# Outlook on Existing and Future Grid Programming Paradigms
## EchoGrid Workshop Athens 2008

Bernhard Schott          bschott@platform.com
Guillaume Mecheneau   gmecheneau@platform.com
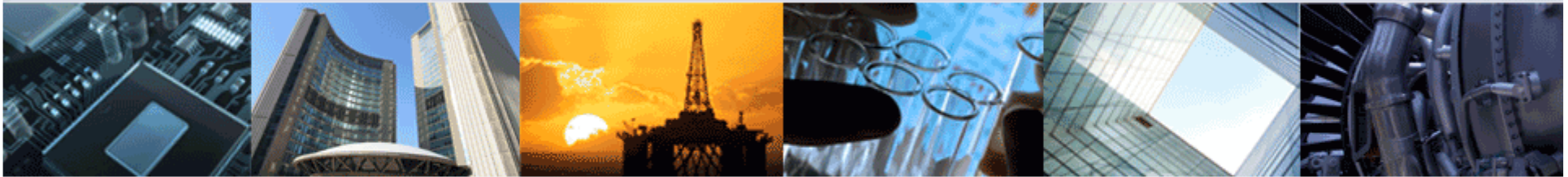Platform EU-Research Team, Paris

Date: 9th June 2008

Agenda

- Platform Computing

- (Grid) Programming Paradigms

- State of the Art and possible extensions
  3 Examples:

  1. Grid embracing the Application: Platform LSF & EGO

  2. SOA application integration: Platform Symphony

  3. On-demand resource acquisition and formation of application specific Grids: QosCosGrid

- Summary: Trends in Grid Programming

# Platform™ Introduction

# Platform Computing

*Platform is a pioneer and the global leader in High Performance Computing infrastructure software, delivering integrated software solutions that enable organizations to improve time-to-results and reduce computing costs.*

## Over 2200 Customers Worldwide

- Electronics, Financial Services, Manufacturing, Life Sciences, Oil & Gas, Government, Universities & Research, Telco …

## Recognized Leader in HPC, Cluster and Grid Computing

- 16 years global experience
- Worldwide offices, resellers and partners
- 24x7 follow the sun support and services

## Growing & Profitable since inception in 1992

- Self-funded
- No debt; money in bank

*Founded by Songnian Zhou, Bing Wu, Jingwen Wang*

**Office Locations**

*North America*
Toronto (HQ)
San Jose
Washington
Detroit
Los Angeles
Boston
New York

*International*
China
Japan
Korea
UK
Germany
France

**VARs**

| | |
|---|---|
| U.S. | South Africa |
| Italy | Sweden |
| Israel | Turkey |
| Germany | India |
| Spain | Malaysia |
| Korea | Thailand |
| Taiwan | Australia |
| Singapore | Austria |
| Japan | France |
| U.K. | China |
| The Netherlands | U.A.E. |
| Poland | Portugal |
| Brunei | |

# Big Companies Trust Us

**Platform**

## Financial Services
- BNP Paribas
- Citigroup
- Deutsche Bank
- Fortis
- HSBC
- JP Morgan Chase
- Lehman
- Mizuho Financial
- MUFG
- Prudential
- Société Générale
- Sal Oppenheim

## Electronics
- AMD
- ARM
- ATI
- Broadcom
- Cadence
- Cisco
- HP
- IBM
- Motorola
- NVIDIA
- Qualcomm
- Samsung
- ST Micro
- Synopsys
- TI
- Toshiba

## Industrial Manufacturing
- Airbus
- BMW
- Boeing
- Bombardier
- British Aerospace
- Daimler & Chrysler
- GM
- Lockheed Martin
- Pratt & Whitney
- Toyota
- Volkswagen
- Xi'an Aircraft Design

## Life Sciences
- AstraZeneca
- Bristol Myers-Squibb
- Celera
- Dupont
- GSK
- Johnson & Johnson
- Merck
- Novartis
- Novo-Nordisk
- Pfizer
- Wellcome Trust Sanger Institute
- Wyeth

## Government & Research
- ASCI
- CERN
- CINECA
- DoD, US
- DoE, US
- ENEA
- ETH
- Fleet Numeric
- GSI
- INFN
- MaxPlanck
- SSC, China
- TACC
- TU Dresden
- Univ Tokyo

## Other Business
- Bell Canada
- Cablevision
- Deutsche Telekom
- Ebay
- Starwood Hotels
- Telecom Italia
- Telefonica
- Sprint
- GE
- IRI
- Cadbury Schweppes

# Partite Ecosystem

**Platform**

## Partner Ecosystem

**Strategic Partners**

intel — DELL — Microsoft — redhat. — SAS — hp invent — IBM

**Premier Partners**

GEMSTONE — Sun microsystems — NICE — The MathWorks — Mentor Graphics — macrovision — cadence — ANSYS — Apple — ARROW — AMD — vmware — CSC — Novell — the BioTeam — CALYPSO — BULL — MSC Software — SCS — CISCO — AVNET — DS SIMULIA — sgi — WIPRO Applying Thought — Schlumberger — Vinetech Utility Computing Service Provider — SC science + computing — LSTC Livermore Software Technology Corp. — SYNOPSYS — ORACLE — TANGOSOL DATA GRID SOLUTIONS — NEC — Satyam What Business Demands. — GIGASPACES

**Select Partners**

Apache DESIGN SOLUTIONS — WinterLogic — D'GIPRO DESIGN ACTIVATION & MARKETING (PVT.) LTD. — T···Systems··· — ATopTech — SCHRÖDINGER. — 安託集團 ATOZ GROUP — ATHENA DESIGN — electric cloud — carus — GNS Systems — HPC SYSTEMS — SCS LIVING TECHNOLOGY — ITOCHU — CoWare The ESL Design Leader — PTC Data Management Solutions — KIMOTION TECHNOLOGIES — COMLINK — LINUX NETWORX — QuIC — SiSoft — A SPEED Software Formerly Powertel — immixTechnology — elim — NANGATE — Engineous software — QuIC — a·s — TotalView TECHNOLOGIES — solido DESIGN AUTOMATION — CRAY — Voltaire — Proline Integrated Intelligence — denali — EVERGRID ALWAYS ON. ALWAYS OPTIMIZED. — accelrys — FUJITSU — it peers optimization services — INCAT — BCS — SONNET — gridwise tech — CVIS — RCH SOLUTIONS — redIT — ponte — KEL — WOLFRAM RESEARCH MAKERS OF MATHEMATICA — T2 TECHNOLOGIES — ICE systems — ENGAGE Technology GRID Computing Solutions

# (Grid) Programming Paradigms

- **(Grid) Programming Paradigms**
  - New Programming Paradigms will include non-Grid applications or usage scenarios, e.g. apps running with our without Grid depending on problem sizes.

- **Expected Future Programming Paradigm**
  - Application express their needs and behavior towards infrastructure that re-arranges itself accordingly
  - Infrastructure describes resource availability and properties (e.g. HPC backbone topology) to the application that (re-)compiles or (re-) configures itself accordingly
  - Platform position paper contributed to the Oxford Challengers Workshop on Standards Roadmap to 2020:
    http://www.w3c.rl.ac.uk/pastevents/ChallengersWorkshop/Schott.pdf
    (Bernhard Schott, Chris Smith, Werner Dubitzky)

**Platform**

We assume a world were applications are capable of adjusting themselves to the available environment; at the same time intelligent environments adjust themselves to applications requirements.

Whether the resource (Grid, Storage) managers provide the necessary dynamic configuration services or the complementary launcher components are provided by the applications, this scenario will result in automatically provisioned environments that would go along with computing resources emulating different machine architectures, types.

<...>

It's this notion of application portability that needs the descriptive capabilities on the application side; portable applications that can work in the grid environment and a grid fabric that can adapt to application needs.

Full document: http://www.w3c.rl.ac.uk/pastevents/ChallengersWorkshop/Schott.pdf

# (Grid) Programming Paradigms

- 3 Examples on State of the Art application to Grid integration

  1. Grid embracing the Application: Platform LSF

     - Platform LSF is the industry workhorse for some >2000 customers around the world. In most cases applications are not "written" for the Grid use, do not use Grid API. The Grid is in charge to adapt itself to the application by comfortable, sophisticated and even dynamic "embracement".

  2. SOA application integration: Platform Symphony

     - Platform Symphony implements "real time" execution of some application classes (SOA) in the HPC Grid (=cluster, extended cluster). Keep It Sweet and Simple: Very easy to use APIs allow the application to request resources (=service instances) instantaneously..

  3. On-demand resource acquisition and formation of application specific Grids: QosCosGrid

     - QosCosGrid (= Quasi Opportunistic Supercomputing for Complex Systems on the Grid) Complex Systems applications use the QosCosGrid-Toolbox to parallelize their workload in order to use distributed resources.

     - Towards the QosCosGrid-Broker they express requirements and behavior.

     - The Broker acquires resources on-demand and form an application specific Grid

**Example 1:**
**Grid embracing the**
**Application:**
**Platform LSF & EGO**

Whatever the Grid does for you already, saves you lines of code and more …

- Overview on the common layered architecture and LSF
  - Foundation: unified resource layer – Enterprise Grid Orchestrator
  - Describing Resource Allocation: Policies on supply

- Embracing the Application:
  - Interfaces: GUI, CLI, APIs, BES/JSDL/HPC-Profile++, DRMAA, Workflow (GUI, CLI, XML), ESUB
  - Arbitrary (dynamic) Resource Semantics

- Describing Application behavior: Application Profiles

- Describing Usage: Policies on consumption
  - LSF Modular Scheduler

# Platform Enterprise Grid Orchestrator

## Open & Decoupled Architecture

**Applications**

| LS | MDA | EDA | CAE | FSI | VM's | J2EE | DB's | ERP | CRM | BI |

**SOA**

**Application Workload Management**

| API/CLI | API/CLI | API/CLI | API/CLI | API/CLI | API/CLI |
|---|---|---|---|---|---|
| Platform LSF HPC | Platform LSF | Platform Symphony | Platform VMO | Platform Process Manager | 3rd Party Middleware Integration |

**Platform EGO SDK/API**

**System Resource Orchestration**

**Platform EGO Standard Services**

| Data Cache | | Event Service | Deployment Service | Logging Service | Portal Service | Service Director | | Storage |
| SNMP | | | | | | | | License |
| Security | | | | | | | | e.g. Infiniband |

**Platform EGO Kernel**

**SOI**

| Infrastructure Plug-ins | Manage | Allocate | Execute | Fail-over | Resources Plug-ins |

**Grid Devices**

| Linux | Windows | Aix | Solaris | Servers | Desktops |
|---|---|---|---|---|---|
| H/W | H/W | H/W | H/W | H/W | H/W |

# Platform EGO – Simple Use Case



**Application Orchestrator**

**EGO SDK**

**Platform EGO**

Q. What resources are available for my application based on current usage?

A. Resources available (Host1, Host2, etc) for Consumer

Q. Submit Allocation Request ?

A. Allocate Hosts to Consumer

Q. Run my Services on EGO

A. Start Up required Application infrastructure (Activities and Services)

# Dynamic, flexible, scalable

**Platform**

- By decoupling DRM and Workload Management, multiple Workload-Managers can dynamically flex based on workload and relative priorities – borrowing resources from other clusters and consumers

- Applications can request additional resources - even during runtime

A,A,A  B,B,B,B  C,C,C  D,D,D

B,B,B

| VMO | LSF | Symphony 1 | Symphony 2 |
|-----|-----|------------|------------|

Platform EGO

# Dynamic, flexible, scalable

- By decoupling DRM and Workload Management, multiple Workload-Managers can dynamically flex based on workload and relative priorities – borrowing resources from other clusters and consumers

- Applications can request additional resources - even during runtime

A,A,A          B,B,B,B          C,C,C          D,D,D

B,B,B

| VMO | LSF | Symphony | Symphony |
|-----|-----|----------|----------|

Platform EGO

- **Application Interfaces:**
  - GUI, CLI, APIs,
  - BES/JSDL/HPC-Profile++,
  - DRMAA,
  - Workflow (GUI, CLI, XML),                    (example: SAS)
  - ESUB: post submit / pre-queuing admin controlled modification of parameters
  - ELIM: dynamic resource metric update
  - Scheduler-Plug-In interface

- ## Arbitrary (dynamic) Resource Semantics
  - Add custom resource descriptions, dynamic metrics by ELIMs
  - Boolean, numeric, string, static & dynamic values, consumable
  - Example: resource [tape_drive==tape-ID]
    bsub –R "rusage tape_drive==XYZ1711" will send the job to the host with a tape drive that has loaded tape with ID=XYZ1711. If not (yet) fulfilled, job will be pending with pending reason "waiting for resource tape_drive==XYZ1711"
    Extension: ELIM to "measure" tape requests, forwarding to tape-library controller
  - Example: Application metrics. DB indicating to load balancing infrastructure remaining #connection capacity via ELIM interface.
  - Example: Application license metric. License manager interfaces via ELIM to load balancing infrastructure, indicating possible number of concurrent use..

# Embracing the application: Platform LSF

- **Describing Application behavior and specifics:**
  **Application Profiles**
  - Independent from "per host" or "per queue" definitions
  - pre-execution, post-execution, jobstarter (job environment), queue
  - resource requirements: mem, cpu-type, scratch space, I/O, licenses, access to networks, databases, file systems,… etc.
  - limits, runtime (soft, hard), mem (per job / per process), limits on processes, threads, …
  - check-pointing method, signaling (signal sequence)
  - chunk (size)
- **Controlling application runtime behavior**
  - Exit codes (lists of -"success" exit codes, rerun exit codes, rerun but not on same host exit codes, max #reruns / #requeues)
  - Trigger on run too long / too short
  - Trigger on lack of activity & mem growth. (Not mem limit!)
  - Triggers optional actions like (default: email to admin) bstop –u $username all

- **Describing Usage: Policies on consumption**
  - Fairshare methods: host partition, cross-queues, inter-queues, hierarchical,
  - Limits matrix: per queue, per user, per host (queue based, host based), per system
  - Policies: priorities, preemption, dynamic priority, parallel-greedy (job slot, mem), backfill, infinite backfill, absolute priority, checkpointing, migration, advance reservation …
  - SLA / goal based scheduling: deadline, speed, throughput, resource-SLA, combinations, … (more next slides)
  - Scheduling target host/CPU/core (per host definition)
  - LSF Modular Scheduler: custom policy plug-ins – use your own policy without need to rewrite the whole scheduler.

- **System reliability:**
  - self-healing, recovery from incidents,
  - policy driven proactive problem containment, "black hole" isolation
  - **no** job loss during operation or in error condition, reconfiguration or failover

- LSF implements and executes SLAs for workload submitted to a "service class" (e.g. deadline, throughput, speed, …)
- SLA is proactively managed by allocating more resources or starting more jobs to achieve scheduling goal
- Handing over estimated run times for the jobs improve scheduler precision – otherwise, scheduler will learn about the runtime
- To resources, this translates into "least impact scheduling"
  - A given set of resources (= a (set of) cluster(s)) can serve more "happy" scientists at the same time
  - = progress in more projects at the same time
  - = faster turn-over, shorter time to results for more projects a the same time.

Classical
opportunistic
scheduling

I need to work now!

100% —

= 8 Job-slots

My jobs

Cluster filled to 100%

now

time

LSF-SLAs : SLA "Deadline"

100% = 8 Job-slots

Classical opportunistic scheduling

I need to work now!

Cluster filled to 100%

now — time

Early enough for me

100%

Free resources for dialog users, real-time requests, online sessions, training sessions, other workload

SLA 1 "deadline"

SLA 1 consumes 50% of cluster

now — time — "deadline"

- How are those scheduling features related to a Grid Programming Paradigm? (question from EchoGrid Athens'08)

- By embracing the application for its complete lifecycle in the Grid, Platform LSF simplifies Grid integration for applications

  – Application profiles implement a simple information exchange between application and infrastructure.

  – Distinct application, consumption and resource policies allow for dynamic adaptation and balancing between application and infrastructure requirements

  – Self-management functions for automated reaction on not-intended application or infrastructure behaviour care for incidents

- These rich infrastructure features take away the need to implement respective functions in the application code

- Build your application to your application needs

- For using the Grid, just supply an application profile giving instructions to the infrastructure: "This is how I want to be executed"

**Platform**™

Example 2:
HPC-SOA application
integration
Platform Symphony

# Performance, at any Scale

| | Symphony at IBM DCCoD |
|---|---|
| **Scalability** <br> 1,000 concurrent clients, 100 applications | **20,000+ CPU's** simulated on 1,000 physical CPUs in one cluster |
| **CPU Utilization** <br> 1-100 clients, 1 sec task, 1KB message, 2,000 CPU | **98%** |
| **Task Throughput** <br> 1KB Message | **2,700 messages/sec** |
| **Single Task Round Trip** | **2.4 ms** |
| **Single Session Round Trip** <br> 100KB common data, 10 second 1KB Task | **11.8 ms** |

Cluster partition usage over the last week

Report period: from 11-19 00:00 to 11-26 00:00

Symphony Client Applications

register
application

SSM

EGO API

execute
<allocation,container>

Resource
Sharing Plan

100 cpus

A

B          C          D

50 cpus    25 cpus    25 cpus

Session
Director

allocate
<consumer,resreq>

return
allocation

EGO Master

Security
Plug-ins

Scheduling
Plug-ins

Resource
Plug-ins

EGO Agent

LIM          LIM          LIM          LIM

PEM          PEM          PEM          PEM

Lending/Borrowing

Consumer

SIM          SIM          SIM          SIM

SI          SI          SI          SI

**Platform**

- ## Symphony: parameter sweep

  1. decompose the problem

  2. execute parameter sweep

  3. merge results / aggregation of results

- ## Keep It Sweet and Simple application integration:

```
static double GridMonteCarloPi(int simulations) {
    // connect to the grid and create a session
    Connection connection = SoamFactory.Connect("computePi");
    Session session = connection.CreateSession(…);




}
```

Ask the grid to allocate resources for the "computePi" application, and make them available in a session

- ## Symphony: parameter sweep

  1. decompose the problem

  2. execute parameter sweep

  3. merge results / aggregation of results

- ## Keep It Sweet and Simple application integration:

```
static double GridMonteCarloPi(int simulations) {
    // connect to the grid and create a session
    Connection connection = SoamFactory.Connect("computePi");
    Session session = connection.CreateSession(…);

    // send the tasks on the grid
    int numTasksToSend = 10;
    double numberOfSimulationsPerTask = simulations/numTasksToSend;
    for (int taskCount = 0; taskCount < numTasksToSend; taskCount++) {
        session.SendTaskInput(numberOfSimulationsPerTask);
    }


}
```

Send parameters of sweep as input messages to compute nodes

**Platform**

- Symphony: parameter sweep

  1. decompose the problem

  2. execute parameter sweep

  3. merge results / aggregation of results

- Keep It Sweet and Simple application integration:

```
static double GridMonteCarloPi(int simulations) {
    // connect to the grid and create a
    Connection connection = SoamFa
    Session session = connection.Cre

    // send the tasks on the grid
    int numTasksToSend = 10;
    double numberOfSimulationsPerT
    for (int taskCount = 0; taskCount
        session.SendTaskInput(numberOfSimulationsPerTask);
    }

                            ;
}
```

Upon message reception each compute node executes the service code and returns a result

```
public override void OnInvoke(TaskContext taskContext)
    {
        double hits = 0; double x, y;
        double simulations = (double) taskContext.GetTaskInput();
        // seed
        Random Rnd = new Random(1234);
        // random throw for each simulation
        for (int i = 1; i < simulations; i++)
        {
            x = Rnd.NextDouble();y = Rnd.NextDouble();
            hits += ((x * x + y * y) <= 1) ? 1 : 0;
        }
        // return the number of hits
        taskContext.SetTaskOutput(hits);
    }
```

- Symphony: parameter sweep
    1. decompose the problem
    2. execute parameter sweep
    3. merge results / aggregation of results

- Keep It Sweet and Simple application integration:

```
static double GridMonteCarloPi(int simulations) {
    // connect to the grid and create a session
    Connection connection = SoamFactory.Connect("computePi");
    Session session = connection.CreateSession(…);

    // send the tasks on the grid
    int numTasksToSend = 10;
    double numberOfSimulationsPerTask = simulations/numTasksToSend;
    for (int taskCount = 0; taskCount < numTasksToSend; taskCount++)  {
            session.SendTaskInput(numberOfSimulationsPerTask);
    }

    // get results back and aggregate them
    EnumItems enumItems = session.FetchTaskOutput((ulong)numTasksToSend);
    foreach (TaskOutputHandle output in enumItems) {
        hits += (double) output.GetTaskOutput();
    }
    return 4 * hits / simulations;
}
```

Client program receives the messages and aggregates the results

On the grid, bigger is not necessarily better

- ## Handle reliability in large scale computation
  - rerunning a task is a drop in the ocean
  - rerunning tasks costs nothing compared to the length of the entire job
- ## Single task failure:
  - automated isolation of black holes - application(-session) specific
  - other application may reuse this host and run perfectly
- ## Potential failure reasons:
  - One (or some) corrupted data records
  - in case of parameter sweep, only that one task will fail, that tries to compute the corrupted record

- XML Application profile (metadata) defines
  - Local action on failure
    - (e.g. "this machine's DB driver is not working properly, restart my task elsewhere")
  - Global action on failure
    - (e.g. "this job's data is corrupted, fail entire job ")
- Client application defines
  - action upon exception or result
    - (e.g. "this input didn't work, let me try another one")
- Overall, more flexibility for failure control
  - Yet simple to define
  - and not always necessary to recompile entire application to extend failure control

- Task is expendable: may be used to steer the application
  - Probe tasks can be used to measure application-specific properties,
    - Send a few montecarlo tasks, get their results, and use the elapsed time to predict overall computation time and reduce/increase precision of simulation
  - Or to measure grid properties
    - In effect, failure control and black listing is a sanity check on resources
- Grid as an open book
  - APIs to get resource availabily allow the application to fine-tune its own problem decomposition
    - No use to launch 100 tasks on 10 cpus, if 20 will do the job…

```
static float GridPortfolioPricing() {
...
    // send the tasks on the grid
    for (each option in Portfolio)  {
        session.SendTaskInput(option);
    }
}
```

```
class Option {
    float price();

}
```

class Vanilla

class Exotic

Client program prices a
diversified portfolio by
sending options on the grid

They may be Vanilla (very
quick pricing , ~1s) or
Exotic (long pricing, ~1h)

Problem: Exotic options finish pricing long after vanillas
More generally, how can a grid cater for
complex object models ?

**Platform**

```
static float GridPortfolioPricing() {
…
    // send the tasks on the grid
    for (each option in Portfolio) {
        session.SendTaskInput(option);
    }
}
```

```
class Option {
    float price();
}
```

V   E   V

class Vanilla

class Exotic

```
float price() {
…
    // send montecarlo tasks on the grid
    for (each trajectory) {
        session.SendTaskInput(trajectory);
    }
}
```

Exotic option reuses the grid by launching a montecarlo job

Solution: Exotic objects reuse the grid to speed up computation!
# Grid Oriented Object, a.k.a. GOO, is fluid…

- If tasks are expendable, preemption becomes possible

- Preemption means reallocation of an application's resource to another application

  - Moving resources according to priorities becomes less costly (in terms of lost cpu time)

- Consequence is: to take full advantage of application fluidity, the underlying container of application must also be quick and flexible

- How do those Platform Symphony features and applied methods relate to a Grid Programming Paradigm? (question from EchoGrid Athens'08)
- By delivering high performance instantiation services to SOA applications, and provide for its messaging and complete lifecycle in the Grid, Platform Symphony simplifies Grid integration for SOA applications
  - Easy to use API provide for short time to results. Low number of additional code lines (see example)
  - Detailed policies and resource SLAs allow for dynamic adaptation and balancing between application and infrastructure requirements. (lending & borrowing between application domains based on controlled policies)
- Application steering looks for the optimal balance between resource allocation size and problem decomposition – automatically achieving best possible performance AND best resource utilization even for applications with dramatically varying runtimes per problem.
  - In the systematic approach towards improved usability of Grids, we investigate whether application steering – up to now done application specific – will become a general infrastructure feature
- These rich infrastructure features take away the need to implement respective functions in the application code

# Example 3: QosCosGrid

# Grid ≈ virtual supercomputer

# Current Implementation



Legend

| Use cases | Parallel application communication channels | Data and file channels | Resource reservation and remote job management | Network and topology monitoring |

Applications

Interfaces, Services, Tools

Middleware

Grid Fabric

Use case 1 | Use case 2 | Use case n

QCG OpenMPI | SSH/GSISSH
QCG ProActive | GridFTP

Remote Grid Interfaces
(OpenDSP based on
DRMAA + Reservation APIs)

Local Resource Management System
(LSF Cluster)

Cluster and Network Monitoring System
(LSF PERF and Ganglia)

Local Network

6/12/2008    48

# Target Implementation

# QosCosGrid workflow

- On-demand resource acquisition and formation of application specific Grids: QosCosGrid
  - QosCosGrid (= Quasi Opportunistic Supercomputing for Complex Systems on the Grid) Complex Systems applications use the QosCosGrid-Toolbox to parallelize their workload in order to use distributed resources.
  - Towards the QosCosGrid-Broker they express requirements and behavior by xml "job profile"
  - The Broker acquires resources on-demand and form an application specific Grid, creates RTG = resource topology graph
  - RTG is used to map the application to the resources, placing MPI communicators at the right place

# QCG Job profile

- The QCG Job Profile document is inspired by an existing XML-based job description language supported by one of the main components of the QCG middleware called GRMS.
- End users can describe topology and resource requirements, in particular:
  - required aggregations and hierarchies of resources (computing nodes, clusters, sub-clusters, storage elements etc.),
  - required resource properties (operating system, memory, number of CPUs, speed of the CPU on a resource),
  - required network and connection properties (bandwidth, latency and capacity),
  - required applications and licenses available at destination computing resources.
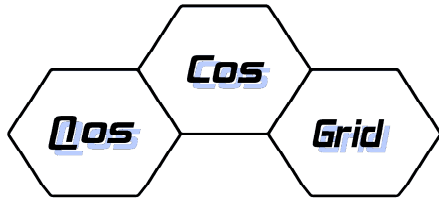
# QCG Resource Description Model

- RTG (Resource Topology Graph)
- A common XML resource description language
- Provide description of:
  - Resources, tasks, processes
  - Topology
  - Communication properties
- Serves as a "bridge" between the various system components
- Used to describe, publish, evaluate, reserve and monitor heterogeneous resources across the QosCos Grid

- Supplementary Java implementation:
  - Functional behavior and logic
  - XML to Java objects marshaling/un-marshaling
  - Specialized types of RTG objects, according to the middleware requirements (i.e. Resource advertisement, Meta-scheduling, SLA's, Monitoring, and User requirements. )
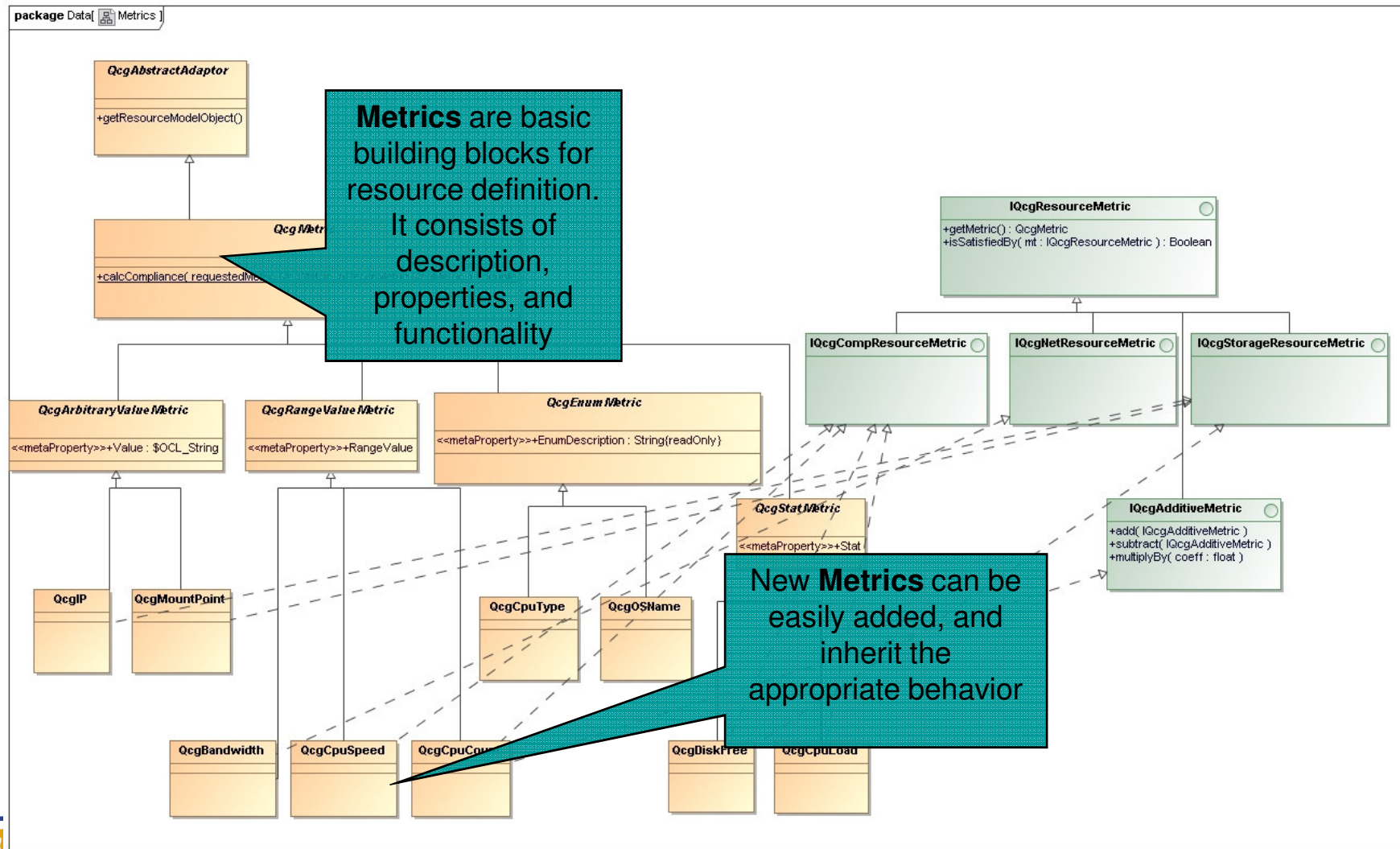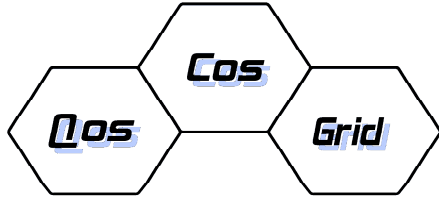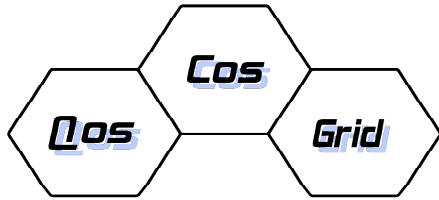
QCG RTG – Resource Description View

**Network Resource**
Define end-to-end network properties between *resources* and *Resource Communication Groups*

**Resource Communication Groups**
Define any combination of <u>nested</u> aggregations of homogeneous/Heterogeneous resources/clusters etc.

**Resource Availability**
Describes availability and other state related information along time intervals

Each **Resource** (computing, storage, network) is defined by a set of **Metrics**:
cpuCount
cpuSpeed
cpuType
memoryTotal
diskTotal
osName
…

AD PSNC
CR4
CR1  CR2  CR3  SR1

AD INFRA
CR1  CR2  CR3
CR4  CR5  CR6  CR7  CR8

AD UPF
CR2  CR3  SR1
CR4  CR5

# Resource Co-Allocation Protocol

- A simple, robust 2-phase commit alike protocol for planning based co-allocation of resource across multiple ADs
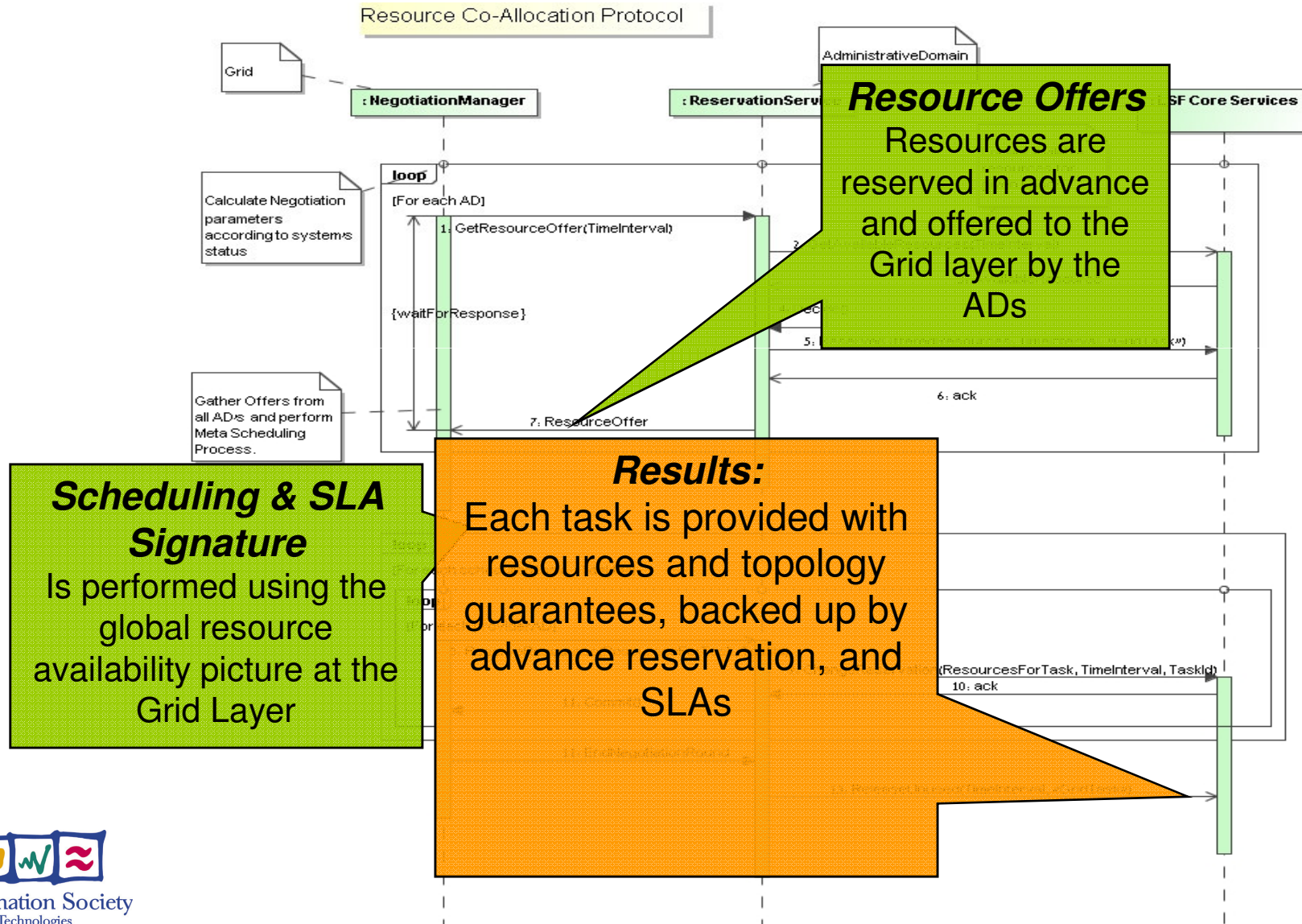
Highlights:

- Simple negotiation strategy, prevent "bargaining" which would prolong the number of negotiation rounds.

- Global guarantied resource pricing methodology, allows each AD to estimate theirs potential profit.

- Each AD decides what to contribute, and provides a guaranteed "Resource Offer" accordingly.

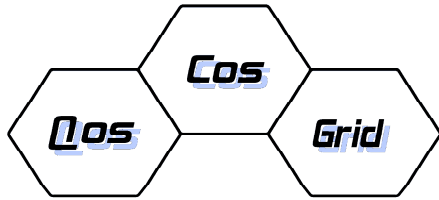- Support co-allocation scenarios, requiring a coordinated booking from different ADs.
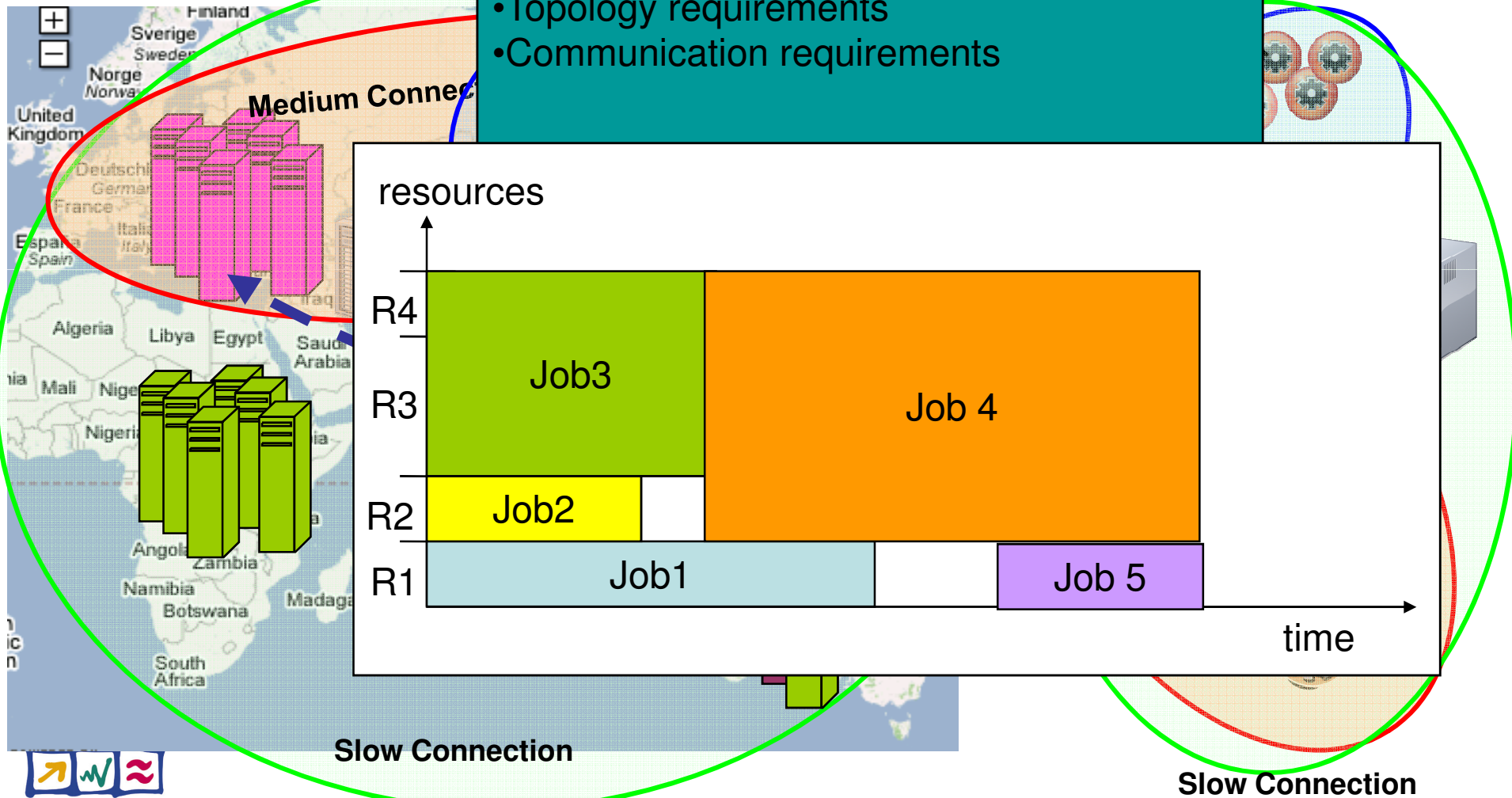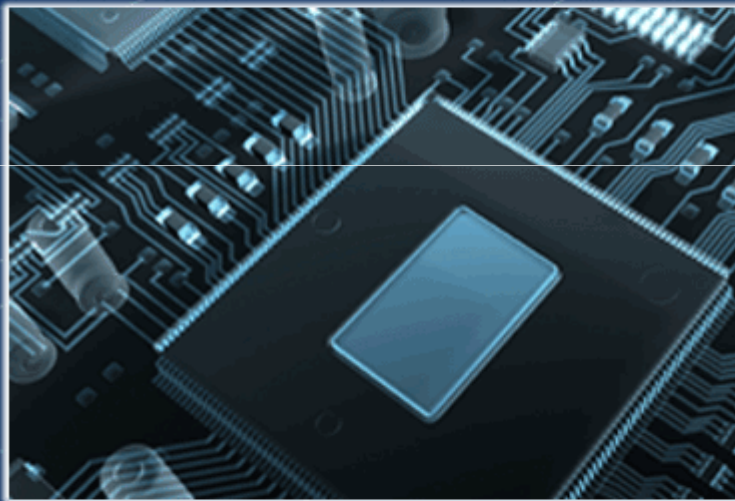
# Resource Co-Allocation Protocol



**Resource Offers**
Resources are reserved in advance and offered to the Grid layer by the ADs

**Scheduling & SLA Signature**
Is performed using the global resource availability picture at the Grid Layer

**Results:**
Each task is provided with resources and topology guarantees, backed up by advance reservation, and SLAs

# Meta Scheduler - The Objective

To provide a resource allocation for the user's tasks, constrained by:
- Resource requirements
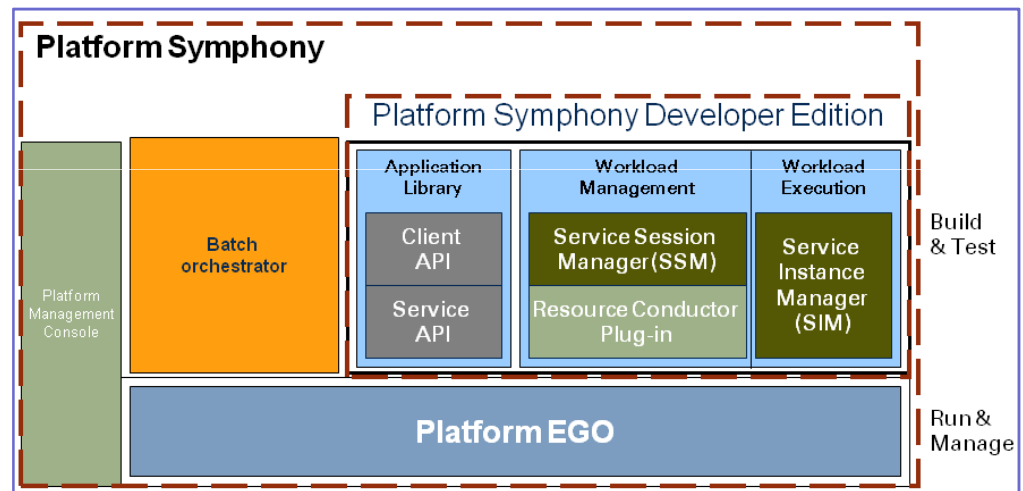- Topology requirements
- Communication requirements

**Medium Connection**

**Slow Connection**

**Slow Connection**



resources

| | |
|---|---|
| R4 | Job3 |
| R3 | |
| R2 | Job2 |
| R1 | Job1 |

Job 4

Job 5

time

Information Society Technologies

# Summary

## Conclusion:

- Less code due to infrastructure services

- Auto-adaption of apps & infrastructure to each other

  - Advantage: overcoming the versioning problem, legacy code problem application tells infrastructure which version to provide.

  - Advanced Grid services: changing app configuration, optimizing problem decomposition

- Maintain application autonomy: run on single host must be possible

  - Example:
    Symphony Developer Edition: single host Grid

- Multi-Components applications / multi-instance (MPI, HPC-SOA) / SOA:

  - production stability depends on total availability (error propagation, see next page!)

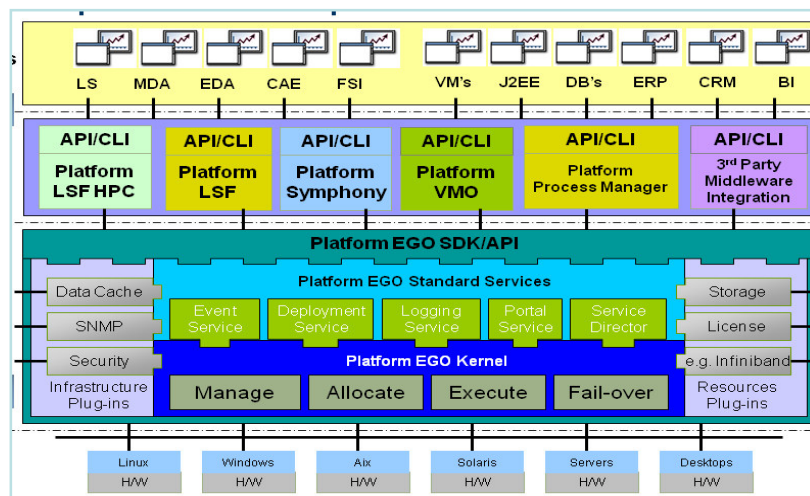  - Auto-response Grid infrastructure to improve total availability



Platform Symphony

Platform Symphony Developer Edition

Platform Management Console | Batch orchestrator | Application Library — Client API, Service API | Workload Management — Service Session Manager (SSM), Resource Conductor Plug-in | Workload Execution — Service Instance Manager (SIM) | Build & Test

Platform EGO | Run & Manage

| Component and System availability | | | Comments |
|---|---|---|---|
| p-component | # component | p-total | Probability per component and total system |
| 0,9 | 2 | 0,81000 | |
| 0,99 | 2 | 0,98010 | 98% good = 175 hours bad every year |
| 0,999 | 2 | 0,99800 | |
| 0,9999 | 2 | 0,99980 | |
| 0,99999 | 2 | 0,99998 | production stability: four 9's minimum |
| 0,9 | 4 | 0,65610 | |
| 0,99 | 4 | 0,96060 | 96% good = 350 hours bad every year |
| 0,999 | 4 | 0,99601 | |
| 0,9999 | 4 | 0,99960 | 99,96% good = 210 minutes bad every year |
| 0,99999 | 4 | 0,99996 | production stability: four 9's minimum |
| 0,9 | 12 | 0,28243 | |
| 0,99 | 12 | 0,88638 | |
| 0,999 | 12 | 0,98807 | |
| 0,9999 | 12 | 0,99880 | |
| 0,99999 | 12 | 0,99988 | will not reach production stability |

| | | | |
|---|---|---|---|
| probability worse than 0,9 | => | more than | 876 hours disfunct per annum |
| probability better than 0,9 | => | less than | 876 hours disfunct per annum |
| probability better than 0,9999 | => | less than | 53 minutes disfunct per annum |

Caution:
be aware –
even 52
minutes per
year could
mean 52
breakdowns
!!
(1 minute
duration
each)

# Reliability → Usability → Value

- **Performance Platform LSF** (by 2006/07):
  - 10millions jobs per day throughput with >95% job-slot utilization
    based on EDA job mix, max 5min for failover. (EDA job mix: 5min, 15min, 30min job-runtime, 120Hz job submission and dispatch rate)

- Performance and Scalability translates into Reliability

- Reliability can be measured as "MTBF" - Mean Transactions (=Jobs) Between Failure

- Reliability is achieved by proactive incident management

  - self-healing, recovery from incidents, policy driven proactive problem containment (both for resources and applications), **no** job loss during operation or in error condition, reconfiguration or failover

- Grid moves toward ease of use, plug&play app integration.
- Applications and infrastructure will describe themselves and adapt to each others requirements resp. offerings
- Infrastructure will cover up on disruptions, incidents, both on resource as application side
- First steps already done! Move on towards flexibility && reliability. Need for standards.
- Looking forward to discuss and share with you!



Bernhard Schott
Dipl. Phys.
EU-Research Program Manager

Platform Computing GmbH

Frankfurt Office
Direct        +49 (0) 69 348 123 35
Mobile       +49 (0) 171 6915 405
Email:       bschott@platform.com
Skype:      bernhard_schott
Web:         http://www.platform.com/