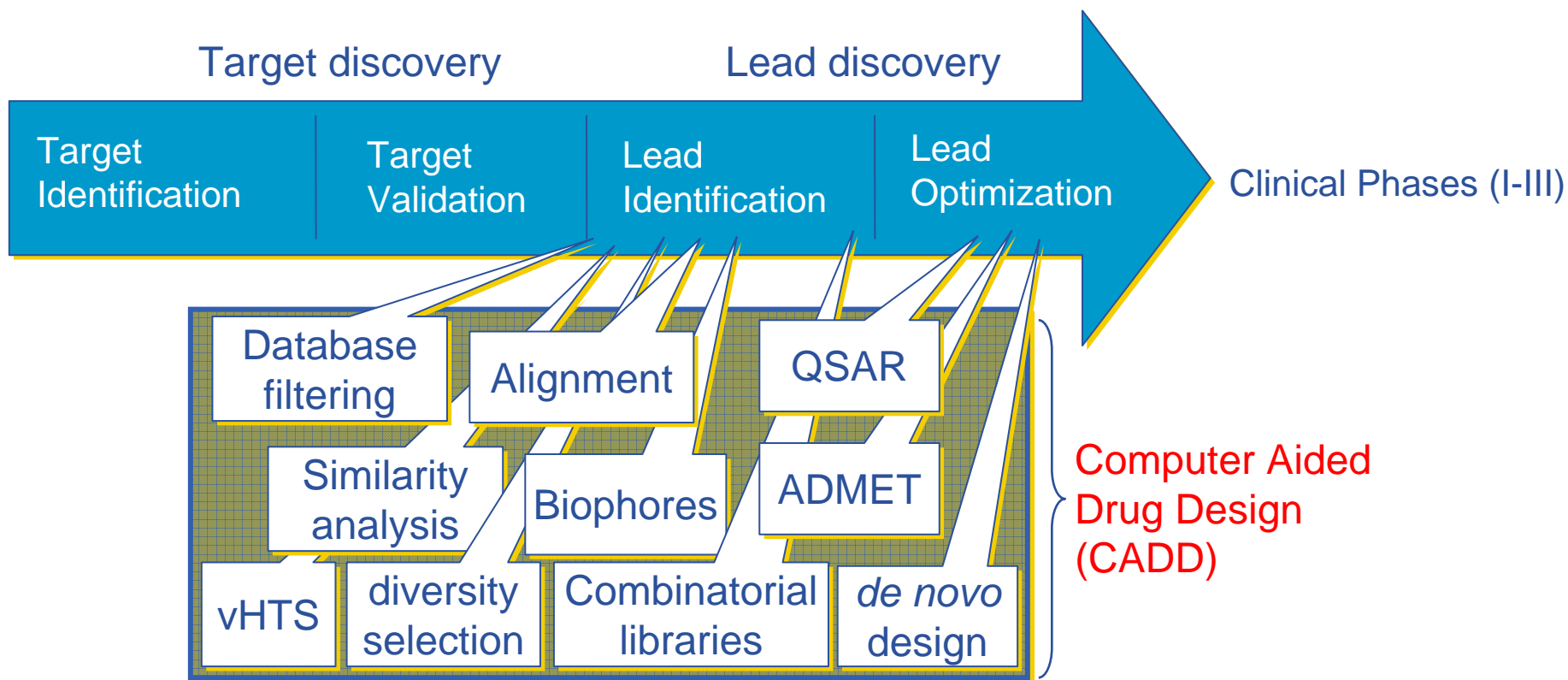


# Grid enabled *in silico* drug discovery

*Vincent Breton*

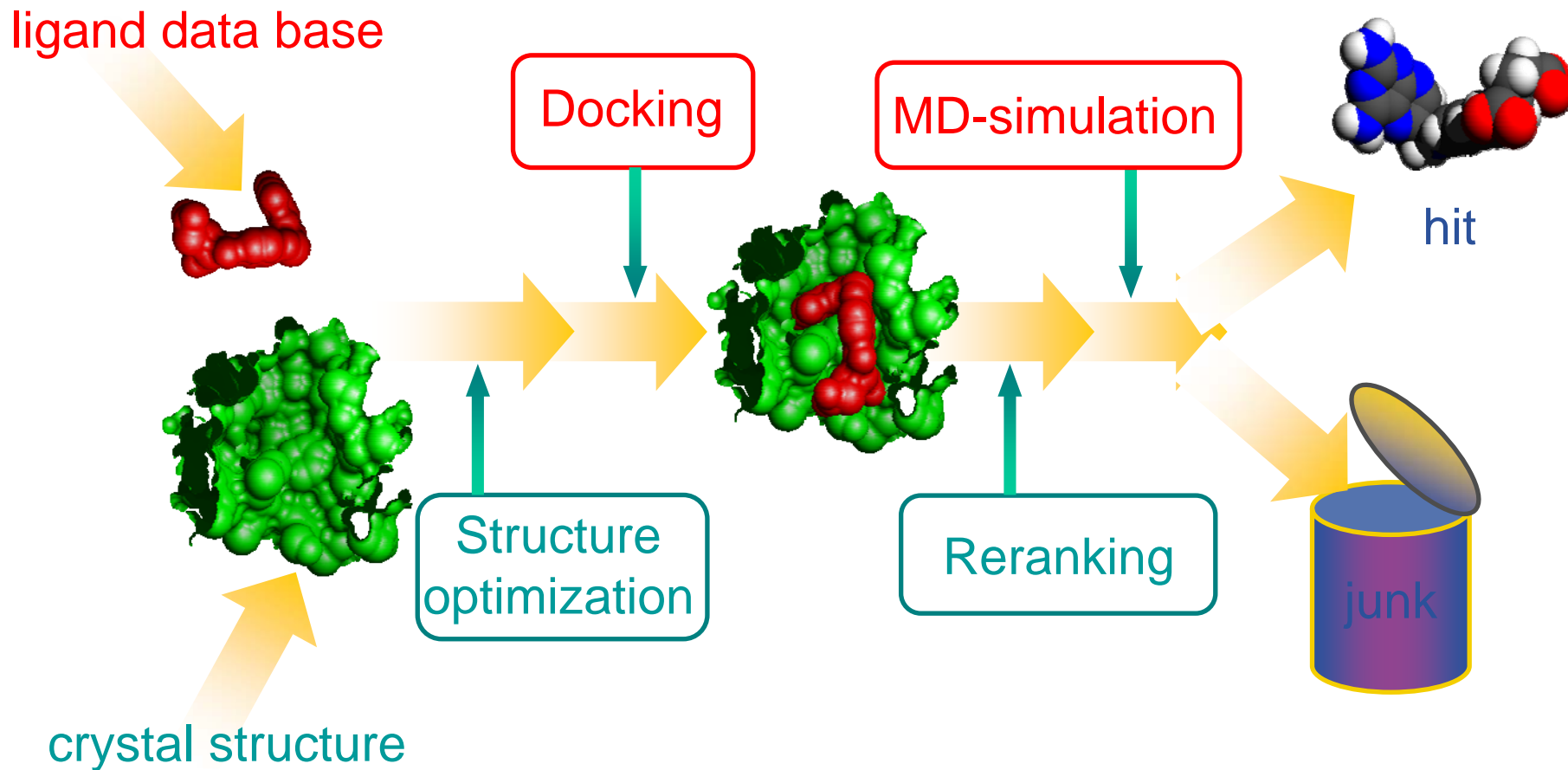
*CNRS/IN2P3*


*Credit for the slides: N. Jacq*



Duration: 12 – 15 years, Costs: 500 - 800 million US \$

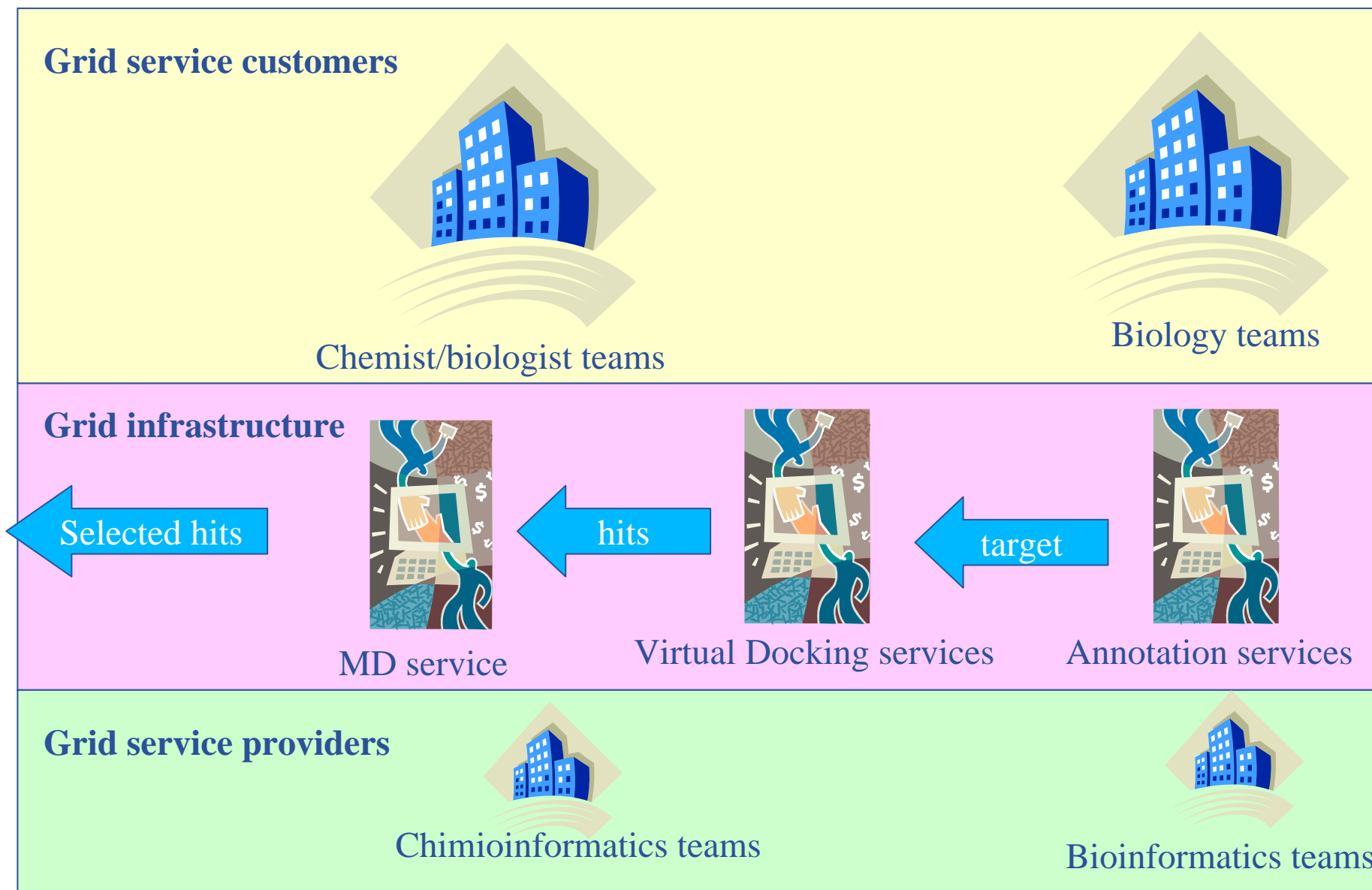
- **28 million compounds currently known**
- **Drug company biologists screen up to 1 million compounds against target using ultra-high throughput technology**
- **Chemists select 50-100 compounds for follow-up**
- **Chemists work on these compounds, developing new, more potent compounds**
- **Pharmacologists test compounds for pharmacokinetic and toxicological profiles**
- **1-2 compounds are selected as potential drugs**



- **Enable scientists to quickly and easily find ligands binding to a particular target protein**
  - growth of targets number
  - growth of 3D structures determination (PDB database)
  - growth of computing power
  - growth of prediction quality of protein-compound interactions
  
- **Experimental screening very expensive : difficult for academic or small companies**
  
- **Enrichment =**  $\frac{\text{Actives molecules}}{\text{Tested molecules}}$  

- **Target identification and validation**
  - Volume of molecular biology data is exponentially increasing
  - Grid added value: interoperability, sharing of data content and tools
  
- **Large scale virtual screening to select the most promising compounds**
  - Distributed computing
  - output data management
  
- **Molecular dynamics to further assess selected compounds**
  - Parallel computing

- **A grid infrastructure uses an identified set of resources properly administered behind firewalls**
- **Grid infrastructures vs pervasive grids**
  - Large scale docking on pervasive grid already achieved (Grid.org, Decryphon, World Community Grid)
    - Centralized job submission and data management
    - Limited security model
    - No output data distribution (web portal)
    - Limited quality of service (no user support)
- **Grid infrastructures vs clusters**
  - Sharing of computing resources
  - Data management: distribution/replication of data
  - Sharing of services (participating groups bring their expertise)





- **Scientific objectives**

- start enabling *in silico* drug discovery in a grid environment to address the deadliest infectious disease on earth: malaria
- Demonstrate to the research communities active in the area of drug discovery the relevance of grid infrastructures

- **Goals of the first “data challenge” (July - September 2005)**

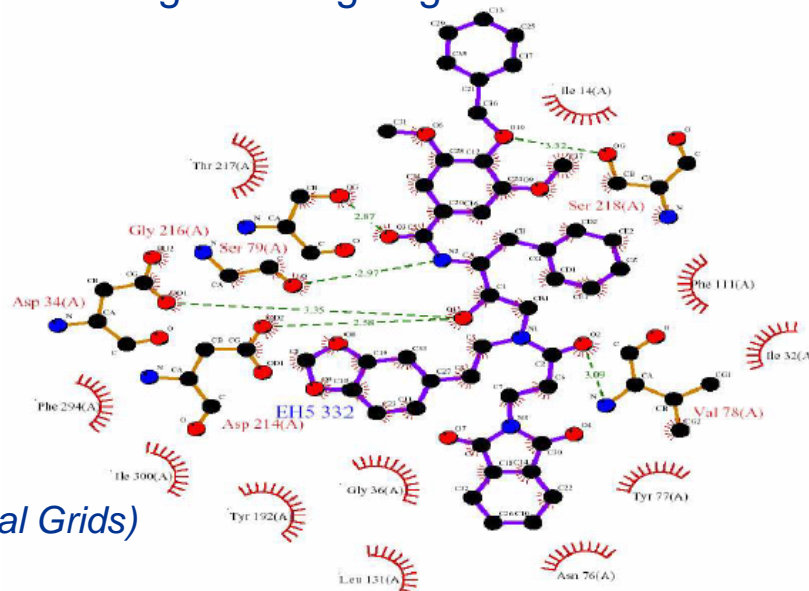
- Biological goal : Proposition of new inhibitors for a family of proteins produced by *plasmodium falciparum*
- Biomedical informatics goal : Deployment of *in silico* virtual screening on the grid
- Grid goal : Deployment of a CPU consuming application generating large data flows to test the grid infrastructure and services.

- **Partners**

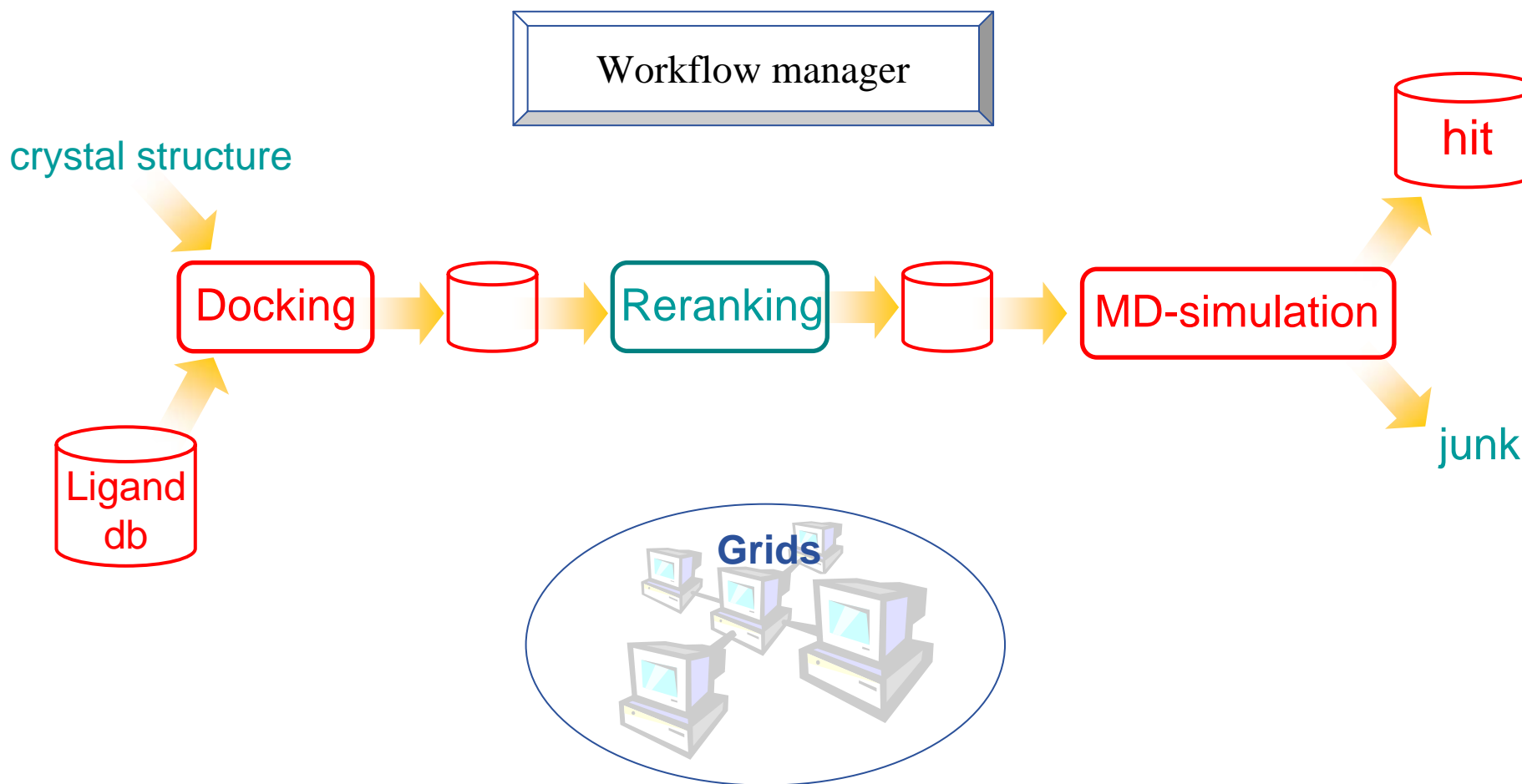
- Fraunhofer SCAI
- CNRS/IN2P3
- CMBA (*Center for Bio-Active Molecules screening*)

**representing different projects:**

- EGEE (EU FP6)
- Simdat (EU FP6)
- Instruire and Campus Grid (French and German Regional Grids)
- Accamba project (french ACI project)

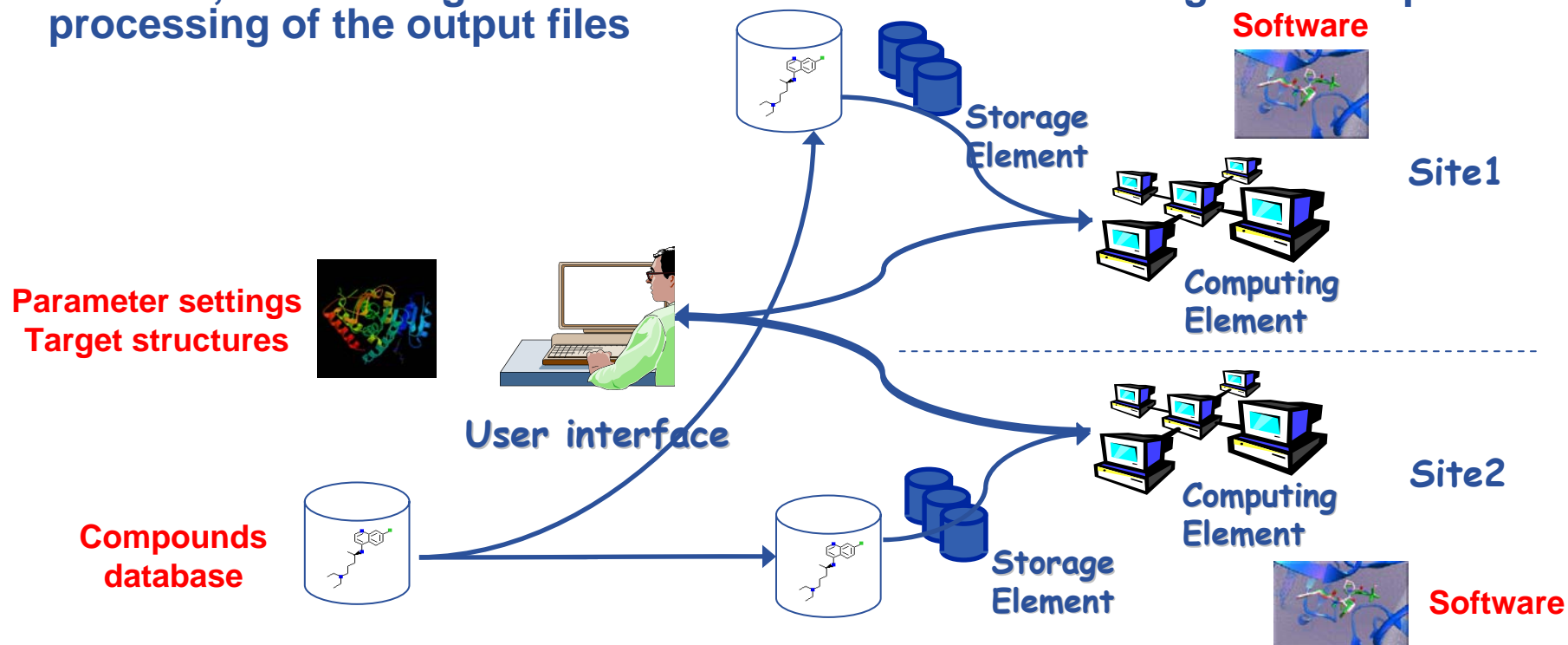


- Deployment of a virtual screening workflow on grid infrastructures



- **Biological information**
  - **Plasmepsin** is a promising aspartic protease target involved in the hemoglobin degradation of *P. falciparum*. 5 different structures are prepared (PDB source)
  - **ZINC** is an open source library of 3,3 millions selected compounds. They are made available by chemistry companies and are ready to be used
- **Biomedical informatics tools**
  - **Autodock** is free for academic, with grid based empirical potential and flexible docking via MC search and incremental construction
  - **FlexX** is licensed required, available for this data challenge during 1 week, with Boehm potential and fragment assembly energy function
- **Grid tools**
  - `wisdom_env` is an environment for an automatic, optimized and fault tolerance workflow using the grid resources and services
  - The biomedical VO will be the infrastructure with dedicated/no-dedicated resources

- Docking is easily distributed once the compound database is available on the grid nodes. Each computing element computes docking probability for a different sample of ligands
- In a first step, docking scores are returned to the user and compared on its local machine.
- Later on, data management services can handle the storage and the post-processing of the output files



# Results of the preliminary tests

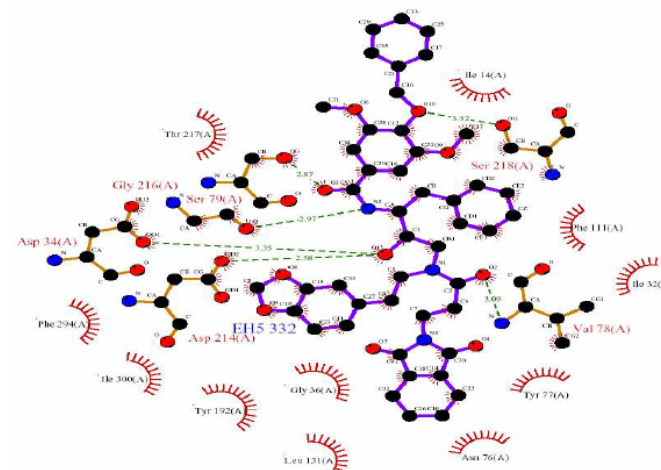
- Docking application deployed since the summer 2004
- +30,000 jobs since January 2005
- Tests performed with the software Autodock on the biomedical VO

	<b>100,000 compounds 500 jobs</b>
<b>Total CPU time for jobs</b>	<b>6 months CPU</b>
<b>User script time</b>	<b>40 h</b>
<b>Gain of time for the user</b>	<b>150</b>
<b><i>CPU time for 1 job</i></b>	<b><i>9h</i></b>
<b><i>Input and output transfer time between SE and CE for 1 job</i></b>	<b><i>2.5 mn</i></b>
<b><i>Waiting time for 1 job due to the grid</i></b>	<b><i>30 mn</i></b>
<b><i>Resubmitted Jobs</i></b>	<b><i>16</i></b>
<b><i>Aborted jobs %</i></b>	<b><i>3%</i></b>

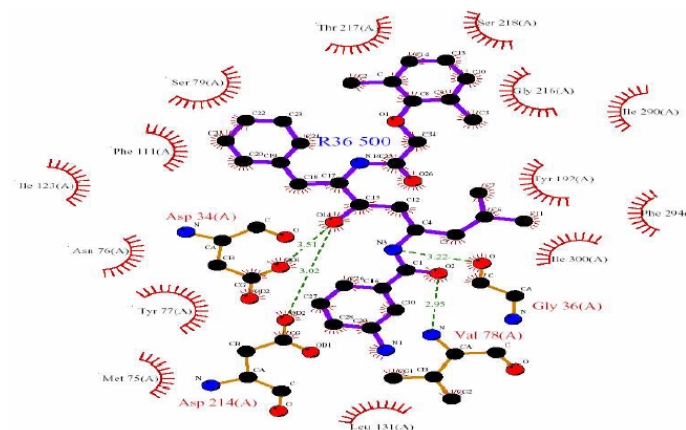
	<b>Scenario 1</b>
<b>Duration</b>	<b>3 weeks</b>
<b>CPU time</b>	<b>80 years CPU</b>
<b>Grid performance</b>	<b>70%</b>
<b>Number of CPU</b>	<b>2,000</b>
<b>Number of grid jobs (20h)</b>	<b>30,000</b>
<b>Storage</b>	<b>2*6 TB</b>
<b>Docking workflow description</b>	
<b>Number of compounds</b>	<b>500,000</b>
<b>Number of parameters settings</b>	<b>4</b>
<b>Objective</b>	<b>Selection of the best hits with short analysis</b>

- FlexX running time : 1 mn
- F. output size : 1MB
- F. job output size : 1.2GB
- F. job compressed output size : 250MB
- Autodock running time : 2.5 mn
- A. output size : 1MB
- A. job output size : 0,5GB
- A. job compressed output size : 100MB

- Post filtering
- Clustering of similar conformations
- Checking pharmacophoric points of each conformation
- Doing statistics on the score distribution
- Re-ranking for interesting compounds
- Sorting and assembly of data



Ligand plot of 1LF3 (plasmepsin II) with inhibitor EH5 332



Ligand plot of 1LEE (Plasmepsin II) with inhibitor R36 500

- **The best hits found by post-treatment will be published and available on a permanent grid storage via a portal**
  - Experimental screening of the most promising hits
- **A knowledge space will be progressively build around these results**
  - to extract and process the most interesting information
  - to enrich the data with the results found later by other *in silico* drug discovery processes
- **The *in silico* drug discovery will be further extend**
  - to include more precise molecular dynamics computations using quantum chemistry software like NAMD



# From drug discovery to drug delivery

- **Drug discovery is about finding new drugs**
- **However, the best drugs are useful provided they are made available to the sick**
- **Drug delivery is a huge challenge for developing countries**
  - Lack of healthcare infrastructures
  - Lack of resources to buy drugs
  - Lack of education to deliver them
  - Lack of information on drug efficiency
- **For drug delivery, grids have a real added value**
  - To collect data in endemic areas
  - To provide data and tools to endemic areas (local research, training)

# Grids for neglected diseases of the developing world

In silico drug discovery process  
(EGEE, SwissBioGRID, ...)

Clermont-Ferrand

SCAI Fraunhofer

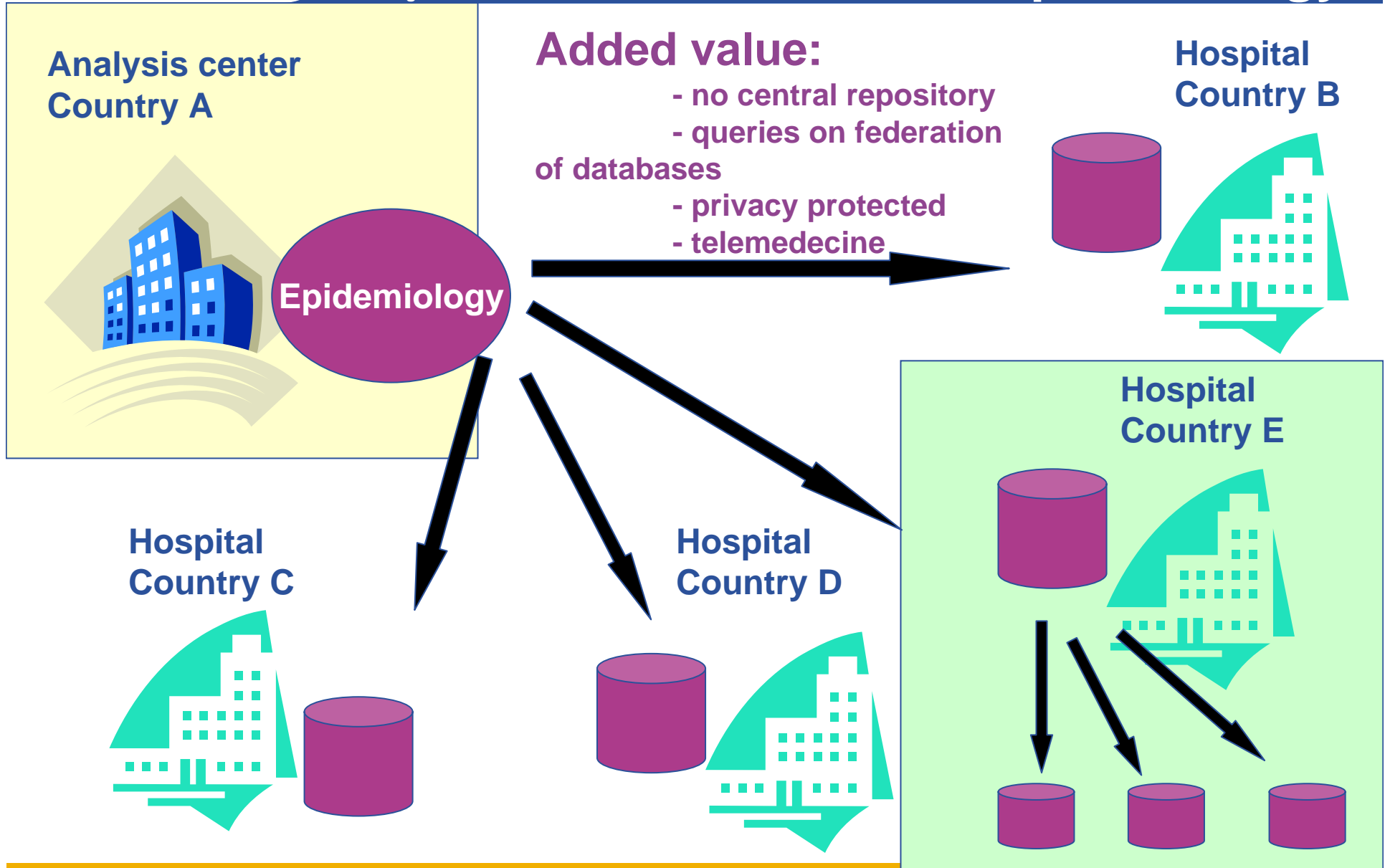
Swiss Biogrid consortium

Support to local  
centres in plagued  
areas (data collection,  
genomics research,  
clinical trials and  
vector control)

Local research centres  
In plagued areas

## The grid impact :

- Computing and storage resources for genomics research and in silico drug discovery
- cross-organizational collaboration space to progress research work
- Federation of patient databases for clinical trials and epidemiology in developing countries



Pharmaceutical laboratory /  
International organization  
Country A

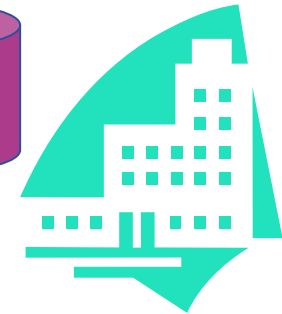
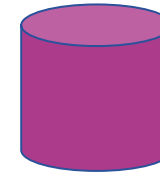


Drug / Vaccine  
assessment

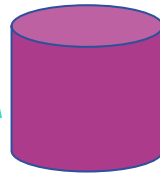
**Added value:**

- no central repository
- queries on federation of databases
- privacy protected-

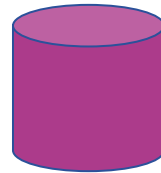
Hospital  
Country B



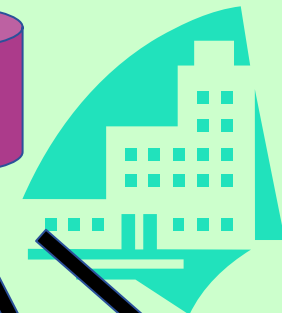
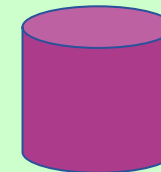
Hospital  
Country C



Hospital  
Country D



Hospital  
Country E



- **Grid enabled telemedecine for medical development**
  - Development of neurosurgery in poverty regions of western China
  - Ophthalmology in Burkina-Faso
    - Collaboration with Schiphra dispensary (Ouagadougou, Burkina Faso)

**Collaboration: NPO Chain of Hope, n°9 Hospital Shanghai (neurosurgery unit), Chuxiong Hospital (Yunnan), CNRS-IN2P3, Clermont-Ferrand hospitals**

Goal: improve patient follow-up by french clinicians

Method: grid-enabled telemedicine web application



- **Grid technologies promise to change the way organizations tackle complex problems by offering unprecedented opportunities for resource sharing and collaboration**
- **Grids should provide the services needed for *in silico* drug discovery**
- **Applied to world health development, grids should also**
  - Help monitor epidemics
  - Strengthen R&D on neglected diseases
  - Grant easier access to eHealth
  -
- **We are looking for joint pilot projects with a pharmaceutical lab**
  - Develop a grid-enabled drug discovery pipeline for malaria
  - Build a federation of databases to address 1 infectious disease (epidemiology, clinical trials, vector control)
  - Study grid added value for drug delivery