Grid-based processing of high-volume meteorological data sets





BI

Outline

Introduction

- Energy meteorology, WISENT
- Challenges
 - Parallel processing, Data transfers

Utilizing Grid technologies

- Condor, Globus Toolkit 4, UNICORE
- Future Work



Energy meteorology



- Research on the influence of weather and climate on the transformation, transport, and utilization of energy from renewable energy sources
 - Forecast models of energy production
 - Finding optimal locations for power plants
- Interdisciplinary field of research (meteorology, physics, engineering, ...)
- Large and heterogeneous data sources (satellites, earth stations, ...)
- Compute-intensive applications on high-volume data sets



wisent.d-grid.de



- German e-Science project in the domain of energy meteorology
- Associated with D-Grid (German Grid initiative)
- Started in October 2005 (duration for 3 years)
- Distributed resources (CPUs, data storages)
- ~ 1 TB new data per month (increasing)
 - Mostly raw or post-processed satellite images
 - Archived in the "Data and Information Management System" (DIMS)
 - ~ 300 TB; planned extension to 3 PB
- Objective: Build Grid infrastructure based on these resources to support (large) data transfers and distributed processing





Parallelization



Status:

- Most applications run on one single machine
- Parallelization is achieved with Parallel virtual machine (PVM) / ppmake
- Most applications can be parallelized at data level

Objective:

- Parallelization of each application utilizing existing CPU resources (desktop PCs, clusters, etc)
- Recognizing user activity on desktop PCs
- Authentication and authorization
- Easy access to computing resources via the Grid infrastructure

Data transfers



Status:

- Multiple (~ 100) data-transfers per day (periodically/on demand)
- Size ranges from a few kilobytes up to several hundred megabytes
- Number and size will increase in future
- Often FTP-based transport with manual error recovery in case of failures

• Objective:

- Security
 - Encrypted data transfers
 - Authentication and authorization
- Reliable data transfers with automatic recovery
- Monitoring for accounting and billing
- Easy initiation of data transfers within the Grid infrastructure

Bottom-up approach







Utilizing Grid technologies

S

Page 8

Parallelization scenario





Utilizing Grid technologies

Page 9

Guido Scherp <guido.scherp@offis.de> Business Information Management

Condor



Parallelization of each application utilizing existing CPU resources

- "Cycle scavenging" using idle-periods of computational resources
- Very suitable for applications using data parallelization
- Rudimentary support for MPI
- Scheduling strategy sufficient to adequately utilize resources

Authentication and authorization

- Possible but not tested
- Not necessarily needed at Intra-Grid level
- Sandbox approach is sufficient

Recognizing user activity on desktop PCs

- Tracking user's activity (mouse, keyboard, etc)
- Migration of jobs to other nodes on user's demand

Easy access to computing resources via the Grid infrastructure

Text-based interface is not very user-friendly?



Pros:

- Good approach for pooling CPUs at Intra-Grid level
- Capable of construction of "desktop-Grids"
- Solid documentation and long development history
- Wide user base

Cons:

- No Open Source project
- Large number of configuration settings (but well-documented)
- Demands on network connectivity

Data transfer scenario







Utilizing Grid technologies

Page 12

Guido Scherp <guido.scherp@offis.de> Business Information Management

Globus Toolkit 4



Security

- Certificates based on X.509 (SSO) for authentication
- Community Authorization Service (CAS) for authorization (not evaluated)

(Reliable) Data transfers

- GridFTP as enhanced FTP
 - GridFTP uses no data channel encryption per default
 - Encryption modes "Safe" and "Private" reduce data throughput
 - New port assigned for each data channel
- Reliable File Transfer (RFT)
 - Does not support GridFTP-based encryption
- Monitoring for accounting and billing
 - Monitoring and discovering service (MDS) (not evaluated)
- Easy initiation of data transfers within the Grid infrastructure
 - Text-based interface is not very user-friendly?

Globus Toolkit 4



Pros:

- Support of WSRF and most services proposed in OGSA
- Comprehensive data services
- Interoperability with Condor
- Cons:
 - Dynamic port assignment conflicts with current firewall policies
 - Open a whole port range (only temporary solution)
 - DLR develops Application Level Gateway (Proxy in the DMZ)
 - Use of encryption in RFT is currently not possible



UNICORE

Currently little experience with UNICORE 5 and evaluation is ongoing

Pros:

- Complete encryption of communication
- Gateway uses only one port
 - Location of Gateway in DMZ could possible
- Graphical client for job submission and monitoring
 - Possible automation?
- Workflow language for modelling process chains
 - Which capabilities?
- UNICORE 6 has promising extensions
 - Support of Grid standards

Cons:

- Execution of a workflow element must be assigned to a specific VSite
 - Optional dynamic VSite-selection
- UPL-based data transfers are not sufficient for large data sets

Future work



- UNICORE (6) seems very interesting for fulfilling many requirements in WISENT
- We still need more experience with UNICORE
- Evaluation of other Grid middleware as gLite, Sun N1 Grid Engine, etc
- More investigation of interoperability





Questions?

