# A Gentle Introduction to Grid Computing

**Borja Sotomayor**

University of Chicago

Department of Computer Science

borja@cs.uchicago.edu

February 18, 2008

---

THE UNIVERSITY OF

CHICAGO

- Website/Wiki: **http://acm.cs.uchicago.edu/**
- Get announcements about ACM events on our mailing list.
  - ◆ https://mailman.cs.uchicago.edu/mailman/listinfo/acm
    (or follow convenient link on ACM website)

# Acknowledgements

- Food provided courtesy of Grid.org and the Department of Computer Science.



- Some slides borrowed from Ian Foster.

# Index

- Introduction
- How does this work?
- More than just a "virtual supercomputer"
- Applications on the Grid
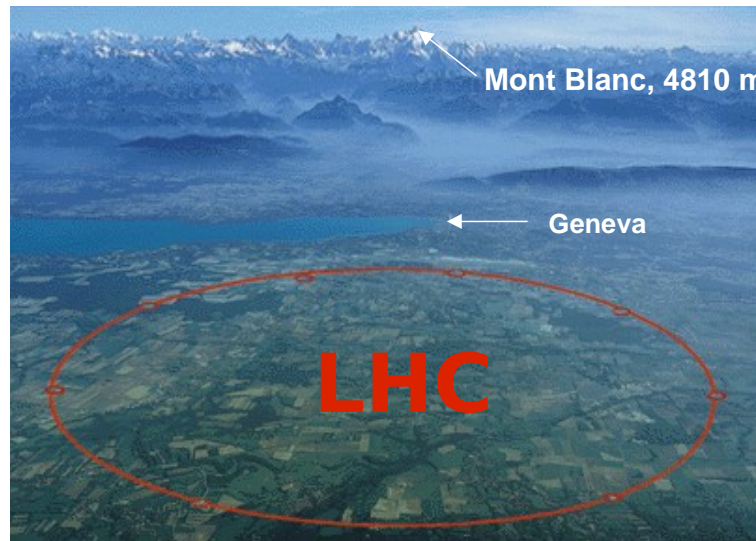- Globus Toolkit 4
- I want to know more!

# Index

# A problem... (I)

# A problem... (II)

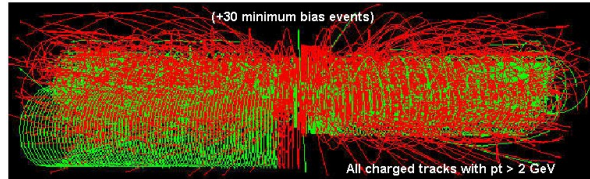

Mont Blanc, 4810 m

Geneva

LHC

---

# A problem... (III)

- The LHC (Large Hadron Collider), which is being built in CERN, is a particle accelerator/collider with a circumference of 27km (16.7mi).
- Will answer many interesting questions, specially: Does the Higgs boson exist?
- When it starts to work this year, it will produce *huge* amounts of information.
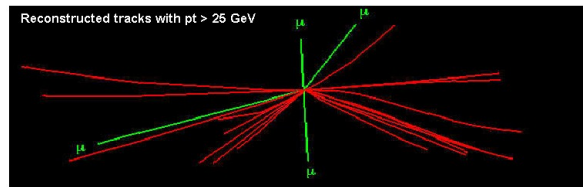
# A problem... (IV)

*From this event (1 event = 1 collision)...*



*(+30 minimum bias events)*

*All charged tracks with pt > 2 GeV*

*We're searching for this characteristic signature:*



*Reconstructed tracks with pt > 25 GeV*

**1 in $10^{13}$**

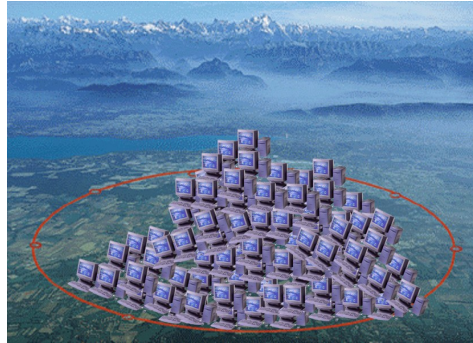**Like looking for one person in a thousand world populations.**

---

# A problem... (V)

- 40 million collisions per second (300 GB/s of data)
- These collisions are filtered, and only "interesting events" are kept.
  - ~300 MB/s (~25 TB/day!). This information requires a (non-trivial) processing, and must be stored for future reference and study.
- LHC will produce ~15 Petabytes of information per year.
  - 1000x of information published in books in a year.
  - 1% of all human-produced information in a year

# A problem... (VII)

- Using current technology, processing and storing all that data in a single site is impossible.
- An estimated 100,000 processors would be needed to deal with the LHC's computational needs.

# A solution

- Problem: A single node can't handle all that work.
  - ◆ But the combined power of *several* sites might be able to handle it.
- Solution: Achieving greater performance and throughput by pooling together resources from different organizations
  - ◆ Informally, this is what Grid Computing is all about (better definition coming up)

# Virtual Organizations



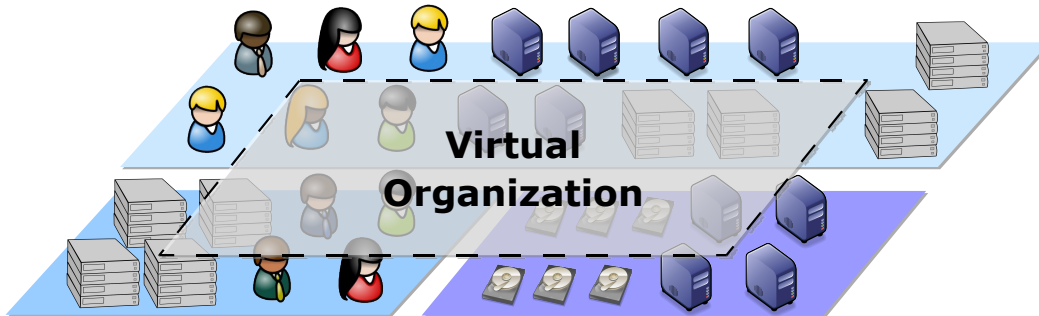- An organization can only use the resources it has under its control.
- Grid Computing involves multiple organizations.
- Resources from these organizations are dynamically pooled together, creating *virtual organizations*, to solve specific problems.

# Easier said than done!

- How do we decide what resources are part of each virtual organization?
- Given a computational task, how do we decide what resources will be allocated to deal with that task? For how long?
- How do we get the resources to communicate amongst themselves? Take into account that these are *heterogeneous* resources from *different* organizations!
- If I want to "split up" a task so that it can be performed in parallel by several computers in different organization, how to I actually "split up" the program?
- A lot of security challenges. For example, how can an organization make sure its resources are only being used by trusted users and that they are not being abused by malicious users?

# Why is this hard?

- No centralized control
  - We *cannot* and *do not* control the decisions of individual organizations. Grid Computing does not impose an all-powerful master that overrides local decisions.
  - We have to reconcile all the different policies in each site.
- Shared heterogeneous resources
  - We cannot assume that all resources are exactly alike (as we could assume in a computing cluster)
  - Resources appear and disappear on the grid.
- Communication and coordination
  - Different sysadmins, users, geo-political restrictions, etc.

# A definition

- Ian Foster provides an (open) definition in the paper *What is the Grid? A Three Point Checklist.*
- A grid is a system that:
  - coordinates resources that are not subject to centralized control...
  - ...using standard, open, general-purpose protocols and interfaces...
  - ...to deliver nontrivial qualities of service
- There is no "The Grid", but there are many production "grids" around the world that support a wide variety of applications.

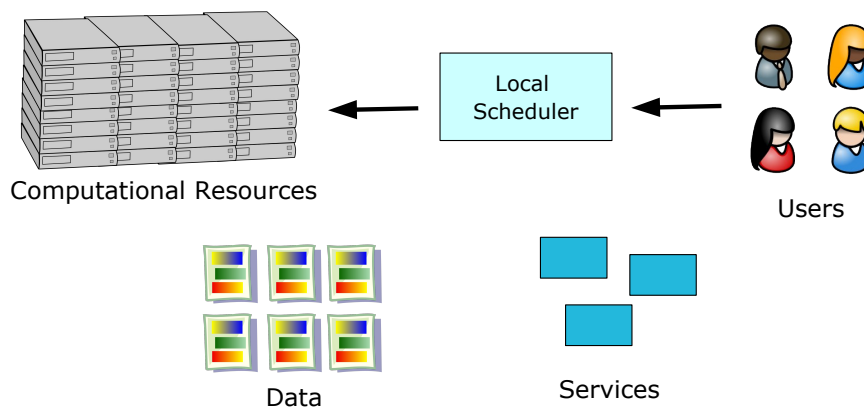# Index

---

# What a typical site could look like



Computational Resources

Local Scheduler

Users

Data

Services

- Example: UChicago has a cluster called Teraport (http://teraport.uchicago.edu/) that uses Torque/Maui as its local scheduler.

# A typical interaction



**Springfield University** — **Shelbyville National Laboratory** — **UofCC (University of Capital City)**

# The nuts and bolts

- The nuts and bolts of a grid system are mainly *middleware*: networked services, protocols, sofware toolkits, etc.
- Standards for grid services:
    - Open Grid Services Architecture, developed by the Open Grid Forum (http://www.ogf.org/)
    - Standards from OASIS, W3C, IETF, ...
    - Many grid standards are still under development.
- ... and implementations:
    - Globus Toolkit (http://www.globus.org/)
    - gLite (http://cern.ch/glite)
    - UNICORE (http://www.unicore.eu/)

# Index

# More than just a "virtual supercomputer"

- Grid Computing is sometimes characterized as creating a "virtual supercomputer", capable of handling huge computations that cannot possibly be run on a single site (e.g., LHC)

- This is only one of the use cases for Grid Computing. There are many more.

- Grid Computing is not just about sharing computational power. A virtual organization also includes *services*, users, instruments, etc.

  - Service-Oriented Science (Ian Foster, Steve Tuecke, *The Many Faces of IT as Service;* Ian Foster, *Service-Oriented Science)*

# First Generation Grids

Focus on aggregating large amounts of resources for massively parallel applications

EGEE

Open Science Grid

# The original "electric grid" analogy

At first, we had to be close to electrical generators

Now: the *electric grid* carries electricity from large producers to small consumers.

Can we do the same with computation?

# Enabling e-Science

- Grid Computing can provide a "cyberinfrastructure" to projects that require...
    - ... analysis of large quantities of data, but are currently confined to running on a single site (or even a single computer!)
    - ... interdisciplinary and/or interorganizational collaboration.
    - ... dynamic reorganization of resources (possibly including instrumentation) to explore new experiment configurations.
- Allowing...
    - ... composition of new applications based on existing services.
    - ... existing systems to scale and test new hypothesis, explore more experiment configurations, etc.
    - ... users to focus on their applications, while grid software takes care of planning execution on remote sites (e.g., Swift: http://www.ci.uchicago.edu/swift/)

25

# Scaling up Social Science: Citation Network Analysis



*Work of James Evans, University of Chicago, Department of Sociology*

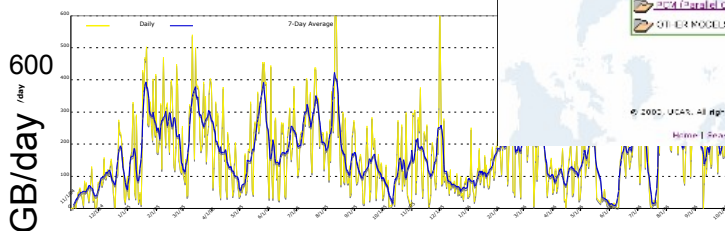http://www.ci.uchicago.edu/research/detail_socialinformatics.php

26 26

# Index

- Introduction
- How does this work?
- More than just a "virtual supercomputer"
- **Applications on the Grid**
- Globus Toolkit 4
- I want to know more!

---

# Earth System Grid

- Provides access to all IPCC data
- >150 TB data downloaded
- >300 scientific papers written

Social Informatics Data Grid

Bennett Berthenthal et al., sidgrid.ci.uchicago.edu


NIH's
Cancer Biomedical Informatics Grid

*caBIG: sharing of infrastructure, applications, and data.*

Cancer Center (8)
Clinical Cancer Center (14)
Comprehensive Cancer Center (39)
Planning Grant (7)

caBIG   cancer Biomedical Informatics Grid

https://cabig.nci.nih.gov/

Integrating Data and Computing, on Demand

**Public PUMA Knowledge Base**

Information about proteins analyzed against ~2 million gene sequences

**Back Office Analysis on Grid**

Millions of BLAST, BLOCKS, etc., on OSG and TeraGrid

Natalia Maltsev et al., http://compbio.mcs.anl.gov/puma2



Global Communities

# Applications

- Type
  - Computation
  - Large volumes of data
  - Distributed collaboration
- Common aspects
  - Size or complexity of the problem
  - Inter-organizational collaboration
  - Sharing of computational resources, data, and instruments.

33

# Index

- Introduction
- How does this work?
- More than just a "virtual supercomputer"
- Applications on the Grid
- Globus Toolkit 4
- I want to know more!

34

# What is the Globus Toolkit?

- A collection of solutions to problems that come up frequently when building collaborative distributed applications
  - Job management, data management, information services, metascheduing, etc.
- Does not provide turnkey solutions. It provides building blocks for software developers and system integrators.
- http://www.globus.org/toolkit/

# Globus Philosophy

- Globus was first established as an open source project in 1996
- The Globus Toolkit is open source to:
  - Allow for inspection
    - for consideration in standardization processes
  - Encourage adoption
    - in pursuit of ubiquity and interoperability
  - Encourage contributions
    - harness the expertise of the community
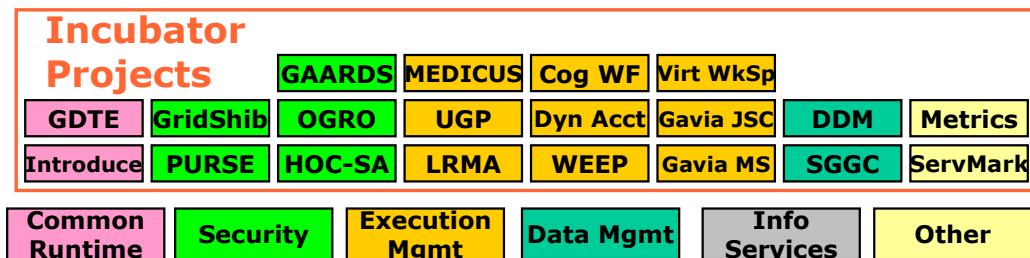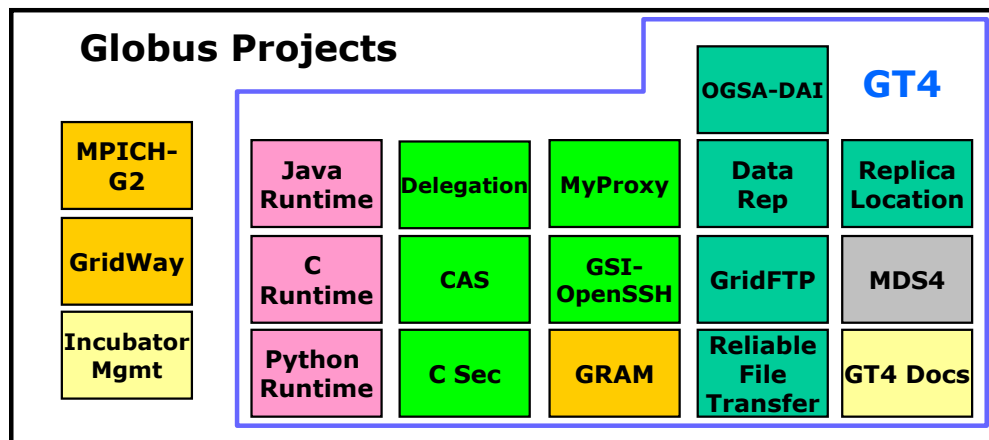- The Globus Toolkit is distributed under the (BSD-style) Apache License version 2

# dev.globus

- Governance model based on Apache Jakarta
  - Consensus based decision making
- Globus software is organized as several dozen "Globus Projects"
  - Each project has its own "Committers" responsible for their products
  - Cross-project coordination through shared interactions and committers meetings
  - New projects can be proposed by anyone through an incubation process.
- A "Globus Management Committee"
  - Overall guidance and conflict resolution

# Globus Software: dev.globus.org

**Globus Projects**

**GT4**

| MPICH-G2 | Java Runtime | Delegation | MyProxy | OGSA-DAI | |
|---|---|---|---|---|---|
| | | | | Data Rep | Replica Location |
| GridWay | C Runtime | CAS | GSI-OpenSSH | GridFTP | MDS4 |
| Incubator Mgmt | Python Runtime | C Sec | GRAM | Reliable File Transfer | GT4 Docs |

**Incubator Projects**

| | GAARDS | MEDICUS | Cog WF | Virt WkSp | | |
|---|---|---|---|---|---|---|
| GDTE | GridShib | OGRO | UGP | Dyn Acct | Gavia JSC | DDM | Metrics |
| Introduce | PURSE | HOC-SA | LRMA | WEEP | Gavia MS | SGGC | ServMark |

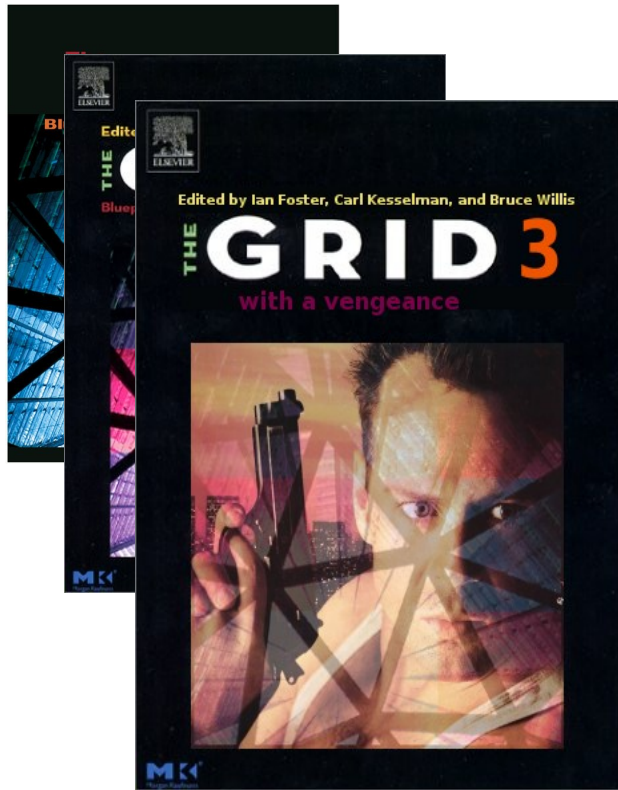| Common Runtime | Security | Execution Mgmt | Data Mgmt | Info Services | Other |
|---|---|---|---|---|---|

# Index

- Introduction
- How does this work?
- More than just a "virtual supercomputer"
- Applications on the Grid
- Globus Toolkit 4
- I want to know more!

# I want to know more!

- GridCafé provides a good introduction to Grid Computing
  - http://gridcafe.web.cern.ch/
- Books...

... d 2". Edited by ...ster and Carl ...man. Morgan ...ann, 2003.

---

# I want to know more!

- Books
  - "Grid Computing: The Savvy Manager's Guide". Pawel Plaszczak, Richard Wellner, Jr. Morgan Kaufmann, 2005.
  - "Globus Toolkit 4: Programming Java Services". Borja Sotomayor, Lisa Childers. Morgan Kaufmann, 2005.
- Websites
  - The Grid Index: http://www.gridindex.org/
  - Grid Gurus blog: http://gridgurus.typepad.com/
  - Ian Foster's blog: http://ianfoster.typepad.com/

# Questions?

**Borja Sotomayor**

Department of Computer Science

University of Chicago

borja@cs.uchicago.edu