# JMEA
## Job Manager Enterprise Application

Thomas Soddemann, RZG

# Overview

- RZG and DEISA
- DEISA and its resources
- Access to Resources in DEISA (s.a. next talk)
- Material Science and Plasma Physics Portal requirements
- JMEA

# History of supercomputing at the RZG

1962:  IBM 7090

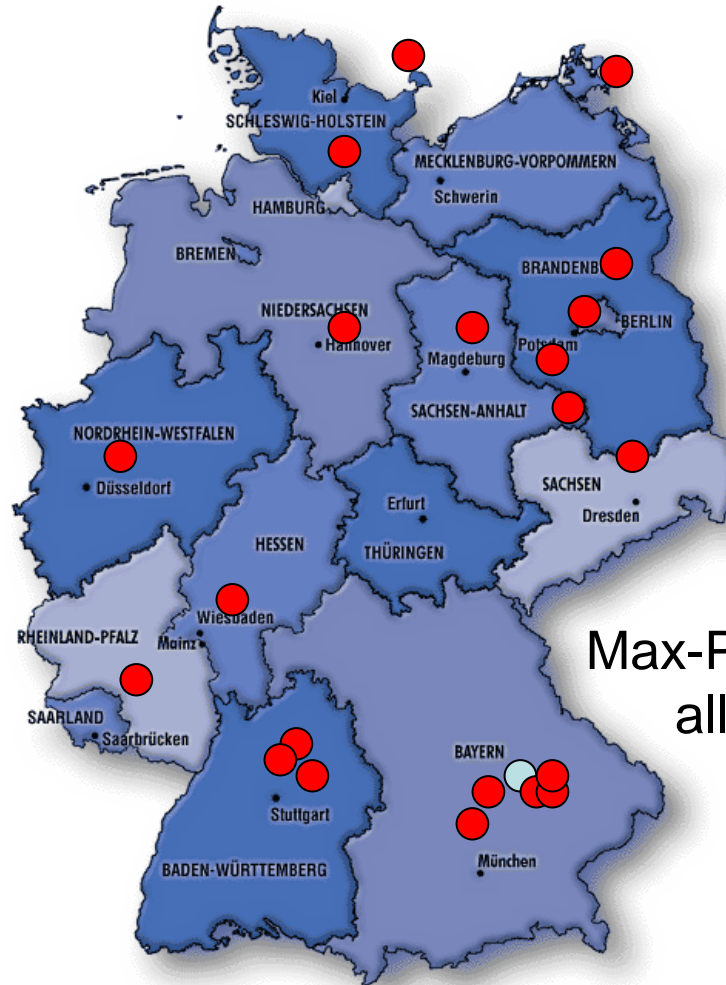1969:  IBM 360/91

1979:   Cray-1

1998: Cray T3E/816

1999: NEC SX-5/3C

2002/2003: IBM p690

# User community

- Supercomputing and Application Support

- Data management and long-term archives

- Data acquisition systems for fusion experiments

- Bioinformatics platform

Users from Max-Planck-Institutes all over Germany, Italy, and the Netherlands

# DEISA – Distributed European Infrastructure for Supercomputing Applications

# DEISA



- DEISA is an European Supercomputing Service built on top of existing national services.

- DEISA deploys and operates a persistent, production quality, distributed supercomputing environment with continental scope

# DEISA



**AIX distributed super-cluster**

## *THE DEISA SUPERCOMPUTING GRID*

**Linux systems (SGI, IBM, …)**

**Vector systems (NEC, …)**

# Deisa Network Status



Dedicated network infrastructure using „Premium IP"

**Participating NRENs**

| | |
|---|---|
| ES – RedIRIS, Spain | IT  – GARR, Italy |
| FI – FUNET, Finland | NL – SURFNET, The Netherlands |
| FR – RENATER, France | UK – UKERNA/JANET, UK |
| DE - DFN, Germany | |

# MC-GPFS Multiple Network Streams



Client

Client

File

I/O- Server

File

I/O- Server

SAN File System

File

I/O-Server

Multiple Streams

• Use of many I/O Server implies high disk performance

• Saturation of a 1 Gbit/s or 10 Gbit/s line is easier to achieve

# The DEISA Super Cluster in 2005/2006

**AIX IBM domain**

ECMWF (UK)

RZG (DE

IDRIS (FR)

CSC (FI)

**HPC Common Global File System**
similar architectures / operation systems
High bandwidth (10 Gbit/s)

CINECA (IT)

Jülich (DE)

Linux

SARA (N

LRZ (D

BSC (

**HPC Common Global File System**
**various architectures / operating systems**
High bandwidth (10 Gbit/s)
**More than 100 TFlop/s, 50 TB memory**

LINUX P

# Ways of accessing resources: CLI

ssh, qsub/llsubmit, qstat/llq, …

# Local Resource Management

- Load Leveler
- LSF
- OpenPBS/PBSpro
- Sun Grid Engine
- Torque

- MC-LoadLeveler

Obvious disadvantages:

- Separate batch script for each environment
- No job rerouting

# Ways of accessing resources: Rich Client Solution

# Ways of accessing resources: Web Portal Solution

# DEISA UNICORE infrastructure

**CINECA user**

**IDRIS user**

job

| Gateway CINECA | Gateway RZG | Gateway FZJ | Gateway IDRIS | Gateway CSC |

| NJS CINECA IBM SP5 | NJS RZG IBM SP4 | NJS FZJ IBM SP4 | NJS IDRIS IBM SP4 | NJS CSC IBM SP4 |

**NJS CINECA IBM SP5**
- uid mapping
- cluster properties

**NJS RZG IBM SP4**
- uid mapping (UUDB)
- spooler
- cluster properties (IDB)

**NJS FZJ IBM SP4**
- uid mapping
- cluster properties

**NJS IDRIS IBM SP4**
- uid mapping
- cluster properties

**NJS CSC IBM SP4**
- uid mapping
- cluster properties

Unicore Summit 2006, Dresden, Germany     August, 31st

# UNICORE deployment in DEISA

# UNICORE Configuration (RZG)



UNICORE Gateway in DMZ

NJS hosted separately

RZG Gateway

DMZ

remote gateway

SSL

SSL socket

NJS

SSL socket

SSL

SSL socket

RZG firewall

Standard socket

TSI

INTRANET

RZG IBM frames

# DEISA Research Activities

**JRA1 – Material Sciences**

CPMD

CP2K

**JRA3 – Plasma Physics**

TORB

# Requirements for a Portal Solution

- Compute Job Handling
  (submit, cancel, hold, status, …)
  => components holding job information

- Session management

- File staging support
- Remote file system access
- Database access

- User Administration (auth*)

Job Manager

Session Manager

Persistence Manager

Identity Manager

# Advantages of Portal Applications

- Give the possibility to hide the complexity of Grid Infrastructures (sensible simplifications vs. mystification MS approach)
- Can give the impression of direct use of an application

*Application Service Provider -> ASP*

- Can be accessed from almost everywhere

# Architecture of a Web Based Enterprise Solution



**Portal Application**

# Architecture of a Web Based Enterprise Solution

| Web Browser |
|:---:|

| Web Application |
|:---:|

| Other EA | JMEA |
|:---:|:---:|

| Other Middleware | UNICORE server side |
|:---:|:---:|

**UNICORE Client**

# JobManager Methods

- submit (submitting the job request),
- cancel (canceling a job request which is not being executed)
- delete (delete a finished job request),
- kill (kill a running job request)
- halt (halt a job which is being executed)
- resume (resume a previously halted job).

- Publish information about
  - resources
  - status of jobs
  - fetch console output.

- In addition
  - Support of Proxy Certificates
  - Support of Explicit Trust delegation

# The JobManager Interface



«interface»
**JobManager**

- getVirtualSites(): Collection
- setGatewayURI(in gateway: String)
- getGatewayURI(): String
- getResources(in vsite: Object, in cert: X509Certificate, in signature: Signature): Collection
- getRunningJobs(in vsite: Object, in user: X509Certificate, in signature: Signature): Collection
- getJobStatus(in jobId: Object, in vsite: Object, in user: X509Certificate, in signature: Signature): JobStatus
- cancelJob(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature)
- killJob(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature)
- haltJob(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature)
- resumeJob(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature)
- deleteJob(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature)
- getResults(in jobId: Object, in vsite: Object, in cert: X509Certificate, in signature: Signature): JobResult

Allows different implementations:
- UNICORE (primary target)
- Globus
- …

# Arcon Library disadvantages (Multi User Application)

In order to avoid race conditions in multi threaded applications,
one should

- omit static variables

  unless they are used for communication between the threads and
  their access is synchronized.

- synchronized access must not lead to a performance bottlenecks

- Keep in mind that the application might be clustered

# Arcon Library disadvantages (Multi User Application)

Arcon Library defines:

- `outcome_dir` which specifies the directory, where streamed files will be stored
- `buffer_size` which reflects the buffer size for connections
- `always_poll` which tells if request are always asynchronous or not.

The abstract class Connection implements three static variables:

- `keep_open` which defines if the next connection is kept open after use.
- `compression` which tells if the transmission should be compressed or not.
- `encrypt` which defines whether the next retrieved connection should be encrypted communication or not.

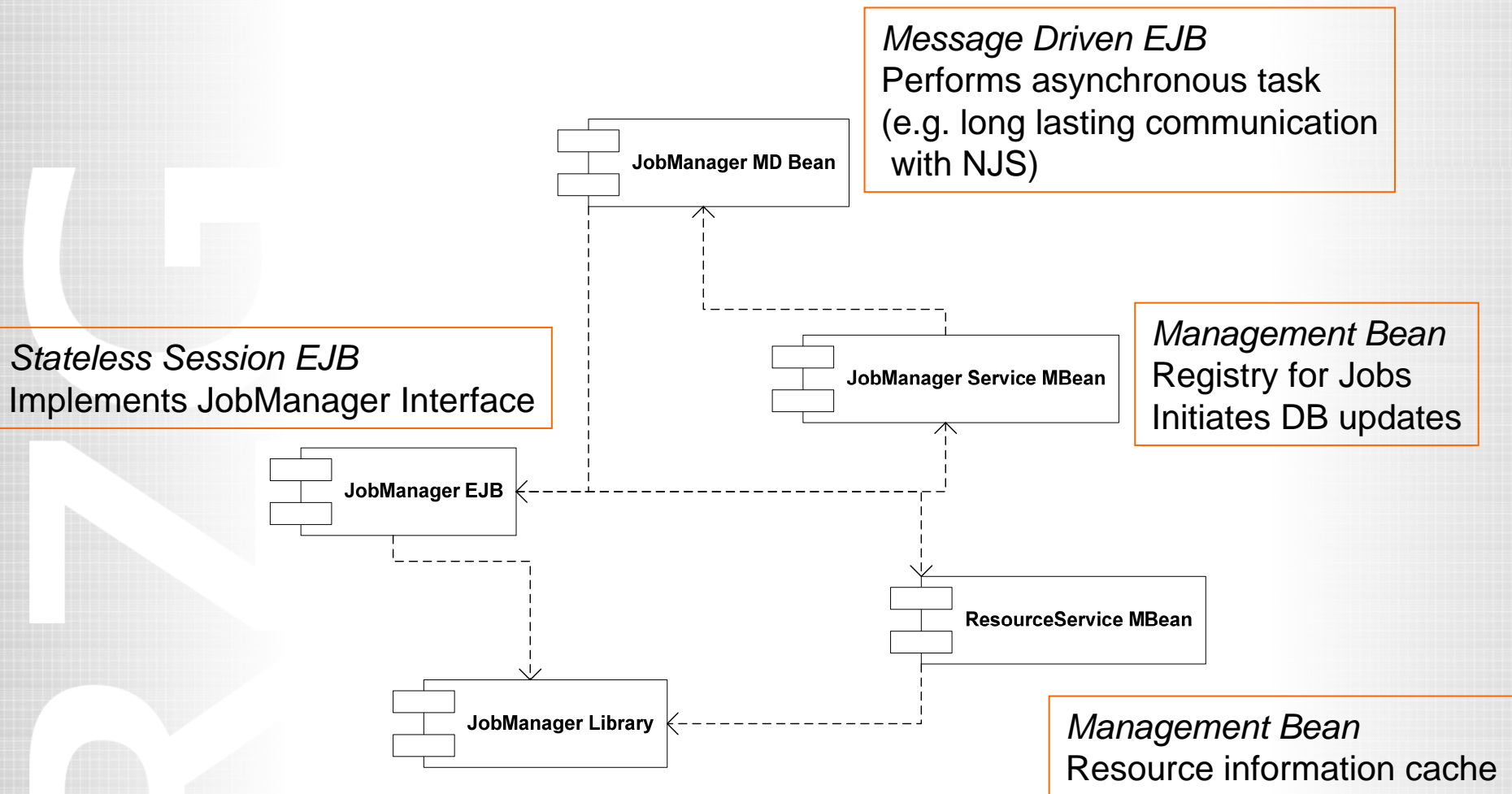# Arcon Library disadvantages (Multi User Application)

Further disadvantages

- Proprietary Logging

- Exceptions used for control flow rather than for error handling

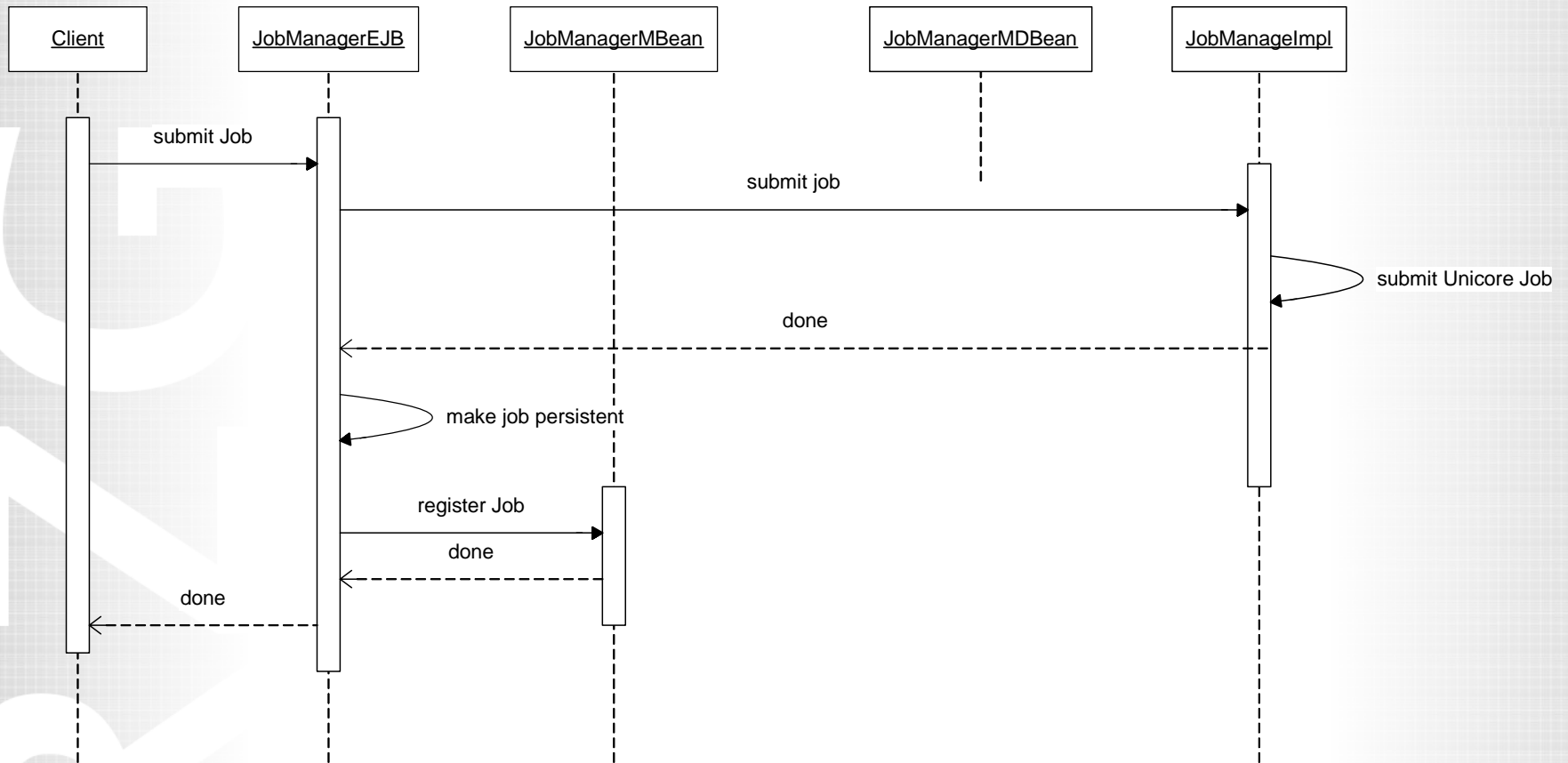- Missing support for ETD (at least in the official release)

# New UNICORE JobManager library

- Need for a new job management library implementation for UNICORE
- Arcon partly code reused/refactored
- Reimplementation of problematic parts
- ETD support
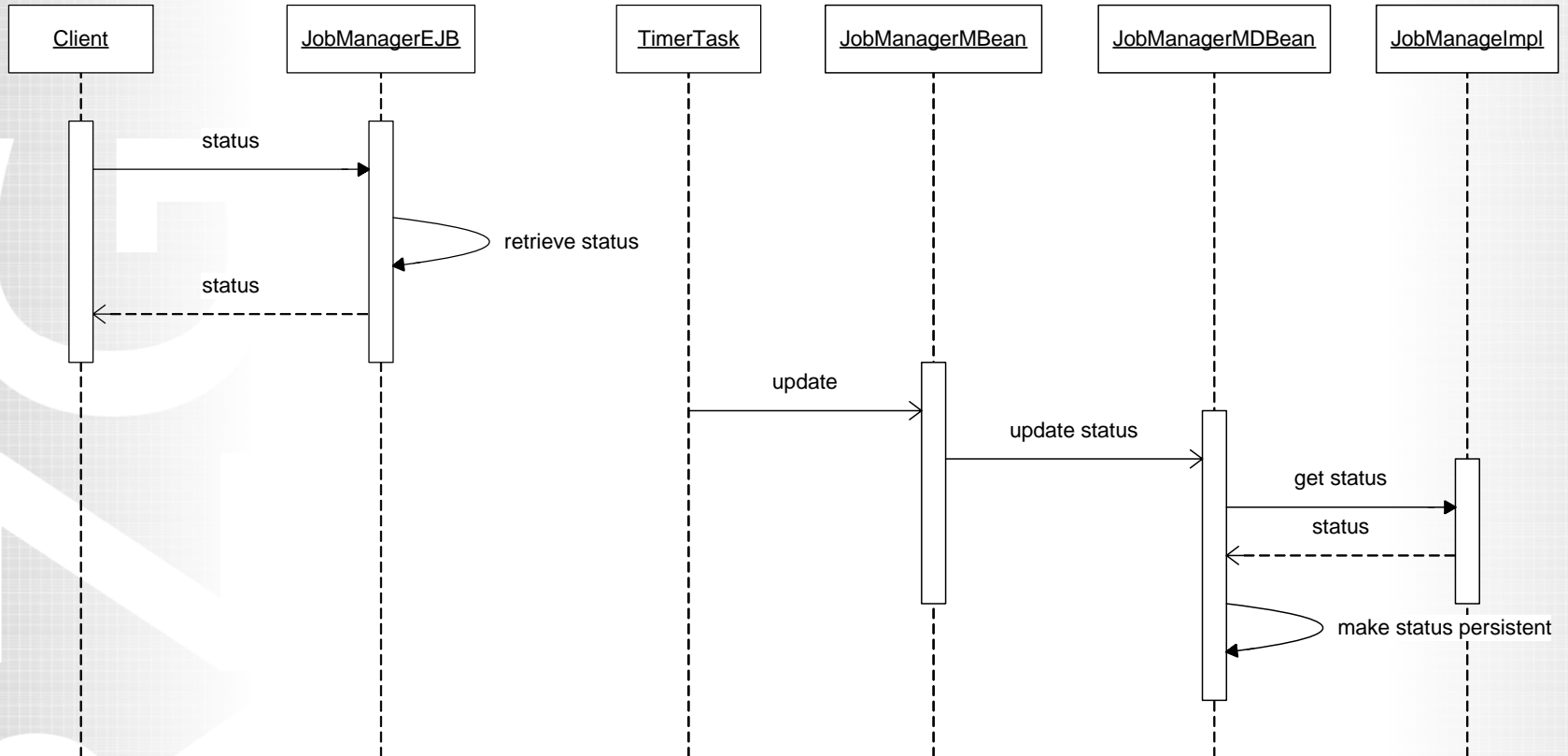- Proxy Certificate Support
- Thread safe

# Job Management Enterprise Application - JMEA

Message Driven EJB
Performs asynchronous task
(e.g. long lasting communication
 with NJS)

JobManager MD Bean

Management Bean
Registry for Jobs
Initiates DB updates

JobManager Service MBean

Stateless Session EJB
Implements JobManager Interface

JobManager EJB

ResourceService MBean

JobManager Library

Management Bean
Resource information cache

# Job Submission Sequence
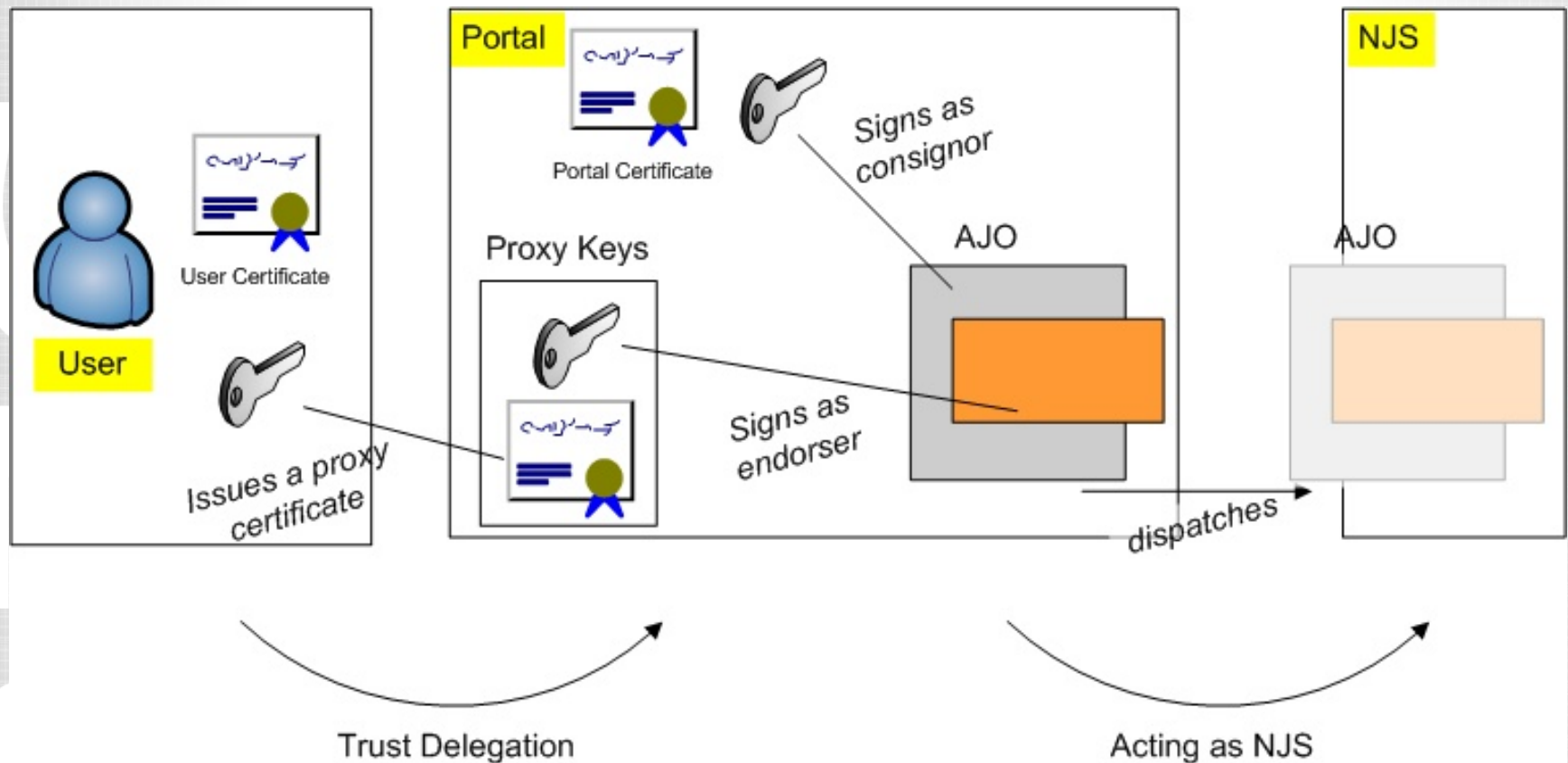
# Job Status Request Sequence

# Facts on JMEA

Advantages

- Responding fast to client requests
- Scalable (number of client requests)
  - Implies scalable database and container infrastructure
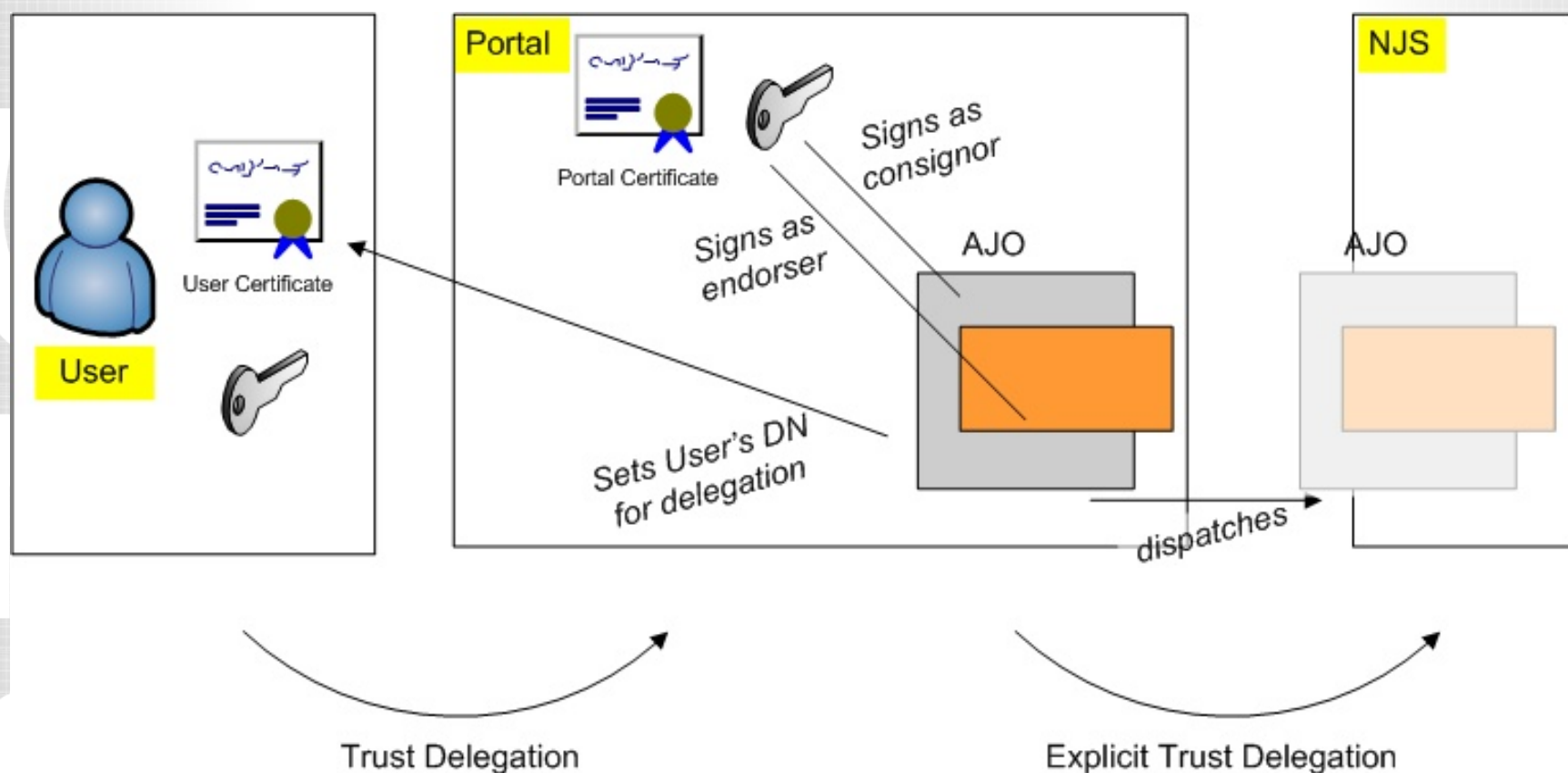- fault tolerant (to some extend)

Disadvantages

- Risk of delivering outdated data
- No support UNICORE alternate file transfer
  The web application has files transferred independently (GPFS, CIFS)

# Security Solutions: Proxy Certificate Approach

# Security Solutions: Explicit Trust Delegation

# Possible UNICORE5 improvements

AJO

- Default constructors for all classes (unless it does not make sense)
  - Background: Persistence

ETD

- Use of X500Principals instead of whole X509 certificates as user attribute

- Allowing more "direct" requests for an ETD agent
  - E.g. for Resource information

# Conclusion

- JMEA is an EA which proves to work with UNICORE5 in DEISA
- It has all basic features implemented which are needed for successful job management
- It is designed to work with the JRA1/JRA3 Web application
- It can be used in a different context (OMII/GridSAM)
- But, it does not provide a standards based interface

- WS-GRAM (4.0 and 4.1) support is being developed
- UNICORE6 support is hopefully given with WS-GRAM 4.2 support

# Thanks

Contributors to the JRA1/JRA3 endeavors:

Johannes Reetz, RZG, Garching, Germany
Daniel Frank, RZG, Garching, Germany

Discussion Partners (GridSAM)
Stephen McGough, IC, London, England

Researchers/Beta-Testers
Matthias Krack, ETH, Zürich, Switzerland
Peter Coveney, UCL, London, England