

The Flux framework: overcoming scheduling and management challenges of exascale workflows

ECP Tutorial, Module 1:

Feb 10, 2023, 2-3:30p CST

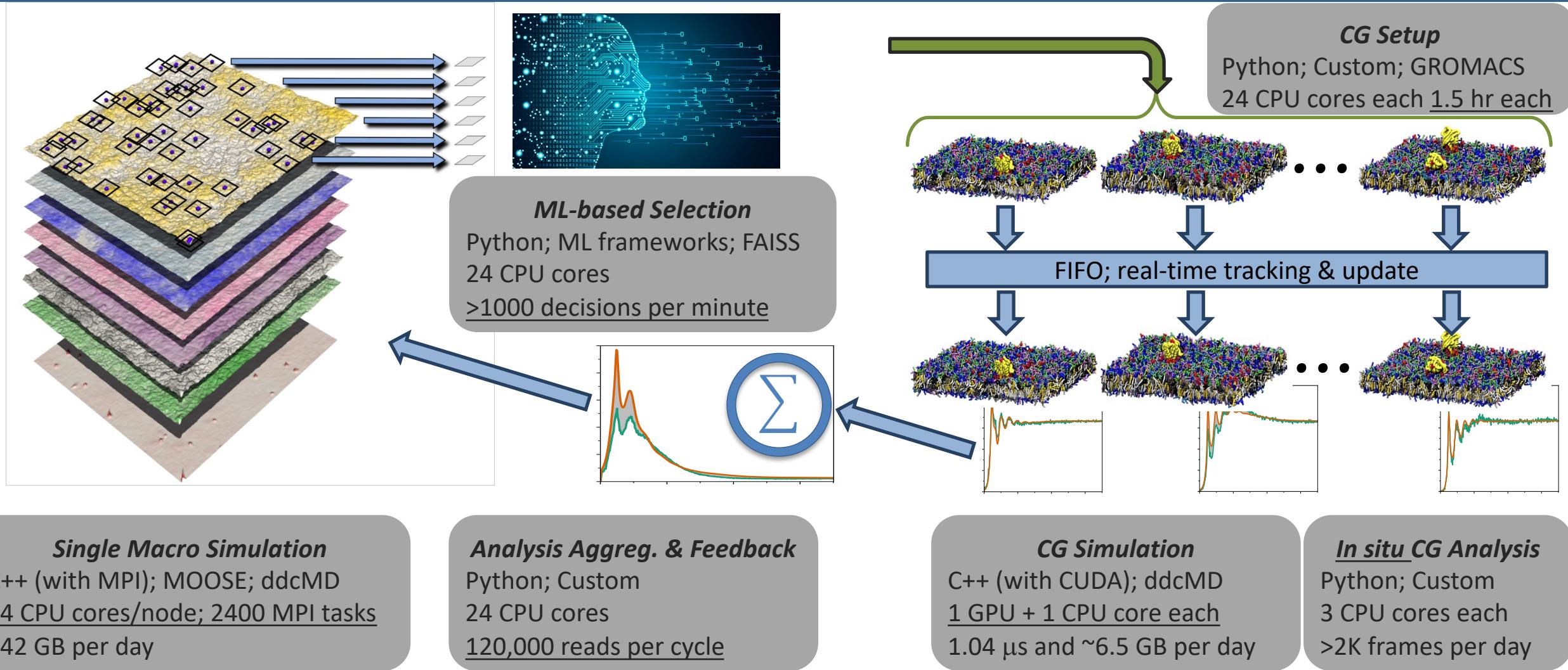
Al Chu, James Corbett, Ryan Day, Jim Garlick, Giorgis Georgakoudis,
Mark Grondona, **Dan Milroy**, Zeke Morton, Chris Moussa, **Tapasya Patki**,
Barry Rountree, Abhik Sarkar, **Tom Scogland**, Vanessa Sochat,
Becky Springmeyer, **Jae-Seung Yeom**



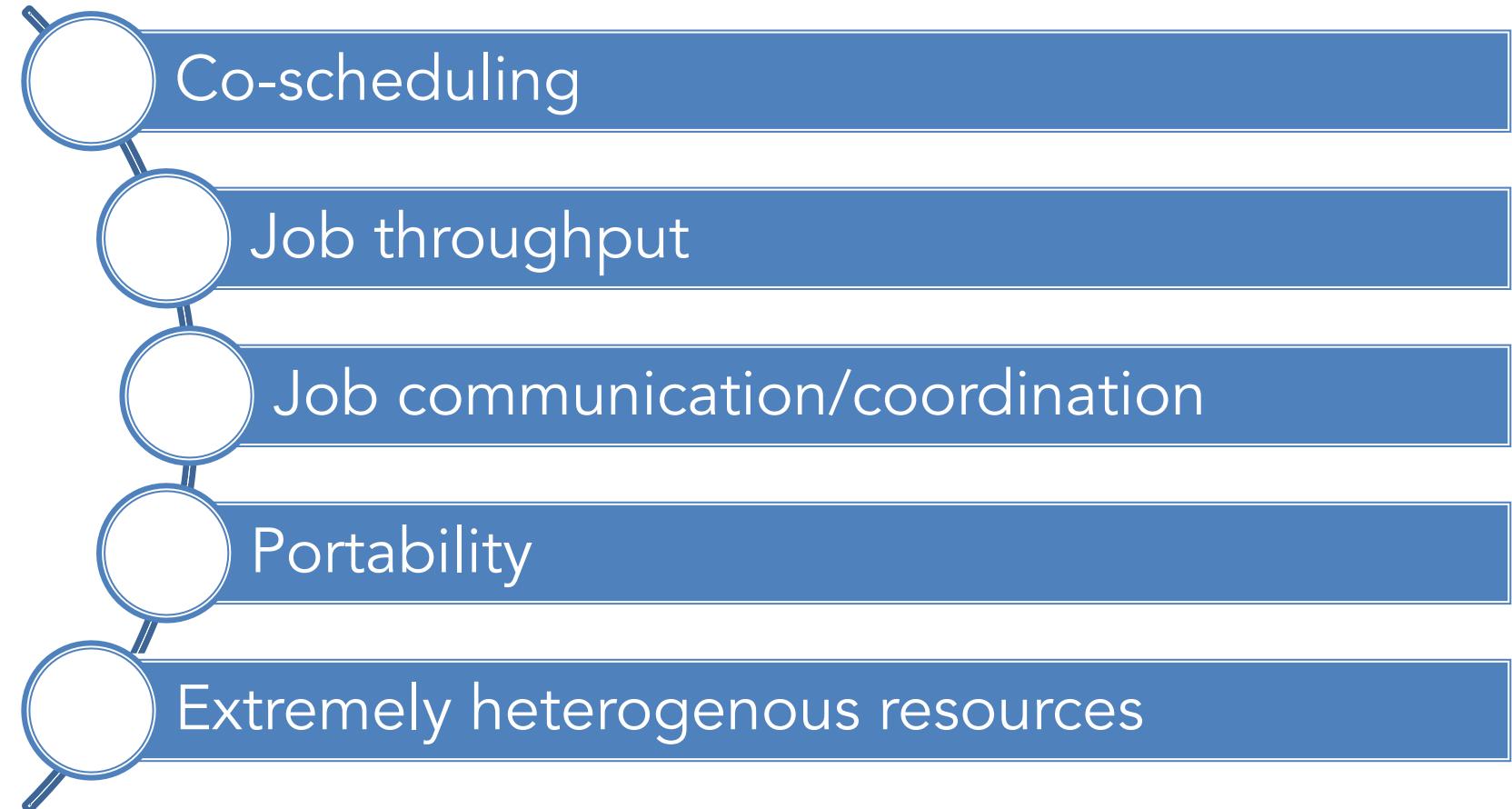
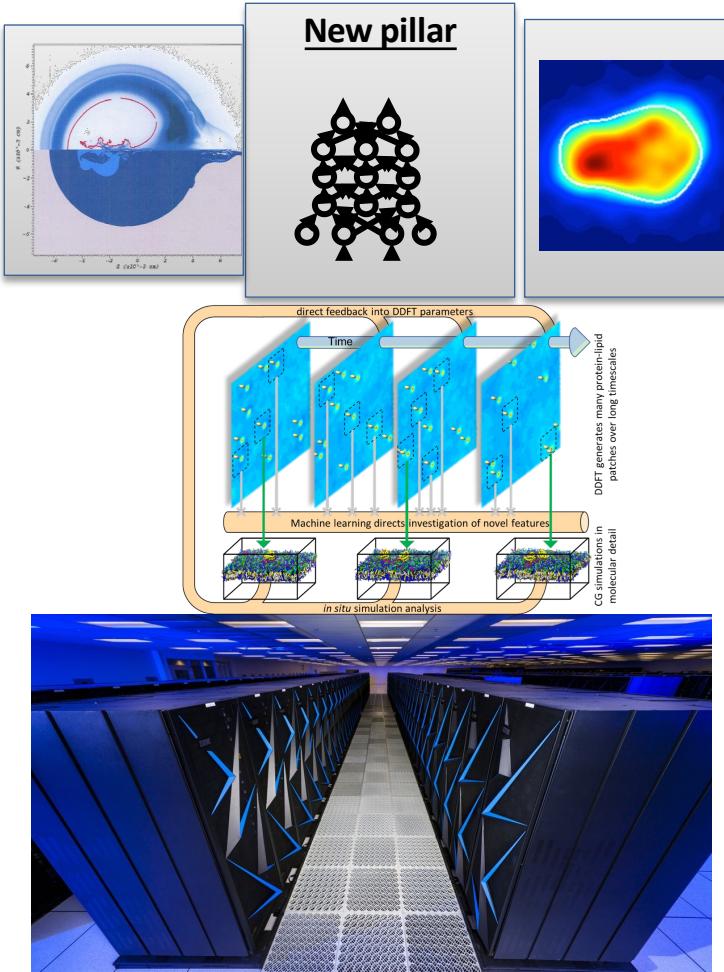
Agenda for Module 1

Introduction to Flux <ul style="list-style-type: none">• Overview and Roadmap• Workflow and Resource Management Challenges at Exascale• Hierarchical and Graph-based Scheduling• Building and using Flux APIs	Tom Scogland	30 min
Hands-on exercises with AWS and Jupyter <ul style="list-style-type: none">• Getting to know Flux and launching jobs• Process and job utilities• Bulk job submission• Flux Python API• Deeper dive	Dan Milroy, Flux team	60 min

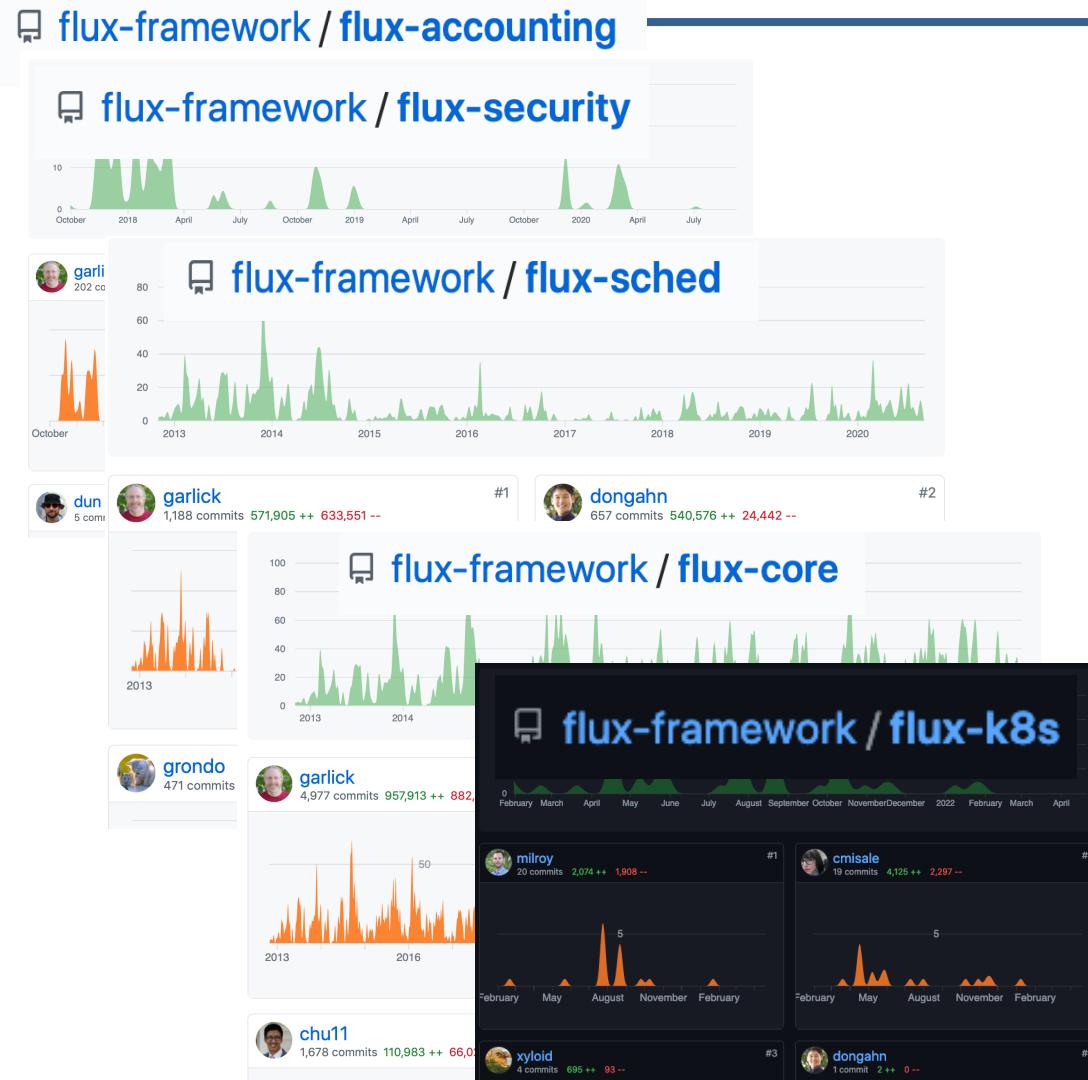
Sierra pre-exascale system is a wakeup call (MuMMI).



Trends towards complex workflows, extreme resource heterogeneity, and converged computing render traditional workload managers increasingly ineffective.



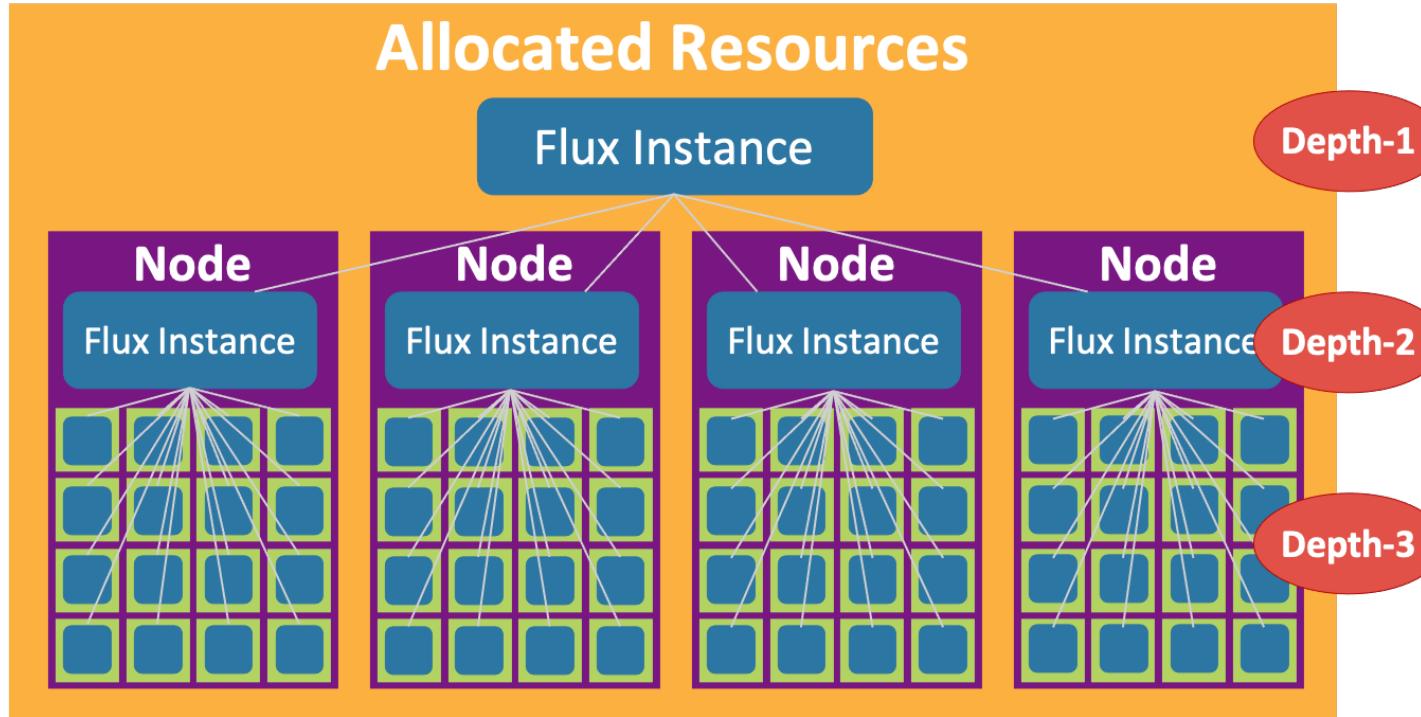
Flux solves key technical problems that emerge from these trends.



- Open-source project in active development at flux-framework GitHub organization
 - Multiple projects: flux-core, -sched, -security, -accounting, -k8s etc.
 - Over 15 contributors including some principal engineers behind Slurm
- Single-user and System instance modes
 - Single-user mode in production for about 4 years
 - Multi-user mode debuting on LLNL Linux clusters
- Plan of record for **LLNL El Capitan** exascale system



Flux's fully hierarchical resource management solves primary deficiencies of the conventional approach.



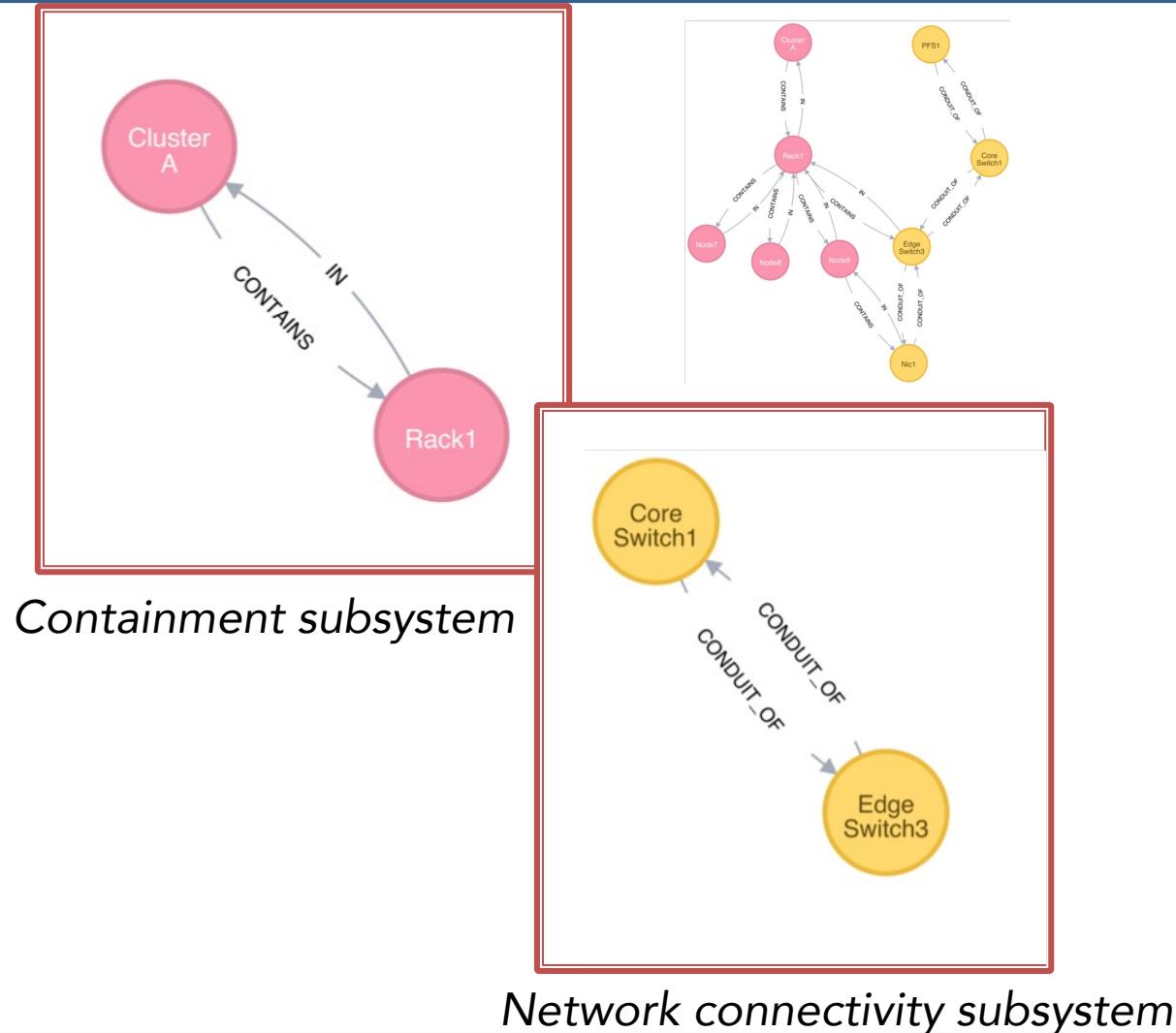
"Fractal scheduling" mitigates centralized scheduler bottleneck

- handles high throughput
- job steps needn't hit central sched

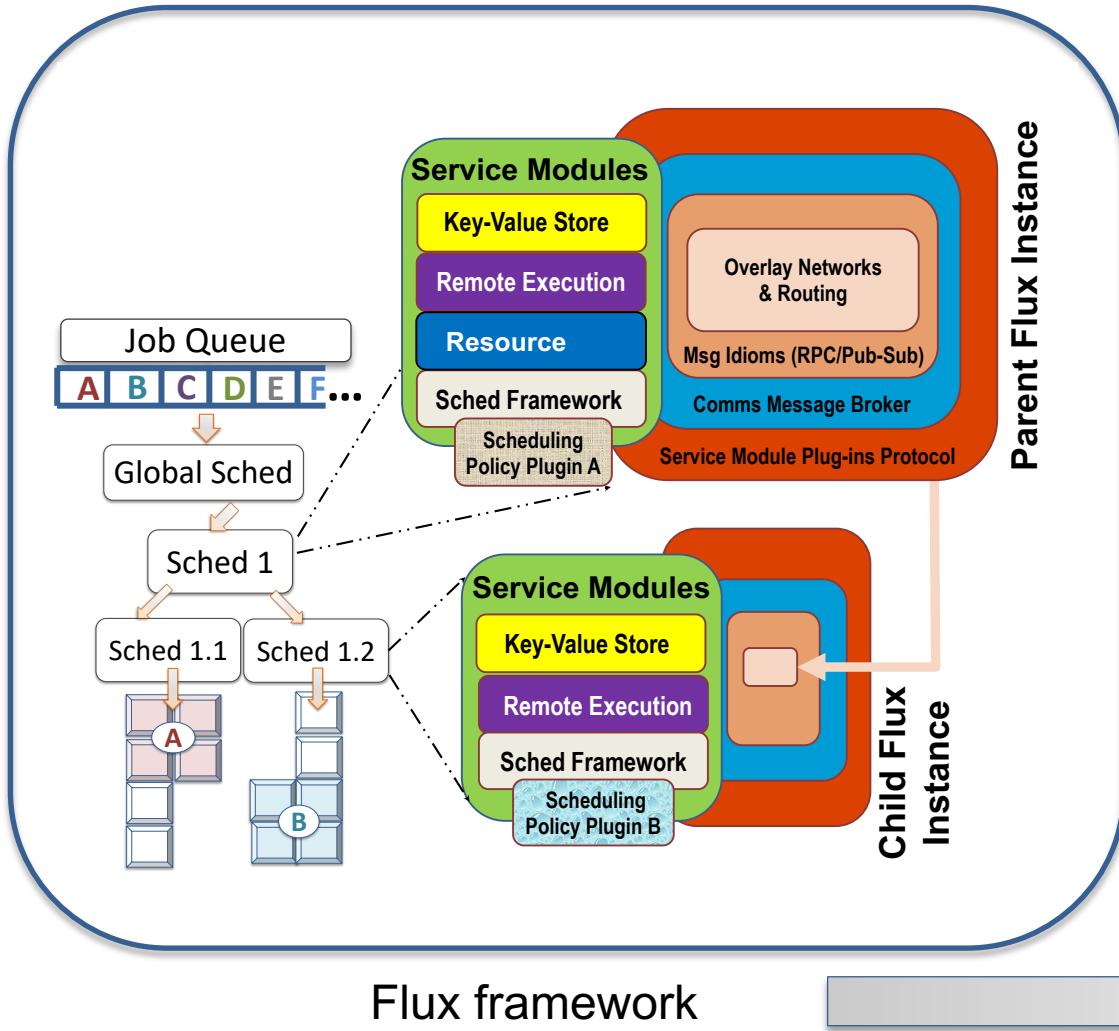
- Full workflow-enablement support
 - Via hierarchical resource subdivision
 - Sub-resource manager per subdivision with service specialization
 - E.g., at LLNL: MuMMI, AHA MoleS, UQP
- Capable of managing resources from almost anywhere
 - Bare metal resources, virtual machines in the cloud, HPC resources allocated by another workload manager
 - Workflows only need to program to Flux
- Provide rich and well-defined interfaces
 - Facilitate communications and coordination among tasks within a workflow
 - CLI, Python, C

Flux pioneers and uses graph-based scheduling to manage complex combinations of extremely heterogenous resources.

- Traditional resource data models are largely ineffective for resource heterogeneity
 - designed when systems were simpler
 - node-centric models
- Elevate resource relationships (edges) to an equal footing with resources (vertices)
- Complex scheduling can be expressed without changing the scheduler code
- Rich and well-defined C and C++ API for graph traversal and allocation



Flux is architected to embody our fully hierarchical scheduling model.



Techniques

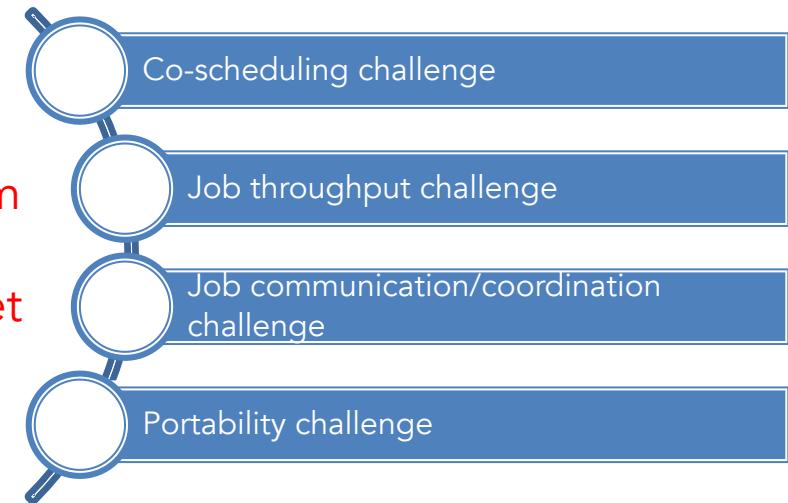
Scheduler Specialization

Scheduler Parallelism

Rich API set

Consistent API set

Challenges



A rich set of well-defined APIs enables easy job coordination and communication.

- Jobs in ensemble-based simulations often require close coordination and communication with the scheduler as well as among them.
 - Traditional CLI-based approach can be slow and cumbersome.
 - Ad hoc approaches (e.g., many empty files) can lead to many side effects.
- Flux provides well-known communication primitives.
 - Pub/sub, request-response, and send-recv patterns
- High-level services
 - Key-value store (KVS) API
 - Job API (submit, wait, state change notification, etc)
- Flux's APIs are consistent across different platforms

Docs » man3

 Edit on GitHub

man3

C Library Functions

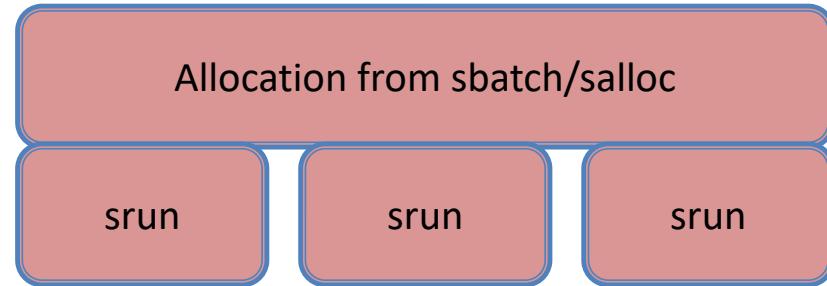
- [flux_attr_get\(3\)](#)
- [flux_aux_set\(3\)](#)
- [flux_child_watcher_create\(3\)](#)
- [flux_content_load\(3\)](#)
- [flux_core_version\(3\)](#)
- [flux_event_decode\(3\)](#)
- [flux_event_publish\(3\)](#)
- [flux_event_subscribe\(3\)](#)
- [flux_fatal_set\(3\)](#)
- [flux_fd_watcher_create\(3\)](#)
- [flux_flags_set\(3\)](#)
- [flux_future_and_then\(3\)](#)
- [flux_future_create\(3\)](#)
- [flux_future_get\(3\)](#)
- [flux_future_wait_all_create\(3\)](#)
- [flux_get_rank\(3\)](#)
- [flux_get_reactor\(3\)](#)
- [flux_handle_watcher_create\(3\)](#)
- [flux_idle_watcher_create\(3\)](#)
- [flux_kvs_commit\(3\)](#)
- [flux_kvs_copy\(3\)](#)
- [flux_kvs_getroot\(3\)](#)
- [flux_kvs_lookup\(3\)](#)
- [flux_kvs_namespace_create\(3\)](#)
- [flux_kvs_txn_create\(3\)](#)
- [flux_log\(3\)](#)
- [flux_msg_cmp\(3\)](#)



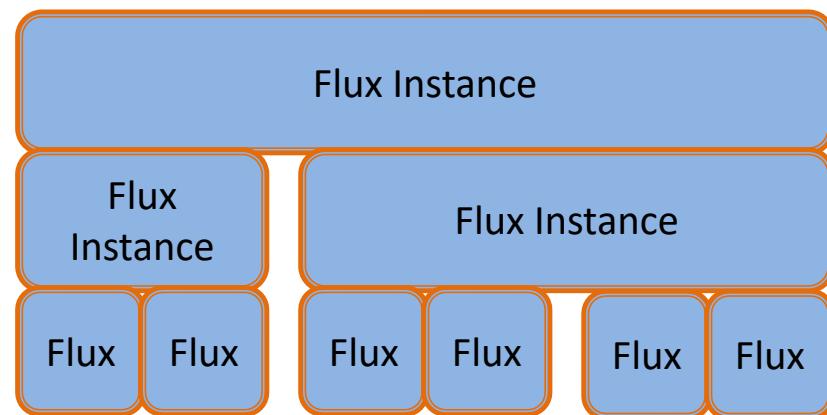
<https://flux-framework.readthedocs.io/projects/flux-core/en/latest/index.html>

Flux's hierarchical scheduling makes managing and scheduling complex allocations easier.

- Significant start-up cost for running on new resource manager (portability)
 - Learn to use the manager
 - Translate batch script logic
- Traditional resource managers can't subdivide allocations (hierarchy)
 - Slurm has two levels: **sbatch** creates allocation, **srun** launches parallel task
 - **sbatch** from inside an allocation starts a separate allocation
- Flux's hierarchical scheduling provides a natural solution
 - Request one large allocation and partition schedulable resource subgroups



The Slurm hierarchy

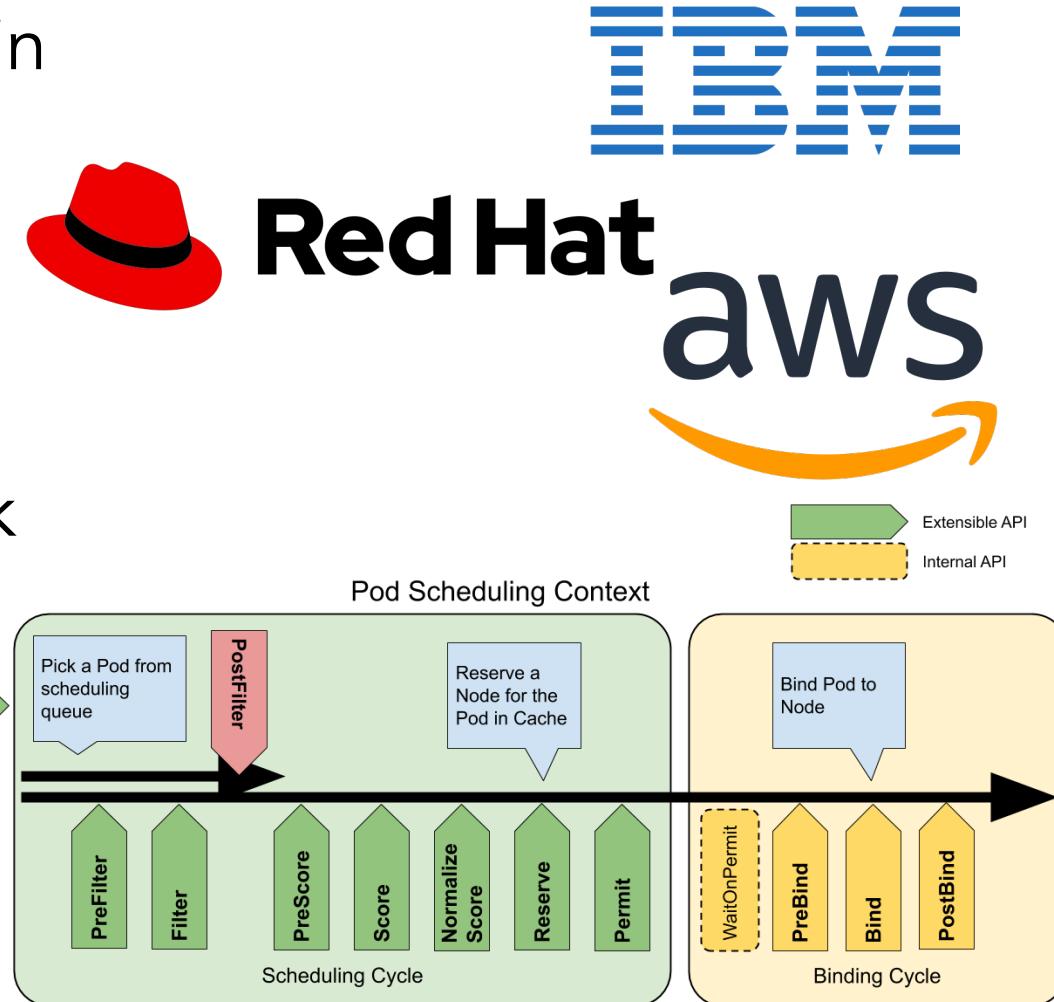


One possible Flux hierarchy

Team formed industry collaborations to tackle converged computing challenges and contribute to community.

We bring complementary backgrounds in HPC, cloud, and performance-oriented orchestration

- LLNL LDRD team
- IBM T.J. Watson Research Center: T. Elengikal, M. Drocco, C. Misale, Y. Park
- AWS: E. Bollig, M. Hugues, H. Poxon, L. Wofford
- Integrate Fluxion into Kubernetes via **Fluence**



Accessing the hands-on tutorial

- Using our EC2 instance at <https://tutorial.flux-framework.org>
 - Choose a unique username (if not, you'll be sharing with someone else)
 - Password: **ecp-flux-tutorial2023**
 - Double-click **ecp23-flux.ipynb** to start
 - To execute a cell in JupyterLab: **Shift+Enter**
 - It's a shared instance- please don't run computationally demanding tasks in your container
 - The instance will disappear shortly after the tutorial
- Running locally with Docker:
 - **git clone <https://github.com/flux-framework/Tutorials.git>**; README in the 2023-ECP JupyterNotebook directory



**Lawrence Livermore
National Laboratory**

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

The Flux framework: overcoming scheduling and management challenges of exascale workflows

ECP Tutorial, Module 2:

Feb 10, 2023, 3:45 – 5:15p CST

Al Chu, James Corbett, Ryan Day, Jim Garlick, Giorgis Georgakoudis,
Mark Grondona, **Dan Milroy**, Zeke Morton, Chris Moussa, **Tapasya Patki**,
Barry Rountree, Abhik Sarkar, **Tom Scogland**, Vanessa Sochat,
Becky Springmeyer, **Jae-Seung Yeom**



Agenda for Module 2

Fluxion and Graph-Based Scheduling • I/O and Variation-aware scheduling	Tapasya Patki	15 min
Converged Computing and Fluence	Daniel Milroy	15 min
I/O Awareness and Data Management • DYAD • Hands on exercises with DYAD	Jae-Seung Yeom and Ian Lumsden	30 min
AHA Moles: Example of Integrated Workflow Management with Flux	Xiaohua Zhang	20min
Summary and Q&A	Tom Scogland, Flux team	10 min

Fluxion: Flux's Resource & Job Model

ECP Tutorial, Feb 10, 2023

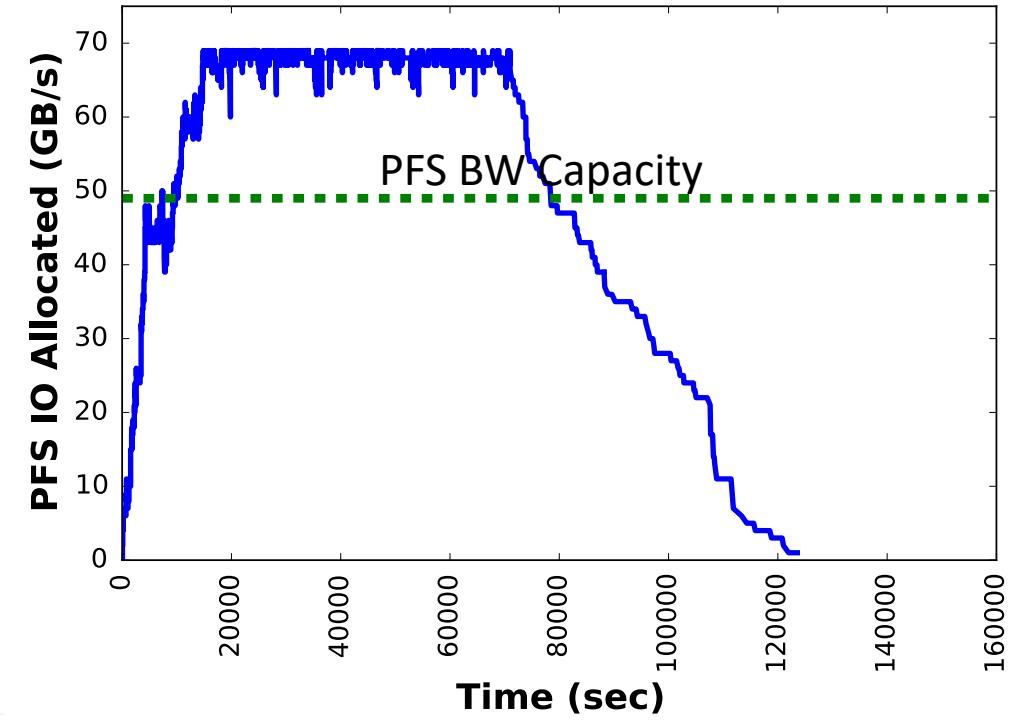
Tapasya Patki

Acknowledgment: Dong H. Ahn, Ned Bass, Subhendu Behera, Albert Chu, Kyle Chard, Brandon Cook, James Corbett, Ryan Day, Franc Di Natalie, David Domyancic, Phil Eckert, Stephen Herbein, Jim Garlick, Elsa Gonsiorowski, Mark Grondona, Helgi Ingolfsson, Shantenu Jha, Zvonko Kaiser, Dan Laney, Carlos Eduardo Gutierrez, Edgar Leon, Don Lipari, Daniel Milroy, Joseph Koning, Claudia Misale, Chris Moussa, Frank Muller, Yoonho Park, Luc Peterson, Barry Rountree, Thomas R. W. Scogland, Becky Springmeyer, Michela Taufer, Jae-Seung Yeom, Veronica Vergara and Xiaohua Zhang



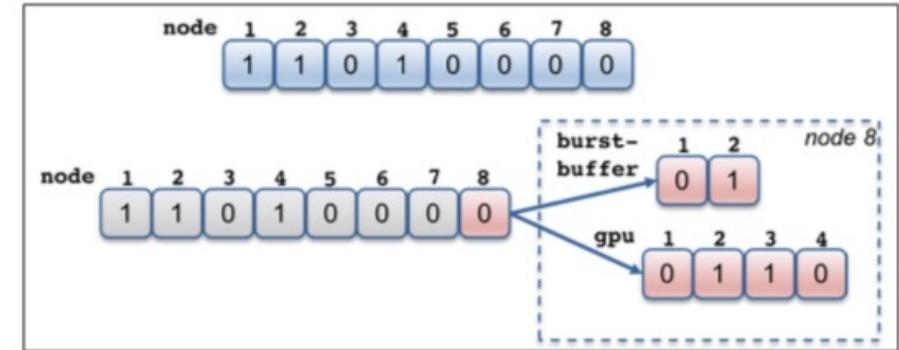
The changes in resource types are challenging.

- Problems are not just confined to the workload/workflow challenge.
- Resource types and their relationships are also becoming increasingly complex.
- Much beyond compute nodes and cores requiring partial occupancy and accounting...
 - GPGPUs
 - Burst buffers
 - I/O and network bandwidth
 - Network locality
 - Power
 - Variation



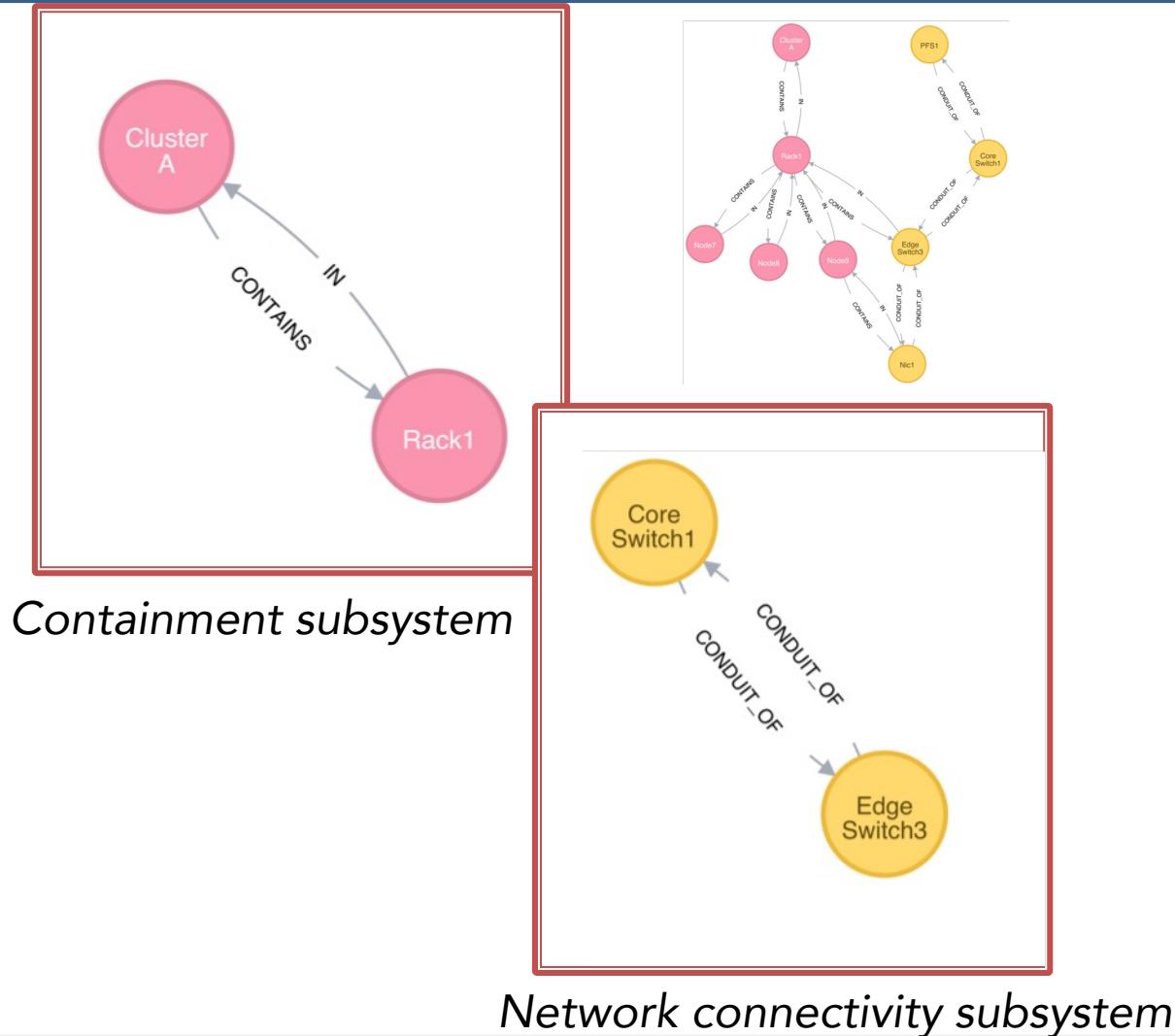
The traditional resource data models are largely ineffective to cope with the resource challenge.

- Designed when the systems are much simpler
 - Node- or core-centric models
 - SLURM: bitmaps to represent a set of compute nodes
 - PBSPro: a linked-list of nodes
- HPC has become far more complex
 - Evolutionary approach to cope with the increased complexity
 - E.g., add auxiliary data structures on top of the node-centric data model
- Can be quickly unwieldy
 - Every new resource type requires new a user-defined type
 - A new relationship requires a complex set of pointers cross-referencing different types.
 - Dynamic updating of flow resources is not supported



Flux pioneers and uses graph-based scheduling to manage complex combinations of extremely heterogenous resources.

- Traditional resource data models are largely ineffective for resource heterogeneity
- Elevate resource relationships (edges) to an equal footing with resources (vertices)
- Complex scheduling can be expressed without changing the scheduler code
- Rich and well-defined C and C++ API for graph traversal and allocation
- High expressibility and level of detail control



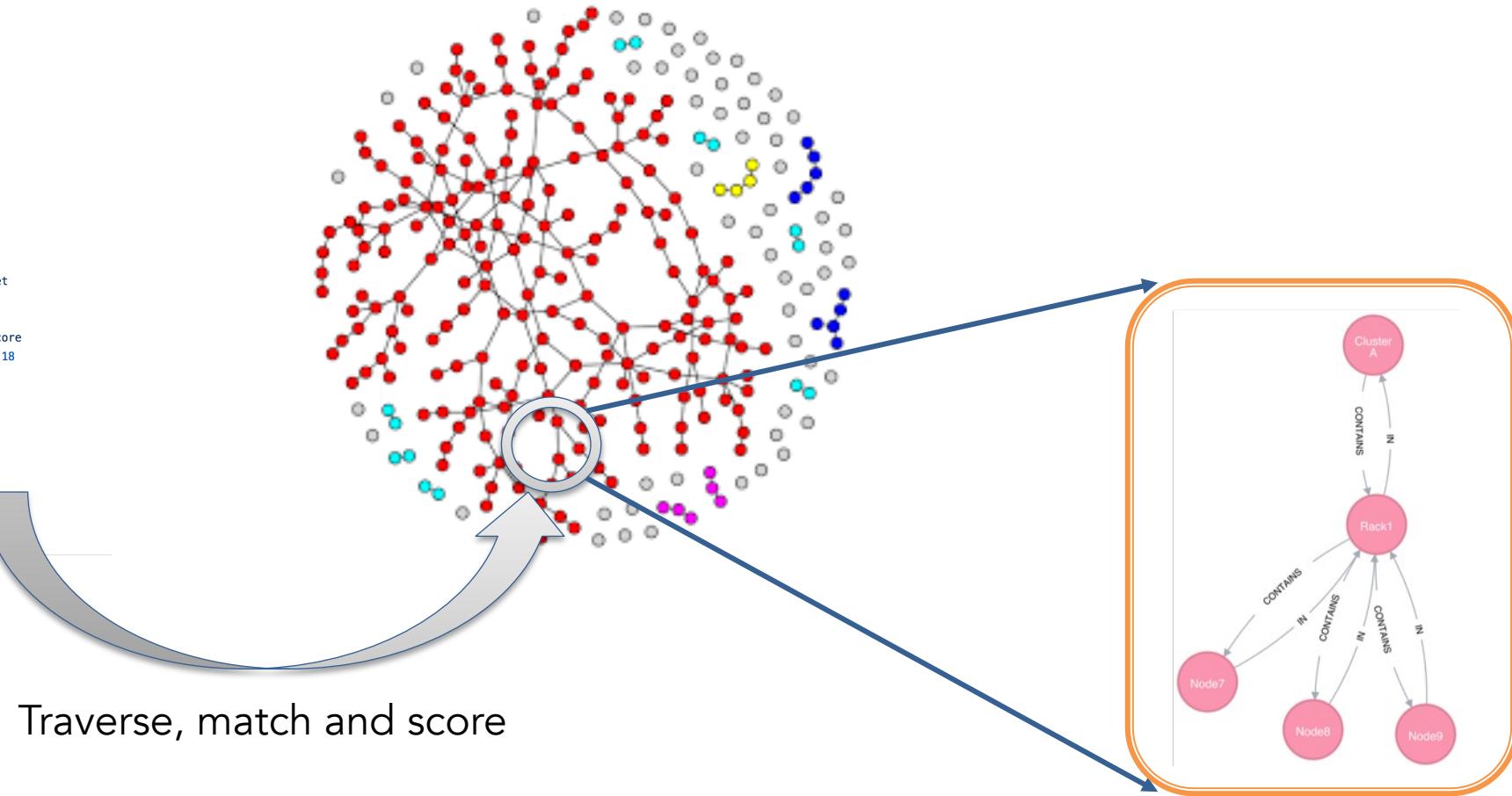
Flux's graph-oriented canonical job-spec allows for a highly expressive resource requests specification.

- Graph-oriented resource requirements
 - Express the resource requirements of a program to the scheduler
 - Express program attributes such as arguments, run time, and task layout, to be considered by the execution service
- cluster->racks[2]->slot[3]->node[1]->sockets[2]->core[18]
- **slot** is the only non-physical resource type
 - Represent a schedulable place where program process or processes will be spawned and contained
- Referenced from the tasks section

```
1  version: 1
2  resources:
3    - type: cluster
4      count: 1
5      with:
6        - type: rack
7          count: 2
8          with:
9            - type: slot
10           label: myslot
11           count: 3
12           with:
13             - type: node
14               count: 1
15               with:
16                 - type: socket
17                   count: 2
18                   with:
19                     - type: core
20                     count: 18
21
22 # a comment
23 attributes:
24   system:
25     duration: 3600
26 tasks:
27   - command: app
28     slot: myslot
29     count:
30       per_slot: 1
```

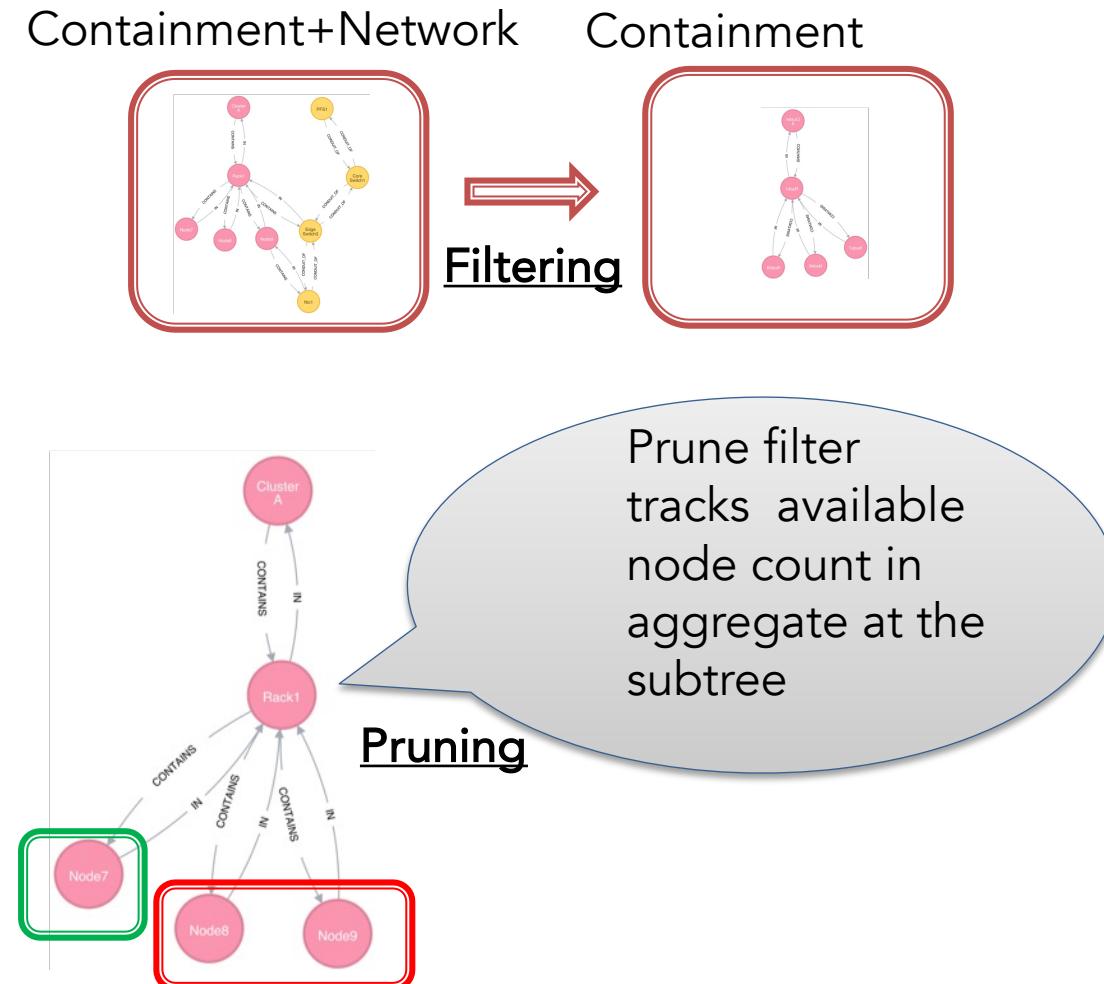
Flux maps our complex scheduling problems into graph matching problems.

```
1 version: 1
2 resources:
3   - type: cluster
4     count: 1
5     with:
6       - type: rack
7         count: 2
8         with:
9           - type: slot
10          label: myslot
11          count: 3
12          with:
13            - type: node
14              count: 1
15              with:
16                - type: socket
17                  count: 2
18                  with:
19                    - type: core
20                      count: 18
21
22 # a comment
23 attributes:
24 system:
25 duration: 3600
26 tasks:
27 - command: app
28   slot: myslot
29   count:
30     per_slot: 1
```

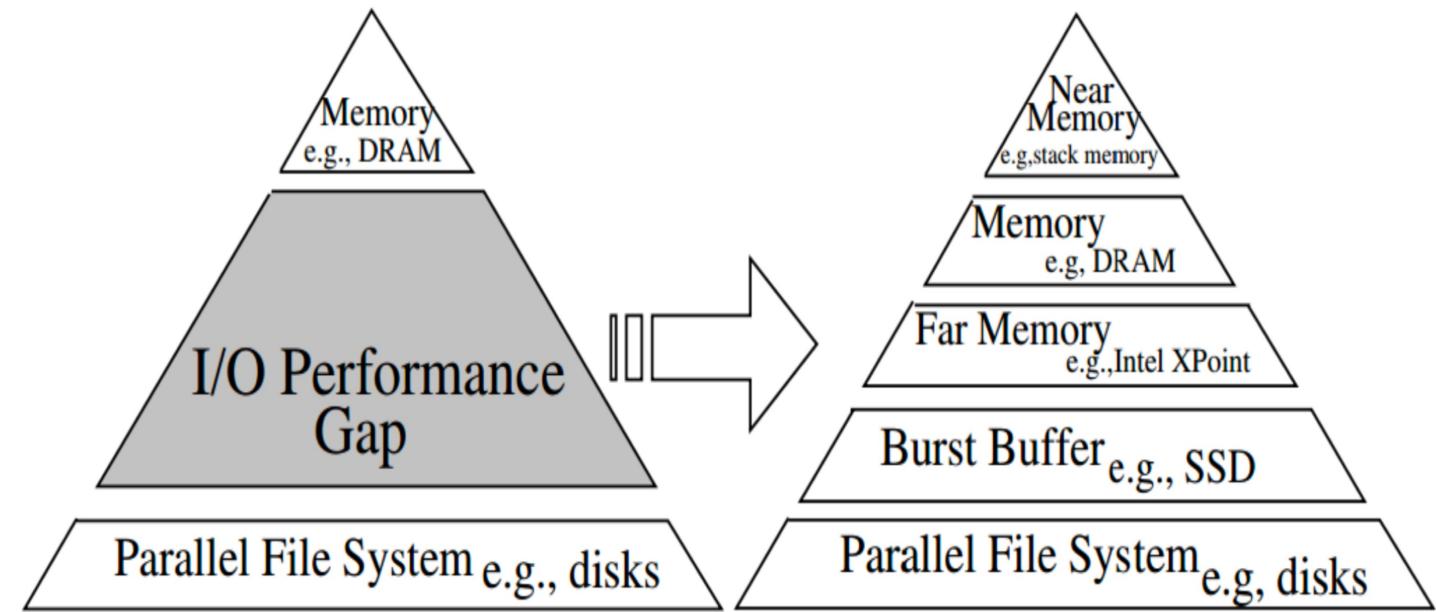
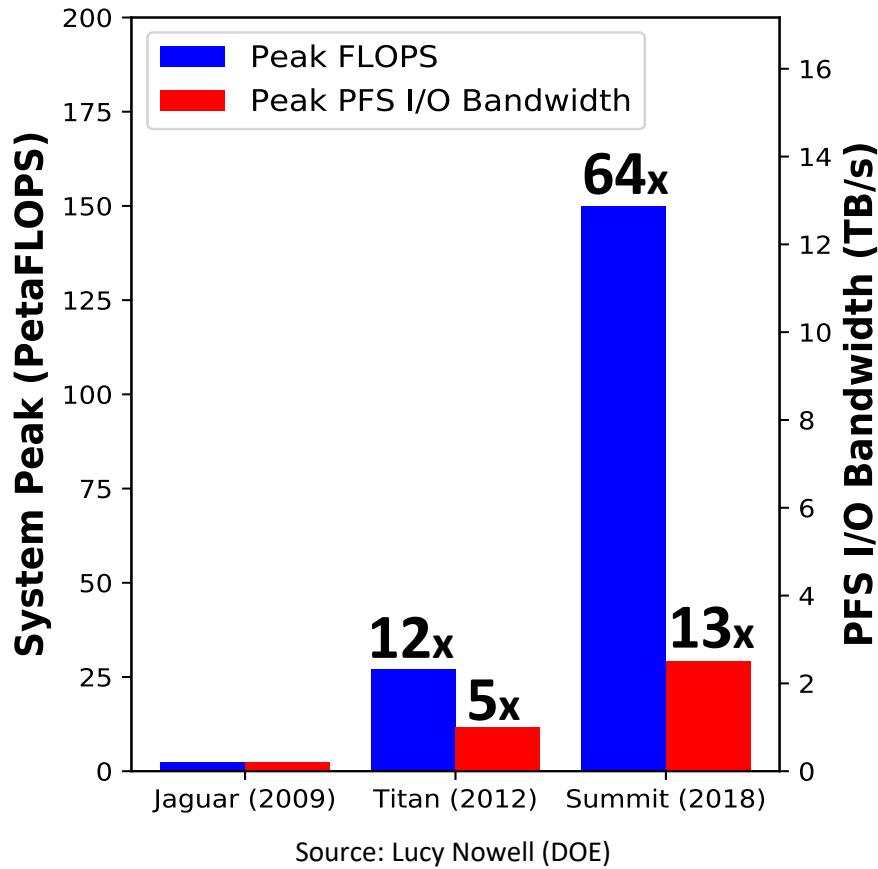


We use graph filtering and pruned searching to manage the graph complexity and optimize our graph search.

- The total graph can be quite complex
 - Two techniques to manage the graph complexity and scalability
- Filtering reduces graph complexity
 - The graph model needs to support schedulers with different complexity
 - Provide a mechanism by which to filter the graph based on what subsystems to use
- Pruned search increases scalability
 - Fast RB tree-based planner is used to implement a pruning filter per each vertex.
 - Pruning filter keeps track of summary information (e.g., aggregates) about subtree resources.
 - Scheduler-driven pruning filter update



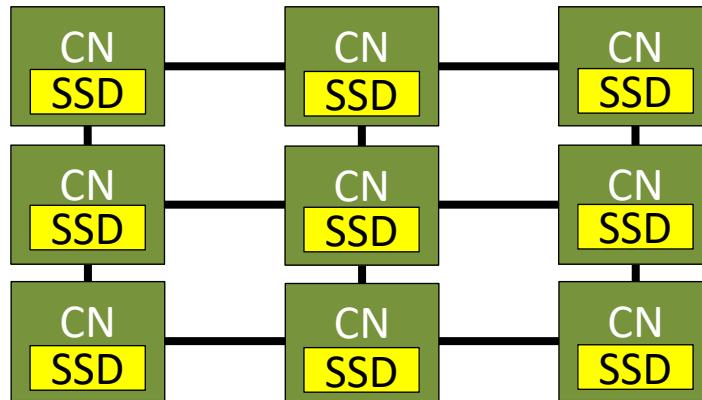
Tiered Storage in HPC



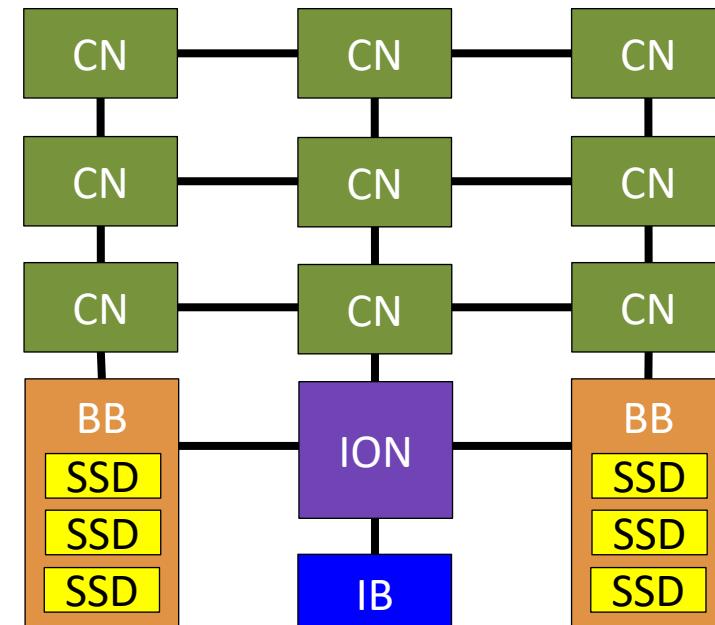
How can we use the Flux graph-based scheduler to allocate these new storage tiers?

Burst Buffer Architectures

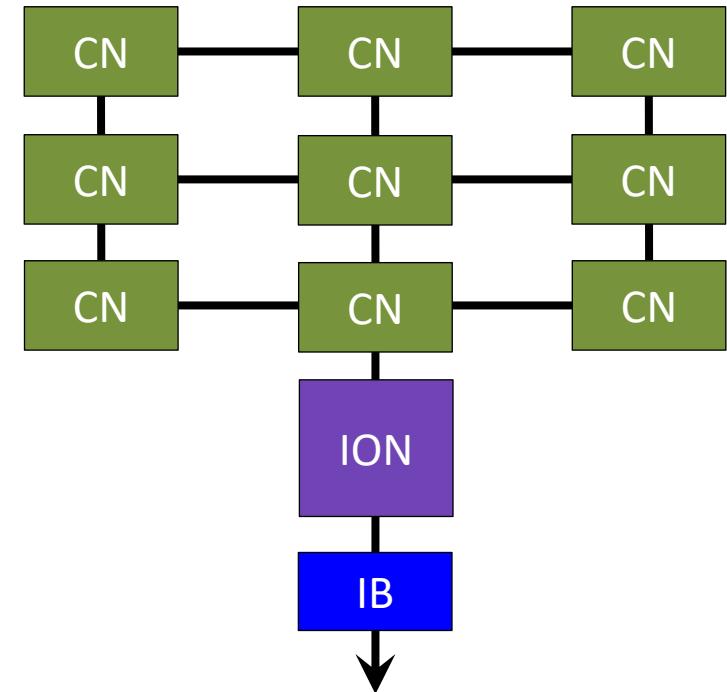
Node-local BB



Remote, shared BB



Filesystem BB



Parallel File System

Parallel File System

Parallel File System

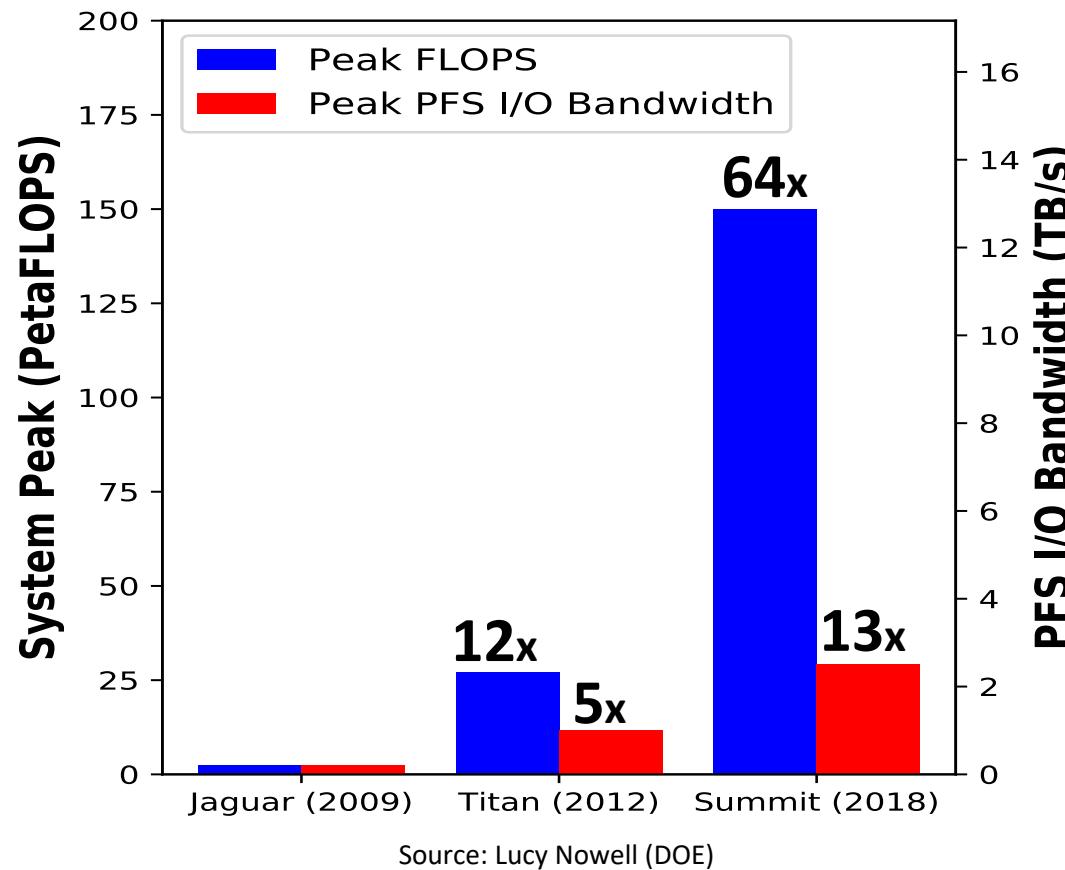
SSD SSD SSD SSD SSD



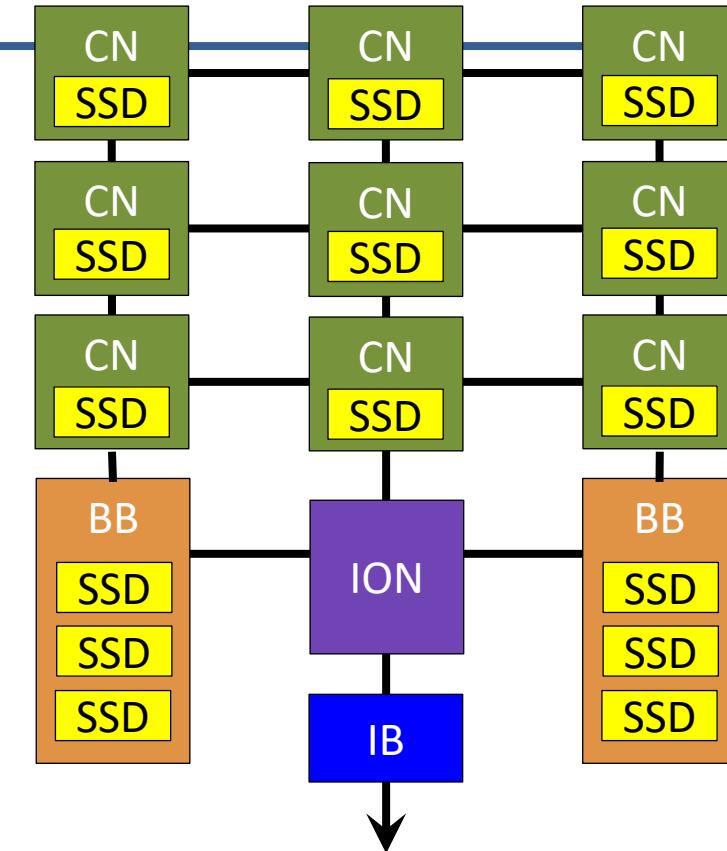
Burst Buffers Functionality “Classes”

- Cache
 - acts like an extended OS page cache for a given job or parallel filesystem
 - all data movement between the BB and the PFS occurs implicitly
- Scratch
 - functions like a /tmp file system for the given node, job, or workflow
- Staging
 - all data movement occurs explicitly for a given job
 - data movement can occur before and after the job is running

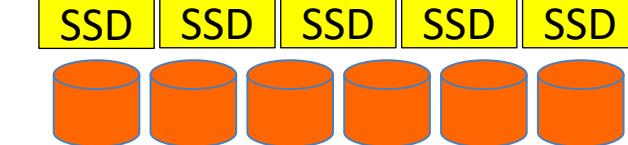
Tiered Storage in HPC



We can use the Flux graph scheduler to allocate these new storage tiers with 0 code changes



Parallel File System

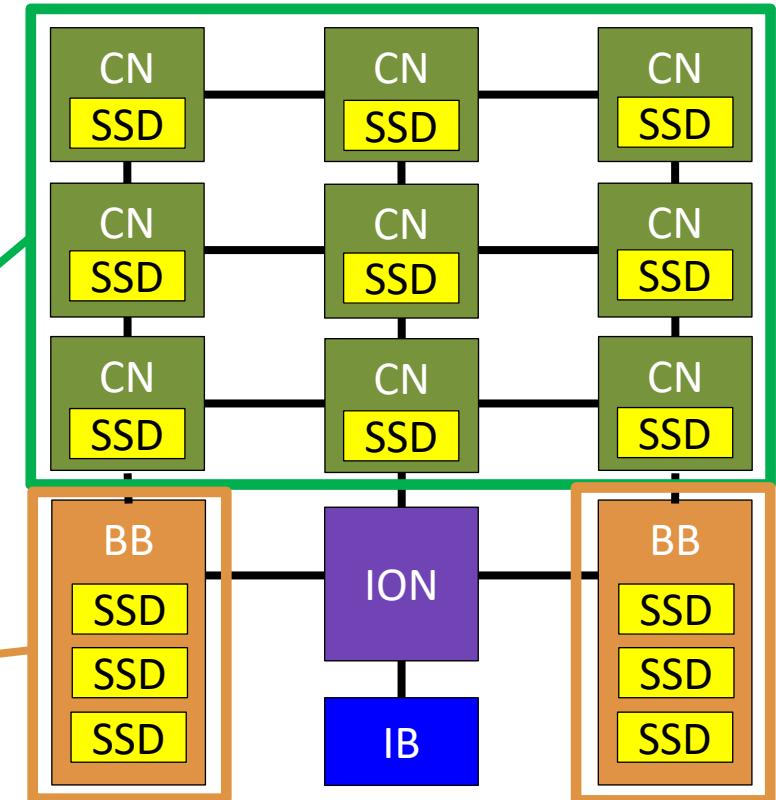


Example Tiered Storage Request in Flux Jobspec

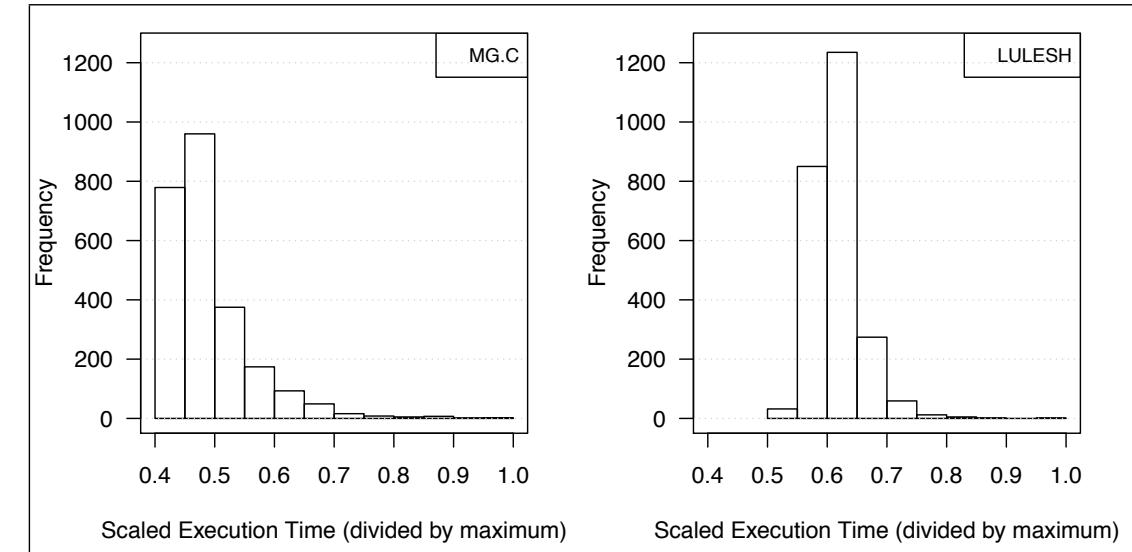
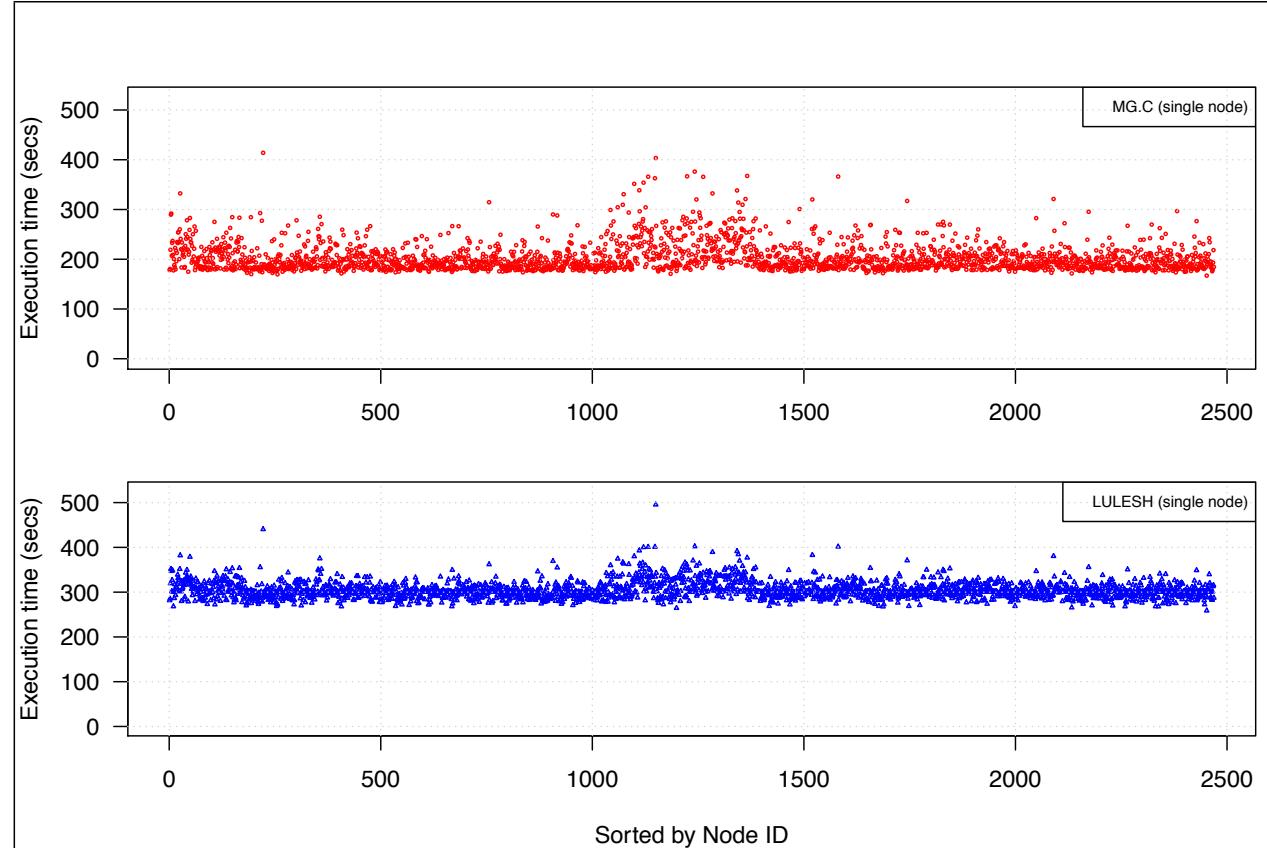
```
resources:
  - type: node
    count: 9
    with:
      - type: slot
        count: 1
        label: default
        with:
          - type: core
            count: 2
          - type: storage
            count: 1
            unit: terabytes
            label: node-local-scratch
  - type: storage
    count: 4
    unit: terabytes
    label: PFS-cache
```

attributes:
storage:

- label: node-local-scratch
mode: scratch
granularity: per-node
stage-in:
 list: /path/to/stage-in-listing
- label: PFS-cache
data-layout: striped
mode: cache
stage-in:
 directory: /path/to/PFS



Variation-aware scheduling with Flux: Addressing Manufacturing Variability, Processor Aging, and inherent heterogeneity

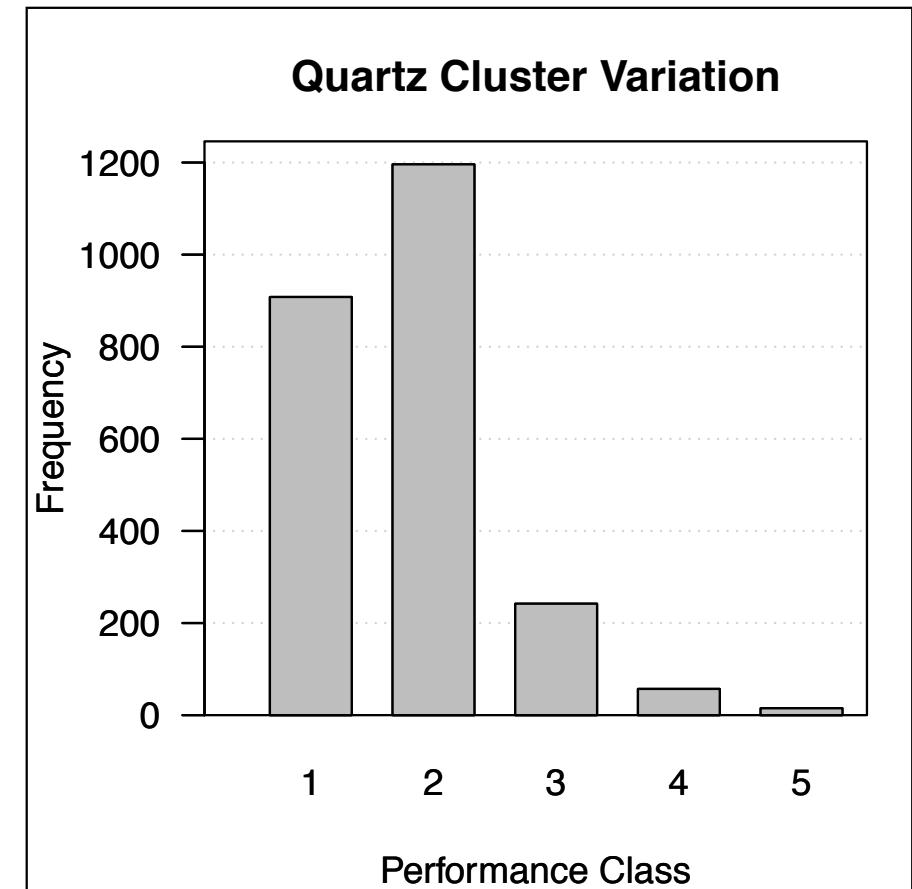


- Real world example under power constraints: Quartz cluster, 2469 nodes, 50 W CPU cap
- 2.47x difference between the slowest and the fastest node for MG
- 1.91x difference for LULESH.

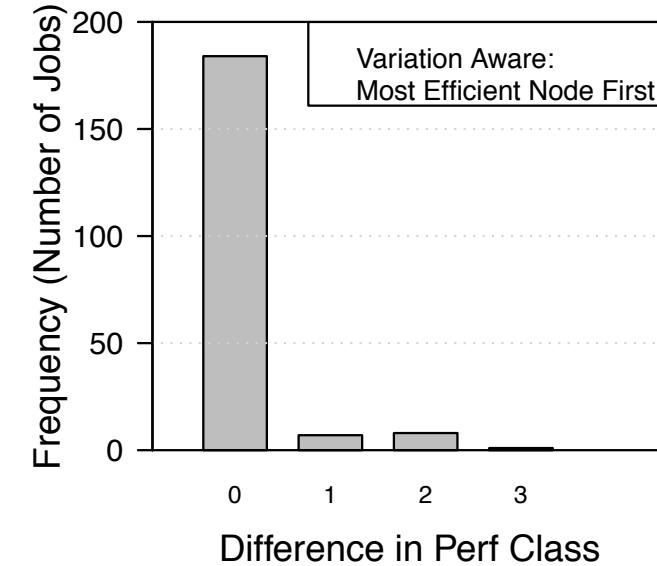
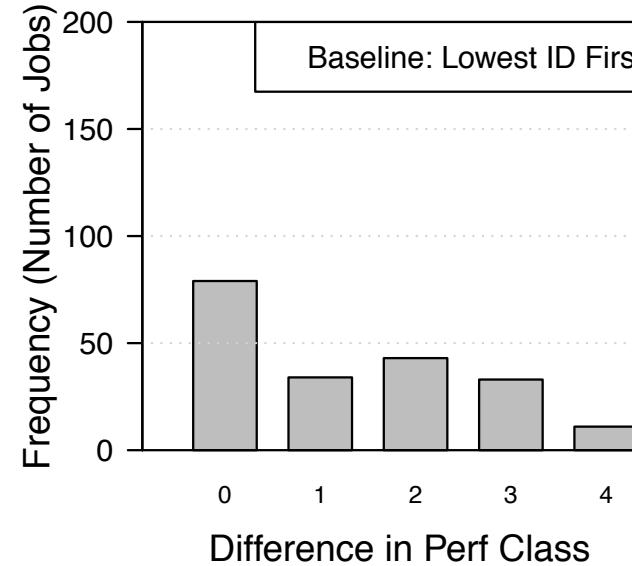
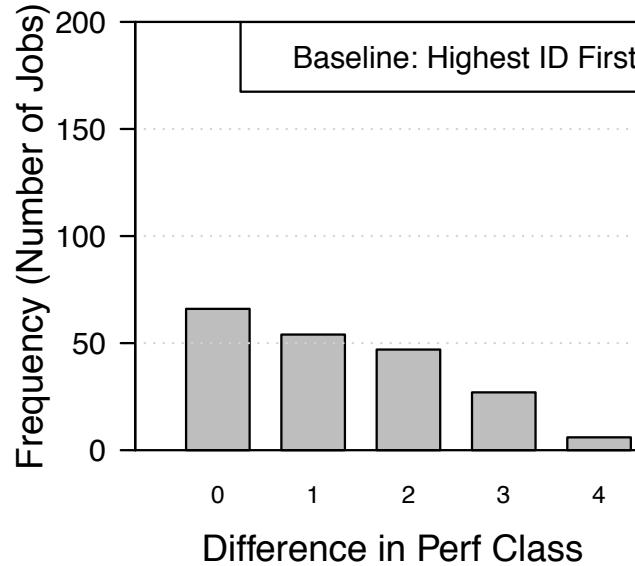
<https://github.com/flux-framework/flux-sched/tree/master/resource/policies>

Example: Statically determining node performance classes

- Ranking every processor is not feasible
- Statically create **bins** of processors with similar performance instead
 - Techniques for this can be simple or complex
 - How many classes to create, which benchmarks to use, which parameters to tweak
 - Our choice: 5 classes, LULESH and MG, 50 W cap
- **Mitigation**
 - **Rank-to-rank:** minimize spreading application across multiple performance classes
 - **Run-to-run:** allocate nodes from same set performance classes to similar applications



Variation-aware scheduling results in 2.4x reduction in rank-to-rank variation in applications with Flux



Flux's graph-based resource model easily and effectively enables this variation-aware scheduler optimization

The Flux Framework and HPC+cloud Converged Computing

Dan Milroy



LLNL-PRES-844974

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344. Lawrence Livermore National Security, LLC

 Lawrence Livermore
National Laboratory

Pre-exascale scientific workflows strain the capabilities of traditional HPC resource managers and schedulers.

Co-scheduling:

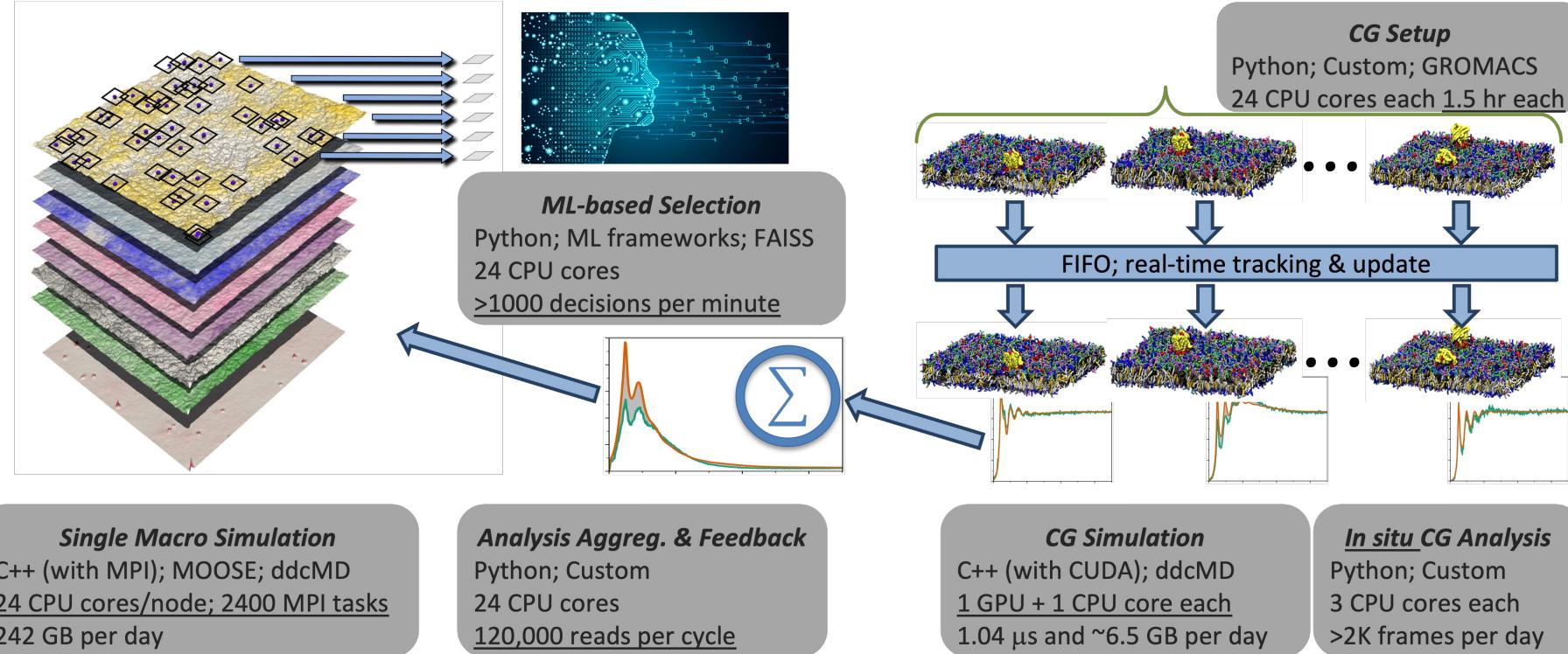
CG, analysis bound to cores
nearest PCIe buses

Job comms/coordination:

36,000 concurrent tasks;
176,000 cores, 16,000 GPUs

Portability:

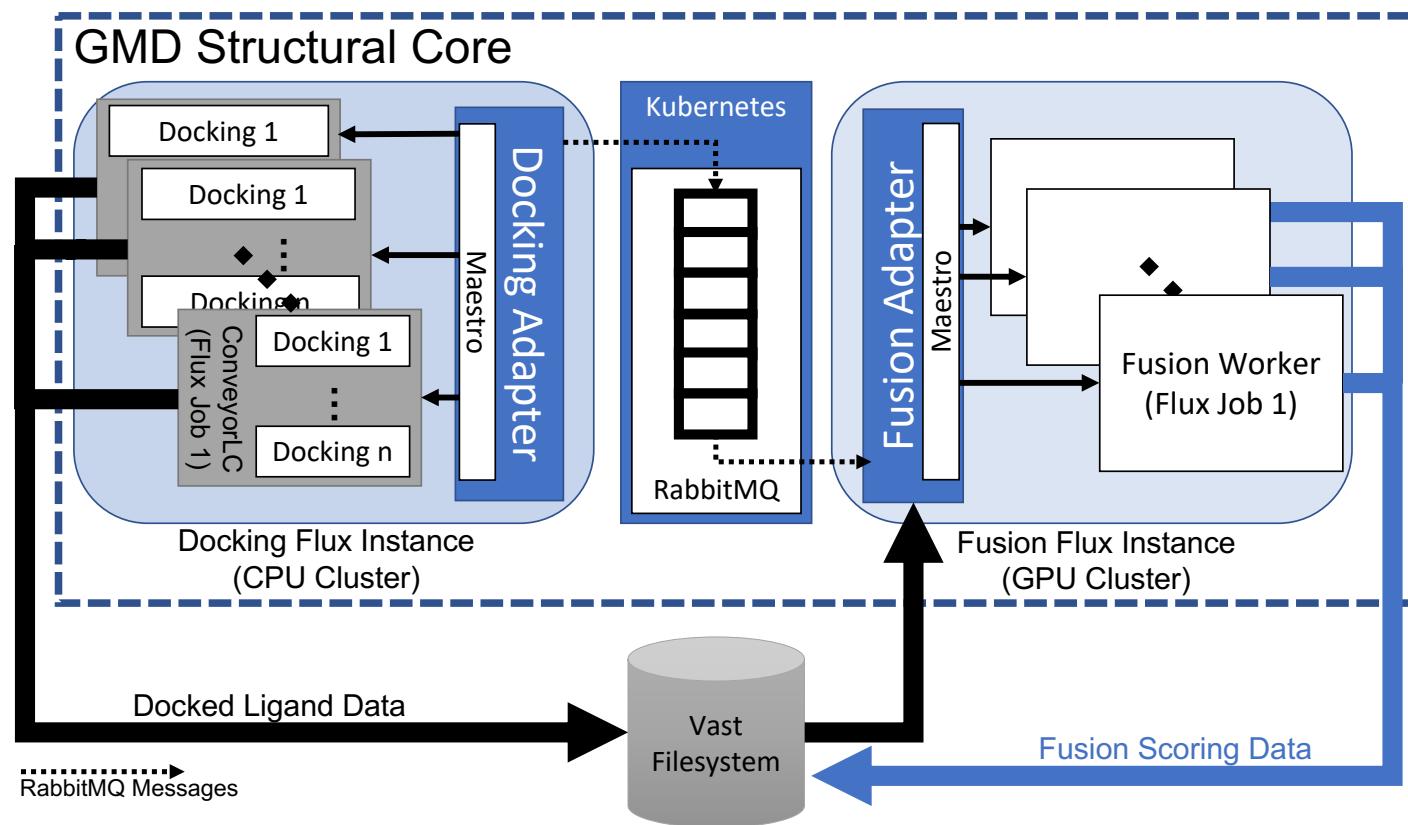
adapt tasks to different
Schedulers/managers



MuMMI: SC'19 best paper, SC'21 paper

MPI-based simulation with in-situ analysis plus ML

Next-generation, cross-cluster scientific workflows are demanding portability and cloud integration.



AHA MoLeS: eScience'22 best paper

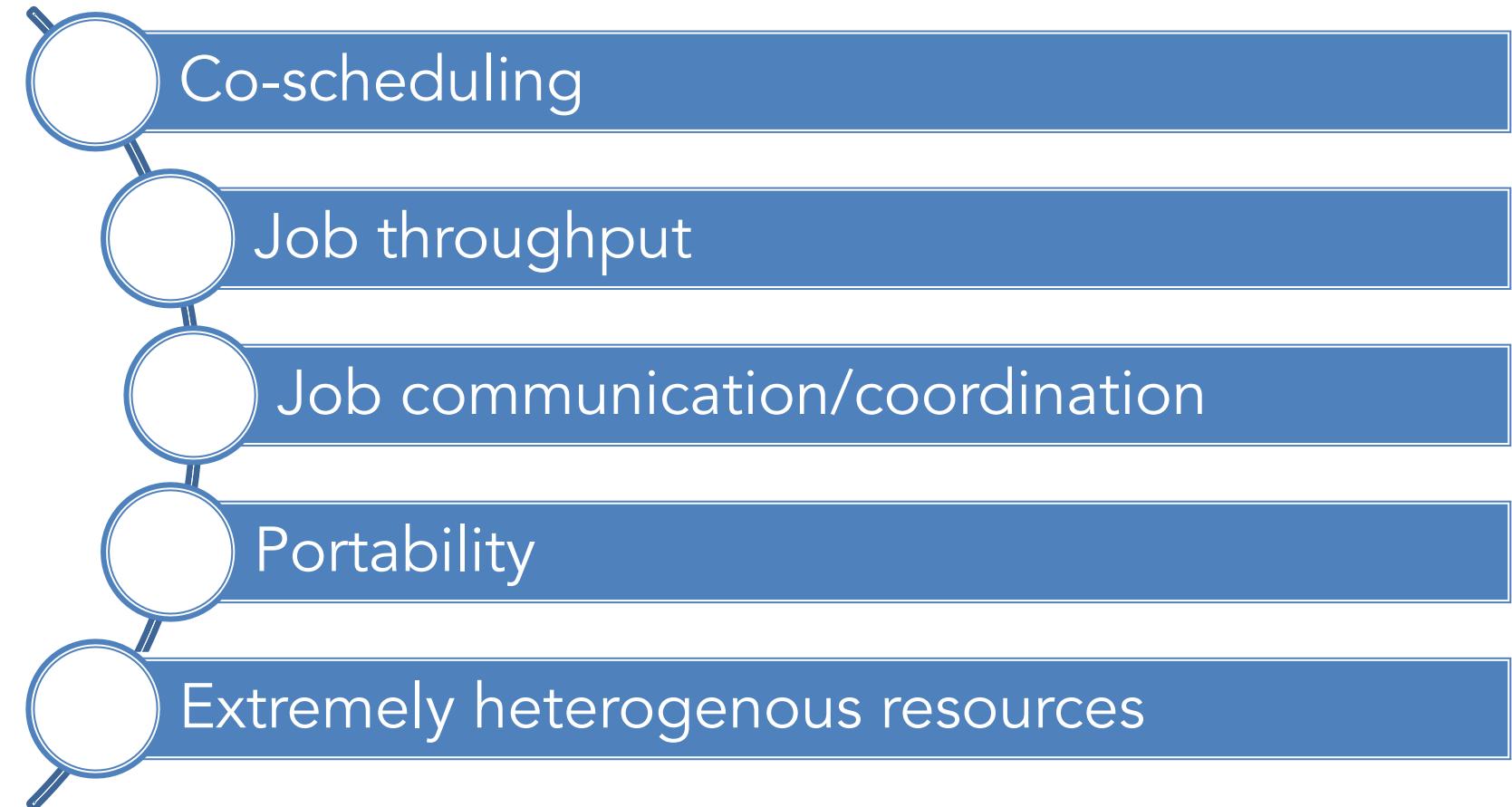
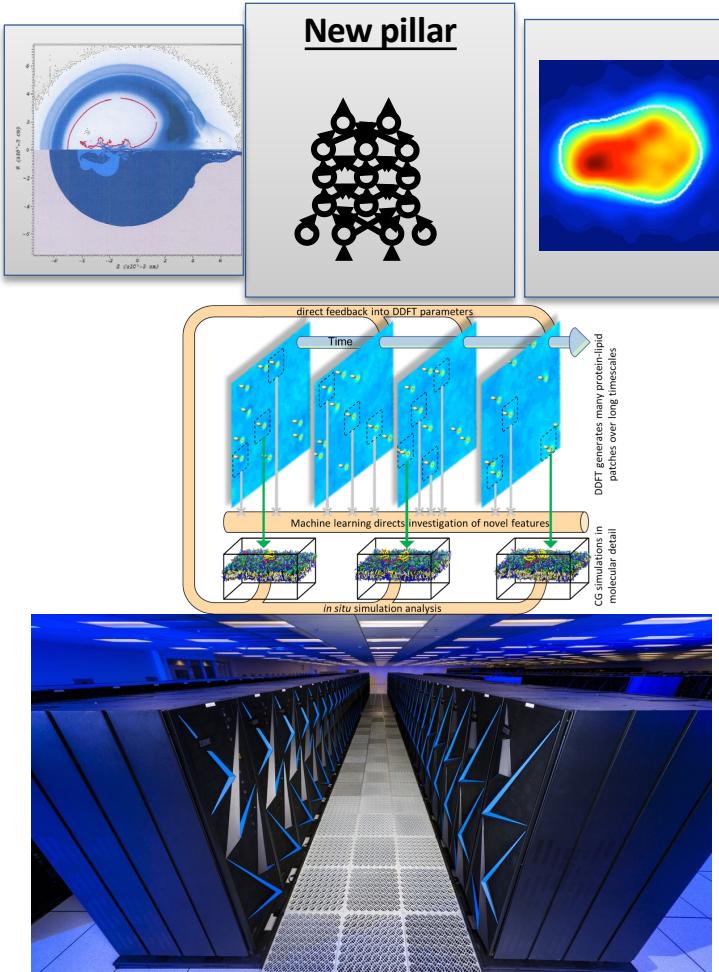
Complex workflows integrating cloud technologies at LLNL and beyond

- Scalable message broker couples MPI-based tasks, analysis, workflow runs anywhere (**AHA MoLeS**)
- HPC simulation with AI/ML surrogates, orchestrated databases (**AMS**)
- Many other examples: ATOM, AMPL, GMD, etc.

2020 lab survey found that 73% of LLNL workflows interested in cloud integration

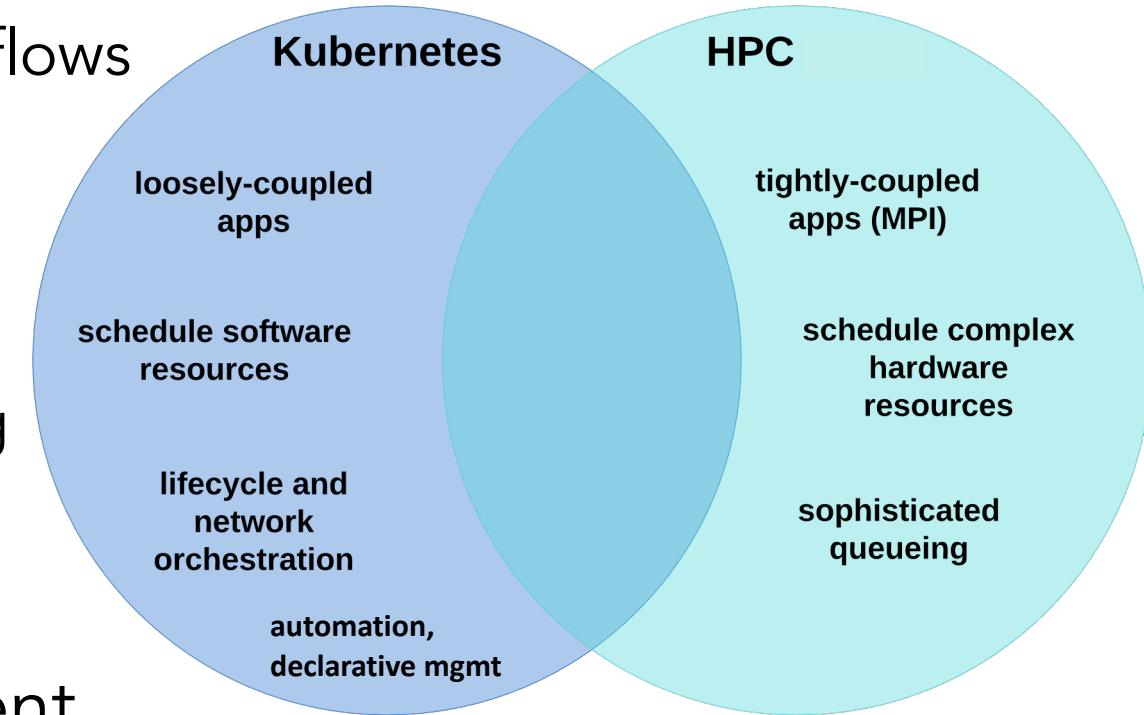
MPI-based simulation with analysis, ML, and containerized components

Trends towards complex workflows, extreme resource heterogeneity, and converged computing render traditional workload managers increasingly ineffective.



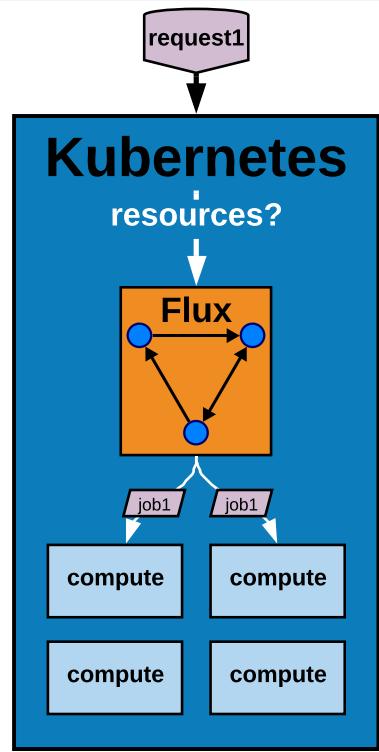
HPC Management and Scheduling and Kubernetes (K8s) are complementary technologies

- HPC limitations:
 - becoming very difficult to manage workflows
 - not designed for dynamism/elasticity
- K8s limitations:
 - designed for loosely coupled apps
 - not focused on performance (scheduling limitations, resource expression, throughput...)

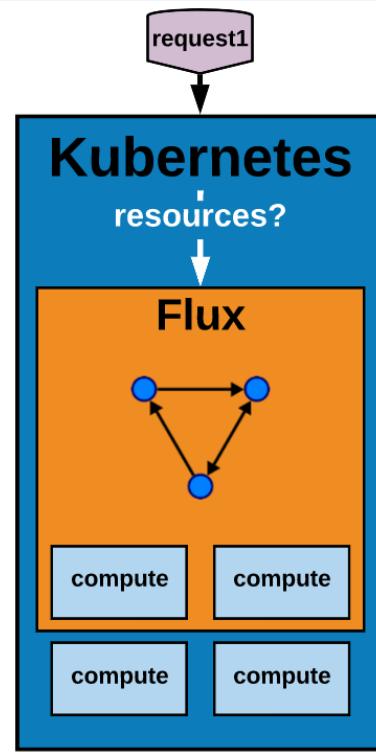


How to create a converged environment that combines the best of both worlds?

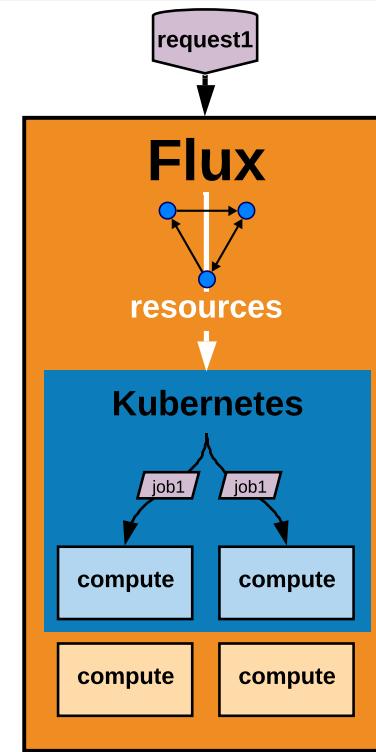
We are contributing to convergence by implementing Flux-based models.



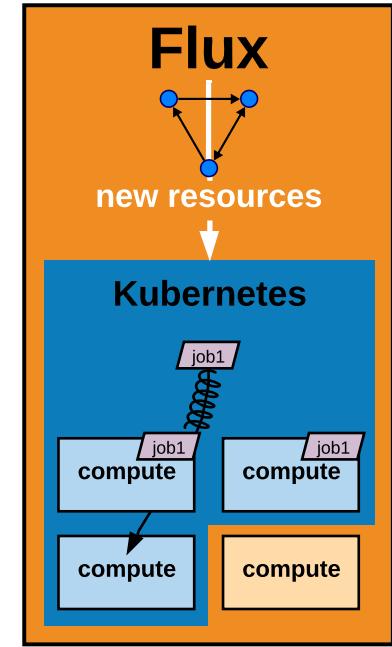
L1: standard APIs, resource representations, sched plugin. HPC nesting model enables portable converged workflows



L1.5: Flux manages, schedules cloud.
(portable converged workflows)



L2: resource co-management (requires L1). **Cloud nesting** model reduces software complexity, increases automation, performance

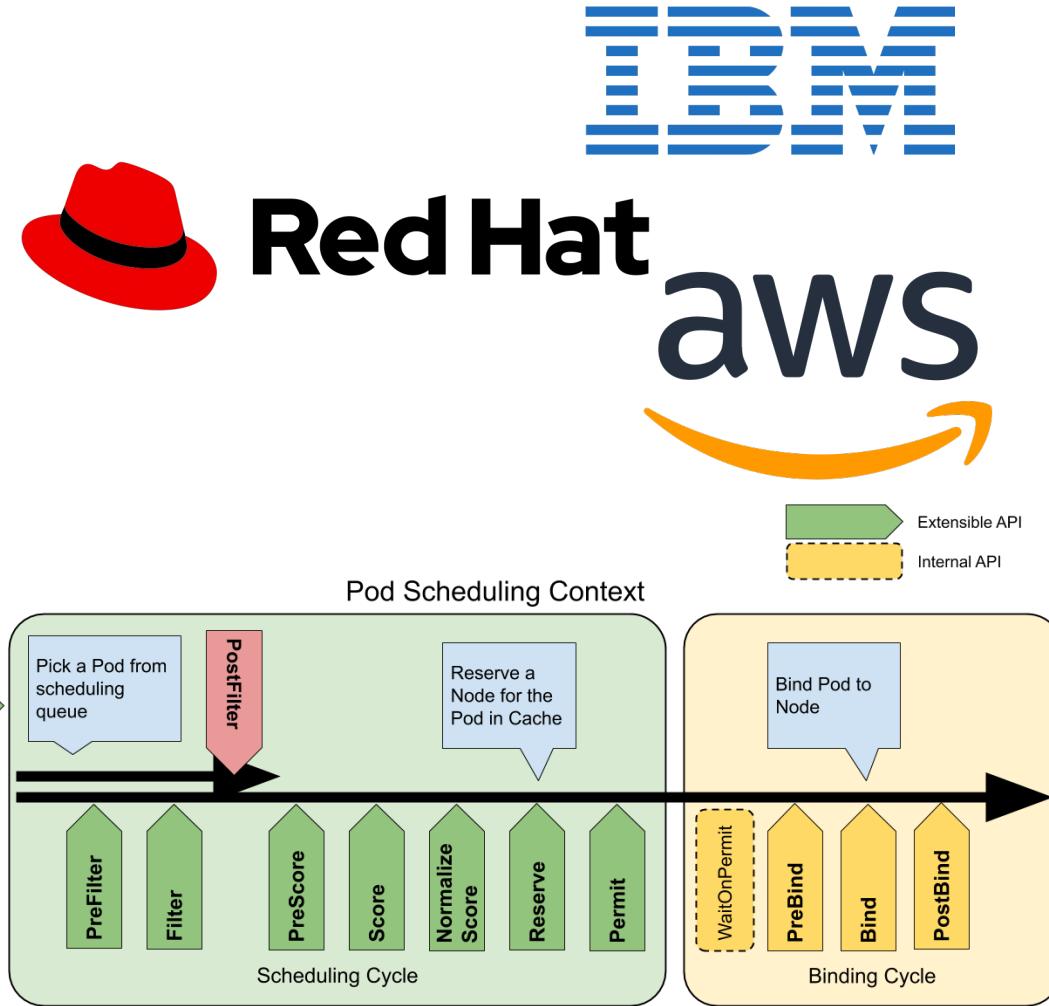


L3: elastic resource co-management (L1+L2)
Elastic cloud nesting model enables autoscaling and dynamism

Team formed industry collaborations to tackle converged computing challenges and contribute to community.

We bring complementary backgrounds in HPC, cloud, and performance-oriented orchestration

- LLNL LDRD team
- IBM T.J. Watson Research Center: T. Elengikal, M. Drocco, C. Misale, Y. Park
- AWS: E. Bollig, M. Hugues, H. Poxon, L. Wofford
- Integrate Fluxion into K8s via [Fluence](#)
- Declarative, automated Flux deployment in K8s via [Flux Operator](#)



The Fluence (FKA KubeFlux) plugin brings HPC-grade scheduling and improved performance to Kubernetes.

K8s Scheduling Framework plugin based on Fluxion scheduler.

Architectural change from monolithic to gRPC-based

- Improves maintainability, separation of concerns

More placement control and functionality

- Gang scheduling
- GPU support
- Topology awareness of Availability Zones (AZs)

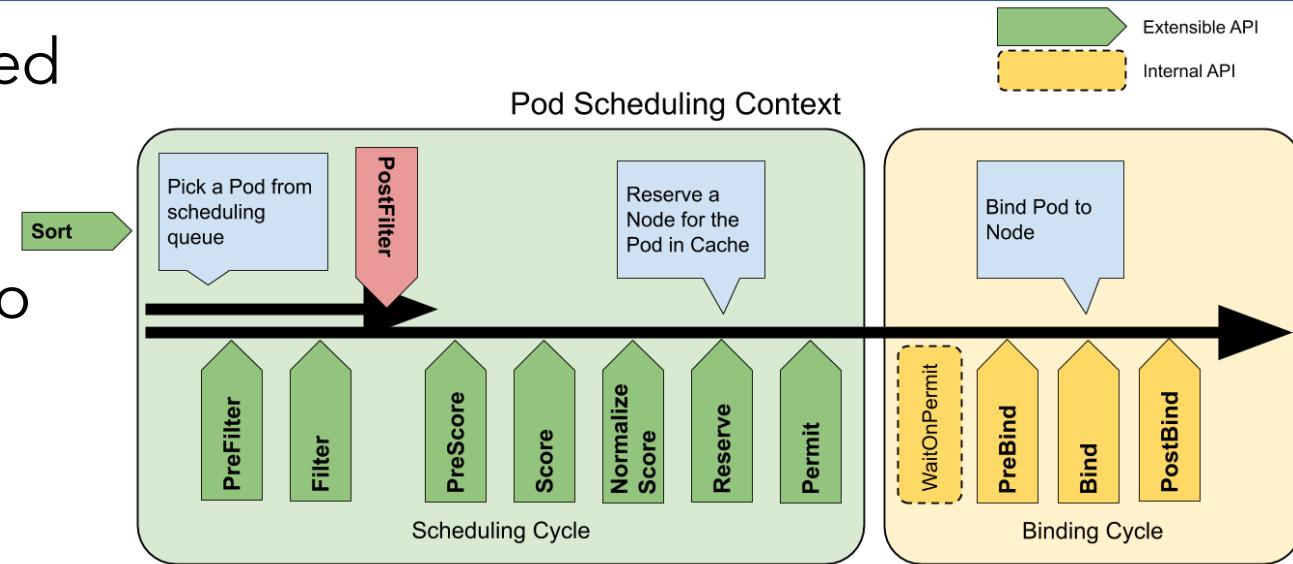
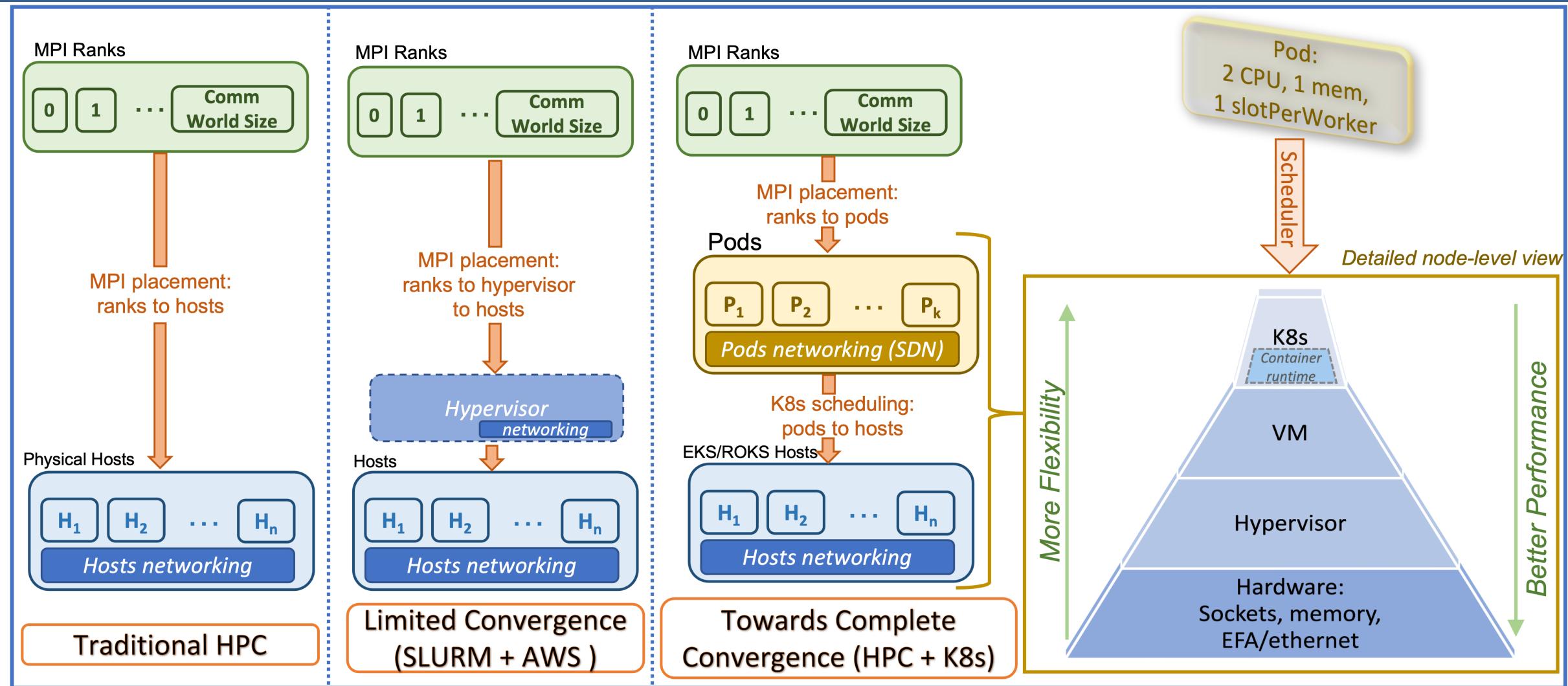


image: <https://kubernetes.io/docs/concepts/scheduling-eviction/scheduling-framework/>

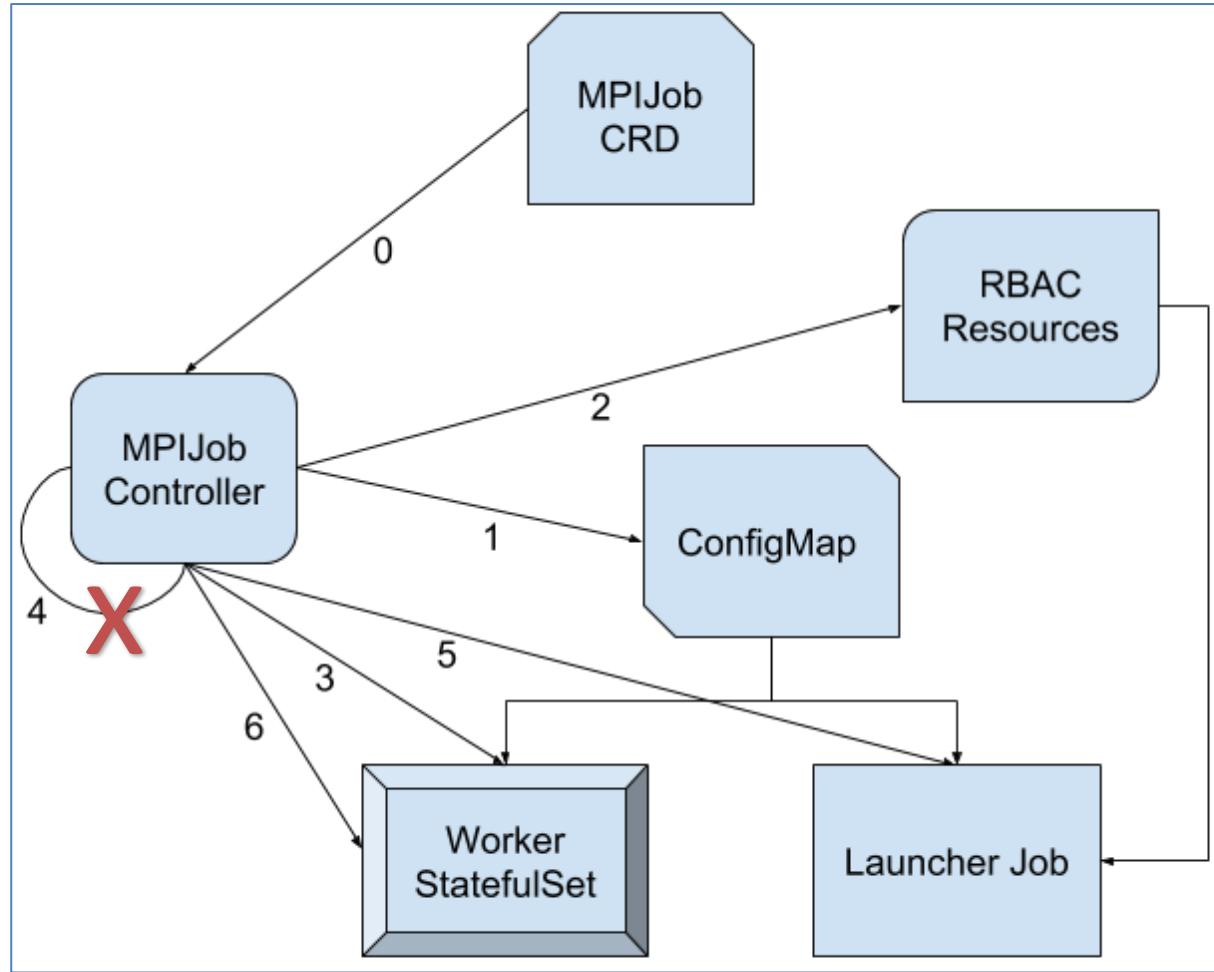
Easier deployment

- Automation through Helm
- Export of Golang modules for easier distribution

L1-1.5: creating HPC cloud “slice” is complicated and requires scalable MPI bootstrap and pod placement.



We enabled cloud-native, declarative MPI to scale three orders of magnitude higher.



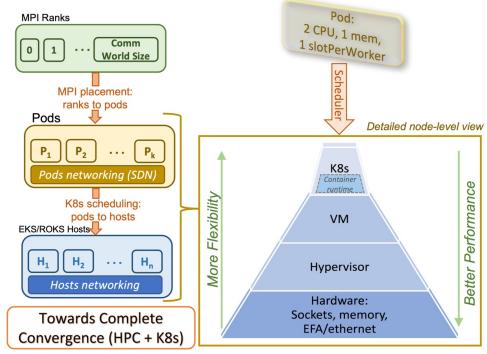
KubeFlow MPI Operator created to manage MPI Jobs; only scaled to around 80 ranks

- Race conditions in Controller (steps 1-4): Launcher doesn't wait for Workers, Workers report **READY** before ssh available
- Fixed DNS flood

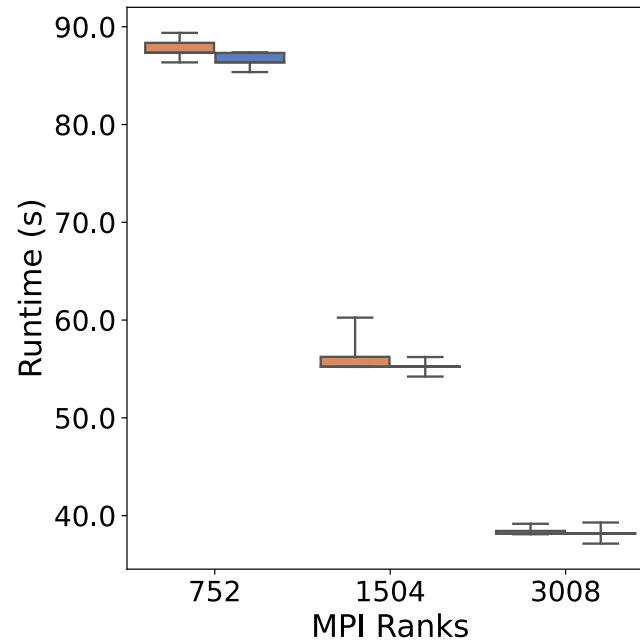
Tested fixed MPI Operator with *Hello World* at 16,384 ranks on AWS EKS

image source Kubeflow: <https://medium.com/kubeflow/introduction-to-kubeflow-mpi-operator-and-industry-adoption-296d5f2e6edc>

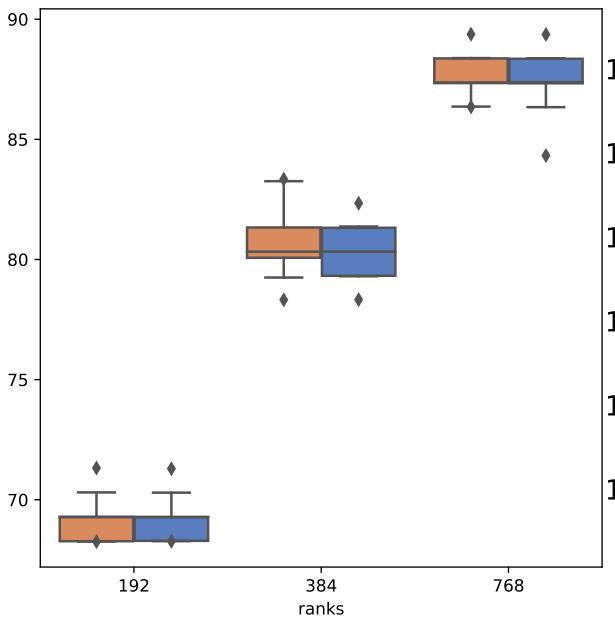
We ran CORAL-2 benchmarks on AWS Elastic Kubernetes Service with the improved MPI Operator.



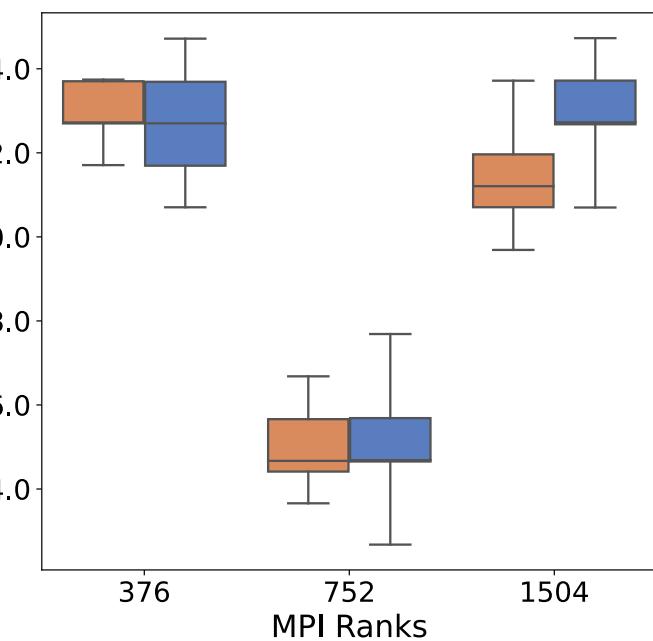
Control pod shape via resource request, MPI shape via ***slotsPerWorker***, ***mpirun*** options



LAMMPS



AMG



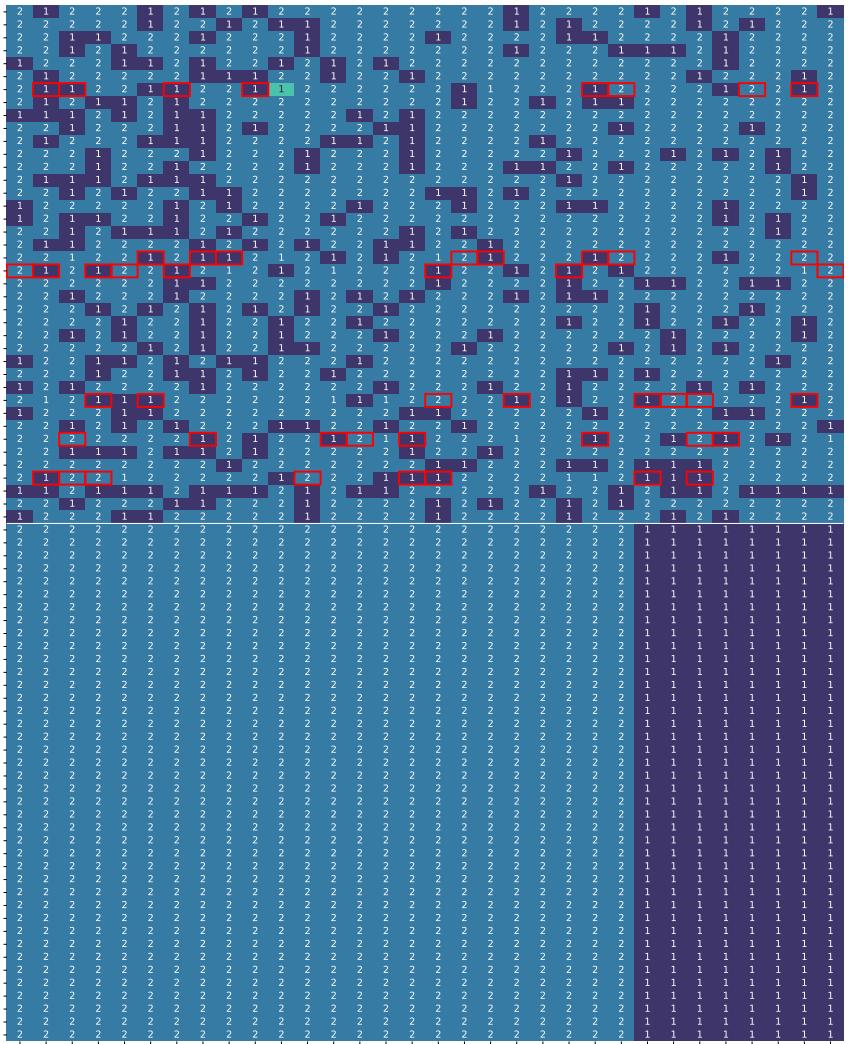
QCMPACK

32 Hpc6a instances with 94 usable CPUs, 360 GB memory, Elastic Fabric Adapter (EFA)

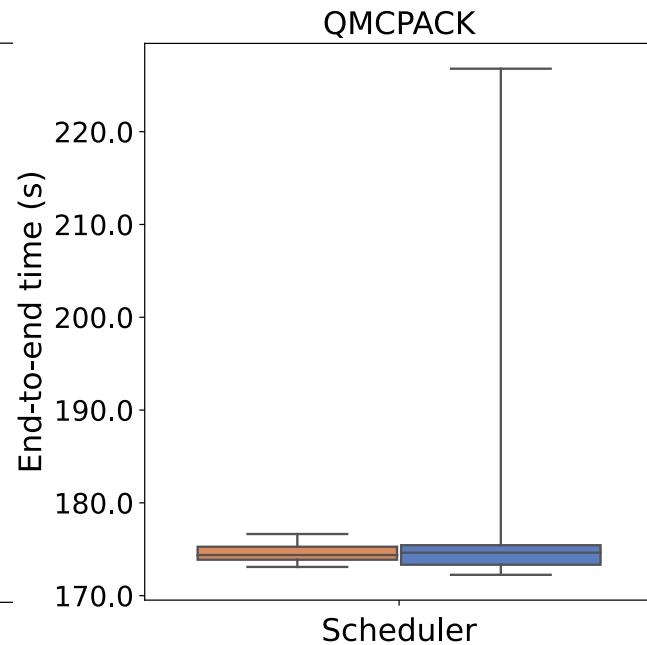
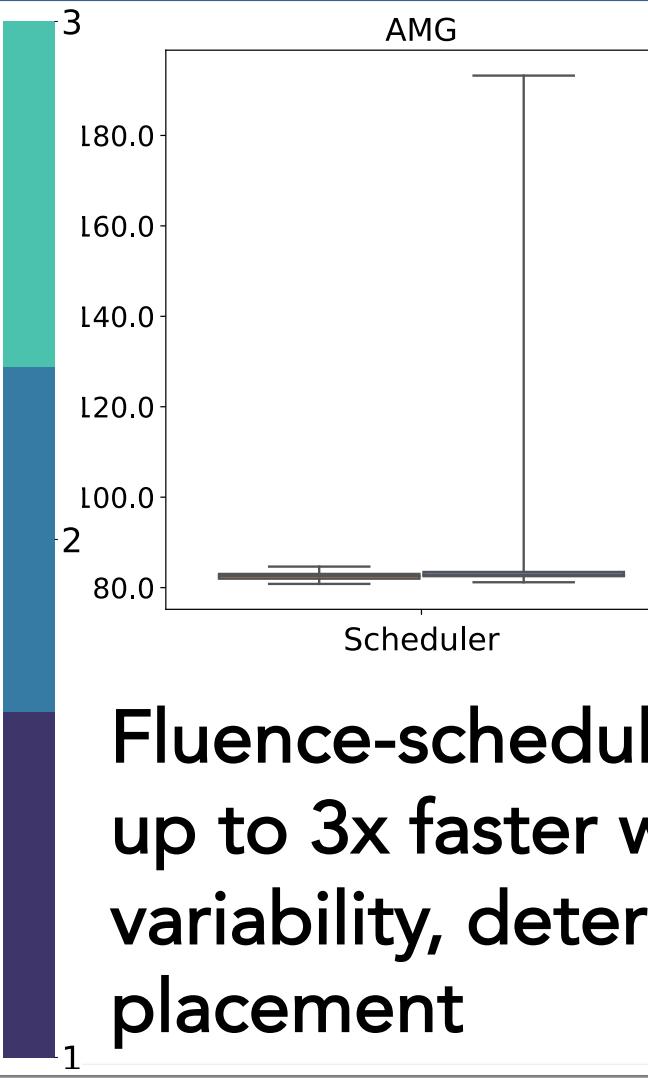
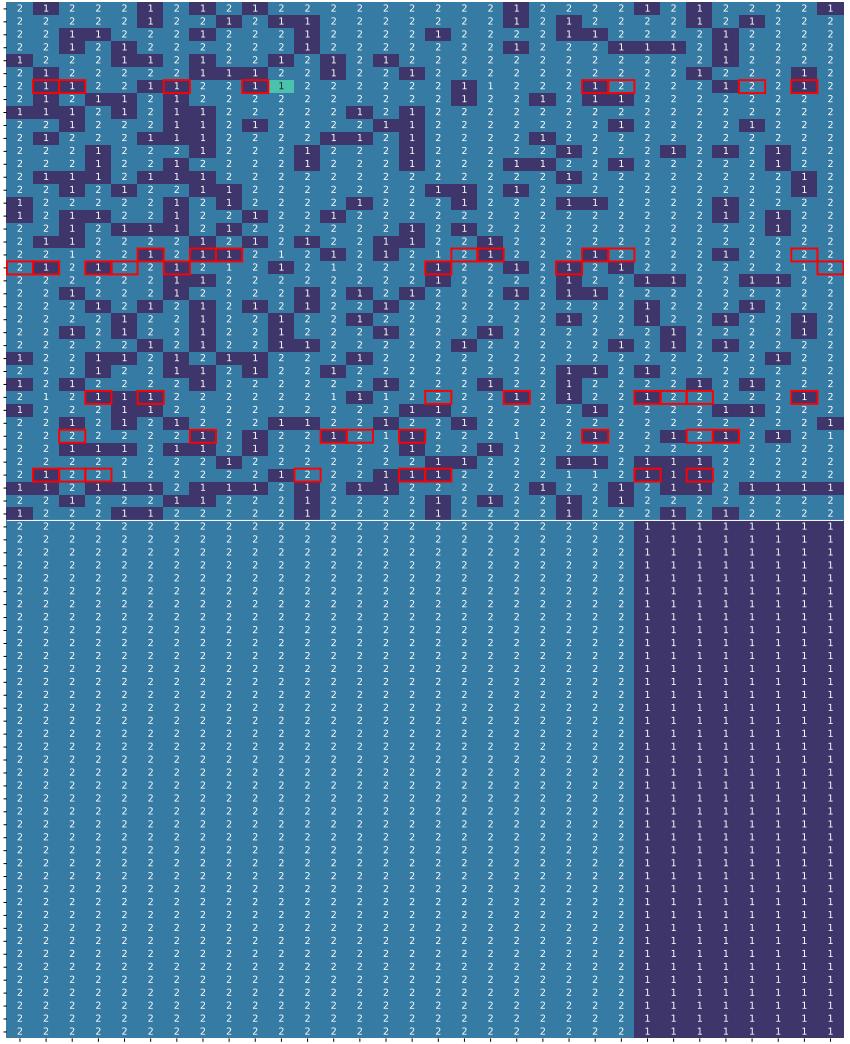
- LAMMPS strong scales 94 ranks/node with EFA up to 3008 ranks
- AMG exhibits decent weak scaling 24 ranks/node with and without EFA
- QCMPACK exhibits inverse (!) scalability after 8 nodes

We compared multi-app simulated workflow performance scheduled by Fluence vs. Kube-scheduler.

Kube-scheduler



Fluence

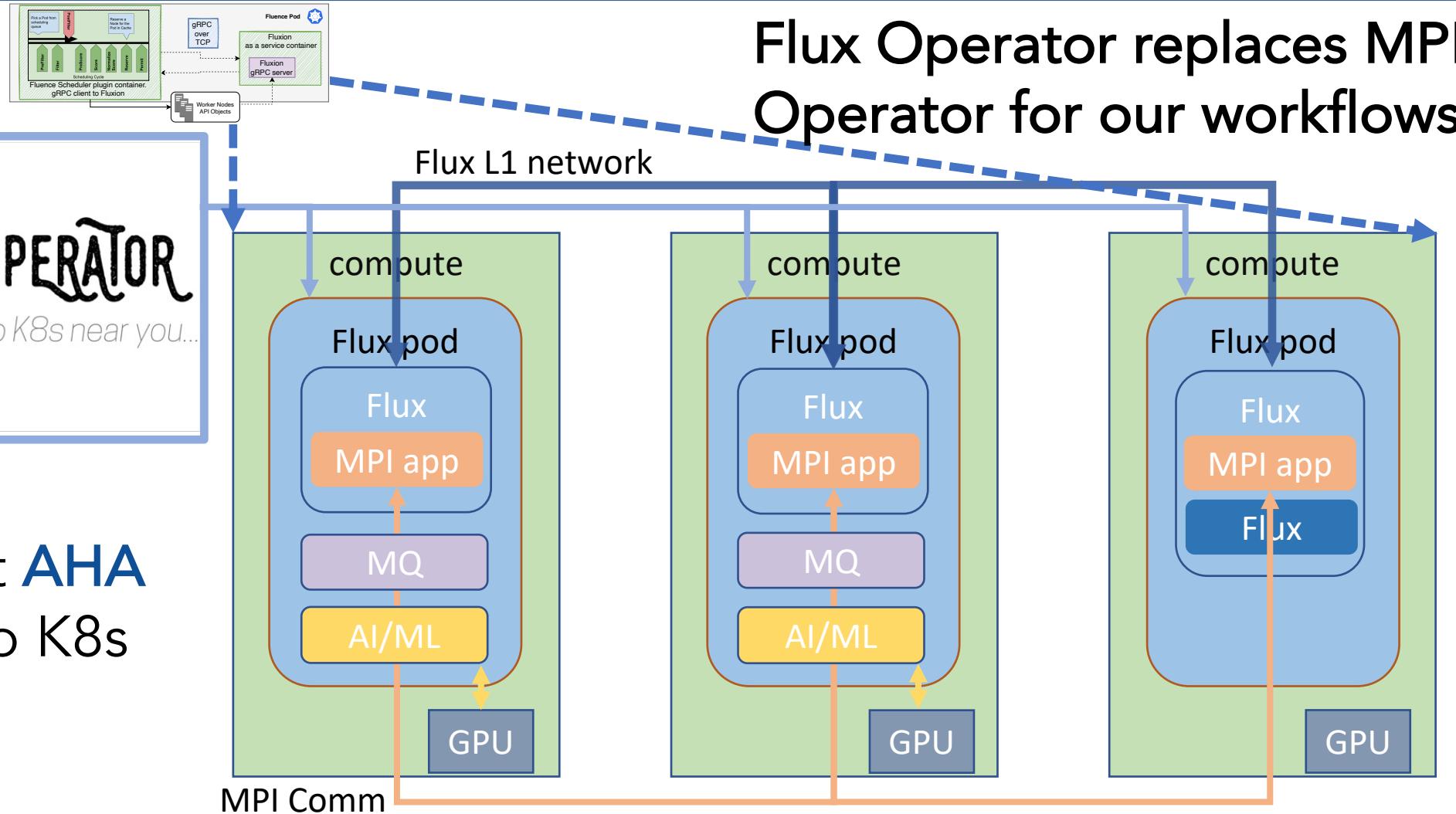


Fluence-scheduled workflows run up to 3x faster with low variability, deterministic placement

L1.5: Fluence + Flux operator will enable hierarchical management with HPC-grade pod scheduling.



Next up: port **AHA**
MoleS, etc. to K8s

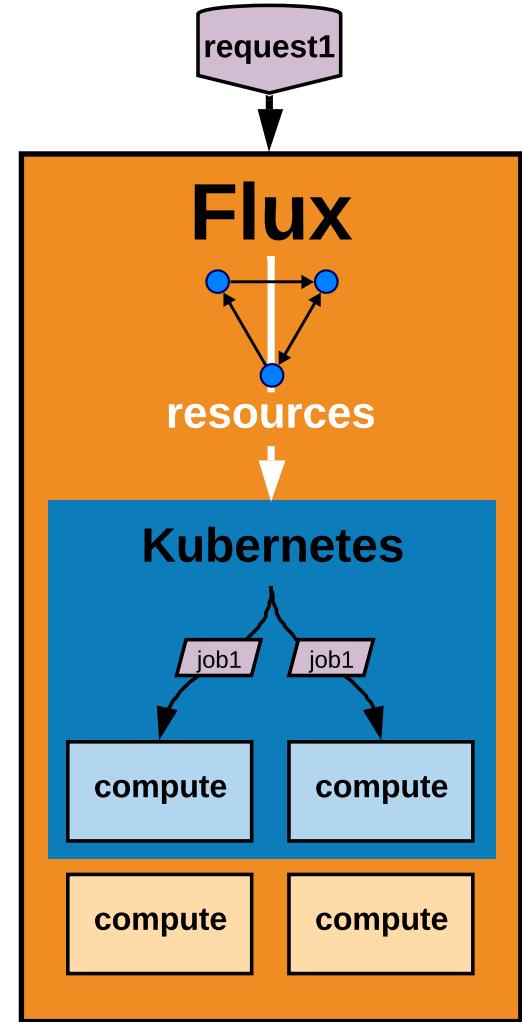


L2: Develop an effective and efficient resource co-management framework

Challenges:

- which manager has authority over resources; how?
- how to minimize/mitigate contention?
- consistency model? (potentially millions of vertices, edges, and thousands of jobs)
- how to deliver high performance?

Benefit: enables tighter coupling between HPC and K8s components, greater automation and reduced complexity via, e.g. Operators



There is still much research to be done to make resource dynamism and elasticity useful as embodied by Flux.

- Flux needs capability to reshape/rebalance its ranks and network
- Automatically resize Fluence resource graph when new K8s resources added; grow and shrink allocations
- Algorithmic enhancements to Fluxion to support elastic jobs
 - Scheduler driven grow/shrink: malleability
 - Application driven grow/shrink

Thank you!
Questions?

DYAD (DYnamic and Asynchronous Data-streamliner)

Jae-Seung Yeom*, Dong H. Ahn[†], Ian Lumsden[‡], Jakob Luettgau[‡], Silvina Caino-Lores[‡], Michela Taufer[‡]

*Lawrence Livermore National Laboratory,

[†]NVIDIA Corporation,

[‡]Department of Electrical Engineering and Computer Science, University of Tennessee, Knoxville



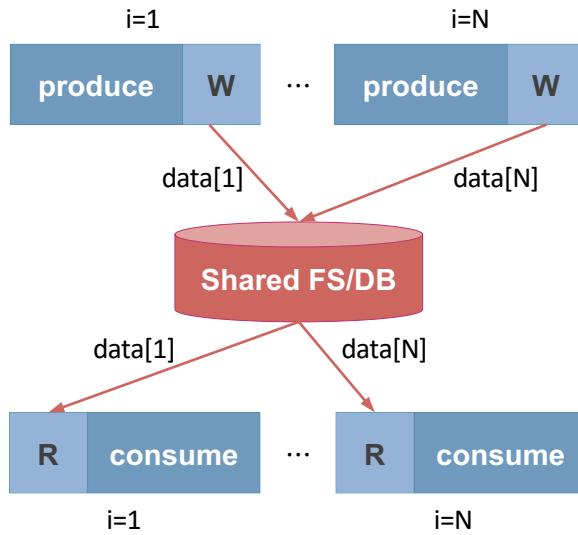
Resolving inter-task data dependence in a workflow

- Existing models
 - *in situ, in-transit, shared file system and so on*
 - Varying degree of performance, user productivity, portability and flexibility
- Challenges
 - Lack of synchronization support at the file/data object level
 - Conflict with code change requirements
 - Poor temporal/spatial locality
 - Low file metadata-operation performance

Benefits we aim to bring in

- No/little code change
 - Productivity (Fast construction of workflows with less effort)
 - Portability
 - Data sharing based on universal file abstraction
 - Open-sourced and proprietary applications
 - widely used POSIX IO instead of specific API
- Performance
 - Use of local storage enable faster accesses
 - Fine-grain file level synchronization exposes further parallelism

Producer-consumer patterns are common in modern workflows



```
void producer() {  
    for (i=1; i <= N; i++) {  
        produce(data[i])  
        write(data[i]);  
    }  
}
```

```
void consumer() {  
    for (i=1; i <= N; i++) {  
        read(data[i]);  
        consume(data[i]);  
    }  
}
```

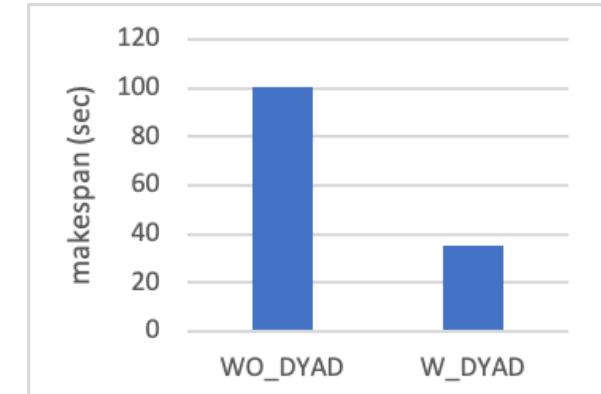
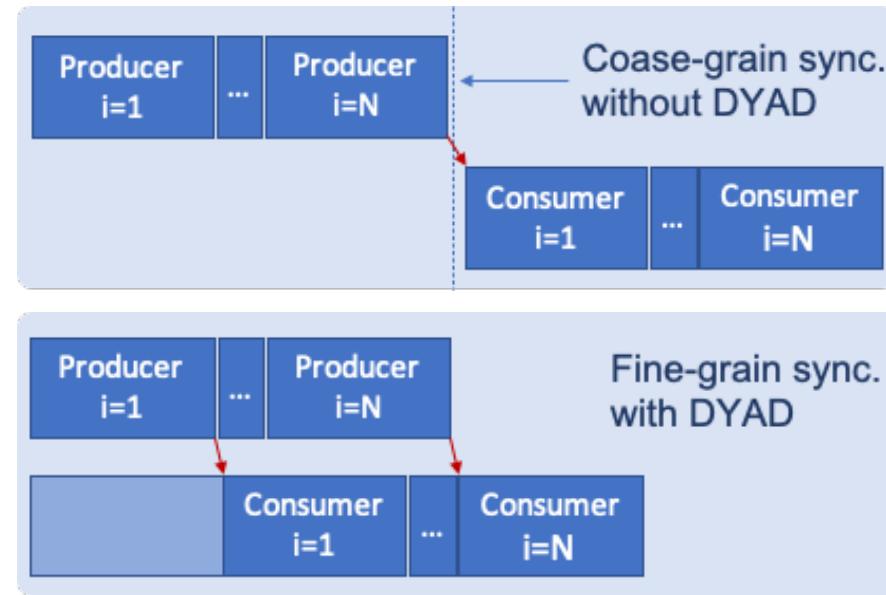
(a) Producer-consumer example

- description: producer task
name: run-producer
run:
cmd: producer
- description: consumer task
name: run-consumer
run:
cmd: consumer
depends: [run-producer]

(b) Maestro Specification in YAML

Fine-grain synchronization with DYAD

```
producer() {  
    for (i=1; i<=N; i++) {  
        produce(data[i])  
        write(data[i]);  
    }  
}  
  
consumer() {  
    for (i=1; i<=N; i++) {  
        read(data[i]);  
        consume(data[i])  
    }  
}
```



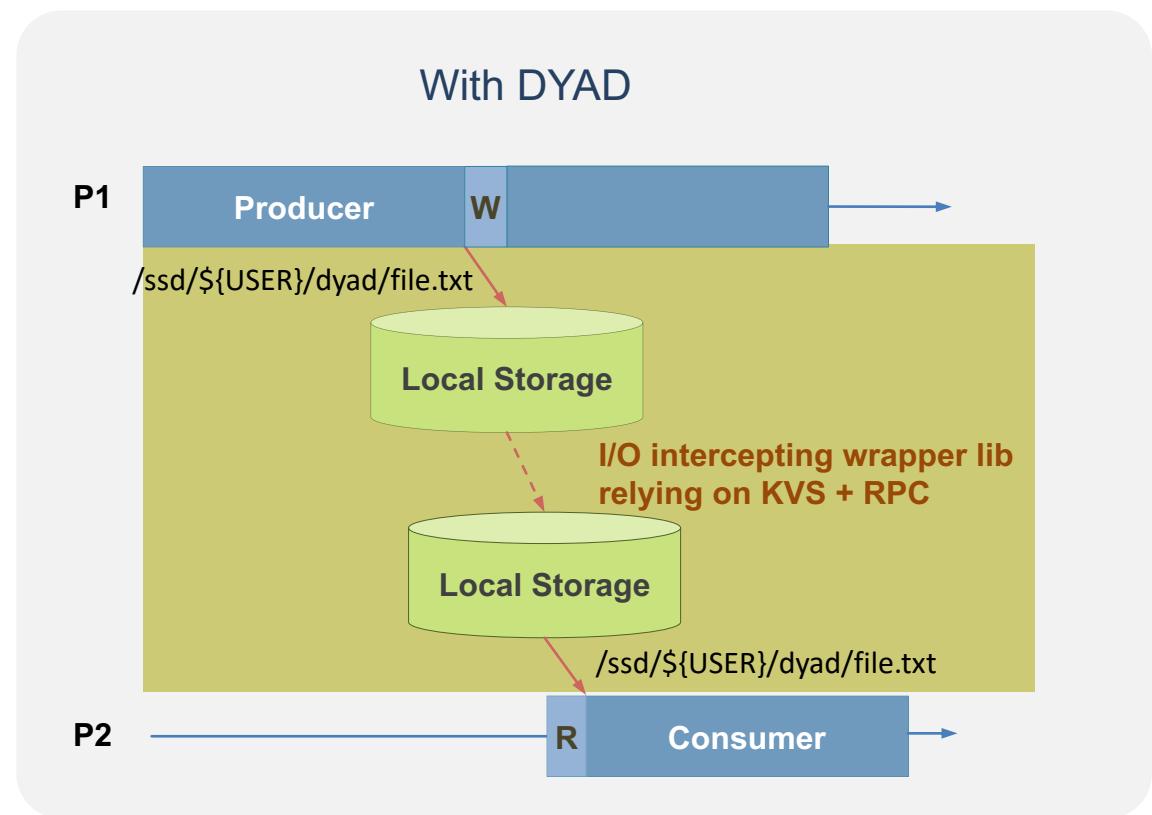
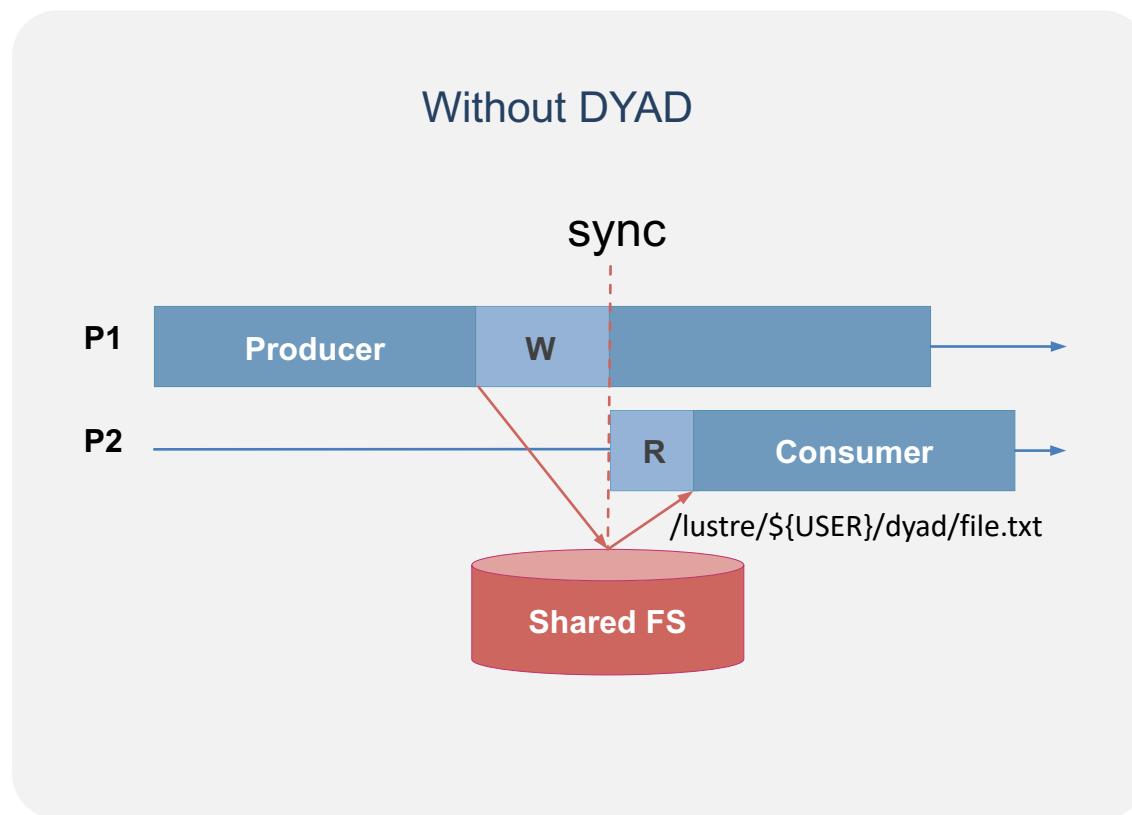
- Explore the effectiveness of the scheduling to minimize the temporal and the spatial distance between producers and consumers

DYAD operation

Explicit synchronization between read and write

- e.g., dependent jobs, API, or polling

- Relies on local storages
- Transparent data transfer between storages and synchronization per shared file



Minimal example

- Running DYAD service with FLUX

```
flux exec -r all flux module load dyad.so /ssd/${USER}/dyad
```

```
flux exec -r all flux module load dyad.so /ssd/${USER}/dyad
```

- Configure DYAD via env variables

```
DYAD_KVS_NAMESPACE  
DYAD_PATH_PRODUCER  
DYAD_PATH_CONSUMER  
DYAD_SHARED_STORAGE
```

- Running user app written in C with DYAD

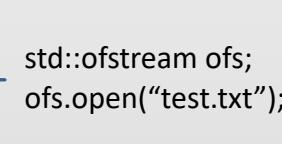
```
LD_PRELOAD=libdyad_sync.so:${LD_PRELOAD} ./app
```

- Using DYAD stream library in C++

```
#include <dyad_stream_api.hpp>

void producer(dyad::dyad_stream_core& dyad)
{
    dyad::ofstream_dyad ofs_dyad;
    ofs_dyad.init(dyad);
    ofs_dyad.open("test.txt");
    std::ofstream& ofs = ofs_dyad.get_stream();
    ofs << "test" << std::endl;
    ofs_dyad.close();
}

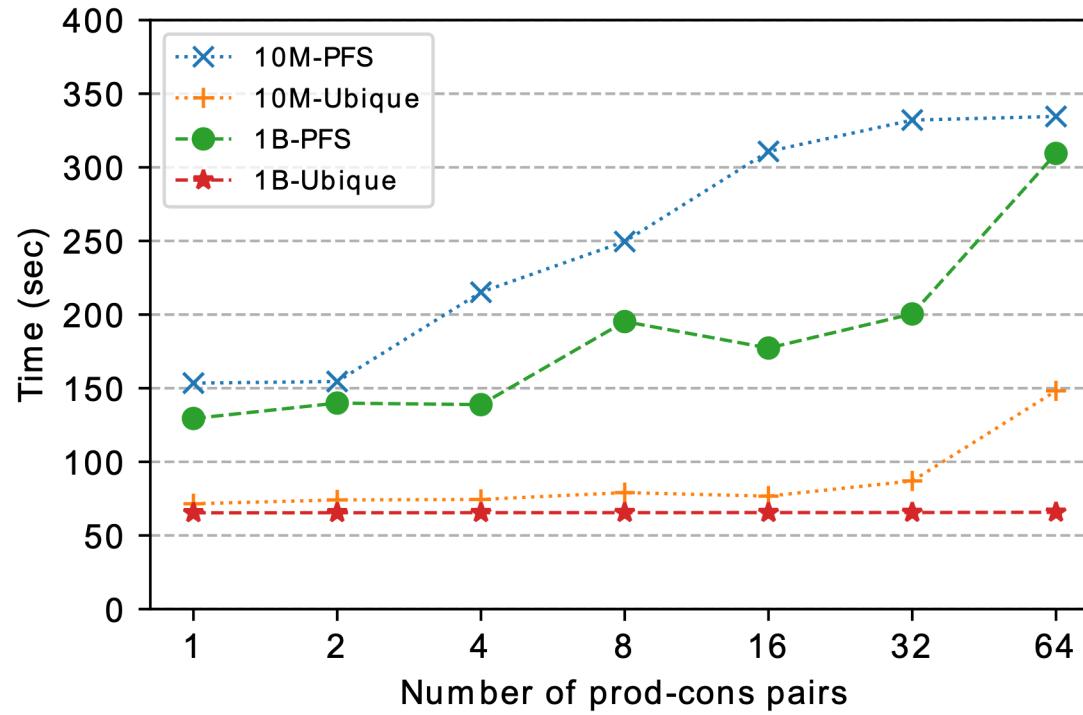
void consumer(dyad::dyad_stream_core& dyad)
{
    dyad::ifstream_dyad ifs_dyad;
    ifs_dyad.init(dyad);
    ifs_dyad.open("test.txt");
    std::ifstream& ifs = ifs_dyad.get_stream();
    std::string line;
    ifs >> line;
    std::cout << line << std::endl;
    ifs_dyad.close();
}
```



DYAD: DYnamic and Asynchronous Data-streamliner

- Handles producer-consumer I/O pattern in a workflow
- Take advantage of local storage
 - Transfer data transfer directly from producers to consumers
- Fine-grain synchronization
 - At file open/close granularity
- Enable dynamic sharing (dynamic dependency)
- Flux API-based tool (KVS, RPC)

Performance of producer-consumer benchmark

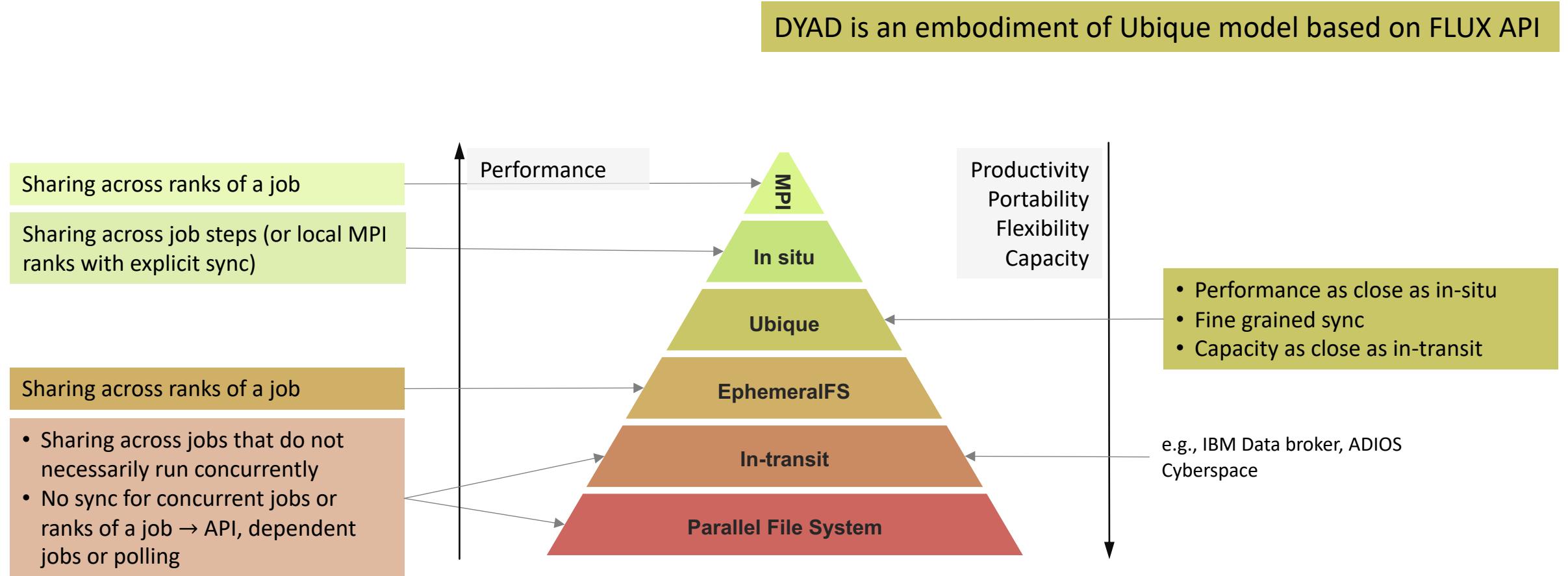


File sharing by DYAD vs parallel file system (Lustre/quartz)

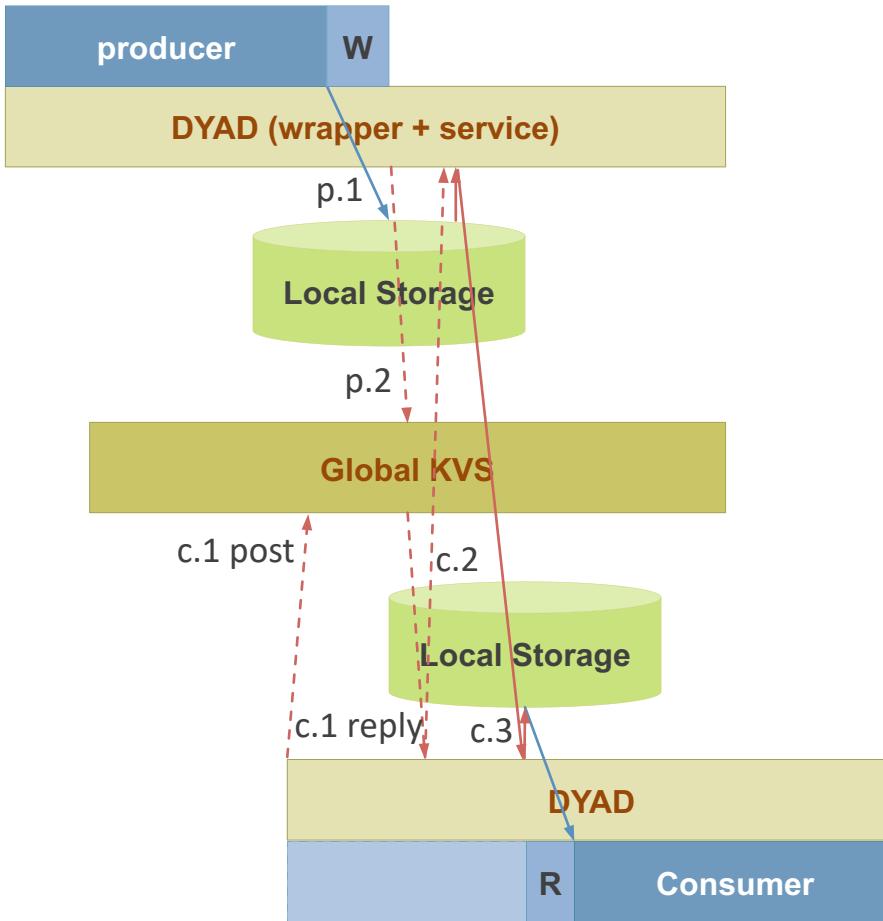
```
void producer() {  
    for (i=1; i <= N; i++) {  
        produce(data[i]);  
        write(data[i]);  
    }  
}
```

```
void consumer() {  
    for (i=1; i <= N;i++) {  
        read(data[i]);  
        consume(data[i])  
    }  
}
```

Design trade-off space for data transfer model



Synchronization and file transfer



DYAD

- DYAD module (or service) runs on each node.
- DYAD wrapper only intercepts I/O for files on the directory it manages
 - If managed dir is on local storage (LS), sync accesses and transfer files
 - If it is on shared storage, only sync accesses.

Producer

- p.1 *write(manged_dir/filename)*
- p.2 *publish(<filename, prod_rank>)*
 - If a file is written into **managed_dir** or its subdirectory, DYAD registers the filename into the global key-value-store (KVS) managed by FLUX
 - KVS entry is a pair of **filename** and **rank**, where rank is the FLUX rank which owns the file and sees it locally

Consumer

- c.1 *query(filename) → prod_rank*
 - Consumer queries KVS to obtain the rank of the file owner (producer). Then, blocking wait.
- c.2 *rpc_get(prod_rank, filename)*
 - Ask the owner rank to transfer the file. Once received, store it on LS
- c.3 *read(manged_dir/filename)*

Scaling the performance of key value store (KVS)

- KVS
 - FLUX module (service)
 - The larger the number of keys, the slower it gets to complete a search/update
- Namespace
 - Within a namespace, users can get/put KVS values completely independent of other KVS namespaces.
 - DYAD generates queries under the namespace given by DYAD_KVS_NAMESPACE env variable
 - Needs to execute `flux kvs namespace create \${DYAD_KVS_NAMESPACE}` in advance
- Hierarchy
 - DYAD creates a hierarchical key string for a given filename
 - env variables DYAD_KEY_DEPTH and DYAD_KEY_BINS to control
 - Murmur3 hash vectorized for 128 bit
- Remove keys that are no longer used
 - `flux kvs namespace remove \${DYAD_KVS_NAMESPACE}`

$key \leftarrow bin_1.bin_2.\dots.bin_d.filename$

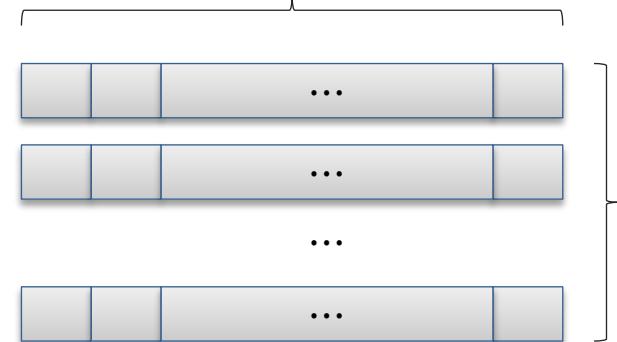
$bin_1 \leftarrow \text{hash}(filename, seed_1)$

$bin_2 \leftarrow \text{hash}(filename, seed_2)$

...

$bin_d \leftarrow \text{hash}(filename, seed_d)$

DYAD_KEY_BINS



DYAD_KEY_DEPTH

Future work

- Generalize the supporting environment
 - Leveraging external data transport layers
 - for HPC and Cloud
 - Leveraging external meta-data management
- Taking advantage of deep storage hierarchy
- Data-locality-aware scheduling
- Data lifespan and ownership management

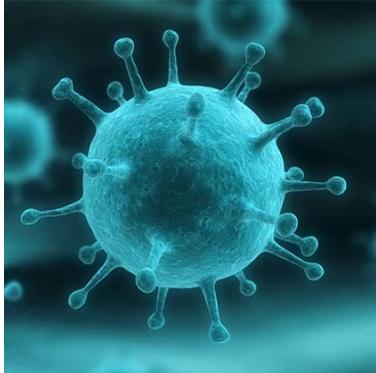
AHA MoleS: Accelerate Drug Discovery by American Heart Association Molecular Screening Workflow

Xiaohua Zhang
BBTD, PLS, LLNL



We need to be prepared for new emerging threats

HIV



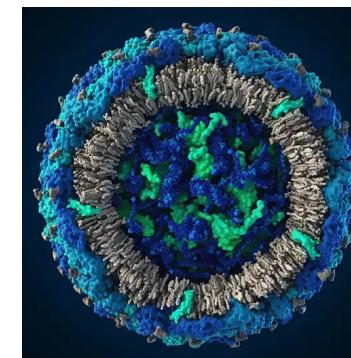
Superbug



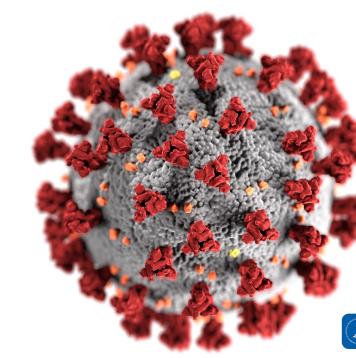
Ebola



Zika



Covid-19

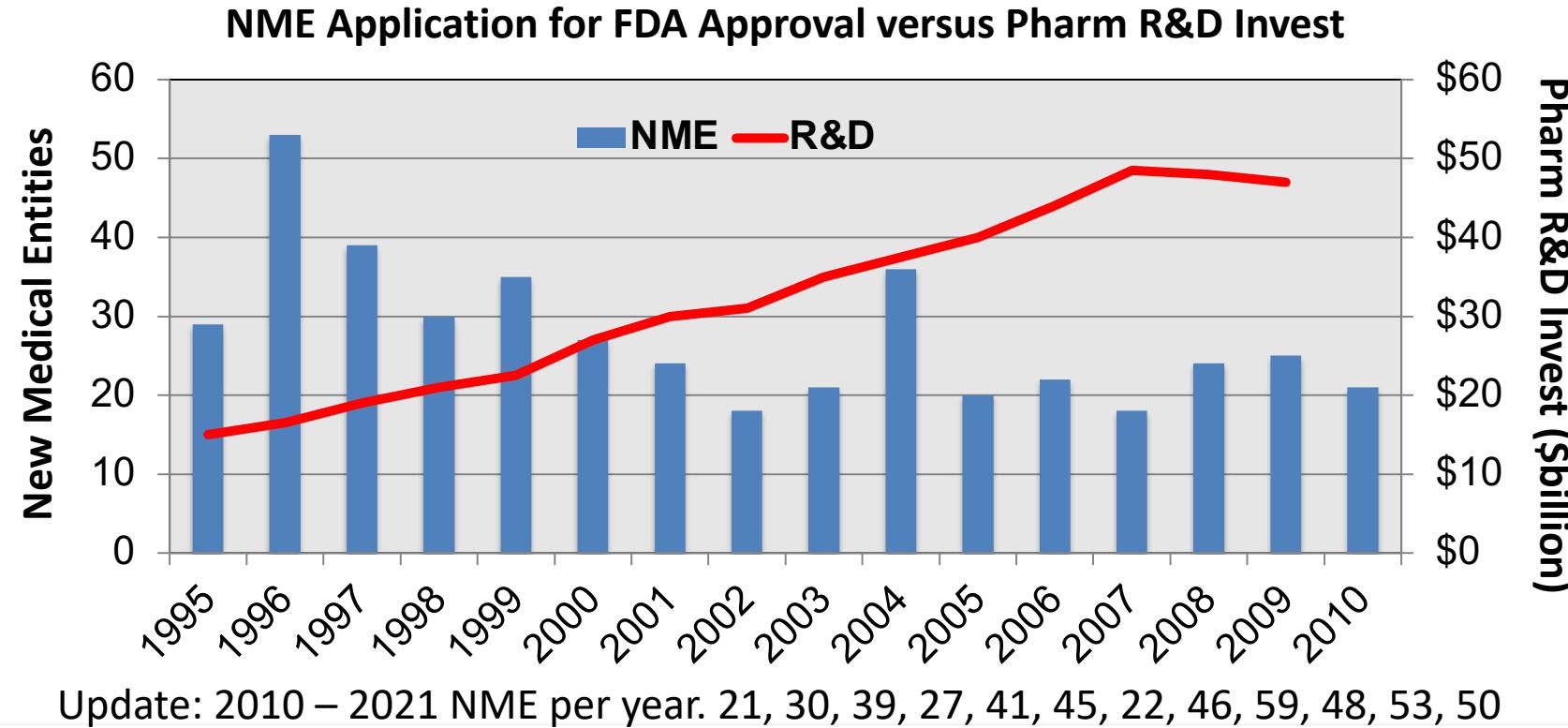


Next?

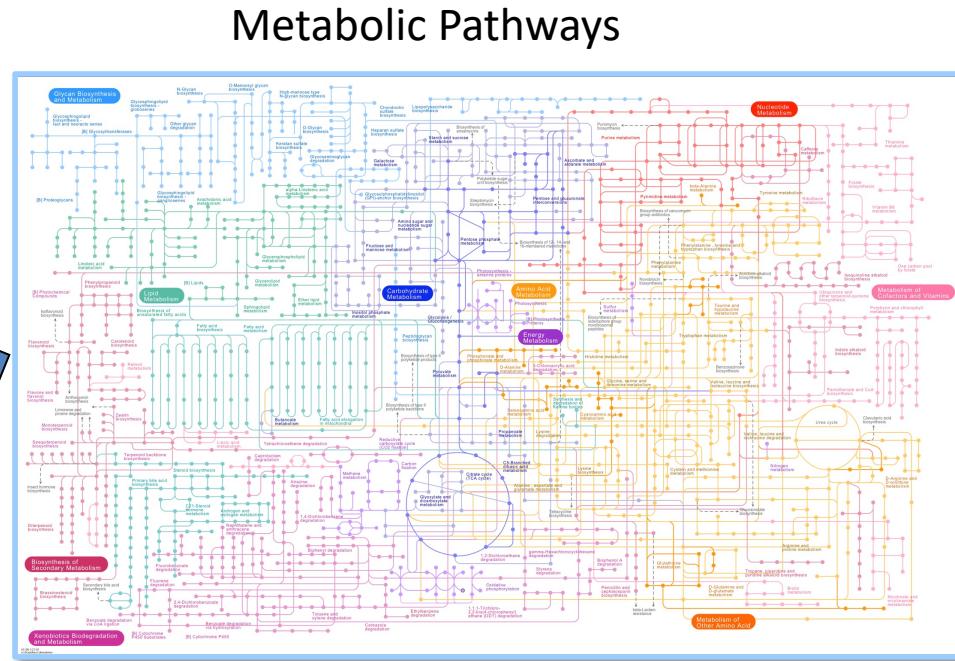
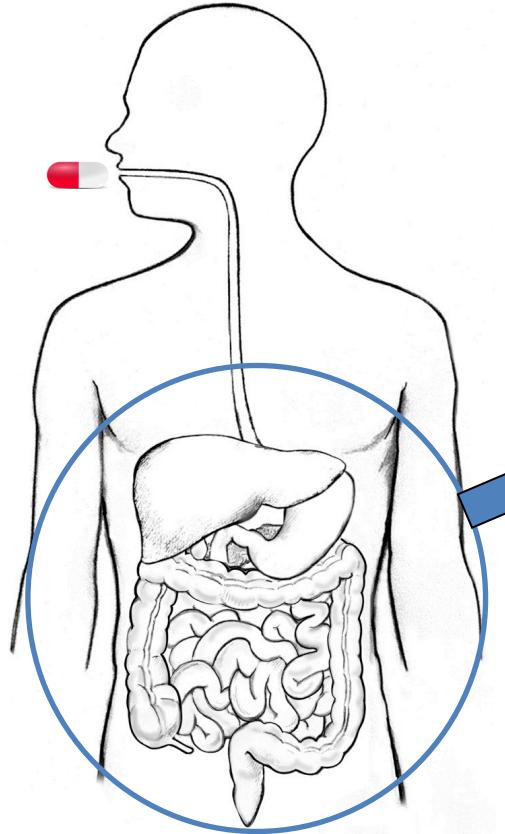


Urgent need for better *in-silico* prediction tools to reduce the cost and failure rate of drug discovery

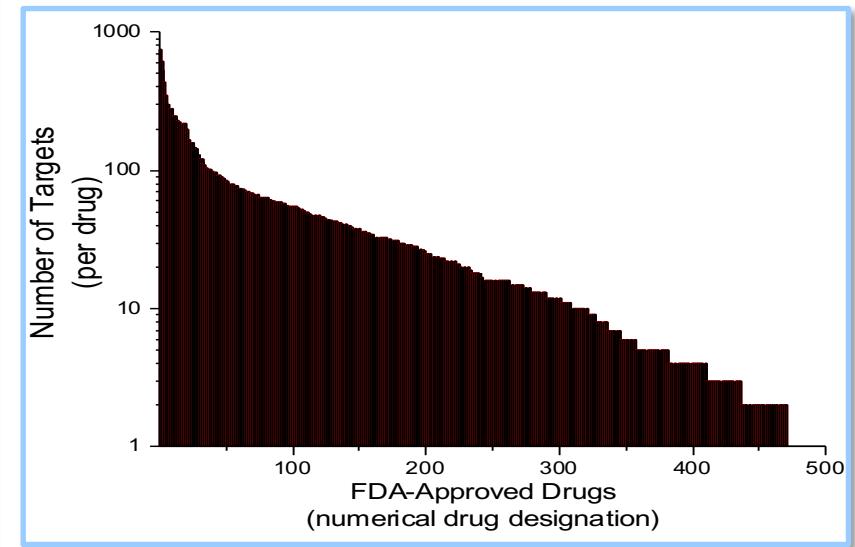
- 7-15 years and > 1 Billion USD per drug
- >90% failure at late stage. 42% of approved drugs (25/59) in 2018 were based on a single pivotal trial.
- Post-market failure due to adverse drug effect (e.g. Vioxx, Lipobay).



Drug on and off target binding along its metabolic pathway

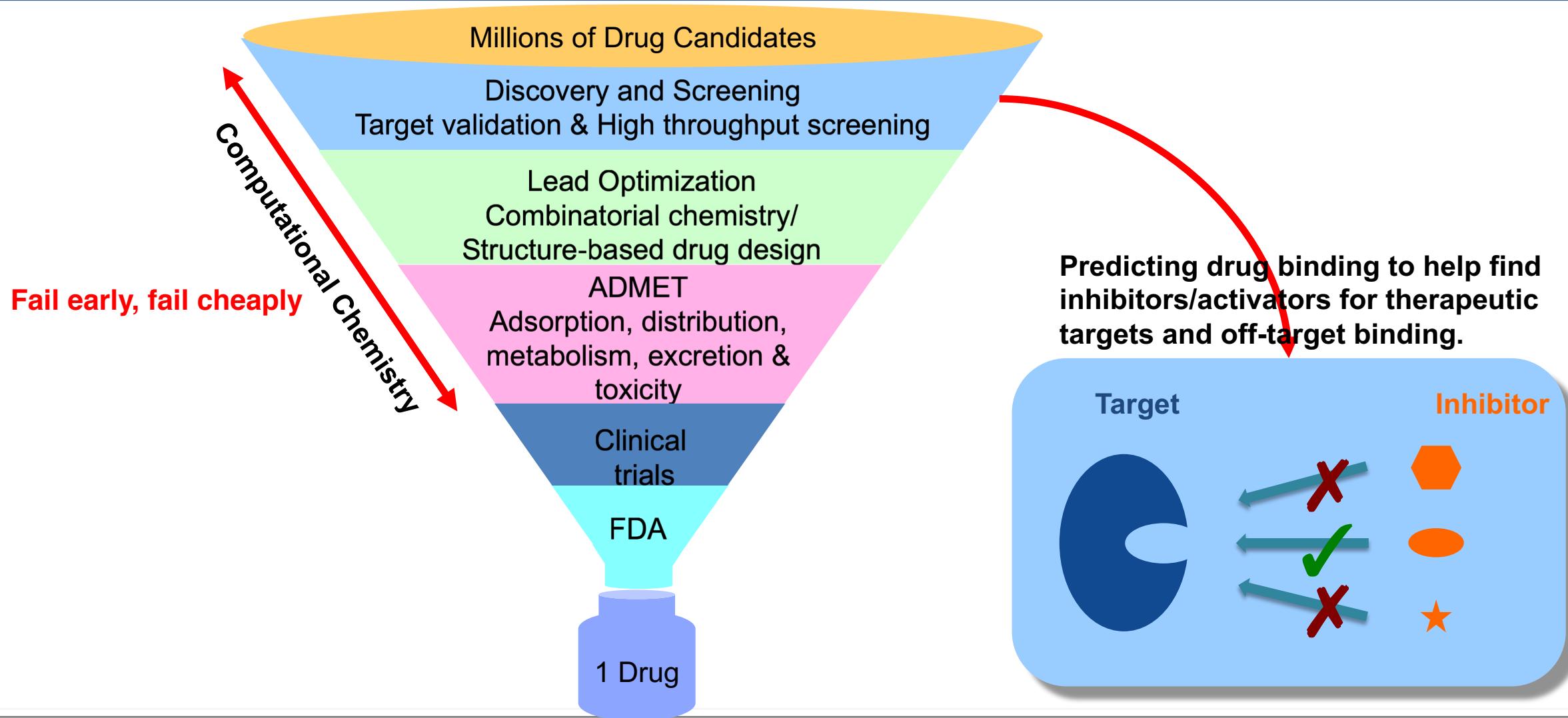


Drug could hit unexpected off-targets and cause severe side effects

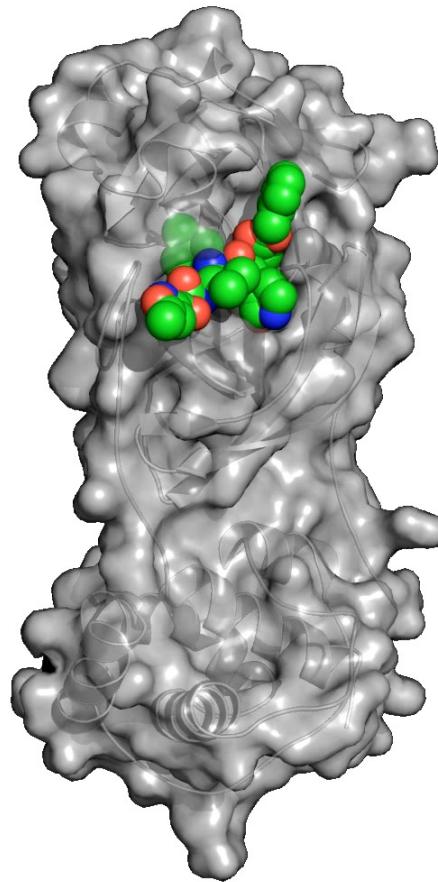


Average: 32.4 targets/drug

What can computational chemistry do for drug design?



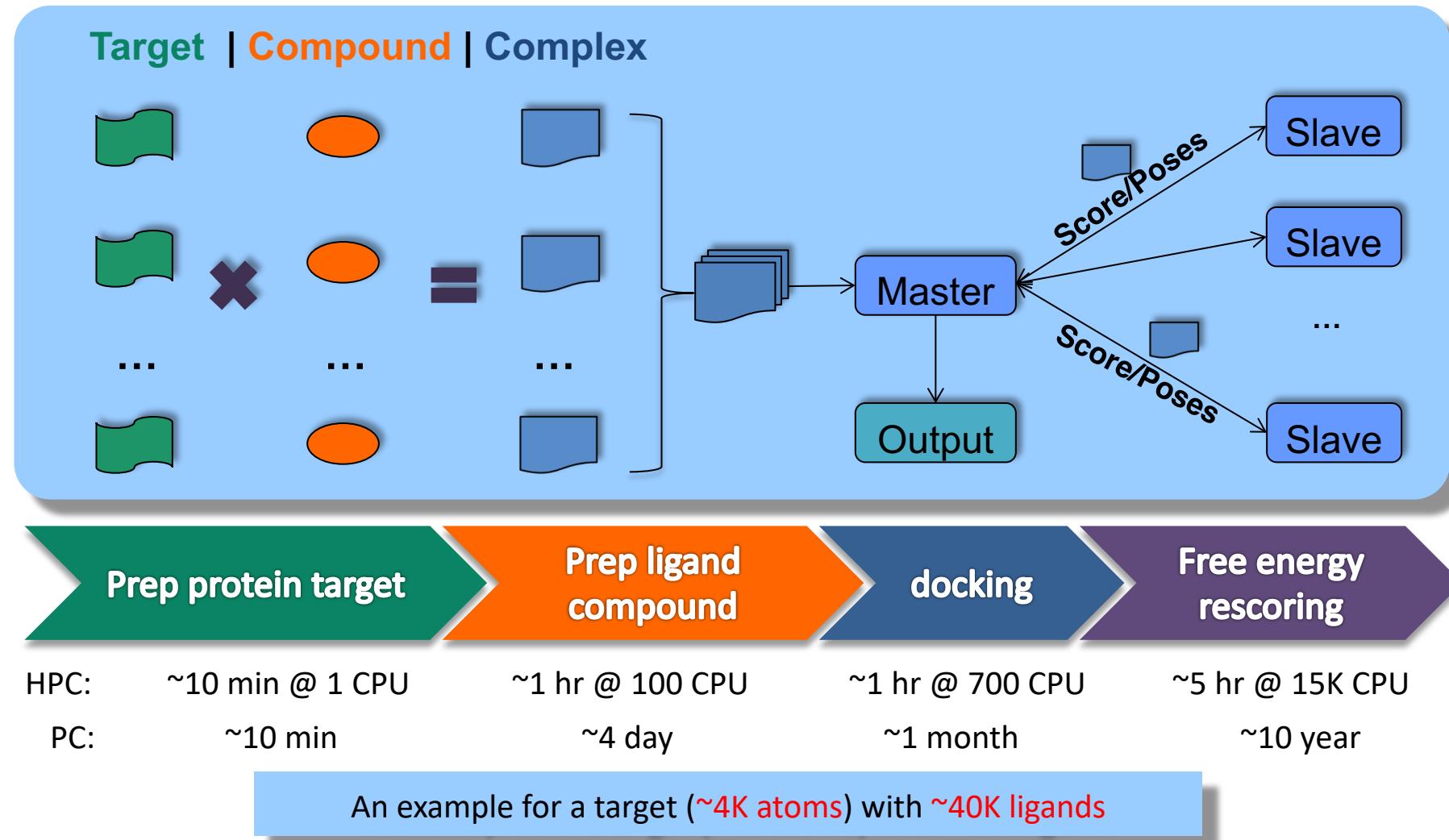
All drugs must bind their protein target



Predict the protein and small molecule binding affinity

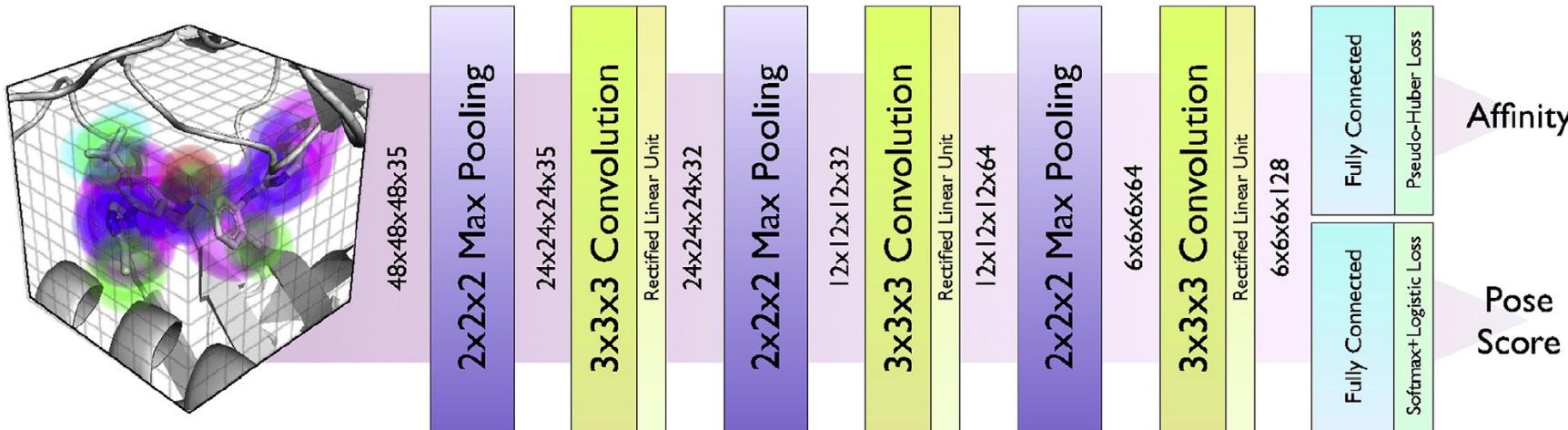
- Physical based prediction
 - Molecular Docking
 - Free energy (MM/GBSA) rescoring
- Machine learning based prediction
 - Fusion model
- Docking is required to obtain the complex structure, which will be used in MM/GBSA rescoring and fusion calculation

ConveyorLC pipeline for physical-based binding prediction



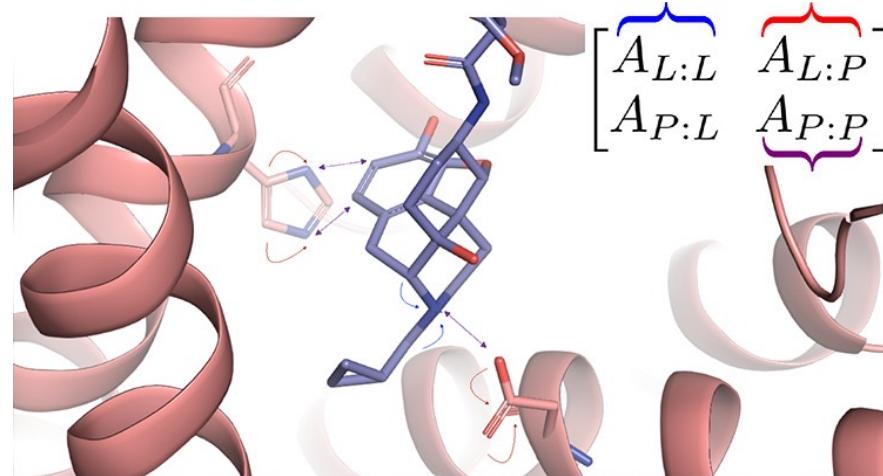
Binding affinity prediction by 3D convolutional neural network

- 3D-CNN discretizes 3D space into 3D-voxels, apply convolution operations to progressively down sample feature maps
- Select features including element type, hybridization, number of heavy atom bonds, partial charges, van der Waals, hydrophobic, ring, donor/acceptor, etc.
- Training on the PDBBind general set and use the Ki and KD values only (13K crystal structures)

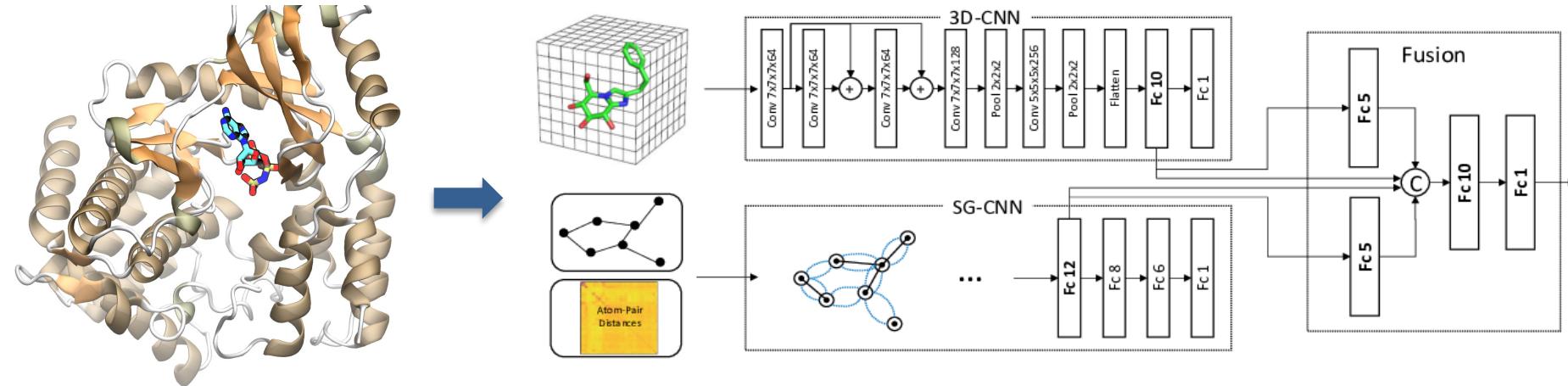


Binding affinity prediction by spatial graphs convolutional neural network

- SG-CNN computes distance matrices for atoms in ligand (L:L), protein (P:P), and complex (P:L,:L:P)
- Form graph where nodes are atoms and edges are the covalent/nonbonded interactions, label edges with distances



Combine 3D-CNN and SG-CNN for better prediction



Performance on PDDBind Core Set — Docking Poses (**Not Xray Structures**)

Model	Pearson r	Spearman r	MAE	RMSE
3D-CNN	0.523	0.503	1.843	2.358
SG-CNN	0.656	0.635	1.343	1.649
Vina	0.616	0.618	-	-
MM/GBSA	0.629	0.641	-	-
Mid Fusion	0.677	0.647	1.351	1.715
Fusion	0.685	0.668	1.340	1.701

Use the **second-** and **fourth-last** layer's output activations from 3D-CNN and SG-CNN, respectively

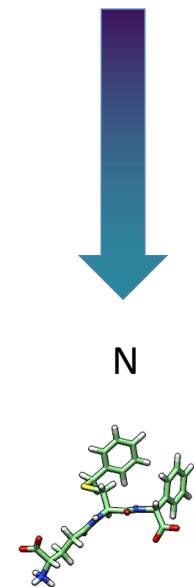
We are creating a queryable protein-small molecule atlas

A composite science workflow (AHA MoleS) is developed for the docking and fusion calculations

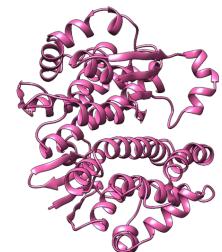
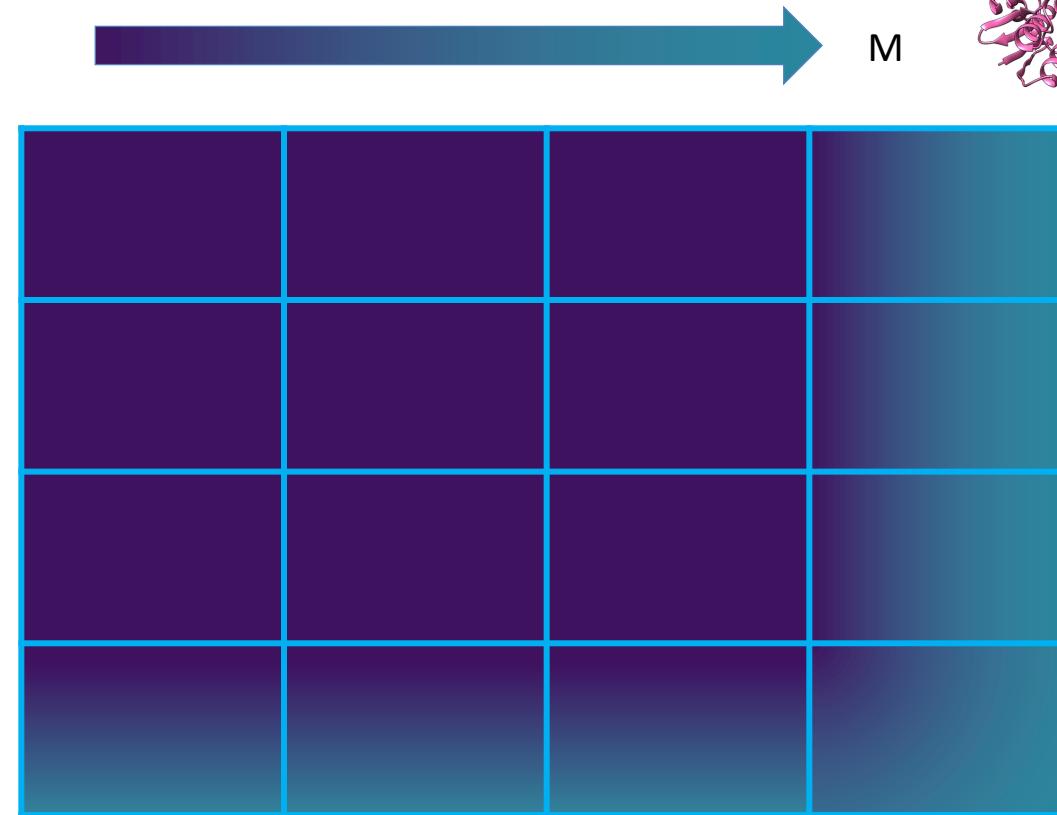
Small Molecule Compounds: ~2.1M

ChEMBL28: ~1.9M
Broad Drugs: 6550
Foodome: 24,114
NCI60: 244,800
PDBBind: 14,744

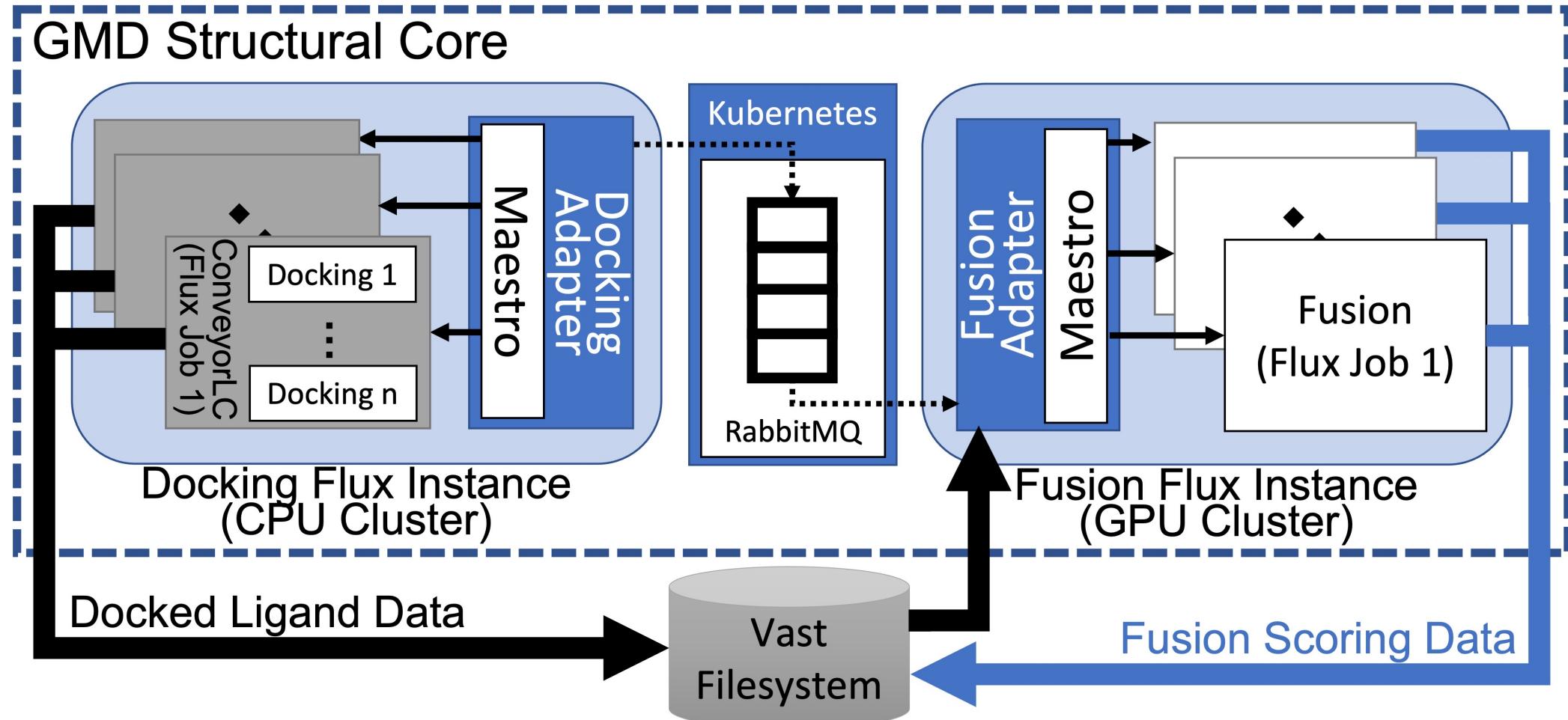
Total docking calculations: ~30 B
Total fusion calculations: ~300 B



Human Proteins: 28,658 models for 14,813 proteins
Experimental structures (>35%): 8,283
Experimental structures (>10%): 12,894



AHA MoleS workflow architecture



Demo of CPU only workflow

Balance of docking and fusion jobs

Beginning of workflow

```
(cpu-workflow-r1) ahashare@ruby964$ fluxj
```

JOBID	USER	NAME	ST	URG	NTASKS	NNODES	R
fAx1yUX9	ahashare	main	R	16	1	1	9
fAJKDjrX	ahashare	batch_completion	R	16	1	1	1
f9XghMpT	ahashare	docking_adapter	R	16	1	1	1
f8eQiCWj	ahashare	fusion_adapter	R	16	1	1	1

- Running on 24 Ruby nodes.
- Compound batch size: ~1500
- 1 docking adapter
- Each docking job uses 3 nodes running 21 MPI tasks
- 36 pre-fetch docking jobs
- 1 fusion adapter

Launching of docking jobs

```
(cpu-workflow-r1) ahashare@ruby964$ fluxj
```

JOBID	USER	NAME	ST	URG	NTASKS	NNODES	R
fGVtBFjd	ahashare	docking	PD	16	21	-	- -
fGdomq6X	ahashare	docking	PD	16	21	-	- -
fGkMjfEo	ahashare	docking	PD	16	21	-	- -
fGtrzY5d	ahashare	docking	PD	16	21	-	- -
fH47g46T	ahashare	docking	PD	16	21	-	- -
fH6ppjy1	ahashare	docking	PD	16	21	-	- -
fHHjsTJF	ahashare	docking	PD	16	21	-	- -
fHyDLqM	ahashare	docking	PD	16	21	-	- -
fHJ2BK03	ahashare	docking	PD	16	21	-	- -
fHX2o1qZ	ahashare	docking	PD	16	21	-	- -
fHdkVjyAP	ahashare	docking	PD	16	21	-	- -
fHdgFo7D	ahashare	docking	PD	16	21	-	- -
fHsJTCK1	ahashare	docking	PD	16	21	-	- -
fHWpfDQB	ahashare	docking	PD	16	21	-	- -
fJ1iD8ZT	ahashare	docking	PD	16	21	-	- -
fJSBqSMM	ahashare	docking	PD	16	21	-	- -
fJBaULoZ	ahashare	docking	PD	16	21	-	- -
fJBbxL5u	ahashare	docking	PD	16	21	-	- -
fJJtUoQ8	ahashare	docking	PD	16	21	-	- -
fJQ8QGAs	ahashare	docking	PD	16	21	-	- -
fJSKuCQX	ahashare	docking	PD	16	21	-	- -
fJSdh3nf	ahashare	docking	PD	16	21	-	- -
fJVHsm6X	ahashare	docking	PD	16	21	-	- -
fJWYX9uR	ahashare	docking	PD	16	21	-	- -
fJcbtDh9	ahashare	docking	PD	16	21	-	- -
fJdkPsao	ahashare	docking	PD	16	21	-	- -
fjesq7mq	ahashare	docking	PD	16	21	-	- -
fJicHjvB	ahashare	docking	PD	16	21	-	- -
fJjMGx6B	ahashare	docking	PD	16	21	-	- -
fJmehqTD	ahashare	docking	PD	16	21	-	- -
fGRnVfN7	ahashare	docking	R	16	3	7.746s	ruby[:]
fGGZn7uy	ahashare	docking	R	16	3	8.106s	ruby[:]
fGSuZHQC	ahashare	docking	R	16	3	8.500s	ruby[:]
fGSuZHQB	ahashare	docking	R	16	3	8.502s	ruby[:]
fG5t5J7r	ahashare	docking	R	16	3	8.504s	ruby[:]
fG5t5J7q	ahashare	docking	R	16	3	8.506s	ruby[:]
fAx1yUX9	ahashare	main	R	16	1	20.19s	ruby[:]
fAJKDjrX	ahashare	batch_completion	R	16	1	21.66s	ruby[:]
f9XghMpT	ahashare	docking_adapter	R	16	1	23.41s	ruby[:]
f8eQiCWj	ahashare	fusion_adapter	R	16	1	25.41s	ruby[:]

```
(cpu-workflow-r1) ahashare@ruby964$ fluxj
```

JOBJID	USER	NAME	ST	URG	NTASKS	NNODES	RUNTIME	NODELIST
fJdKPsao	ahashare	docking	PD	16	21	-	- -	
fjesq7mq	ahashare	docking	PD	16	21	-	- -	
fJicHjvB	ahashare	docking	PD	16	21	-	- -	
fJjMGx6B	ahashare	docking	PD	16	21	-	- -	
fJmehqTD	ahashare	docking	PD	16	21	-	- -	
f2tLcX2ym	ahashare	docking	PD	16	21	-	- -	
f3vQADVxu	ahashare	docking	PD	16	21	-	- -	
f6J5TKAA3	ahashare	docking	PD	16	21	-	- -	
f7GSxBJrb	ahashare	docking	PD	16	21	-	- -	
fAUNh0K2o	ahashare	docking	PD	16	21	-	- -	
fBzXmkNoh	ahashare	docking	PD	16	21	-	- -	
fC78gGtHD	ahashare	docking	PD	16	21	-	- -	
fGWzoRrk9	ahashare	docking	PD	16	21	-	- -	
fLdF1luryD	ahashare	docking	PD	16	21	-	- -	
fP9zhLjlS	ahashare	docking	PD	16	21	-	- -	
fQ6ivUqFm	ahashare	docking	PD	16	21	-	- -	
fQXKx1Xzo	ahashare	docking	PD	16	21	-	- -	
fQgSXZUVV	ahashare	docking	PD	16	21	-	- -	
fRBDD57SX	ahashare	docking	PD	16	21	-	- -	
fUBULEtto	ahashare	docking	PD	16	21	-	- -	
fVv7LqGiu	ahashare	docking	PD	16	21	-	- -	
fWZKbXjRR	ahashare	docking	PD	16	21	-	- -	
fYJQ6dkWk	ahashare	docking	PD	16	21	-	- -	
fbqmuXHQ7	ahashare	docking	PD	16	21	-	- -	
ffhTcCQRm	ahashare	docking	PD	16	21	-	- -	
fg7VqczzP	ahashare	docking	PD	16	21	-	- -	
fgFQQUrEb	ahashare	docking	PD	16	21	-	- -	
fnqL87zb9	ahashare	docking	PD	16	21	-	- -	
finXkq26B	ahashare	docking	PD	16	21	-	- -	
fmYGYcrs1	ahashare	docking	PD	16	21	-	- -	
fmYP8hfsd	ahashare	fusion	R	31	1	1	4.794m	ruby1453
fJcbtDh9	ahashare	docking	R	16	21	3	5.112m	ruby[1239,1269,1286]
fJWYX9uR	ahashare	docking	R	16	21	3	11.14m	ruby[1182,1190,1213]
fhq6rgRdq	ahashare	fusion	R	31	1	1	12.94m	ruby1481
fJVHsm6X	ahashare	docking	R	16	21	3	13.28m	ruby[227,407,491]
fgFaJE4bh	ahashare	fusion	R	31	1	1	16.4m	ruby1455
fJSdh3nf	ahashare	docking	R	16	21	3	16.75m	ruby[944,999,1048]
fJJtUoQ8	ahashare	docking	R	16	21	3	26.36m	ruby[1300,1352,1383]
fYJNdiTK	ahashare	fusion	R	31	1	1	33.84m	ruby1391
fJBaULoZ	ahashare	docking	R	16	21	3	37.97m	ruby[520,757,819]
fUBdmJLBh	ahashare	fusion	R	31	1	1	42.87m	ruby1453
fQgB1ezjq	ahashare	fusion	R	31	1	1	50.56m	ruby1481
fAx1yUX9	ahashare	main	R	16	1	1	1.702h	ruby1391
fAJKDjrX	ahashare	batch_completion	R	16	1	1	1.702h	ruby1453
f9XghMpT	ahashare	docking_adapter	R	16	1	1	1.703h	ruby1455
f8eQiCWj	ahashare	fusion_adapter	R	16	1	1	1.703h	ruby1481

Scaling on hybrid CPU and GPU clusters

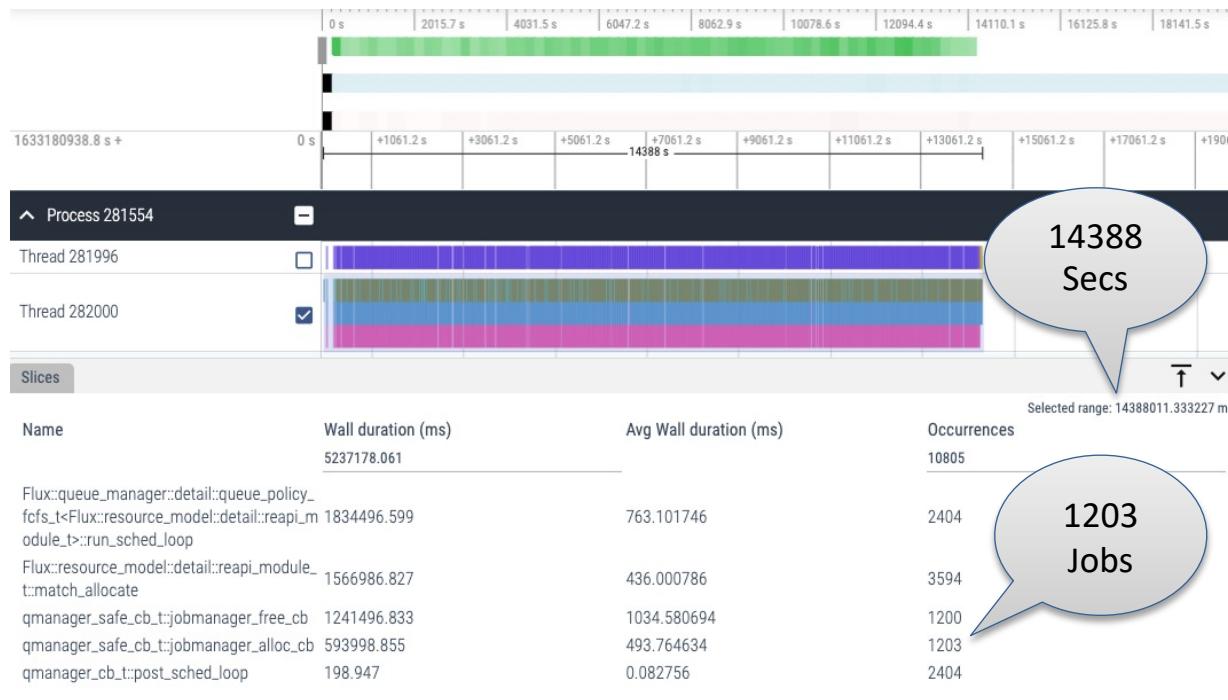
- Workflow consists of just Docking and Fusion currently scaling to ~1000- Ruby nodes and 140 or 220 Lassen nodes respectively (uses Flux)
- **Compound batch size:** ~1500
- 2 Docking adapters process each batch via maestro/flux/MPI sub-workflow:
 - **per batch:** Use 145 nodes or 5 nodes.
 - Pre-fetch set to 36 batches for each adapter (2 adapters * 36 prefetch * 11 nodes = 792 nodes)
- Each Fusion worker processes each batch on lassen with:
 - **per batch:** 1 task using 1 GPU and 10 cores-per-task (~30 minutes per batch).
 - Each node runs 4 fusion works (560 or 880 workers)
- **hybrid cluster workflow** requires DAT to acquire resources simultaneously on both clusters

Performance of different configurations using PerfFlowApect

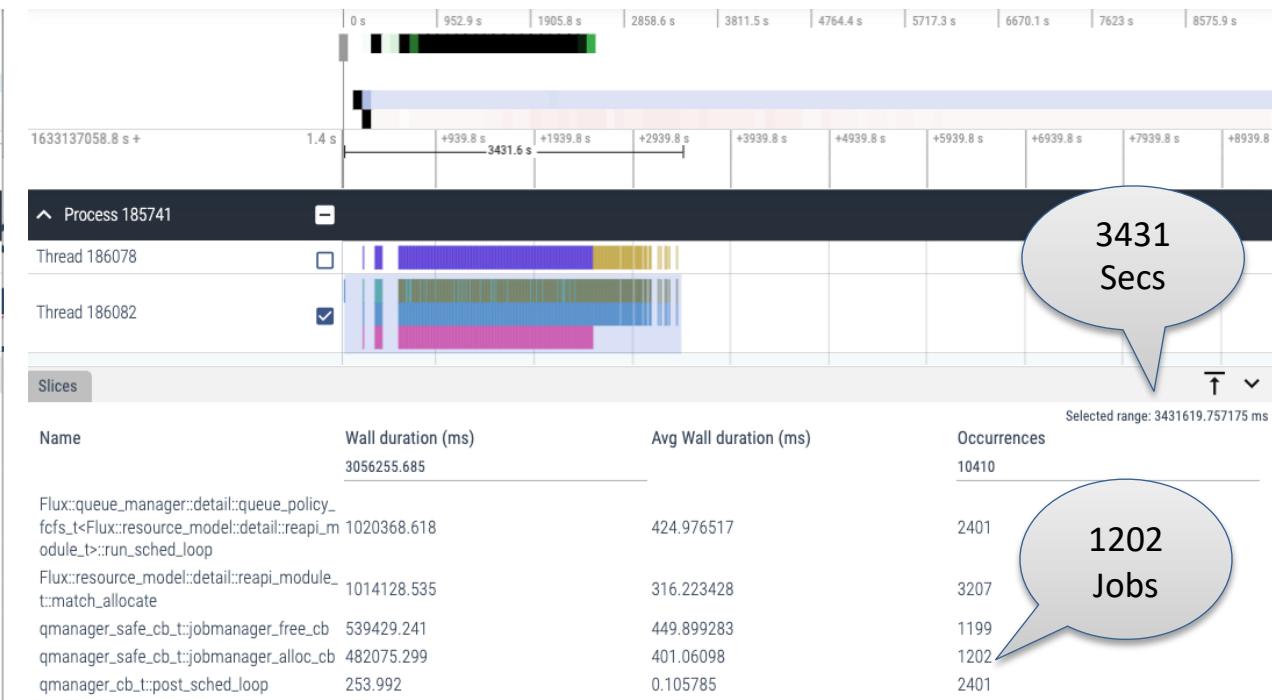
Flux Version	Docking Nodes/Job	Fusion Total Nodes	Makespan (second)	Time Per Mol. (ms)
0.26.0	145	220	15208	7.67
0.26.0	5	220	7841	3.95
0.26.0	145	140	15257	7.69
0.26.0	5	140	9717	4.9
0.29.0	145	220	19049	9.6
0.29.0	5	220	8041	4.05
0.29.0	145	140	19236	9.7
0.29.0	5	140	9666	4.87

Note: Use the PerfFlowApect data at the top-level Flux instance on Ruby

Performance impacted by docking job size



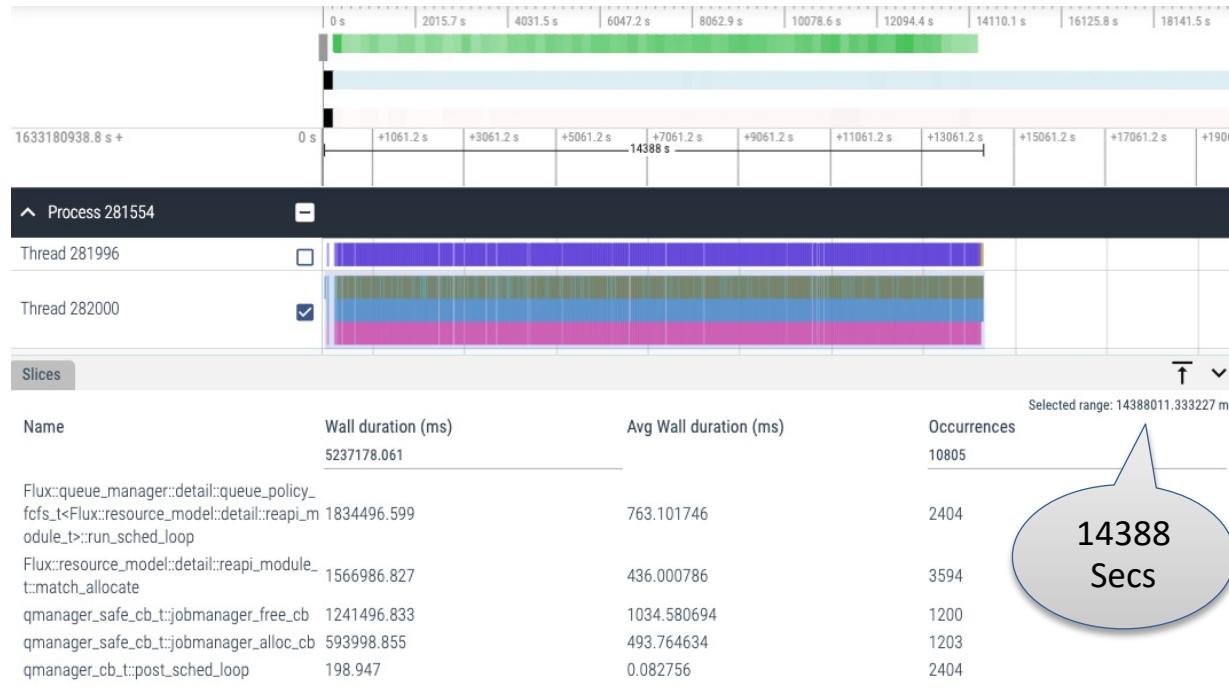
Flux=0.26.0 Dock=145 nodes/Job Fusion=220



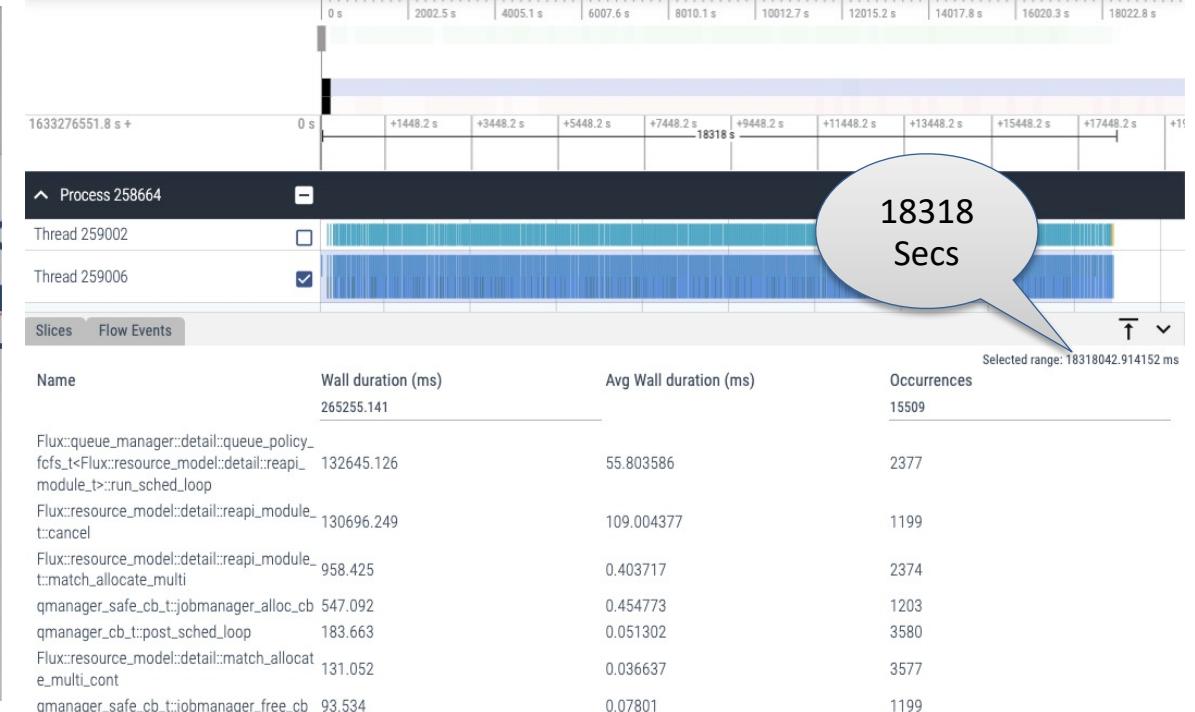
Flux=0.26.0 Dock=5 nodes/Job Fusion=220

Using larger number of nodes per docking job slow down the calculation by **1.93x**

Performance affected by flux overhead



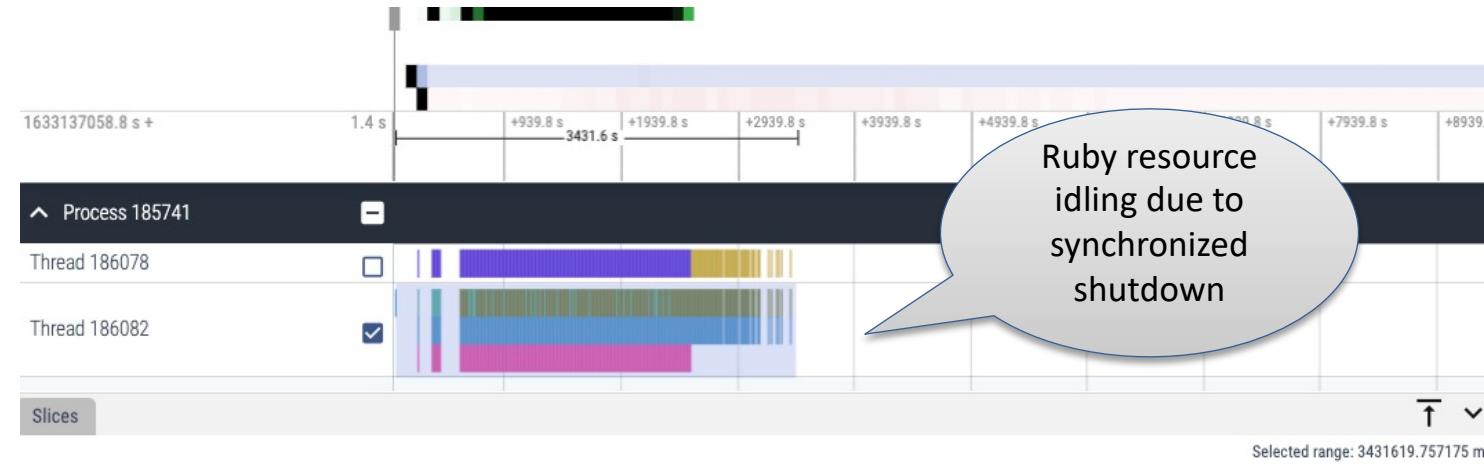
Flux=0.26.0 Dock=145 nodes/job Fusion=220 nodes



Flux=0.29.0 Dock=145 nodes/job Fusion=220 nodes

The higher scheduling overhead with Flux Version 0.29.0 slows down the calculations by **1.27x**.

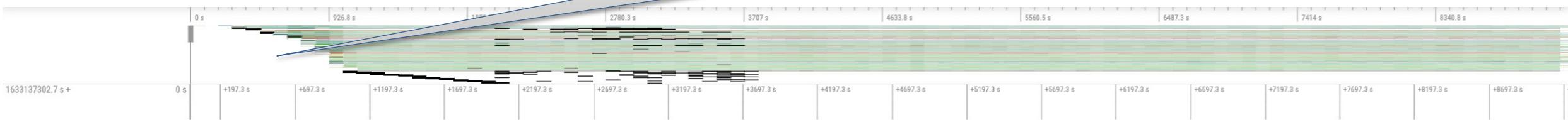
How PerfFlow analysis can help understand other future choices of performance variables



Flux=0.26.0 Dock=5 nodes/job Fusion=220 nodes

Use the PerfFlowApect data at the top-level Flux instance on Ruby

Lassen resource idling due to synchronized startup



Use the PerfFlowApect data at the top-level Flux instance on Lassen

Conclusions

- AHA MoleS workflow supports scaling from a small number of nodes to 1000+ nodes, potentially spanning multiple heterogeneous clusters depending on the application and workload resource requirements
- Computer resource is idling due to synchronization of docking and fusion
- For Ruby resources, the more imbalanced Docking and fusion throughputs are the worse
- For Lassen resources, the larger the batch size is the more lead time idling. It cannot be determined by flux-level perfflow, Fusion adapter and maestro level perfflows are required.
- In the future, workflow needs to loose/remove synchronization between docking and fusion for the hybrid computer resource.

Acknowledgements

- LLNL PLS/BBTD
 - Felice Lightstone
 - Drew Bennett
 - Brian Bennion
 - Dan Kirshner
 - Edmond Lau
 - Ali Navid
 - Sergio Wong
 - Fangqiang Zhu
- LLNL GS
 - Jonathan Allen
 - Aidan Epstein
 - Stewart He
 - Derek Jones
 - Jeff Mast
 - Garrett Stevenson
 - Marisa Torres
 - Adam Zemla
- LLNL Computing
 - Dong Ahn
 - Bronis de Supinski
 - Francesco Di Natale
 - Stephen Herbein
 - Sam Ade Jacobs
 - Ian Karlin
 - Hyojin Kim
 - Daniel Milroy
 - Brian Van Essen
- Funding
 - American Heart Association CRADA



Thank you!
Questions?

[https://github.com/flux-
framework/Tutorials.git](https://github.com/flux-framework/Tutorials.git)