



LABIOS: A Label-Based I/O System

Anthony Kougkas, Hariharan Devarajan, Jay Lofstead*, and Xian-He Sun
Illinois Institute of Technology, *Sandia National Laboratories

Wednesday, June 26th, 2019



Highlights

Storage
Malleability

Asynchronous I/O

Resource
Heterogeneity

Data Provisioning

Storage Bridging



“The tools and cultures of HPC and BigData analytics have diverged, to the detriment of both; **unification** is essential to address a spectrum of major research domains.”

- D. Reed & J. Dongarra

Diversity of I/O Requirements

Feature	I/O Requirement	HPC	Cloud	Optimizations
Data consistency	Data passed to the I/O system must be consistent between operations.	Strong, POSIX	Eventual, Immutable	Tunable consistency
Data access	Multiple processes must be able to operate on the same data concurrently.	Shared concurrent	Multiple replicas	Collective I/O, Concurrent handlers, Complex locks
Global namespace	Data identifiers must be resolved and recognizable in a global namespace that can be accessed from anywhere.	Hierarchical, Directory, Nesting	Flat	Namespace partitioning, Client-side caching, Decoupling data-metadata, Connectors
Fault tolerance	Data must be protected against faults and errors.	Special hardware, check-pointing	Replication, Data partitioning	Erasur coding
Scale	Support for extreme scale and multi-tenancy	Few large jobs, Batch processing	Many small jobs, Iterative	Job scheduling, I/O buffering, Scale-out
Locality	Jobs are spawned where data reside.	Remote storage	Node local	Data aggregations
Ease of use	Interface, user-friendliness and ease of deployment.	High-level I/O libraries	Simple data formats	Amazon S3, Openstack Swift

LABIOS: Label-Based I/O System



- Distributed, scalable, and adaptive storage solution
- Fully decoupled architecture
- Software defined storage (SDS)
- Energy-aware enabling power-capped I/O
- Reactive storage with tunable I/O performance
- Flexible API
- Intersection of HPC and BigData



7/2/2019

Slide 5

A silhouette of a person pushing a large ball up a hill, symbolizing a challenge or struggle. The person is on the left, pushing the ball towards the right. The background is a solid blue color.

Key Challenges

- How to efficiently utilize and share I/O resources?
- How to leverage storage heterogeneity?
- How to provide an elastic storage system?
- How to support a wide range of I/O interfaces?
- How to balance energy – performance?



I/O is a logistics problem

And people are really good at this.

A simple analogy (before)

- Sending a gift

- Drive to three different retailers
- Purchase items independently
- Decide what package to use
- Decide which delivery provider
- Decide options (priority, etc)

- Performing I/O

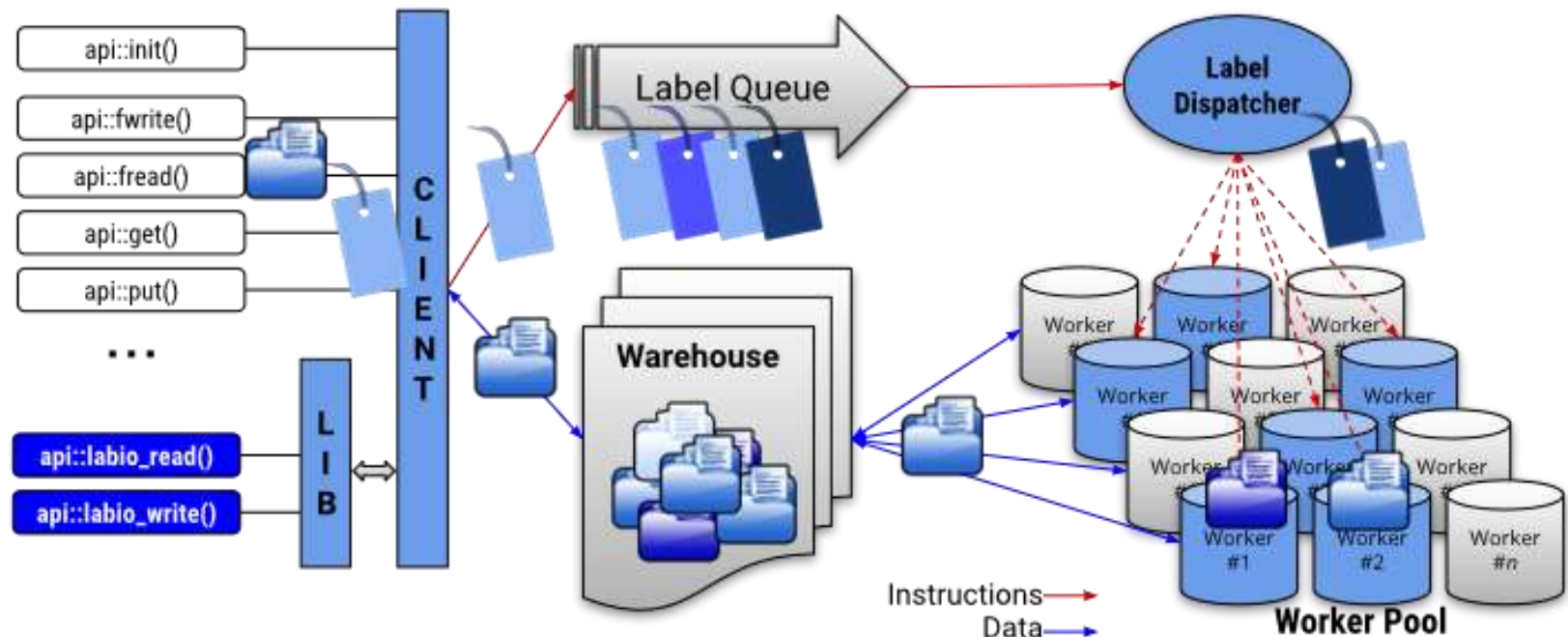
- Three different data sources
- Acquire data elements
- Data representation (file, object, etc)
- Storage device (SSD, HDD, etc)
- Storage semantics

A simple analogy (now)

- Sending a gift via an online retailer
 - Add items to cart
 - Specify details (payment info, etc)
 - Submit order
- Performing I/O with LABIOS
 - Create labels
 - Define label attributes
 - Push labels to queue

Overview

- I/O requests are transformed into a configurable unit, called a (data) **Label**.
 - A label is a tuple of an operation and a pointer to the data.
 - Resembles a shipping label following a Post Office package.
- Labels are pushed to a distributed queue.
- Data or contents are pushed into a warehouse.
- A dispatcher distributes labels to the workers.
- Workers execute labels independently (i.e., fully decoupled).



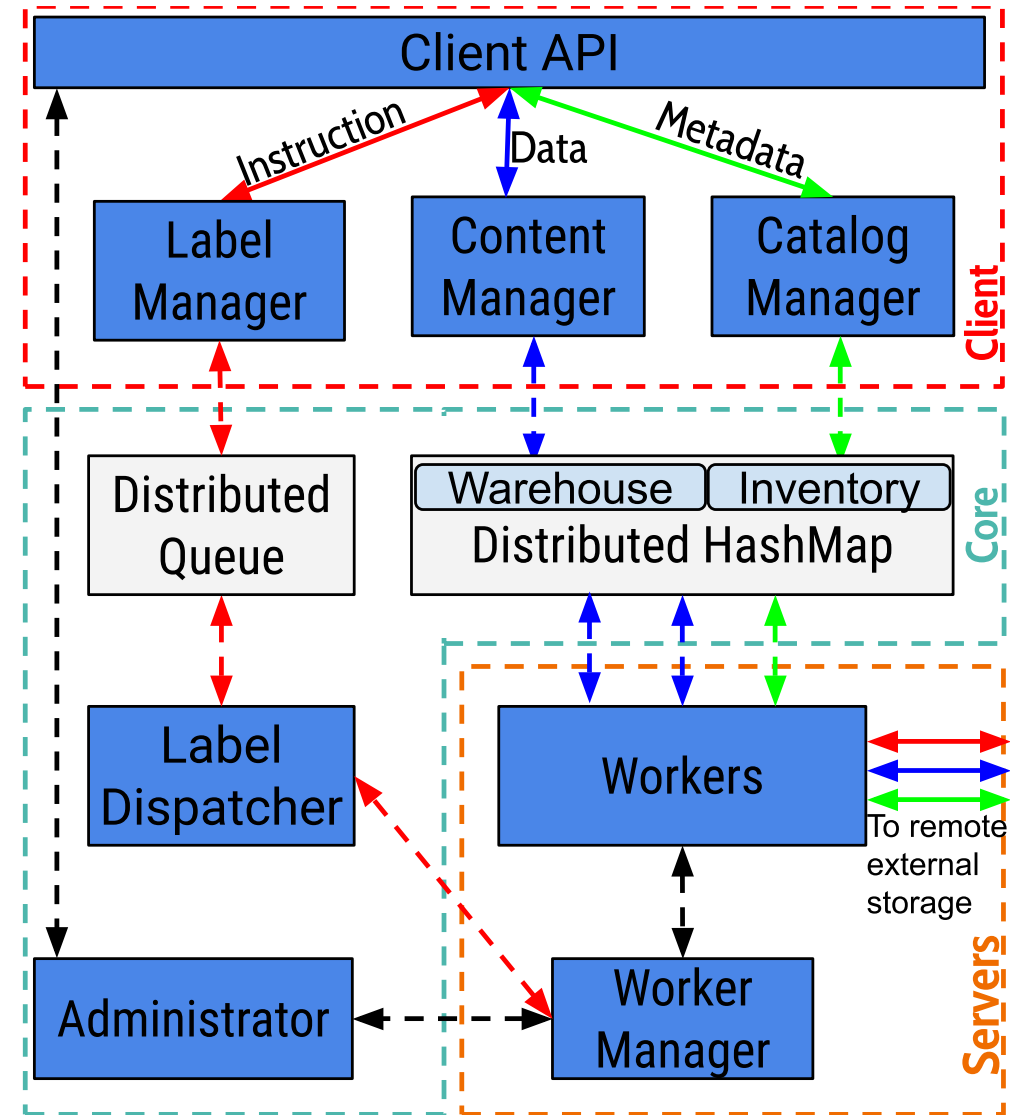
LABIOS Data Model - Labels

- Storage-independent abstraction expressing I/O intent.
- A tuple of one or more operations and a pointer to its input data.
- Exclusive to each application.
- Immutable, independent of one another, and cannot be re-used.
- Label structure includes:
 - Type
 - Unique identifier
 - Source and destination
 - memory address, file path, server IP, network port
 - Function pointer (user-defined or pre-defined)
 - all functions are store in a shared program repository
 - Set of flags indicating label's state
 - queued, scheduled, pending, cached, invalidated, etc.,



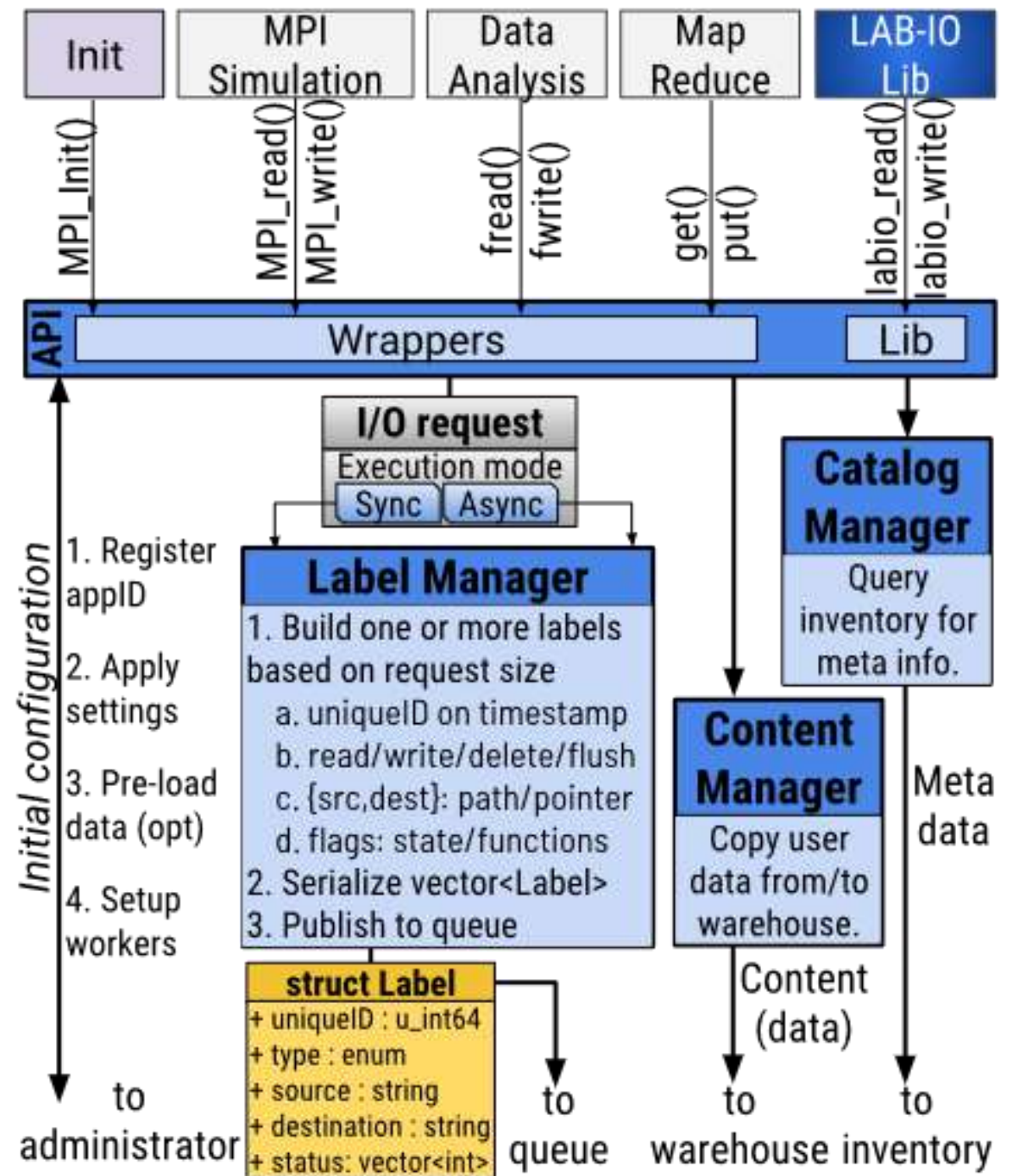
High-level Architecture

- Two main ideas:
 1. Split the data, metadata, and instruction paths.
 2. Decouple storage servers from the application.



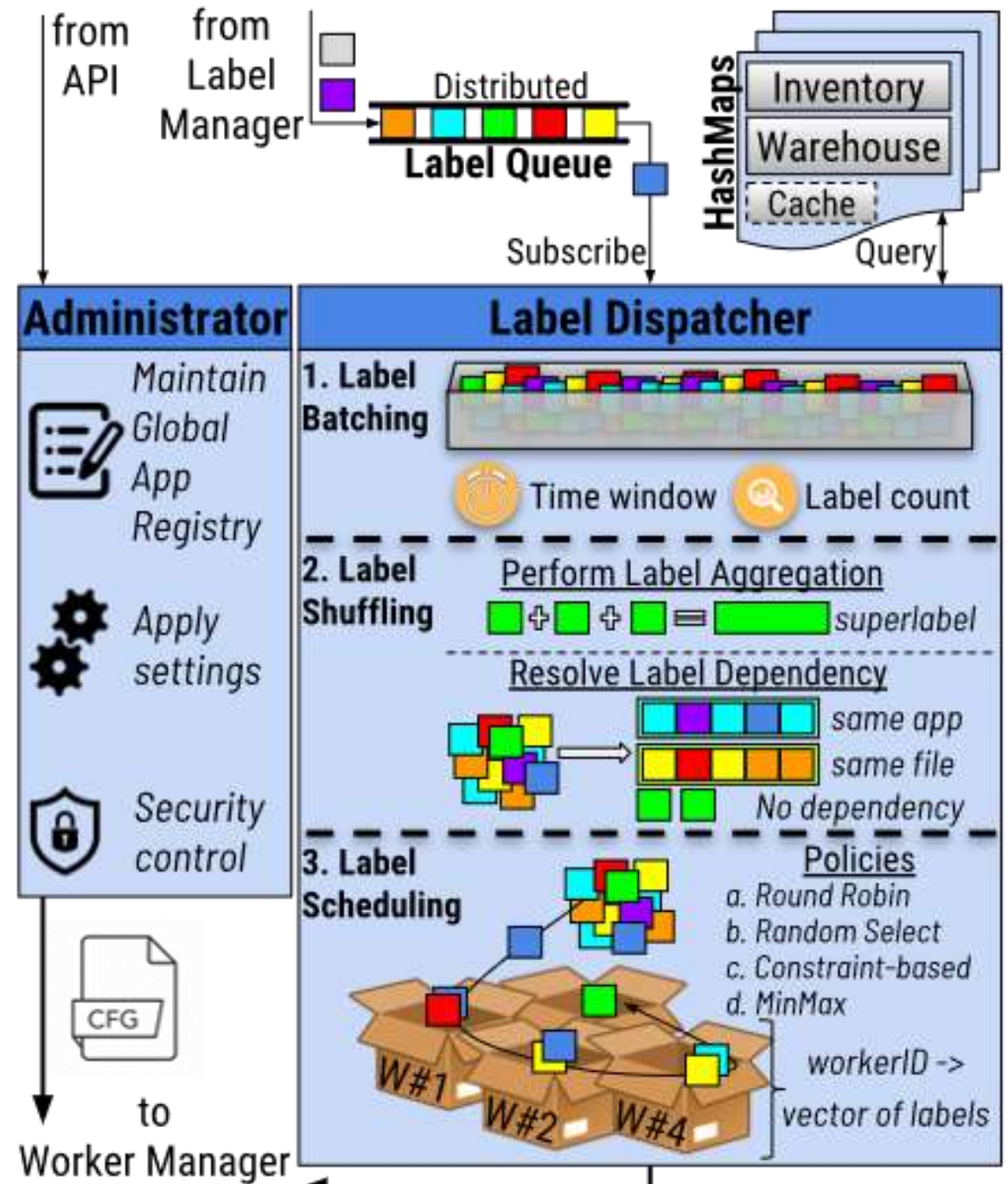
LABIOS Client

- Objectives:
 - Performs system initialization per-application.
 - Accepts application's I/O requests.
 - Builds labels based on the incoming I/O.
- Modules:
 - Label manager
 - Content manager
 - Catalog manager



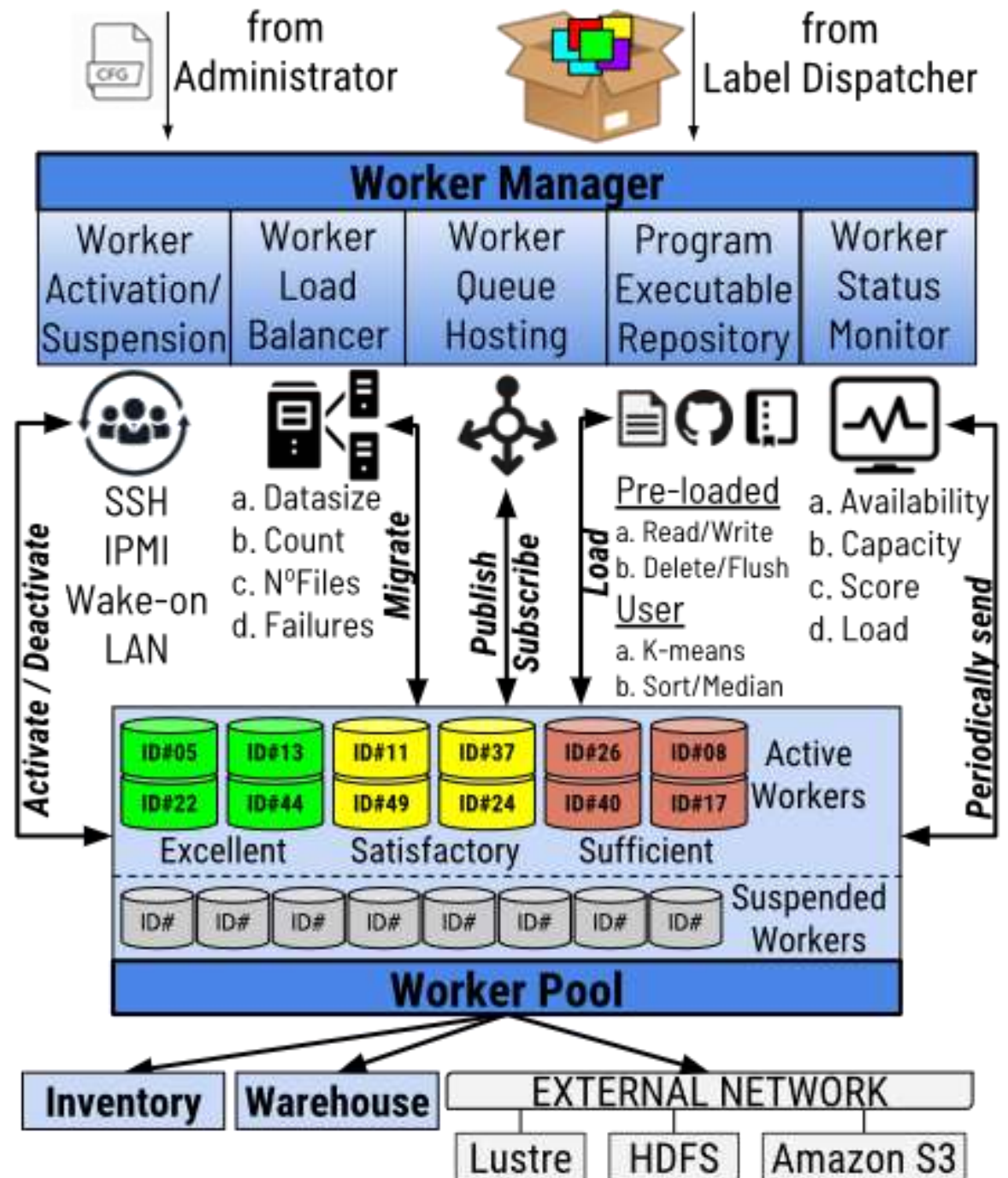
LABIOS Core

- Manages the instruction, data, and metadata flow separately.
- Distributed data structures:
 - Label queue
 - Warehouse
- Modules:
 - Administrator
 - Label Dispatcher



LABIOS Server

- Manages workers (i.e., storage servers)
- Modules:
 - Worker
 - Worker manager



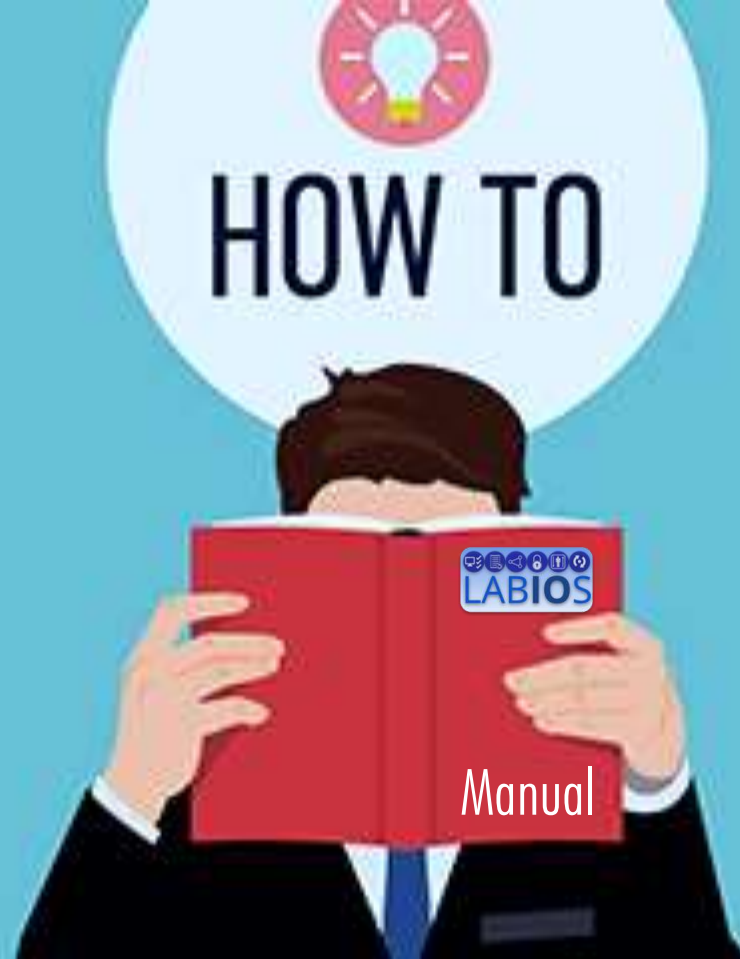
Worker

- The storage server in LABIOS
- Responsibilities:
 - service its own queue
 - execute labels
 - calculate its own worker score and send it to the worker manager periodically
 - auto-suspend itself if there are no labels in its queue for a given time window
 - connect to external storage sources
- Weighting system expresses the scheduling policy
- Final score is a double precision between 0 and 1
 - Higher score -> better worker

Variable	Value	Example
Availability	1-active, 0-suspended	1
Capacity	Double [0,1] (ratio remaining/total)	0.75
Load	Double [0,1] (ratio current/max queue size)	0.50
Speed	Integer [1,5] (grouping)	4
Energy	Integer [1,5] (grouping)	3

$$Score(workerID) = \sum_{n=1}^5 Weight_j \times Variable_j$$

Priority	Availability	Capacity	Load	Speed	Energy
Low latency	0.5	0	0.35	0.15	0
Energy savings	0	0.15	0.2	0.15	0.5
High bandwidth	0	0.15	0.15	0.70	0



LABIOS in depth

- Label Scheduling
- Deployment Models
- LABIOS API example

Label Scheduling

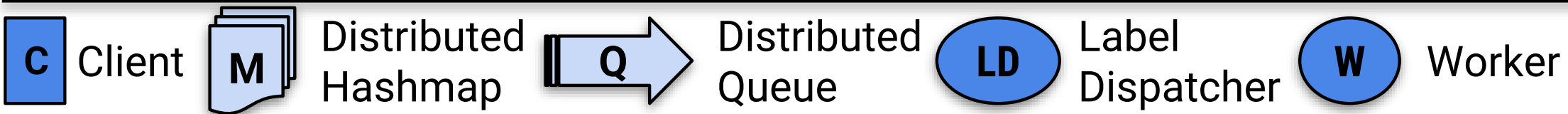
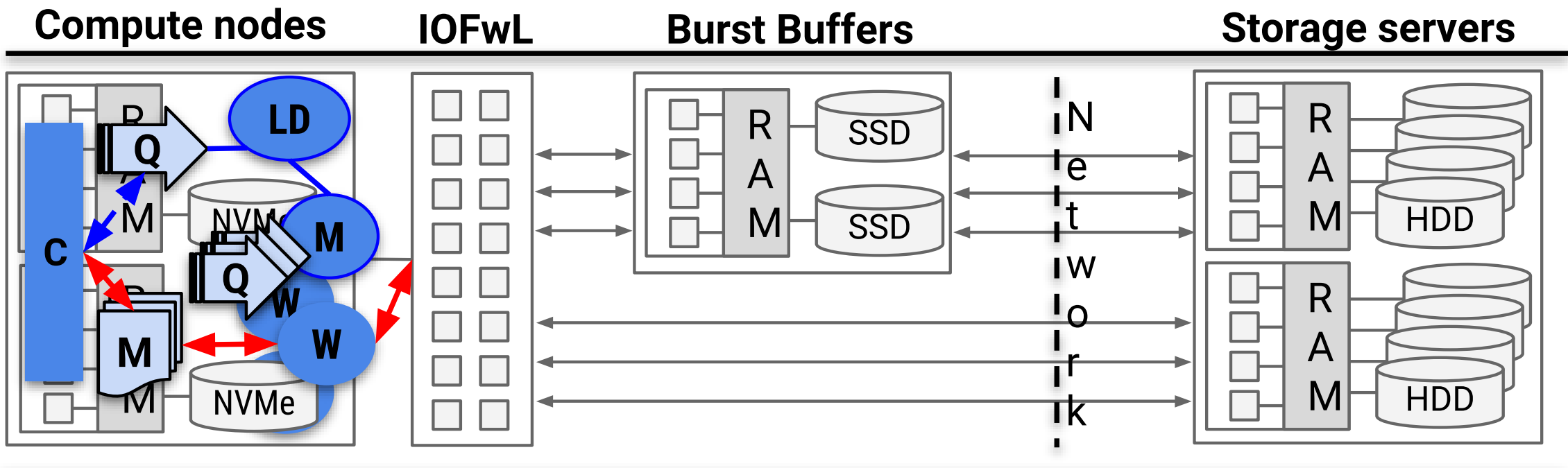
SCHEDULE

- Data distribution via scheduling policies:
 1. Round Robin
 - similar with PFS
 2. Random Select
 - randomly select workers
 3. Constraint-based (i.e., heuristically)
 - Priorities based on higher weight in the worker score
 - Availability, load, capacity, performance
 4. MinMax
 - Minimize energy consumption while maximizing performance
 - Subject to load and capacity
 - NP-hard combinatorial optimization
 - Multidimensional knapsack algorithm



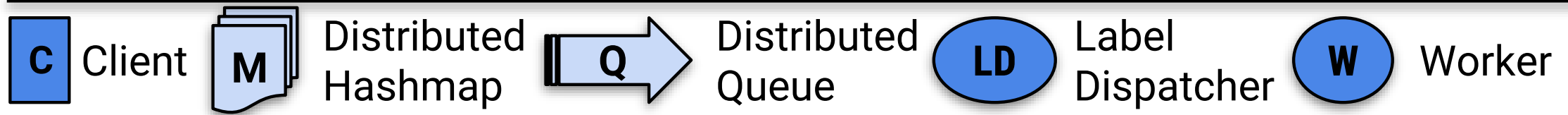
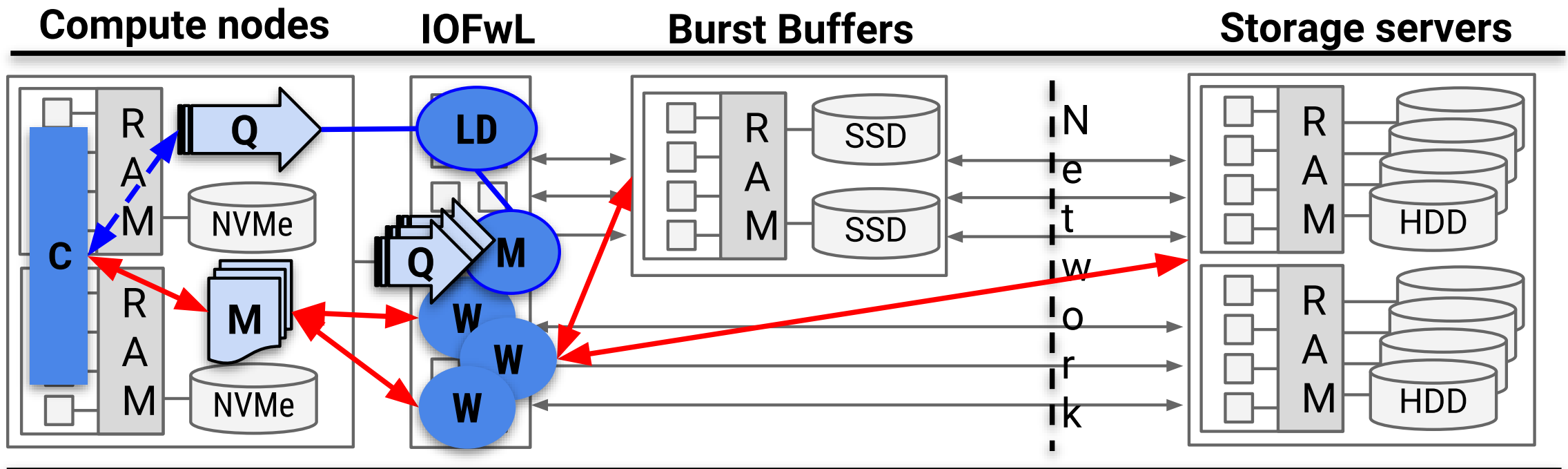
Deployment Model

- LABIOS can:
 1. replace existing distributed storage solutions
 2. be used as I/O accelerator to one or more underlying storage subsystems
- Machine model in use (motivated by the recent machines Summit in ORNL or Cori on LBNL):
 - Compute nodes equipped with a large amount of RAM
 - Local NVMe devices in each compute node
 - An I/O forwarding layer
 - A shared burst buffer installation based on SSD equipped nodes, and
 - A remote PFS installation based on HDDs



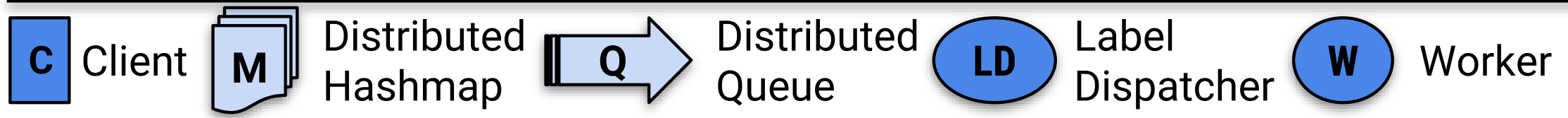
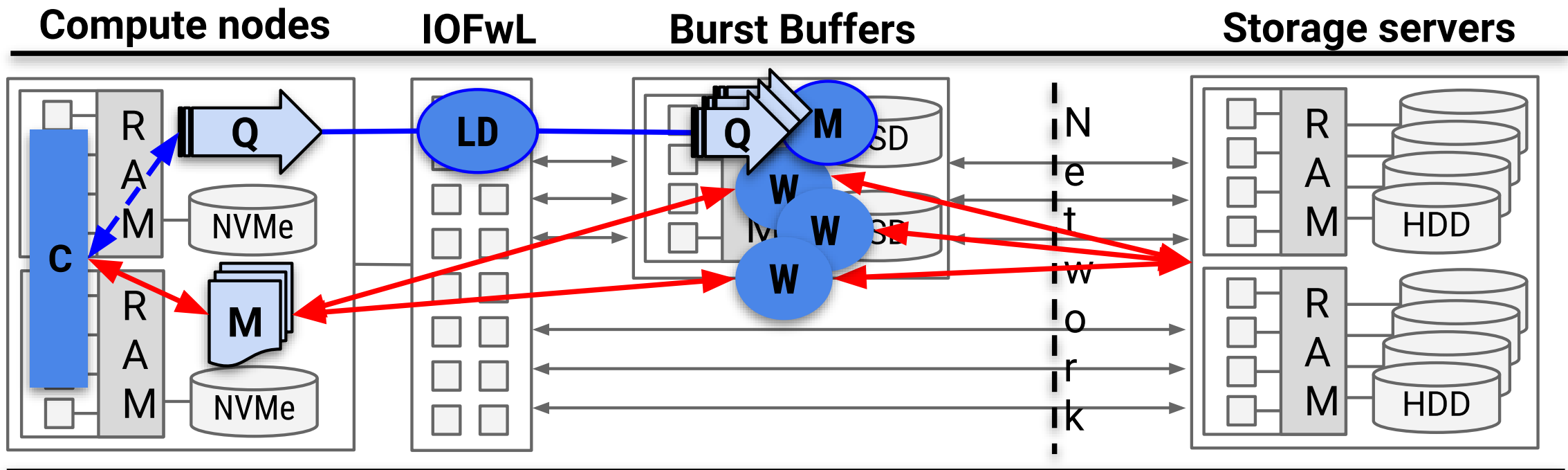
LABIOS as I/O accelerator (in compute nodes)

- Pros
 - Fast distributed cache
 - For temporary I/O
 - On top of external sources
 - Hadoop workloads with node local I/O
- Cons
 - Overheads by using compute cores to run its services
 - I/O traffic mixed with compute network



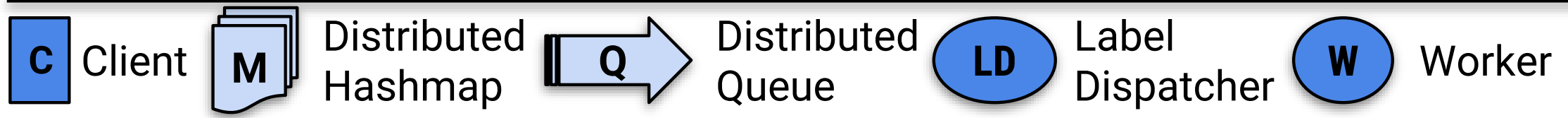
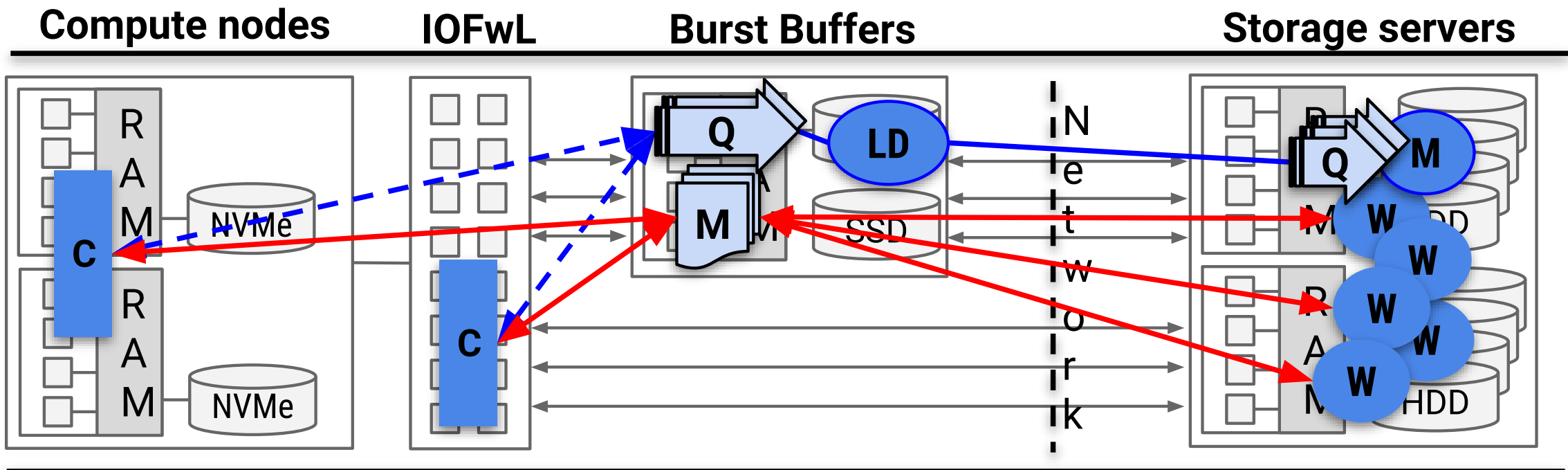
LABIOS as I/O forwarder (in ION)

- Pros
 - Asynchronous I/O
 - Non-blocking data movement
 - Connect to external storage
- Cons
 - Subject to I/O forwarding layer
 - Limited scalability



LABIOS as I/O buffering (in burst buffers)

- Pros
 - Fast scratch space
 - Data sharing between applications
 - In-situ visualization and analysis
- Cons
 - Requires additional resources (e.g., buffers)
 - Storage
 - Network



LABIOS as Remote Distributed Storage

- Pros
 - Scalability
 - Better resource utilization
 - Higher flexibility
- Cons
 - Increased deployment complexity
 - Requires systems admins

LABIOS API Example

```
1  #include <tabios.hpp>
2  ...
3  Client client = InitClient(ip, port, connConfig);
4  std::string path = "pvfs2:/data/integers.dat";
5  LabelSrc src = new LabelSrc(path, src_offset, size);
6  LabelType type = SDS_IN_SITU; //complex type
7  LabelFlag flags = CACHED | MPI_IO; //keep in cache
8  std::function<int(vector<int>)> fn = FindMedian;
9  Label label = client.CreateLabel(type,src,fn,flags);
10 Status status =client.IPublishLabel(label);
11 ... //perform other computations
12 client.WaitLabel(&status);
13 int median = std::static_cast<int>(status.data);
```

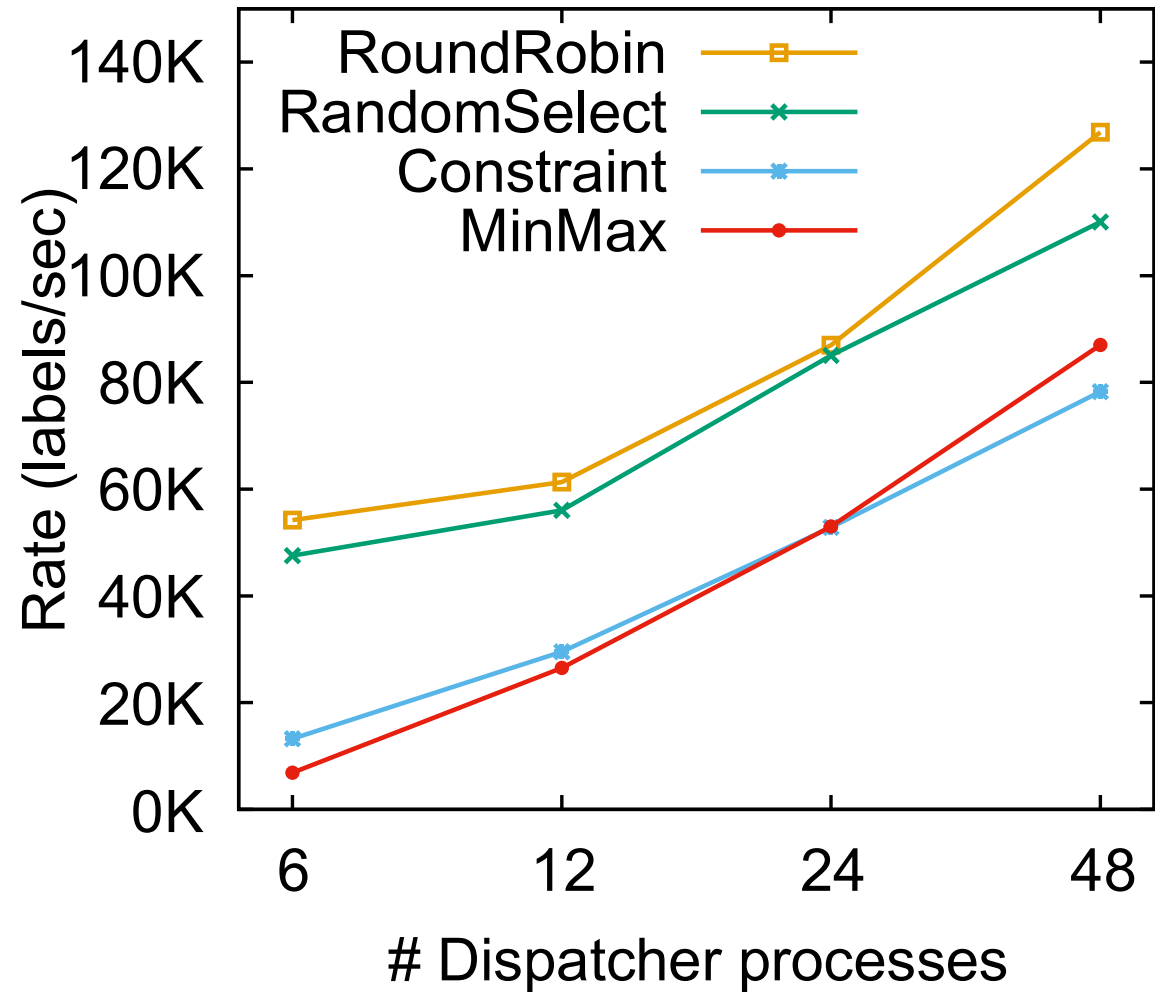



Testbed

- All experiments on bare metal on Chameleon:
 - 64 client nodes
 - 8 burst buffer nodes
 - 32 storage servers
- Cluster OS: CentOS 7.1
- PFS: OrangeFS 2.9.6
- Workloads:
 - CM1 simulation
 - HACC simulation
 - Montage application
 - K-means clustering

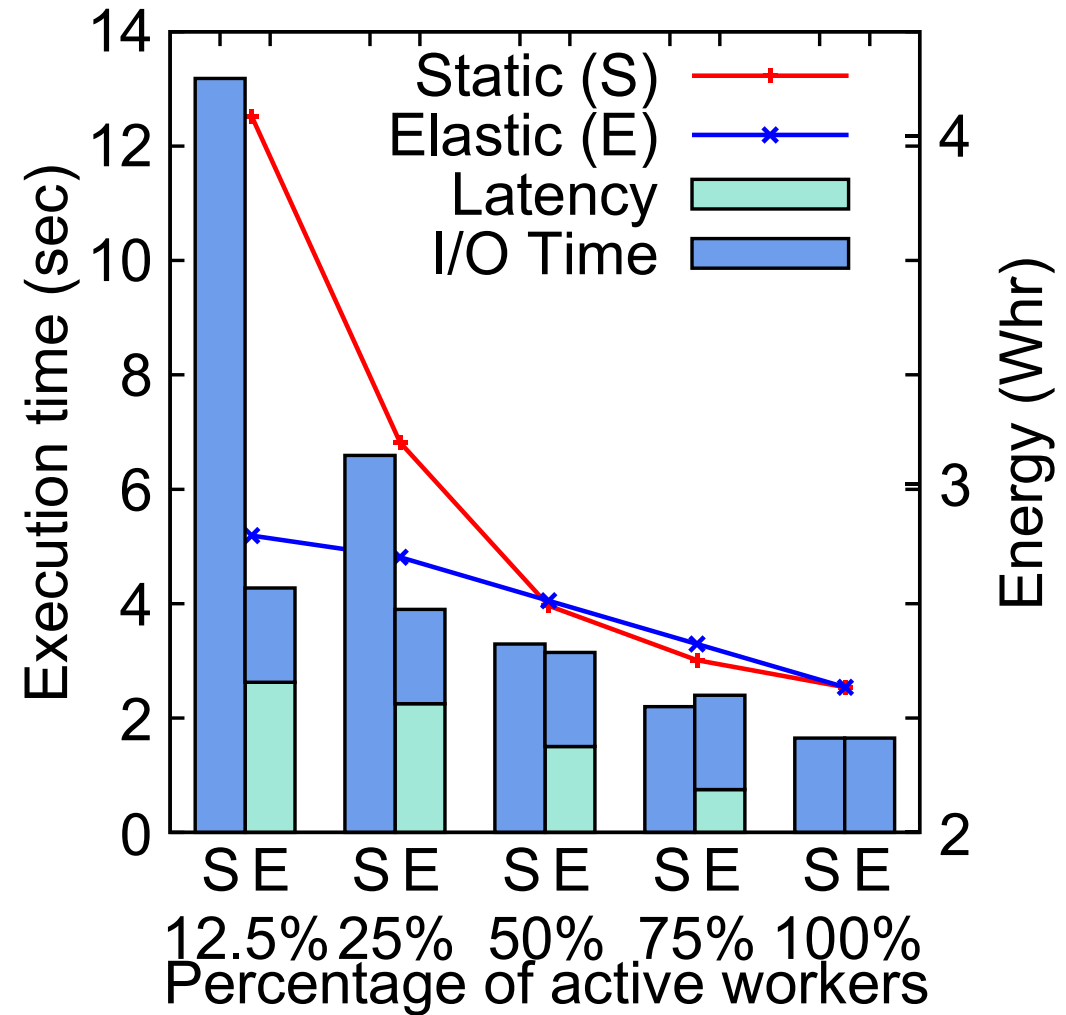
Label Dispatching Throughput

- Metric: Labels per second
- Dispatcher runs on a dedicated node
- 100K auto generated labels
 - Mixed read and write
 - Equal size
- Linear scalability
 - Round robin and random select 55-125K
 - Constraint-based more communication intensive
 - MinMax more CPU intensive due to DP approach

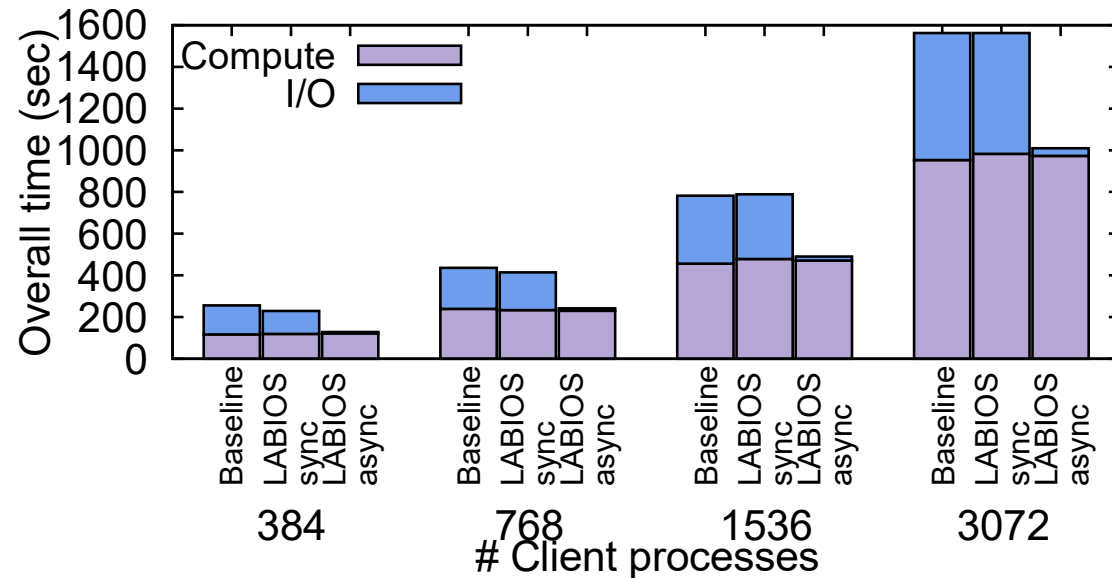


Storage Malleability

- Metric: Total I/O time in sec
- 4096 labels of 1MB each
- Vary the ratio of active – suspended workers
- Worker activation in 3 sec on average
- Worker allocation techniques
 - Static (S): labels only on active workers
 - Elastic (E): labels to all workers (even on suspended paying the penalty of activation)
- When small % of workers are active, elastic boosts performance
- When enough workers are active, activation latency hurts performance

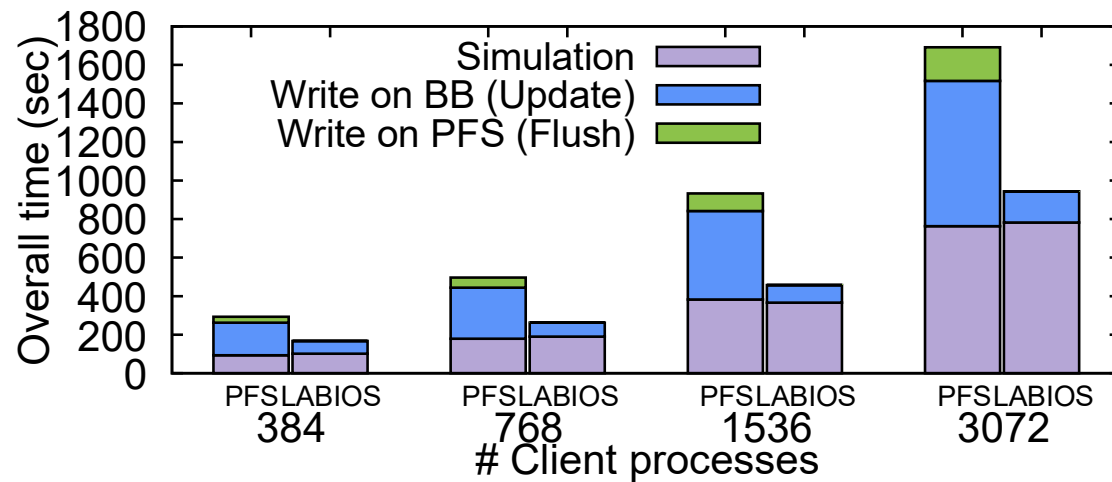


I/O Asynchronicity



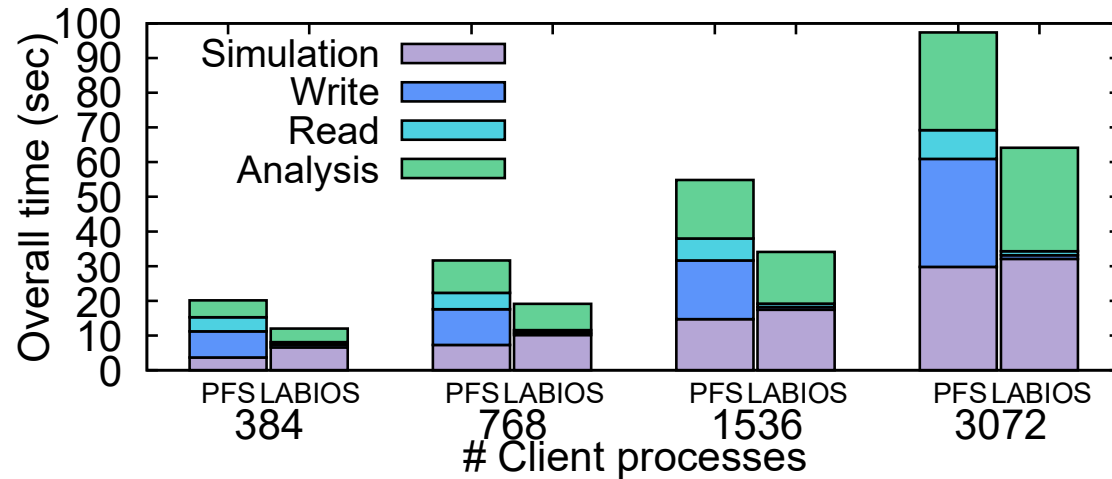
- Metric: Overall execution time in sec
- Support of both sync – async modes
- Label paradigm fits (naturally) in async
- CM1 simulation scaled up to 3072 processes with 16 time steps
- Each process writes 32MB of I/O
 - 100GB per step for the 3072 case
- Sync mode competitive with PFS baseline
- Async mode overlaps label execution with computations
 - **16x boost** in I/O performance
 - **40% reduction** in execution time

Resource Heterogeneity



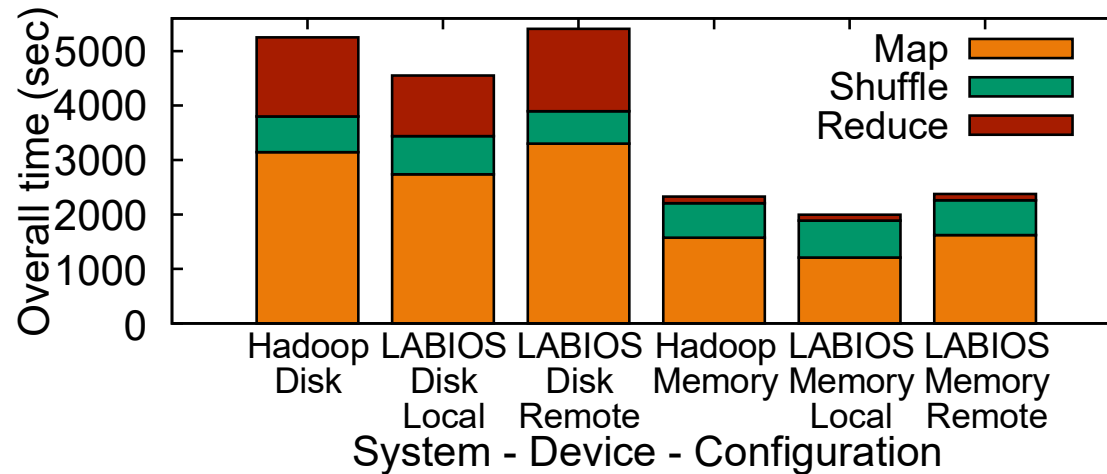
- Metric: Overall execution time in sec
- HACC simulation scaled up to 3072 processes with 16 time steps
- Update-heavy workload
 - Each process updates 32MB of I/O
 - Checkpoint in burst buffers
 - Final flush of last checkpoint data to PFS
- **6x improvement** in I/O performance
- Flushing in the background from workers

Data Provisioning



- Metric: Overall execution time in sec
- Montage application
 - Multiple executables that share data
- 50GB of intermediate results in temporary files in PFS
- LABIOS shares data via the warehouse (i.e., in-memory)
 - Label destination is analysis compute nodes
- Performance acceleration
 - No temporary files are created in remote storage
 - Simulation and analysis can be pipelined
- **17x boost** in I/O performance
- **65% reduction** in execution time

Storage Bridging



- Metric: Overall execution time in sec
- Two modes for LABIOS:
 - Node-local I/O (similar to HDFS)
 - Remote external I/O (similar to HPC)
- Map processes read 32MB each and then write them back to storage
- Reduce processes read 32MB each
- Shuffle sends 32MB through network
- Hadoop-memory optimized version
 - No disk I/O for intermediate results
- LABIOS employs collective I/O to perform data aggregations
- LABIOS successfully integrates MapReduce with HPC

Conclusions

- Supporting a wide range of workflows with different, often conflicting, I/O requirements under a single platform is challenging.
- A new way to perform I/O is required. Desired features include:
 - Storage malleability
 - Asynchronous I/O
 - Resource Heterogeneity
 - Data Provisioning
 - Storage Bridging
- LABIOS provides storage flexibility, versatility, and agility due to a new data model, the (data) labels and its decoupled data-centric architecture.
- LABIOS can boost I/O performance on certain workloads by up to **17x** and reduce overall execution time by **40-60%**.

Thank you
Any questions?

LABIOS: A Label-Based I/O System

Anthony Kougkas, Hariharan Devarajan, Jay Lofstead*, and Xian-He Sun
Illinois Institute of Technology, *Sandia National Laboratories

Find more at:

www.cs.iit.edu/~scs

www.akougkas.com/research/labios

Please come to our poster Slot #30 tonight at 6:30pm in Room 301A

We would like to thank
our sponsors the
National Science
Foundation

