



# **Introduction to the Grid**

**Peter Kacsuk**  
**MTA SZTAKI**  
**[www.lpds.sztaki.hu](http://www.lpds.sztaki.hu)**



# *Agenda*

- From Metacomputers to the Grid
- Grid Applications
- Job Managers in the Grid - Condor
- Grid Middleware – Globus
- Grid Application Environments



# Grid Computing in the News

BusinessWeek

online

GRID COMPUTING PLANET.COM

EARTHWEB

BW MAGAZINE

DAILY BRIEFING

INVESTING

GLOBAL BUSINESS

TECHNOLOGY

SMALL BUSINESS [Internet.com](#)

management Networking & Communications Web Development Hardware & Systems Software Development

JUNE 3, 2002

[ridcomputingplanet.com](#) : News

FREE TECH NEWSLETTERS

YOUR EMAIL

Sign Up

Grid Computing Planet Text

INFORMATION TECHNOLOGY

## Who Needs Supercomputers?

Grid software lets companies tap other machines in their network for

## Researchers Achieve Production Grid Breakthrough

By [Paul Shread](#)



BW MAGAZINE

U.S. EDITION

[Full Table of Contents](#)

[Cover Story](#)

[Up Front](#)

[Readers Report](#)

AN MIT ENTERPRISE TECHNOLOGY REVIEW

THE TR PATENT SCORECARD 2

Search the top patentees in major high-tech indu

Emerging Technologies and Their Impact

Home

Browse by Topic

Focus On

Current Issue

Archive

Opinions/Forums

Events

TR STORE

PHYSICS TODAY.org

Gradscl

Search

FEATURE ARTICLES PHYSICS UPDATE LETTERS SEARCH & DISCOVER

Feature Article

TABLE OF CONTENTS

PAST CONTENTS

LINKS TO PHYSICS TODAY ADVERTISERS

PLACE AN AD

BUYER'S GUIDE

ABOUT US

CONTACT US

## The Grid: A New Infrastructure for 21st Century Science

As computer networks become cheaper and more powerful, a new computing paradigm is poised to revolutionize science and engineering.

[Ian Foster](#)

Grid Computing

By M. Mitchell Waldrop May 2002

Hook enough computers together and what

ZDNet News

Technology News Now

Page One

Hardware

## Grid computing earns a living

By Stephen Shankland  
Special to ZDNet News  
February 21, 2002, 4:00 AM PT

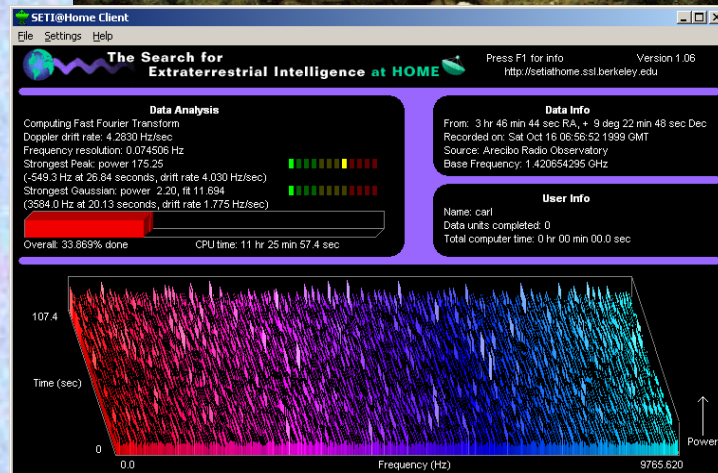
Is grid computing in your company's future?



# Real World Distributed Applications

## • SETI@home

- 3.8M users in 226 countries
- 1200 CPU years/day
- 38 TF sustained (Japanese Earth Simulator is 40 TF peak)
- 1.7 ZETAflop over last 3 years (10<sup>21</sup>, beyond peta and exa ...)
- Highly heterogeneous: >77 different processor types





# Progress in Grid Systems

**Supercomputing  
(PVM/MPI)**

**Network  
Computing (sockets)**

**Clusters**

**Cluster  
computing**

**High-throughput  
computing**

**High-performance  
computing**

**Web Computing  
(scripts)**

**OO Computing  
(CORBA)**

**Client/server**

**Object Web**

**Condor**

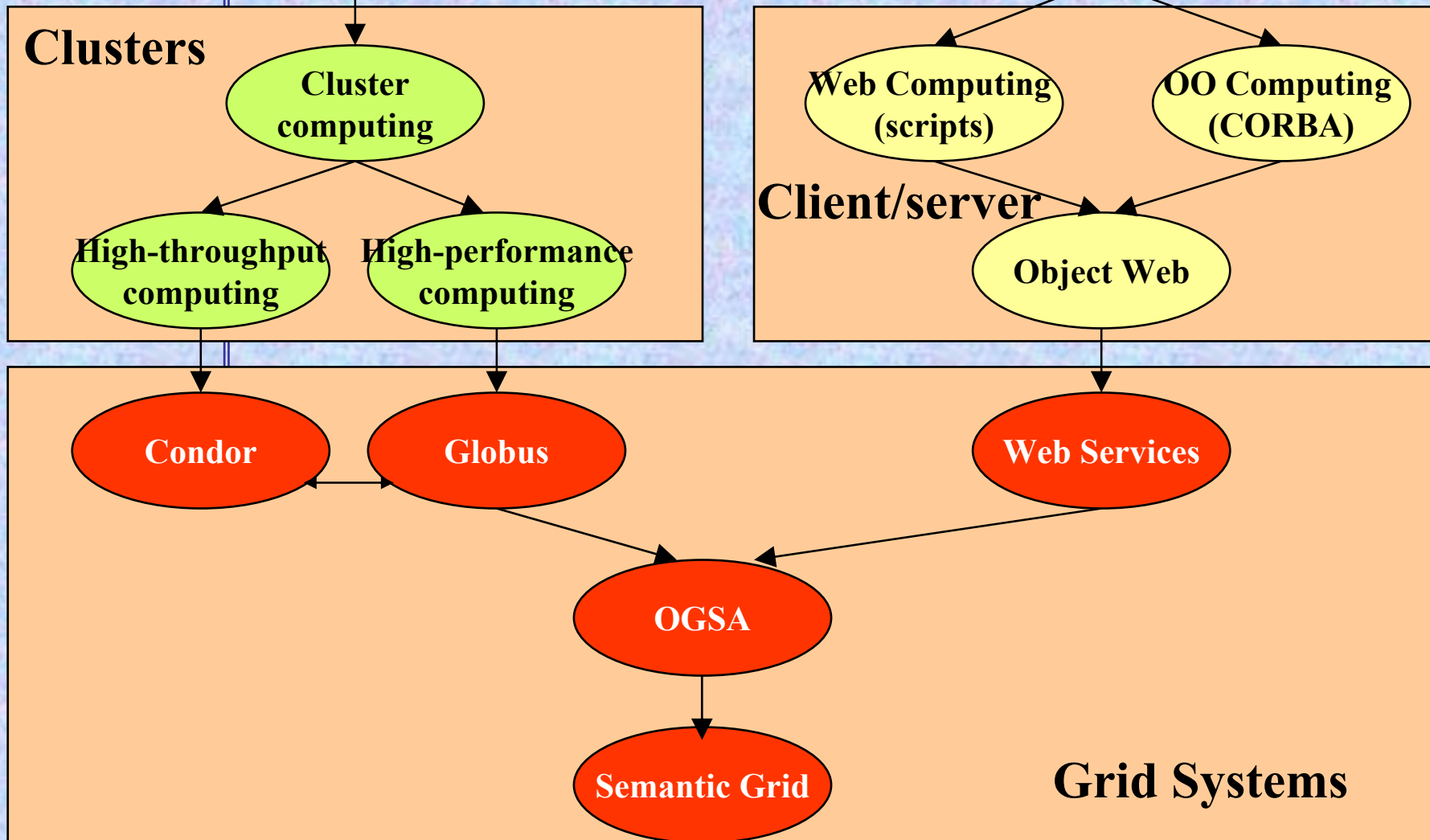
**Globus**

**Web Services**

**OGSA**

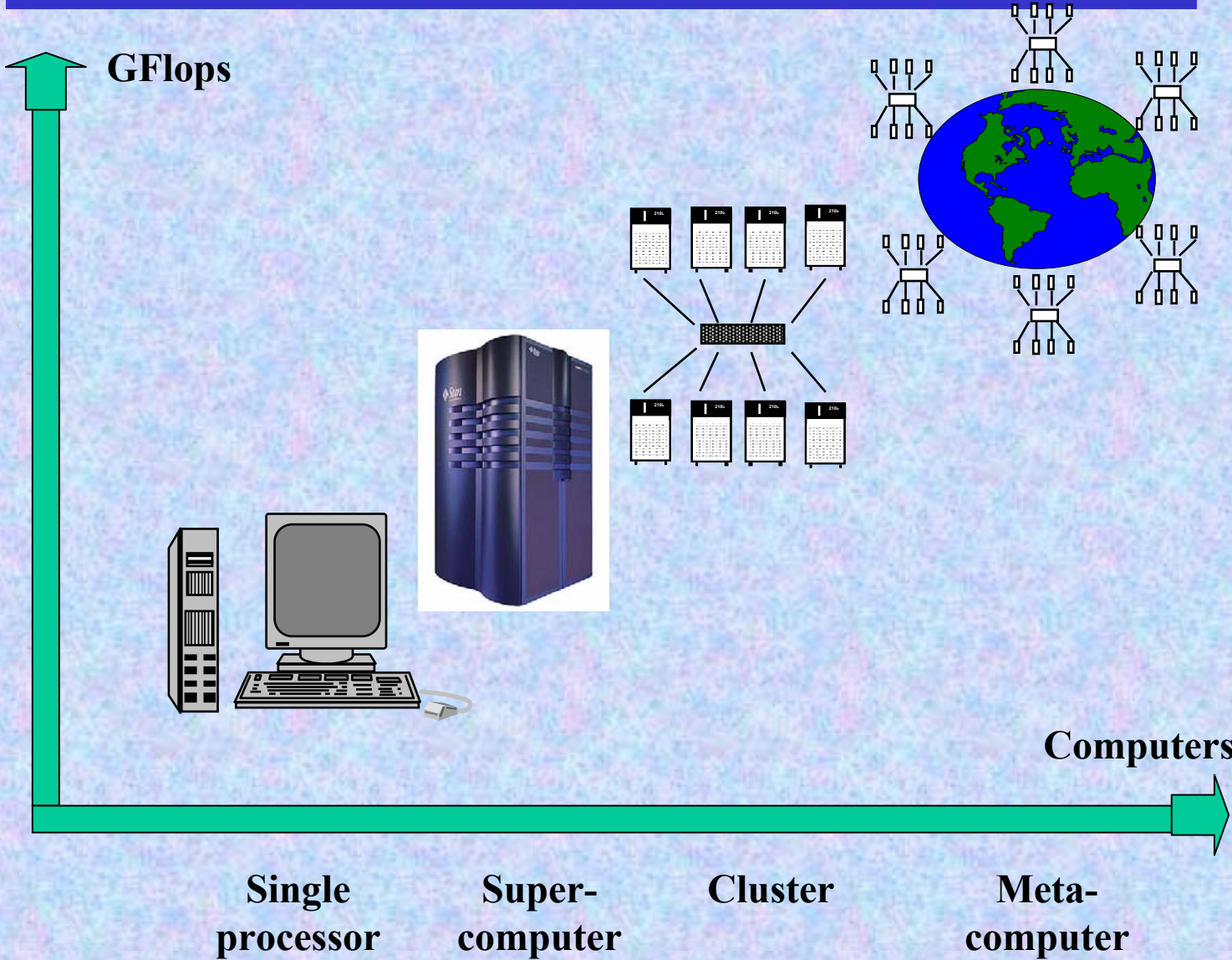
**Semantic Grid**

**Grid Systems**





# *Progress to the Grid*



# *A metaszámítógépek megalkotásának eredeti motivációi*

- **Az un. nagy kihívást jelentő problémák** megoldása heteket sőt hónapokat vesz igénybe még a szuperszámítógépeken is



- **Különböző szuperszámítógépeket és klasztereket** kellett összekapcsolni **távolsági hálózatokkal** annak érdekében, hogy a fenti problémákat ésszerű időn belül meg lehessen oldani

## *A metaszámítógép eredeti jelentése*

Metaszámítógép

=

Szuperszám.  
technológia

+

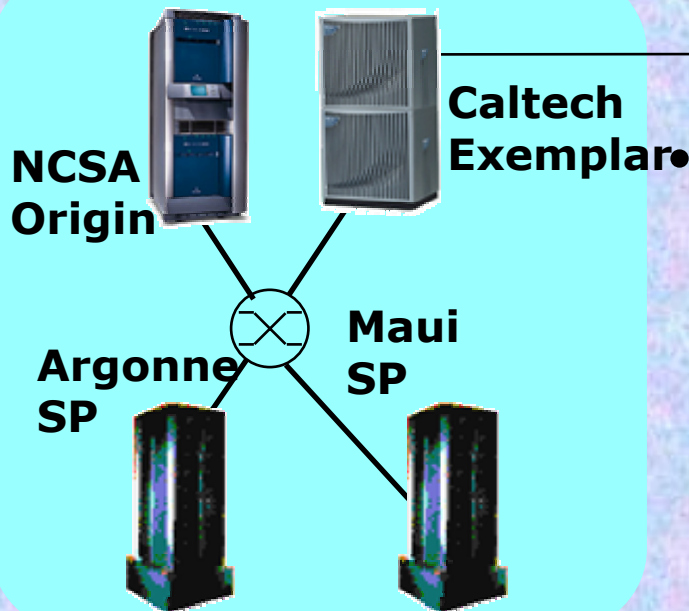
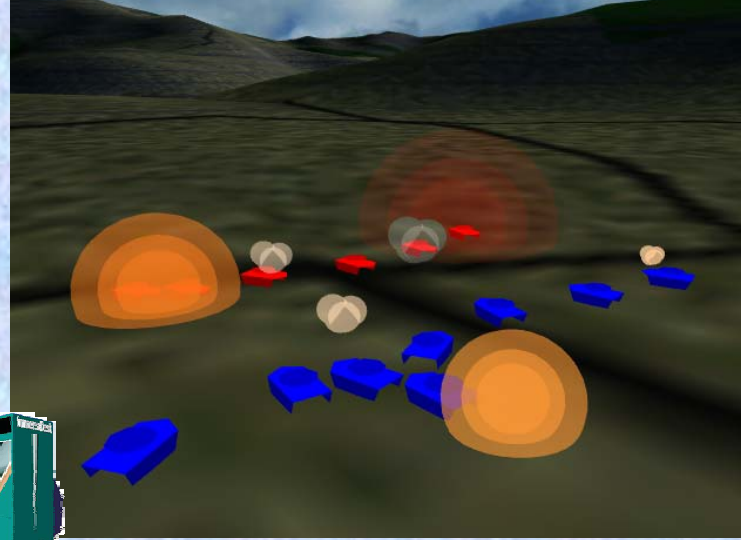
Távolsági  
hálózat

## A metaszámítógép eredeti célja

- **Nagyobb teljesítményt** elérni, mint az egyedi szuperszámítógépek/klaszterek tudnak biztosítani



# *Distributed Supercomputing*

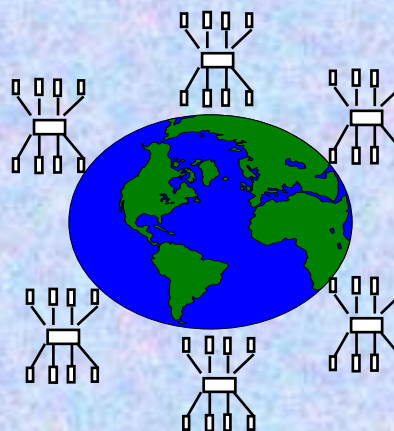
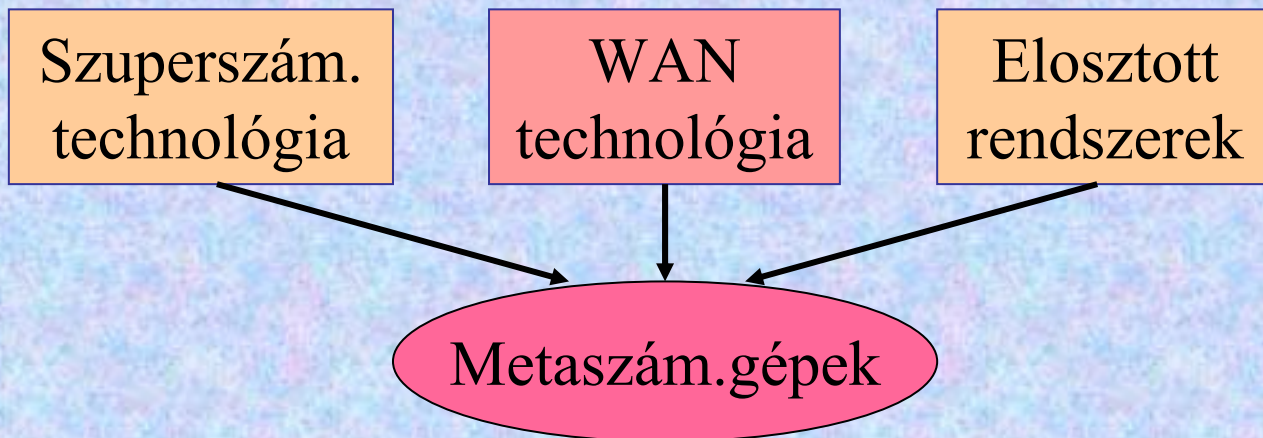


## Issues:

- Resource discovery, scheduling
- Configuration
- Multiple comm methods
- Message passing (MPI)
- Scalability
- Fault tolerance

**SF-Express Distributed Interactive Simulation**  
Caltech, USC/TST

# *Második lépés a metaszámítógépek felé*





## *Mi is a metaszámítógép?*

- A metaszámítógép olyan számítógépek együttese, amelyek
  - **heterogének** minden szempontból
  - földrajzilag **elosztottak**
  - **távolsági hálózattal** vannak összekötve
  - **egyetlen komputer képét** alkotják (**SSI**)
- Metaszámítás jelentése:
  - hálózat alapú
  - elosztott **szuperszámítógép technológia**

## *További motivációk a metaszámítógépek megalkotására*

- A távolsági hálózattal elérhető számítási és egyéb **erőforrások hatékonyabb kihasználása**



- **Különböző számítógépeket** kell távolsági hálózattal összekötni a **szabad ciklusok kihasználására**
- **Különböző speciális készülékeket** kell távolsági hálózattal összekötni **kollaboratív munka biztosításához**

## *Motivációk a GRID infrastruktúra megalkotására*

- Olyan számítási és adatfeldolgozási **GRID**-et létrehozni, amely hasonlóan széleskörű, mint az információ elérése a weben

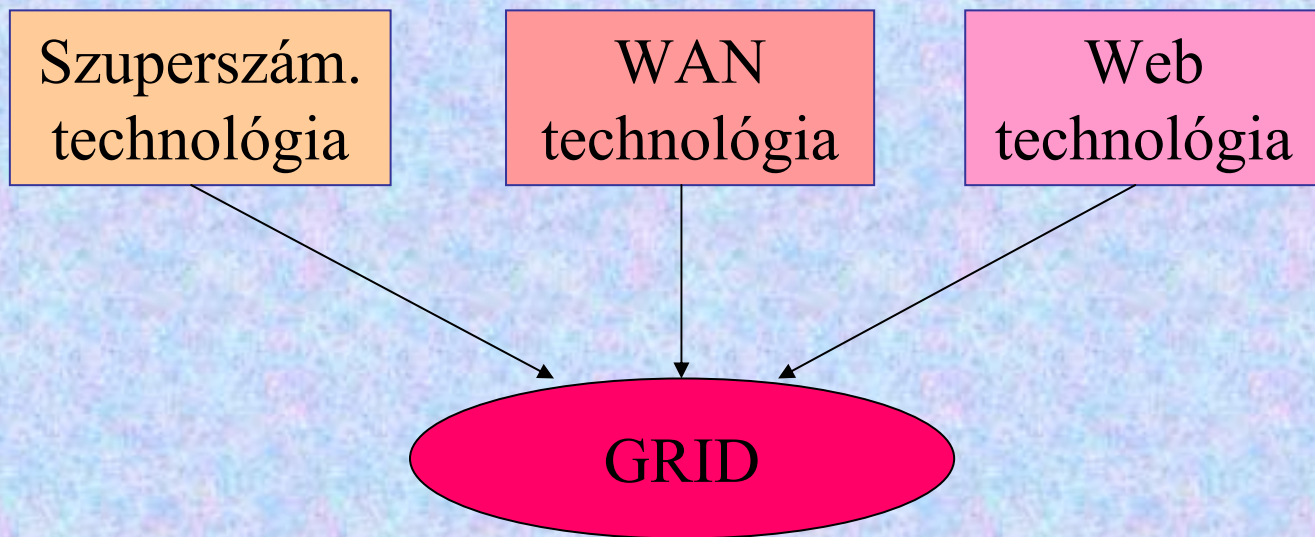


- **Bármely** számítógépet/készüléket célszerű összekötni távolsági hálózattal, hogy **univerzális feldolgozó kapacitást** nyerjünk



- **Grid** = általánosított metaszámítógép technológia

# *A GRID megalkotásához vezető technológiák*





## *Mi is a GRID?*

- A GRID olyan számítógépek, tárolóegységek és egyéb készülékek együttese, amelyek
  - **heterogének** minden szempontból
  - földrajzilag **elosztottak**
  - **távolsági hálózattal** vannak összekötve
  - **egyetlen komputer képét** alkotják (SSI)
- Az általánosított metaszámítási technológia (GRID) jelentése:
  - hálózat alapú
  - elosztott **adatfeldolgozási technológia**

# *A GRID alkalmazási területei*

- High throughput computing
- Virtuális laboratórium
  - Kollaboratív tervezés
- Adat intenzív alkalmazások
  - adatbányászat, részecskefizika
- Földrajzi információs rendszerek

- Távoli jelenlét (tele-immersion)
- Vállalati rendszerek

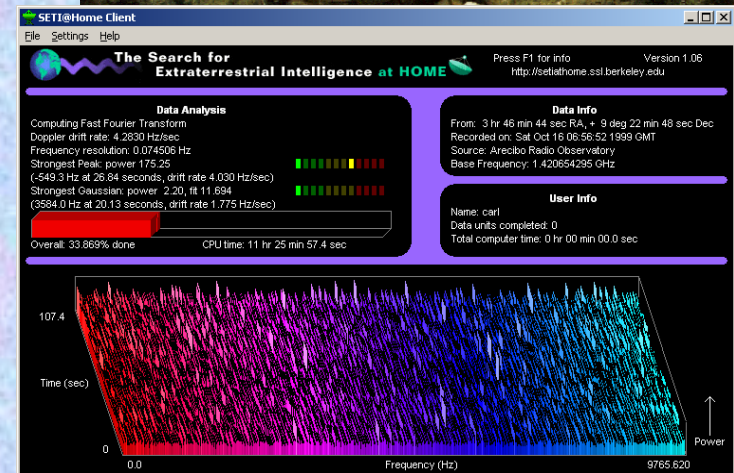




# Real World Distributed Applications

## • SETI@home

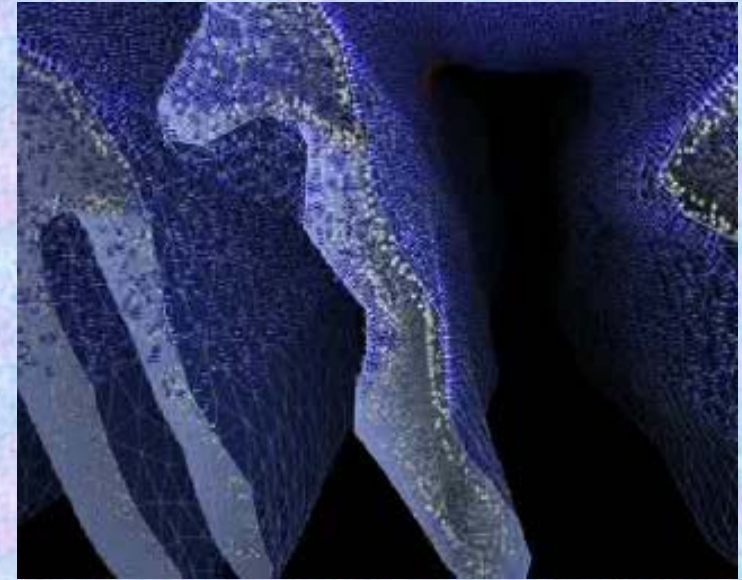
- 3.8M users in 226 countries
- 1200 CPU years/day
- 38 TF sustained (Japanese Earth Simulator is 40 TF peak)
- 1.7 ZETAflop over last 3 years (10<sup>21</sup>, beyond peta and exa ...)
- Highly heterogeneous: >77 different processor types



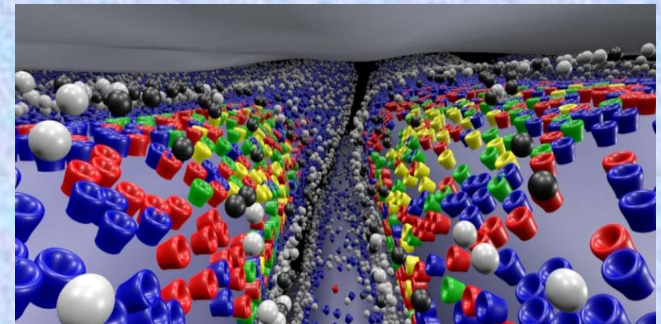


## **MCell -- Simulation of neuromuscular synaptic transmission**

- Uses Monte Carlo diffusion and chemical reaction algorithm in 3D to simulate complex biochemical interactions of molecules
- Molecular environment represented as 3D space in which trajectories of ligands against cell membranes tracked
- Ultimate Goal: A complete molecular model of neuro-transmission at level of entire cell



**MCell Animation**



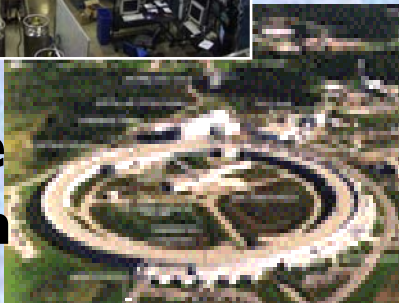


# Online Instruments

Advanced Photon Source

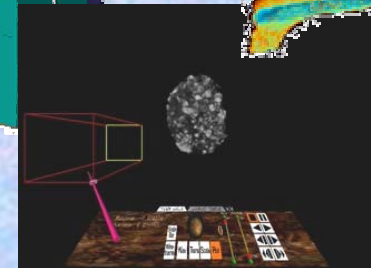
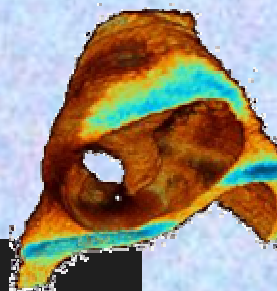
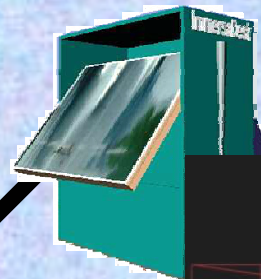


real-time collection



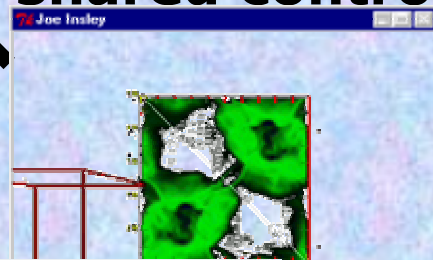
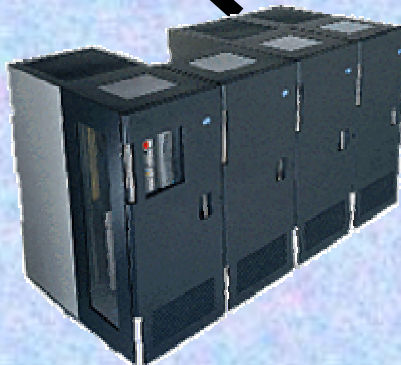
tomographic reconstruction

wide-area dissemination



desktop & VR clients with shared controls

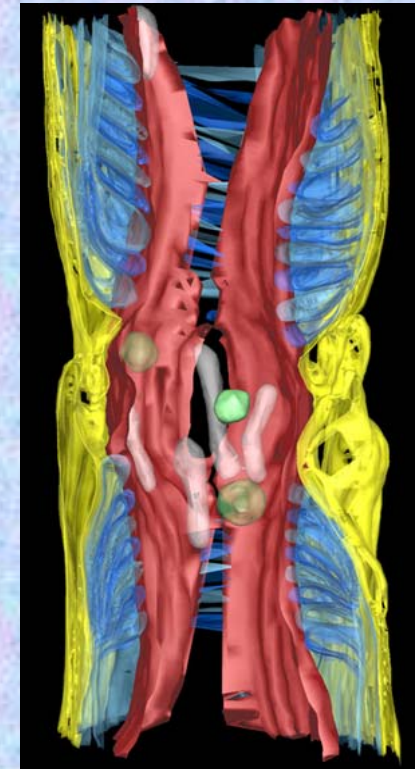
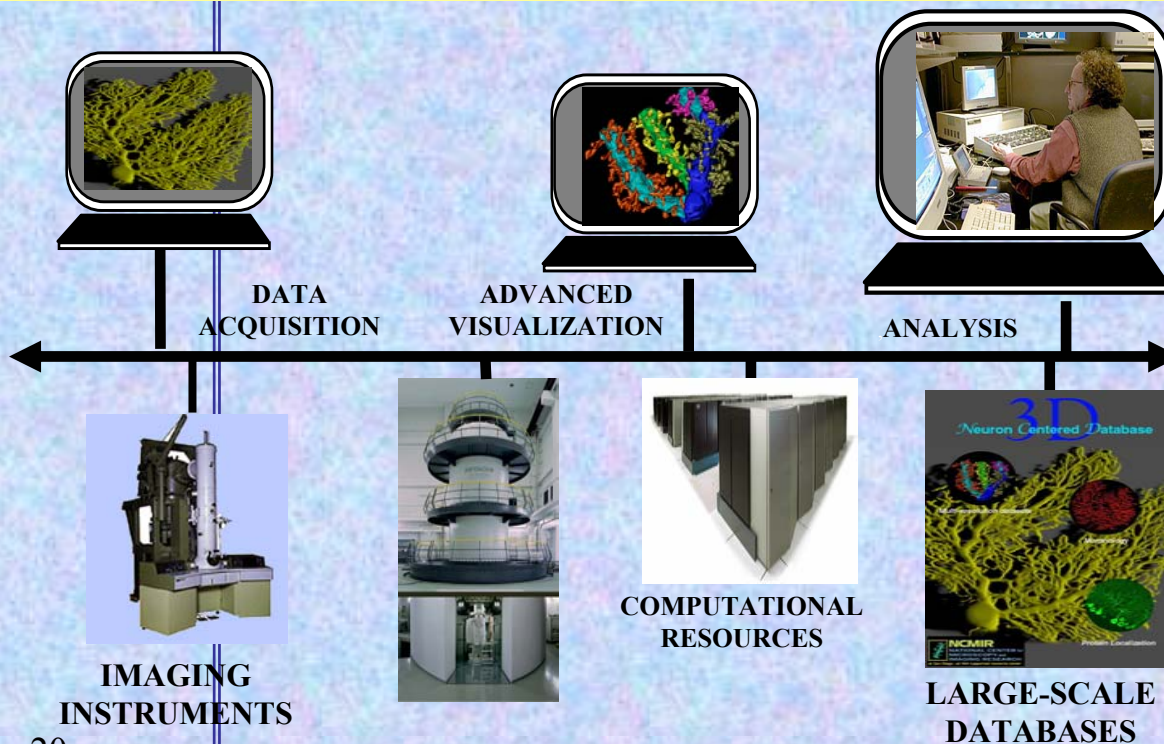
archival storage



DOE X-ray source grand challenge: ANL, USC/ISI, U.Chicago

# Telescience – Collaborative Engineering

- Links computation and data management to unique, expensive instrumentation
- Requires advanced visualization tools for segmentation and analysis of the data
- Provides critical database of biological structure info for neuroscientists



**3D Model of the Node of Ranvier**



# Az "álmom" alkalmazás

Távoli víz. Los Angelesben



Távoli víz. és megfigyelés Berlinből



Távoli megfigyelés és monitorozás Ferihegyről



Előző szimuláció eredményeinek víz.-ja egy bécsi kávéházból



Ekvi-felületek

HDF5

http

T3E: Jülich (Németo.)



Origin: NCSA (USA)



Globus



1. Cactus web portálról indító szimuláció (Róma)

Grid kiterjesztésű Cactus futása elosztott szuperszámítógépeken



# High-Throughput Computing

The screenshot shows the Nimrod-G interface with several key elements highlighted by red circles and labels:

- Deadline:** A red circle highlights the 'Deadline' field, which is set to 'Dec 1998 07 15 : 25 : 00'. A red label 'Deadline' points to this field.
- Cost:** A red circle highlights the 'Current expenditure' and 'Budget' fields, both showing '\$0.00'. A red label 'Cost' points to these fields.
- Available Machines:** A red circle highlights the 'Available Machines' section at the bottom, which lists four resources with their respective task counts: 'mcs.anl.gov - fork' (17/116), 'va.mcs.anl.gov - fork' (2/31), 'p16.mcs.anl.gov - easymcs' (6/44), and 'enali.mcs.anl.gov - fork' (36/65). A red label 'Available Machines' points to this section.

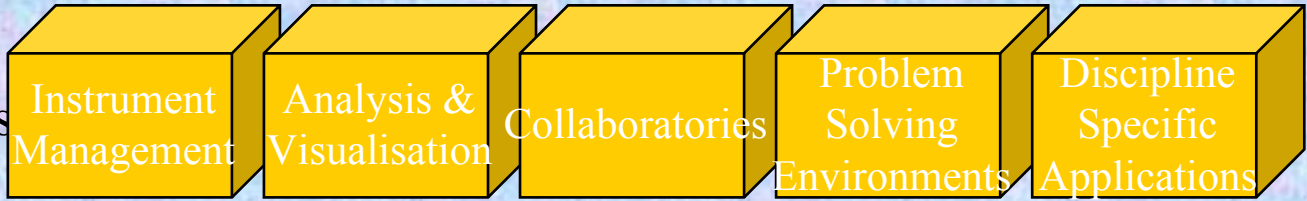
The interface also displays a grid of task status (Waiting, Completed, Failed) and a table of completion statistics.

- **Schedule many independent tasks**
  - Parameter studies
  - Data analysis
- **Issues:**
  - Resource discovery
  - Data Access
  - Scheduling
  - Reservation
  - Security
  - Accounting
  - Code management

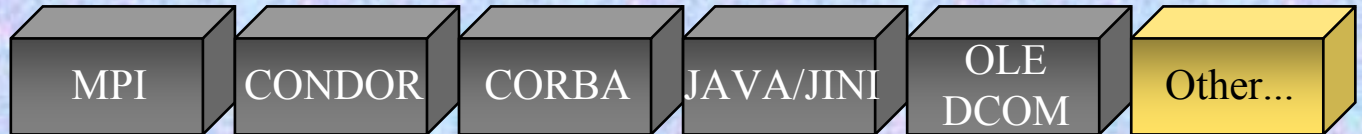


# Generic Grid Architecture

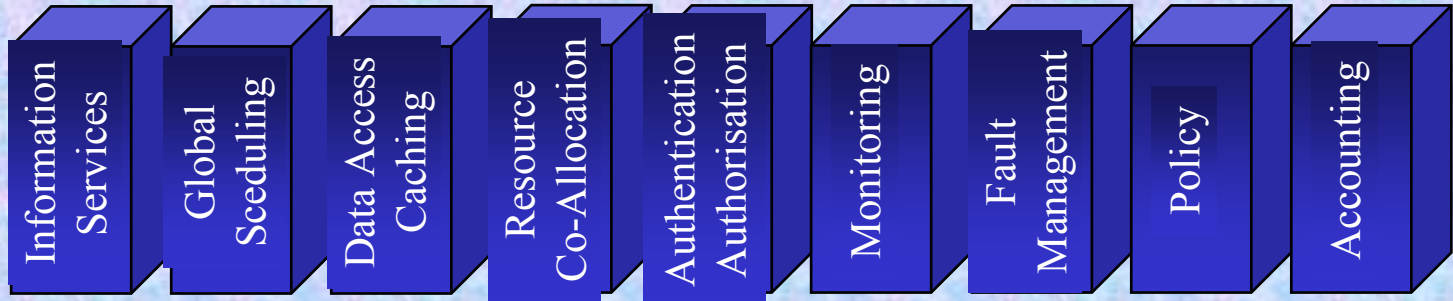
Application Environments



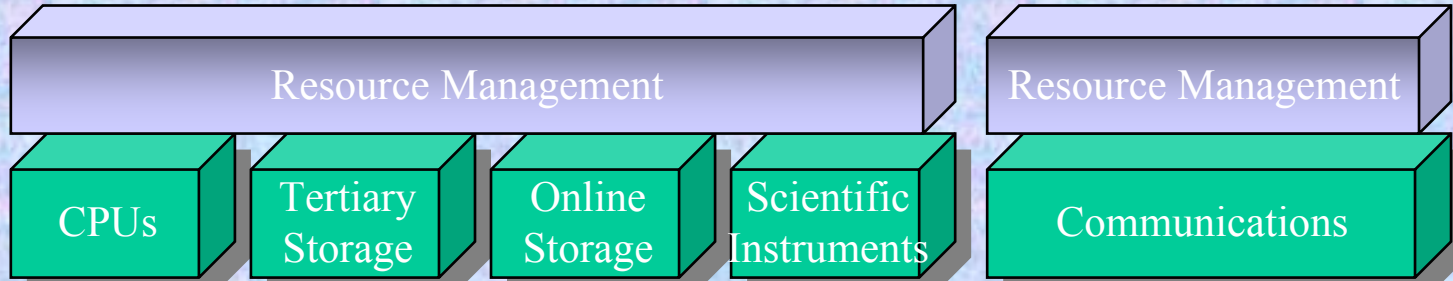
Application Support



Grid Common Services



Grid Fabric - local resources





# Problem Solving Environments

## Examples:

- Problem solving env. for computational chemistry
- Application web portals

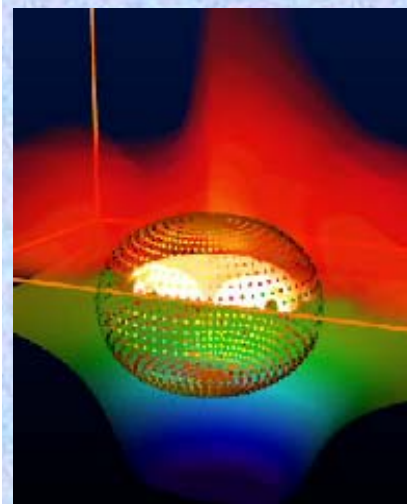
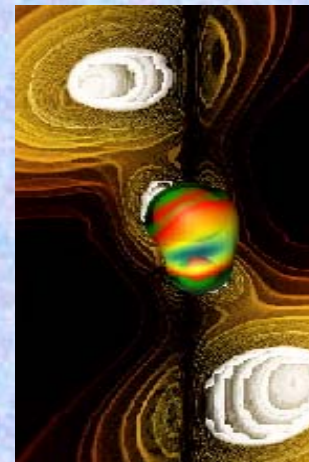
## Issues:

- Remote job submission, monitoring, and control
- Resource discovery
- Distributed data archive
- Security
- Accounting

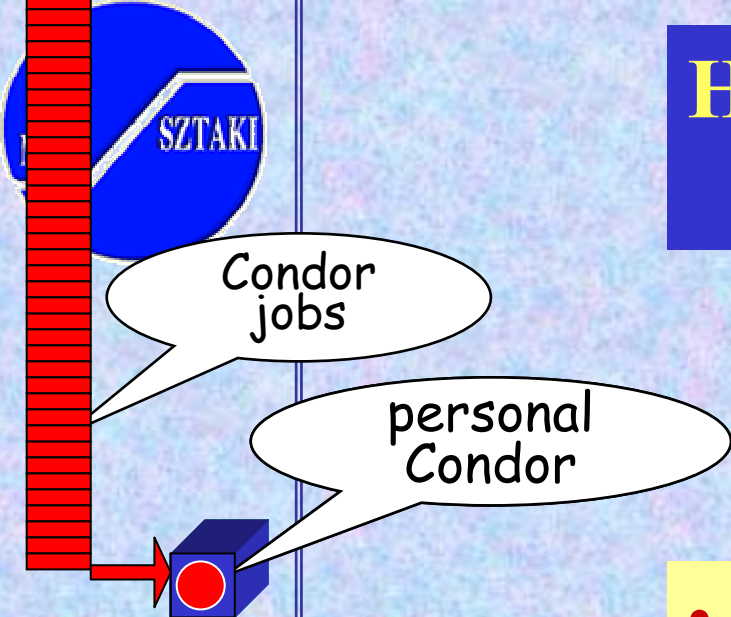


# *Cactus*

- ❏ **Black Holes (prime source for GW) – from Misner BH collisions to grazing BH collisions with initial spin and momentum (Brandt-Bruegmann initial data)**
- ❏ **Gravitational Waves (Evolution of Brill Waves, collapse of pure GW, investigation of critical amplitude, i.e. when do black holes form)**
- ❏ **Neutron Stars (NASA Neutron Star Grand Challenge, GR hydrodynamics, neutron stars colliding to black holes,...)**



# High throughput Computing: Condor



- **Cél:** A gridben lévő számítógépek szabad ciklusainak kihasználása
- **Megvalósítási lépések (1):** A személyes PC v. munkaállomás átalakítása személyes Condor géppé

# High throughput Computing: Condor



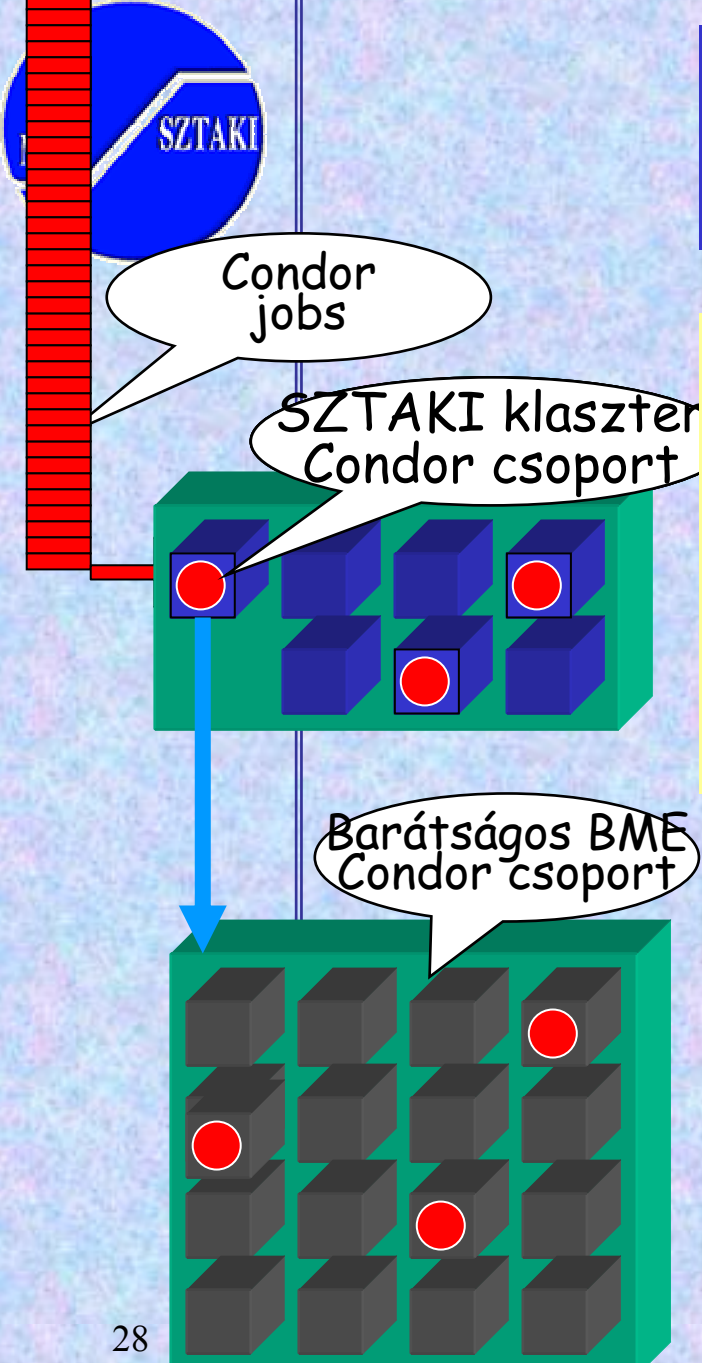
Condor jobs

SZTAKI klaszter  
Condor csoport

- **Megvalósítási lépések (2):**  
Intézeti Condor csoport létrehozása

# High throughput Computing: Condor

- **Megvalósítási lépések (3):**  
Intézeti Condor csoport összekapcsolása más “barátságos” Condor csoportokkal.



- **Megvalósítási lépések (4): Grid erőforrások ideiglenes kihasználása**

Magyaro.-i Grid

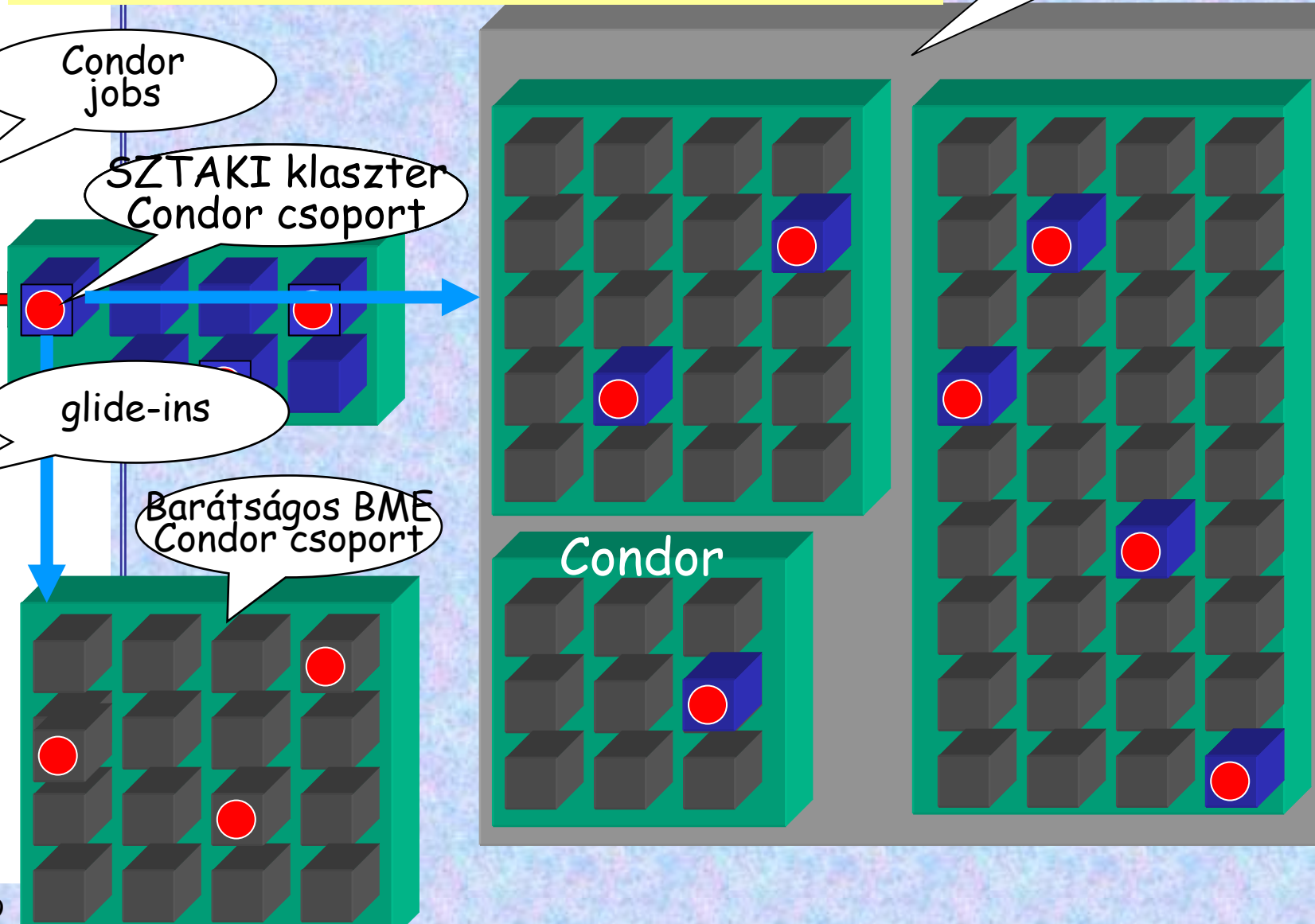
Condor jobs

SZTAKI klaszter  
Condor csoport

glide-ins

Barátságos BME  
Condor csoport

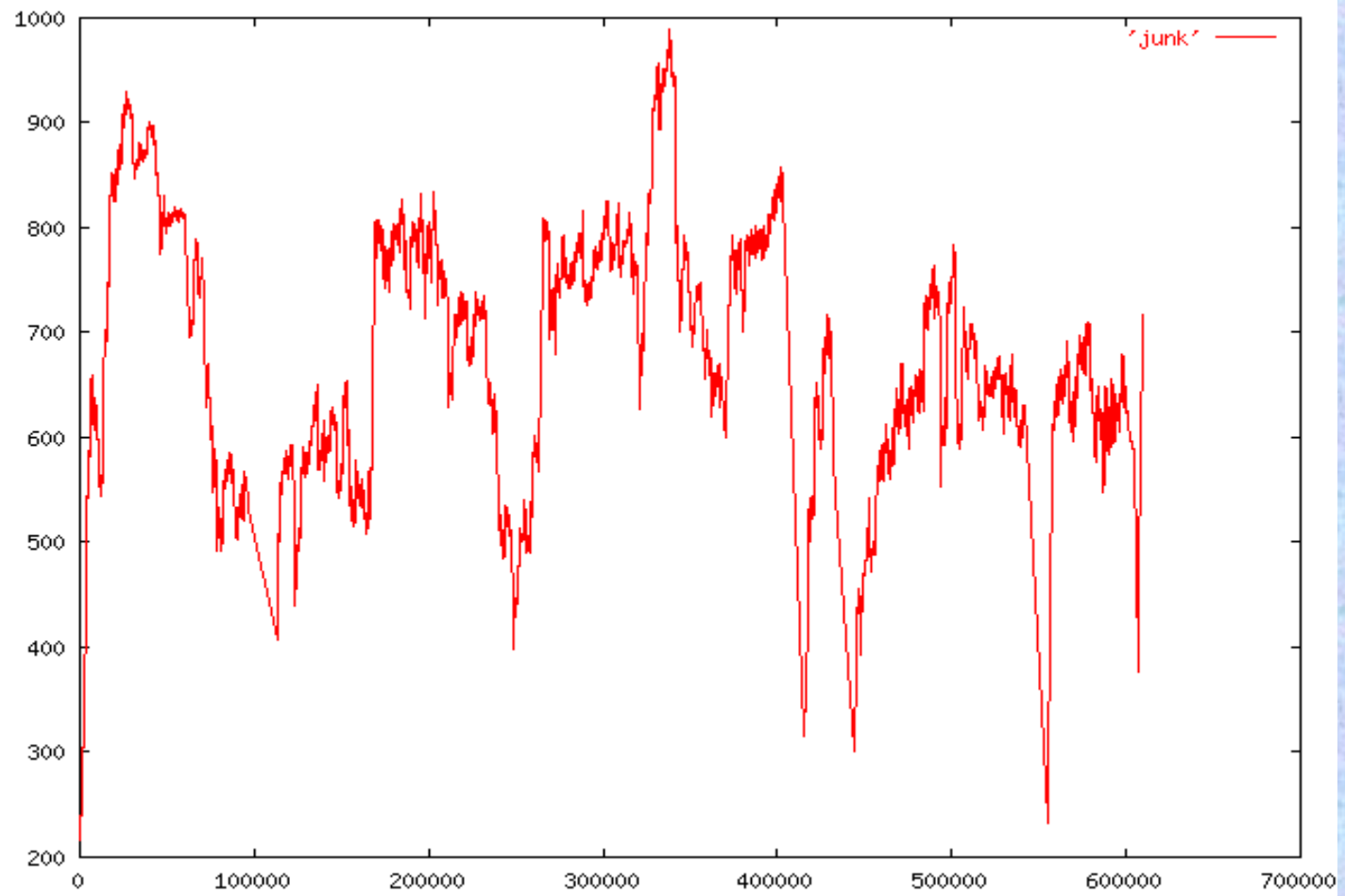
Condor



# *NUG30 kvadratikus allokálási probléma*

- Megoldva 7 nap alatt 10.9 év helyett
- Az első 600K másodperc ...

Processzorok  
száma





# *A NUG30 kísérletben alkalmazott számítógépek*

## **Barátságos Condor csoportok:**

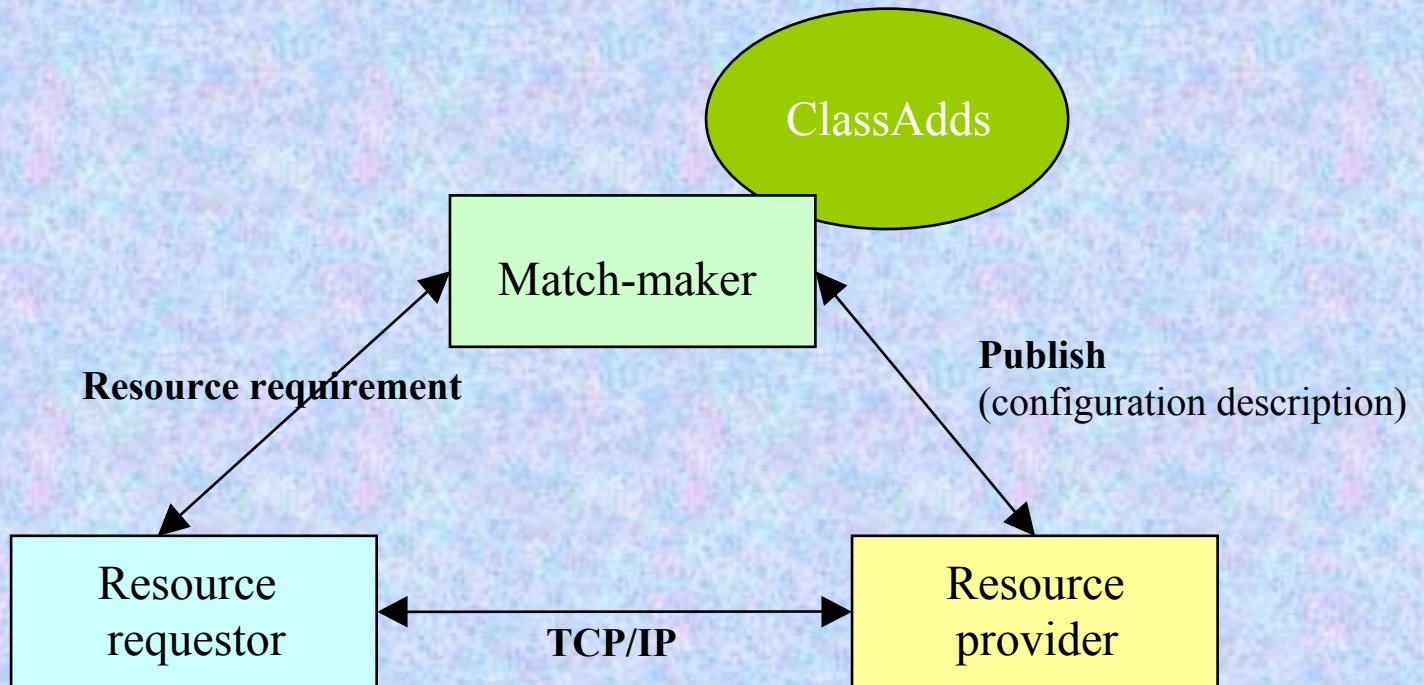
- the main Condor pool at U. of Wisconsin (600 processors)
- the Condor pool at Georgia Tech. U. (190 Linux boxes)
- the Condor pool at UNM (40 processors)
- the Condor pool at U. of Columbia (16 processors)
- the Condor pool at Northwestern U. (12 processors)
- the Condor pool at NCSA (65 processors)
- the Condor pool at INFN (Olaszo.) (200 processors)

## **Grid erőforrások:**

- Origin 2000 at NCSA
- Origin 2000 at Argonne National Lab



# The Condor model



**Your program moves to resource(s)**



**Security is a serious problem!**

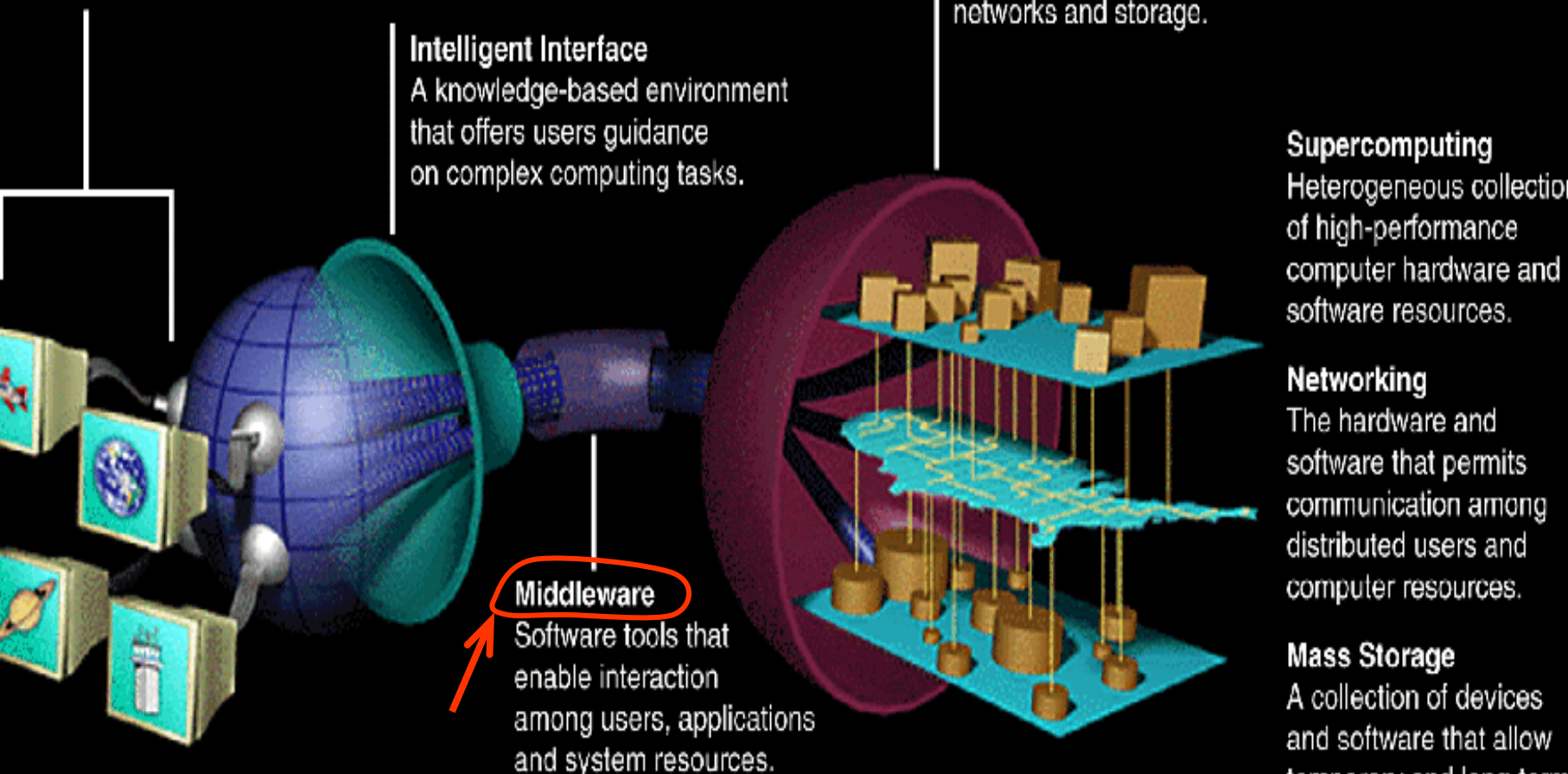


# *GRID/metaszámítási rendszerek létrehozásának komponensei*

- Programozási modell
- Architektúrális komponensek
- Middleware
- Programozási környezet
- Ütemezés és erőforrás kezelés
- Kommunikációs rendszerek és protokollok
- Rendszer problémák megoldása (pl. biztonsági kérdések)

## System Users

Scientists and engineers  
using computation to  
accomplish Lab missions.



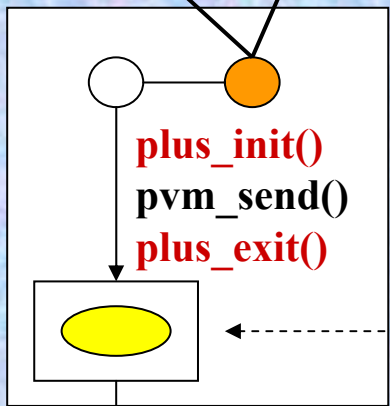
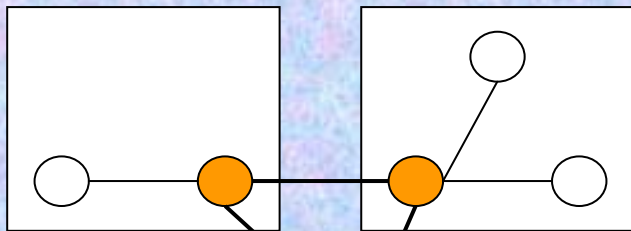


# *Programozási modellek*

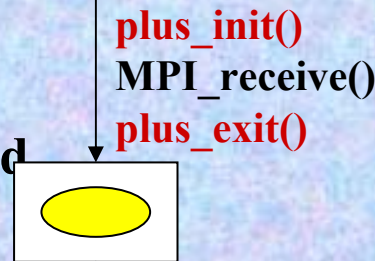
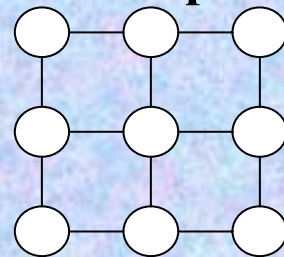
- **Üzenet alapú**
  - PVM, MPI, PVMPI, Globus, PLUS
- **Makro adatfolyam (data-flow)**
  - Webflow, Legion
- **Objektum-orientált modellek**
  - Java, CORBA, DCOM, Legion, Globe
- **Intelligens mobil agensek**
  - applet, servlet

# Processz kommunikáció PVM és MPI között PLUS segítségével

## C+PVM klaszteren



## C+MPI szuperszámítógépen



Frontend

hálózat



PVM démon



PLUS démon



## *Middleware koncepciók*

- Hol helyezkedik el a middleware?
- A middleware célja:
  - Egy radikálisan **heterogén** környezetet virtuálisan **homogén** rendszerré alakítani
- Három fő koncepció:
  - Toolkit (mix-and-match) modell
    - Globus
  - **Objektum-orientált modell**
    - Legion, Globe
  - **Internet-www modell**
    - Webflow



# *Globus Layered Architecture*

## Applications

**Application Toolkits**

Global components: GlobusView, Testbed Status

Toolkits: DUROC, MPI, Condor-G, HPC++, Nimrod/G, globusrun

**Grid Services**

Global components: Nexus, GRAM

Services: I/O, MDS-2, GSI, GSI-FTP, HBM, GASS

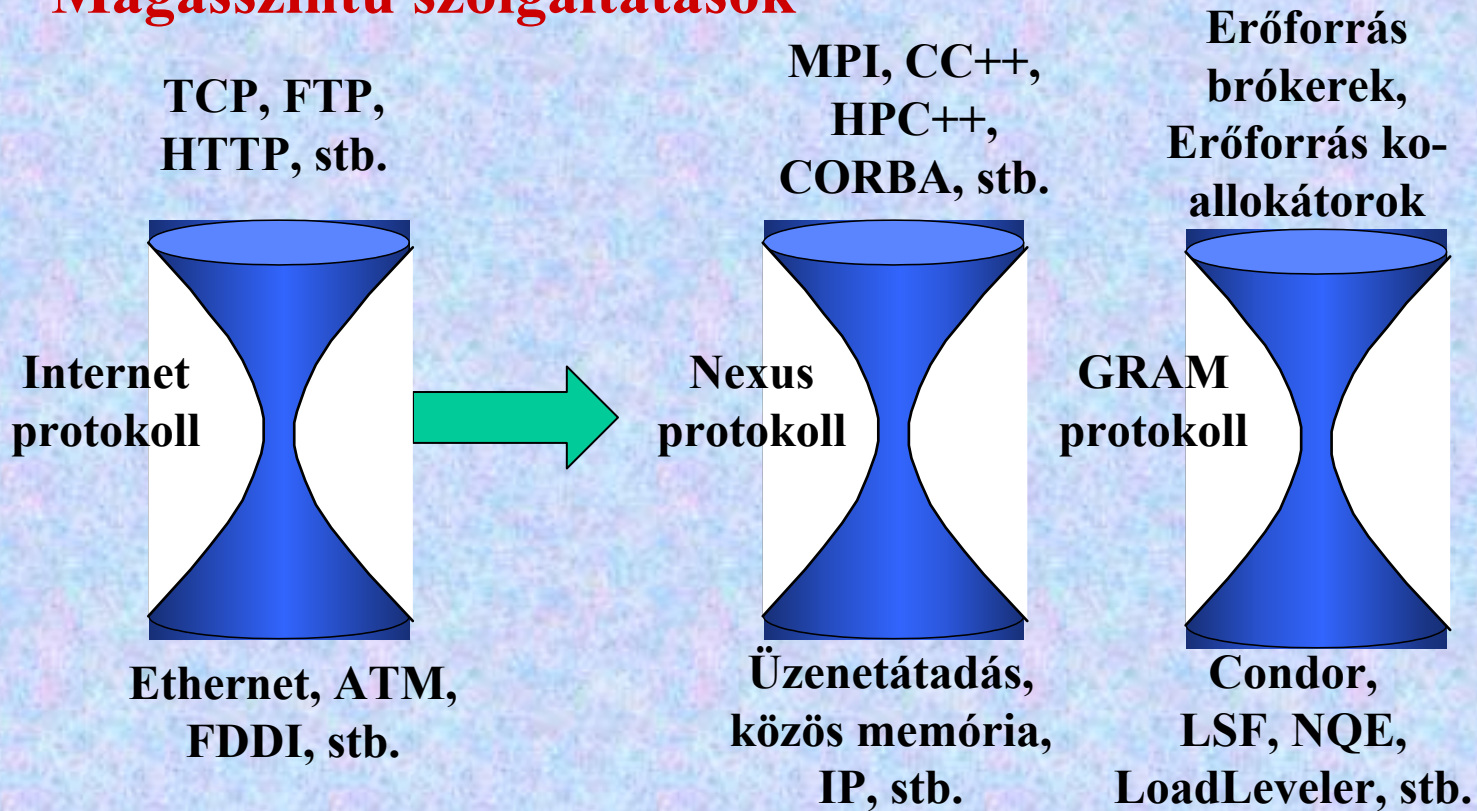
**Grid Fabric**

Global components: Condor, MPI, TCP, UDP

OS/Software: LSF, PBS, NQE, Linux, NT, Solaris, DiffServ

# A Globus homokóra koncepciója

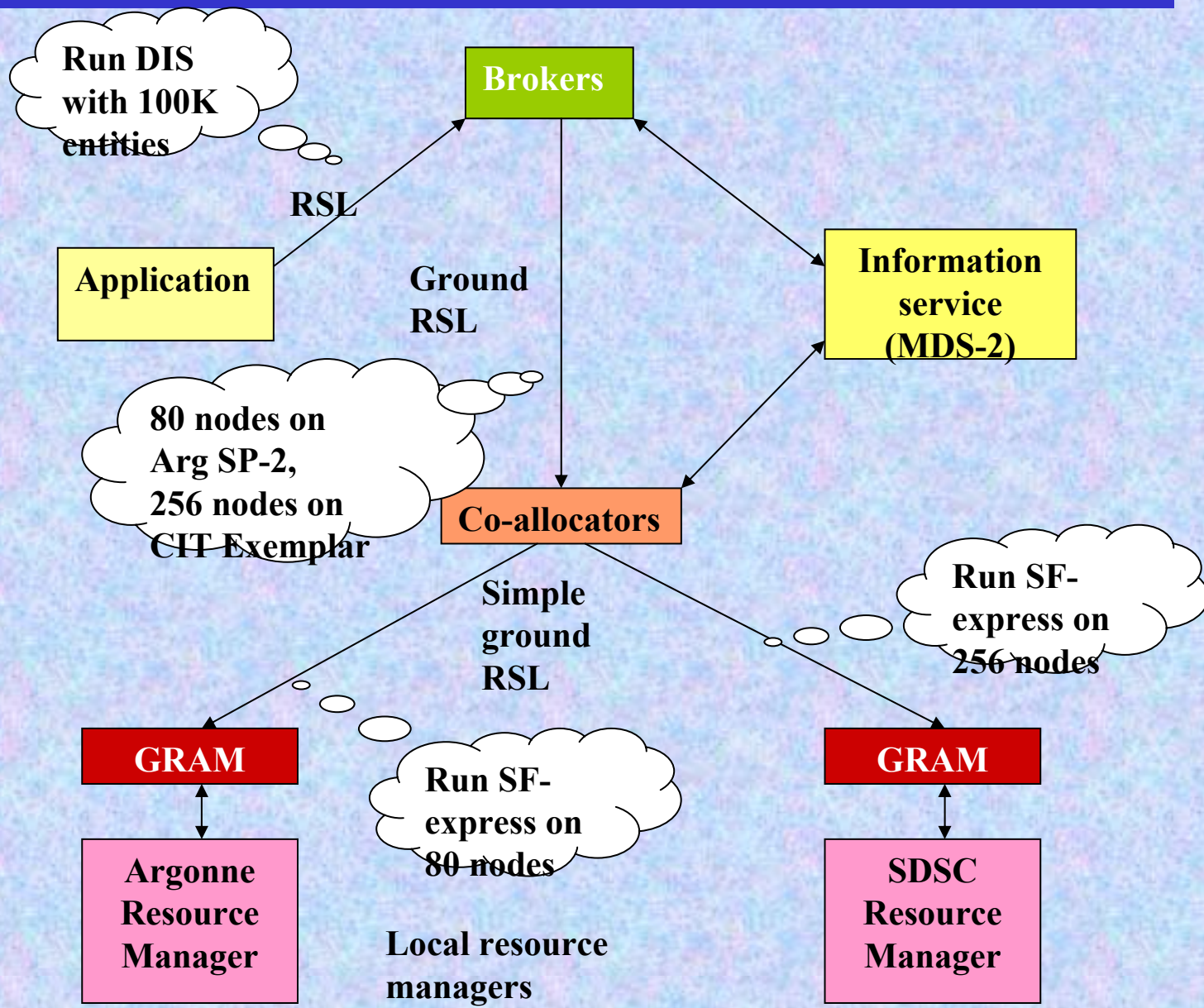
## Magasszintű szolgáltatások



## Alacsonyszintű eszközök

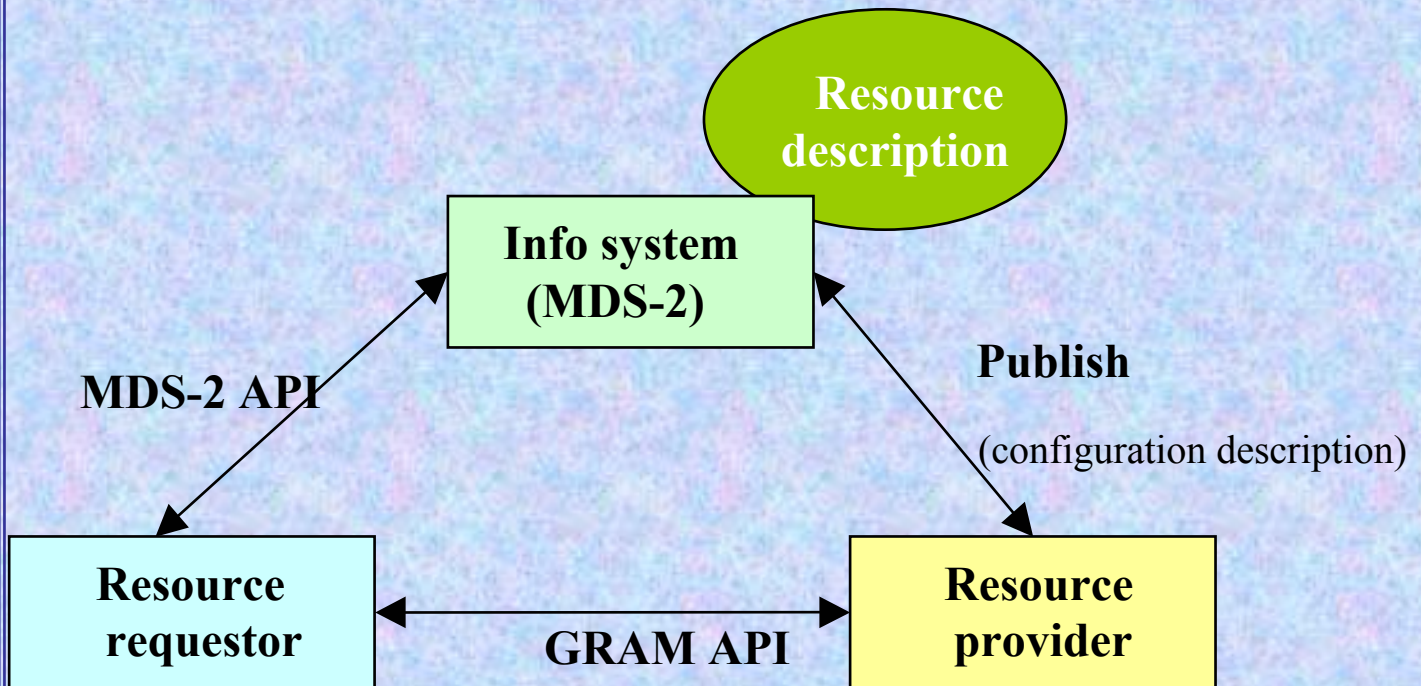


# Globus hierarchical resource management architecture





# The Globus Model



**Your program moves to resource(s)**

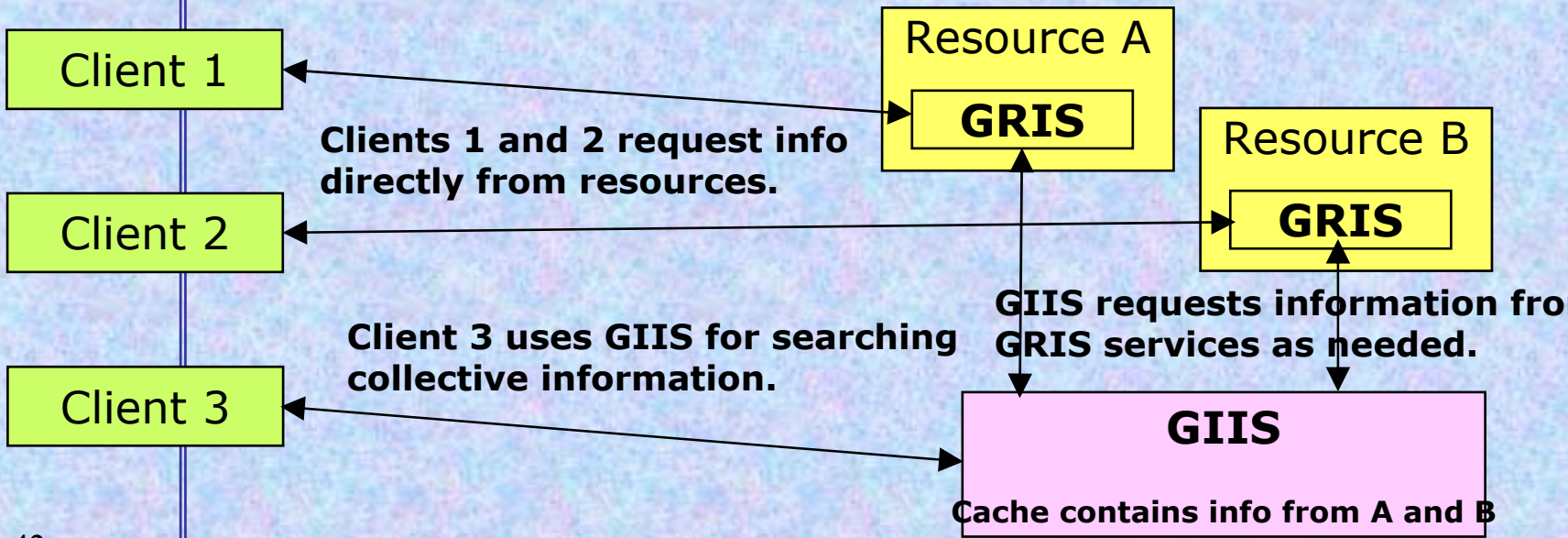


**Security is a serious problem!**



# “Standard” MDS Architecture (MDS-2)

- Resources run a standard information service (**GRIS**) which speaks LDAP and provides information about the resource (no searching).
- GIIS** provides a “caching” service much like a **web search engine**. Resources register with GIIS and GIIS pulls information from them when requested by a client and the cache is expired.
- GIIS** provides the collective-level indexing/searching function.





# *Grid Security Infrastructure (GSI)*

PKI for  
credentials

Proxies and delegation (GSI  
Extensions) for secure **single  
Sign-on**

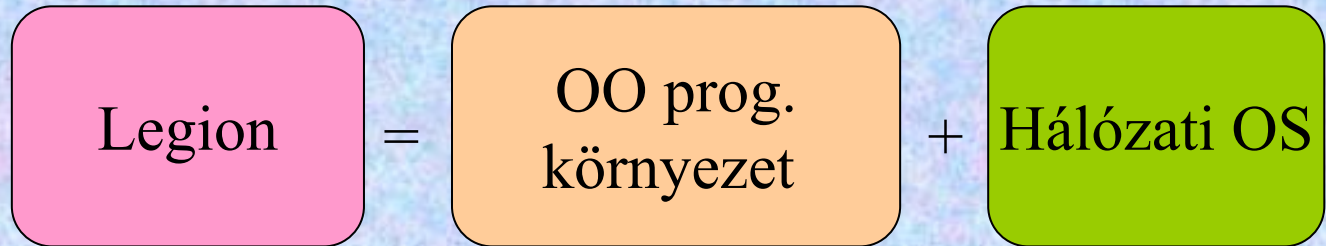
Proxies and Delegation

PKI  
(CAs and  
Certificates)

SSL

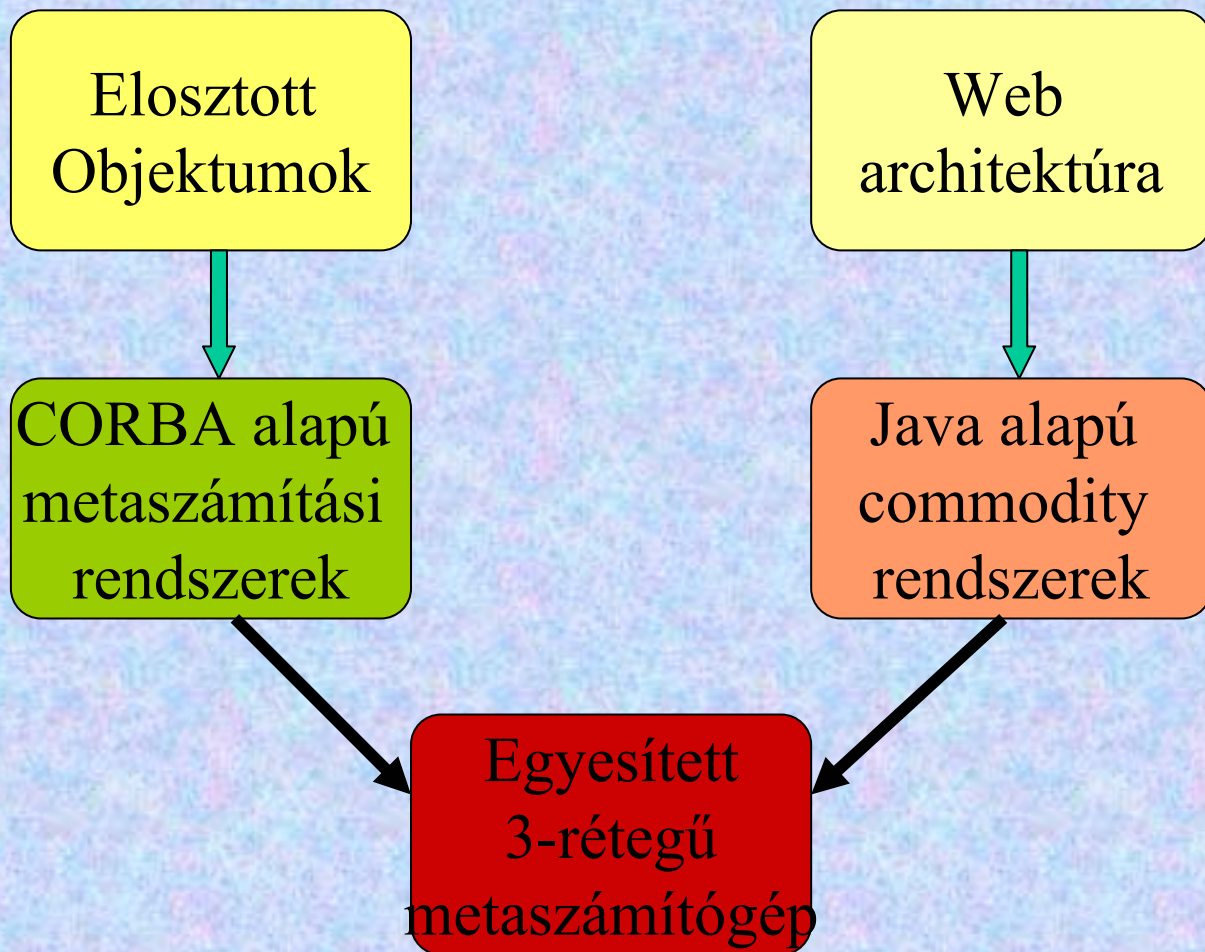
SSL (Secure  
Socket Layer)  
for  
Authentication  
and message  
protection

# Legion

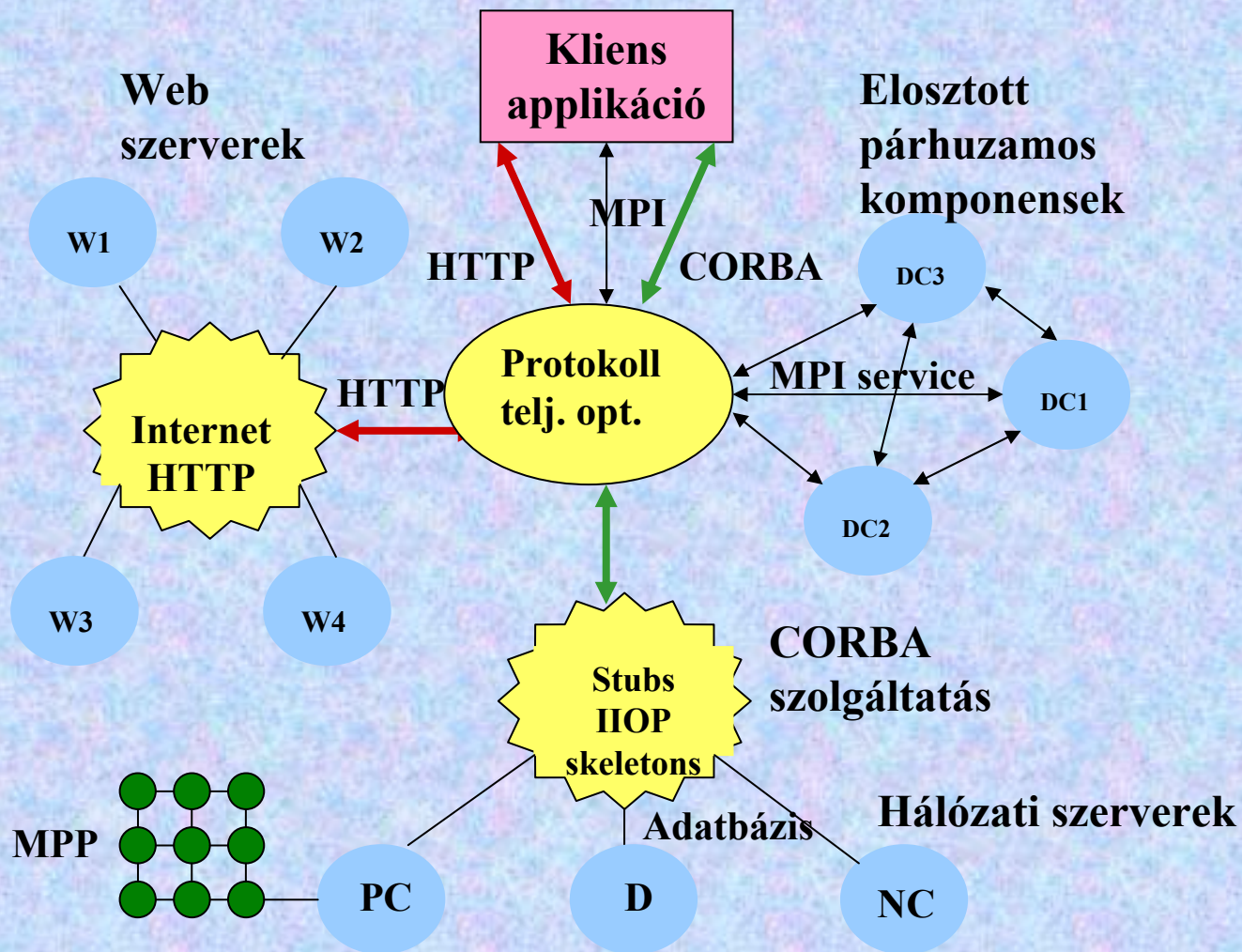


- Magasszintű funkciókat épít egy egyszerű **unifikált objektum modellre** alapozva.
  - Az objektumok processzek, amelyek a munkaelosztás, ütemezés és erőforráskezelés alanyai
- A software IC koncepciót valósítja meg a **makro-dataflow** technikára alapozva

# *Commodity 3-rétegű architektúrák*



# Egyesített 3-rétegű metaszámítógép



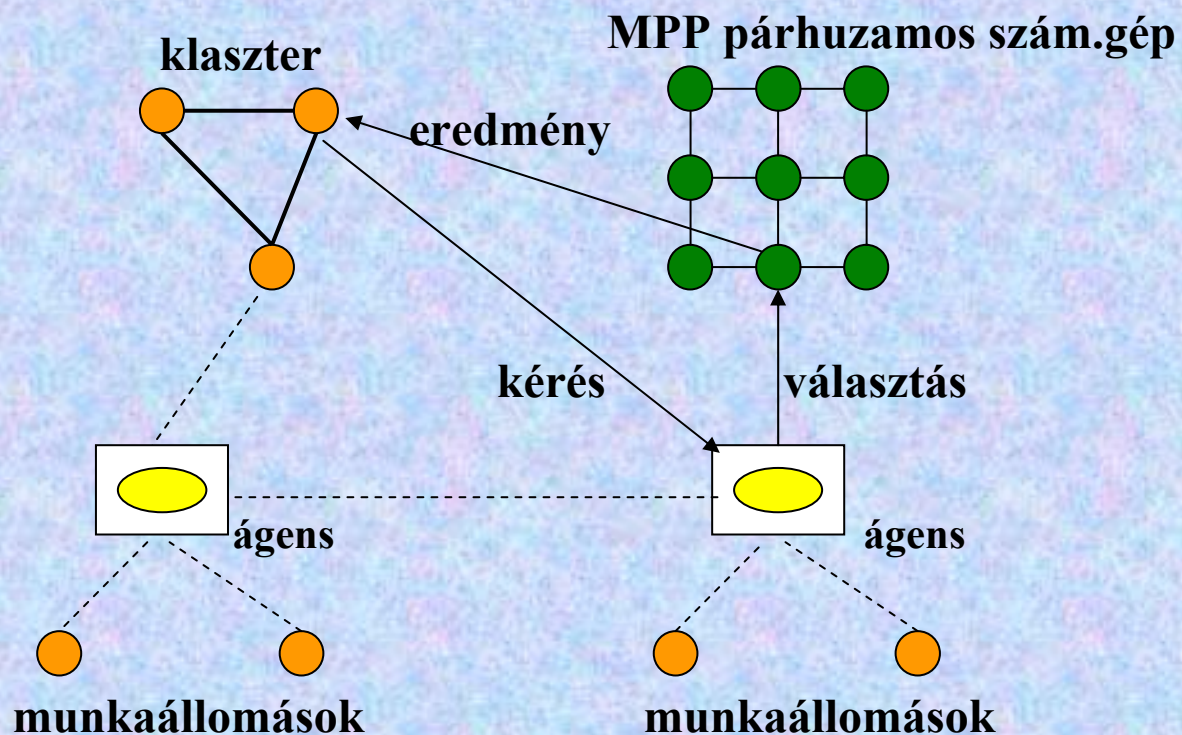


# *Programozási környezetek*

- **Toolkit alapú**
  - Cactus
  - Paderborn toolkit:
    - PLUS kommunikációs környezet
    - Resource and Service Description (RSD)
- **Integrált környezetek**
  - P-GRADE GRID verziója
- **Applikáció specifikus környezetek**
  - NetSolve

A GRID hálózatban rendelkezésre álló erőforrások kihasználása **numerikus számításokhoz**

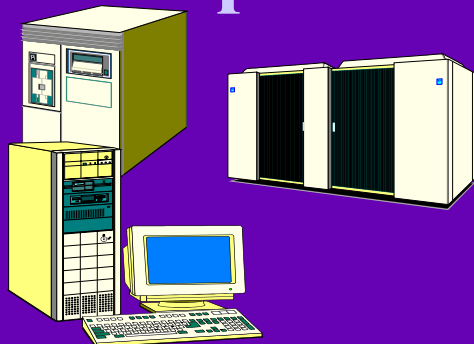
Fő koncepció: kliens - szerver - **ágens**





# NetSolve

## Computational Resources



<b>Hardware:</b>	<b>Software:</b>
Clusters	Routines
MPP	Libraries
Workstations	Applications
Globus, Condor, MPI, PVM	

**Choice**

**Reply**

**Matlab**  
**Mathematica**  
**C, Fortran**  
**Java, Excel**  
**Java GUI**

**Request**

**Agent**

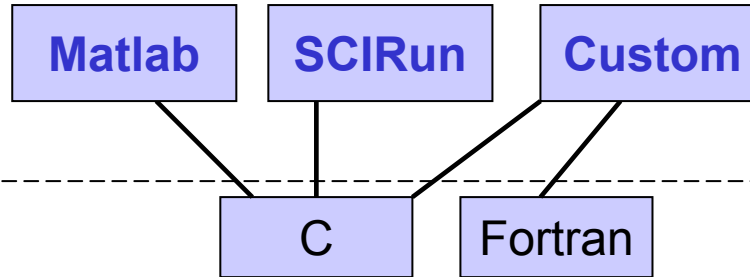
**Scheduler**  
**Database**



**Client - RPC like**

# NetSolve Infrastructure

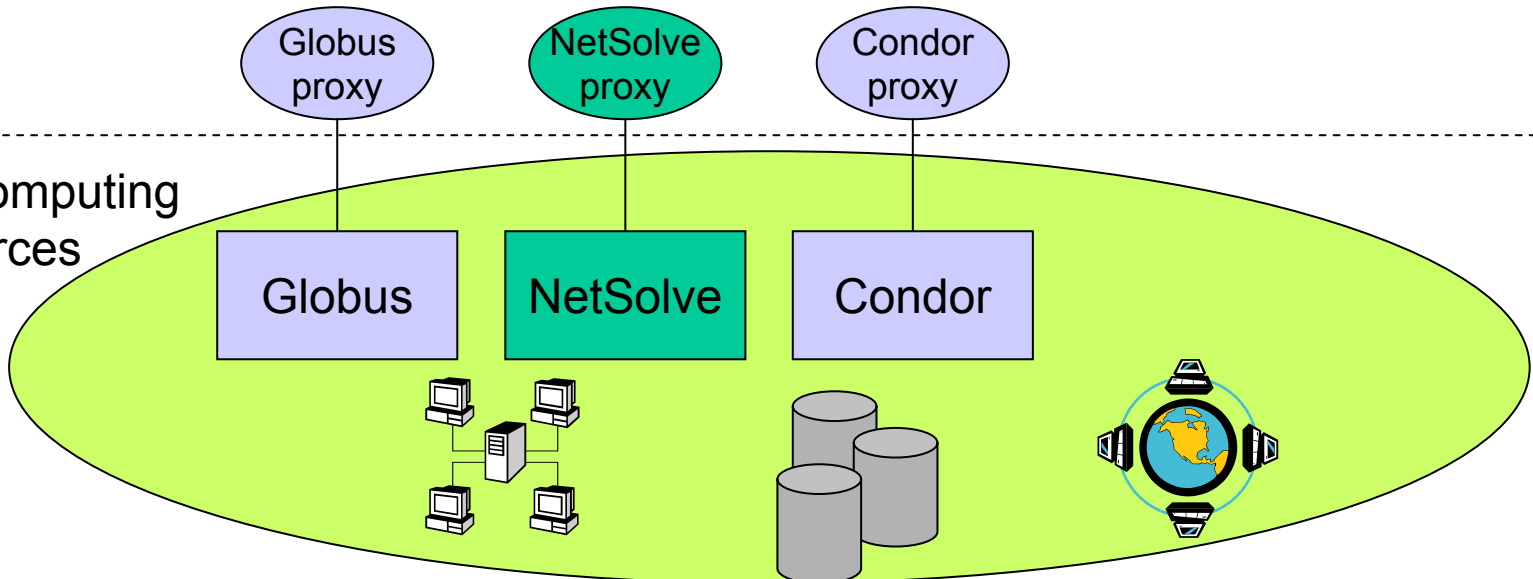
PSEs and  
Applications



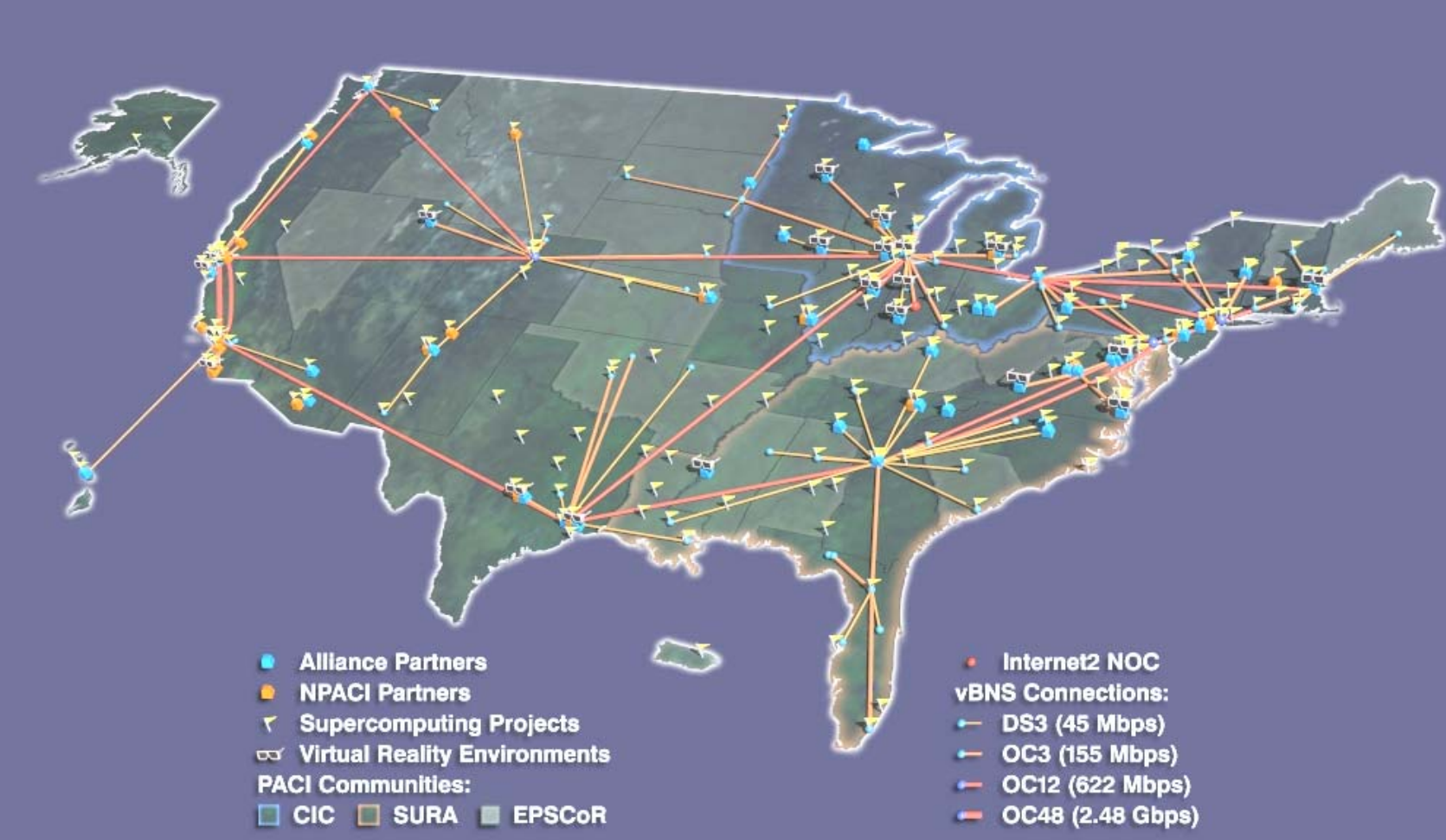
Middleware



Metacomputing  
Resources

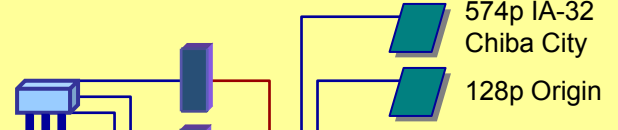
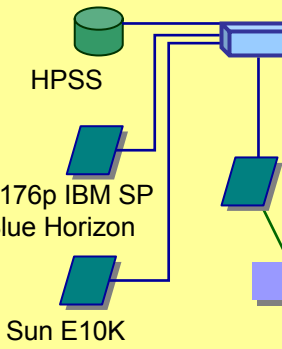
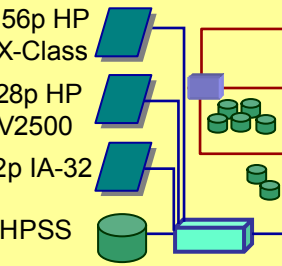


# Hol tartott az USA 2000-ben?

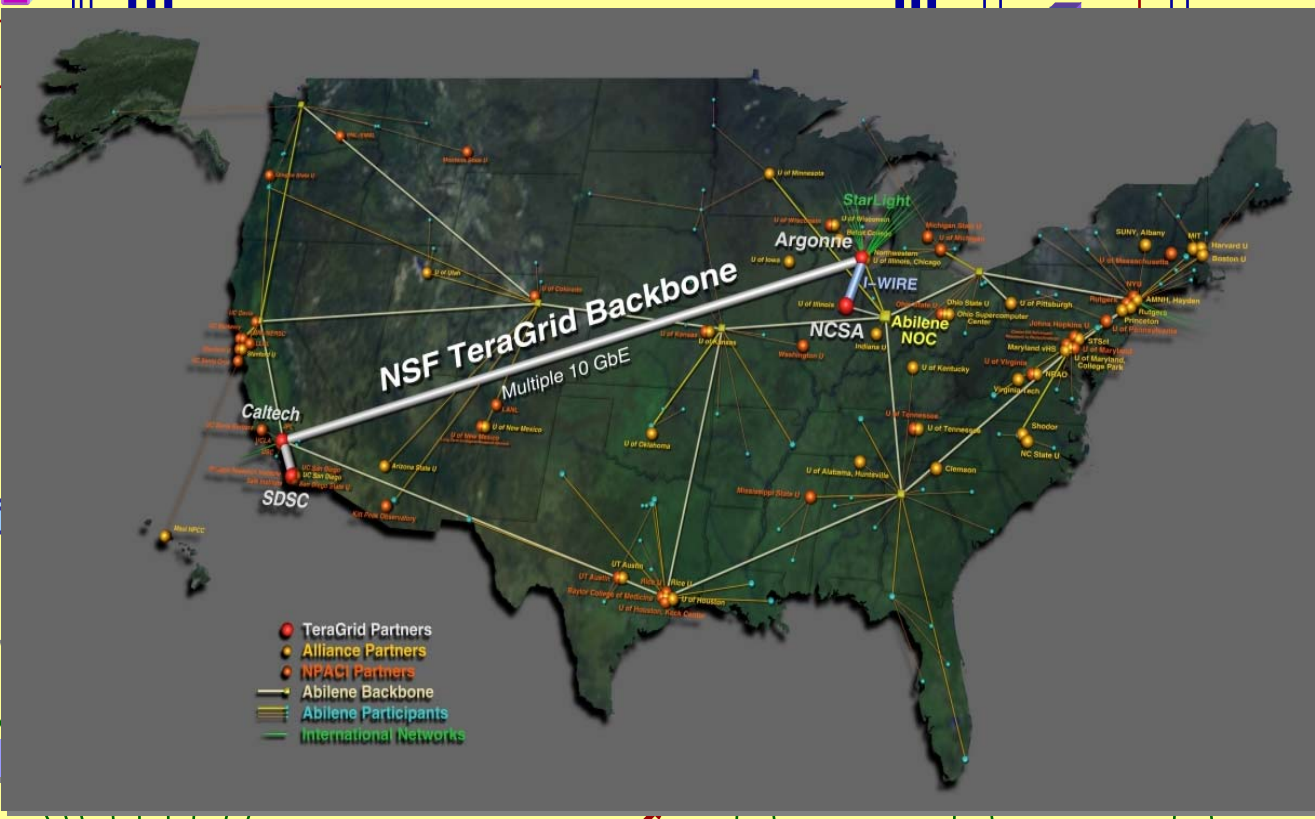


NSF National Technology Grid

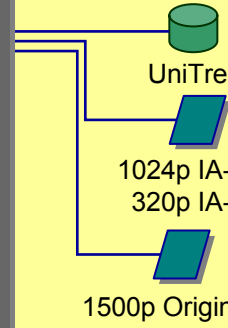
# Hol tart ma az USA? - TeraGrid



HR Display & VR Facilities  
HPSS



ization



**NCSA: Compute-Intensive**



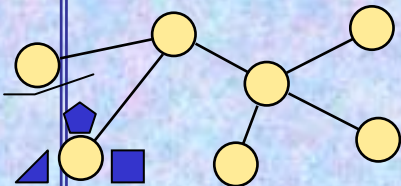
# *Fundamental grid functionalities*

- The essential grid functionalities are:
  - **Resource abstraction**
    - Physical resources can be assigned to virtual resource needs (matched by properties)
    - Grid provides a mapping between virtual resource needs and physical resources
  - **User abstraction**
    - User of the physical machine may be different from the user of the virtual machine
    - Grid provides a temporal mapping between virtual and physical users

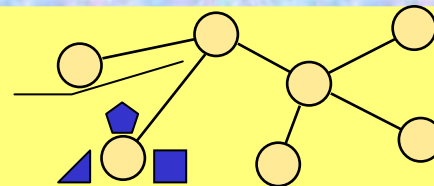


# Conventional distributed environments and grids

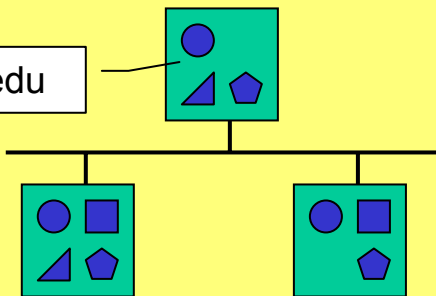
Smith  
4 nodes



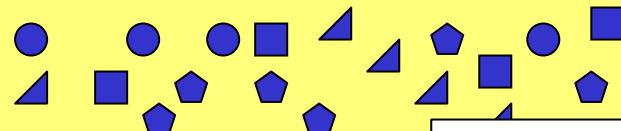
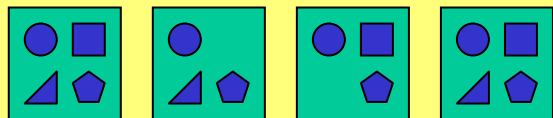
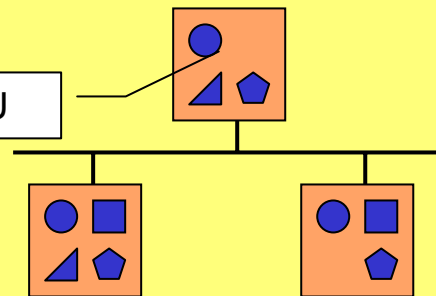
Smith  
4 CPU,  
memory,  
storage



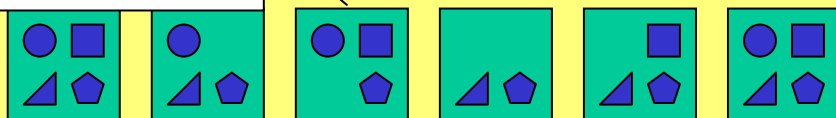
smith@n1.edu



Smith 1 CPU



smith@n1.edu



griduser@n1.edu

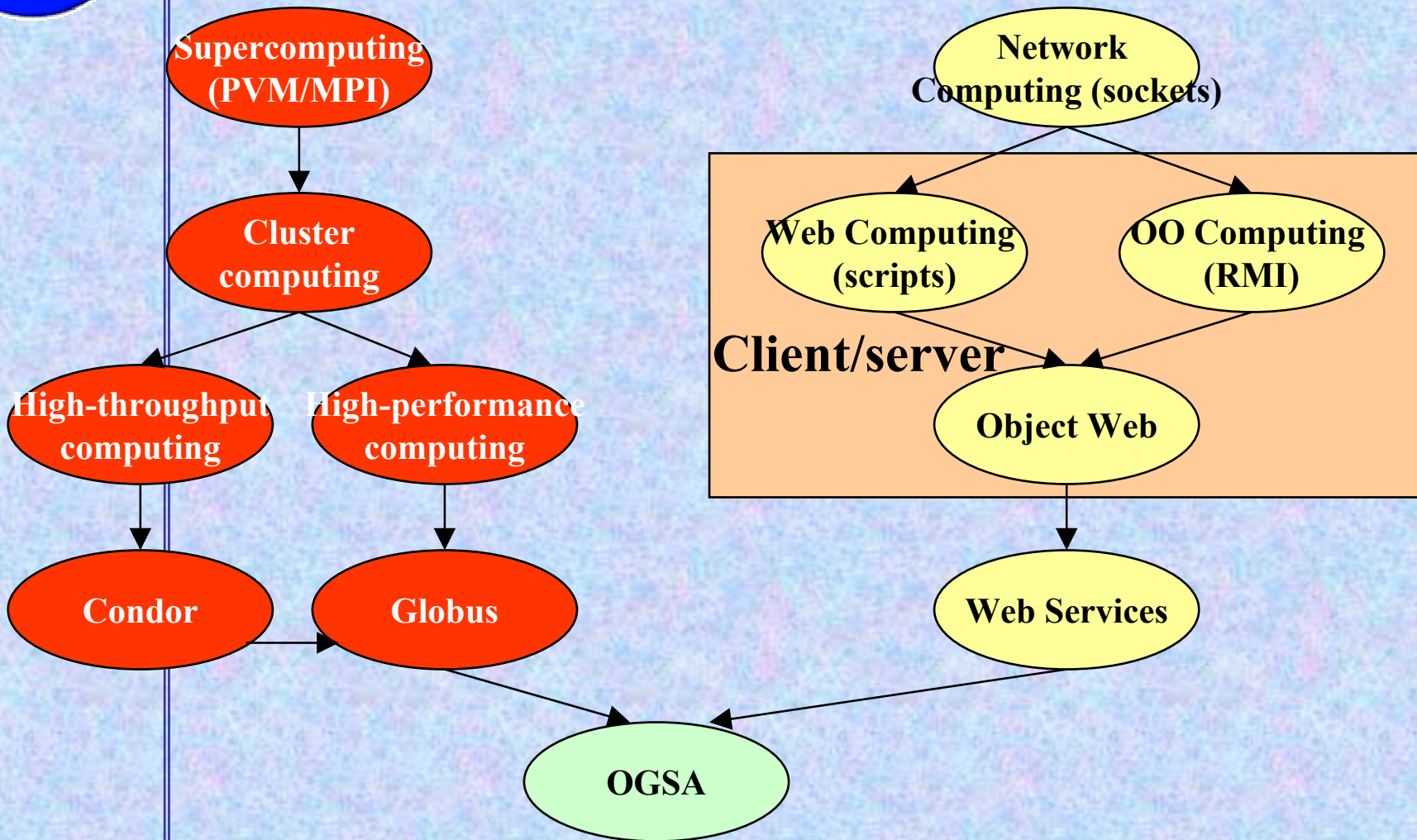


smith@n2.edu

p12@n2.edu

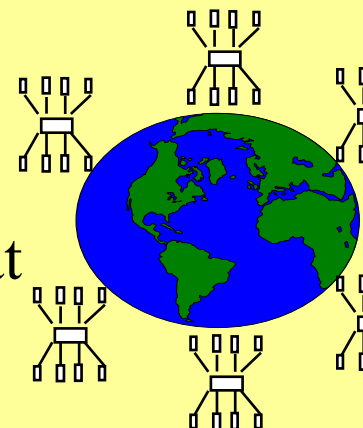


# Grid rendszerek fejlődése



# Összefoglalás

- A GRID/metaszámítógép egy új technológia, amely **integrálja**:
  - a szuperszámítógép technológiát
  - a távolsági hálózatok technológiáját
  - a WWW technológiát



- **A Grid egy új, az elektromos áramot elosztó hálózathoz hasonló infrastruktúra kialakulásához fog vezetni**
- **Ez az új infrastruktúra a www technológiához hasonlóan óriási hatással lesz az informatikai társadalom továbbfejlődésére**



# *Szuperszámítógépek, klaszterek és metaszámítógépek összehasonlítása I.*

	Supercomputer	Cluster	NOW	Metacomputing system
Processing units (nodes)	Microprocessors	PCs, workstations	PCs, workstations	Supercomputers, clusters, PCs, workstations
Number of nodes	100 - 1000	10 - 100	10 - 100	100 - 10000
Communication network	Buses, switches	LAN	LAN	Internet
Node OS	Homogeneous	Typically homogeneous	Typically heterogeneous	Heterogeneous
Inter-node security	Nonexistent	Rarely required	Necessary	Necessary



# Szuperszámítógépek, klastterek és metaszámítógépek összehasonlítása II.

	Supercomputer	Cluster	Metacomputing system
Programming models	<ul style="list-style-type: none"> <li>• Shared memory</li> <li>• Message passing</li> </ul>	<ul style="list-style-type: none"> <li>• Shared memory</li> <li>• Message passing</li> <li>• Peer-to-peer</li> <li>• Client-server</li> </ul>	<ul style="list-style-type: none"> <li>• Message passing</li> <li>• Client-server</li> <li>• Code shipping</li> <li>• Proxy computing</li> <li>• Intelligent mobile agents</li> </ul>
Programming language	<ul style="list-style-type: none"> <li>• HPF</li> <li>• (C/Fortran)+MPI</li> </ul>	<ul style="list-style-type: none"> <li>• HPF</li> <li>• (C/Fortran)+MPI</li> </ul>	<ul style="list-style-type: none"> <li>• HPF</li> <li>• (C/Fortran)+MPI</li> <li>• Java/CORBA</li> </ul>
Middleware	<ul style="list-style-type: none"> <li>• No</li> </ul>	<ul style="list-style-type: none"> <li>• Limited forms</li> </ul>	<ul style="list-style-type: none"> <li>• Toolkit approach</li> <li>• Three-tier commodity (Java/CORBA)</li> <li>• Object-oriented</li> </ul>
Programming environment	<ul style="list-style-type: none"> <li>• Toolkit approach</li> <li>• Integrated environment</li> </ul>	<ul style="list-style-type: none"> <li>• Toolkit approach</li> <li>• Integrated environment</li> </ul>	<ul style="list-style-type: none"> <li>• Toolkit based</li> <li>• Application specific</li> <li>• Integrated environment</li> </ul>
Resource allocation	<ul style="list-style-type: none"> <li>• Mapping</li> <li>• Load balancing</li> </ul>	<ul style="list-style-type: none"> <li>• Mapping</li> <li>• Load balancing</li> </ul>	<ul style="list-style-type: none"> <li>• Resource manager</li> </ul>
QoS	No	No	Yes
Security	No	No	Yes

*Köszönöm a figyelmüket*



További információ: [www.lpds.sztaki.hu](http://www.lpds.sztaki.hu)