

Or
What is SRB Matrix?

Data Grid Automation

**Arun Jagatheesan et al.,
San Diego Supercomputer Center
University of California, San Diego**

VLDB Workshop on Data Management in Grids
Trondheim, Norway, 2-3 September 2005



Talk Outline

- **Data grid Landscape**
- **Long-run data management processes**
 - Data Grid ILM
 - Data Grid Triggers
 - Dataflow Pipelines
- **Execution Logic – Data Grid Language**



Data Grid Landscape

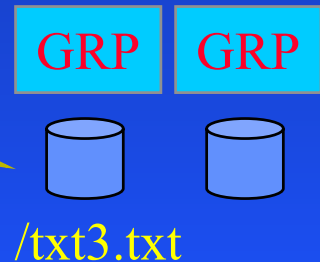


The "Grid" Vision



Data Grid Resource Providers

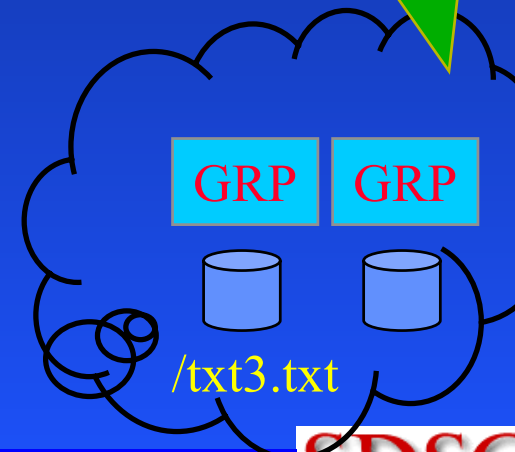
Grid Resource Providers
(GRP) providing content
and/or storage



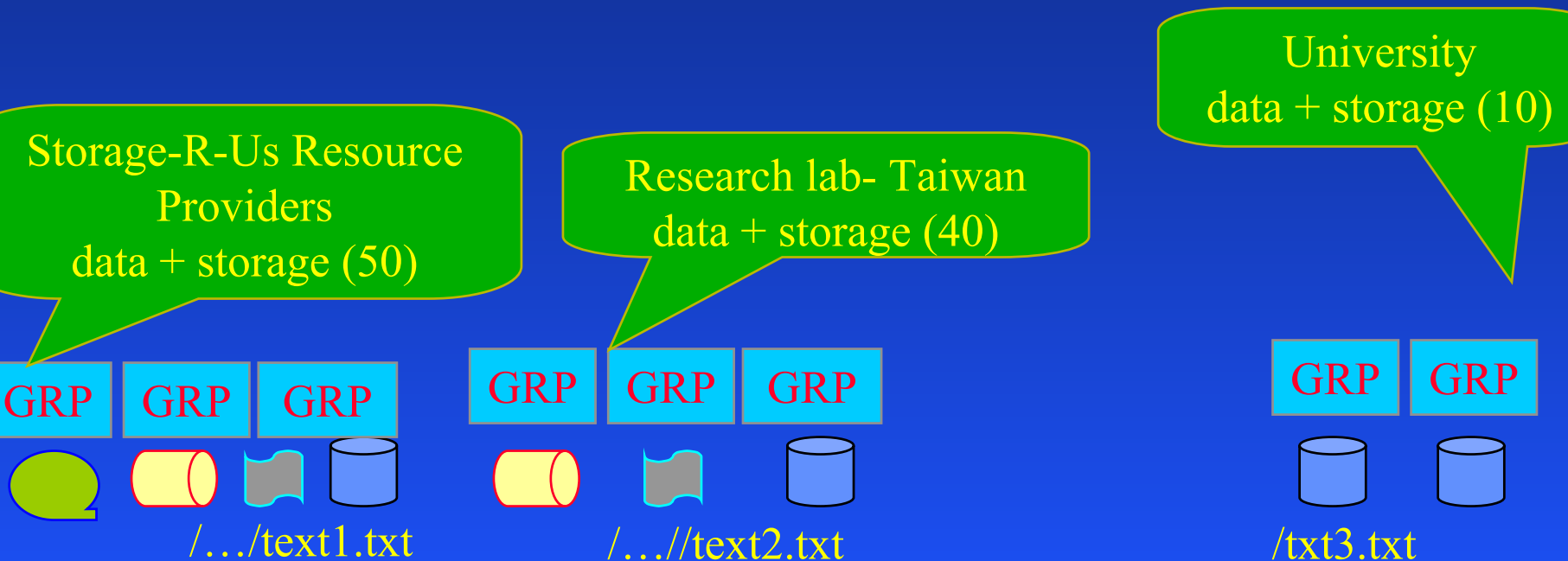
Data Grid Administrative Domain

- Single administrative domain with one or more GRP Resource Providers
- Example - building a shared collection

Research Lab



Data Grid Administrative domains



Data Grid (Enterprise Utility)

Physical Resources managed by autonomous administrative domains of the same enterprise (ABCZ.com)

IT Department US



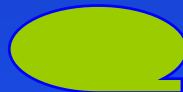
ABCZ.com US

IT Department Asia



ABCZ.com Asia

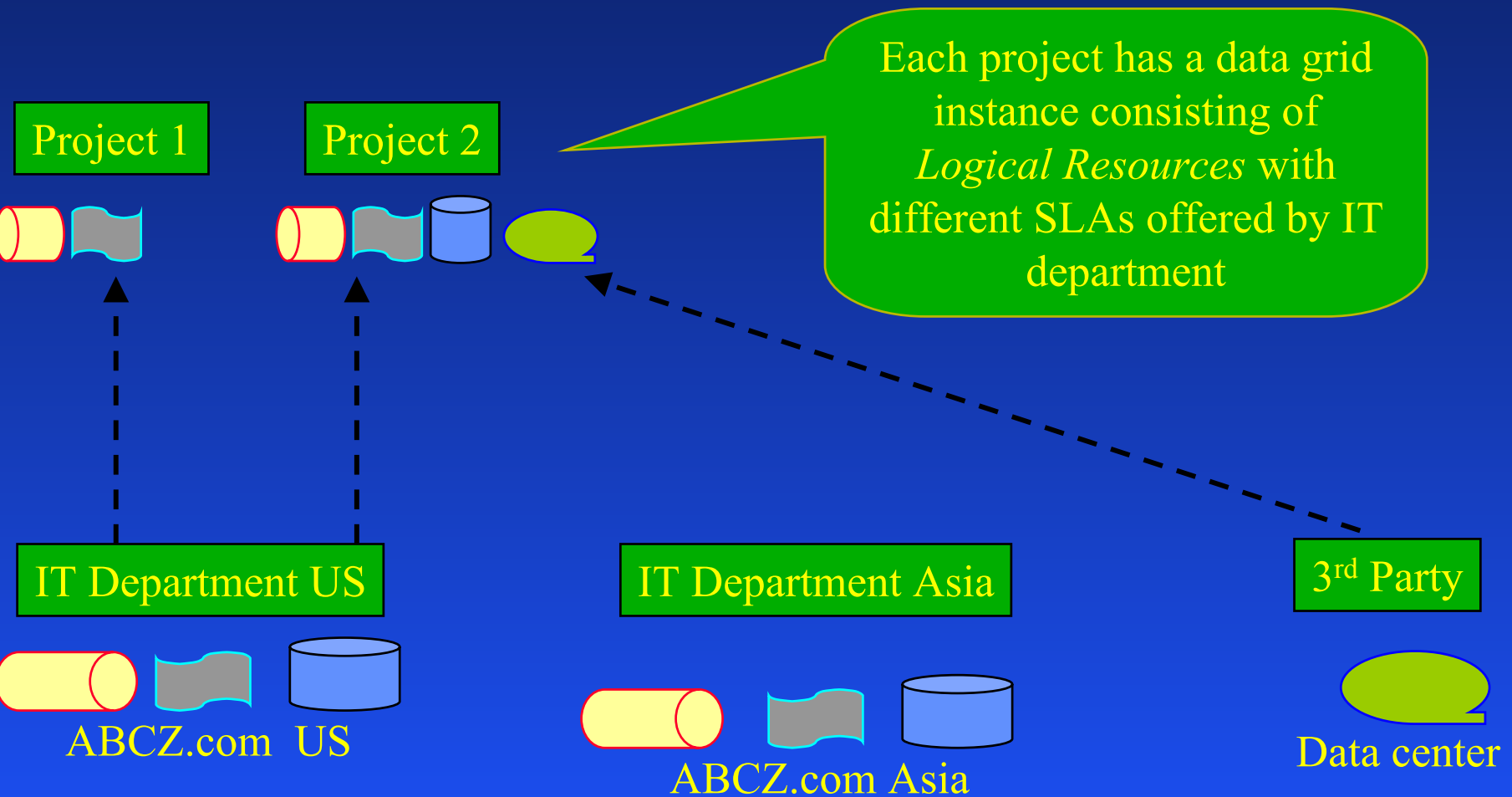
3rd Party



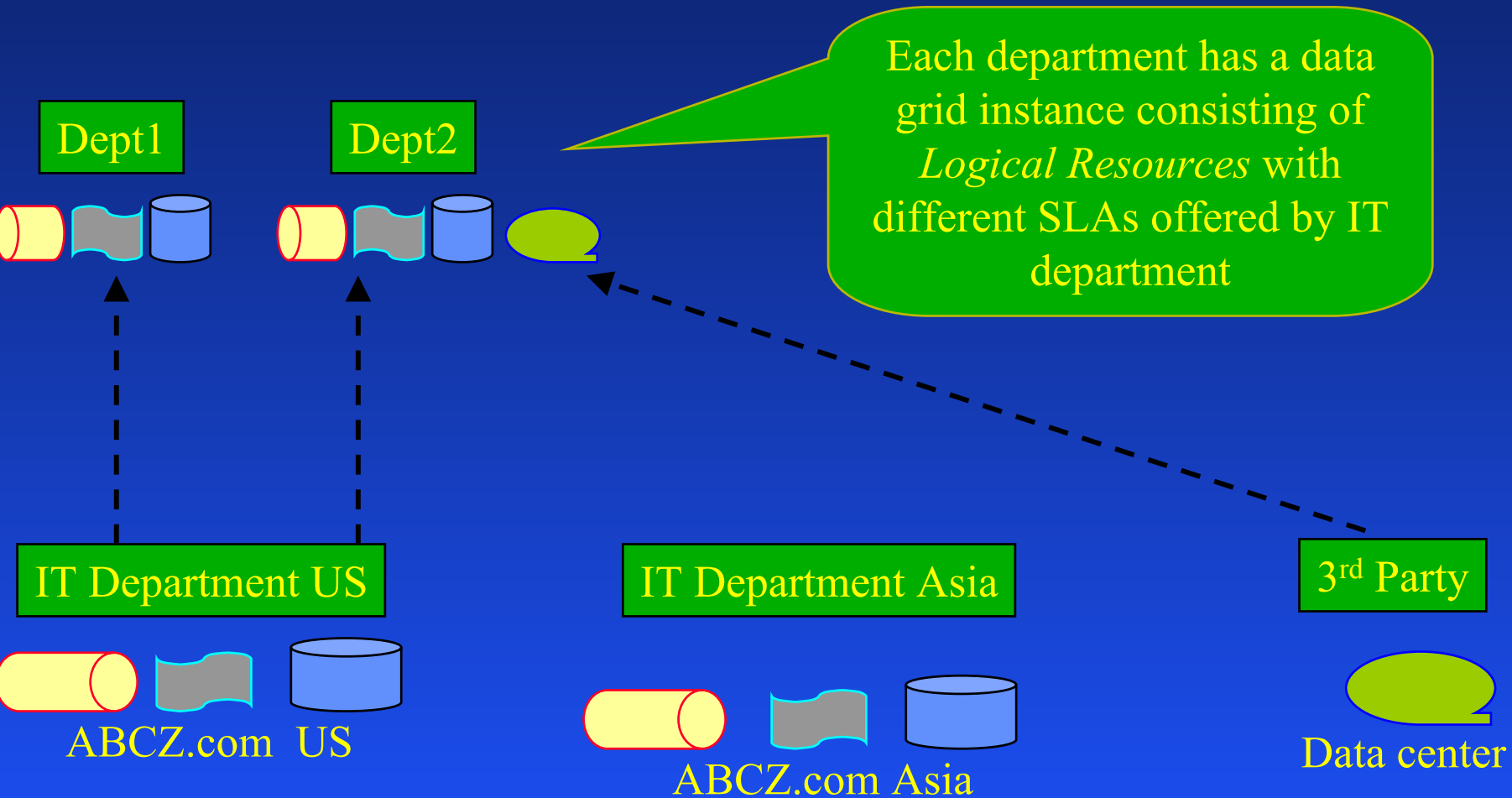
Data center



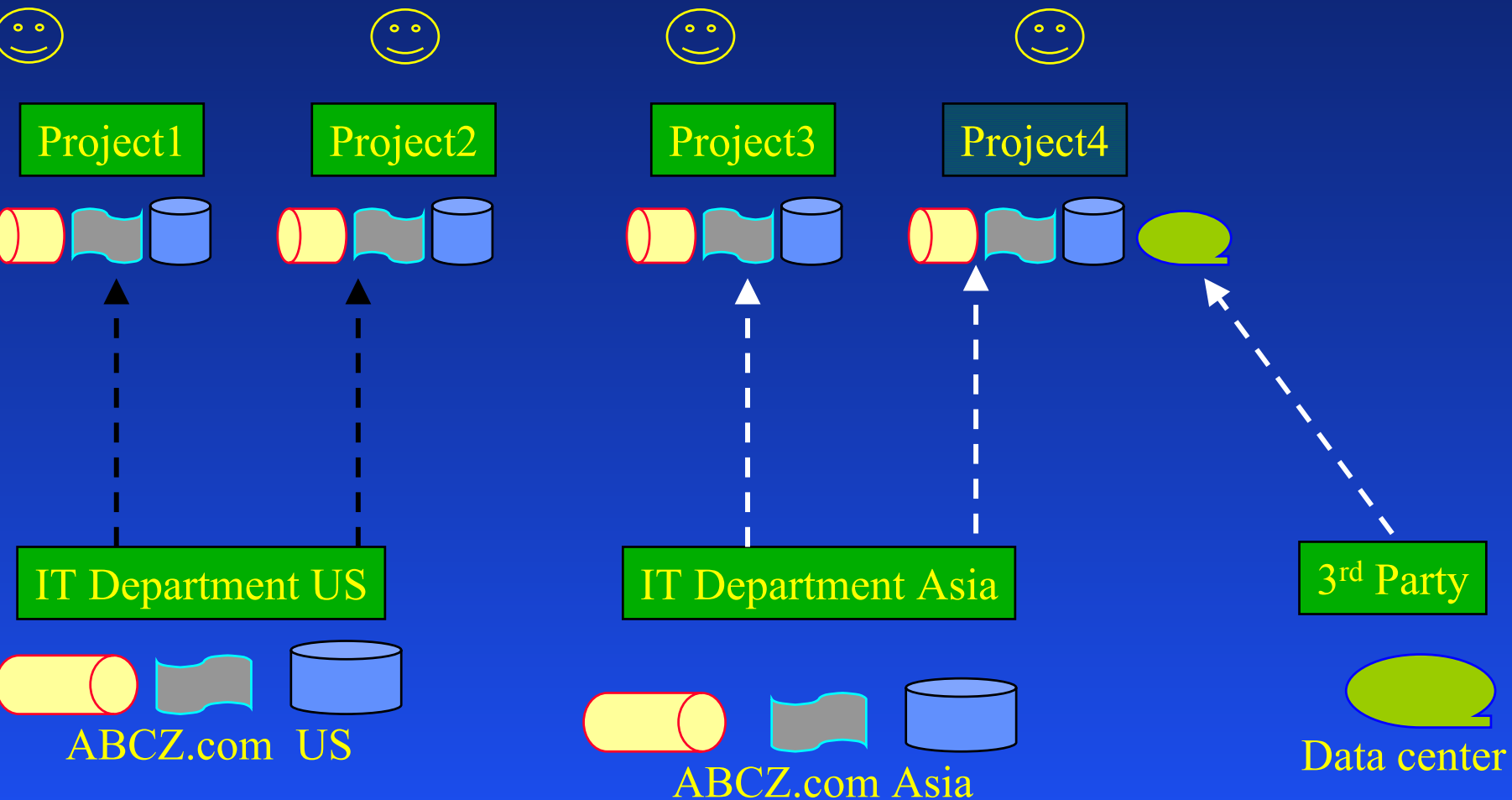
Data Grid (Enterprise Utility)



Data Grid (Enterprise Utility)



Data Grid (Enterprise Utility)



Long-run Processes in Data Grid

- Data Grid ILM
- Data Grid Triggers
- Data Gridflows



Data Grid ILM

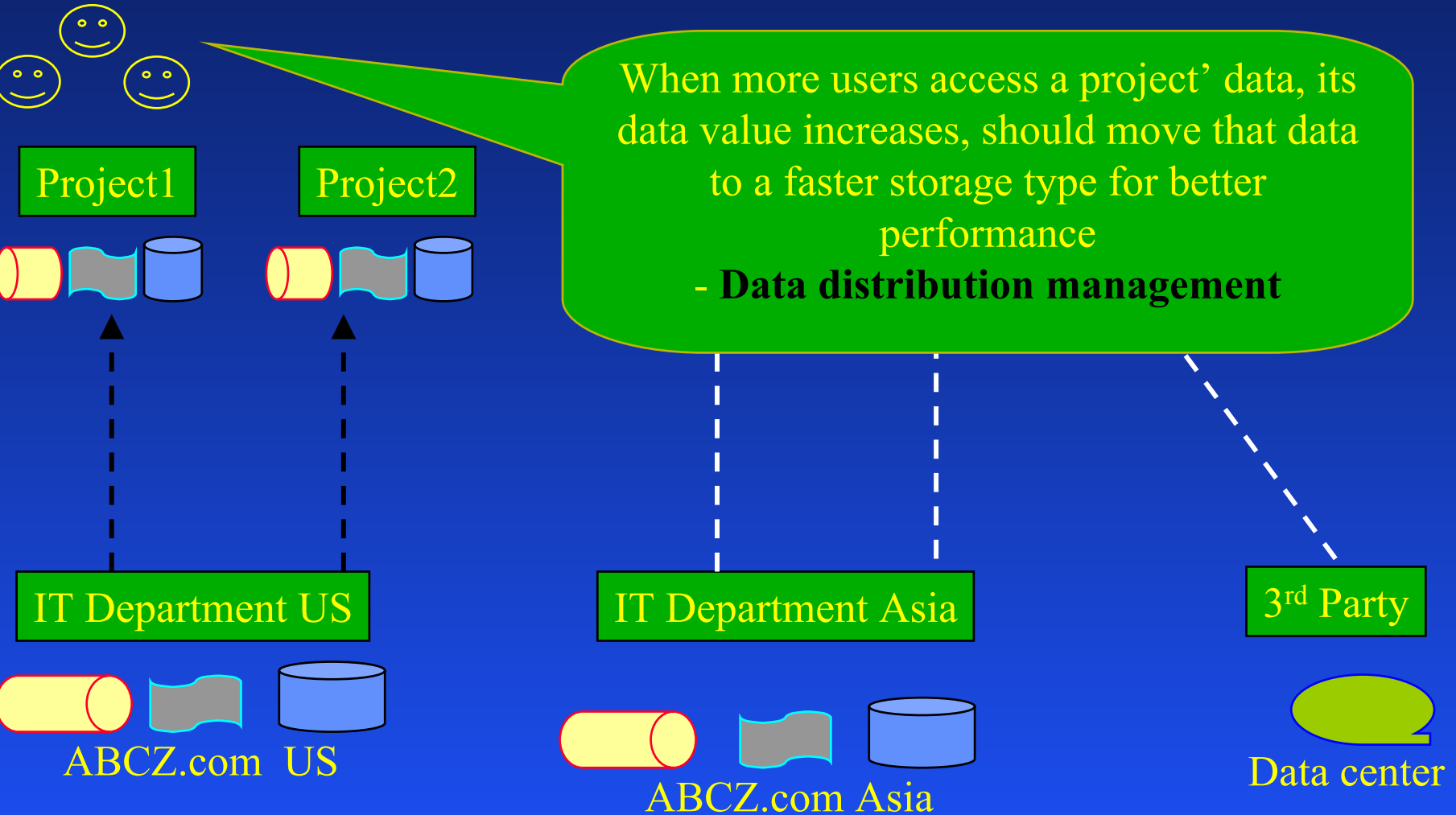


Change is Constant

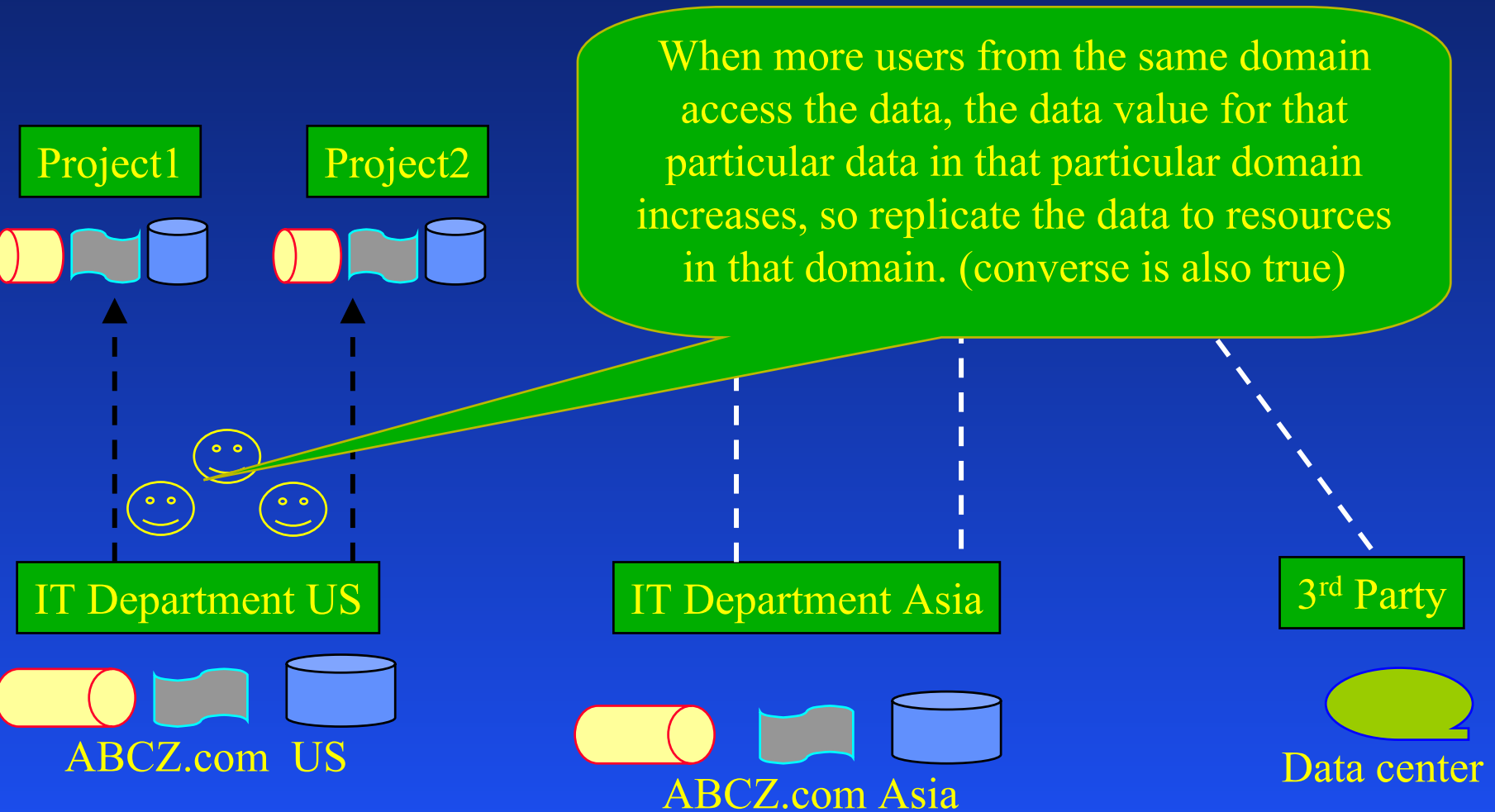
- **Changes in access patterns**
 - Based on number of users accessing a data
 - Domains which want to access data
- **Data Value**
 - The value of a data set (collections) for a particular domain is based on its business model and users' access patterns
 - Each domain will have a different importance based on its users and its role in a data grid



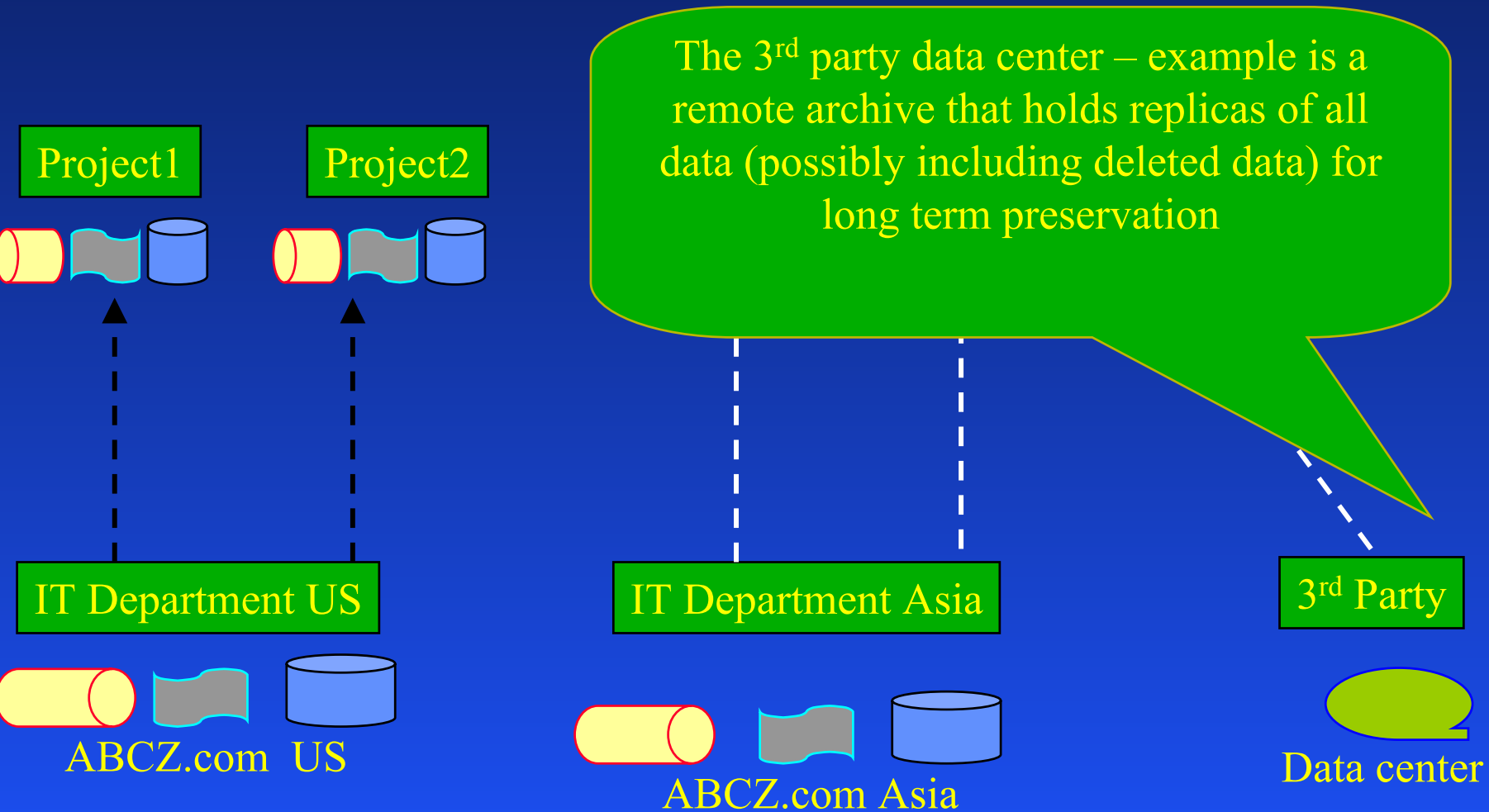
“Data Value” based on users



“Data Value” based on domain



“Data Value” based on role



Data Grid ILM

- ILM = Information Lifecycle Management
- Dynamic application of data placement and data retention policies (rules)
- Based on characterization of “business value of data” and storage cost
- HSM = Hierarchical Storage Management, based on “data time stamps”. ILM goes one step further to evaluate “data value”
- Applying this concept on Data Grid requires managing different business rules on different autonomous domains



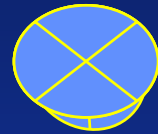
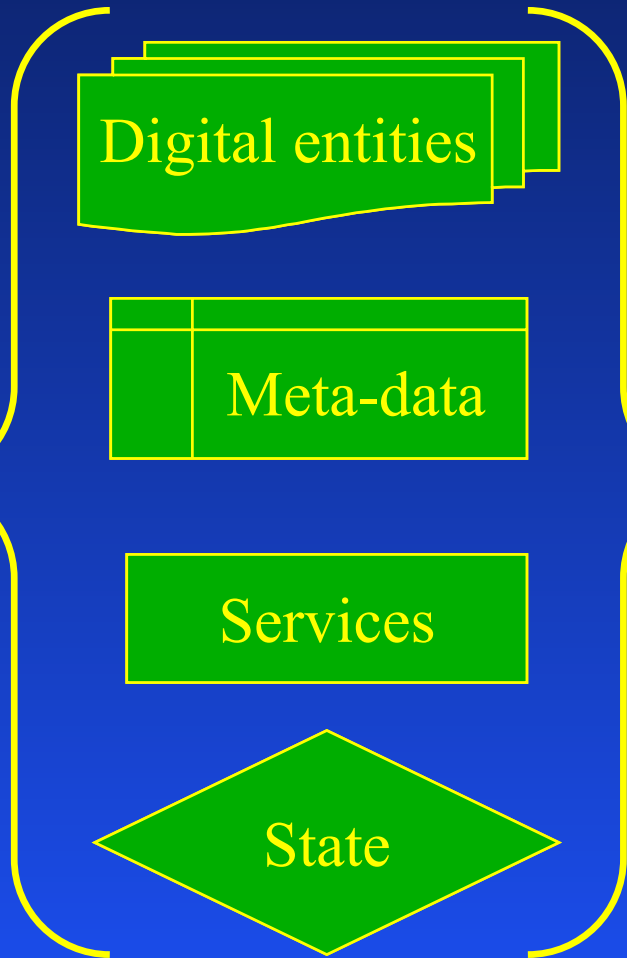
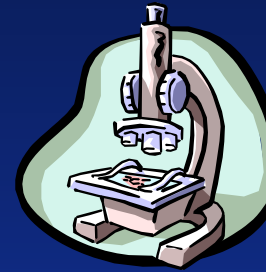
Data Grid Triggers



Data Grid Triggers

- **Similar to triggers in databases**
- **Based on ECA concepts**
 - Event
 - Condition
 - Action
- **Example**
 - Event = Insert new file in collection (“/ourProject/data”)
 - Condition = (color= “blue” && galaxy = “Andromedia”)
 - Action = Run (selectiveDataReplicator.dgl)

Data ⇔ Discovery

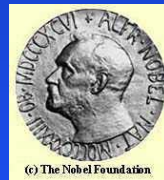


New data

updates relationships among data in collections

Services invoked to analyze new relationships

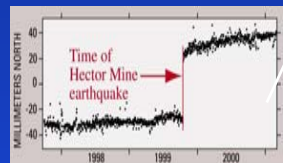
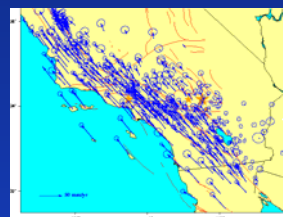
DGMS applications get notified of state updates



Data Gridflows



Gridflow in SCEC (data → information pipeline)



Metadata derivation

Ingest Data

Ingest Metadata

Determine analysis pipeline

Initiate automated analysis

Organize result data into distributed data grid collections

Use the optimal set of resources based on the task – on demand

Pipeline could be triggered by input at data source or by a data request from user

All gridflow activities stored for data flow provenance



Data Grid Language (DGL)



Data Grid Language

- **Requirement**

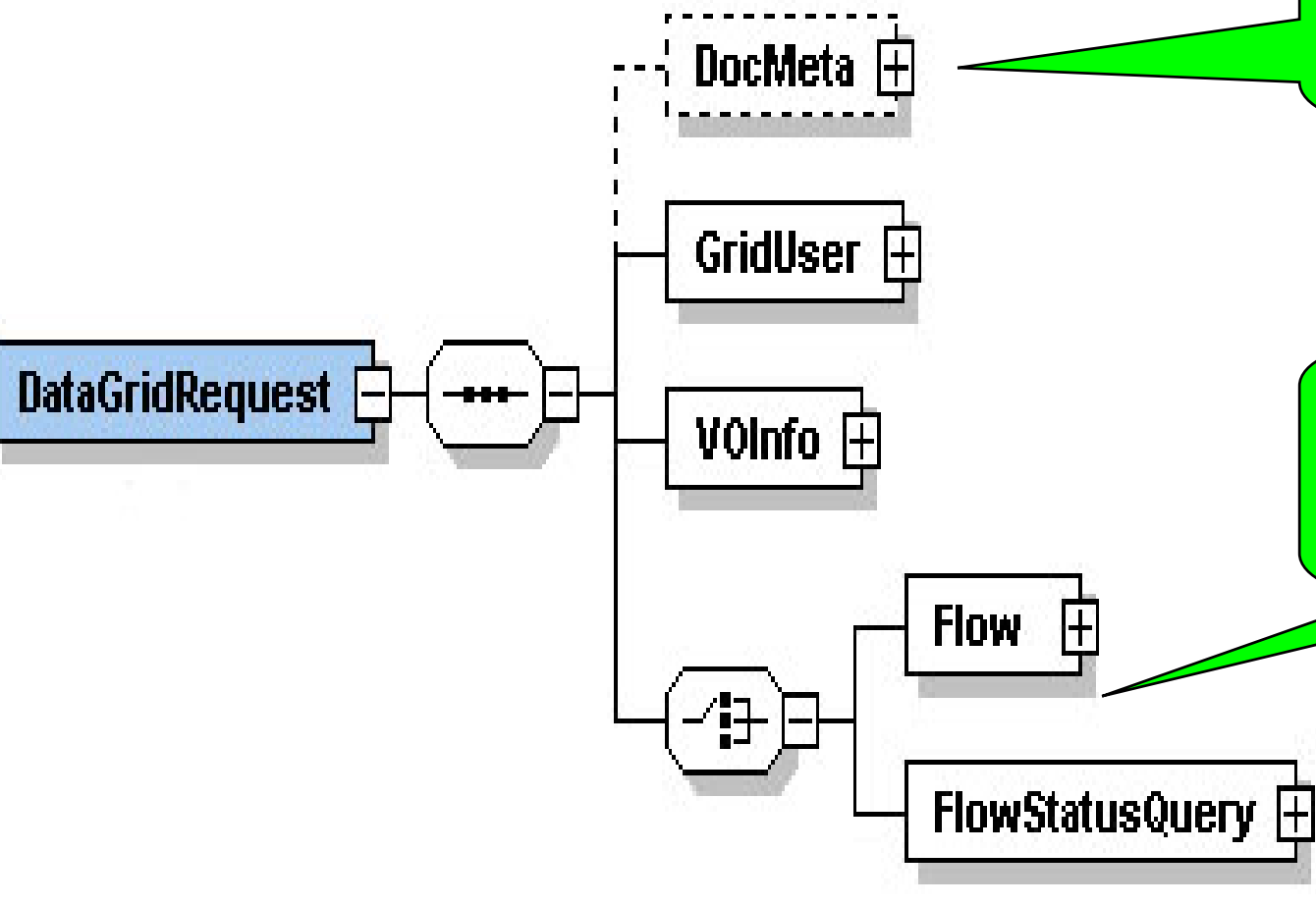
- Data Grid ILM process
 - The long run process that has to be run is described in DGL
- Data Grid Triggers
 - Action part of the ECA (Event-Condition-Action) logic
- Data Gridflows
 - Step by step execution of long run process on Data Grid

- **Analogy of SQL in relational databases**

- Long-run process procedures stored and executed in Data Grid it self
- Captures the “Infrastructure Execution Logic”



DGL Request



Annotations about the Data Grid Request

Can be either a Flow or a Status Query

DGL Requests (2 types)

- **Data Grid Flow**

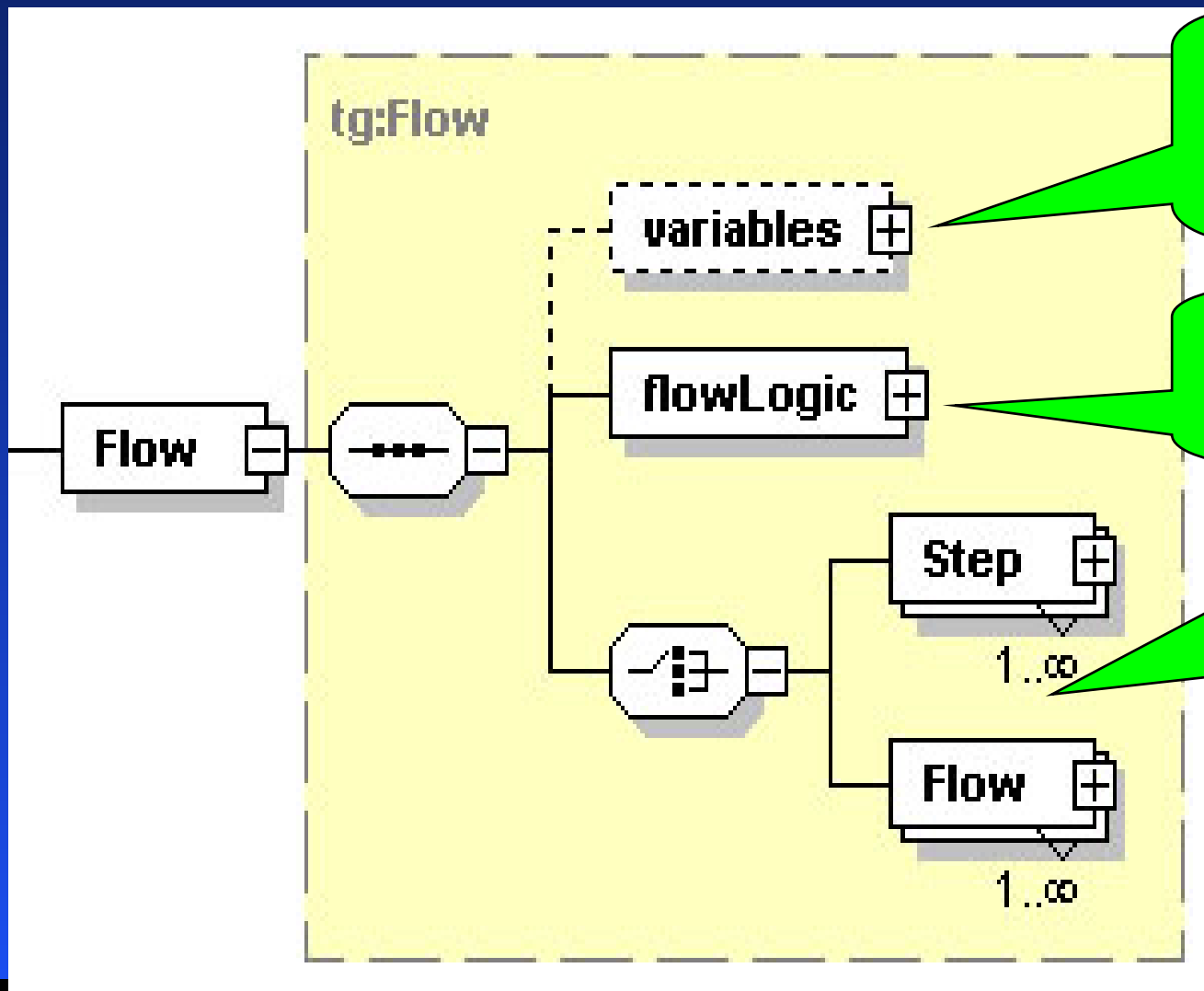
- An XML Structure that describes the execution logic, associated procedural rules and DGL variables. Can be synchronous or asynchronous flow

- **Status Query**

- An XML Structure used to query the execution status any gridflow or a sub-flow at any granular level. Status Queries can be made for both synchronous and asynchronous flows



Flow

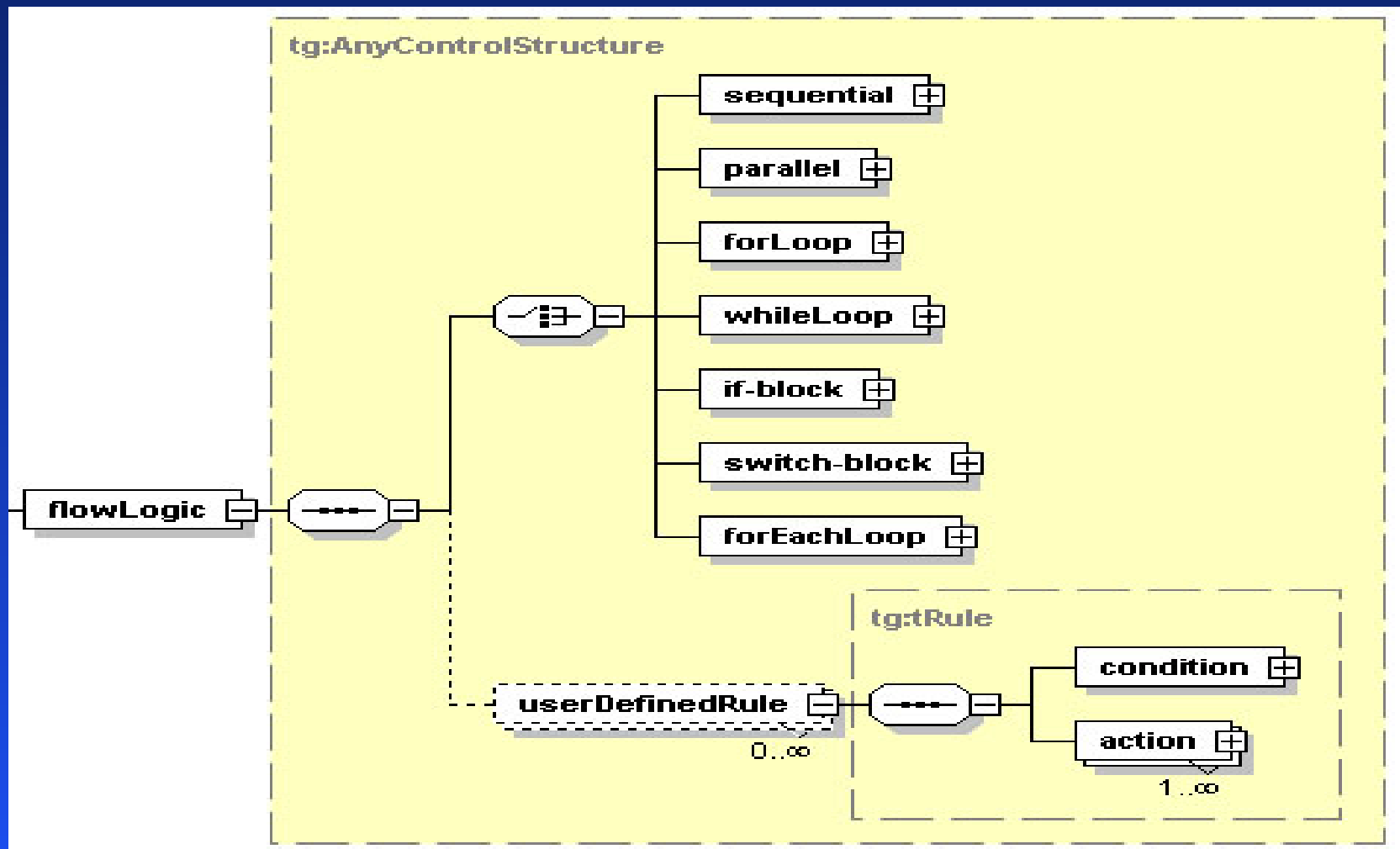


Scoped Variables that can control the flow

Logic used by the sub-members

Sub-members that are the real execution statements

Flow Logic (How a flow executes)



```
...
<userDefinedRule name="beforeEntry">
<condition>
<simpleQuery>$numVar == 1</simpleQuery>
</condition>

<action name="true">
<actionString>SET var1 = 1</actionString>
</action>
<action name="true">
<actionString>SET var2 = "foo"</actionString>
</action>
<action name="false">
<actionString>SET var1 = 0</actionString>
</action>
</userDefinedRule>
...
```

What is SRB Matrix?

- **Matrix provides a Web Service interface to the SRB**
 - Web Service expressed in Data Grid Language, communicated using SOAP
- **Matrix provides a Service Oriented Architecture for Data Grid or Digital Library Clients**
 - Asynchronous end-user applications
 - Long run operations presented to users as portlets
- **Data Grid Automation and ILM**
 - File Triggers on unstructured data
 - Automated movement or management of data



Matrix Gridflow Server Architecture

JAXM Wrapper

WSDL Description

Event Publish
Subscribe,
Notification

JMS Messaging
Interface

SOAP Service for Matrix Clients

Matrix Data Grid Request Processor

Sangam P2P Gridflow Broker and Protocols

Transaction Handler

Status Query Handler

Workflow Query Processor

Flow Handler and
Execution Manager

XQuery
Processor

ECA rules
Handler

Gridflow Meta data
Manager

Matrix Agent Abstraction

Persistence (Store)
Abstraction

SDSC SRB
Agents

Other SDSC
Data Services

Agents for java,
WSDL and other
grid executables

JDBC

In Memory
Store



Conclusion

- **Data Grids are evolving**
- **Data Grid Automation of long-run processes essential**
- **Need a language for Data Grid Automation**
- **Data Grid Language is one such effort as part of the SRB Matrix Project**
- **Open source project for anyone to use (or join)**
- **talk2matrix@sdsc.edu (or arun@sdsc.edu)**

