



# Fast Forward Storage & I/O

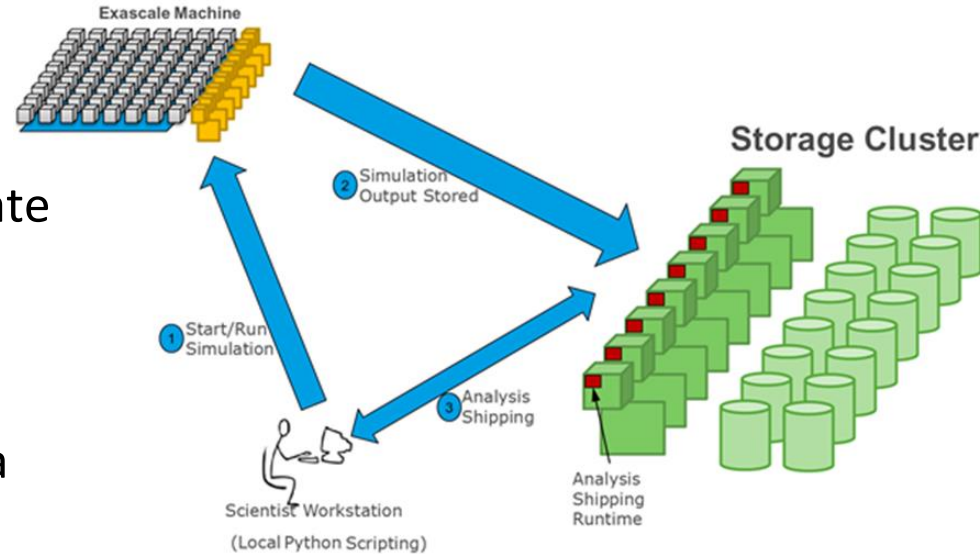
Jeff Layton (Eric Barton)

# DOE Fast Forward IO and Storage

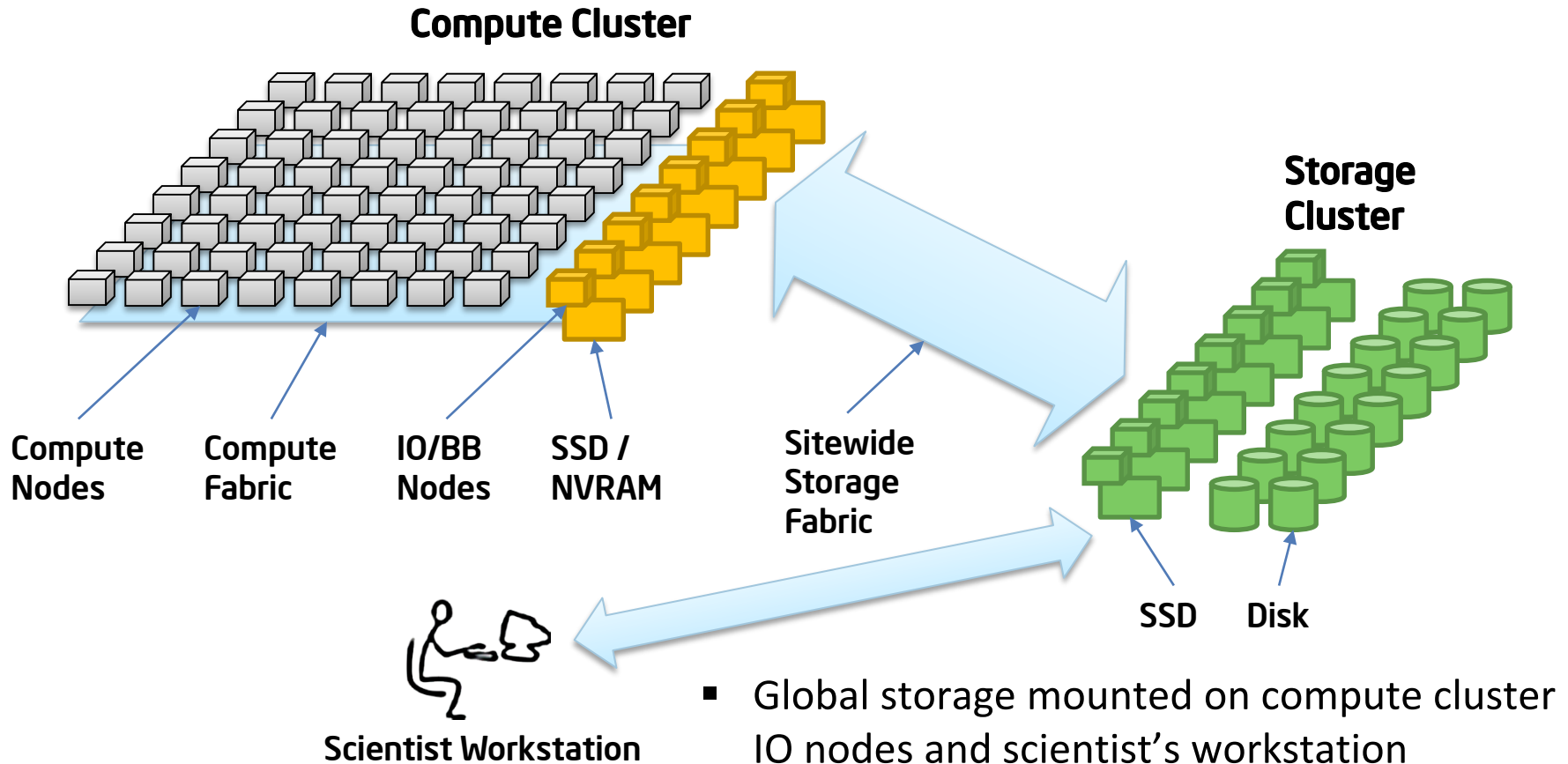
- Exascale R&D sponsored by 7 leading US national labs
  - Solutions to currently intractable problems of Exascale required to meet the 2020 goal of an Exascale system
- Whamcloud & partners won the IO & Storage contract
  - Proposal to rethink the whole HPC IO stack from top to bottom
    - Develop a working prototype
    - Demonstrate utility of prototype in HPC and Big Data
  - HDF Group – HDF5 modifications and extensions
  - EMC – Burst Buffer manager & I/O Dispatcher
  - Whamcloud – Distributed Application Object Storage (DAOS)
  - Cray – Scale out test
- Contract renegotiated on Intel acquisition of Whamcloud
  - Intel – Arbitrary Connected Graph computation
  - DDN – Versioning Object Storage Device

# Fast Forward Storage & IO Project Goals

- Make Exascale storage a tool of the Scientist
  - Tractable data management
  - Comprehensive interaction
  - Move compute to data or data to compute as appropriate
- Overcome today's IO limits
  - Multi-petabyte datasets
  - Explosive growth of metadata
  - Horizontal scaling & jitter
- Support unprecedented fault tolerance
  - Deterministic interactions with failing hardware & software
  - Fast & scalable recovery
  - Enable multiple redundancy & integrity schemes

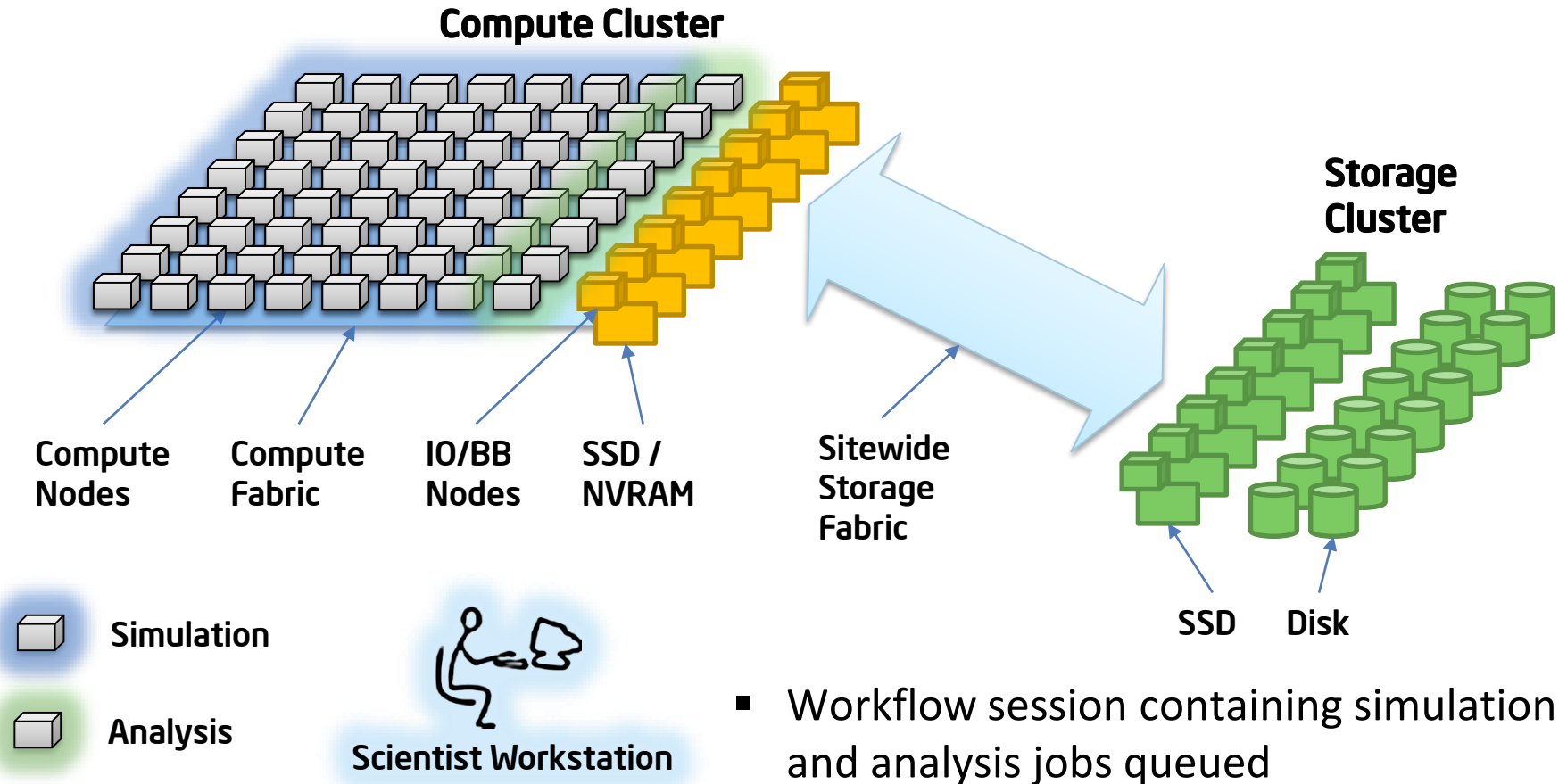


# Fast Forward I/O Architecture

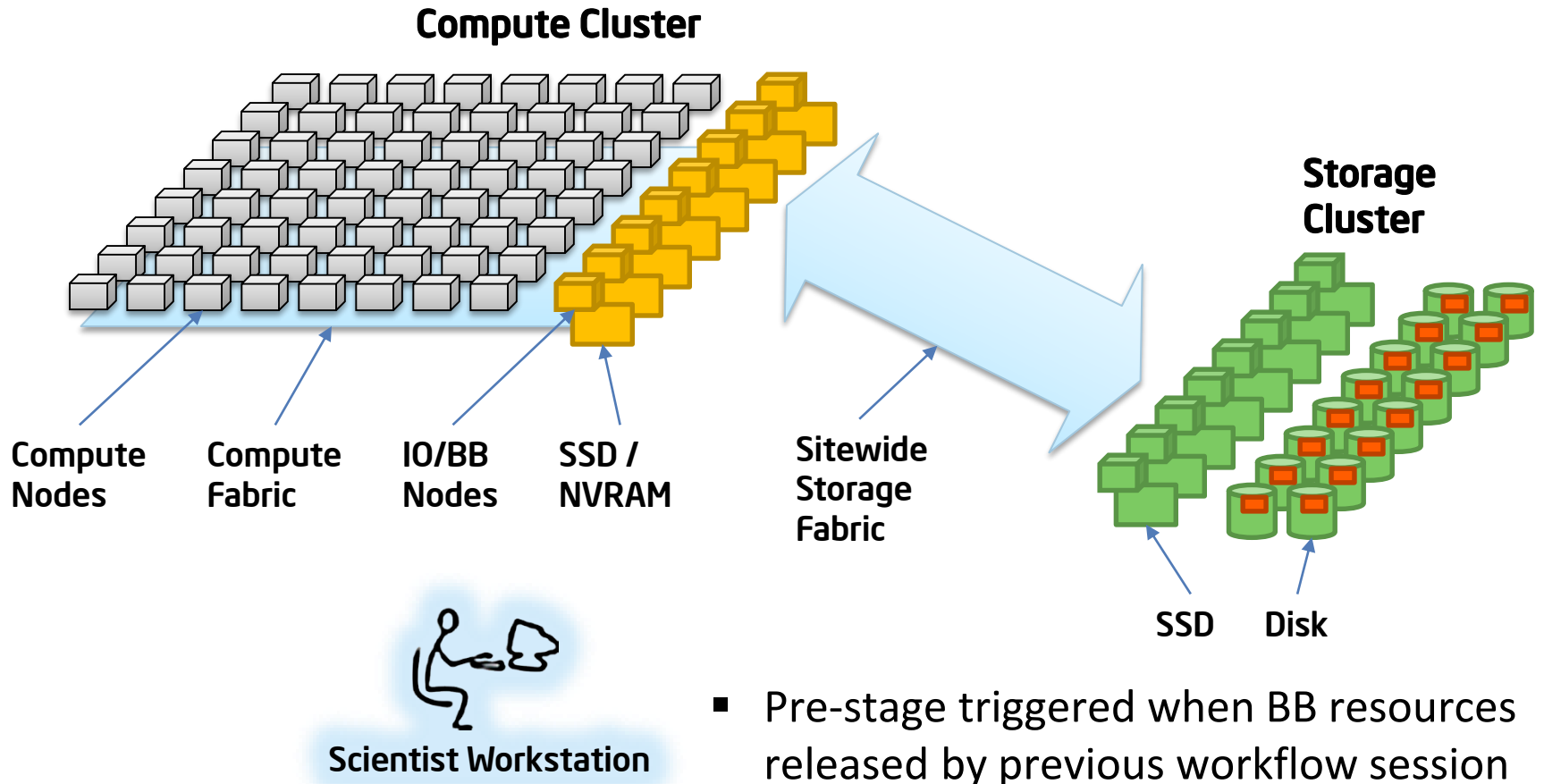


- Global storage mounted on compute cluster IO nodes and scientist's workstation
- I/O forwarding from compute to IO nodes

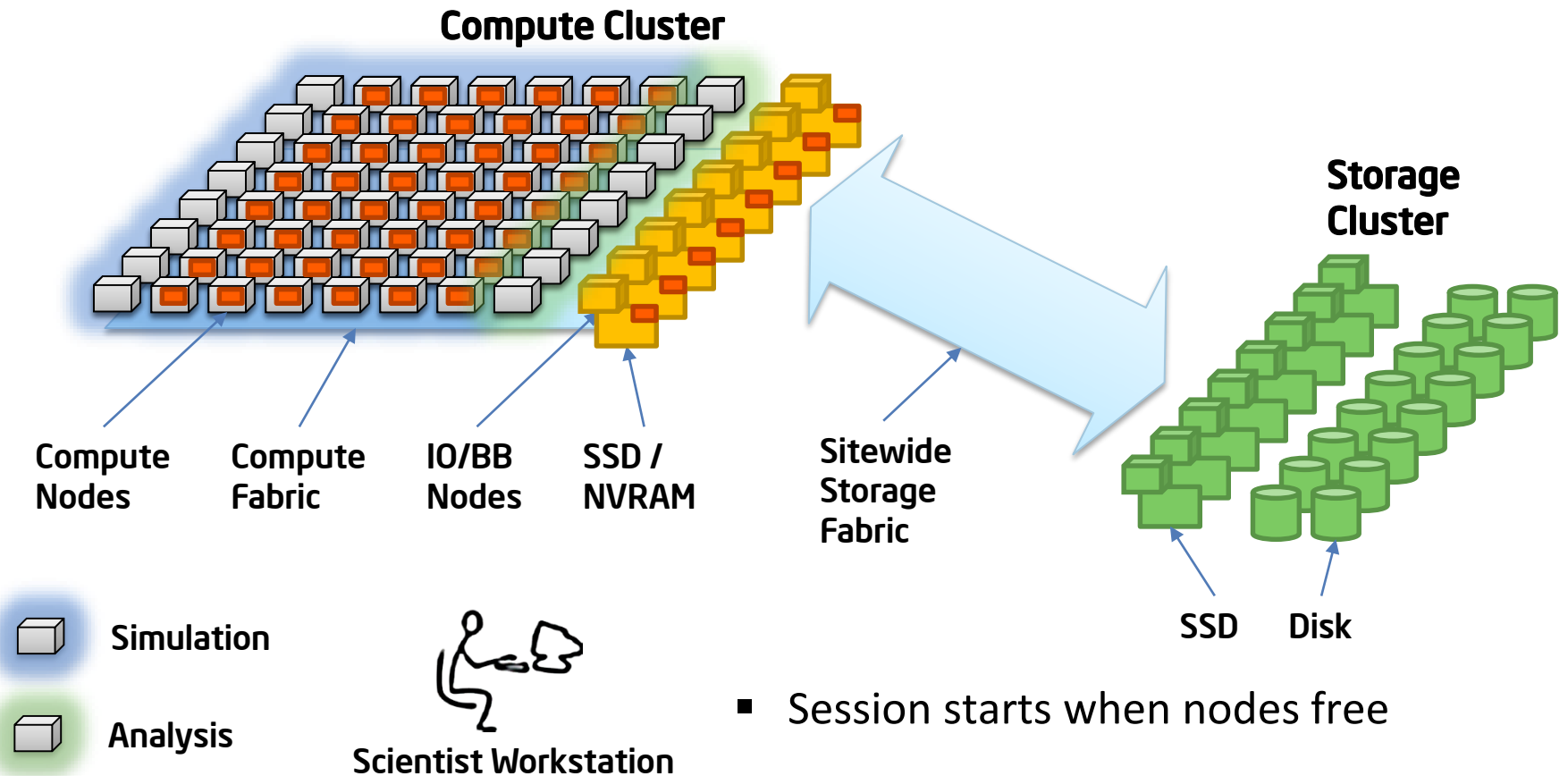
# Workflow: Simulation + In-transit Analysis



# Workflow: Pre-stage to Burst Buffer

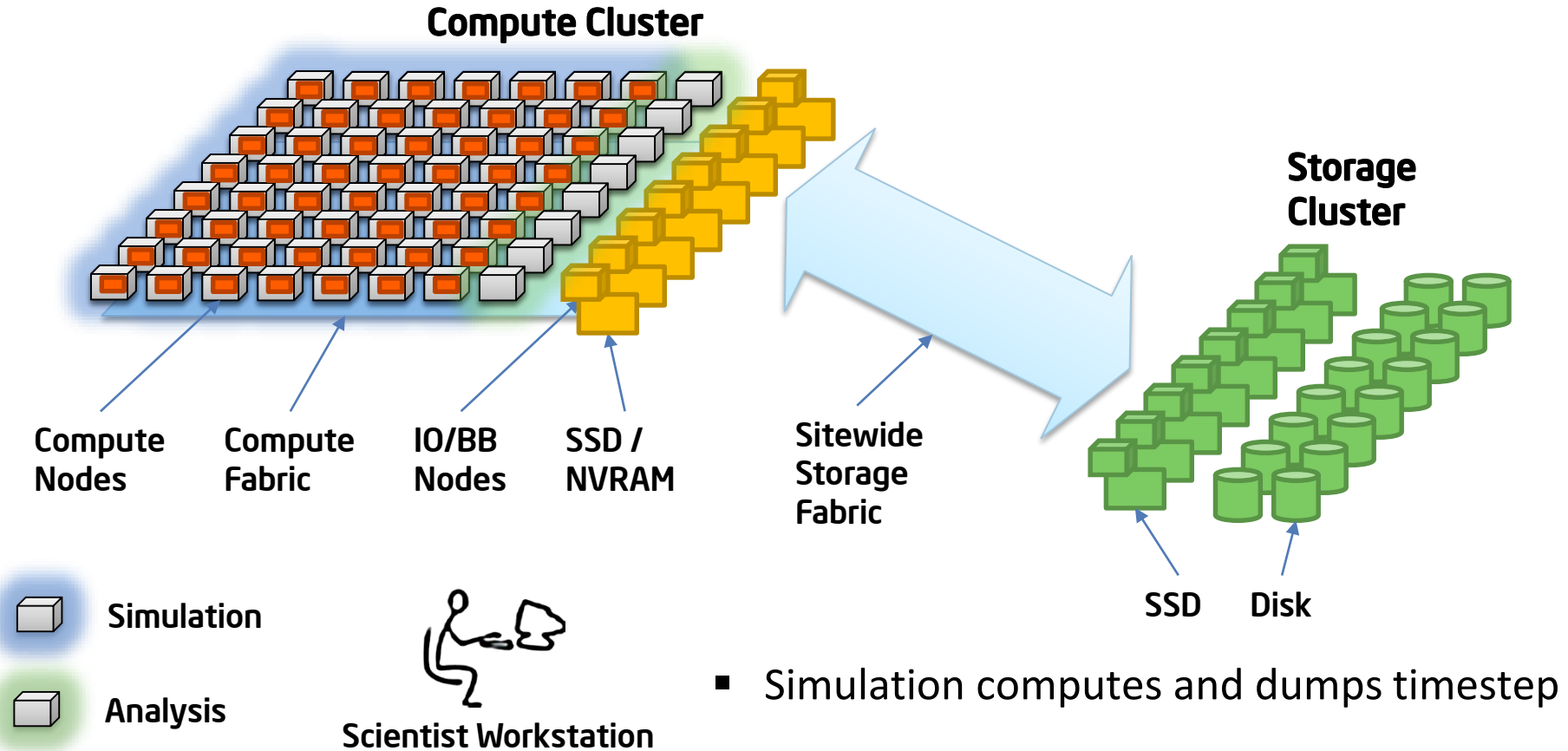


# Workflow: Start Session



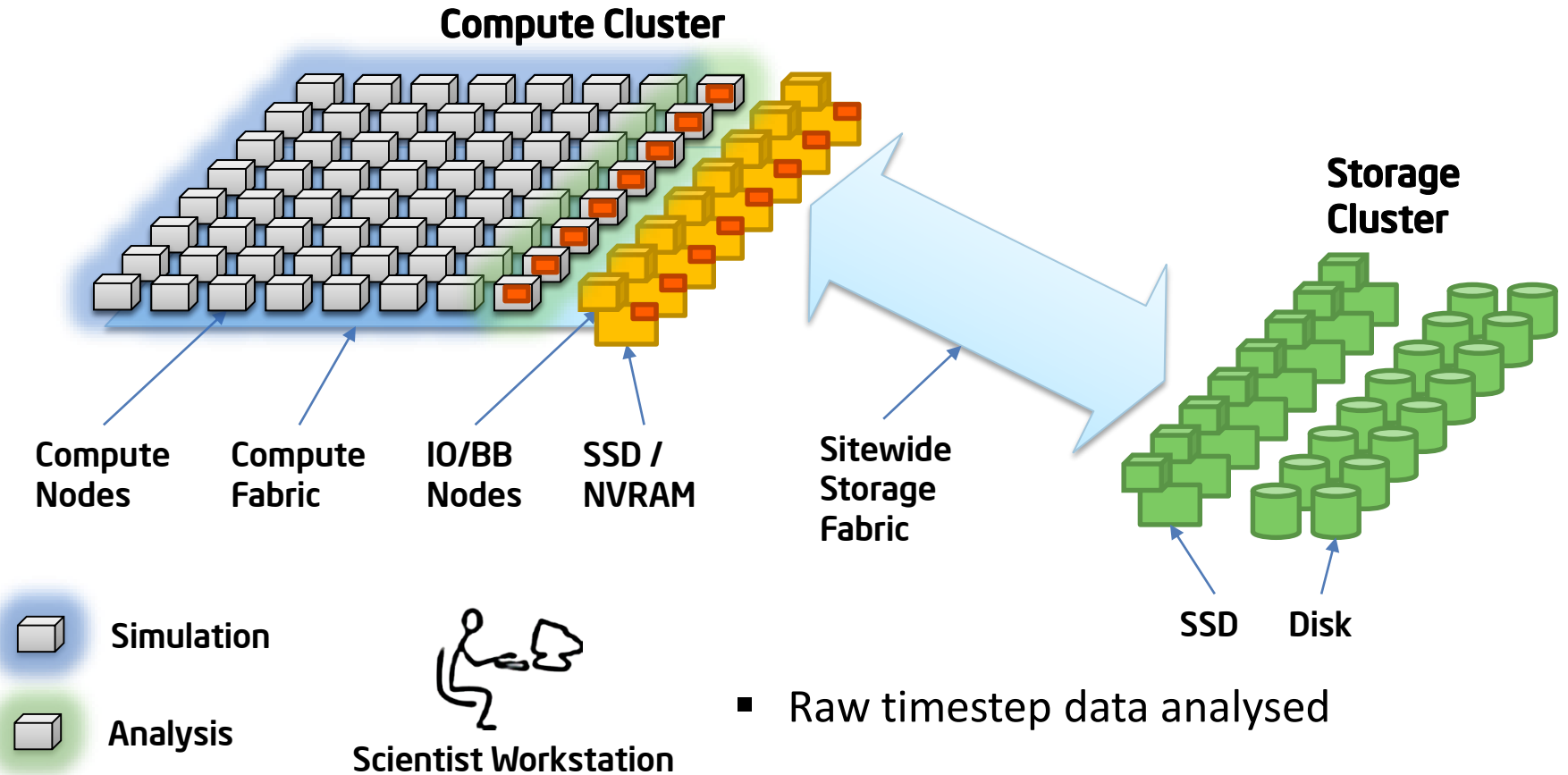
- Session starts when nodes free
- Previous session may still be persisting data from BB to global storage

# Workflow: Dump Timestep



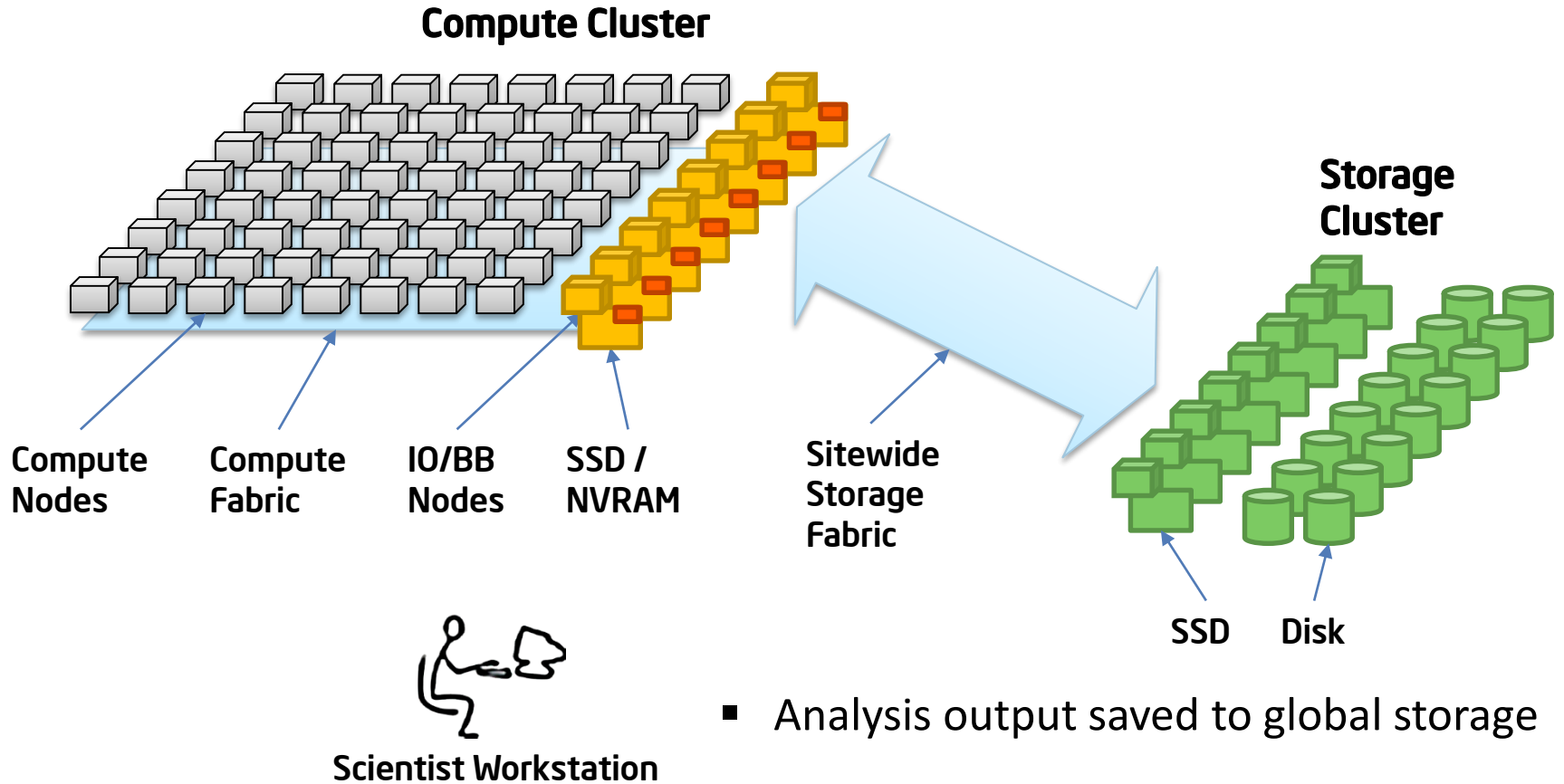


# Workflow: In-transit Analysis



- Raw timestep data analysed
- Analysis data saved to BB
- Raw timestep may be discarded

# Workflow: Persist to Global Storage

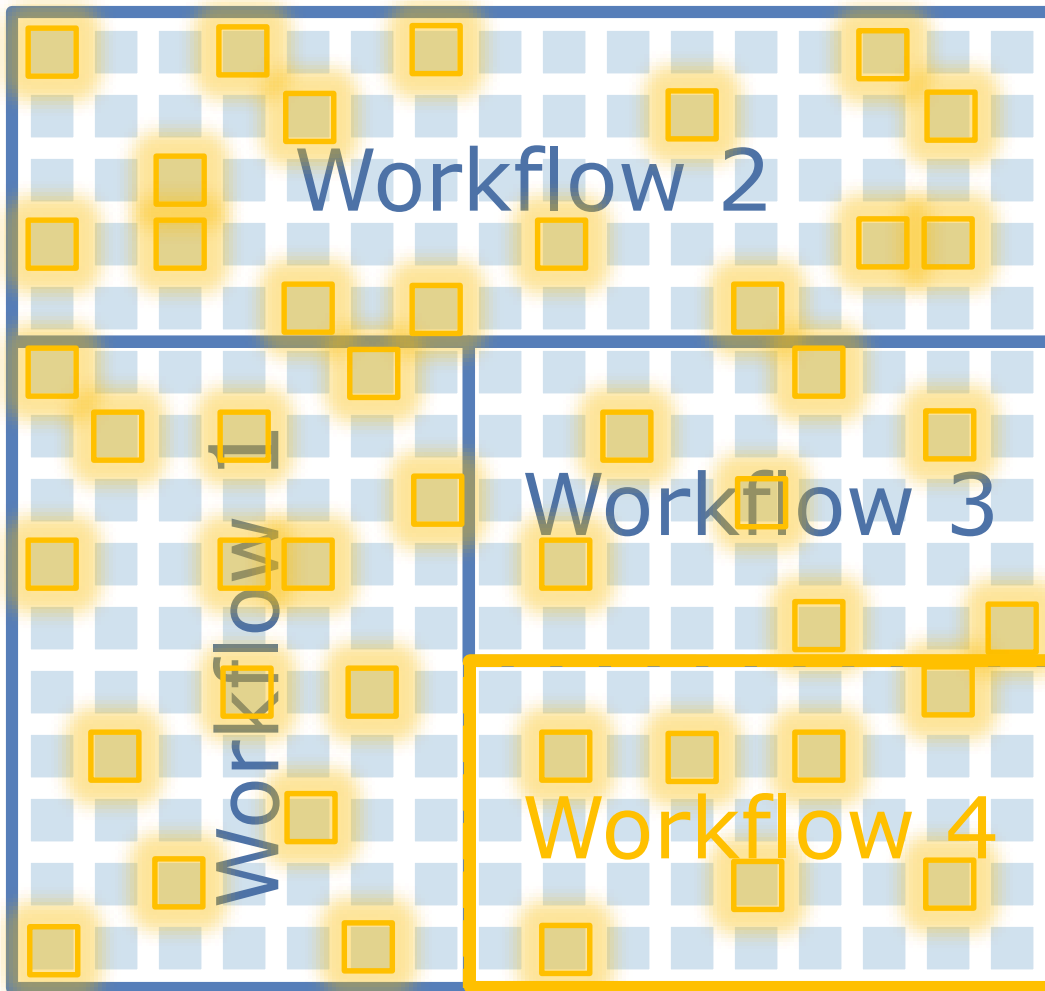


# Ubiquitous NVRAM

- O(1TB) compute node-local storage
- Instant-on
  - 0 power standby
- Load-store byte-granular access
  - Invites Distributed Persistent Memory programming models
  - Order of magnitude larger in-core working sets
- Storage fully leverages fabric

	Disk	Edge BB	NVRAM
Checkpoint / Search	1 hour	6 minutes	6 seconds
Capacity (# datasets)	30	3-5	10-30

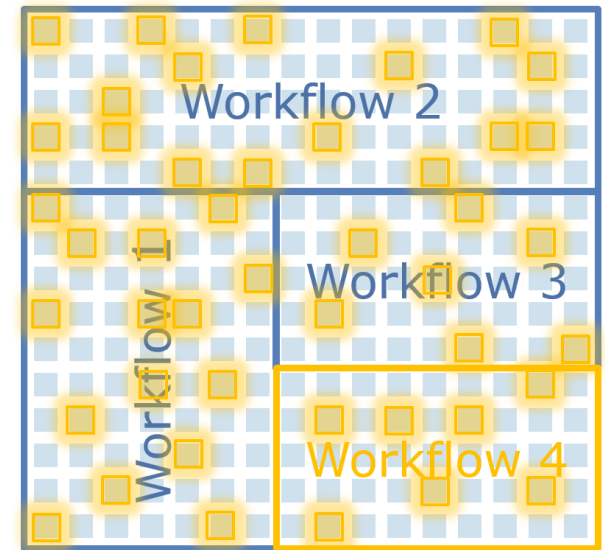
# Scheduling Persistent Memory



- Issues
  - Space management
  - Interference & jitter
- Workflow Session 4 ready to run
- Data not local
- Migrate
- Workflow Session 4 started

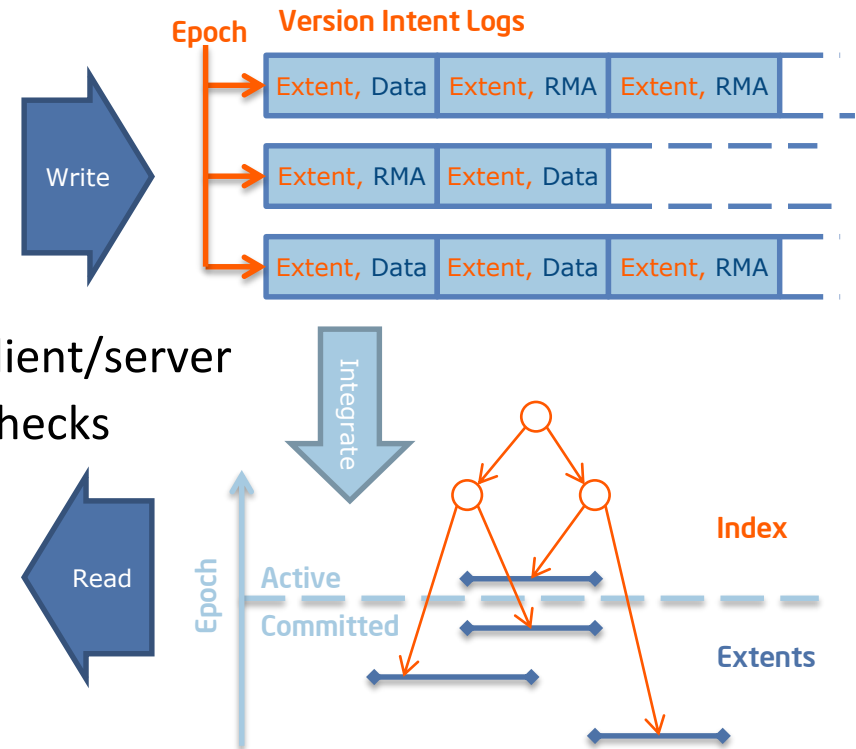
# Persistent Memory v. Storage

- Persistent Memory is fast but it's...
  - Local to the process using it
  - Inaccessible on node failure
  - Fixed schema
- Storage may be slower but it's...
  - Globally accessible
  - Consistent & durable
  - Snapshotable / Cloneable / Migrateable
- APIs required to...
  - Convert PM ↔ Storage
    - Persist / Instantiate Distributed Persistent Memory images
    - PM schema conversion
  - Support workflow scheduler integration
    - Data-aware process instantiation
    - Process-aware data migration



# DAOS-M

- Client & Server OS bypass
- Connectionless
  - Peer-to-peer connectivity = ~100x client/server
  - Heavyweight security / ownership checks once on container open
- Memory VOSD
  - PM programming model
    - No block I/O stack latency
    - Byte granular
  - Read
  - Extremely low latency
  - committed writes integrated on index traversal
  - Write
    - Incoming data and metadata logged
    - Integration processes inserts into index



# Summary

- Ubiquitous NVRAM changes the game
- 3 order of magnitude step change in performance from disk
  - Terabytes/s -> Petabytes/s
  - mS latency ->  $\mu$ S latency
- Workflows will change to exploit
  - Persistent Memory programming models
  - Data aware workflow scheduling
- Storage software must change to exploit
  - Same transactional guarantees required
  - End-to-end OS bypass required
  - Scalable comms/security context establishment
  - More I/O stack configuration flexibility

