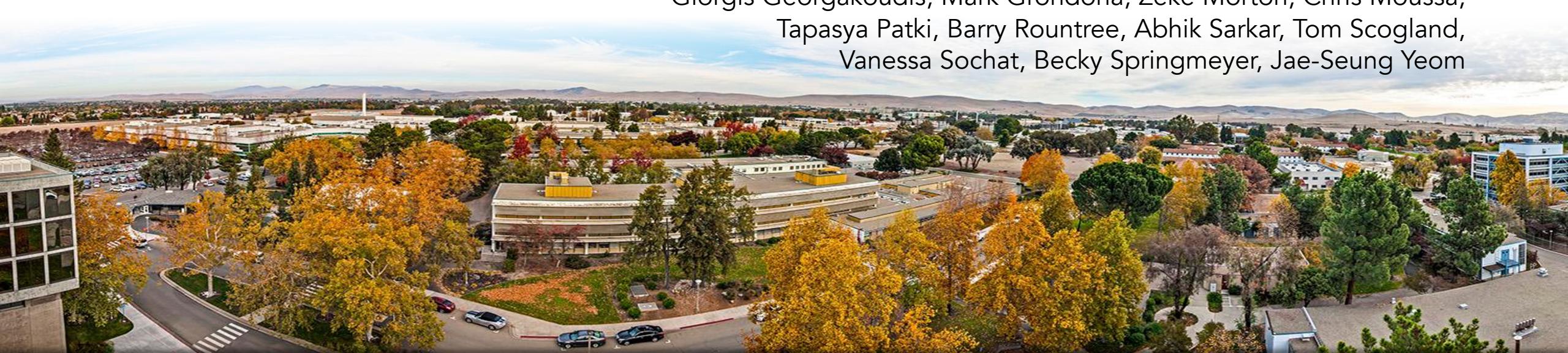


The Flux framework: overcoming scheduling and management challenges of exascale workflows

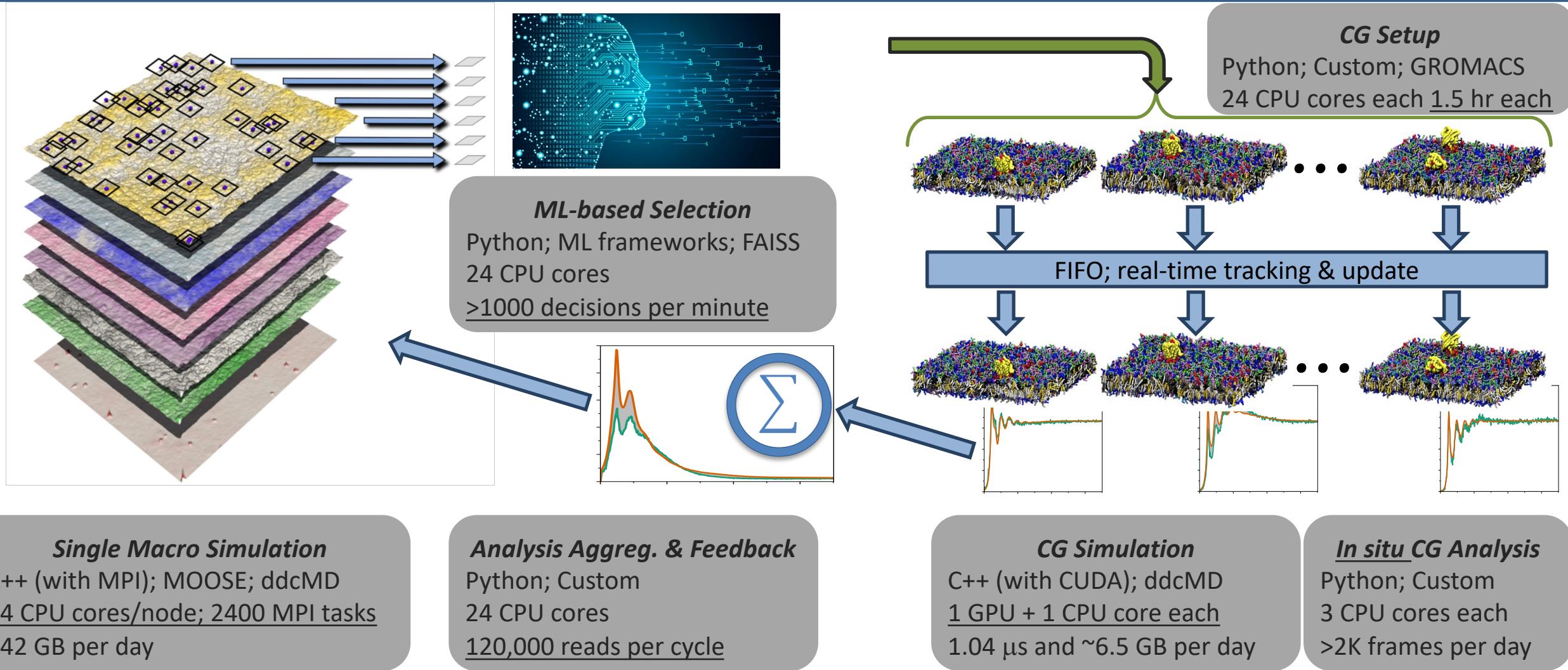
RADIUSS Tutorial on AWS
August 17, 2022

Dan Milroy

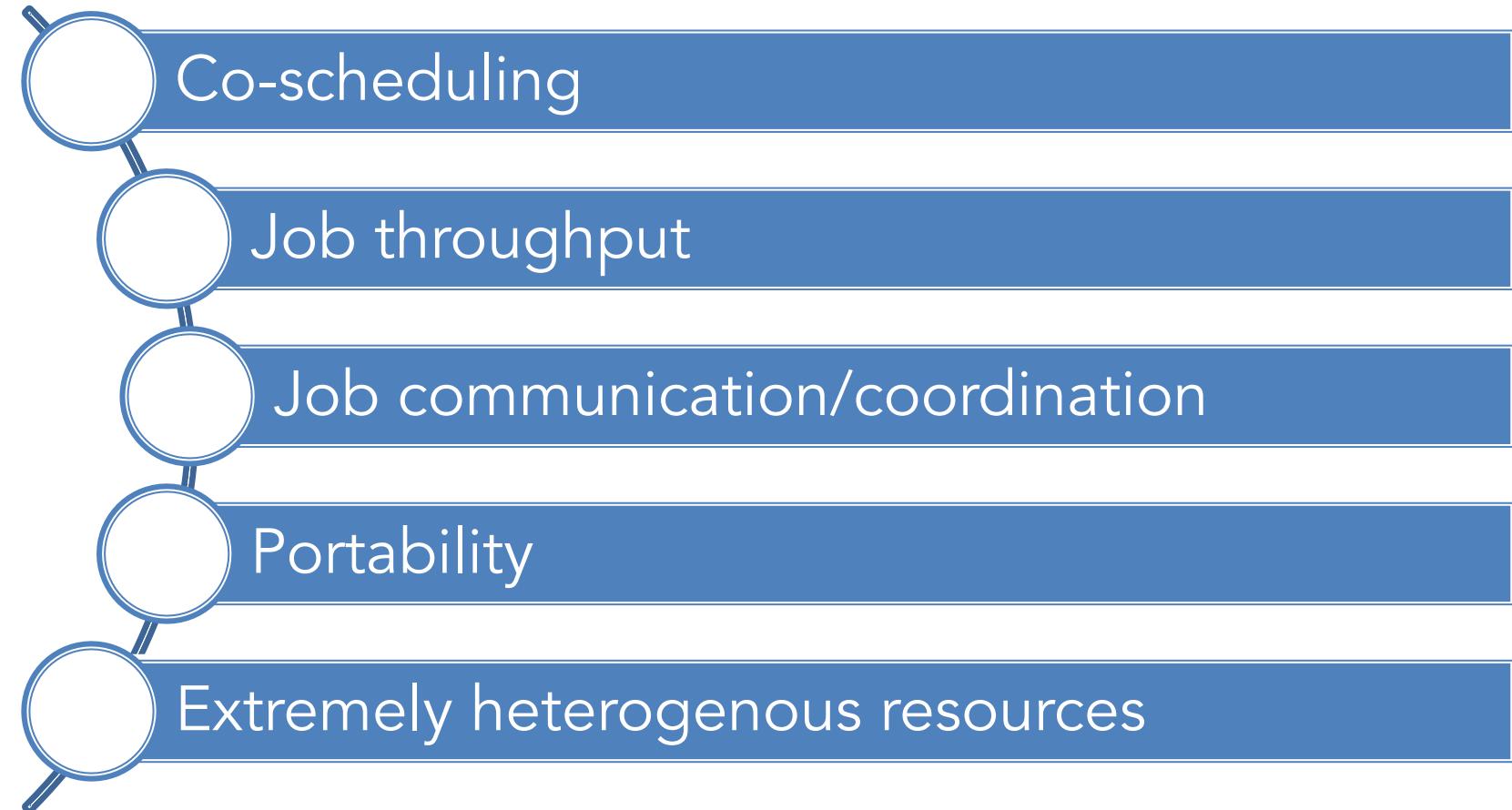
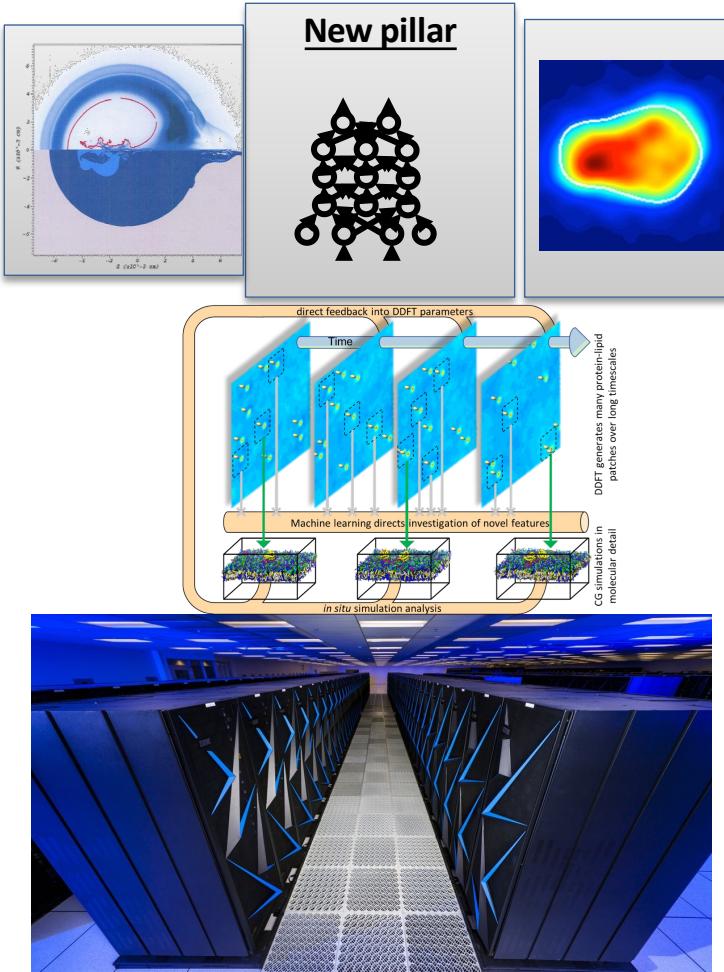
Al Chu, James Corbett, Ryan Day, Jim Garlick,
Giorgis Georgakoudis, Mark Grondona, Zeke Morton, Chris Moussa,
Tapasya Patki, Barry Rountree, Abhik Sarkar, Tom Scogland,
Vanessa Sochat, Becky Springmeyer, Jae-Seung Yeom



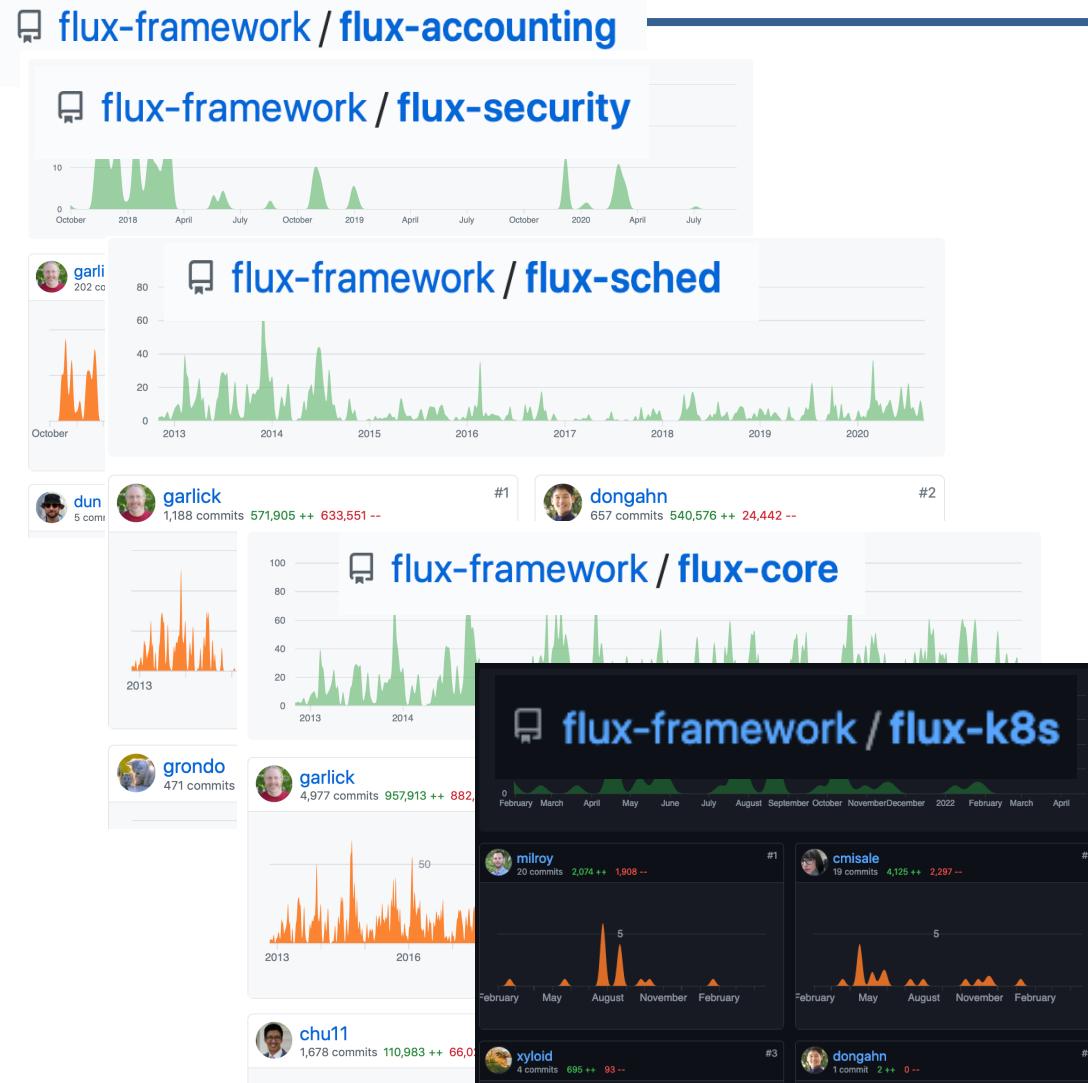
Sierra pre-exascale system is a wakeup call (MuMMI).



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



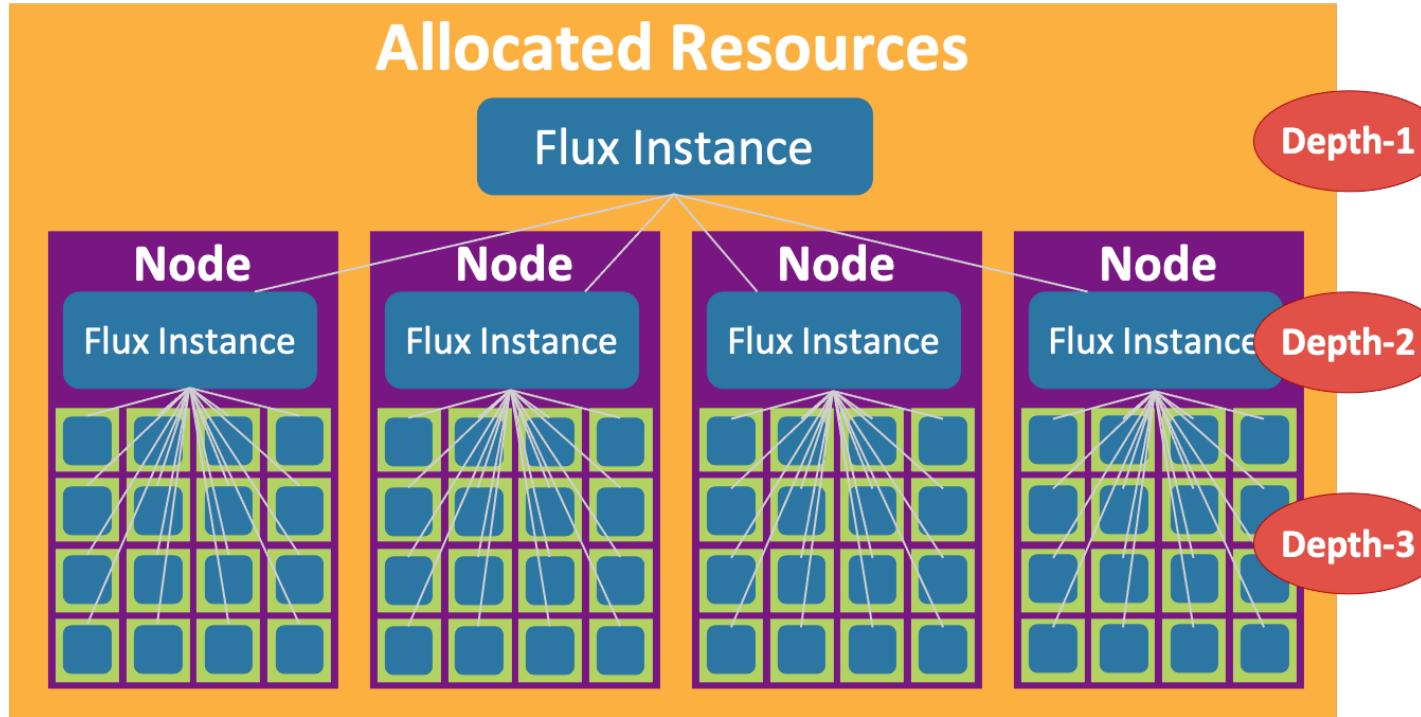
Flux solves key technical problems that emerge from these trends.



- 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.

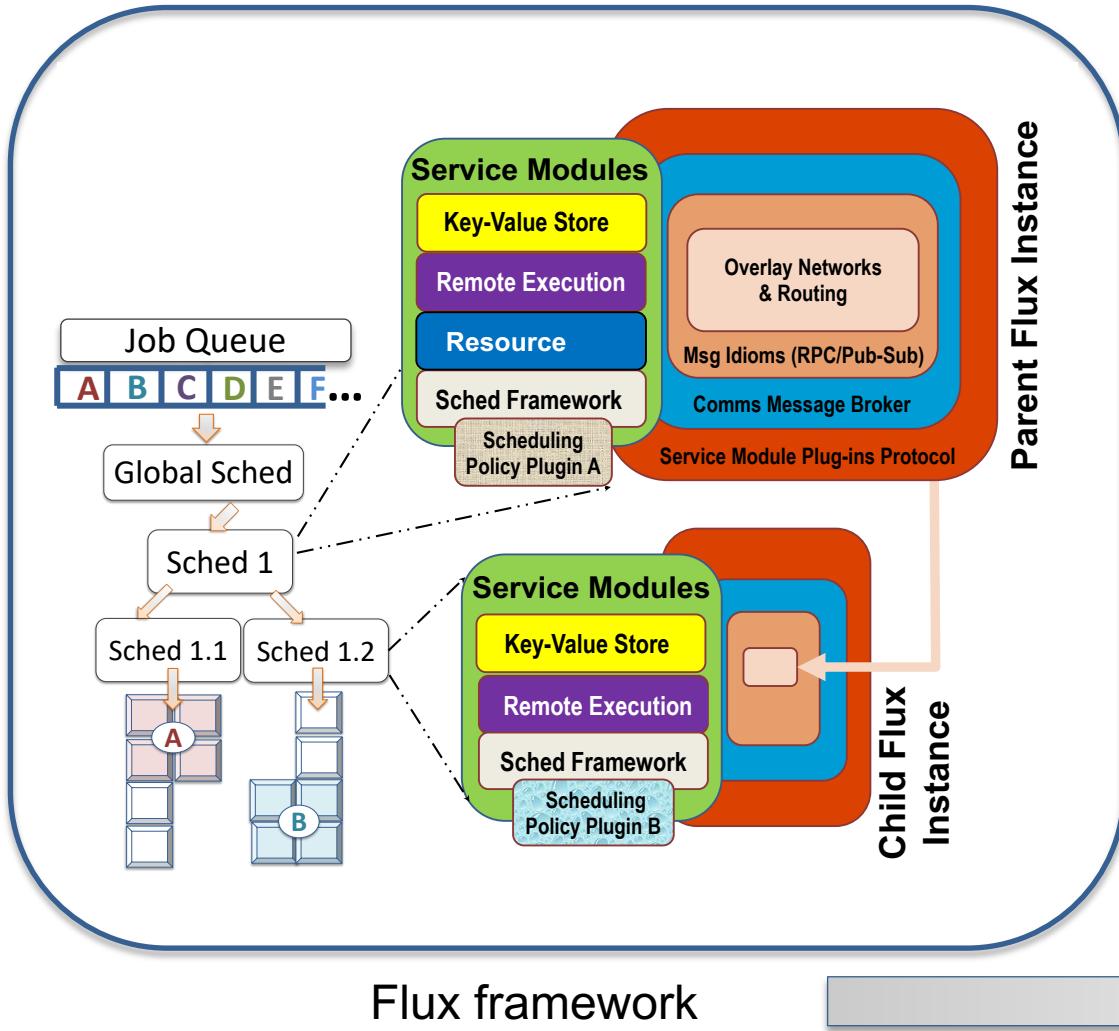


"Fractal scheduling" mitigates centralized scheduler bottleneck

- 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 is architected to embody our fully hierarchical scheduling model.



Techniques

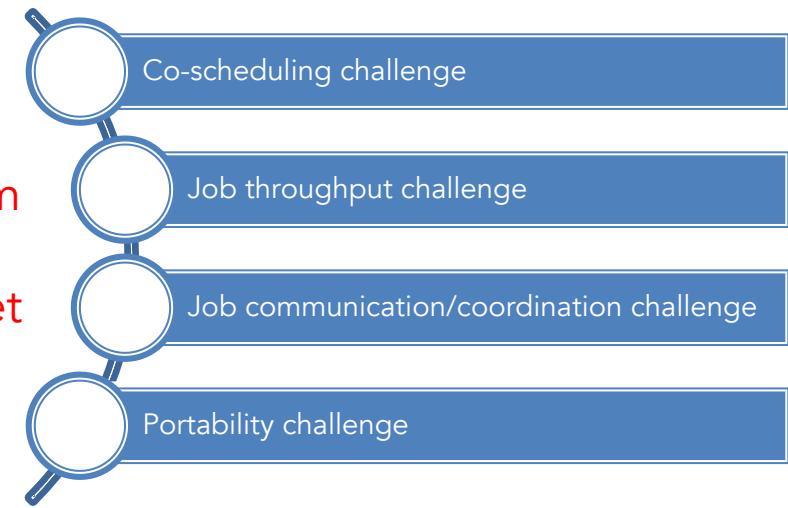
Scheduler Specialization

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

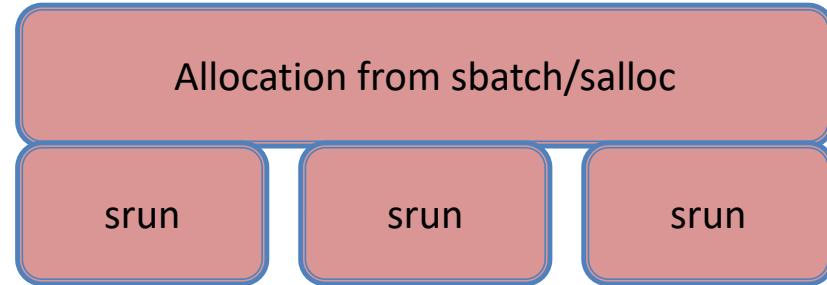
- [flux_attr_get\(3\)](#)
- [flux_aux_set\(3\)](#)
- [flux_child_watcher_create\(3\)](#)
- [flux_content_load\(3\)](#)
- [flux_core_version\(3\)](#)
- [flux_event_decode\(3\)](#)
- [flux_event_publish\(3\)](#)
- [flux_event_subscribe\(3\)](#)
- [flux_fatal_set\(3\)](#)
- [flux_fd_watcher_create\(3\)](#)
- [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\(3\)](#)



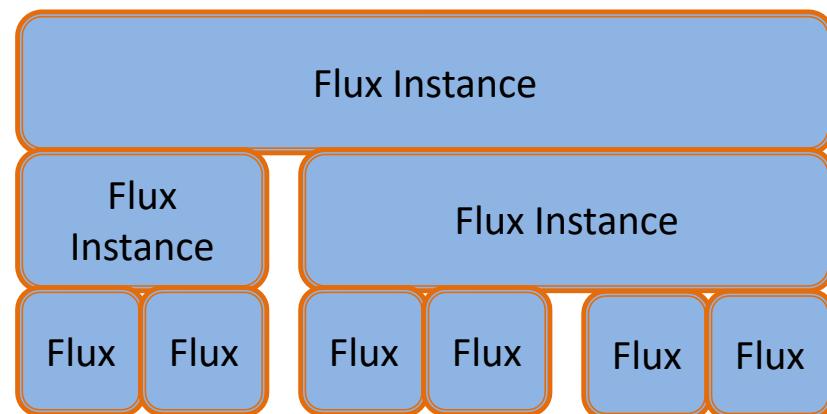
<https://flux-framework.readthedocs.io/projects/flux-core/en/latest/index.html>

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



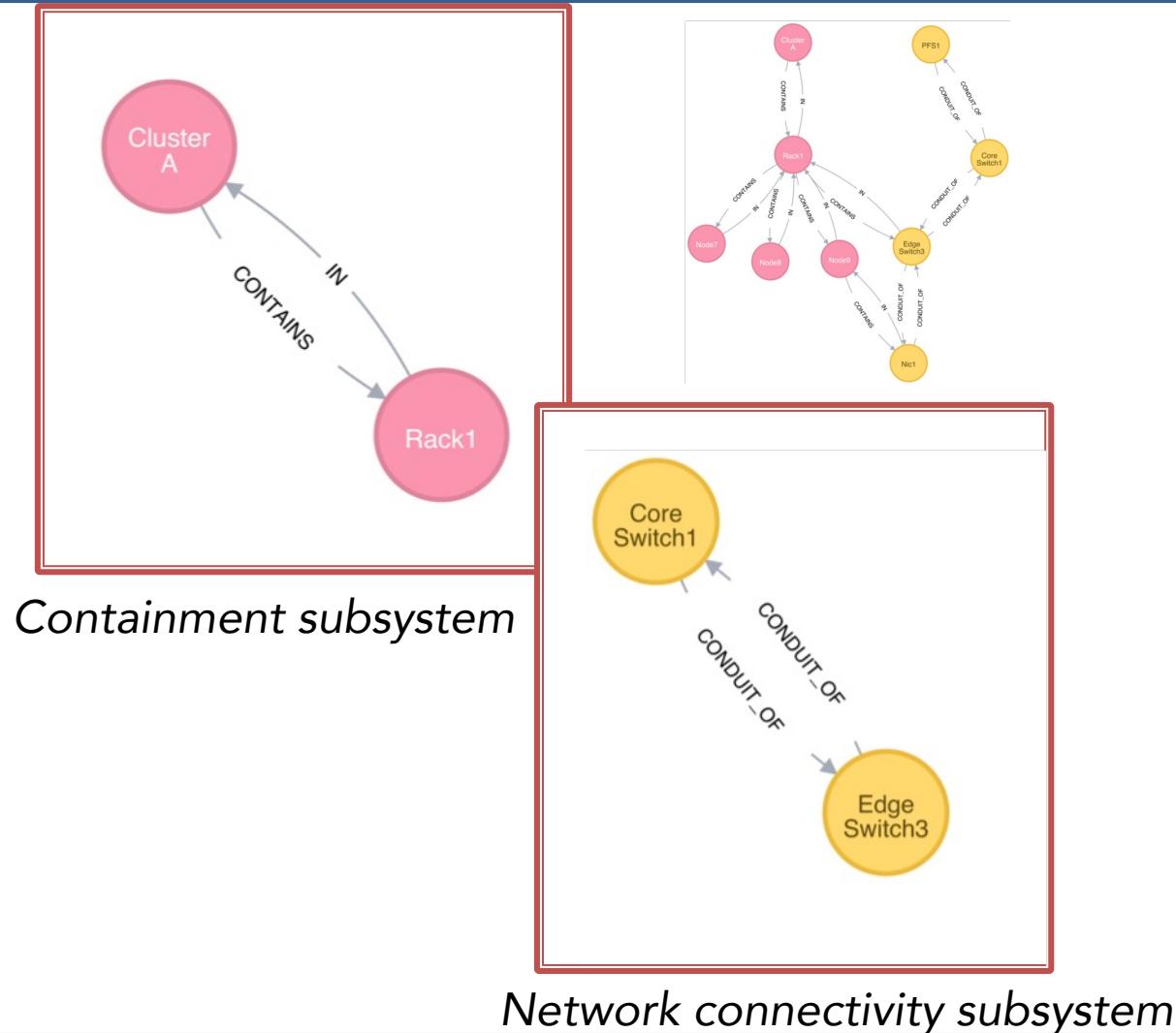
The Slurm hierarchy



One possible Flux hierarchy

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

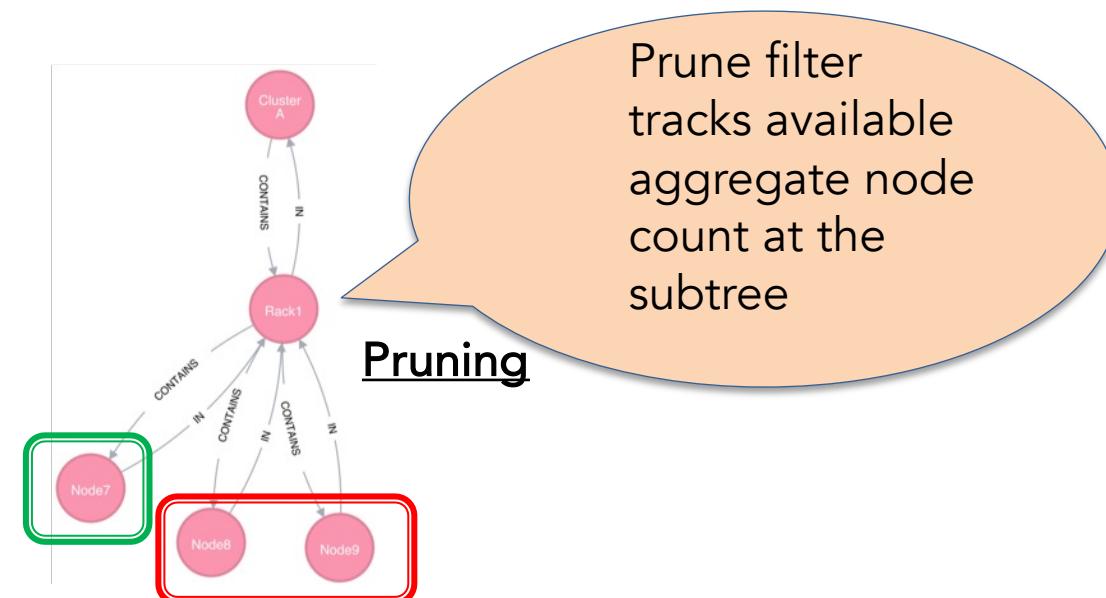
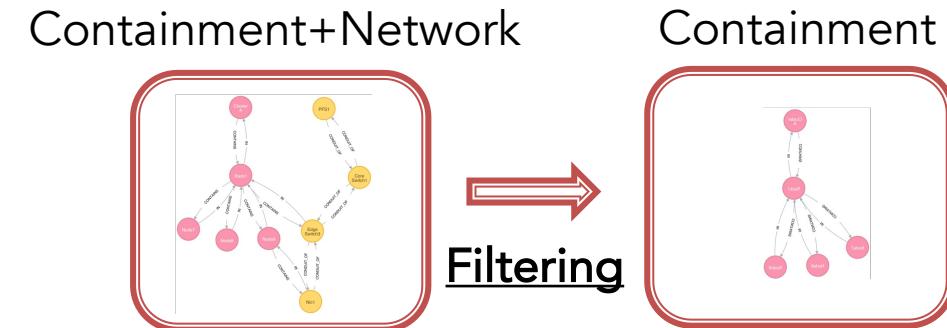
- Traditional resource data models are largely ineffective for resource heterogeneity
 - designed when systems were simpler
 - node-centric models
- Elevate resource relationships (edges) to an equal footing with resources (vertices)
- Complex scheduling can be expressed without changing the scheduler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



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

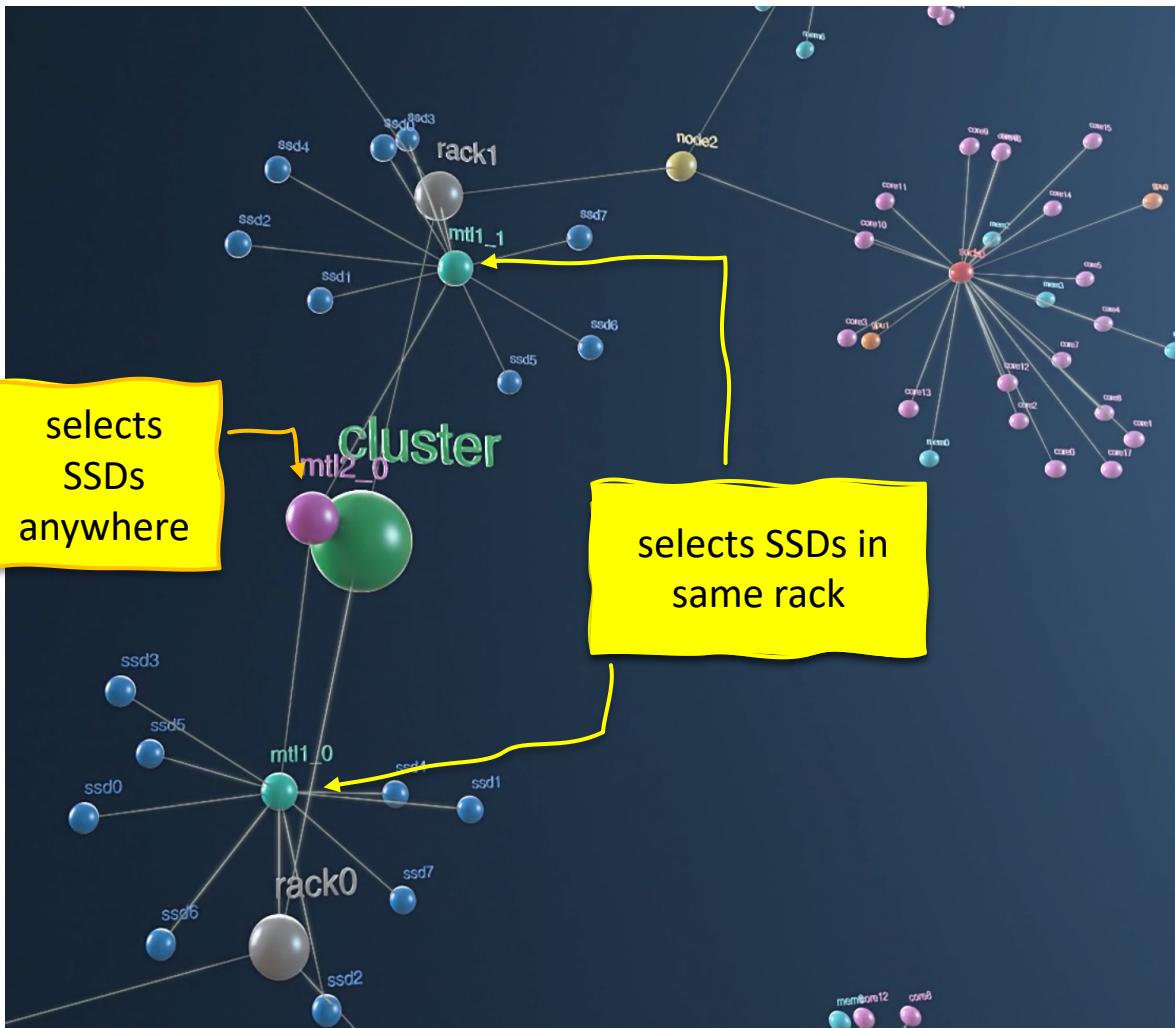


Flux's graph-oriented canonical jobspec allows for a highly expressive resource request specification.

- Graph-oriented resource requirements
 - Express the resource requirements of a program to the scheduler
 - Express program attributes such as arguments, runtime, and task layout to the execution service
- cluster → rack[2] → slot[3] → node[1] → socket[2] → core[18]
- **slot** is the only non-physical resource type
 - Represent a schedulable place where program processes will be spawned and contained
- Tasks section references slot and defines command

```
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
```

Fluxion graph approach solves the critical scheduling need for El Capitan's Rabbit multi-tiered storage.

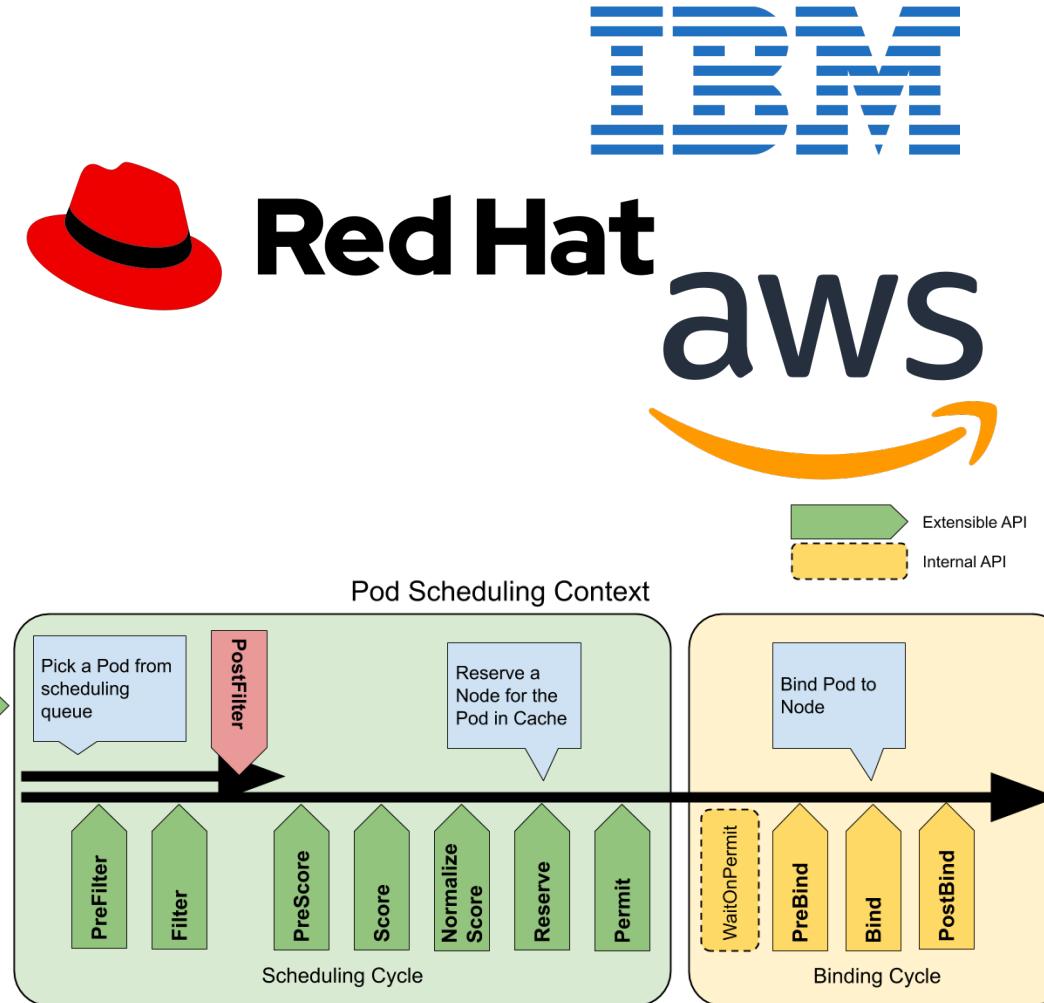


- Schedule SSDs in each rack
 - mount as node-local storage to compute in same rack
 - used to build ephemeral Lustre
 - deemed too difficult for traditional schedulers
- Easily enabled by Fluxion; no code change
- End-to-end job state transition
 - DataWarp/Rabbit mock service containers by HPE
 - Uses Kubernetes CRD

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
- Red Hat: E. Arango Gutierrez
- Newly added AWS: E. Bollig, M. Hugues, H. Poxon, L. Wofford
- Integrate Fluxion into Kubernetes via KubeFlux



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: <password>
 - Double-click **radiuss-aws-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 2022-RADIUSS-AWS 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.