

DESY Features on Top of HTCondor

The talk provides some details of special DESY configurations. It focuses on features we need for user registry integration, node maintenance operations and fair share / quota handling.



With the help of job transformations defining job classes and proper job duration and memory setting, we setup a smooth and transparent operating model.

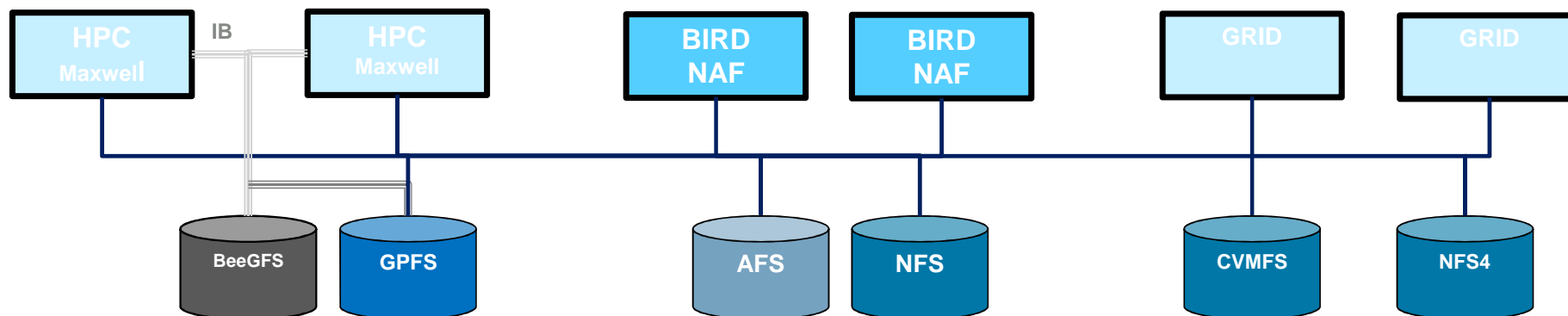
DESY/IT-Systems:
Thomas Finnern
Martin Flemming
Christoph Beyer
Yves Kemp



Overview DESY Batch Infrastructure !



High Performance Computing	BIRD Computing	Test Cluster	Grid Computing	CE
				ARC CE



Outline of Talk

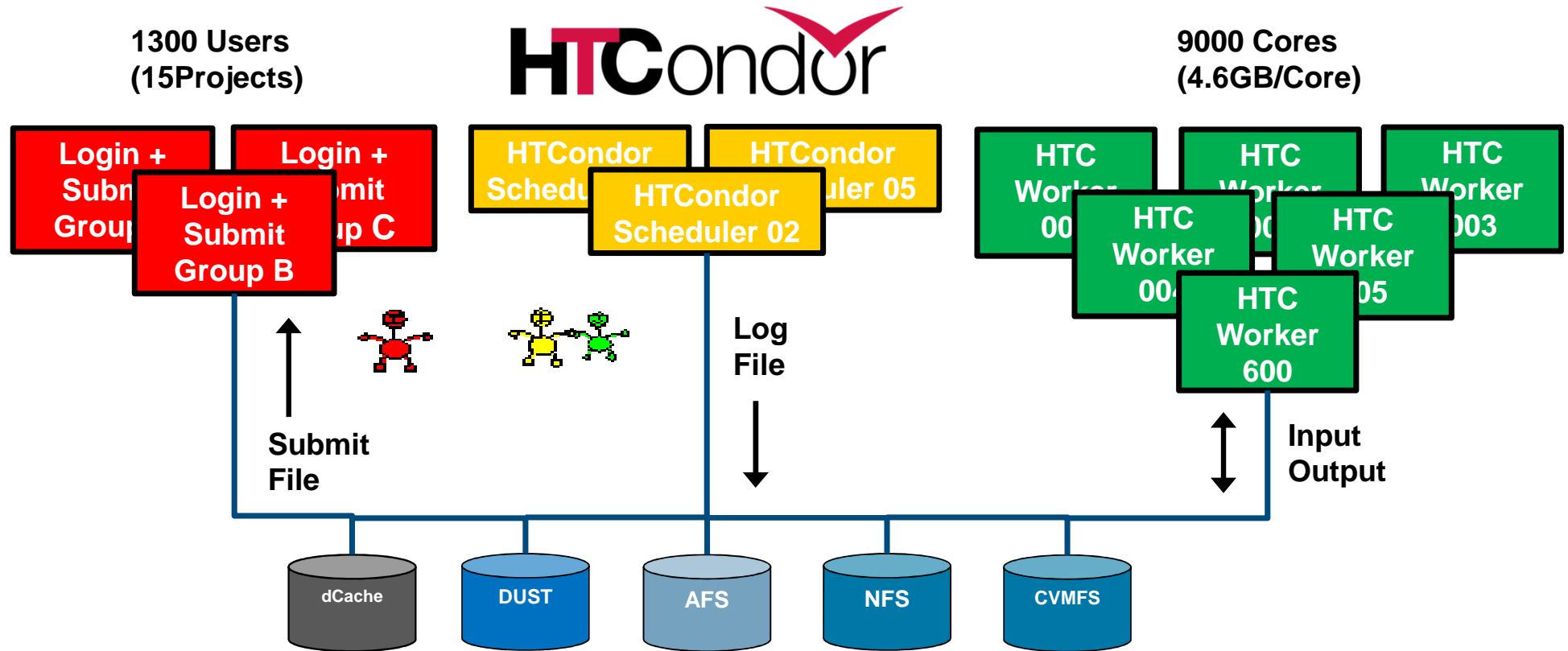
Seamless Integration of HTCondor into the DESY Environment

- Main Focus on BIRD Facility
 - BIRD/NAF Overview
- The Base: Job Classes
- Implementing Dynamic Fair Share
- Node Automation and Control
- User Registry Integration
 - Creating User.Map and Share Groups
 - Adding Blacklists and Maintenance
- Outlook and Conclusions



BIRD, NAF, HTC and HPC:
Batch Infrastructure Resource at DESY
National Analysis Facility
High Throughput Computing
High Performance Computing

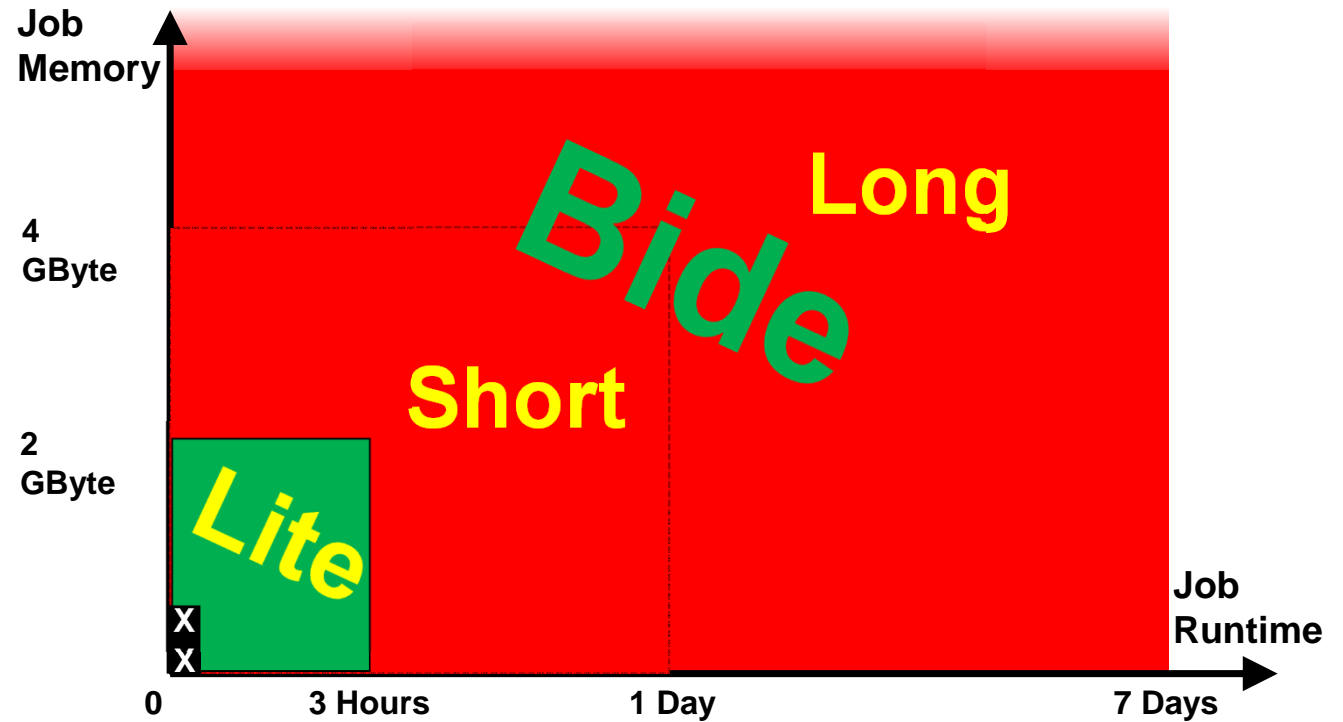
BIRD/NAF Simple Block View



The Base: Job Classes



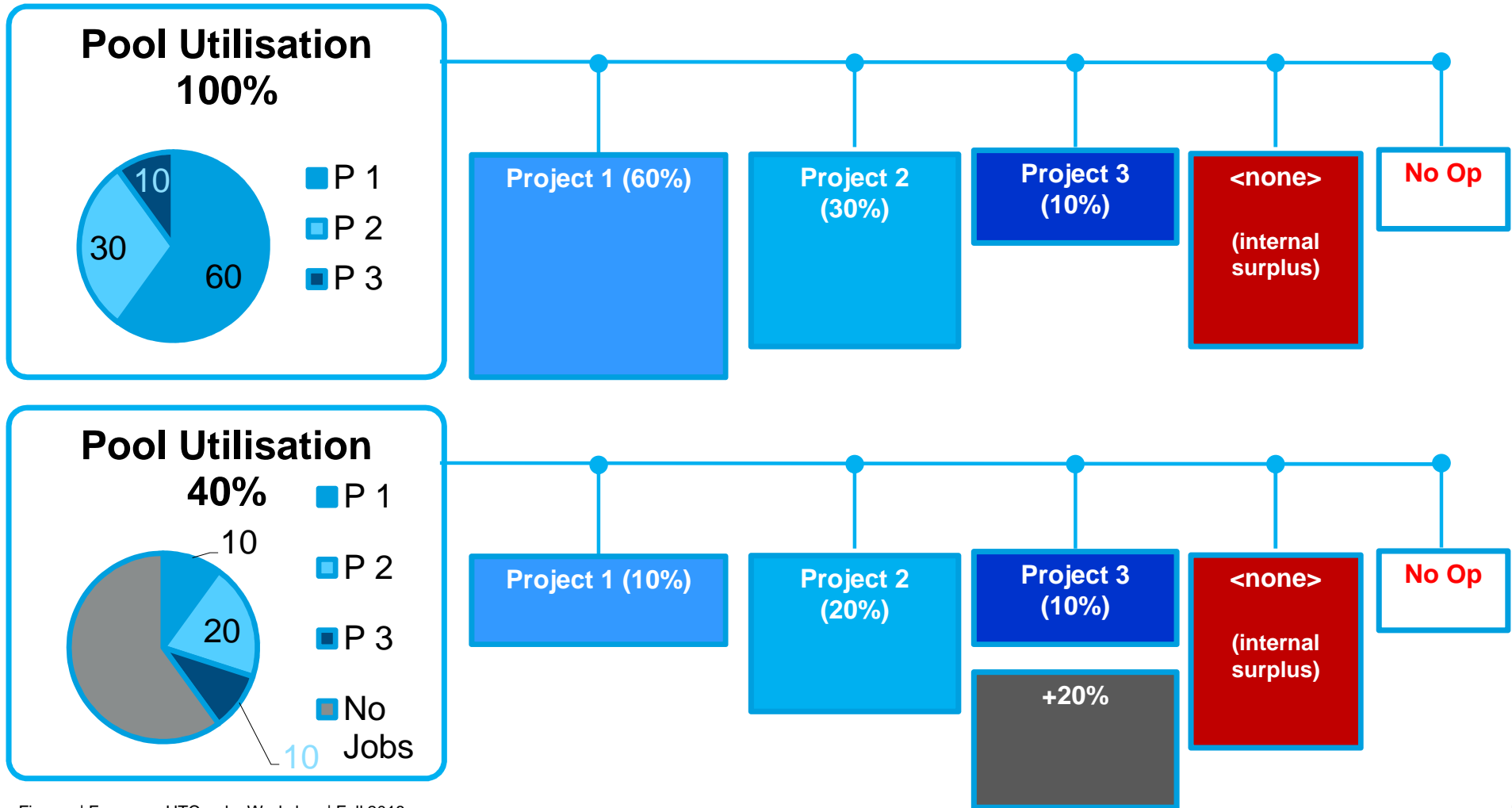
- Job Classes
 - Lite, Short and Long
 - For Quotas and Shares
 - Lite and Bide
- Job Types
 - Single, Array, Multicore, Multiarray
 - For Informational Purpose
- XX No-Go-Area
 - Minimum Memory Setting
 - Minimum Runtime Out of Control



Implementing Dynamic Fair Share (1)



Two different utilisations

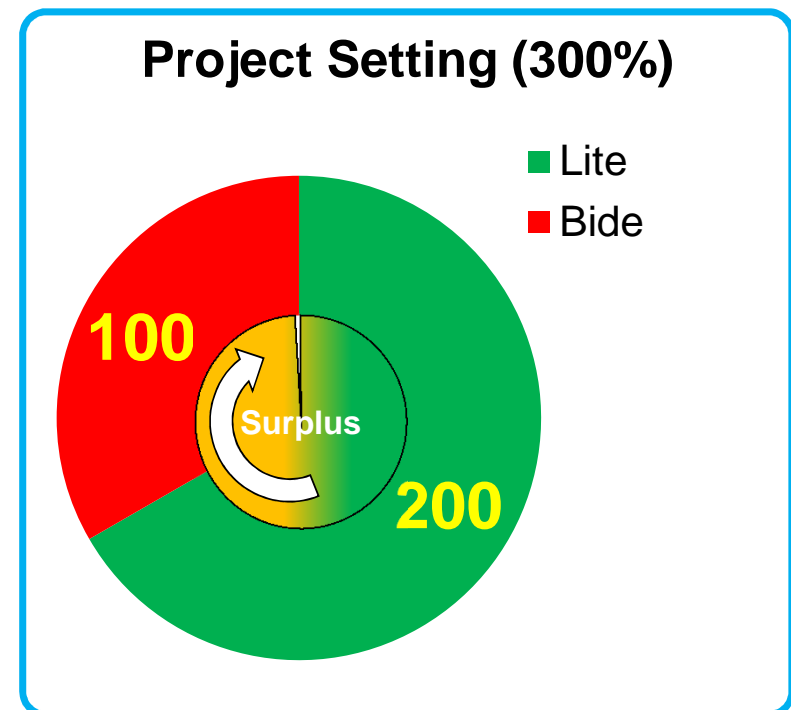


Implementing Dynamic Fair Share (2)

Static hierarchical Quotas with partial Surplus



- Job Classes and Accounting Groups
 - `NEGOTIATOR_ALLOW_QUOTA_OVERSUBSCRIPTION = True`
 - `ProjectX = NCores * Share(ProjectX) * 300%`
 - `ProjectX.lite = NCores * Share(ProjectX) * 200%`
 - `ProjectX.bide = NCores * Share(ProjectX) * 100%`
 - `GROUP_ACCEPT_SURPLUS_ProjectX.lite = true`
 - `GROUP_ACCEPT_SURPLUS = False`
 - `GROUP_AUTOREGGROUP = False`
- Hierarchical Quota with Overcommitment
 - Fairshare over time is proportional to Group Quota Ratio
 - Partial surplus for lite jobs for fast node fill
- 300 % Overcommitment
 - Allows flooding the Cluster with lite Jobs
 - Forbids flooding the Cluster with bide Jobs

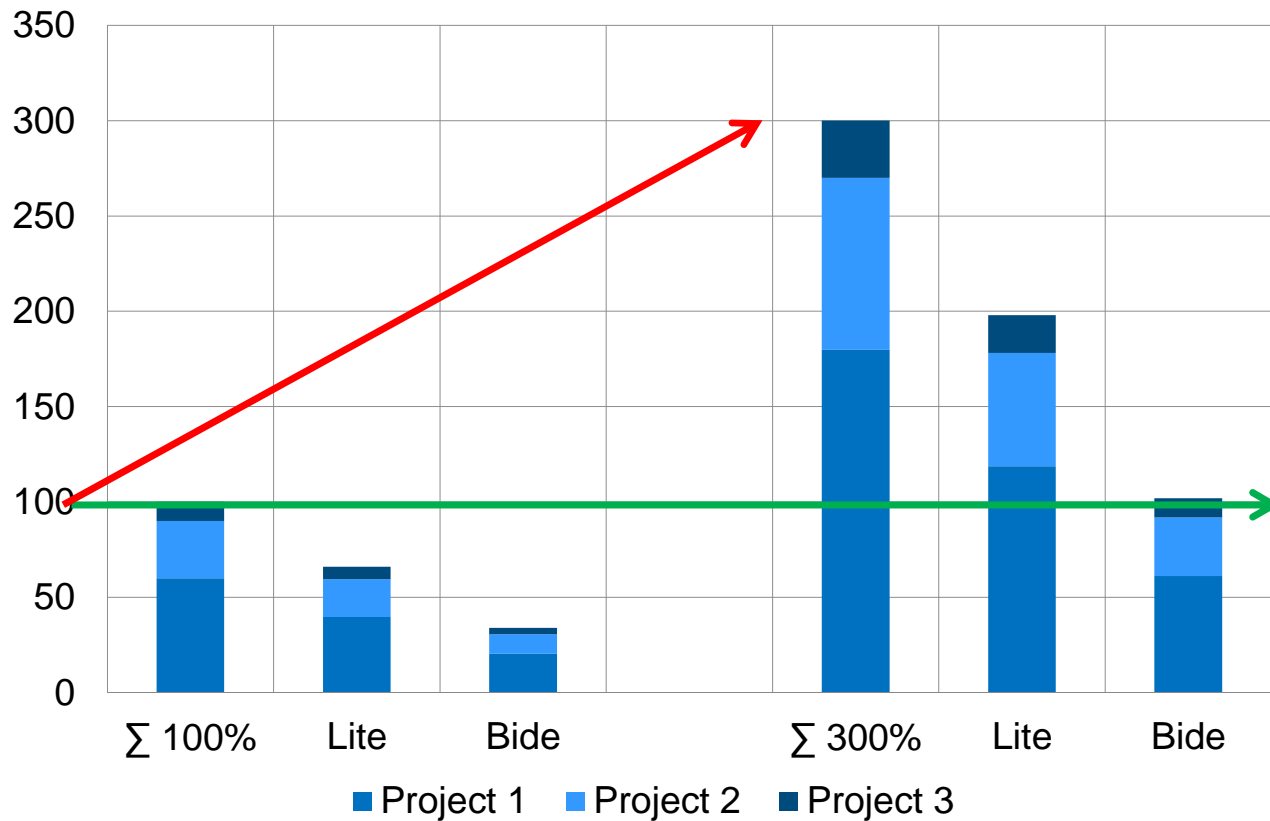


Quota and Fairshare Handling (3)

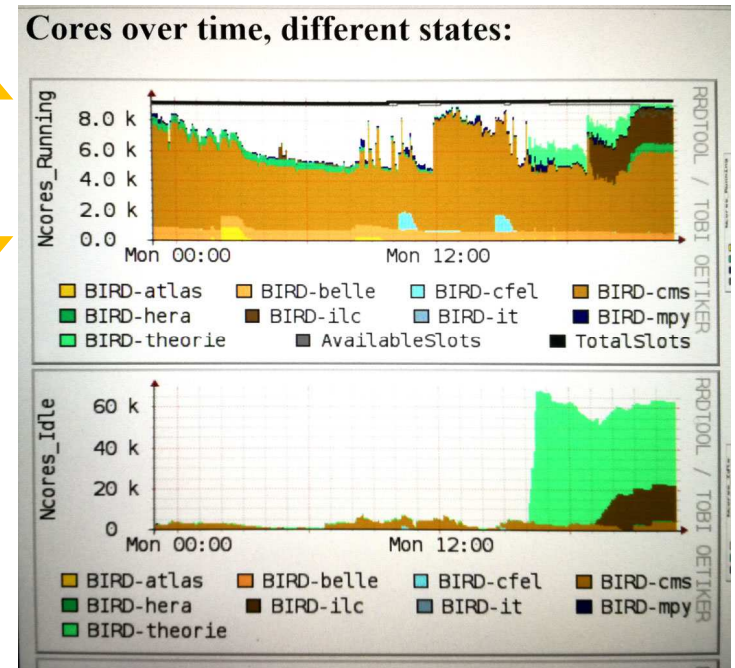
300 % Overcommitment



Quota and Fair Share (X=3)



Head Room



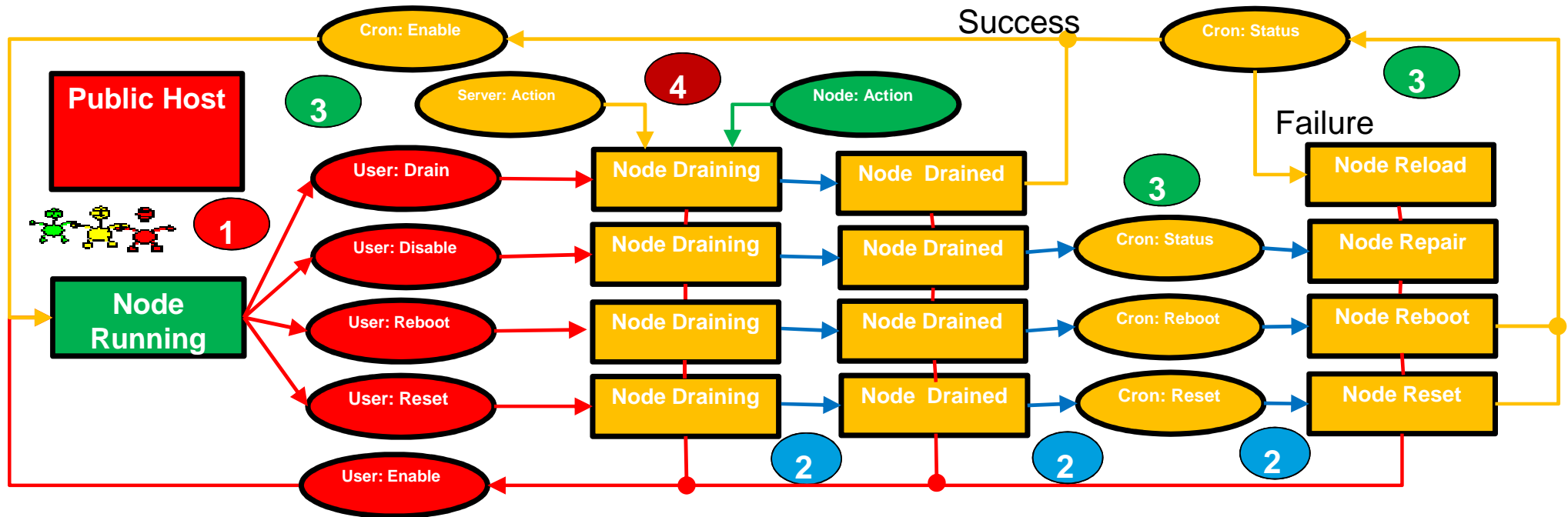
Node Automation and Control



- Automated Operation of Nodes
 - For Problems (e.g. Node Failures)
 - For Service (e.g. Cluster Kernel Update)
 - Manually/CLI or by Scripting
- Disable, drain, reboot and reset Nodes
 - No Preemption or Job Killing
 - No specific Operator Knowledge needed
- Authentication
 - Based on Authenticated Remote Command (`arc`)
 - User based (Operators, Admins) `Batchnode.sh`
 - Server based (Scripts, Cluster Reboot, Kernel Upgrades, ...)
 - Node based (local Monitoring, ...)
- Transparent
 - All States in one View
 - Hourly Status Update
 - Sets/resets exact Icinga Downtimes
 - Works for all Pools
 - GRID, BIRD, TEST, ...

```
[finnern@pal43 ~]$ batchnode show all
Check all
Asking Cluster of bm-test:
  Hosts: Is-Weg-Test
  Date:Time      Host                State: Reason
20170224:1205   Is-Weg-Test.desy.de  Gone: Wo bin ich ?
Asking Cluster of bird-htc-master02:
  Hosts: reference wn4-test bird777 bird781
  Date:Time      Host                State: Reason
20170620:2325   reference.desy.de    Repair: Wer bin ich ?
20170921:2321   wn4-test.desy.de     Gone: afs write problems to user output 2635.8 2635.9 26
20171009:1739   bird777.desy.de      Off/Draining: AutoReboot:bird781 HEPiX example
20171009:1739   bird781.desy.de      Off/Draining: HEPiX example
Asking Cluster of birdsrv1:
  Hosts: bird700 bird666 bird196 bird400 bird588 weg bird298 bird630 bird428 bird436 bird337
  Date:Time      Host                State: Reason
20171009:1734   bird700.desy.de      Draining: AutoReset:hanging dr jobs
20171009:1739   bird666.desy.de      Draining: HEPiX example
20170901:1538   bird196.desy.de      Repair: Maintenance
20170920:1846   bird400.desy.de      Repair: rt753219 Frage: dr jobs on bird400 and bird630
20171006:1333   bird588.desy.de      Booting: AutoReboot:kernel-Auto-Maintenance
20170712:1432   weg.desy.de          Gone: logging test
20170815:1426   bird298.desy.de      Repair: BIOS Update wegen YERR Errors
20170920:1847   bird630.desy.de      Repair: rt753219 Frage: dr jobs on bird400 and bird630
20171005:0824   bird428.desy.de      Reload: rt762558_r4todo_am_11-10
20171003:0114   bird436.desy.de      Reload: AutoReboot:kernel-Auto-Maintenance on birdsrv1.desy.de
20170930:1544   bird337.desy.de      Reload: AutoReset:hanging dr jobs
Asking Cluster of condor01:
  Hosts: batch0236 wn5-test
  Date:Time      Host                State: Reason
20170920:1624   batch0236.desy.de    Booting: AutoReboot:Reinstall EL7
```

Node Automation and Control



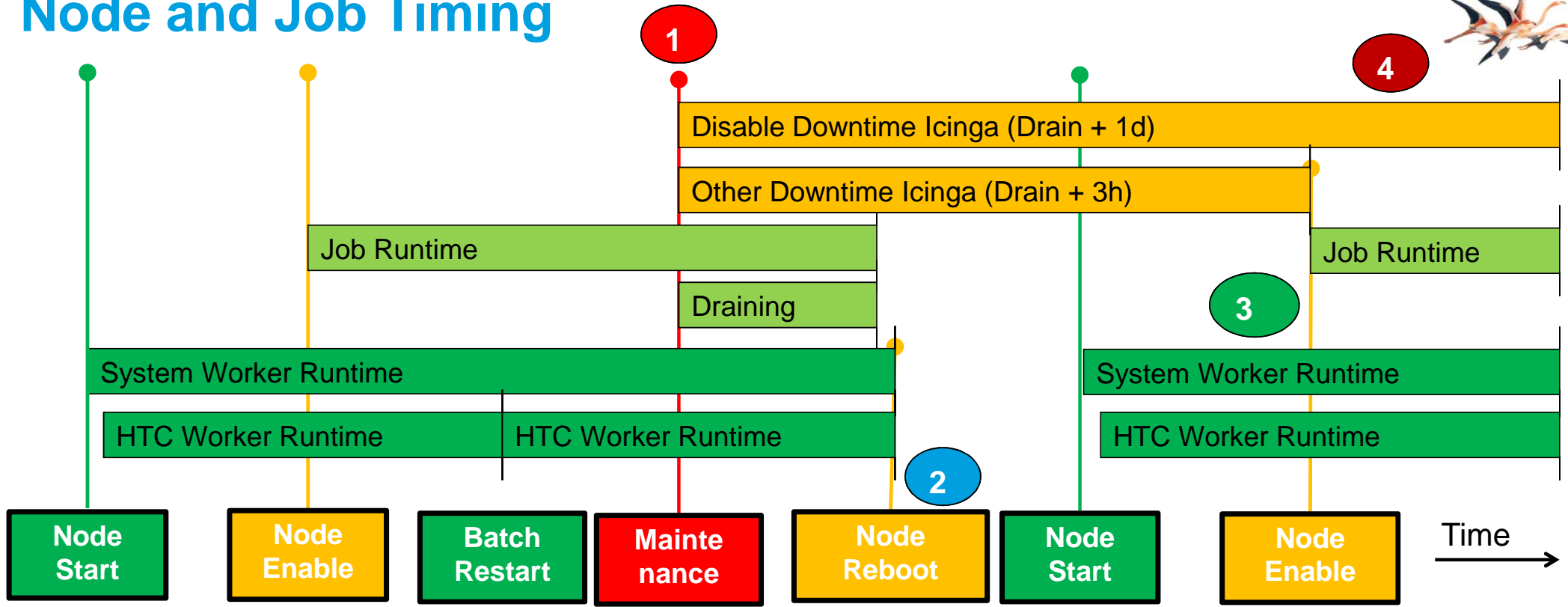
1	2	3	4
Batchnode.sh	waiting	Node.cron	Server and Node Actions
	Time	Cron.hourly	Scripts and Tools

Node and Job Timing



- Policies set on Scheduler
 - Maximum, requested and default Job Runtimes
 - Resulting Drain Times
 - BIRD and GRID are different
- Values
 - Default Runtime 3 hours
 - Runtimes requestable from 1 Week to a few Minutes
 - Vacate Time is Part of Job Runtime
 - Vacate Time 5 Minutes
- Implementation
 - HTC Feature: Periodic Remove
 - Runtime Calculation within HTC Interface
 - Node Runtime essential in Process Management

Node and Job Timing

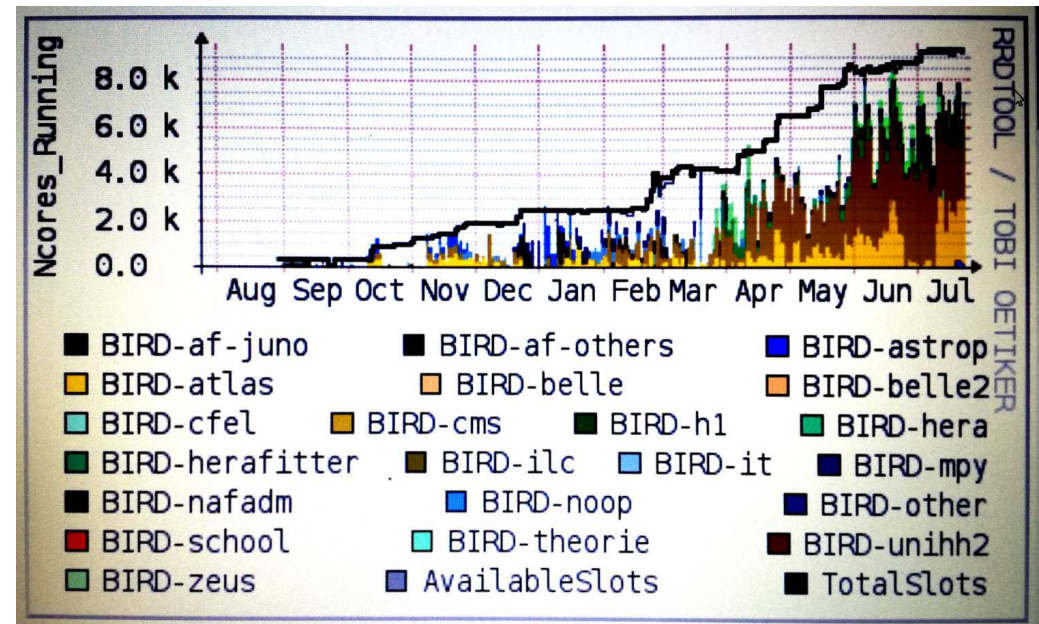


1	2	3	4
Batchnode.sh	Node.cron (Sys reboot)	Node.cron (node enable)	Node.cron (node status)
	Cron.hourly	Cron.hourly	Cron.hourly

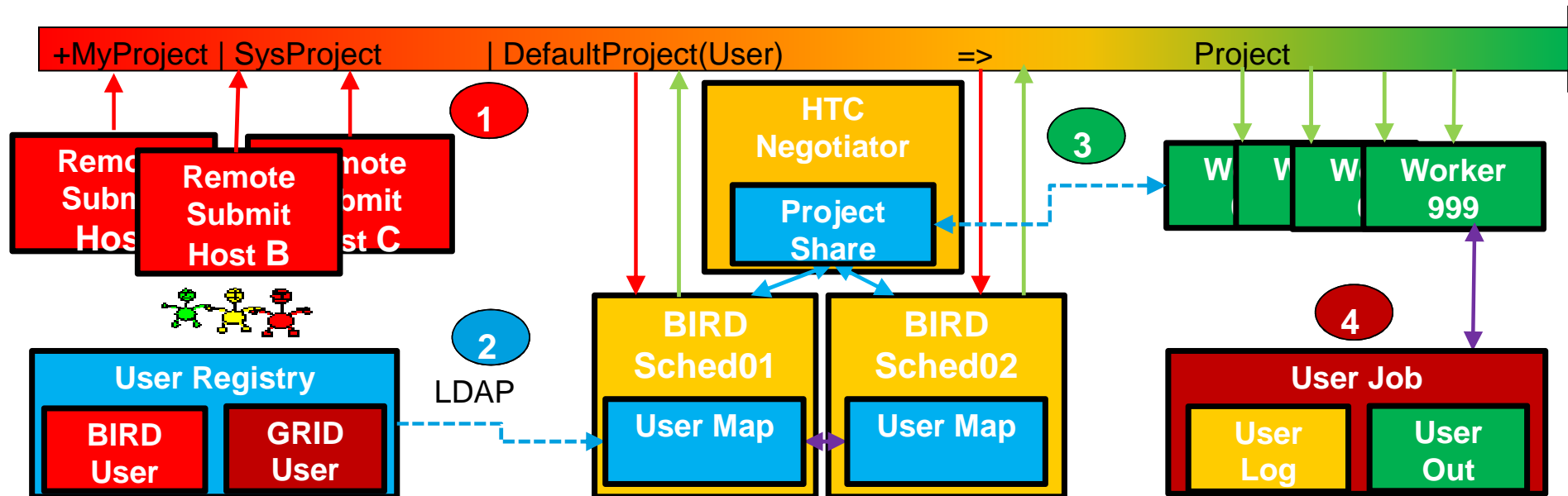
User Registry Integration



- Request adequate Project on Submit
- User can select another Project
- Project defaults to primary Registry Group
- Resulting Project will be checked against Registry
- Resulting project will define Fair Share
- Resulting Project is primary Project on Worker
- Jobs with invalid Project do not run
- Reuse Database für User Blacklisting
- Reuse Database for Maintenance Control

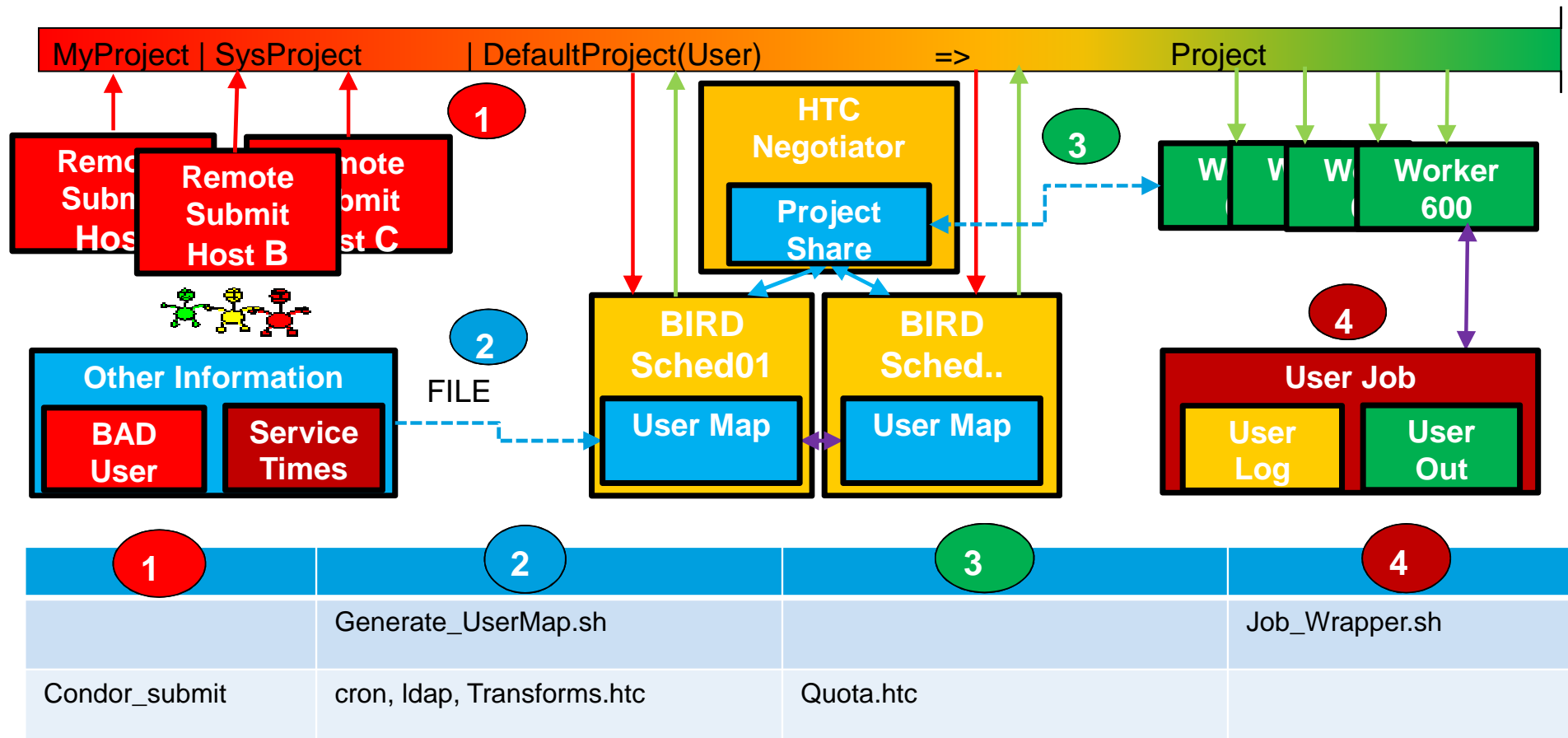


User Registry Integration



1	2	3	4
	Generate_UserMap.sh		Job_Wrapper.sh
Condor_submit	cron, ldap, Transforms.htc	Quota.htc	

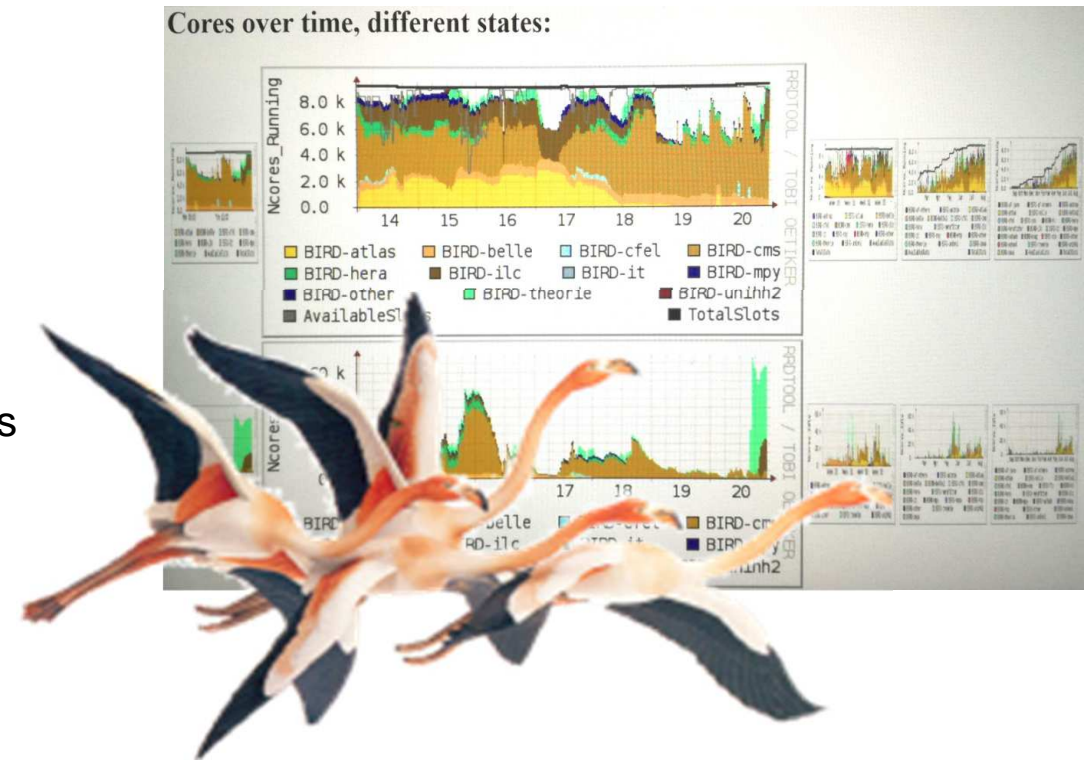
User Blacklisting and Maintenance Control



Outlook and Conclusions

- BIRD/NAF
 - „Proof of Concept“ for new features done
- GRID and BIRD/NAF
 - Some common operating tools running
- Next Steps may be ...
 - More BIRD and GRID Integration
 - Singularity for different operating system flavours
 - Get rid of SL6
 - Backfill of HPC resources with HTCondor
- Working on ...
 - Smarter Config
- Waiting for ...
 - Full AFS/Kerberos Integration

Cores over time, different states:



Thank you for listening



Questions ?

Answers !

