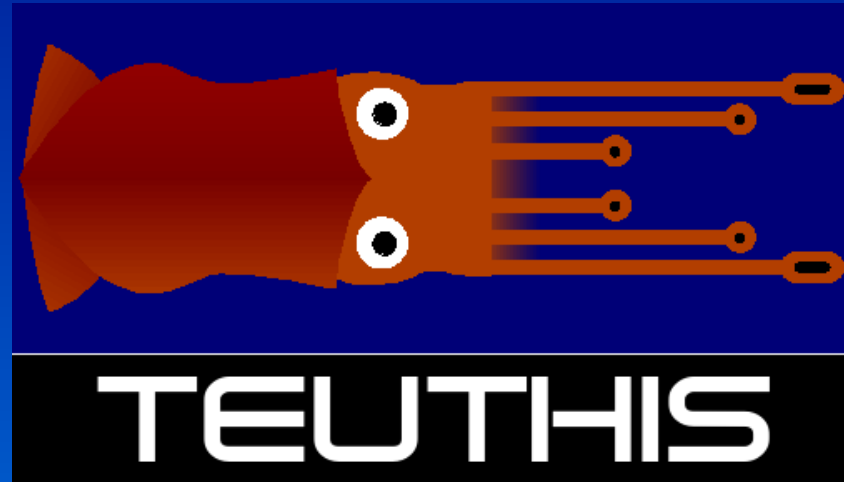
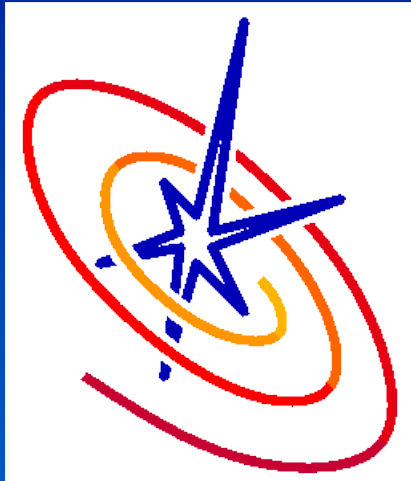


Galaxy Cluster Simulations at NCSA using FLASH and Teuthis



Paul Ricker

National Center for Supercomputing Applications
University of Illinois at Urbana-Champaign

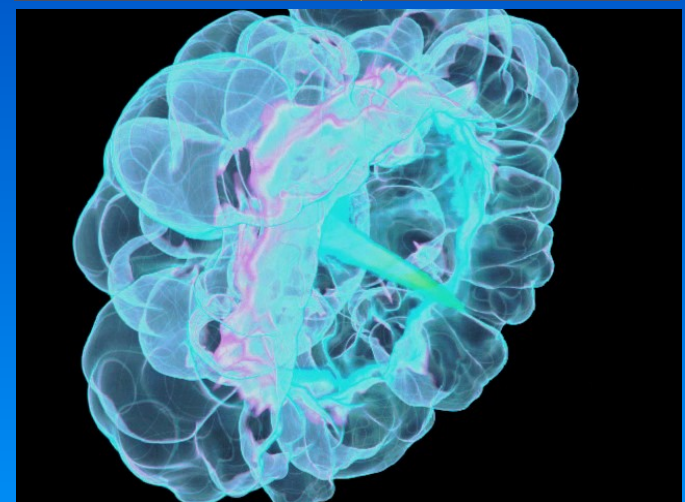
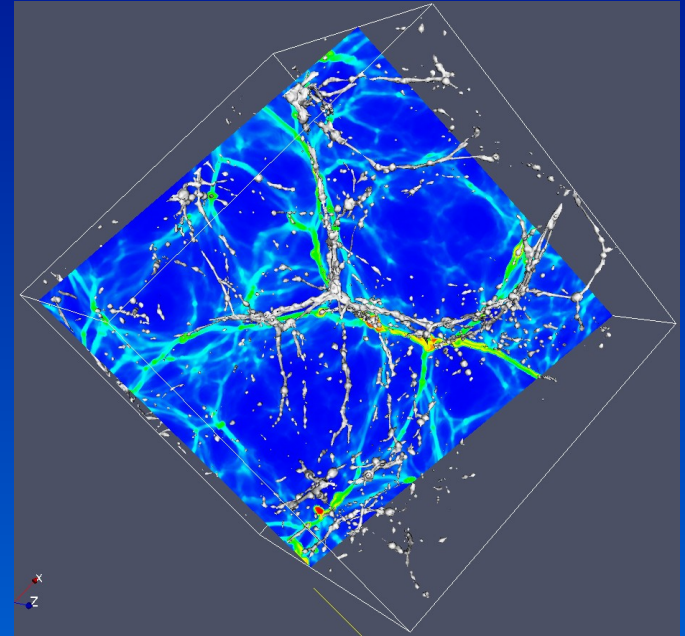
SC06

November 14, 2006



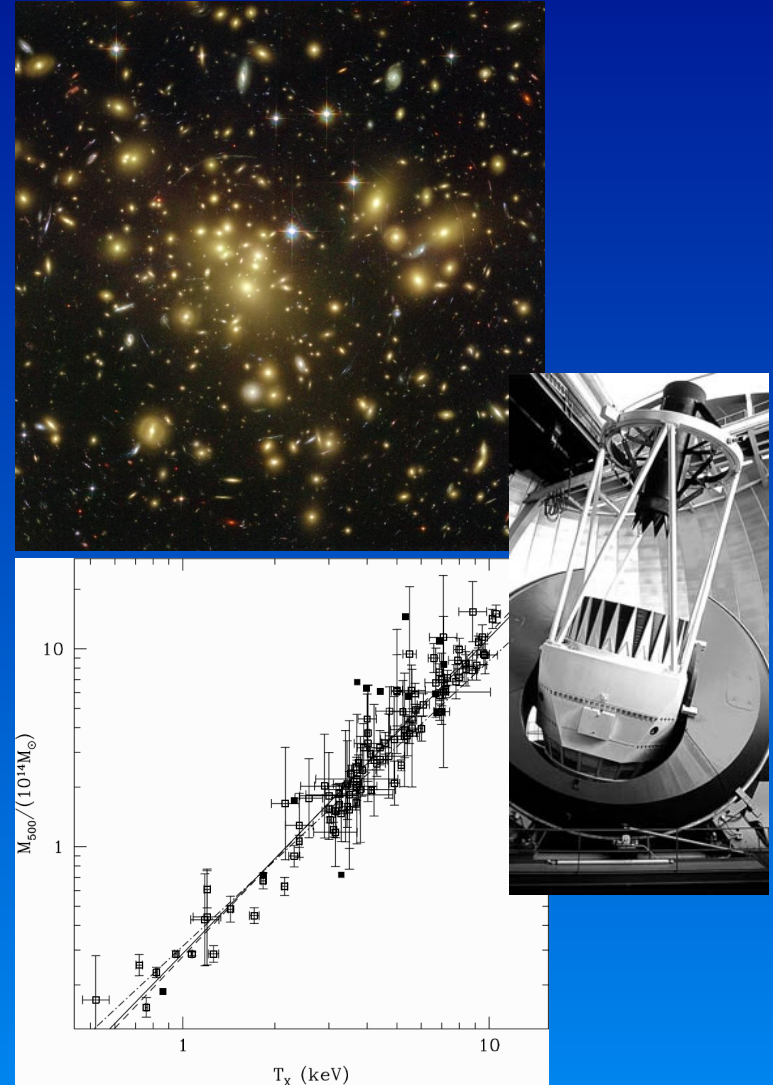
Astrophysical simulations

- **Typical problems**
 - Cosmological structure formation
 - Supernova explosions
 - Planetary disks
- **Algorithms**
 - Partial differential equations
 - Feature extraction
 - Statistical analysis
 - Ray casting
- **Data types**
 - Large checkpoint files (100 GB+)
 - Data subsets for plotting
 - Object/feature catalogs
 - Images and movies



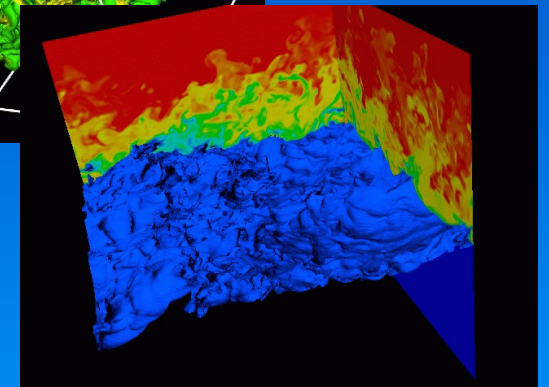
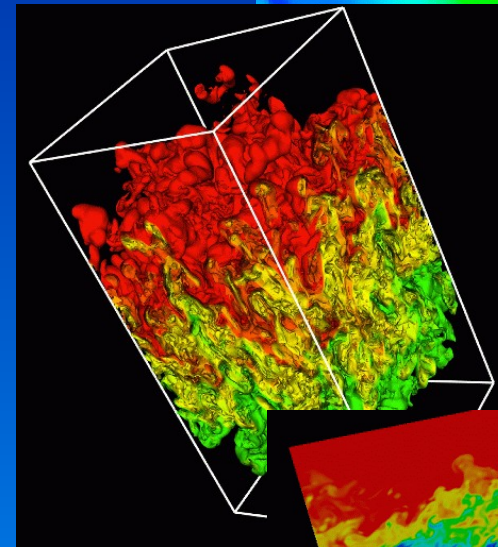
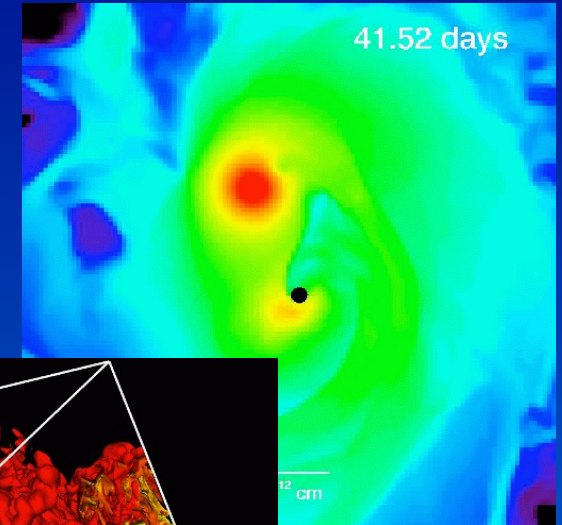
Galaxy clusters

- Largest gravitationally bound objects in the Universe
 - ~ few to > 1000 galaxies
 - $\sim 10^{13-15}$ solar masses
 - $\sim \text{few} \times 10^6$ light-years across
- Simulations
 - Form through gravitational instability
 - Create many to study statistics
 - Simulated observations to compare with cluster surveys



What is FLASH?

- Astrophysical simulation code developed under DOE ASCI program at University of Chicago
 - Adaptive mesh refinement (AMR)
 - Particles and gravity
 - Nuclear reactions
- **Community code**
 - Free: <http://flash.uchicago.edu>
 - Modular framework (560,000 lines)
 - MPI parallel
 - Validated against experiments
 - 200+ users





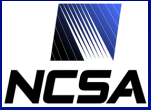
What is Teuthis?

- **A control panel**
 - Remotely configure and build applications
 - Submit and track remote jobs
 - Painlessly create parameter studies and restart jobs
- **A data manager**
 - Stage and archive data
 - Keep track of where datasets are stored
- **A notebook**
 - Organize job metadata by purpose and disposition
- **An aid to collaboration**
 - Share notebook files with collaborators

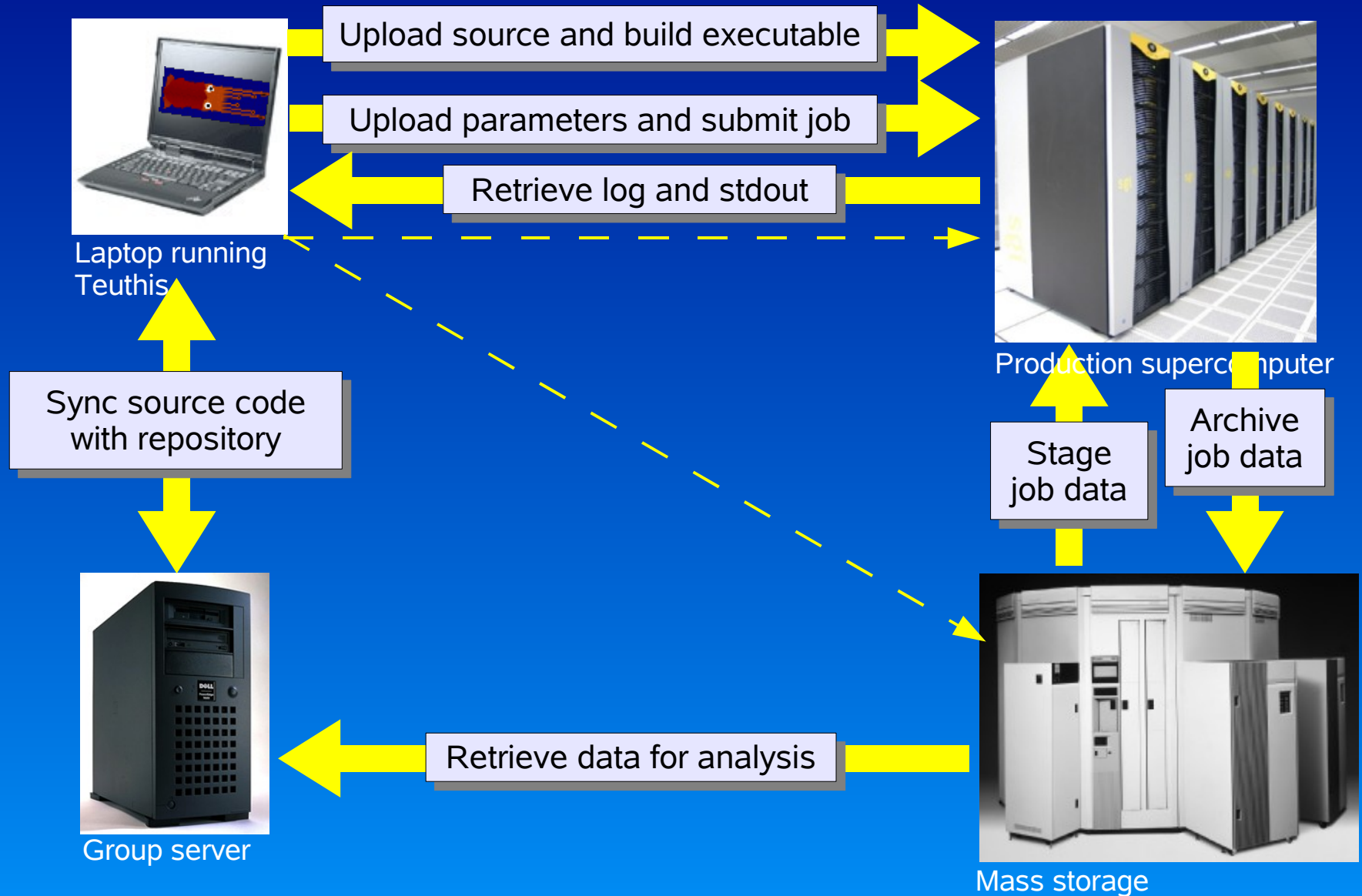
Design philosophy

- **Small is beautiful**
 - Small footprint – run on Tungsten someday?
 - Minimal prerequisites – avoid dependency hell
- **Exploit others' expertise**
 - Use external tools when possible...
 - ... but only for “extra” functionality
- **Think locally**
 - Single point of interaction with one's work
 - Local metadata store complete and authoritative
- **Don't get tied down**
 - Separate GUI from backend
 - Everything open-source and cross-platform





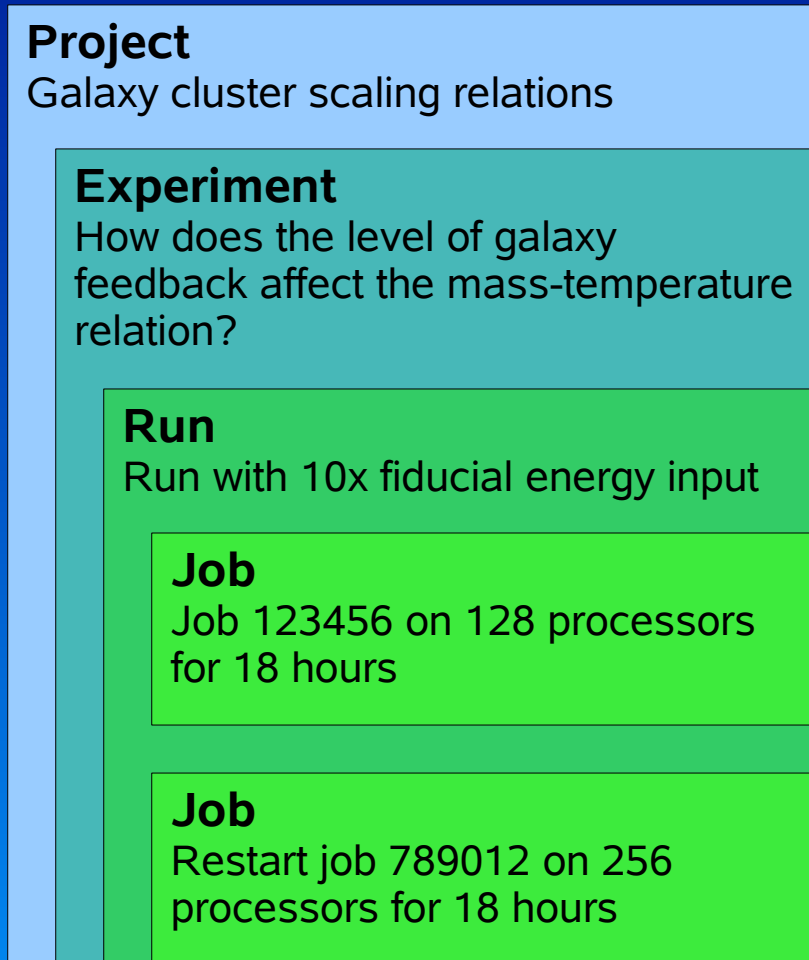
Running simulations with Teuthis





Objects manipulated by Teuthis

Workflow hierarchy



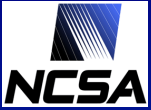
Resources

Application

- Accepts text parameter file
- Executes noninteractively
- May need to be compiled
- Produces log file, screen output, data files

Machine

- Login host
- Access method
- Queuing system
- Paths and commands



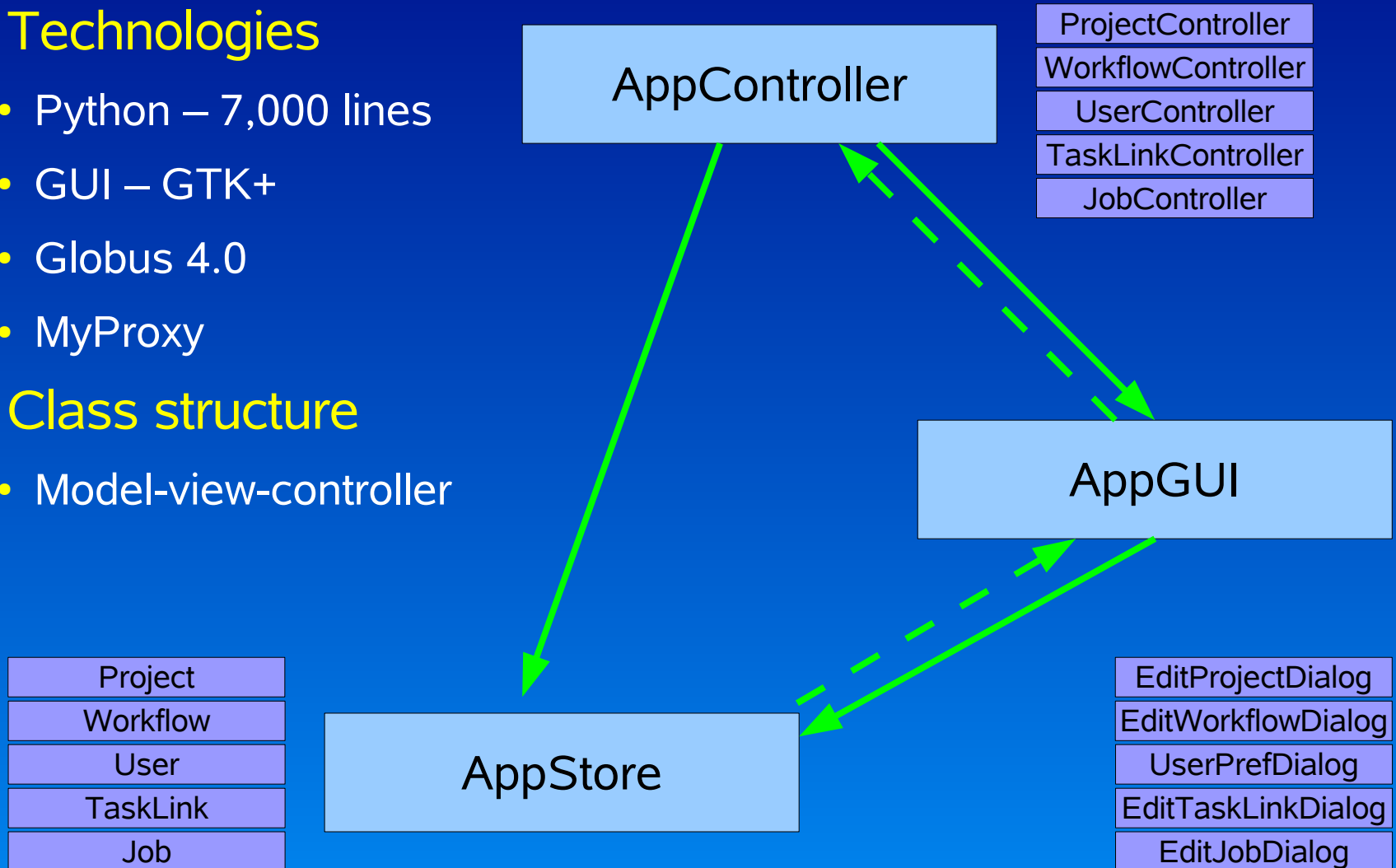
Under the hood

Technologies

- Python – 7,000 lines
- GUI – GTK+
- Globus 4.0
- MyProxy

Class structure

- Model-view-controller

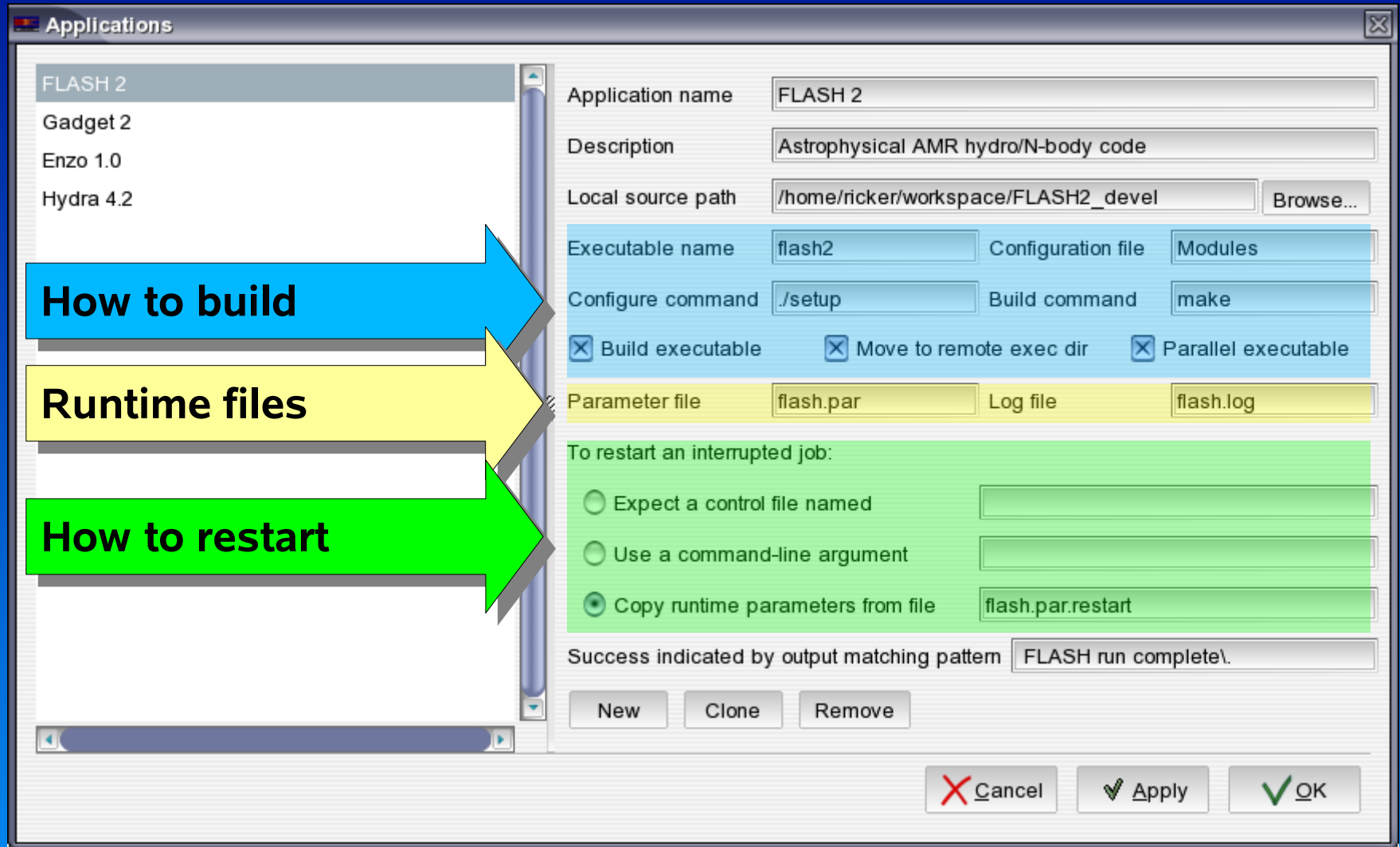




Project view

Simulation Manager 1.0			
File View Settings Help			
Name	Description	Status	Date last modified
▼ FLASH testing	Testing Simulation Manager using FLASH		Tue Oct 4 17:1
▼ Basic Sedov test (local)	Test of local jobs		Tue Oct 4 19:14:5
Run A			Tue Oct 4 19:15:0
▼ Basic Sedov test (cobalt)	Test of jobs on a machine with PBS queuing and using :		Wed Oct 5 01:43:
▼ Run A	Single run using default parameters to test job submission	Complete	Wed Oct 5 01:17:
Job A0004	Original	28909 [21:26 10/04/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 01:50:2
Job A0004	Restart of 28909	28910 [21:27 10/04/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 01:16:7
▼ Sedov scaling test (cobalt)	Test with varying number of processors		Wed Oct 5 01:43:
▼ Run A1		Complete	Wed Oct 5 02:07:
Job A10001	Original	28924 [01:44 10/05/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 02:05:3
▼ Run A2		Complete	Wed Oct 5 02:07:
Job A20001	Original	28925 [01:45 10/05/2005] 2 CPUs/00:10 (Complete) Successful completion	Wed Oct 5 02:05:2
▼ Run A4		Complete	Wed Oct 5 02:07:
Job A40001	Original	28926 [01:45 10/05/2005] 4 CPUs/00:10 (Complete) Successful completion	Wed Oct 5 02:05:7
▼ Run A8		Complete	Wed Oct 5 02:07:
Job A80001	Original	28927 [01:45 10/05/2005] 8 CPUs/00:10 (Complete) Successful completion	Wed Oct 5 02:05:0
▼ Sedov test with varying parameter (cobalt)	Test of jobs with a single varying parameter (lrefine_mi		Wed Oct 5 02:09:
▼ Run A		Complete	Wed Oct 5 02:08:
Job A0001	Original	28928 [01:52 10/05/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 01:58:3
Job A0001 Copy	Original	29231 [15:01 10/05/2005] 1 CPU/00:10 (Complete) No data	Wed Oct 5 15:01:5
▼ Run B		Complete	Wed Oct 5 02:08:
Job B0001	Original	28929 [01:52 10/05/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 01:58:4
▼ Run C		Complete	Wed Oct 5 02:08:
Job C0001	Original	28930 [01:53 10/05/2005] 1 CPU/00:10 (Complete) Successful completion	Wed Oct 5 01:58:5
▼ Run D		Complete	Wed Oct 5 02:08:
Job D0001	Original	28931 [01:53 10/05/2005] 1 CPU/00:20 (Complete) Successful completion	Wed Oct 5 01:59:0
▼ Run E		In progress	Wed Oct 5 02:08:
Job E0001	Original	28932 [01:53 10/05/2005] 1 CPU/00:20 (Complete) Successful completion; k	Wed Oct 5 02:02:5
▼ Run F		In progress	Wed Oct 5 02:09:
Job F0001	Original	28933 [01:53 10/05/2005] 1 CPU/00:30 (Complete) Exceeded MAXBLOCKS	Wed Oct 5 01:55:2

Configuring applications



The screenshot shows the 'Applications' configuration window for 'FLASH 2'. The window is divided into a left sidebar and a main configuration area. The sidebar lists 'FLASH 2', 'Gadget 2', 'Enzo 1.0', and 'Hydra 4.2'. The main area contains the following fields and options:

- Application name: FLASH 2
- Description: Astrophysical AMR hydro/N-body code
- Local source path: /home/ricker/workspace/FLASH2_devel (with a 'Browse...' button)
- Executable name: flash2
- Configuration file: Modules
- Configure command: ./setup
- Build command: make
- Build executable:
- Move to remote exec dir:
- Parallel executable:
- Parameter file: flash.par
- Log file: flash.log
- To restart an interrupted job:
 - Expect a control file named
 - Use a command-line argument
 - Copy runtime parameters from file (with value: flash.par.restart)
- Success indicated by output matching pattern: FLASH run complete\.

At the bottom of the window are buttons for 'New', 'Clone', 'Remove', 'Cancel', 'Apply', and 'OK'. Three callout arrows are overlaid on the left side of the window:

- A blue arrow labeled 'How to build' points to the 'Build executable', 'Move to remote exec dir', and 'Parallel executable' options.
- A yellow arrow labeled 'Runtime files' points to the 'Parameter file' and 'Log file' fields.
- A green arrow labeled 'How to restart' points to the 'To restart an interrupted job' section.



Tested applications

Application	Exec name and args	Config file	Config command	Build command	Move executable	Parallel	Parameter file	Log file	Auto restart method	Notes
FLASH 2.x	flash2	Modules	/setup	gmake	optional	yes	flash.par	flash.log	Copy from flash.par.restart; or command line argument "-chk_file" followed by manual addition of checkpoint file name	flash.par.restart not available in standard distribution; need patch Need to set up site directory for remote site Leave log_file parameter unset
Gadget 2	Gadget2 gadget.param	Makefile	N/A	gmake	optional	yes	gadget.param	info.txt	Command line argument "1"	Upload custom makefile as your configuration file Use "." for OutputDir parameter Leave InfoFile parameter unset
Enzo 1.0.1	enzo.exe EnzoParms	N/A	/configure -bindir=XX	cd amr_mpi/src; gmake mach-YY; gmake; gmake install	yes	yes	EnzoParms	OutputLevelInformation	Command line argument "-r"; must manually add name of last restart file	XX = absolute path to build directory Need to set up Make.mach.YY file for remote machine; place in config directory
Hydra 4.2	hydra	makeflags	N/A	make clean; make	yes	no	prun.dat	pr0001.log	None; manually edit prun.dat	May need to create a new src/system.YY file for remote machine YY Modify src/dumpdata.F, src/readdata.F, and src/gravsubs.F to read/write to / rather than data/ To change array sizes, edit include/psize.inc on local machine and sync source Upload custom makeflags file as your configuration file; set RUNDIR to "." Use 0001 as run number in prun.dat



Configuring machines

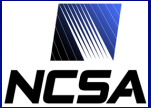
The screenshot shows a configuration window titled "Machines" with a list of machines on the left: cobalt, tungsten, and tungsten (ssh-agent). The "cobalt" machine is selected, and its configuration is shown in the main area. The configuration is divided into several sections:

- Basic information:** Machine name: cobalt, Description: SGI Altix at NCSA, Type: Compute (selected), Online data (unselected).
- Access:** Host: cobalt.ncsa.uiuc.edu, User ID: [empty], Access method: gsissh+uberftp, Realm: [empty].
- Job submission:** Job template: [empty], Queues: standard, OS type: Unix.

Two callout boxes highlight specific configuration areas:

- Access methods:** local, ssh, ssh-agent, Kerberos, GSI
- Queuing methods:** Unix (no queue), PBS, LSF, LoadLeveler, User-specified

At the bottom of the window are buttons for "New", "Clone", "Remove", "Cancel", "Apply", and "OK".



Experiment dialog

2 Set up input and archive data

1 Configure and build application

3 Set up execution host

4 Choose parameters to fix/vary

5 Generate runs

The screenshot shows the 'Experiment properties' dialog box with the following sections highlighted in green:

- Information:** Experiment name: Sedov test with file transfer; Description: Sedov test on a remote machine with two third-party file transfers.
- Data:** Src machine: tungsten; Src files: /u/ac/ricker/input/la128.tar, /u/ac/ricker/input/lcdms64.tar.
- Application:** Application: FLASH 2.4; Configuration file: (empty); Remote build dir: /u/ac/ricker/build/test/FLASH_2.4; Config command: ./setup sedov -auto; Build command: make; mv flash2 /u/ac/ricker/exec; Executable to use: /u/ac/ricker/exec/flash2.
- Execution:** Exec machine: cobalt; Queue: normal; Account: (empty); # of CPUs (range): 1; Tiling: 1; Mem/node (MB): 1000.
- Parameters:** Template: /home/ricker/flash.par; Parameter 1: lrefine_max, Range: 1-6; Parameter 2: (empty), Range: (empty); Parameter 3: (empty), Range: (empty); Parameter 4: (empty), Range: (empty); Parameter 5: (empty), Range: (empty); Parameter 6: (empty), Range: (empty); Parameter 7: (empty), Range: (empty); Parameter 8: (empty), Range: (empty).

Buttons at the bottom: Cancel, Apply, OK.



Job dialog

Job properties

Local job information

Local job ID: C0001
 Comments: Original
 Disposition: No data
 Created: Thu Nov 10 03:06:53 2005
 Last modified: Thu Nov 10 03:06:53 2005

Application

Application: FLASH 2.4
 Executable to use: /u/ac/ricker/exec/flash2
 Comments:

Execution

Exec machine: cobalt Wall time: 00 h 00 m 00 s
 Queue: normal Account:
 No. of CPUs: 1 Tiling: 1 Mem/node (MB): 1000

Data

Src machine: tungsten
 Src files: /u/ac/ricker/input/la128.tar
 /u/ac/ricker/input/lcdms64.tar
 Dest machine: tungsten
 Dest path: /u/ac/ricker/test

Actions

View parameters Submit Status
 View log file View output Continue job Archive data

Remote job information

Remote job ID: 000000
 Submitted: Not yet submitted
 Run status: Unsubmitted

Cancel Apply OK

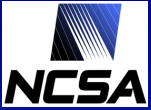
1 Submit

2 Transfer

5 Archive

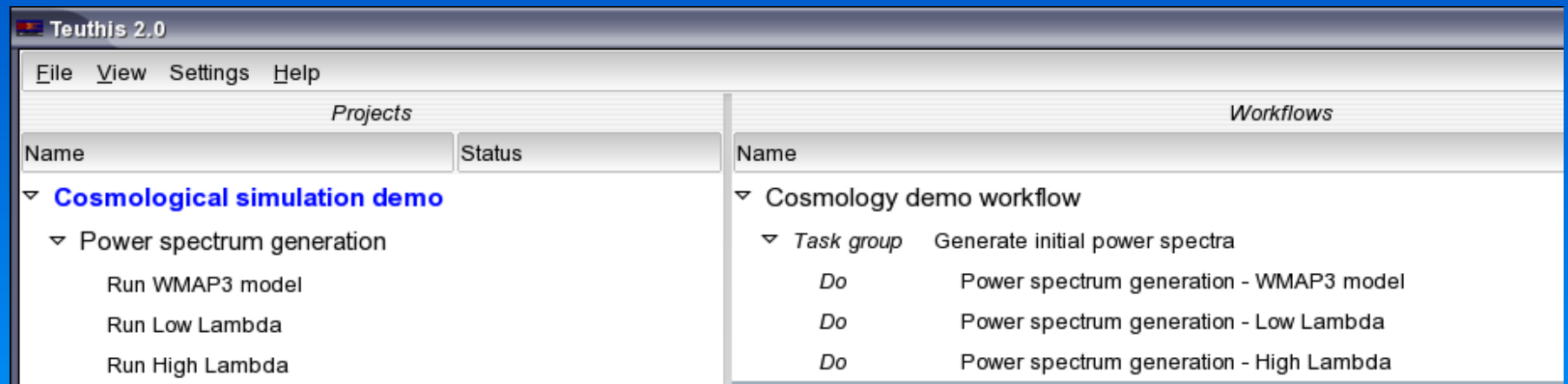
4 View

3 Check



New features in Teuthis 2.0

- Data
 - Background file transfers with retry
- Job submission and monitoring
 - GRAM job submission
 - Workflow management
- User interface
 - Refactoring
 - Usability improvements



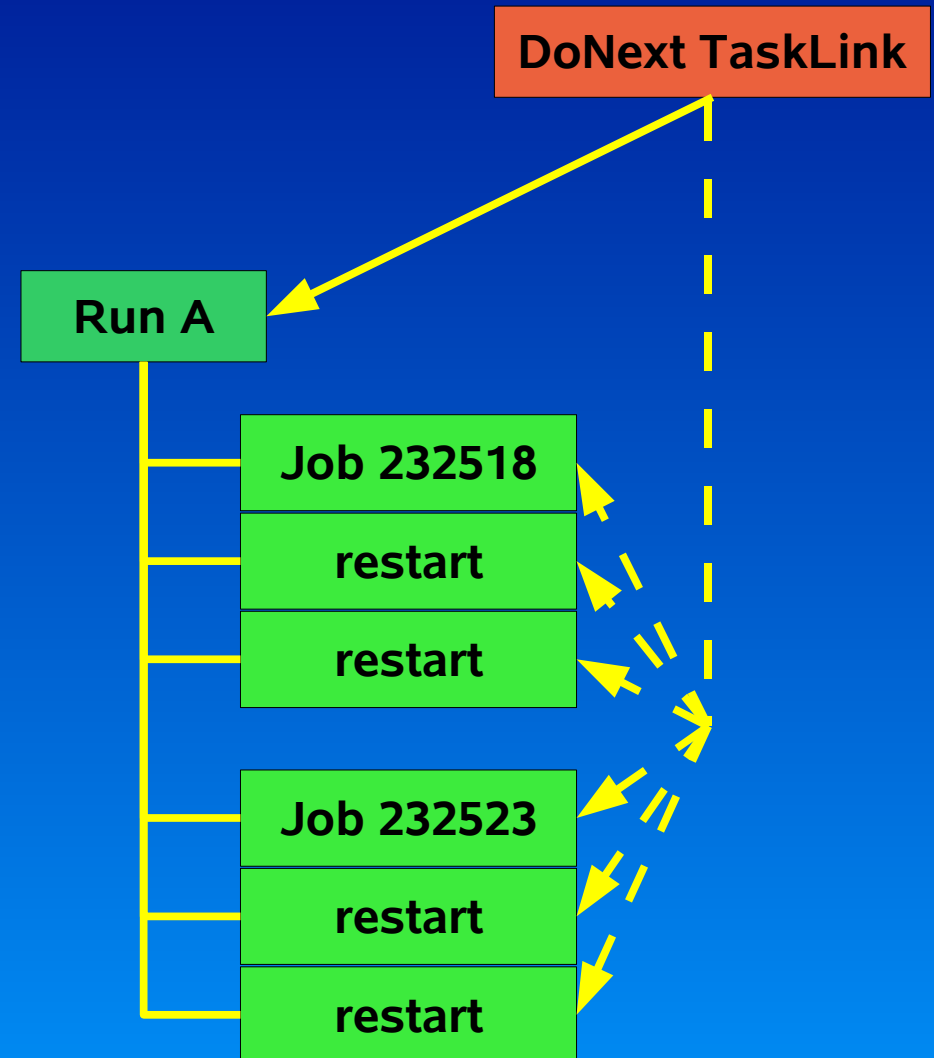
Teuthis workflow management

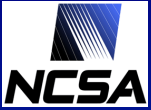
• Task links

- Atomic workflow unit
- Linked to a particular run
- Invocation triggers new job: stage in – exec – stage out
- Automatic job continuation
- Pattern matching conditionals

• Types

- Data source – Static
- Unconditional – DoNext, DoTogether
- Conditional – DoIf, WhileDo
- Grouping – Task Groups





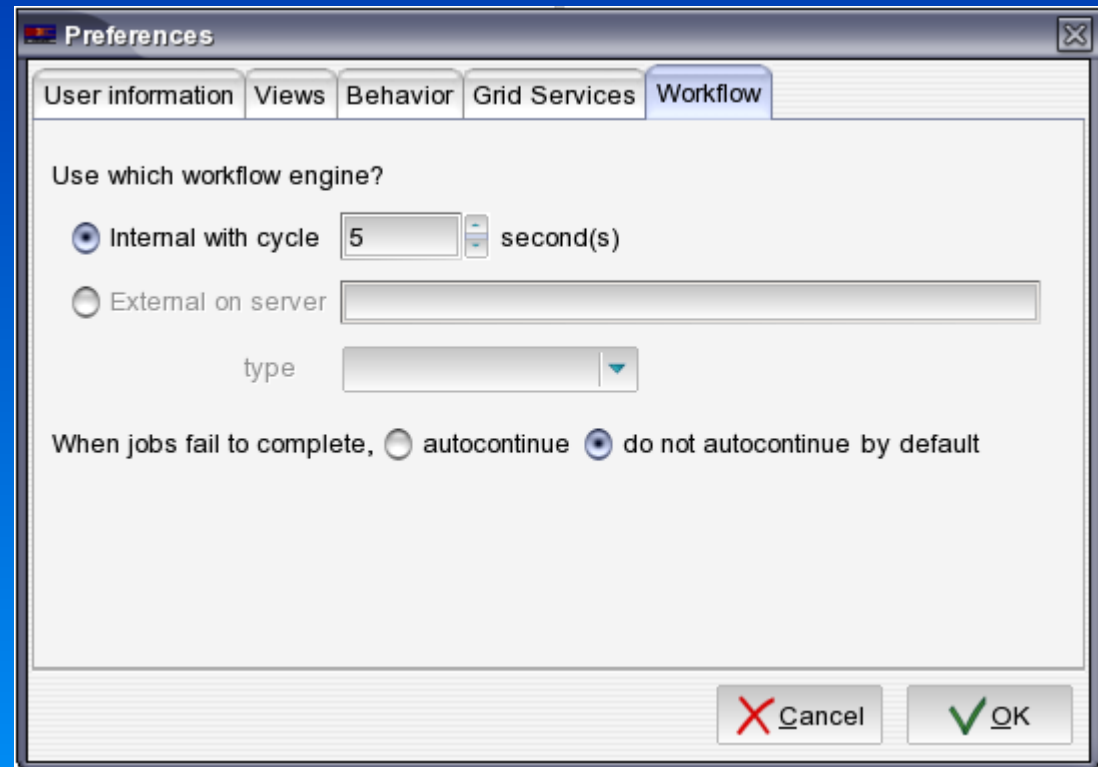
Creating and editing task links

- Creating and editing task links
 - Add from popup menus
 - Drag runs/jobs from Project Pane

Teuthis workflow management

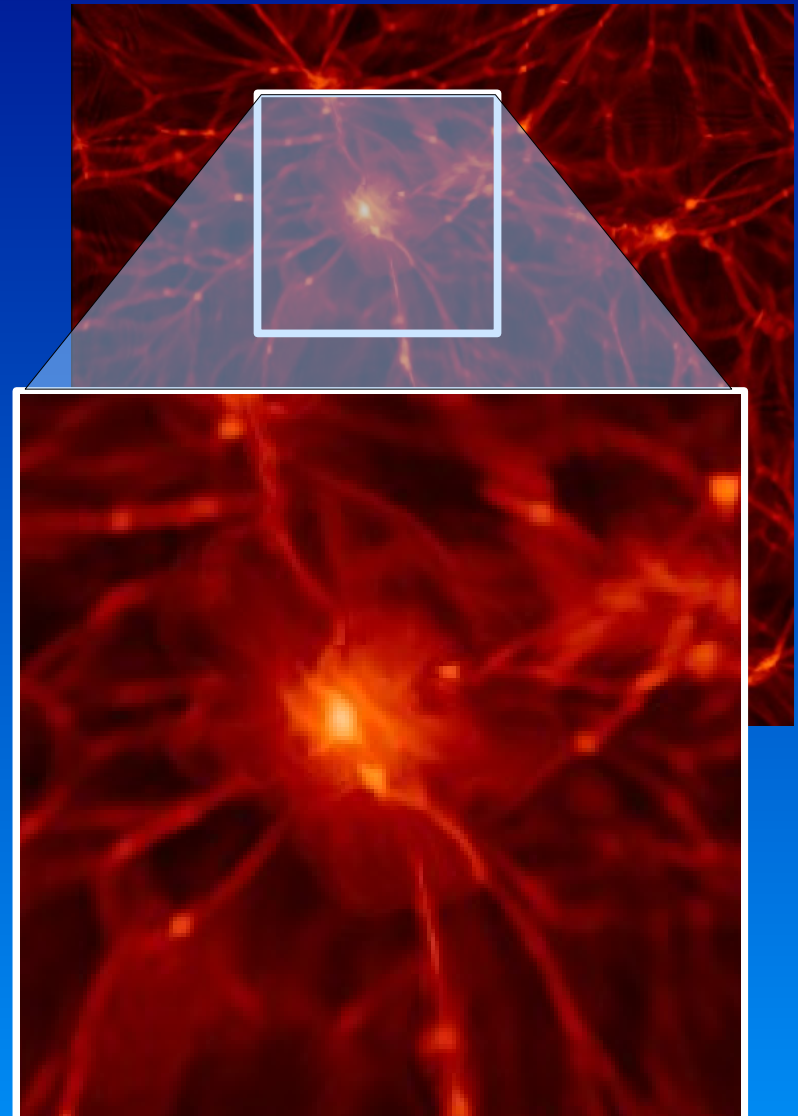
• Workflow execution choices

- Internal engine: use polling to test job status and advance at preset interval
- External engine: hand script off, watch for messages
- Pause/halt/resume



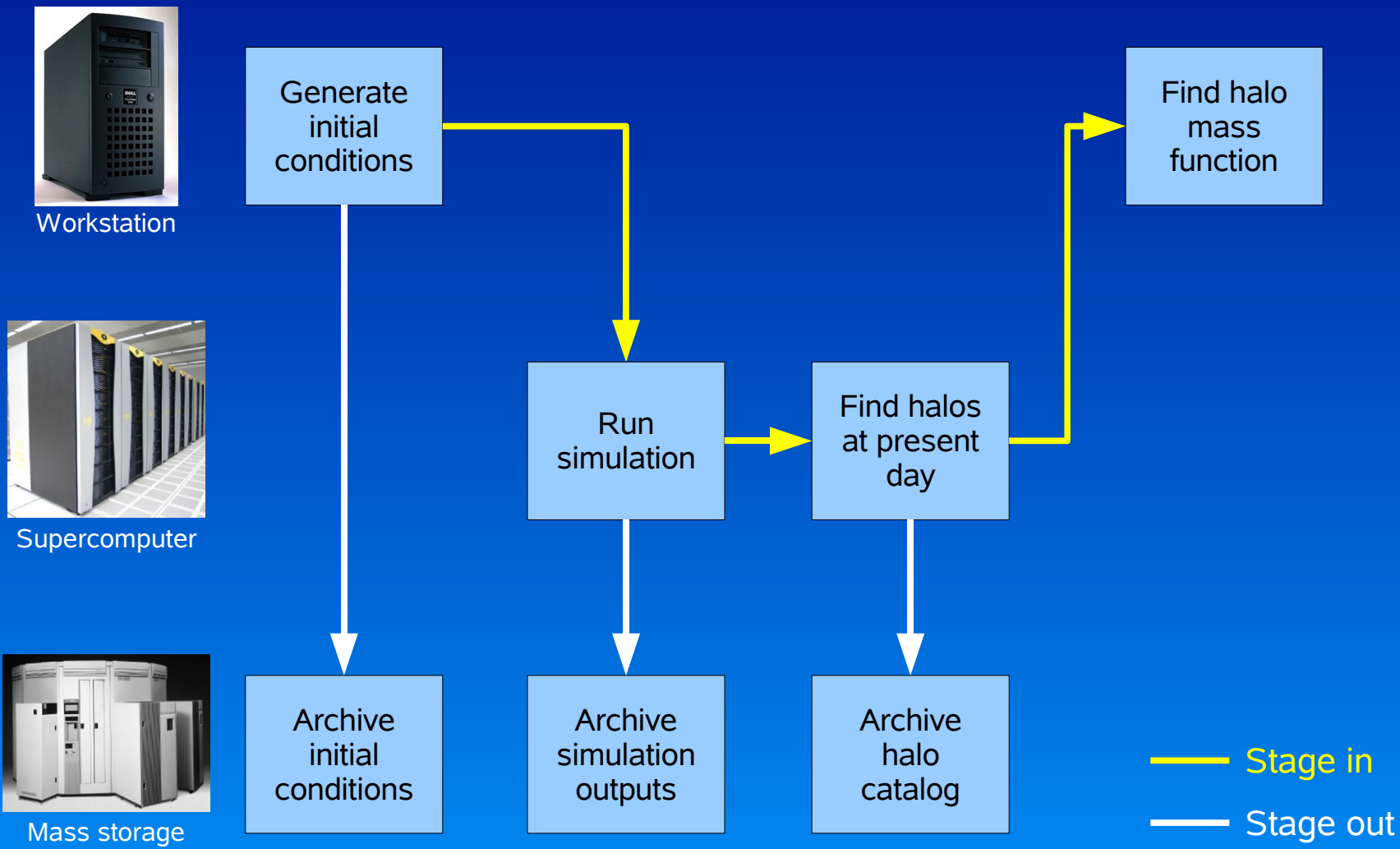
An example workflow

- Galaxy cluster formation
 - Gravity amplifies random initial perturbations
 - Very large dynamic range
- Resimulation approach
 - Large, low-resolution box to locate clusters
 - Resample initial volume containing matter that ends up in cluster
 - Resimulation this volume with surroundings at low resolution to provide boundary conditions





Example workflow #1 – simple





Example workflow #1 – Teuthis implementation

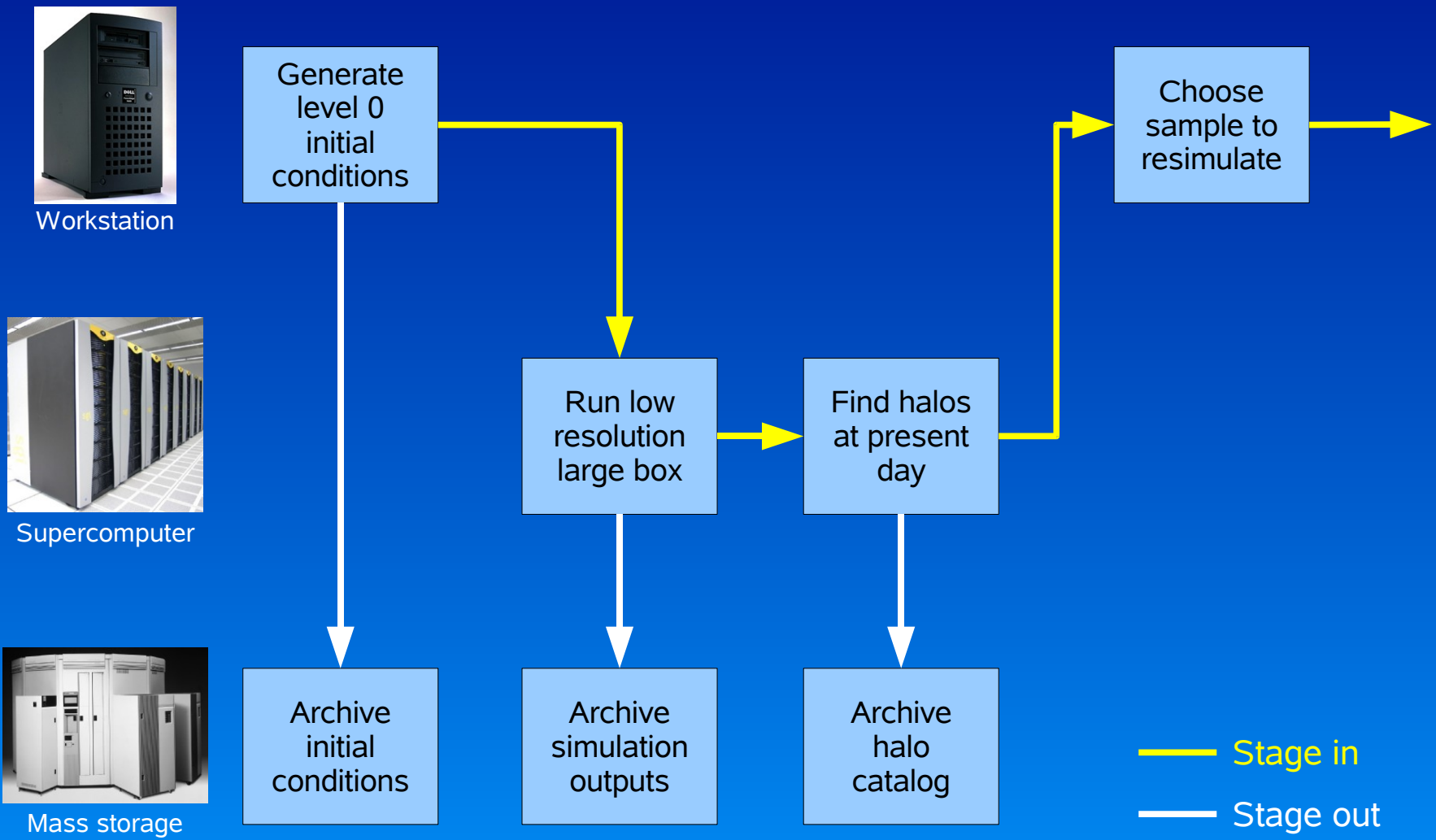
The screenshot shows the Teuthis 2.0 interface with two main panels: Projects and Workflows.

Projects Panel:

- Cosmological simulation demo**
 - Power spectrum generation
 - Run WMAP3 model
 - Job WMA0002: 16399 [02:41 11/14/2006] 1 CPU/00:00 (Complete) No data
 - Run Low Lambda
 - Run High Lambda
 - Initial conditions generation
 - Run A
 - Job A0002: 000000 [02:41 11/14/2006] 1 CPU/00:00 (Complete) No data
 - Run B
 - Job B0001: 16431 [02:42 11/14/2006] 1 CPU/00:00 (Complete) No data
 - Run C
 - Run D
 - Simulation
 - Run WMAP3
 - Job WMA0002: 16454 [02:42 11/14/2006] 1 CPU/00:00 (Complete) No data
 - Job WMA0003: 000000 [Not yet submitted] 1 CPU/00:00 (Unsubmitted) No data
 - Run Low Lambda
 - Run High Lambda
 - Halo catalog construction
 - Run WMAP3 - realization A
 - Run WMAP3 - realization B
 - Run WMAP3 - realization C
 - Run WMAP3 - realization D
 - Plot halos
 - Run WMAP3 - realization A
 - Run WMAP3 - realization B
 - Run WMAP3 - realization C
 - Run WMAP3 - realization D
 - Data movement
 - Run Copy power spectrum data for
 - Run Copy power spectrum data for
 - Run Copy power spectrum data for
 - Run Copy initial conditions for realiz
 - Run Copy initial conditions for realiz
 - Run Copy initial conditions for realiz



Example workflow #2 - resimulation





Example workflow #2 - resimulation

Do together (one for each cluster in sample)



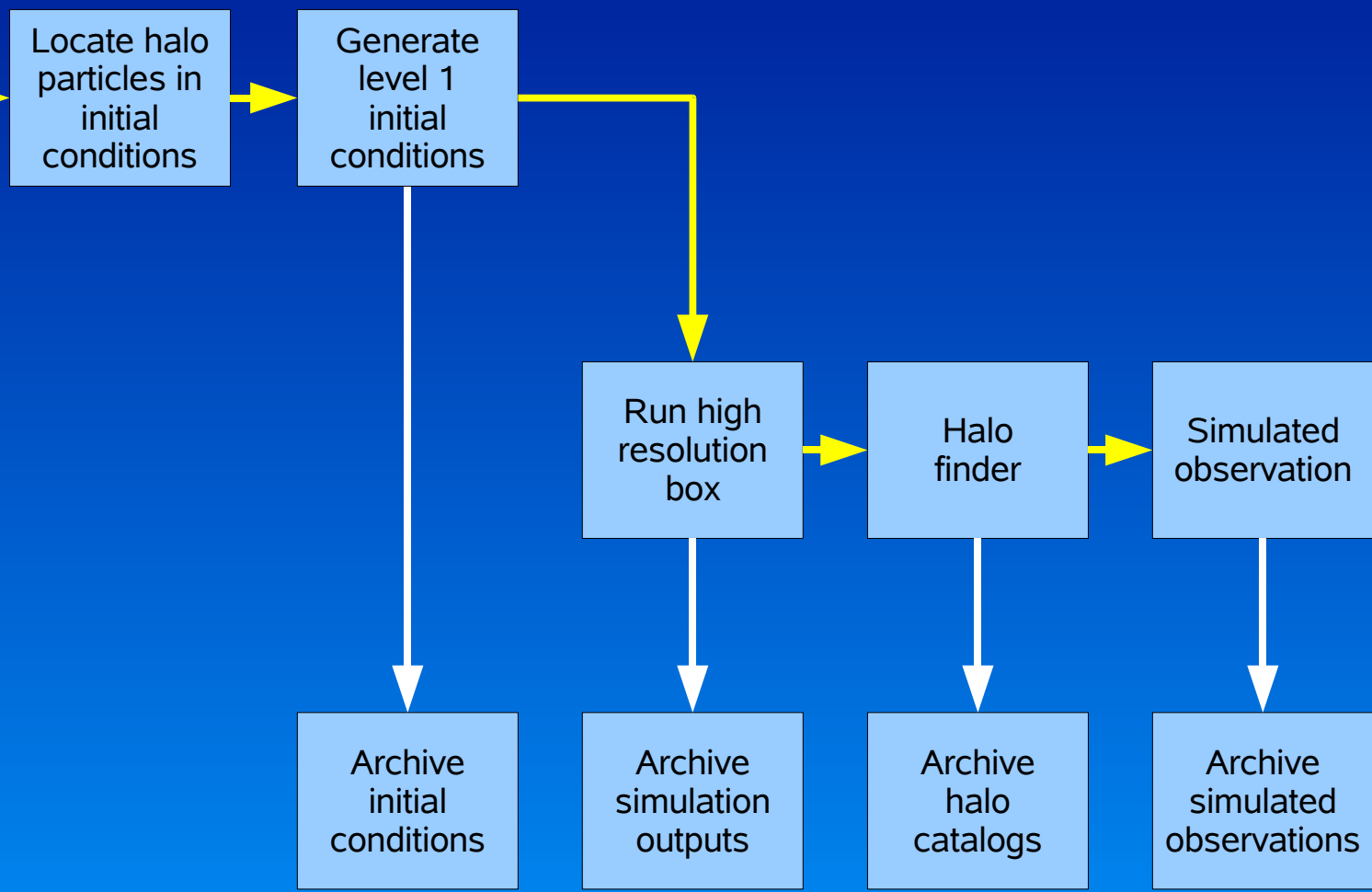
Workstation



Supercomputer



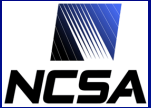
Mass storage





Future plans

- **More sophisticated experiment designs**
 - Latin hypercube, random samples
- **More sophisticated workflows**
 - DoTogether with multiple execution threads
 - DoTogether with linked threads (e.g. data mover)
- **Integration with other workflow engines**
 - Export workflow script and submit
- **More complete data tracking**
 - Store lists of files transferred with types, checksums, sizes
 - Should be able to store other data besides logs/stdout
- **Integration with observational data management tools**
 - Portal version



Getting Teuthis

1.0 release available at

<http://mazama.ncsa.uiuc.edu/projects/teuthis>

TEUTHIS

About

- Documentation
- Download
- Support
- Presentations
- Related projects
- CI home
- Internal pages

Welcome to Teuthis!

Teuthis is a tool intended to improve the efficiency with which computational scientists make use of computing resources, particularly high-performance computers. It is designed especially for the needs of astrophysical simulations, but any computational task that takes a set of input parameters from a file and runs noninteractively can be managed using Teuthis.

With Teuthis you can:

- Remotely configure and build applications from local source code
- Submit and track jobs on remote computing resources
- Painlessly schedule and track multiple restart jobs
- Stage and archive data on different machines
- Create large parameter studies with a few simple operations
- Organize job metadata by purpose and disposition
- Share calculation records with collaborators

Name	Descr
FLASH testing	
Run A	
Job A0001	Orig
Sedov scaling test	
Run A1	
Job A10001	Orig
Job A10001 Copy	Orig
Run A2	
Job A20001	Orig
Run A4	