

Towards automation of computing fabrics using tools from the fabric management workpackage of the EU DataGrid project

Olof Barring
(WP4)

Olof.Barring@cern.ch



<https://edms.cern.ch/document/376367/1>



Talk Outline

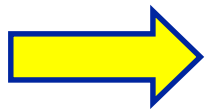
- ◆ Introduction to EU DataGrid workpackage 4
- ◆ Automated management of large clusters
- ◆ Components design and development status
- ◆ Experience with LCFG for system configuration and installation
- ◆ Summary and outlook

Authors

Olof Barring, Maite Barroso Lopez, German Cancio, Sylvain Chapeland, Lionel Cons, Piotr Poznański, Philippe Defert, Jan Iven, Thorsten Kleinwort, Bernd Panzer-Steindel, Jaroslaw Polok, Catherine Rafflin, Alan Silverman, Tim Smith, Jan Van Eldik - CERN
Massimo Biasotto, Andrea Chierici, Luca Dellagnello, Gaetano Maron, Michele Michelotto, Cristine Aiftimiei, Marco Serra, Enrico Ferro – INFN
Thomas Röblitz, Florian Schintke – ZIB
Lord Hess, Volker Lindenstruth, Frank Pister, Timm Morten Steinbeck – EVG UNI HEI
David Groep, Martijn Steenbakkers – NIKHEF/FOM
Paul Anderson, Tim Colles, Alexander Holt, Alastair Scobie, Michael George - PPARC

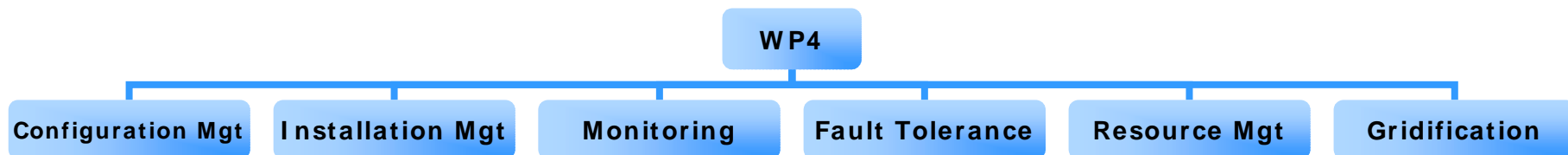
WP4 objective and partners

“To deliver a computing fabric comprised of all the necessary tools to manage a center providing grid services on clusters of thousands of nodes.”

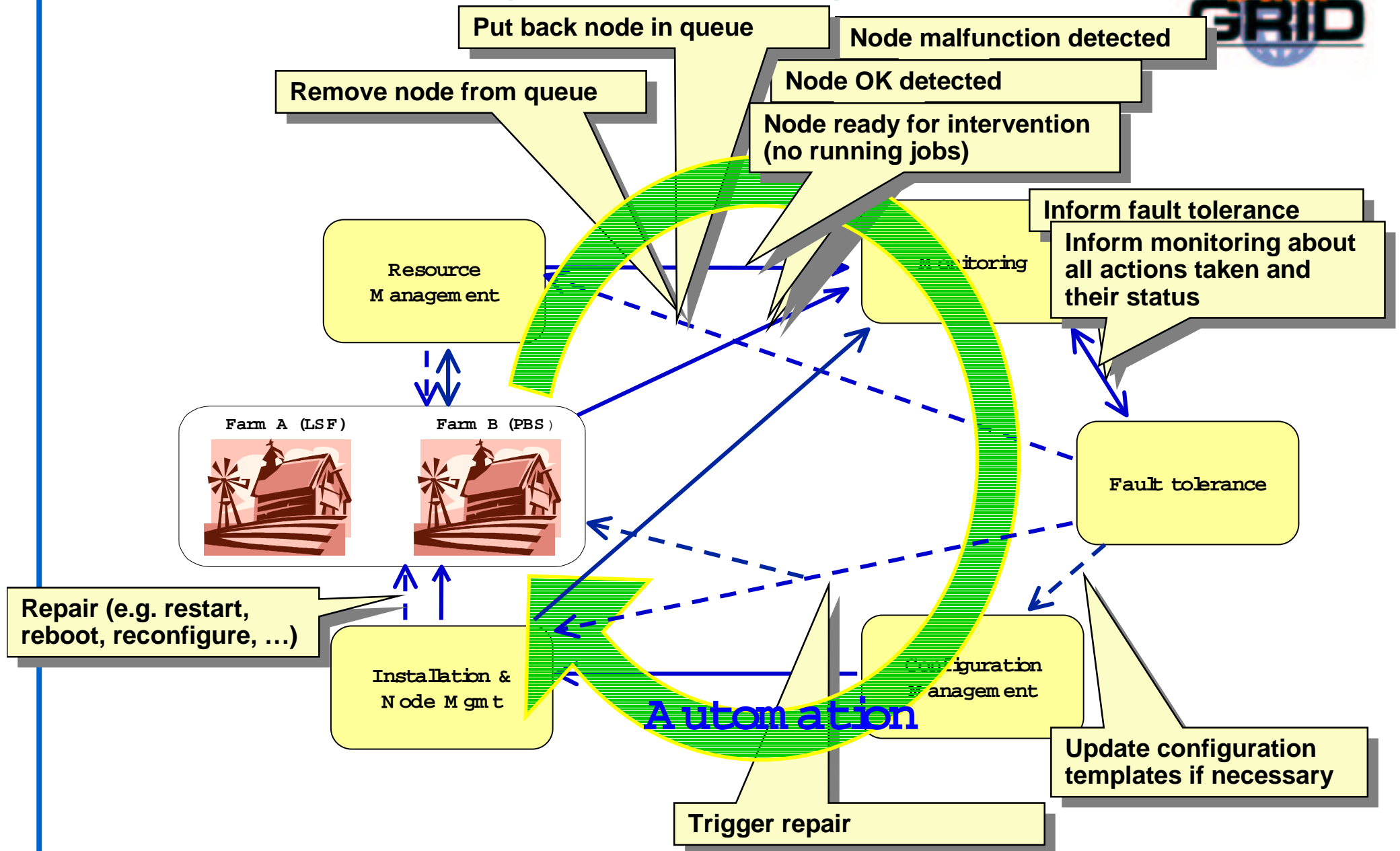


- User job management (Grid and local)
- Automated management of large clusters

- ◆ 6 partners: CERN, NIKHEF, ZIB, KIP, PPARC, INFN.
- ◆ ~ 14 FTEs (6 funded by the EU).
- ◆ The development work divided into 6 subtasks:



Automated management of large clusters



Monitoring subsystem: design

Monitoring Sensor Agent

- Calls plug-in sensors to sample configured metrics
- Stores all collected data in a local disk buffer
- Sends the collected data to the global repository

Plug-in sensors

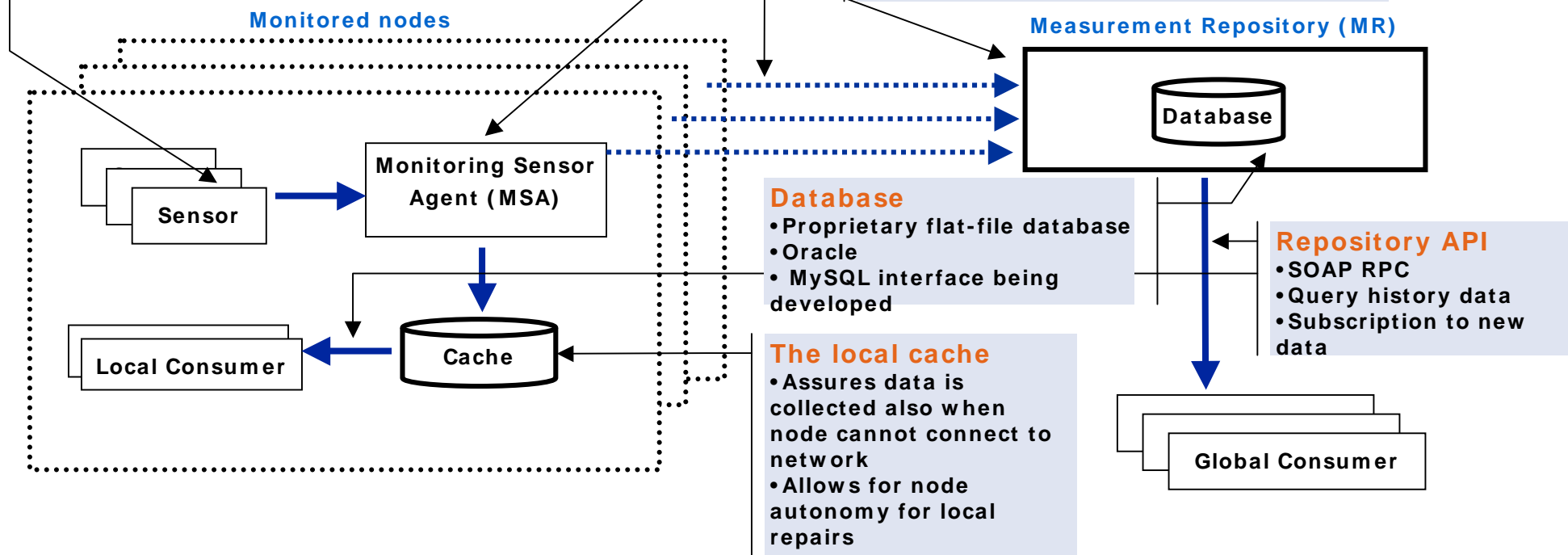
- Programs/ scripts that implements a simple sensor-agent ASCII text protocol
- A C++ interface class is provided on top of the text protocol to facilitate implementation of new sensors

Transport

- Transport is pluggable.
- Two proprietary protocols over UDP and TCP are currently supported where only the latter can guarantee the delivery

Measurement Repository

- The data is stored in a database
- A memory cache guarantees fast access to most recent data, which is normally what is used for fault tolerance correlations.





Monitoring subsystem: status

◆ Local nodes:

- Monitoring Sensor Agent (MSA) and UDP based proprietary protocol are ready and used on CERN production clusters since more than a year
- The TCP based proprietary protocol exists as prototype. More testing and functionality needed to be ready for production use

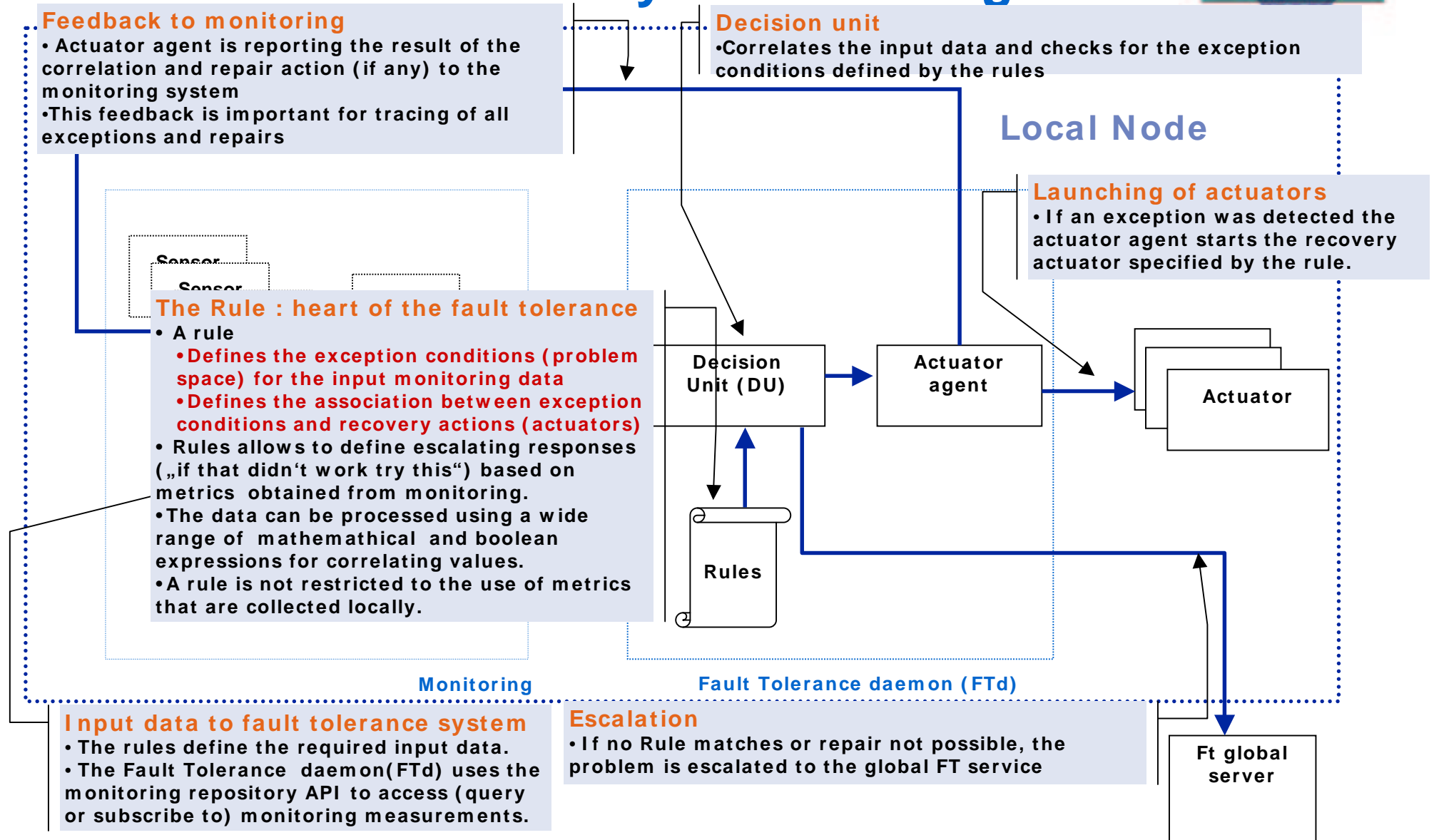
◆ Central services

- Repository server exists with both flatfiles and Oracle database. It is currently being evaluated for production use at CERN. Support for MySQL is planned for later in 2003
- Alarm display: still in early prototype phase.

◆ Repository API for local and global consumers:

- C library implementation of API (same for local and global consumers)
- Bindings for other languages can probably be generated directly from the WSDL

Fault tolerance subsystem: design





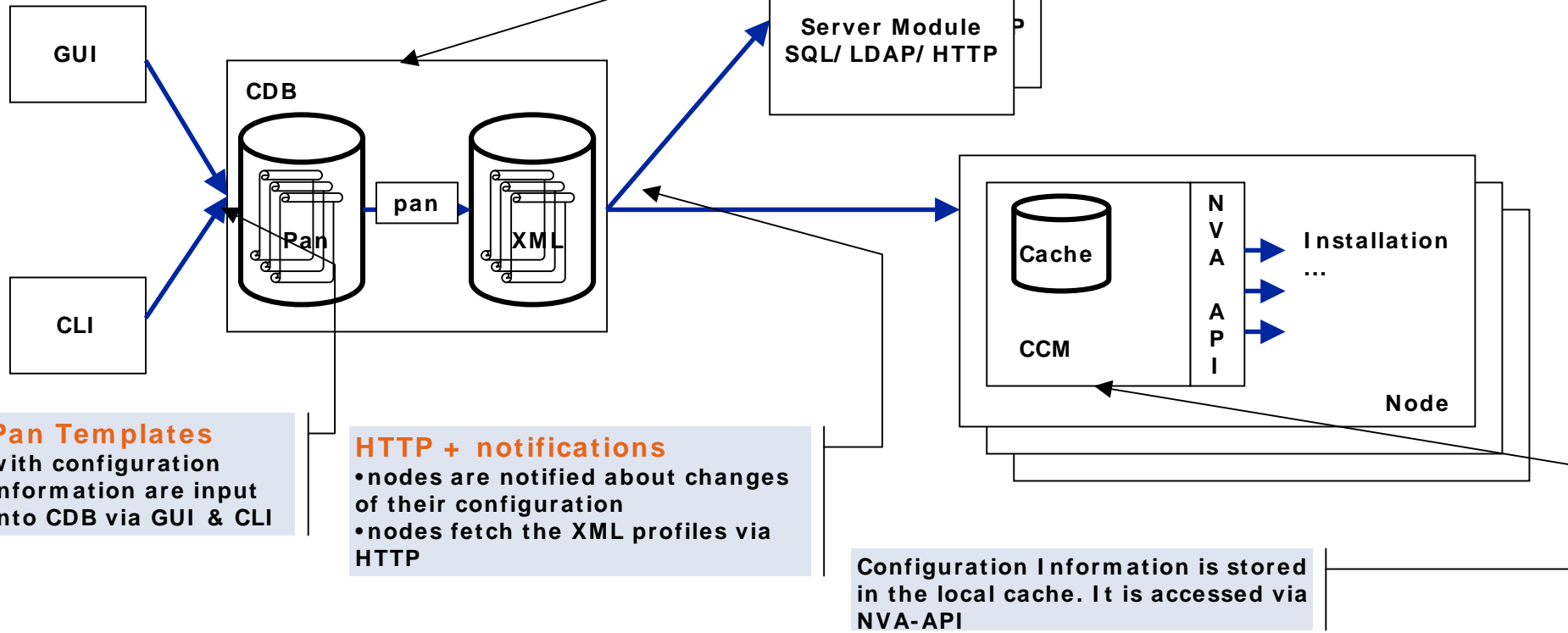
Fault tolerance subsystem: status

- ◆ Not yet ready for production deployment
- ◆ Prototype was demonstrated working together with the fabric monitoring system at EU review in February 2003
 - Web-based rule editor
 - Central Rule repository (MySQL)
 - Local FTd (fault tolerance daemon) that
 - Automatically subscribes to monitoring metrics specified by the rules
 - Launches the associated actuators when the correlation evaluates to an exception
 - Reports back to the monitoring system the recovery actions taken and their status
 - Global correlations not yet supported

Configuration Management subsystem: design

Configuration Data Base (CDB)
 Configuration Information store. The information is updated in transactions, it is validated and versioned. Pan Templates are compiled into XML profiles

Server Modules
 Provide different access patterns to Configuration Information



Pan Templates
 with configuration information are input into CDB via GUI & CLI

HTTP + notifications
 • nodes are notified about changes of their configuration
 • nodes fetch the XML profiles via HTTP

Configuration Information is stored in the local cache. It is accessed via NVA-API

Configuration Management subsystem: status



- ◆ System is implemented (except for CLI and Server Modules), most of the components in 1.0 production version,
- ◆ Pilot deployment of the complete system for LCG 1,

In parallel:

- ◆ System being consolidated,
- ◆ Issues of scalability and security being studied and addressed,
- ◆ Server Modules under development (SQL).

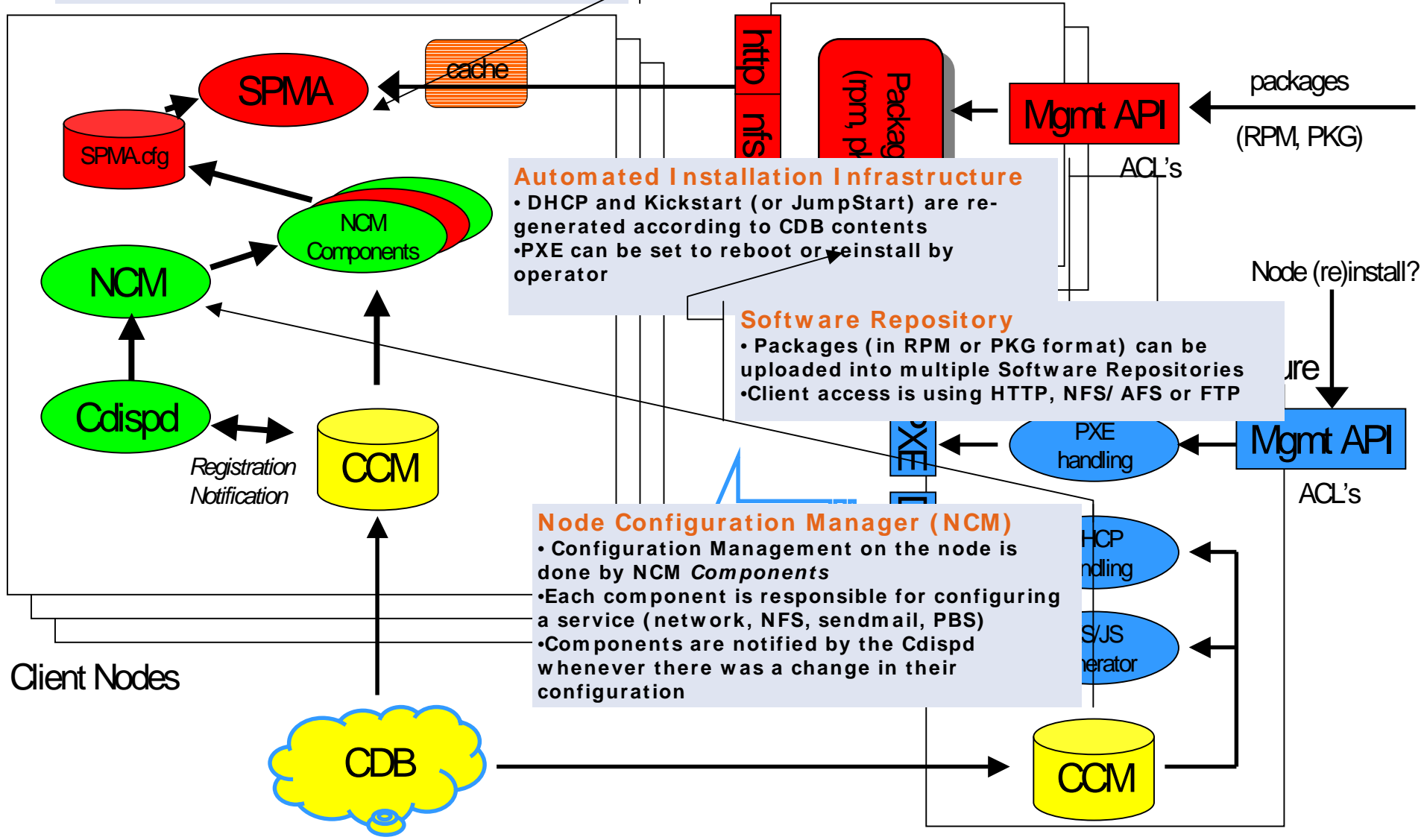
More information:

<http://cern.ch/hep-proj-grid-config/>

Implementation: design

Software Package Mgmt Agent (SPMA)

- SPMA manages the installed packages
- Runs on Linux (RPM) or Solaris (PKG)
- SPMA configuration done via an NCM component
- Can use a local cache for pre-fetching packages (simultaneous upgrades of large farms)



Automated Installation Infrastructure

- DHCP and Kickstart (or JumpStart) are re-generated according to CDB contents
- PXE can be set to reboot or reinstall by operator

Software Repository

- Packages (in RPM or PKG format) can be uploaded into multiple Software Repositories
- Client access is using HTTP, NFS/ AFS or FTP

Node Configuration Manager (NCM)

- Configuration Management on the node is done by NCM Components
- Each component is responsible for configuring a service (network, NFS, sendmail, PBS)
- Components are notified by the Cdispd whenever there was a change in their configuration



Installation subsystem: status

- ◆ Software Repository and SPMA
 - First pilot being deployed on CERN Computer Centre for the central CERN production (batch & interactive) services
- ◆ Node Configuration Manager (NCM)
 - Design phase
 - Implementation available in Q2 2003
- ◆ Automated Installation Infrastructure (AII)
 - Design phase
 - Linux Implementation expected for Q2 2003



LCFG experience

- ◆ LCFG (Local Configuration) tool from Univ. of Edinburgh has been used for fabric installation and configuration management at the EDG testbed since the first project release:
 - Tool modified to be adapted to EDG testbed needs
 - Learned a lot from it to understand what we really want
 - Used at almost all EDG testbed sites → very valuable feedback from a large O(5-10) group of site administrators
- ◆ Disadvantages with LCFG
 - Enforces a private per component configuration schema
 - High level language lacks possibilities to attach compile time validation
 - Maintains proprietary solutions where standards exist (e.g. base installation)



Summary & Future Work

- ◆ Experience and feedback with existing tools and prototypes helped to get requirements and early feedback from users
- ◆ First implementation now ready for all the subsystems
- ◆ Some of them already deployed at CERN and/or EDG testbed. The rest will come during this year.
- ◆ What is still missing:
 - General: Scalability, Security, GUIs
 - Integration between the different fabric subsystems to build a consistent set of fabric management tools
 - From prototype to production quality

Thanks to the EU and our national funding agencies for their support of this work