# Middleware Framework
# for HG2C Project

## Dec 13, 2006

Jaeyoung Choi

Soongsil University, Seoul Korea

choi@ssu.ac.kr

# Middleware Framework for HG2C Project
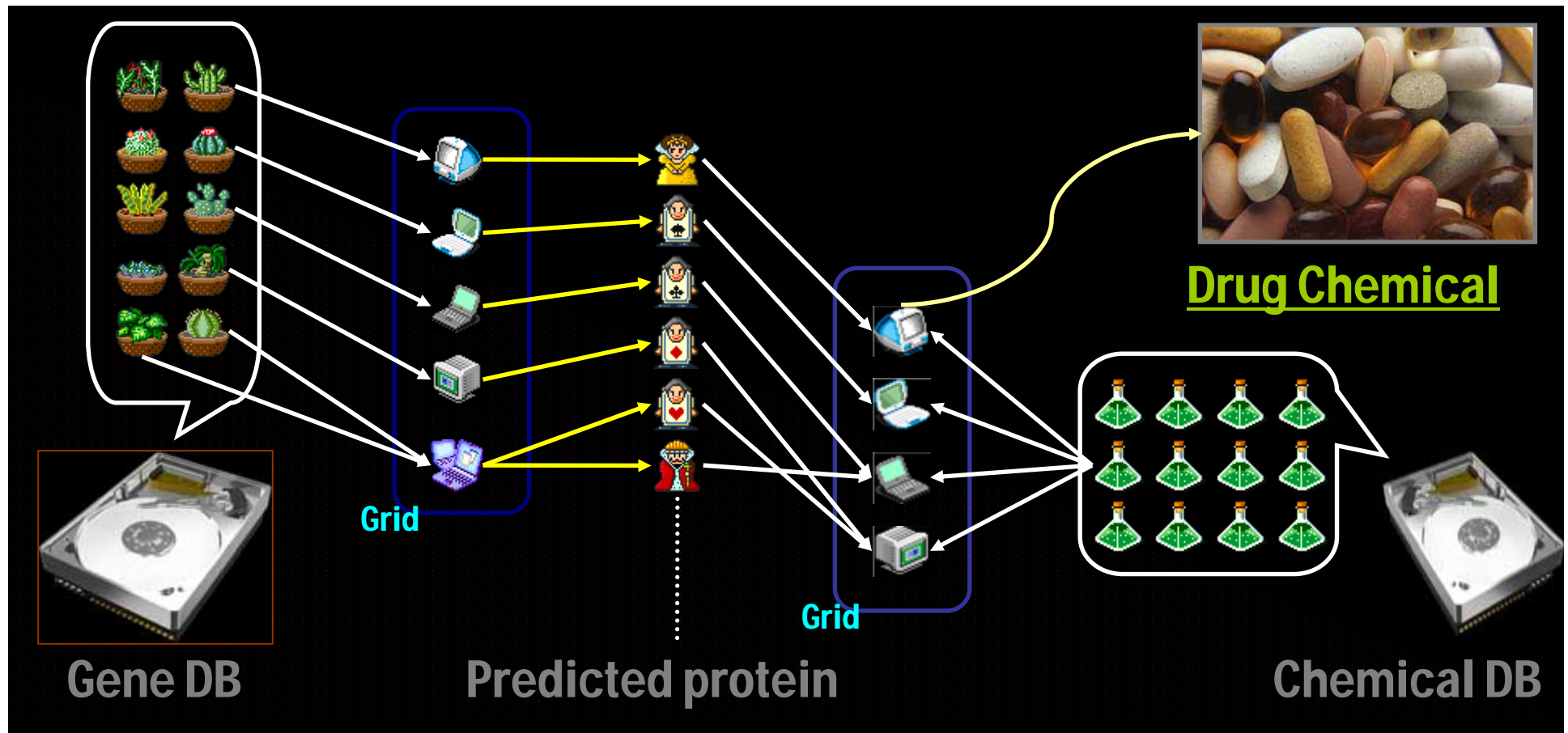
December 2006

Jaeyoung Choi

Soongsil University, Seoul Korea

choi@ssu.ac.kr
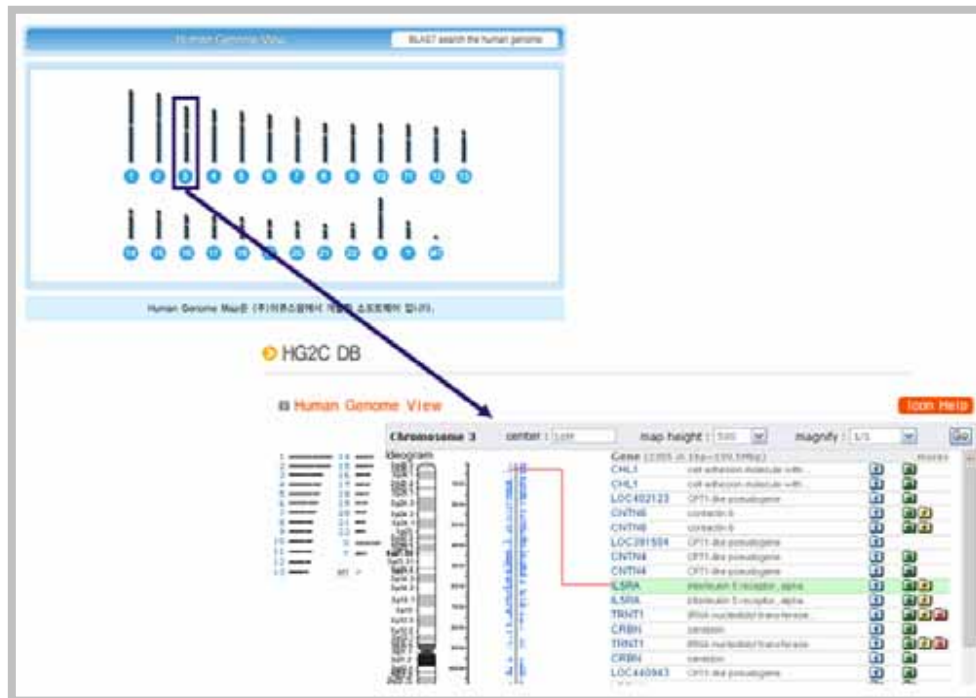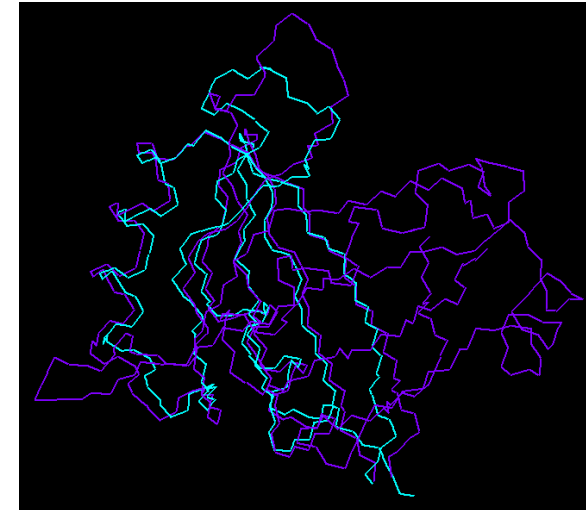
# HG2C (Human Genomes to Chemicals)

# Dream of HG2C



**Drug Chemical**

Grid

Grid

Gene DB          Predicted protein          Chemical DB

# HG2C (1)

- ❑ **Genome**
  - ◆ **Literature Information DB system (PubLink™)**

    **Genes - Proteins - Chemicals**
  - ◆ **Chromosome/Gene/DNA sequence relation**

# HG2C (2)

- ❏ Gene
  - ◆ From DNA sequence
  - ◆ BLAST homology search
  - ◆ Comparative modelling from sequence
  - ◆ Protein structure prediction

# HG2C (3)

- Protein (IDPro$^{TM}$, PharmoMap$^{TM}$)
  - ◆ Active site searching from protein structure
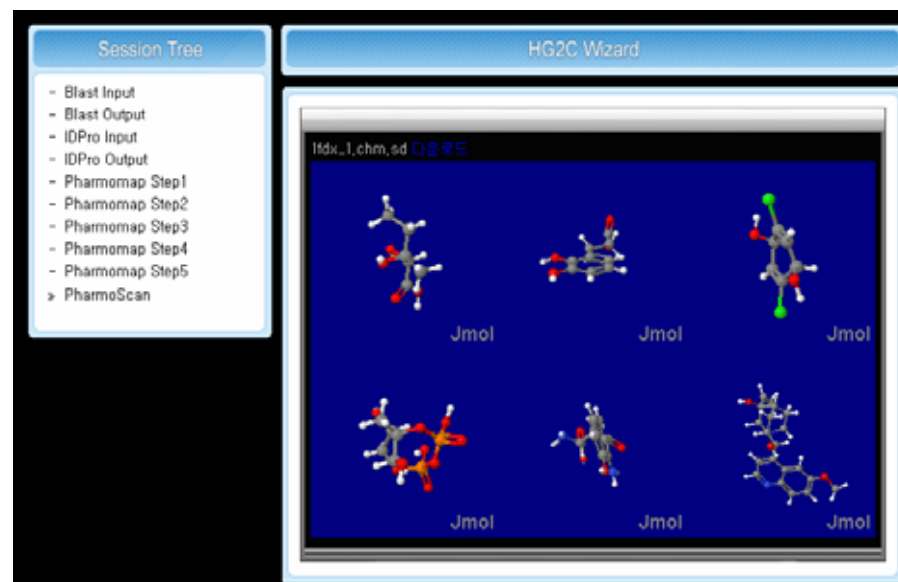  - ◆ Interaction feature searching from active sites

# HG2C (4)

- Chemicals (PharmoScan$^{TM}$)
  - Combinations of interaction features
  - Chemical DB scanning : interaction combination as scanning query
  - Scanned results : virtual / known compounds

# HG2C (5)

- Database
  - Iterative calculations on all the human genes
  - Browsing result DB :
    - Gene-Protein-Chemical relations :
      - ✓ key information for drug discovery

  - Result chemical analysis for drug utility and chemical attributes :
    - IDBManage$^{TM}$
    - JDDW (IDPro$^{TM}$, IDPharmo$^{TM}$, IDChemo$^{TM}$)
    - PubLink$^{TM}$

# HG2C (6)

❑ **Massive Calculation**

- ◆ For all the genomes with a huge chemical DB

- ◆ Choices of :
  - Protein model candidates
  - Active site candidates
  - Virtual screening query candidates

- ◆ Huge size of computing requirement :
  - Grid computing: dynamic computing resource management with grid application utility MAGE support

# HG2C Portal (1)

- Completion of Human Genome Project
- Theoretical prediction of chemical relations among gene functions and chemicals
- Complete independence of the each simulation
- Based on virtual computing:
  - Grid computing technology : MSF, MAGE
- http://www.hg2c.org

# HG2C Portal (2)

# How to explore? - Sequentially

- Sequential processing: gene1, gene2, gene3...
- Automatic batch processing on Grid environment
- Using work flow engine MSF features of batch running, data file transferring and storing

# We have explored … (1)

❑ **Using 1 ~ 21 CPUs (Sep. 1st. 2006)**

  ◆ 7,191 genes calculated

  ◆ 2,264 meaningful protein models achieved (31%)

  ◆ 874 meaningful chemical sets achieved (12%)

# We have explored … (2)

❑ **Chromosome dimension vs. Structural significance**

  ◆ Independent protein homology patterns along the distribution on chromosome (StdDev: 9%)

  ◆ Independent structural complexity patterns along the distribution on chromosome (StdDev: 5%)



Scan rate

# Implementation Features

- HG2C DB opened  http://www.hg2c.org
  - Human chromosome : Gene – Protein – Chemical

- Interactive GUI opened

- MSF integrated in modular usage level

- MAGE being integrated

- MSF package released http://sourceforge.net/

- MAGE package released http://sourceforge.net/

# How to explore? - Directionally

- ❑ **Directional processing**
  - ◆ Interactive GUI for researcher decision :
    - ▪ Protein models, active sites VS queries, etc.
  - ◆ Conditional decision module :
    - ▪ Additional decision program

- ❑ **MSF utility**
  - ◆ User workspace management for multiple researcher
  - ◆ Batch processing management for various procedure

# MSF and MAGE

# HG2C System Architecture

| | |
|---|---|
| **Application Research** | Application SW |
| **Human Genome to Chemome (HG2C)** | |
| **BT Portal (Web)** | High level : Application Support |
| **Service Component Framework (MAGE)** | |
| Visualization — **MSF 2.0** — Steering, Workflow, PSE toolkit | |
| Global scheduler, Adaptive API, Monitoring, Streaming, Text Search, Data Mining | Mid level : Fundamental Service |
| **Globus Toolkit (GT3, GT4)** | Low level : Resource Management |
| **Infrastructure (Cluster, Network, ...)** | Infrastructure |

# Meta Services

- Provide reusable and adaptable workflow environments

- Define a part of a workflow as a new service
  - Workflow instance can be declared as a workflow unit in the service description
  - By overriding some attributes of a workflow unit,
    - Pass parameters of a service to the workflow's attributes
    - Setup service specific information
  - The new service can be wrapped to a Web service or a Grid service, therefore it can be easily reused

- Manage service specific information
  - Restrict resources to allocate a specific service (user's preference and/or organization's policy)
  - Schedule jobs with priority

# Meta Services Framework (1)

- **Meta Services Framework (MSF)**
  - MSF is a workflow system for Bio Grid portal
  - Users can easily compose a DAG-based workflow
    using legacy applications such as a BLAST
  - Schedule user's workflows on Grid environments
  - Provide reusable workflows using Meta Services
  - Users can compose services, flows, and tasks using XML
  - Can be easily installed and configured
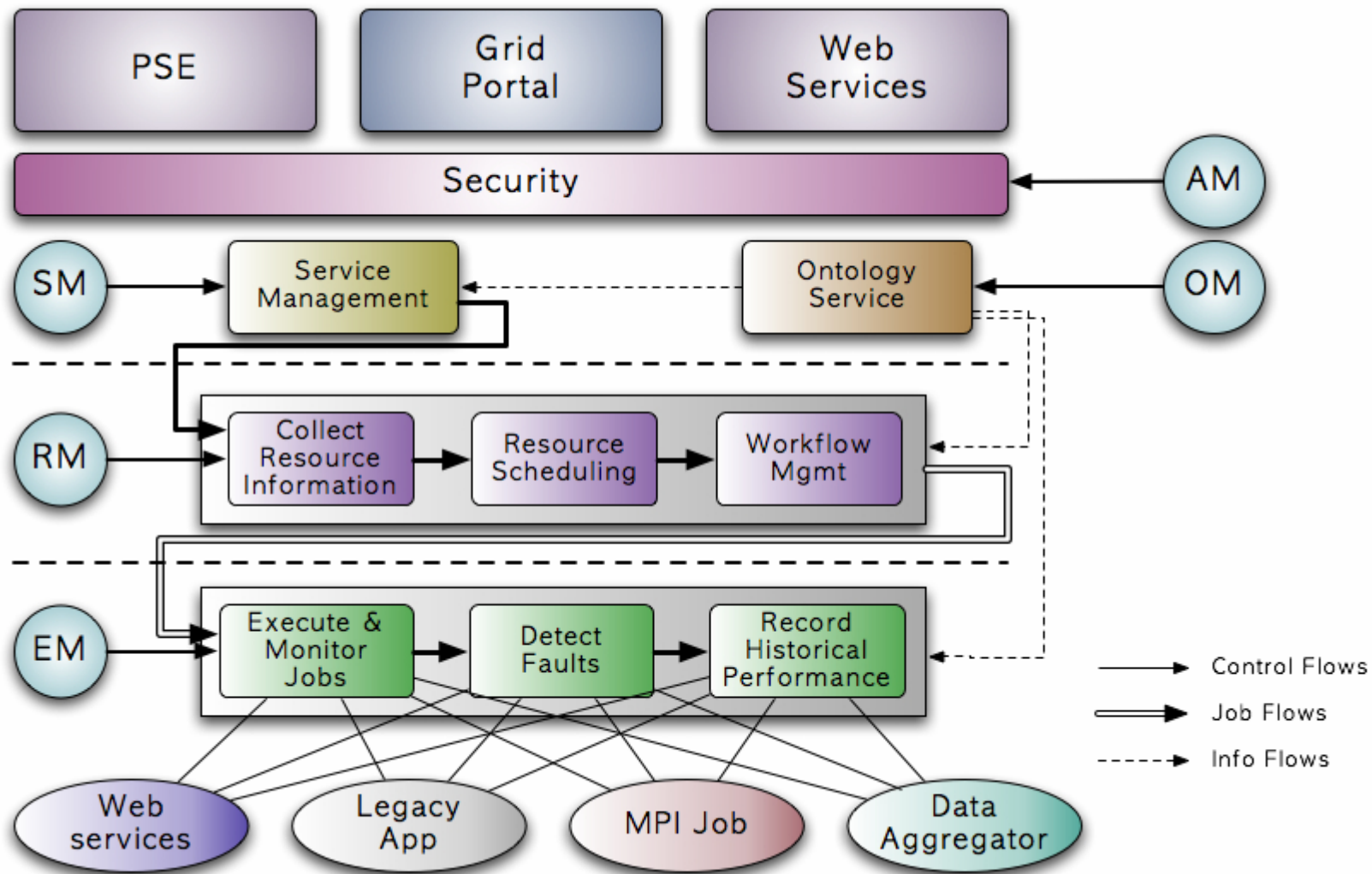
# Meta Services Framework (2)

- ❑ MSF Workflow model
  - ◆ Workflow divided into three layers
    - service layer, flow layer, and task layer
  - ◆ Increase reusability of workflow
  - ◆ Service layer is specially designed using Meta services concept

- ❑ MSF Middleware Architecture
  - ◆ Consist of five modules (SM, RM, EM, AM, OM)
  - ◆ Three agent modules process each layer of the workflow model (SM - service, RM - flow, EM - task)
  - ◆ AM manages authentication and access control
  - ◆ OM maintains XML description

# Five Agents in the Architecture

- **SM (Service Manager) - service**
  - Manage meta services and converts a meta service to a workflow

- **RM (Resource Manager) - flow**
  - Collect resource information, and allocating resources

- **EM (Execution Manager) - task**
  - Launch & monitor (workflow) jobs, detect faults, and collect results and performance data

- **AM (Access Manager)**
  - User authentication, environment setup, and a job submission service

- **OM (Ontology Manager)**
  - Manage ontology of service, flow, and task

# Middleware Architecture

# Operations defined in MSF

| Agent | Type | Operation |
|---|---|---|
| AM | USER | AUTH_USER |
| | | PROXY_INIT |
| | | PROXY_UPDATE |
| | | PROXY_INFO |
| | SYSTEM | CHECK_PRIV |
| | ADMIN | ADD_USER |
| SM | USER | REQUEST_SERVICE |
| | | FORWARD_FLOW |
| | SYSTEM | PROCESS_POLICY |
| | | ANALYZE_META_SERIVCE |
| RM | USER | EXECUTE_FLOW |
| | | CANCEL_FLOW |
| | | GET_FLOW_STATUS |
| | | GET_FLOW_QSTATUS |
| | SYSTEM | ANALYZE_FLOW |
| | ADMIN | REGISTER_RM_INFO |

| Agent | Type | Operation |
|---|---|---|
| EM | USER | EXECUTE_TASK |
| | | CANCEL_TASK |
| | | SET_JOB_PRIORITY |
| | | GET_QSTATUS |
| | | GET_NODE_INFO |
| | SYSTEM | ANALYZE_TASK |
| | ADMIN | REGISTER_TO_RM |
| | | SET_PE |
| OM | USER | RETRIEVE_DESCRIPTION |
| | | STORE_DESCRIPTION |
| | | SEARCH_DESCRIPTION |
| | | GET_SERVICE_LIST |
| | | GET_FLOW_LIST |
| | | GET_TASK_LIST |

AsiaGrid,  Dec 2006

# MAGE features

- Provide API for easy development of Grid application
- Provide transparency to end-users and developers
  - Protocol transparency
  - Running location transparency
  - Message interpret transparency
- Provide layered architecture for easy replacement
- Mobility for each agent

# MAGE architecture



AsiaGrid, Dec 2006

# Communication Layer

- To provide easy replacement of the communication protocol without affecting other layers
  - Administrator can select suitable protocol components before running application
  - Tasks and message interpretation does not affected by changing of communication protocol

# Interpreter Layer

- Interpret received messages and deliver to the appropriate task agents
- Two basic query components implemented
    - Monitor Query: use SQL's SELECT-like statement
    - Table Query: use name=value pair table

# Task Agent Management Layer

- Control the life cycle of task agents
- Provide feature of installation from remote
  - Base function for mobile agents
  - Provide function of stop the job and resume at another node

# MSF Management using MAGE



Grid Management Service

MAGE

Third Party Modules

Data Aggregator

MSF Pkg module

RM module

EM module

AM module

SM module

OM module

MSF Setup module

Distribute & Deploy MSF module    Setup MSF Modules to cooperate

# Conclusion

- MSF contains essential functionalities for Grid portals
  - workflow, service interface, job distribution, and parameter scheduling

- Distributing and deploying MSF modules using MAGE increases reconfigurability and adaptability to MSF modules
  - MSF modules can used as agents
  - MAGE environments can provides various services to MSF modules
    - → Resource information, System environments, ..

# HG2C Portal

# Meta Services Framework Demo (1)

# Meta Services Framework Demo (2)

# HG2C Pubmap

# HG2C References

- HG2C Homepage
  - http://www.hg2c.org


- Open sources of MSF & MAGE
  - http:// www.sourceforge.net/projects/mage4ubi
  - http:// www.sourceforge.net/projects/msf

# Q & A

AsiaGrid, Dec 2006