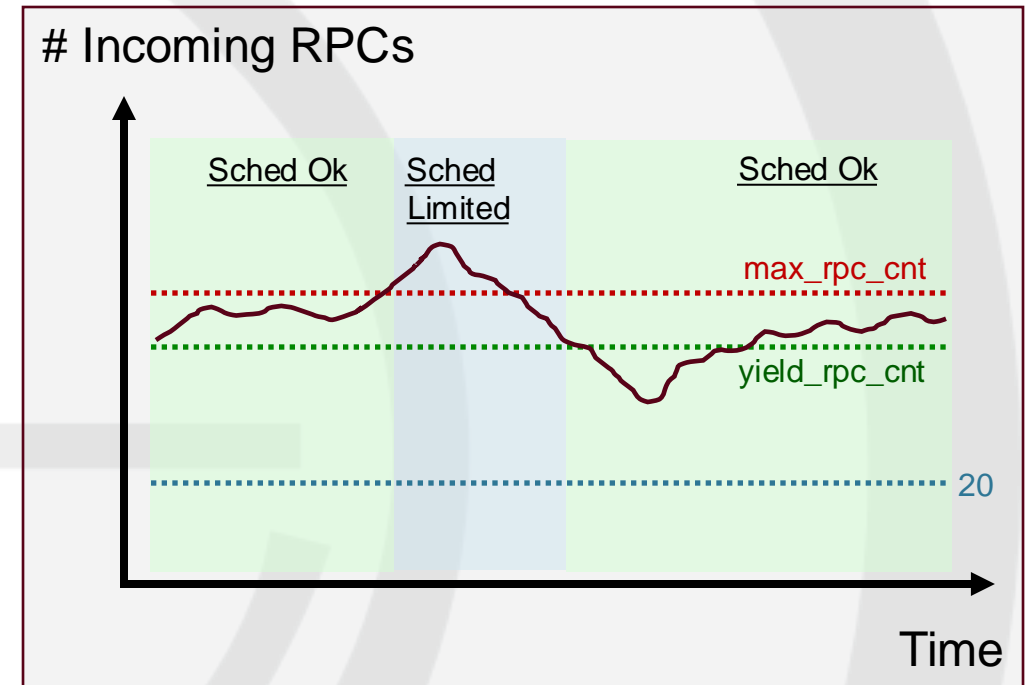
- # of RPCs still used in parallel
  - To help throttle down sched during incoming RPC peaks
    - For non-rpc\_queue aware RPCs
    - For rpc\_queue enabled RPCs while pushing them to their queue

```
slurmctld - enqueued RPCs - max processed count per cycle - by message type
```

message_type	average_microsec_per_rpc	max_per_cycle
REQUEST_JOB_INFO_SINGLE	39	4096
REQUEST_FED_INFO	14	3553
MESSAGE_EPILOG_COMPLETE	19	2493
REQUEST_COMPLETE_BATCH_SCRIPT	151	2284

# Slurm enhancements for High-throughput

- Enhancement of the sched yield / stop thresholds conf
  - Keep sched enabled under heavy yet manageable RPC load
  - Make yield\_rpc\_cnt configurable, example :
    - **max\_rpc\_cnt = 150**
    - **yield\_rpc\_cnt = 100**
      - (greater than the default max of 20)



- Other modifications have been made to help improving sched performances like
  - Reducing the amount of time required per backfill job scheduling attempt taking some shortcuts
    - May not apply to all workloads
  - Reusing free resources as soon as possible even when nodes are completing
    - to better handle large SMP nodes with single-core jobs

# Slurm at Jump Trading

## Agenda

- Jump Trading HPC Overview
- Workload Characteristics and Slurm Migration
- Slurm features sponsored by Jump
- Slurm enhancements for high job throughput
- **What's next ?**

# What's next ?

- SchedMD has integrated Jump **rpc\_queue** enhancements in **24.05**
  - Still considered experimental and hidden
  - A new `rpc_queue.yaml` file enables to configure the various per RPC type tunings
  - Other improvements added and/or ongoing by SchedMD to further improve HT in **24.05** and **24.11**
- SchedMD has reviewed some of our scheduling modifications
  - And is actively working on alternatives to further improve performances in **24.11** (part of the in-progress development sponsored by Jump)
- Additional local patches / features (not discussed in this presentation) maintenance
  - Let's discuss later if some of the following subjects matter to you
    - Automatic association of a default Account to jobs belonging to undefined users
    - Direct step-cancelation in OOM situations, `cgroup(v1)` time-sharing issues between jobs
    - Opt-in memory overcommitment in jobs to use in best-effort more memory than allocated when possible

# Slurm at **Jump Trading**

**Q&A**

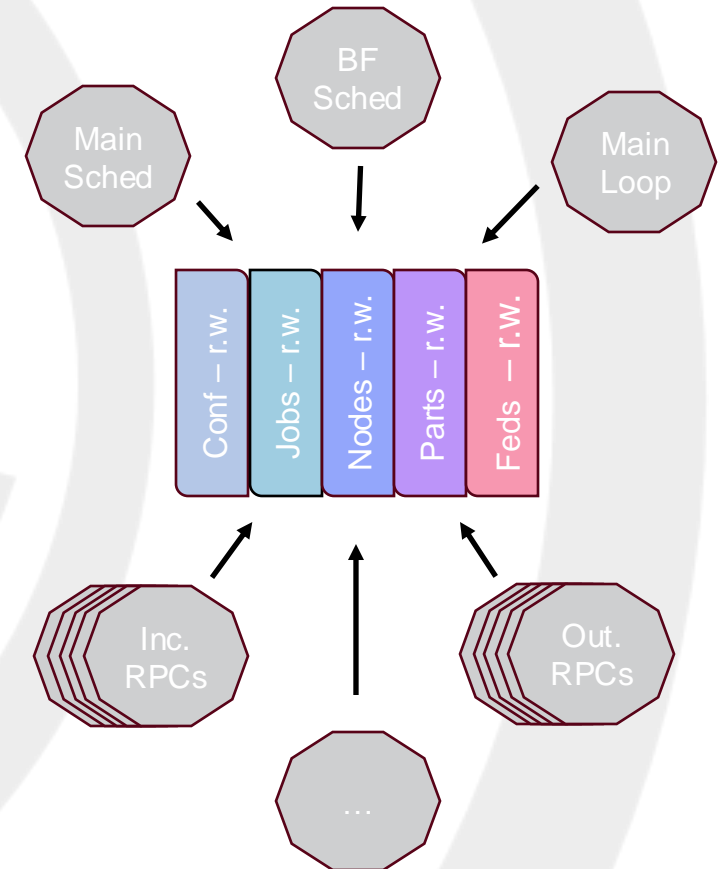
# Slurm at **Jump Trading**

## **Backup Slides**

# Slurm enhancements for High-throughput

- Context

- Slurm Controller is highly/multi-threaded
  - Incoming RPCs threads
  - Outgoing RPCs threads (agents for compute nodes actions requests)
  - Backfill scheduler thread
  - Main sched agent thread
  - Main loop/house-keeping thread
  - ...
- Slurm internal data structures are protected by read/write locks
  - Configuration / Jobs / Nodes / Partitions / Federations
  - Locks ensure that states stay consistent across the variety of handled events
- Slurm threads acquire the read / write locks they need when needed
  - Potentially preventing other concurrent components to do their job
  - Potentially being delayed / starved by other components





# Slurm enhancements for High-throughput

- Context
  - Internal Slurm components acquire the read / write locks they need
    - Potentially preventing other concurrent components to do their job
    - Potentially being delayed / starved by other components
  - Internal Slurm components timing may vary depending on
    - The size and complexity of the cluster
      - amounts of resources / generic resources / GPUs / licenses / ...
    - The number of active/recent jobs
      - in different states including running, pending, recently terminated
    - The specificities of configuration
      - # of partitions, # of various nodes weights, preemption / oversubscription options, ...

# Slurm enhancements for High-throughput

- Context
  - Some potential large locks-time consumers are also natively throttled
    - One-by-one serialization in various RPCs processing
      - REQ\_RESOURCE\_ALLOCATION, REQ\_SUBMIT\_BATCH\_JOB, REQ\_KILL\_JOB, ...
      - \_throttle\_start / \_throttle\_fini with extra usleep to help locks rotation between RPCs
        - Bigger usleep when outgoing RPCs/Agents is high (**LOTS\_OF\_AGENTS**)
    - Tend to increase the # of associated incoming RPC threads in case of burst
    - Thus, triggering the yield/stop of the schedulers
  - Other heavy locks-time consumers appear at some scales
    - RPC\_JOB\_INFO / RPC\_JOB\_USER\_INFO with (dozens of) thousands of jobs, ...

# Slurm enhancements for High-throughput

- Slurm in high-throughput mode is pushing the model to its limits
  - High submission and/or completion rates means lot of incoming RPCs
  - Schedulers yielding almost all the time
  - Not enough locks-time to schedule enough jobs
  - Global usage of resources lower than expected / possible
- Schedulers even with enough cycles, have difficulties to keep the resources fully used under high turn-over pressure

# Slurm enhancements for High-throughput

- **rpc\_queue** enhancements at Jump
  - New logic+tunings to add a **credits|time-based concurrency** model
    - For each RPC type, the following tunings can be defined :
    - **max\_per\_cycle** : # of RPCs that can be processed in a single batch
    - **max\_usec\_per\_cycle** : max amount of microsecs for a single batch
    - **yield\_sleep** : amount of microsecs to yield between 2 batches under pressure
    - **interval** : amount of microsecs to sleep between 2 batches with no pressure
    - **max\_queued** : max # of pending RPCs in the queue before backpressure
    - **hard\_drop** : indicates if backpressure should trigger retries on client side or not
  - Allow **additional RPCs** to be used with **rpc\_queue**
    - To cover observed major elements (exp: epilog completion)

# Slurm enhancements for High-throughput

- Example of sched modification for HT on large SMP nodes
  - Slurm Schedulers avoid scheduling new jobs on completing nodes
  - Large SMP nodes tend to always be completing under HT situations
  - Large SMP nodes are only eligible for jobs when not a single job is currently ending / completing
  - This can reduce global efficiency, as free cores/mem can not be used
  - **bf\_ignore\_cg\_state** to allow the backfill scheduler to schedule jobs on completing nodes