# dCache, a managed storage in grid

Patrick
for the dCache Team

support and funding by

*Patrick Fuhrmann et al.*   *BE-grid, Brussels, BE*   *October 16, 2007*

# Topics

Project Topology

Why do we need storage elements in the grid world ?

The idea behind the LCG (gLite) storage element.

Available Solutions

The dCache implementation

dCache in a nutshell

Weak points and outlook

Usage

Selected Topics

# Project Topology : *The Team*

## Head of dCache.ORG

Patrick Fuhrmann

## Core Team (Desy and Fermi)

Andrew Baranovski
Bjoern Boettscher
Ted Hesselroth
Alex Kulyavtsev
Iryna Koslova
Dmitri Litvintsev
David Melkumyan
Dirk Pleiter
Martin Radicke
Owen Synge
Neha Sharma
Vladimir Podstavkov

## Head of Development FNAL :

Timur Perelmutov

## Head of Development DESY :

Tigran Mkrtchyan

## External

### Development

Gerd Behrmann, NDGF
Jonathan Schaeffer, IN2P3

### Support and Help

Abhishek Singh Rana, SDSC

Greig Cowan, gridPP

Stijn De Weirdt (Quattor)

Maarten Lithmaath, CERN

Flavia Donno, CERN

*We need to serve large amounts of data locally*

- *Access from local Compute Element*
- *Huge amount of simultaneously open files.*
- *Posix like access  (What does that mean ?)*

*We need to exchange large amount of data with remote sites*

- *Streaming protocols.*
- *Optimized for low latency (wide area) links.*
- *Possibly controlling 'link reservation'.*

dCache.ORG

## *We need to allow storage control*

- *Space reservation to guarantee maximum streaming.*

- *Define space properties (TAPE, ONLINE, ...)*

- *Transport protocol negotiation.*

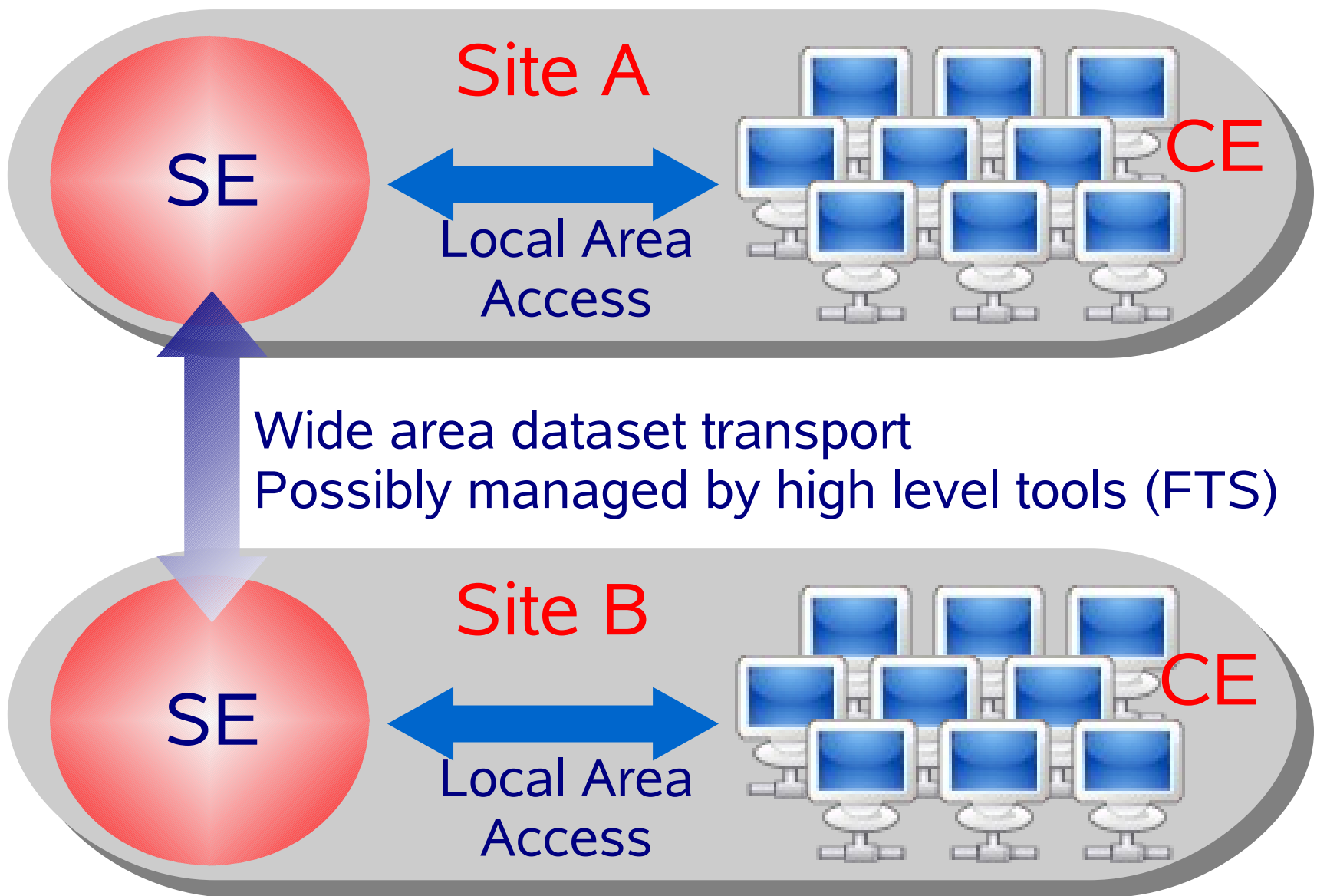## *We need to publish SE specific information*

- *Clients need to select 'best' SE or CE for a job.*
- *Availability*
- *Available Space (max, used, free ...)*
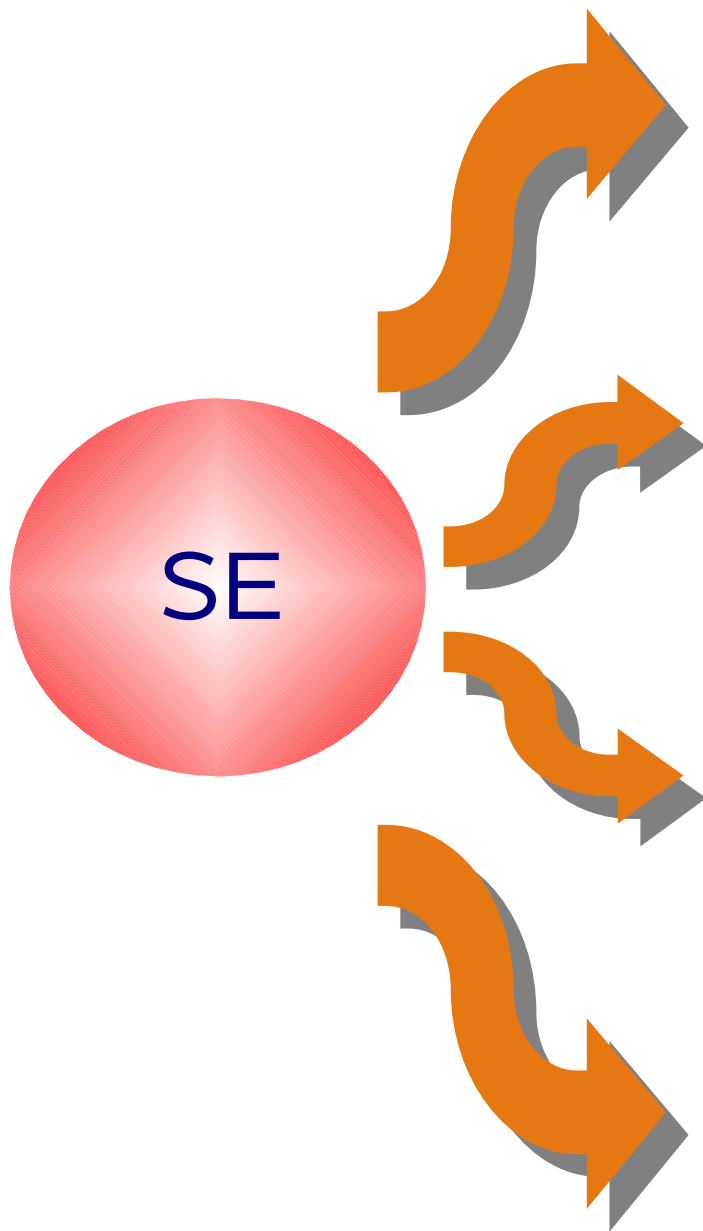- *Supported Spaces (Tape, disk ...)*
- *Which VO owns which space ?*

*dCache.ORG*

**SE**

## Information Publishing
Content : GLUE
Transport : LDAP

## SRM Storage Resource Management
Space/Protocol Management

## Wide Area Transport Protocol
In use : gsiFtp
Discussed : http(s)
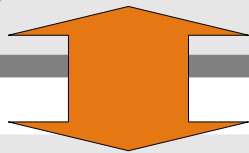
## Local Access Protocol
(gsi)dCap  or rfio and xRoot

dCache.ORG

dCache.ORG

*dCache.ORG*

**Common Protocols**
infoProvider, SRM, gsiFtp, rfio, dCap, xRoot

*dCache*

*CASTOR*

*DPM*

*SToRM*

*GPFS*

*TSM ®*
*HPSS ®*
*DMF ®*
*Enstore, OSM*

# The dCache SE implementation
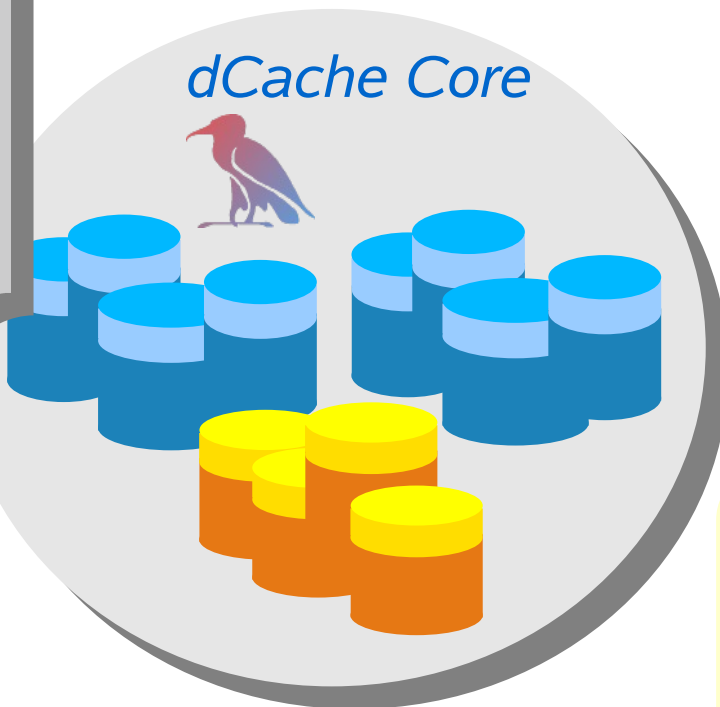# Black Box View

**High Level Services**

Resilient Manager

Admin Module (ssh, jpython)

Maintenance Module

Flush Manager
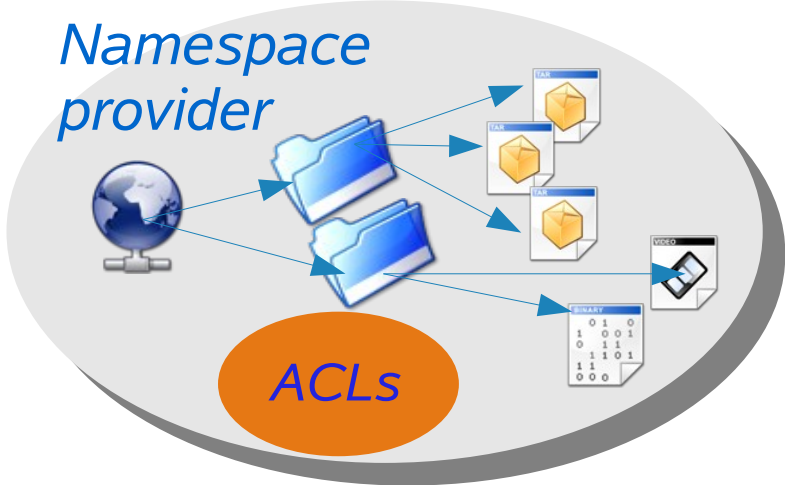
Hopping Manager

**dCache Core**

*Information Protocol(s)*

*Storage Management Protocol(s)*
SRM 1.1  2.2

**Tape Storage**

*OSM, Enstore Tsm, Hpss, DMF*

*Data & Namespace Protocols*

(NFS 4.1)    dCap

ftp (V2)    gsiFtp

xRoot

(http)

**Namespace provider**

*Namespace ONLY*
NFS 2 / 3

ACLs

dCache.ORG

Patrick Fuhrmann et al.    BE-grid, Brussels, BE    October 16, 2007
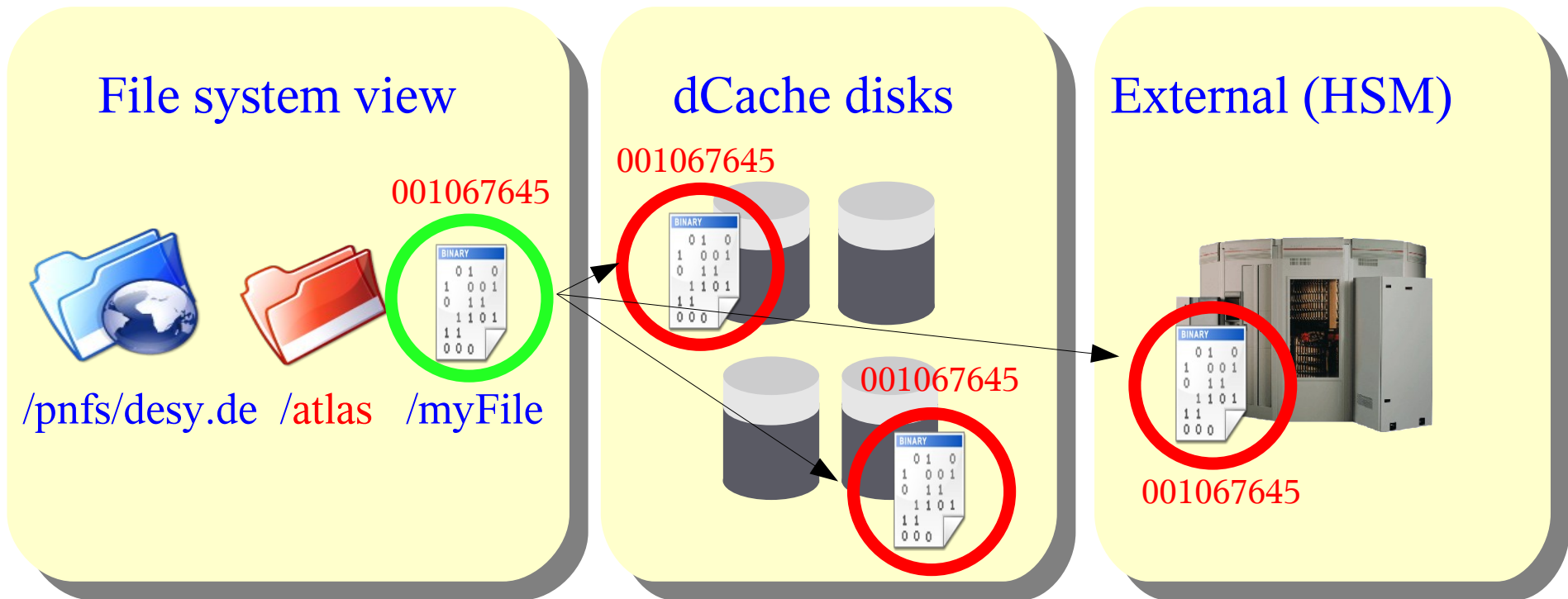
# dCache in a Nutshell

- **Strict name space and data storage separation, allowing**

  - consistent name space operations (mv, rm, mkdir e.t.c)

  - consistent access control per directory resp. file

  - managing multiple internal and external copies of the same file

  - convenient name space management by nfs (or http)



**File system view**

001067645

/pnfs/desy.de   /atlas   /myFile

**dCache disks**

001067645

001067645

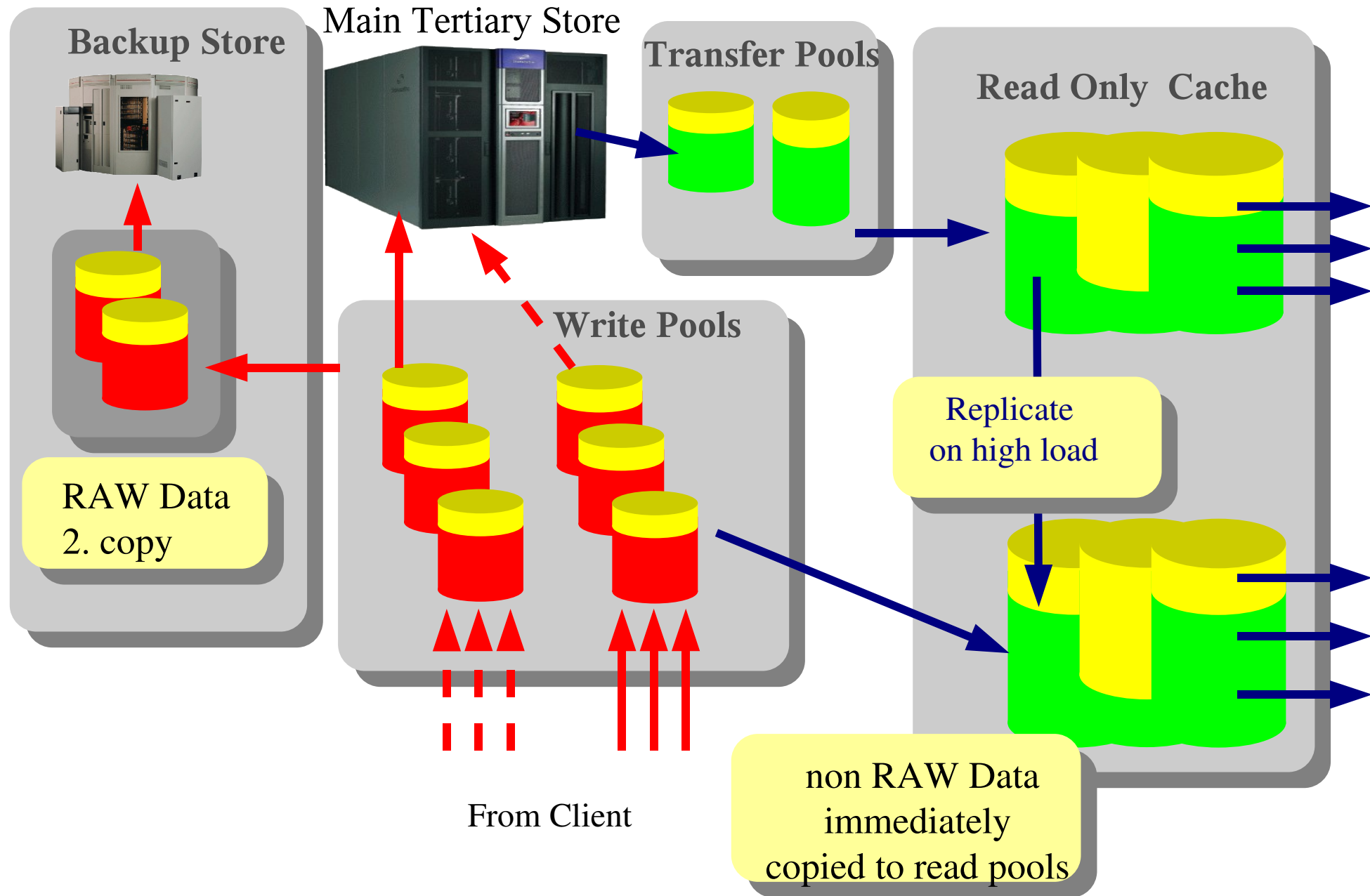**External (HSM)**

001067645

001067645

- **dCache partitioning for very large installations**

- **File hopping on**

    - automated hot spot detection

    - configuration (read only, write only, stage only pools)

    - on arrival (configurable)

# *File Hopping*



Backup Store

Main Tertiary Store

Transfer Pools

Read Only Cache

RAW Data
2. copy

Write Pools

Replicate
on high load

From Client

non RAW Data
immediately
copied to read pools

- **Overload and meltdown protection**

  - Request Scheduler.

  - Separate I/O queues per protocol    (load balancing)

- **Supported protocols : (gsi)ftp , (gsi)dCap, xRoot, SRM, nfs2/3**

- **xRoot support**

  - Vector read

  - Currently working on asyn I/O

dCache.ORG

*I/O Request*

*Space Manager*

*Pool Candidates selected by*
Protocol
Client IP number/net
Data Flow Direction
Name Space Attributes (Directory)
SRM Spaces

*Dispatcher by request Attributes*

dCache.ORG

*List of candidates*

*Dispatcher by Pool Cost*

*Pool Protocol Queues*

*xRoot*    *dCap*    *gsiFtp*

# In the Nutshell

**HSM Support**

- TSM, HPSS, DMF, Enstore, Osm
- Automated migration and restore
- Working on Central Flush facility
- support of multiple, non overlapping HSM systems (NDGF approach)

**Misc**

- Graphical User Interface
- Command line interface
- Jpython interface
- SRM watch
- NEW : Monitoring Plots

*Weak points :*

      *Posix like is NOT posix (file system driver)*

      *Http(s) not really supported*

      *Security might not be sufficient*

# Posix like is NOT posix

*dCache.ORG*

*dCache.ORG*

**Linked Library**

**Preload Library**

App
dCap

App

libC
iNodes
xfs    ext3

SE

preload dCap
libC
iNodes
xfs    ext3

*Application needs to be linked with the dCap library.*

*Application stays unchanged but doesn't work in all cases. (Static linked, Some C++ apps.)*

*App*

*libC*

*iNodes*

*xfs*  *NFS 4.1*  →  SE

*Application doesn't need to be changed.*
*NFS 4.1 driver comes with OS.*

*Solution is on the way....*

We are currently putting significant efforts in the NFS 4.1 protocol

*Deployment Advantages :*

   Clients are coming for free (provided by all major OS vendors).

*Technical Advantages :*

- NFS 4.1 is aware of distributed data

- Faster (optimized) e.g.:
  - Compound RPC calls
  - 'Stat' produces 3 RPC calls in v3 but only one in v4

- GSS authentication
  - Built in mandatory security on file system level

- ACL's

- OPEN / CLOSE semantic (so system can keep track on open files)

- 'DEAD' client discovery (by client to server pings)

# Goal : Industry standards in HEP ?

dCache.ORG

dCache.ORG

## In use at 9 Tier I centers

- fzk (Karlsruhe, GR)
- in2p3 (Lyon, FR)
- RAL (Rutherford, UK)
- BNL (New York. US)
- FERMILab (Chicago, US)
- SARA (Amsterdam. NL)
- PIC (Spain)
- Triumf (Canada)
- NDGF (NorduGrid)

About 40 Tier II's

dCache is part of VDT (OSG)

We are expecting > 20 PB per site > 2011

**dCache will hold the largest share of the LHC data.**

Tier II DESY

Tier 1 fzk

Tier 1 ***

Tier 0 CERN

Tier 1 RAL

Tier 1 IN2P3

300 MB/Sec

*Some more hot topics*

*The wonderful world of*

# SRM 2.2

*Only if there is a lot of time left*
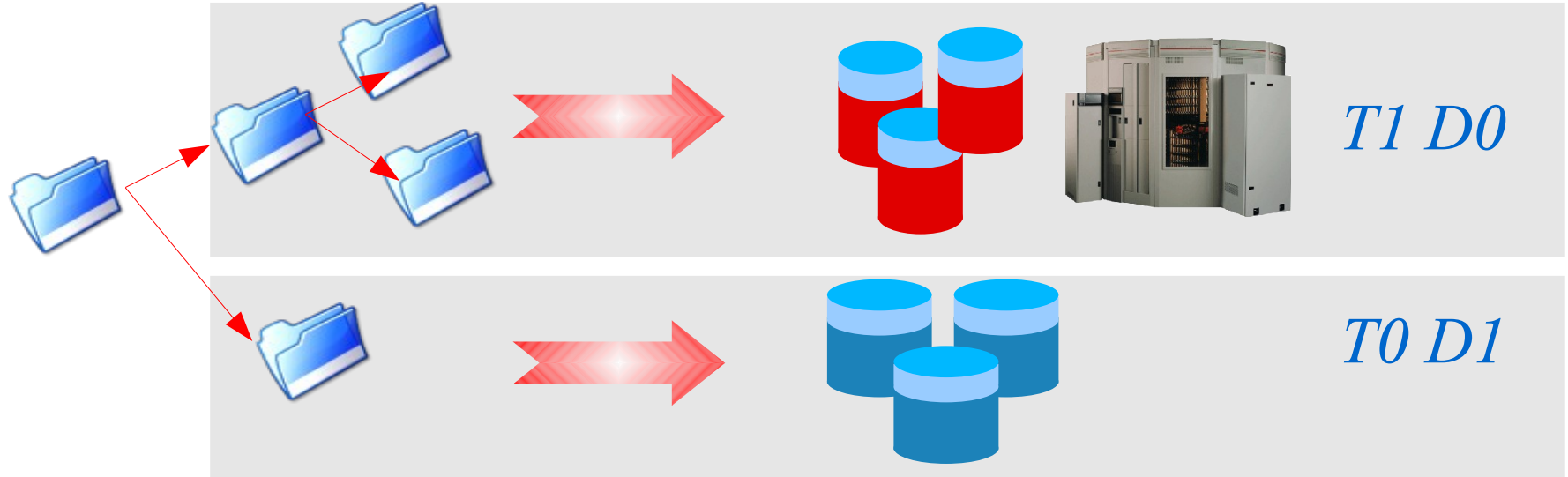
dCache.ORG

dCache.ORG

*The SRM in dCache supports*

- *CUSTODIAL (T1Dx)*
- *NON-CUSTODIAL (T0D1)*
- *Dynamic Space Reservation*
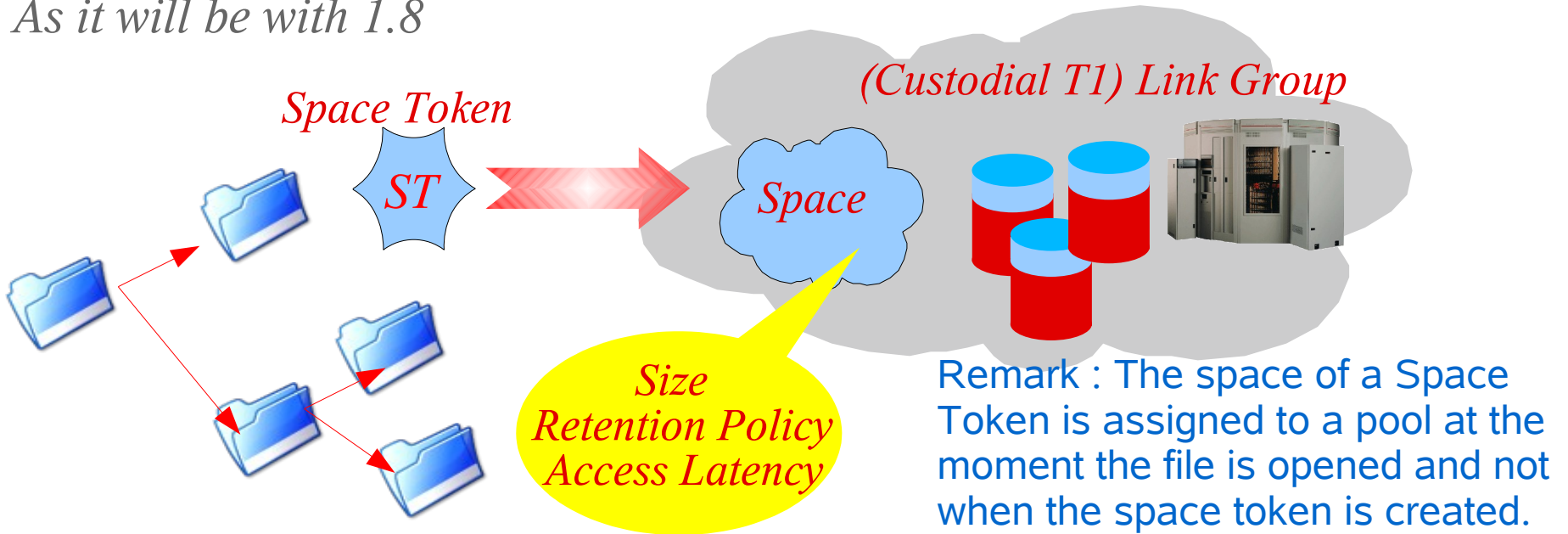- *late pool binding for spaces*
- *and more*

**dCache.ORG**

*As it used to be ( <= 1.7 )*



*T1 D0*

*T0 D1*

**dCache.ORG**

*As it will be with 1.8*

*Space Token*

*ST*

*(Custodial T1) Link Group*

*Space*

*Size
Retention Policy
Access Latency*

Remark : The space of a Space Token is assigned to a pool at the moment the file is opened and not when the space token is created.

# Further reading

## www.dCache.ORG