



## Grid enabled life science applications: status and perspectives

Vincent Breton LPC IN2P3/CNRS Credit for the slides: M. Hofmann, N. Jacq, V. Kasam, J. Montagnat



Institut National de Physique Nucléaire et de Physique des Particules







http://www.itb.cnr.it/bioinfogrid



#### • Introduction

- EGEE biomedical applications
- Focus on virtual screening

#### the Vision

#### **Computing Grid** For data crunching applications

An environment, created through the sharing of resources, in which heterogeneous and dispersed data :

- molecular data (ex. genomics, proteomics)
- cellular data (ex. pathways)
- tissue data (ex. cancer types, wound healing)
- personal data (ex. EHR)
- population (ex. epidemiology)

as well as applications, can be accessed by all users as an tailored information providing system according to their authorisation and without loss of information.

Data Grid Distributed and optimized storage of large amounts of accessible data Knowledge Grid Intelligent use of Data Grid for knowledge creation and tools provisions to all users





# Computing grid applications are being deployed successfully

A few successful data grids (BIRN, BRIDGES, Medical Data Manager)

No knowledge grid yet deployed

Data Grid Distributed and optimized storage of large amounts of accessible data Knowledge Grid Intelligent use of Data Grid for knowledge creation and tools provisions to all users



# **EGEE** biomedical applications

# **Biomed Achievements (I)**

- Goal: demonstrate grid potential for real-scale biomedical applications.
- Start of EGEE in Apr 04: several app. prototypes developed, not yet deployed.
- First year achievements
  - Organization of work
    - Creation of the "Biomed" Virtual Organization
    - Deployment of associated services (guinea pig VO)
      - Definition of application test cases
      - Use of them to test new or updated EGEE components
    - Creation of the "Biomed Task Force"
      - Biomedical user support: GGUS, Data Challenge, tutorials, ...
      - Collaboration with middleware and infrastructure activities
  - Application deployment
    - Successful deployment of applications in the field of bioinformatics and medical imaging
    - ~70k jobs from Biomed users reported at EGEE's first review

# **Biomed Achievements (II)**

- Development of secured data management and complex data flows on the grid
  - Medical Data Management group has demonstrated complete chain for processing medical images on the grid using these services
- First CPU-intensive grid deployments for bioinformatics in the world
  - In silico drug discovery against malaria and bird flu
  - Very large impact in the grid community
  - Biologically-relevant results being processed
- Sustained growth of the "Biomed" VO
  - New apps. interested in joining the VO: 11 in DNA4.4 inventory
  - 3 sub-areas: bioinformatics, medical imaging, drug discovery
  - ~ ~80 users
  - 1000 jobs / day on average



#### **Medical image processing**

#### GATE: Radiotherapy planning

- CNRS-IN2P3
- Monte Carlo simulation
- Parallel execution on different seeds



- Pharmacokinetics: contrast agent diffusion study
  - UPV
  - Medical images registration
  - Distribution of registration pairs



#### Medical image processing

#### • SiMRI3D MRI simulation

- CNRS-CREATIS
- Magnetic Resonance physics simulation (Bloch's equation)
- Parallel processing (MPI)



- gPTM3D: Radiological images segmen
  - CNRS-LRI, CNRS-LAL
  - Deformable-contour based segmentation
  - Interactivity through agent-based scheduling



#### **Bioinformatics**

#### BioinfoGRID

- GPS@: bioinformatics portal
  - CNRS-IBCP
  - http://gpsa.ibcp.fr/ web portal
  - Existing (but overloaded NPSA portal)
  - Tens of bioinformatics legacy code
  - Thousands of potential users

Welcome on GPSA, Grid Genomic Web P	ortal - Mozilla Firefox			- 0			
Eichier Edition Affichage Aller à Ma	rque-pages Ouțils <u>A</u> ide						
• 📦 - 🛃 🙁 😭 🗋 http://gps	a.ibcp.fr/php/to_html.php?method_class	=similarity	👻 🔘 ок	G.			
🗋 Red Hat, Inc. 📄 Red Hat Network 🗀 Sup	port 🗀 Shop 🗀 Products 🗀 Training	9					
Grid Pl Bioinform	rotein Sequence	<b>e @nalysi</b> 10 protein sequence :	<b>3</b> analysis.	e <mark>e</mark> ee			
[GPSA] [n	wSEQ] [HELP] [REFERENCES] [?	NPS@] [PBIL-Gerland	] [ <u>PBIL]</u>				
<u>Accueil</u> <u>PSSP</u>	<u>Patterns</u>	<u>Alignment</u>		Similarity			
Query sequence filename :	/biogrid/vlefort/work/test/prt_test_bk.f	asta	Parcourir				
Sequence databank filename :	Ifn:genomics_gpsa/db/swissprot/spro	genomics_gpsa/db/swissprot/sprot.seg F					
Ssearch William Pearson, 1991 William Pearson [more]      Blastp (only on SwissProt) William William Pearson [more]  Running mode : EGEE      RUN_CLEAR	m Pearson, 1991						
GPS@ grid portal 2005. Contact: Christoy Terminé	he Blanchet	1.	_				

- Electron-microscopic image reconstruction
  - CNB-CSIC
  - Image filtering and noise reduction
  - 3D structure analysis





#### Secure Medical Data Management with EGEE

- 3 SRM-DICOM servers CNRS-I3S; CNRS-CREATIS, CNRS-LAL, CERN
- Distributed secure medical data management services



#### Potential vs current impact of EGEE applications on scientific community

NETTAB, 10/7/06 - 12

Application	dev	user	Potentially impacted community	Most limiting factors						
GATE	8	12	400	Overhead on jobs execution time Middleware stability Storage space						
CDSS	9	9	30 (mental diseases) + 50 (soft tissue tumours) in a short term.	For production: overhead for short jobs. For training the classifiers: computing time (around 1 week)						
gPTM3D	0.5	1	Tens (clinical researchers)	Jobs submission response time (in particular queuing delay) Lack of firewall-proof connectivity solution						
Simri3D	10	10	Several hundreds from the MR physics, medical and image processing communities	Correct handling of MPI jobs (too many errors today). Lack of scheduling time estimation.						
Bronze Std	2	4	Tens to hundreds once a proper interface has been set up.	Capacity to handle lot (hundreds) of jobs concurrently in an efficient manner. Currently the speed-up achieved is far from the expected bound. This is still under investigation but probably related to the bottlenecks of centralized RBs / UIs.						
Pharmcok.	9	10	Hundreds when the tool will prove stable and accurate enough.	Sufficient computing power dedicated to the application. In production it should represent 3 CPU years per year.						
GPS@	3	10	Difficult to estimate as the portal is opened anonymously to the biological community (probably several thousands)	Efficient handling of short jobs. Anonymous users authorization. Automatic replication.						
Xmipp_ML	5	10-15	Will be proposed to a NoE: hundreds.	CPU intensive Reliable MPI support High data throughput (data replication)						
SPLATCHE	4	5	Limited to a community of specialists in a short term	Sufficient CPUs availability (> 70) to compete with local cluster Multi-data jobs submission capability						
WISDOM	8	9	20 in a short term (2006). In the order of 100 later.	Reliability of services, especially WMS Security of data Number of CPUs						
GROCK	4	Tens	Thousands Deference: ECEI	es in the working nodes.						



# Focus on virtual screening



# Addressing neglected and emerging diseases

- Neglected and emerging diseases are major public health concerns in the beginning of the 21st century
  - Neglected diseases keep suffering lack of R&D

World Health

Organization

- Emerging diseases are a growing threat to world public health
- Both emerging and neglected diseases

require:

Early detection

•Emergence, resistance

Epidemiological watch

•Emergence, resistance

- Prevention
- Search for new drugs
- Search for vaccines



on the part of the World Health Organization concerning the legal status of any country, territory, city or area or of its authorities,

or concerning the delimitation of its frontiers or boundaries. Dotted lines on maps represent approximate border lines for which there may not yet be full agreement

© WHO 2006. All rights reserved

Communicable Diseases (CDS) World Health Organization



The grid added value for international collaboration on emerging and neglected diseases

 Grids offer unprecedented opportunities for sharing information and resources world-wide



Grids are unique tools for : -Collecting and sharing information (Epidemiology, Genomics) -Networking experts -Mobilizing resources routinely or in emergency (vaccine & drug discovery)



- Grids open new perspectives to *in silico* drug discovery
  - Reduced cost for R&D against neglected diseases
  - Accelerating factor for R&D against emerging diseases
- EGEE plays a pioneering role in exploring grid impact
  - Data challenge against malaria in the summer 2005
  - Data challenge against bird flu in April-May 2006





#### World wide In Silico Docking On Malaria

**BioinfoGRID** 



Disease	Endemic Countries	People at Risk (million)	Clinical Incidence/yr (million)	Deaths/yr (million)	Disease Burden (DALYs- million)
HIV/AIDS	180	5.900	40	2.8	86
Malaria	101	2.400	300-500	1.2	44.7
TB	211	1.987	8	1.6	35.4
African	36	60	0.3-0.5	0.05	1.5
trypanosomiasis					
Chagas Disease	21	100	16-18	0.01	0.7
Leishmaniasis	88	350	12	0.05	2
Filariasis	80	1.000	120		5.8
Schistosomiasis	76	500-600	140	0.01	1.7
Onchocerciasis	36	120	18		0.5
Leprosy	24		0.8		0.2



UNICEF/UNDP/World Bank/WHO Special Programme for Research and Training in Tropical Diseases (TDR)



### Where grids can help addressing neglected diseases

- Contribute to the development and deployment of new drugs and vaccines
  - Improve collection of epidemiological data for research (modeling, molecular biology)
  - Improve the deployment of clinical trials on plagued areas
  - Speed-up drug discovery process (in silico virtual screening)
- Improve disease monitoring
  - Monitor the impact of policies and programs
  - Monitor drug delivery and vector control
  - Improve epidemics warning and monitoring system
- Improve the ability of developing countries to undertake health innovation
  - Strengthen the integration of life science research laboratories in the world community
  - Provide access to resources
  - Provide access to bioinformatics services



#### First initiative: World-wide In Silico Docking On Malaria (WISDOM)

- Initial partners: Fraunhofer Institute, CNRS – IN2P3
- Significant biological parameters
  - two different molecular docking applications (Autodock and FlexX)
  - about one million virtual ligands selected
  - target proteins from the parasite responsible for malaria
- Significant numbers
  - Total of about 46 million ligands docked in 6 weeks
  - 1TB of data produced
  - Up 1000 computers in 15 countries used simultaneously corresponding to about 80 CPU years
  - Average crunching factor ~600



#### Number of running and waiting jobs vs time



Number of running and waiting jobs vs time



#### Deployment on EGEE infrastructure, wisdom.euegee.fr



Countries with nodes contributing to the data challenge WISDOM





## Strategies in result analysis

#### Results based on Scoring

Results based on match information

Results based on consensus scoring

•Results based on different parameter settings

Results based on knowledge on binding site

					Ter	minal					- 0
Ele	Edt	View	Terminal	Go	Help						
SEL	ECTED 1	MATCH	S: 1lee h2	. ndf	1	lee s	nin				
		+								+	
No	Itig.	Lig.	Ligand	- i	Rec.	Rec	Rec	- i	Rec.	Receptor	Opt.
1	Atom	ANO_	IA-Type	- i	Atom	AA.	Cha	(in)	AANo	IA-Type	Engy.
		+								+	
1.1	1 02	31	h_acc	- 1	water	1			138	h_don	-2.70
	1 N1	14	h_don	- 1	water	1			169	h_acc	-2.70
	1 014	22	h_don		water	1			120	h_acc	1-2.70
	1 C19	34	phenyl_cen	teri	CD2	TY8	IA.	- 1	77	phenyl_ring	1-0.70
	1 C16	9	phenyl_cen	terl	с	THR	IA.		217	amide	1-0.70
	1 C2	8	ch3_phe		CC	TYR.	A.		192	phenyl_center	r -0.70
	1 C	1 1	phenyl_cen	teri	CE2	TY8	A		192	phenyl_ring	-0.70
	1 C20	35	phenyl_rin	e I	CG	TY8	A.		77	phenyl_center	r -0.70
	1 C19	34	phenyl_cen	ter]	CD1	ILE	IA.		123	ch3_phe	1-0.70
	1 C19	34	phenyl_cen	terl	062	ILE	IA.		32	ch3_phe	1-0.70
	1 C13	1 5	phenyl_rin	εI	CC.	PHE	A.	- 1	294	phenyl_center	r -0.70
	1 C10	4	phenyl_rin	¢	CG	PHE	A.		294	[phenyl_center	r -0.70
	1 C	1 1	phenyl_cen	ter	061	VAL.	A.		78	ch3_phe	-0.70
	1 C	1	phenyl_cen	ter	CE1	PHE	1A.		294	phenyl_ring	-0.70
	1 C16	9	phenyl_cen	ter	CE	MET	A.		15	ch3_phe	-0.70
	1 N3	29	h_don		0	CLY	1A.		216	h_acc	1-4.70
	1 014	22	h_don	- 1	OD1	LASP	IA.		214	h_acc	1-4.70
	1 01	16	h_acc	- 1	CHE	TY8	1A.		192	h_don	1-4.70
	1 026	1.9	h_acc		26	VAL.	A.		78	h_don	-4.70
1.1	11N16	1 20	h_don		0	CLY	LA.		36	lh_acc	1-4.70







Credit: V. Kasam

**Fraunhofer Institute** 

#### Top 10 compounds by scoring

BioinfoGRID



Potentially new inhibitors: guanidino compounds

Top scoring, good binding mode, interactions to key residues
NETTAB, 10/7/06 - 23



# Compounds for Molecular Dynamics: Guanidino

#### compounds





Note: Satisfied all criteria, good binding mode, interactions to key residues, good score, appropriate descriptors.

#### Compounds from consensus scoring



**BioinfoGRID** 

#### NETTAB, 10/7/06 - 25





# From virtual docking to virtual screening

**BioinfoGRID** 



NETTAB, 10/7/06 - 27

# The next steps



- Docking step still requires a lot of manual intervention
  - Goal: reduce as much as possible the time needed for experts to analyze the results
  - Task: improve output data collection and post-docking analysis
  - Contribution from CNR-ITB, within the framework of EGEE-II
- The next step after docking is Molecular Dynamics
  - Goal: grid-enable the reranking of the best hits
  - Task: deploy Molecular Dynamics computations on grid infrastructures
  - Contribution from CNRS-IN2P3, within the framework of BioinfoGRID
- Beyond virtual screening, the long term vision: building a grid for malaria
  - To provide services to research labs working on malaria
  - To collect and analyze epidemiological data



# A grid for malaria



Aventis, Hospitals in subsaharian Africa,

NETTAB, 10/7/06 - 29





- WISDOM-II is the second large scale docking deployment against neglected diseases
- Biological goals
  - Validation of virtual vs in vitro screening
  - Virtual docking on new malaria targets
    - New targets from Univ. Pretoria, Univ. Los Andes, CEA Grenoble, Univ. Modena
  - New compound libraries
    - Thai library of compounds
  - Possible extension to other neglected diseases
    - Contacts with Univ. Glasgow
- Grid goals:
  - Improve the user interface, the job submission system and the postprocessing (BioinfoGRID, EGEE, Embrace)
  - Test the infrastructure at a larger scale (100 -> 500 CPU years)
  - Test the deployment on several infrastructures: Auvergrid, EGEE, EELA



WISDOM	MD reranking In vitro t						testii	ng	Further processing										
WISDOM II	Preparation Deploy ment					Analysis					MD reranking								
Avian Flu	Analysis					MD reranking					Further processing								
Avian Flu II						Preparation Deploy ment				Analysis									
Month	6 06	7 06	8 06	9 06	10 06	11 06	12 06	1 07	2 07	3 07	4 07	5 07	6 07	7 07	8 07	9 07	10 07	11 07	1







- Life science applications are currently running on grid infrastructures
  - Example of EGEE biomedical applications
- Most of these applications are compute intensive
  - Emergence of data grid applications
- Very large scale virtual docking achieved
  - More details from H.C. Lee talk
- Exciting perspectives to develop in silico drug discovery
  - Collaboration of partners and EC projects
  - WISDOM-II, further step towards a malaria grid





- Joint tutorial co-organized by Embrace, EGEE and BioinfoGRID
- Location: Clermont-Ferrand
- Dates: 2 days in October 2006
- Goal: get bioinformaticians to understand the grid to the extent that they can install and use it
- Contact: Florence Jacq, fjacq@clermont.in2p3.fr