

Migration from Grid Engine to HTCondor

The talk provides details about the goals and general setup of this migration. It further focuses on Kerberos support, registry integration and node operating automation.

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



大学共同利用機関法人
高エネルギー加速器研究機構



Outline of Talk



- Plan and Status
- Approach and Policies
- Kerberos and AFS Integration
- User Registry Integration
- Node Automation and Control
- Outlook and Conclusions

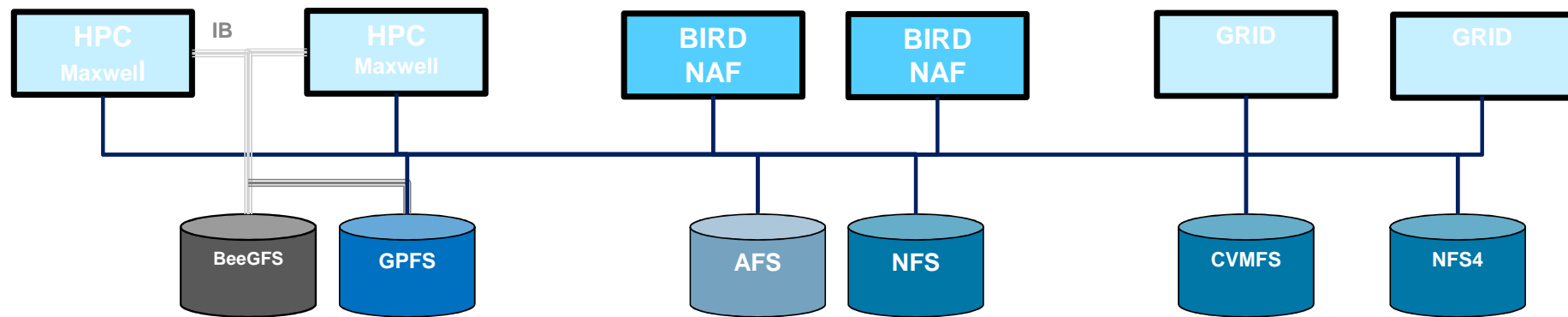
- Main Focus: BIRD

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



We have a Plan !

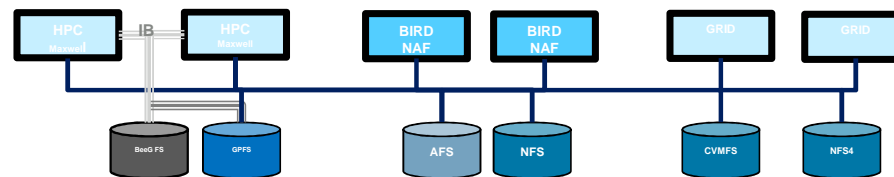
Past	Calendar Reservation Tool	Son Of Grid Engine SoGE		Torque + MySched	Cream CE
Work in Progress		SoGE	HTC Pilot	HTC Grid	
Running		ARC CE			
Future					



Approach and Policies for HTCondor



- User Friendly Migration
 - Kerberos and AFS Integration
 - Growing BIRD Pilot with 500+ Cores
- 1 “Master” Server supporting BIRD
 - Collector and Negotiator
 - Quota/Fairshare configuration
- 2 BIRD-Scheduler
 - For job-friendly service restart
 - With secure token access
 - Policy settings on scheduler
 - Check projects against registry
 - Forward submitter project to worker
- 10 Remote Submit Hosts
 - User Login for job preparation
 - Project specific default settings
 - No dependencies from/to running jobs
- 500 Worker Nodes
 - Common node setup
 - Currently 30+ in Pilot
- Optional GRID Integration
 - Add GRID share to negotiator
 - Add GRID scheduler
 - Add GRID worker resource



Kerberos and AFS Integration

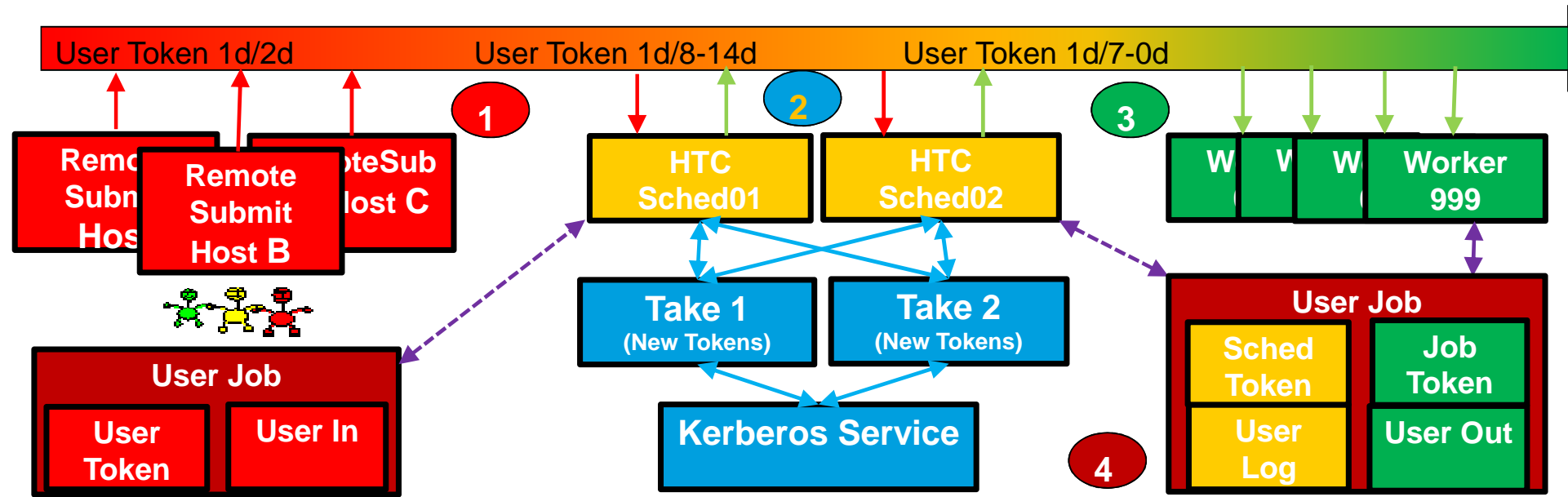


- Kerberos authentication
- AFS as shared filesystem
- Valid tokens during job run time
- 1 week maximal job run time
- Secure token generator on protected servers
- Consistently prolong current token
- Generation of AFS tokens out of Kerberos tokens

Colors

Red: Interactive User
Orange: HTC Server
Green: WorkerNodes
Blue: Authentication

Kerberos and AFS Integration

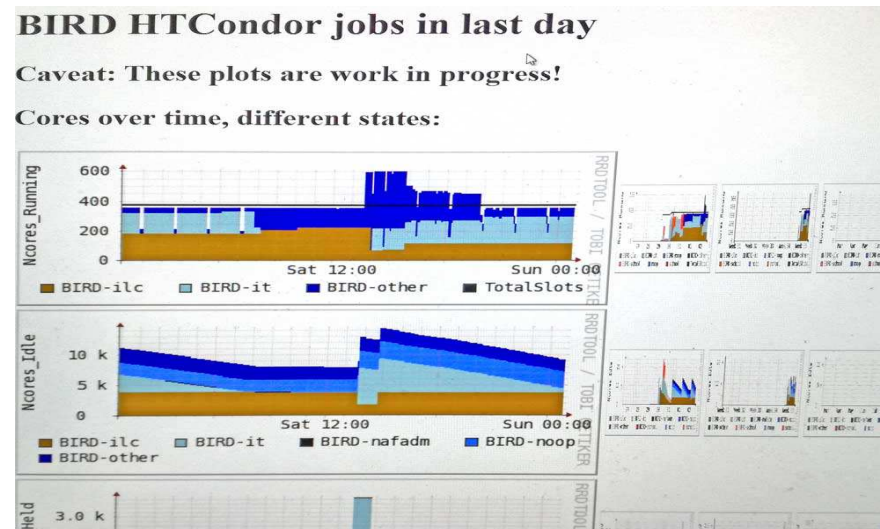


1	2	3	4
Get-Token.sh	Token_Shepherd_Sched.sh	Token_Shepherd_Worker.sh	(Job_Wrapper.sh)
kinit, aklog	arc(Take), kinit, condor_aklog	kinit-prolong, condor_aklog	Token_Update.htc

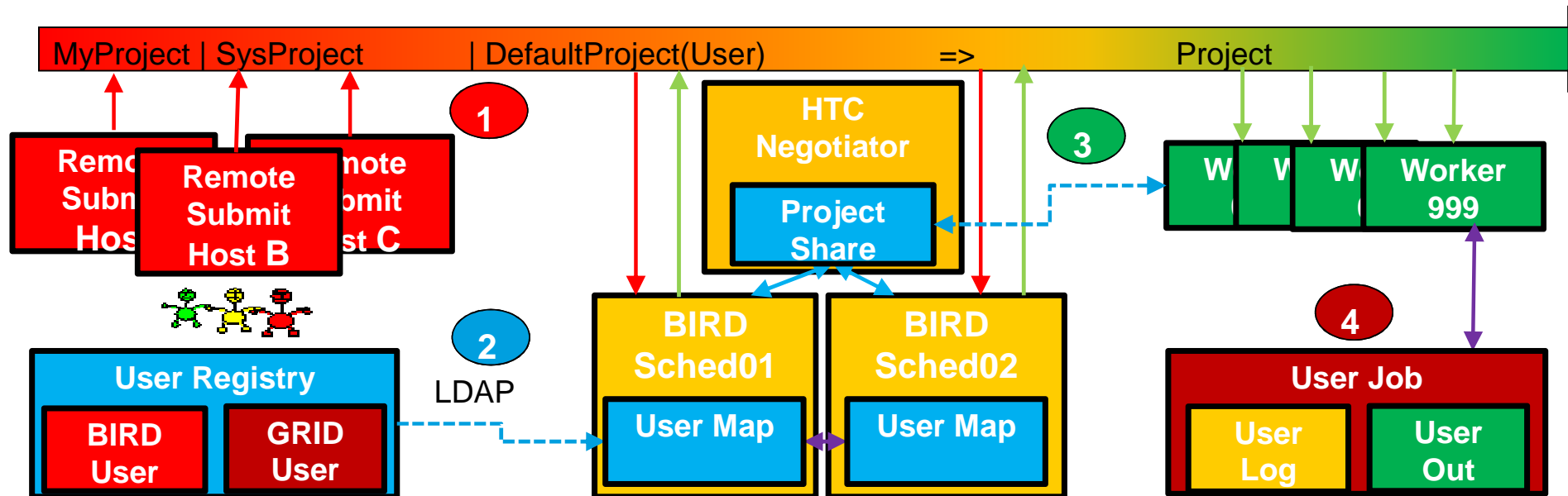
User Registry Integration



- Set adequate project on group submit host
- User may switch to another project
- Project defaults to primary registry group
- Resulting project will be checked against registry
- Resulting project will define fairshare/quota group
- Resulting project will be set on worker as primary group



User Registry Integration



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

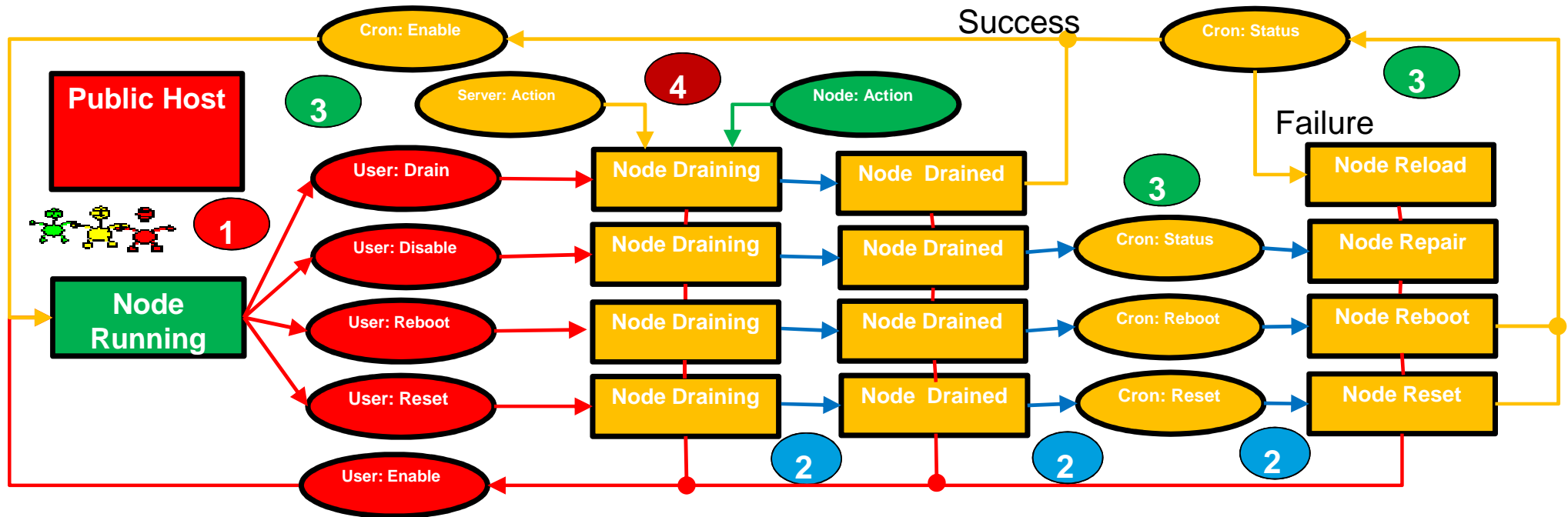


Node Automation and Control

- 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 **A**uthenticated **R**emote **C**ommand (**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
 - SoGE, GRID, PILOT, 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: AutoResetHanging 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: AutoResetHanging 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

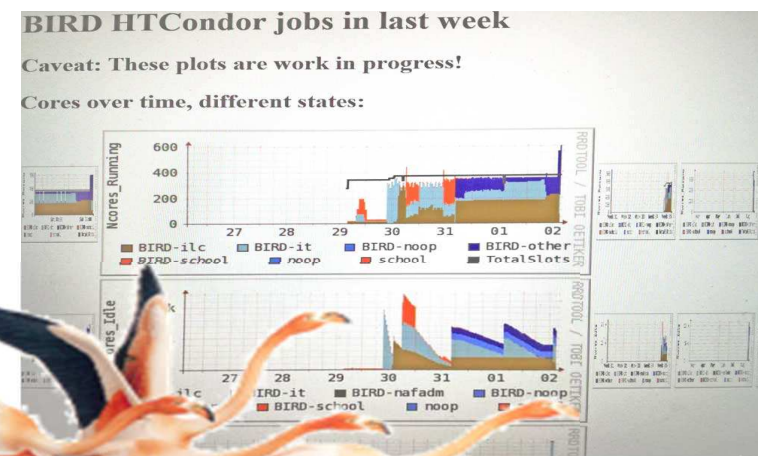


1	2	3	4
Batchnode.sh	waiting	Node.cron	Server and Node Actions
	Time	Cron.hourly	Scripts and Tools

Outlook and Conclusions



- BIRD/NAF
 - „Proof of Concept“ for planned feature done
 - Waited for HTCondor 8.7.3 providing
 - Full Kerberos and AFS support
 - Transforms for Scheduler Policy Settings
 - Pilot running with first users
 - Smooth Transition appears to be possible
- GRID and BIRD/NAF
 - Some common operating tools running
 - Pilot HTC Server are prepared
- Next Steps may be ...
 - More BIRD and GRID Integration
 - Docker for different operating system flavours
 - Backfill of HPC resources with HTCondor



Questions ?

Answers !

