# Data Management
# in Grid Enviroment

# Data Management in Grids

➢ *Safe and reliable file transfer between grid nodes*

➢ *Gathering and maintenance of information about data sources in the grid*

➢ *Data access optimization*

➢ *Providing of metadata services*

➢ *Data sets consistency maintenance*

➢ *Virtual organization support*

➢ *Security (authorization, authentification)*

➢ *Storage systems*

➢ *Grid monitoring services*

**EGEE**
Enabling Grids for
E-science in Europe

## *Based on FTP:*

➢ *Widely implemented and well-understood standard protocol*

➢ *Well-defined architecture for protocol extensions*

➢ *Numerous extensions where already added*

➢ *Transfers between client and server support*
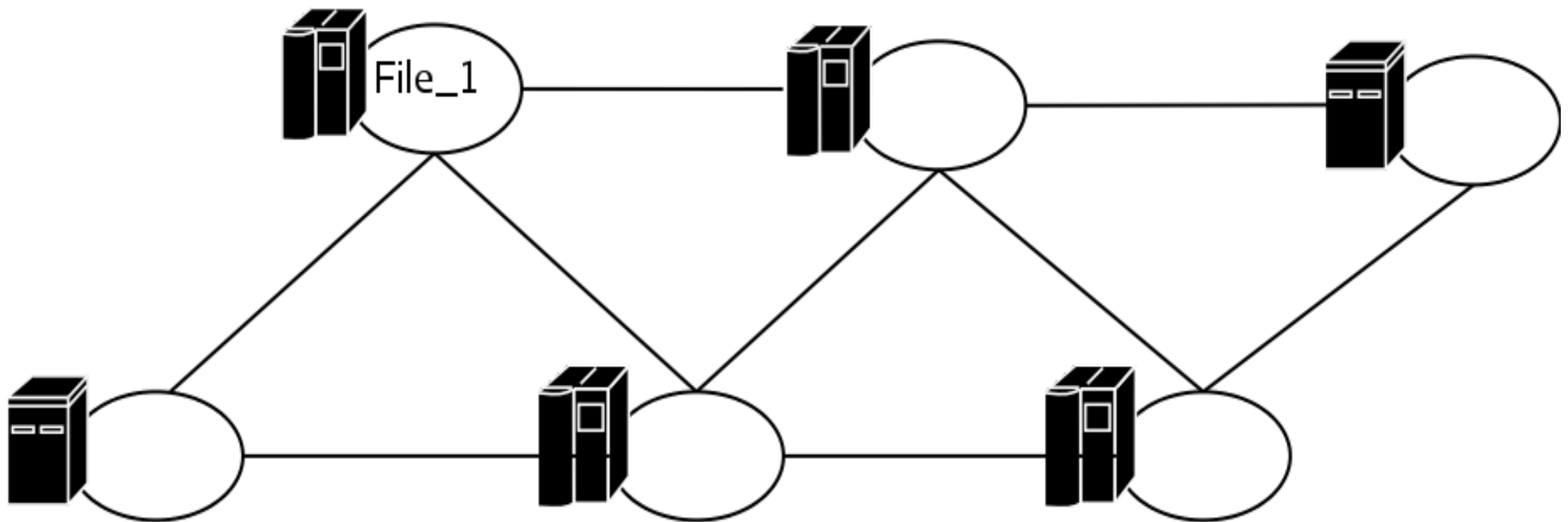
➢ *Third party transfers between two servers support*

## GridFTP - extension of FTP:

➢ **Grid Security Infrastructure (GSI) and Kerberos support**

➢ **Third-party control of data transfer**

➢ **Parallel data transfer**

➢ **Partial file transfer**

➢ **Stripped data transport**

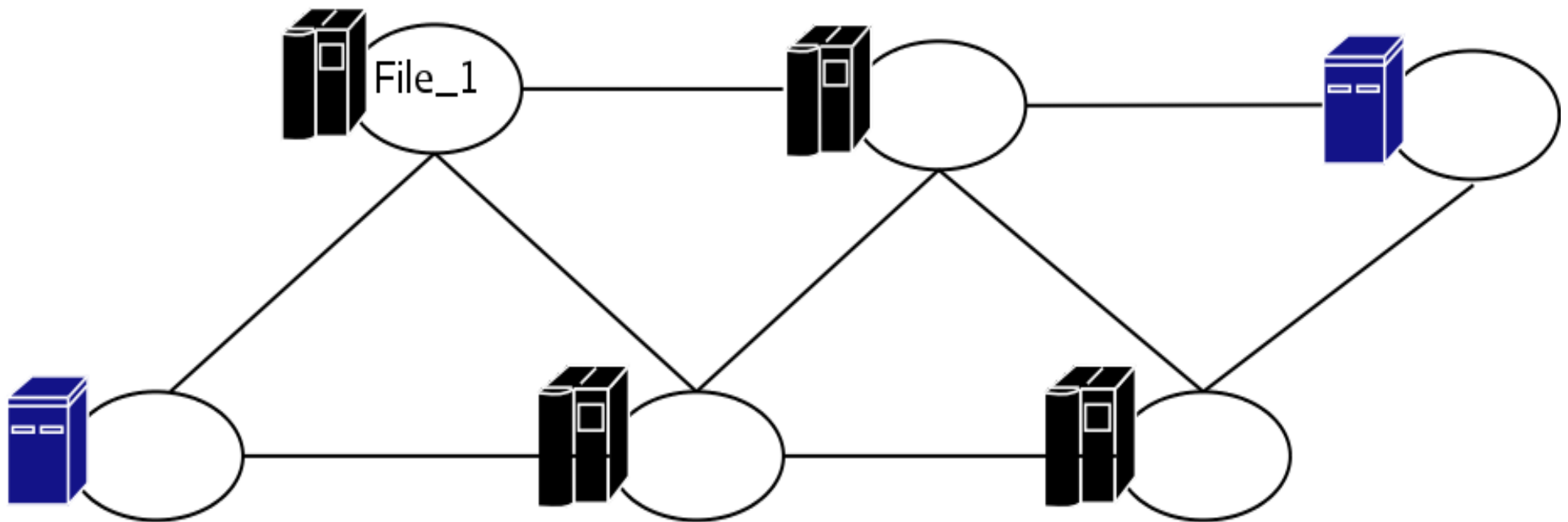➢ **Support for reliable data transfer**

## *Replication of Data Sets - key concept to increase data availability*

## Replication of Data Sets - key concept to increase data availability

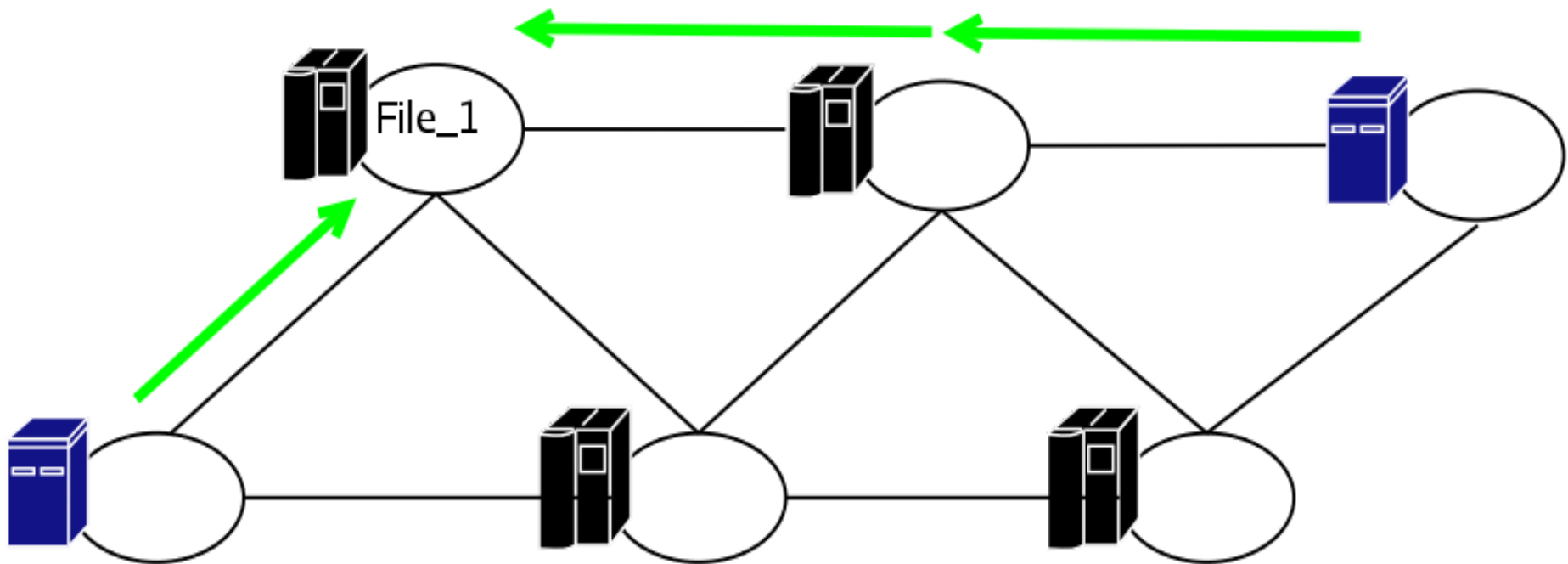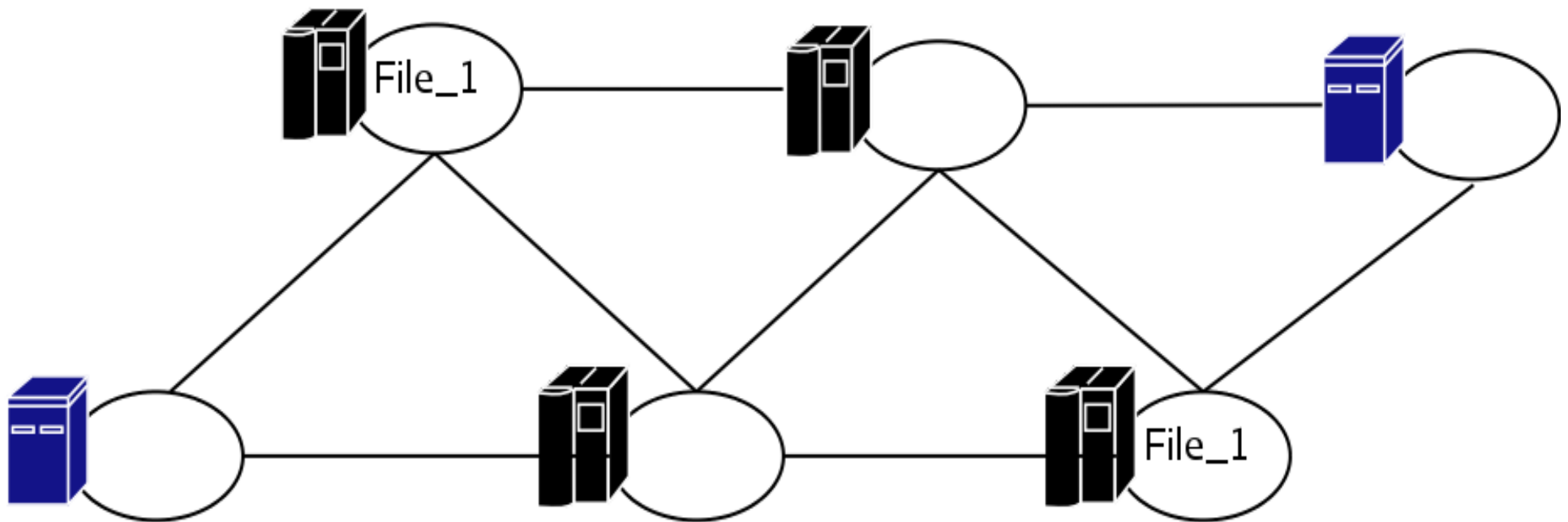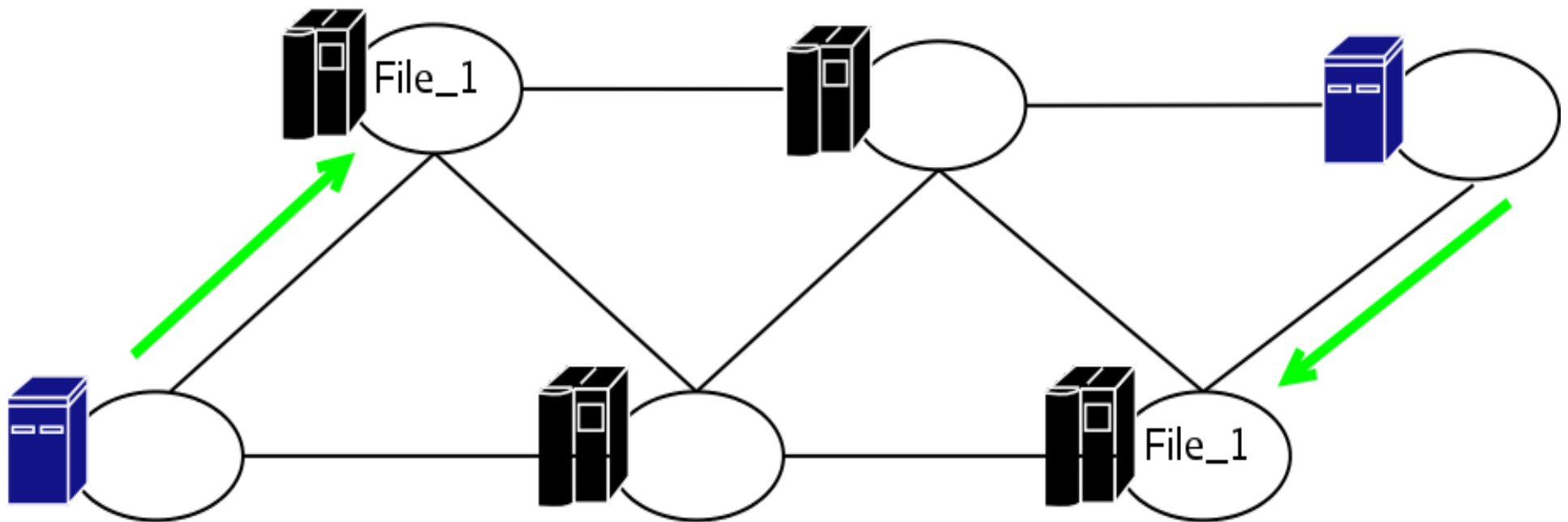## Replication of Data Sets - key concept to increase data availability

## Replication of Data Sets - key concept to increase data availability

**eGee**
Enabling Grids for
E-science in Europe

*Replication of Data Sets -* key concept to increase data availability

# Replica Location Service

➢ *Provide information about replicas location*

➢ *Replicas are identified by LFN and PFN*

➢ *LFN - Logical File Name*

  *- global unique file identifier*

➢ *PFN – Physical File Name*

  *- replica name on concrete Storage Element*

## Replica Location Catalog

➢ **Service run at each grid node**

➢ **Stores information about files stored at the grid node**

➢ **Provides mapping between LFN and PFN**

➢ **Propagate information to Replica Location Index**

➢ **Soft state update mechanism**

**Information Propagation:**
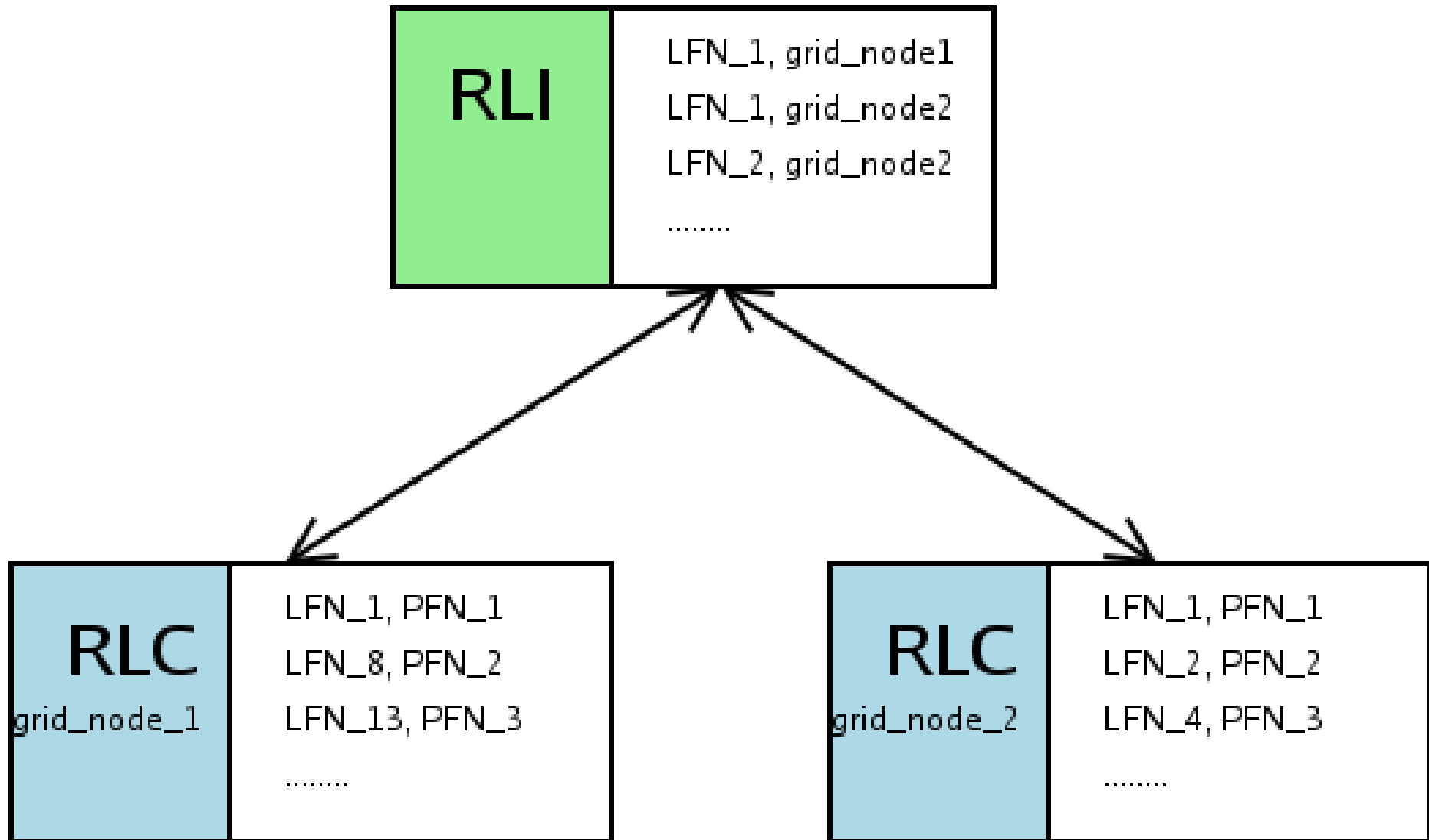
➢ **Full**

➢ **Compressed (Bloom filters)**

## *Replica Location Index*

➢ *Central grid service*

➢ *Gathers information from grid RLCs*

➢ *Provides mapping between LFN and a set of grid nodes containing given file*

# Replica Location Service

| RLI | LFN_1, grid_node1<br>LFN_1, grid_node2<br>LFN_2, grid_node2<br>........ |
| --- | --- |

| RLC<br>grid_node_1 | LFN_1, PFN_1<br>LFN_8, PFN_2<br>LFN_13, PFN_3<br>........ |
| --- | --- |

| RLC<br>grid_node_2 | LFN_1, PFN_1<br>LFN_2, PFN_2<br>LFN_4, PFN_3<br>........ |
| --- | --- |

## Replica Manger Service

**Integrates GritFTP and catalog services**

### RLS Implementations

**EU Data Grid**

   -only full information transfer between RLCs and RLI

**Globus Toolkit**

   - compressed information transfer to RLI

   - bloom filters

# Replica Access Optimization

## Short-term optimization

➢ select best file replica for given grid node an LFN

➢ selection based on current network status & replica location

➢ implemented in EU Data Grid project

## Long-term optimization

➢ automated file replication and deletion

➢ prediction of replica usefulness at a given site

➢ only simulations were done

➢ *descriptive information about data*

➢ *lot of metadata types*
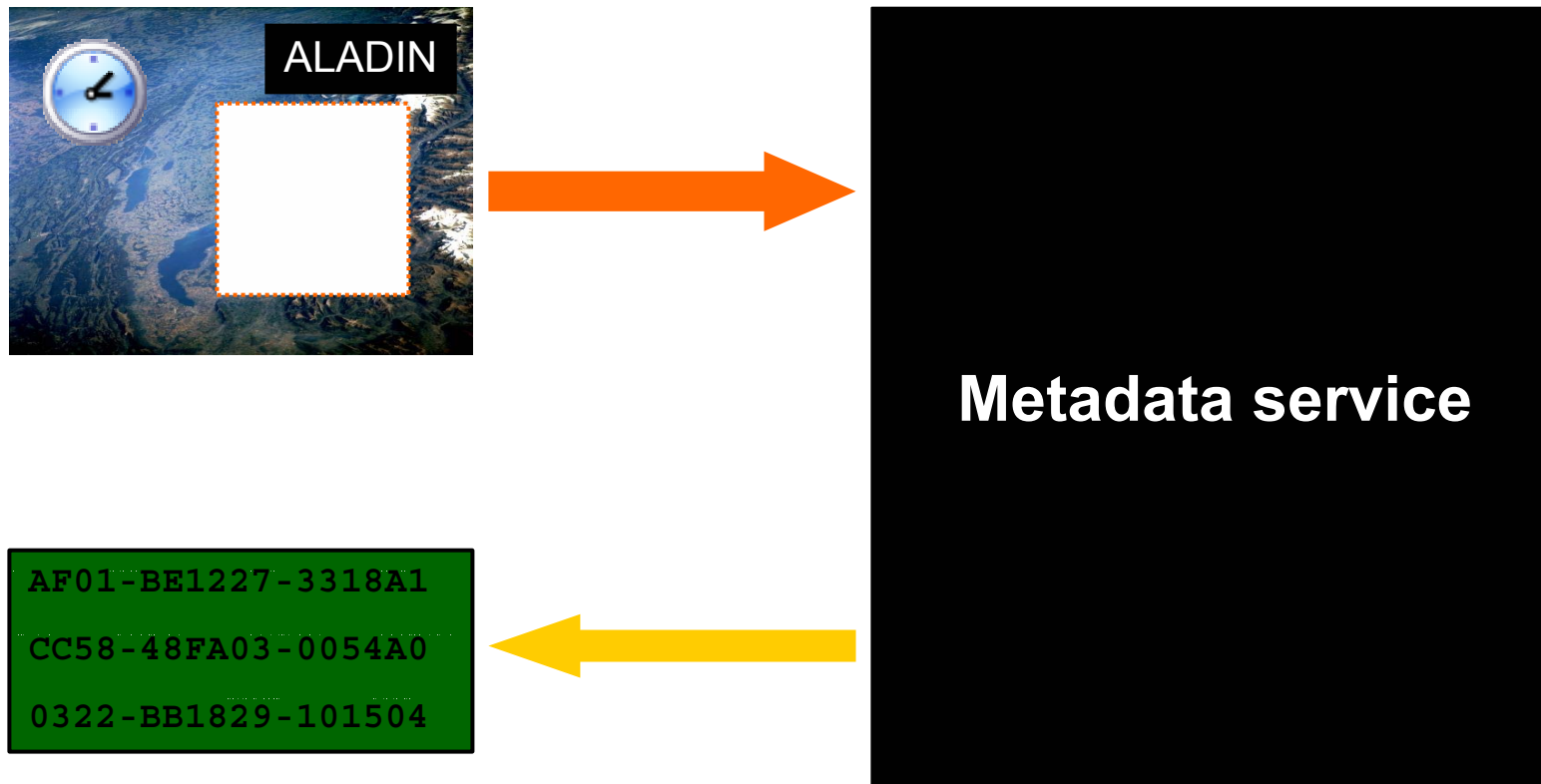
*general metadata (file type, file size, creation date)*

*security metadata*

*application specific metadata*

➢ *metadata description by ontologies*

# Replica Consistency Services

- **difficult task**

- **standard transactional two-phase commit protocol is not suitable for grids**

- **current state: no replica updates are allowed**

- **data update only by file versioning**