



A GRID-based multilayer architecture for bioinformatics

*Ezio Bartocci, Diletta Cacciagrano, Nicola Cannata,
Flavio Corradini, Emanuela Merelli,*

➤ Mathematic and Computer Science Department, University of Camerino, Italy

Luciano Milanesi,

➤ Institute for Biomedical Technologies, CNR, Milan, Italy

Paolo Romano

➤ National Cancer Research Institute, Genova, Italy

A GRID-based multilayer architecture for **bioinformatics**

Which bioinformatics?

Bioinformatics has changed...

```
BLASTN 2.2.4 [Aug-26-2002]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.

Query= gi|8777285|gb|AB044774.1|AB044774 AB044774 Panax ginseng hairy
root Panax ginseng cDNA clone HR36, mRNA sequence
(223 letters)

Database: /local/wwwstud/html/bioinfo3-64/EST-ginseng.fas
21 sequences; 7481 total letters

Searching done

Sequences producing significant alignments:

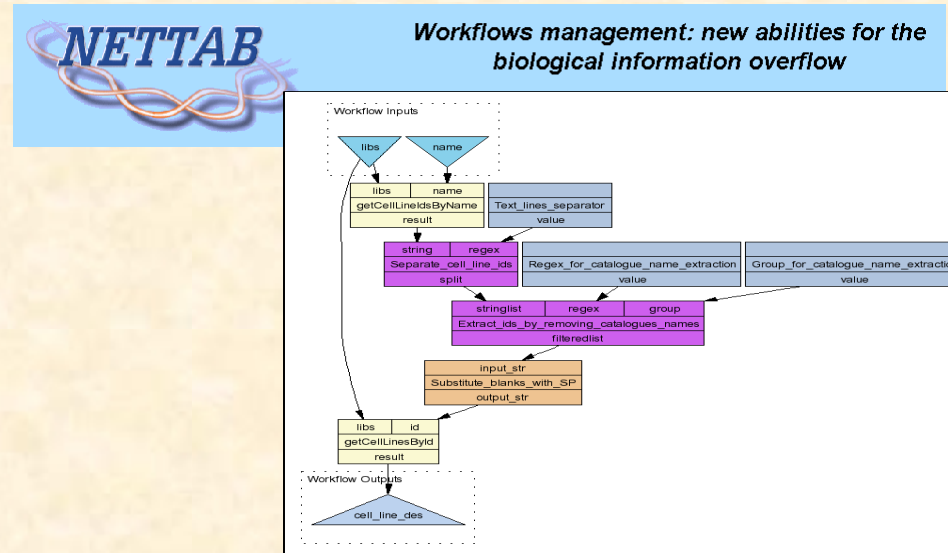
              Score      E
            (bits)  Value
gi|8777285|gb|AB044774.1|AB044774 AB044774 Panax ginseng hairy r... 442  e-128
gi|8777269|gb|AB044758.1|AB044758 AB044758 Panax ginseng hairy r...  24  0.077
gi|8777271|gb|AB044760.1|AB044760 AB044760 Panax ginseng hairy r...  22  0.30

>gi|8777285|gb|AB044774.1|AB044774 AB044774 Panax ginseng hairy root
Panax ginseng cDNA clone HR36, mRNA sequence
Length = 223

Score = 442 bits (223), Expect = e-128
Identities = 223/223 (100%)
Strand = Plus / Plus

Query: 1   ggctcccatgtttgttaaaaaatagtgctccctgctgctgaagattctgaattctatg 60
          |||
Sbjct: 1   ggctcccatgtttgttaaaaaatagtgctccctgctgctgaagattctgaattctatg 60
```

From command line... (till 90s)



NCBI **protein-protein BLAST**

Nucleotide Protein Translations Retrieve results for an RID

[Search](#)

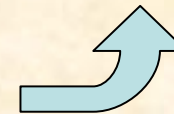
[Set subsequence](#) From: To:

[Choose database](#)

[Do CD-Search](#) ☒

Now: **BLAST!** or [Reset query](#) [Reset all](#)

To web services and workflows (now)



to web interfaces... (with the advent of the WWW)

... very fast ...

TRENDS GUIDE TO BIOINFORMATICS

Biological data, and DNA sequence data in particular, are accumulating at a phenomenal rate. By around 2005, it is likely that the DNA sequence of the complete human genome will have been determined. Although this achievement might seem an end in itself, in reality it is only the beginning. In order to exploit the wealth of DNA sequence and other

biological data, a new science has arisen that fuses biology with mathematics and computer science – 'bioinformatics'.

To find the genes within the genomic sequence is a massive task in itself. Once apparent, otherwise uncharacterized coding regions must be assigned a function. Thereafter, the interactions between genes and gene products must be understood at all levels, not merely in the context of the pathways within and between cells but also in terms of the evolution of gene families within and between species.

These questions can all be addressed using bioinformatics.

Bioinformatics touches all of biology, and straightforward access to data via the Internet means that a wealth of information

over, the because mathematical language these nes the ic disciplines clearly ples of technical in the trievig function plecular Guide to

COPYRIGHT INFORMATION

©1998 Elsevier Science. All rights reserved. This supplement and the individual contributions contained in it are protected under copyright by Elsevier Science. See the box in the accompanying Trends journal for further terms and conditions that apply to the copyright. Except as outlined in the terms and conditions, no part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without prior written permission of the publisher.



trends guide to bioinformatics

Bioinformatics – a new era
Mark Boguski

3

Text-based database searching
Fran Lewitter

7

Fundamentals of database searching
Stephen Altschul

9

Practical database searching
Steven Brenner

12

Computational genefinding
David Haussler

15

Multiple-alignment & -sequence searches
Sean Eddy

18

Protein classification & functional assignment
Kay Hofmann

22

Phylogenetic analysis & comparative genomics
James Lake and Jonathan Moore

24

Databases of biological information
Minoru Kanehisa

27

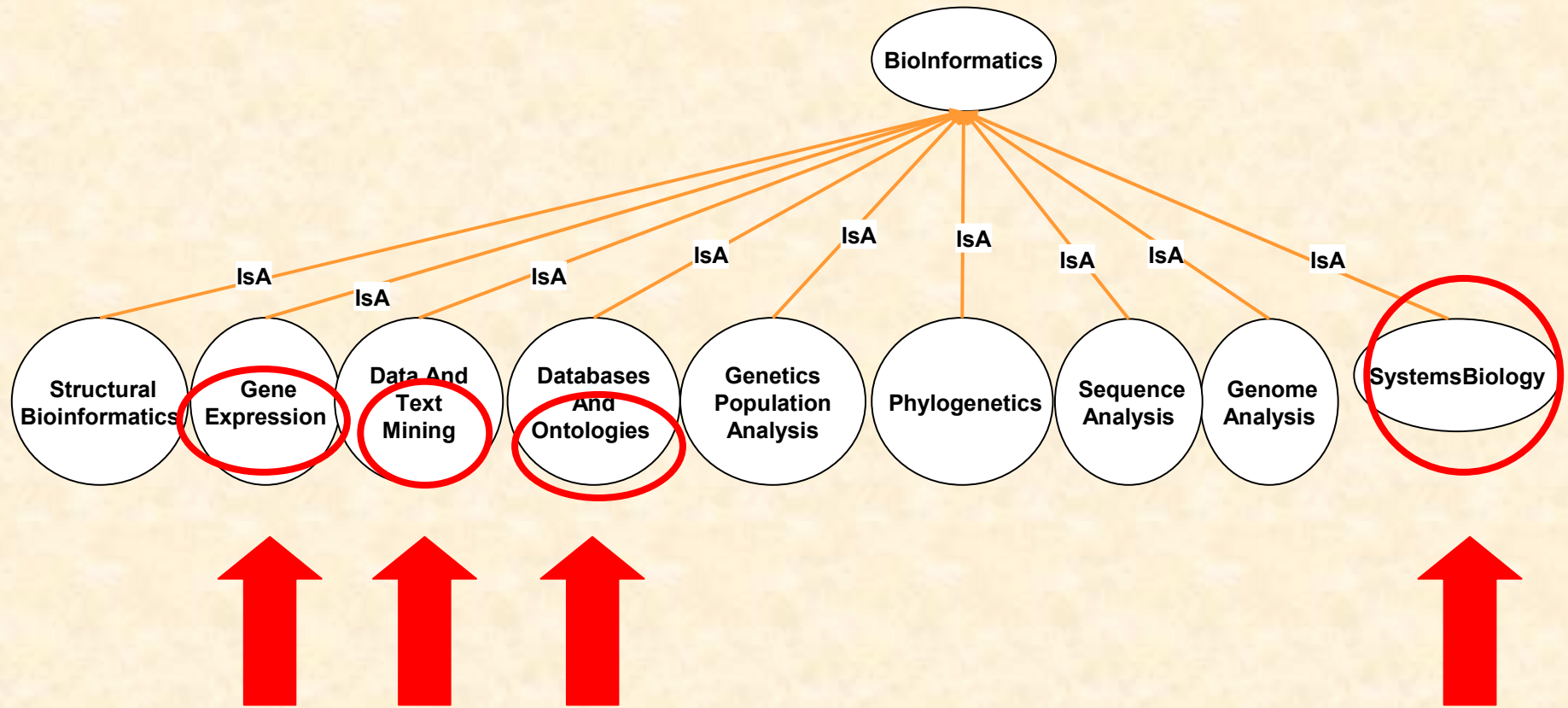
Functional genomics
Michael Brownstein, Jeffrey Trent and Mark Boguski

30

The future of bioinformatics
Janet Thornton

32

Glossary



The classification introduced for articles of “Oxford’s Bioinformatics” in 2005

...and grown

(Or better burst?.....)

PSTT Vol. 2, No. 9 September 1999

| editorial

Information **overflow** from discovery to development



The pharm
tools to d
section of

Ian Shaw
Momentum Heat
Temple Court
Cathedral Road
Cardiff
UK CF1 9HA
tel: +44 1222 78
fax: +44 1222 7
e-mail: shaw@m



Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Journal of Biomedical Informatics xxx (2006) xxx-xxx

Journal of
Biomedical
Informatics

www.elsevier.com/locate/yjbin

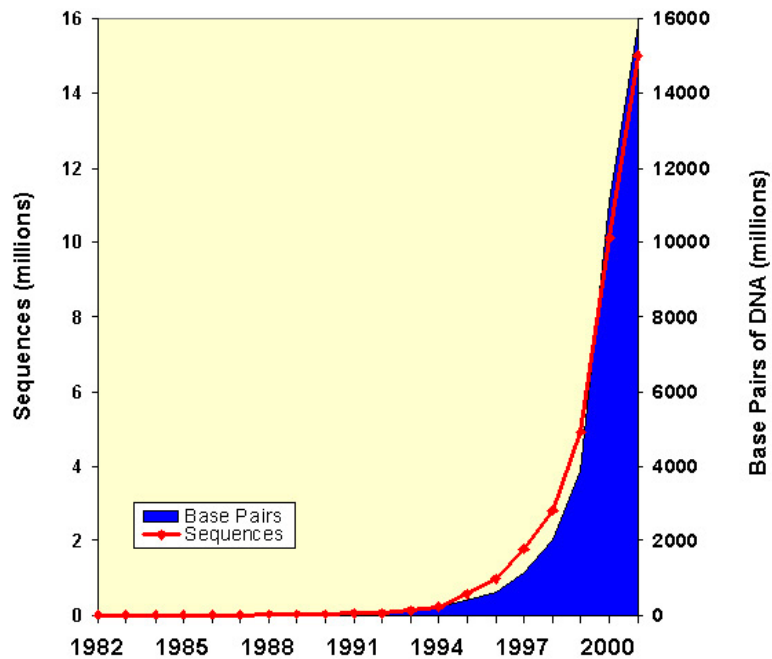
Beyond the **data deluge**: Data integration and bio-ontologies

Judith A. Blake *, Carol J. Bult

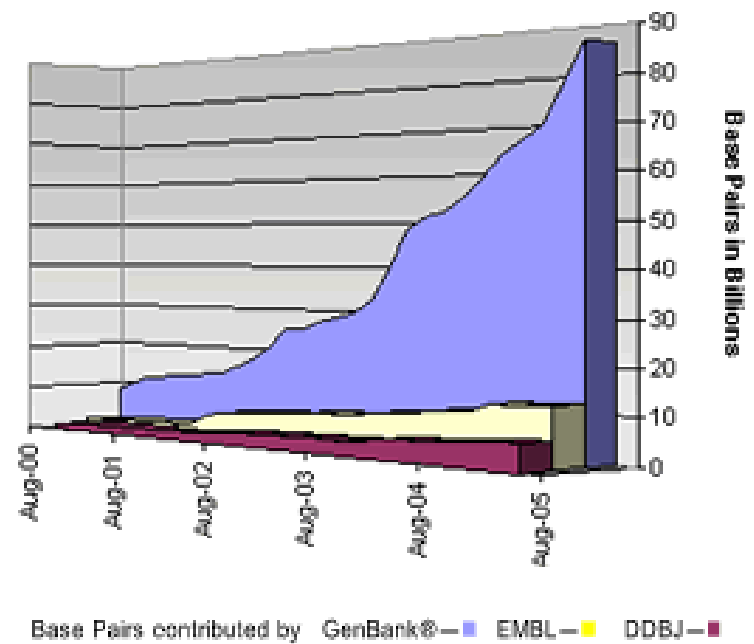
The Jackson Laboratory, Bar Harbor, ME, USA

Received 31 August 2005

Growth of GenBank

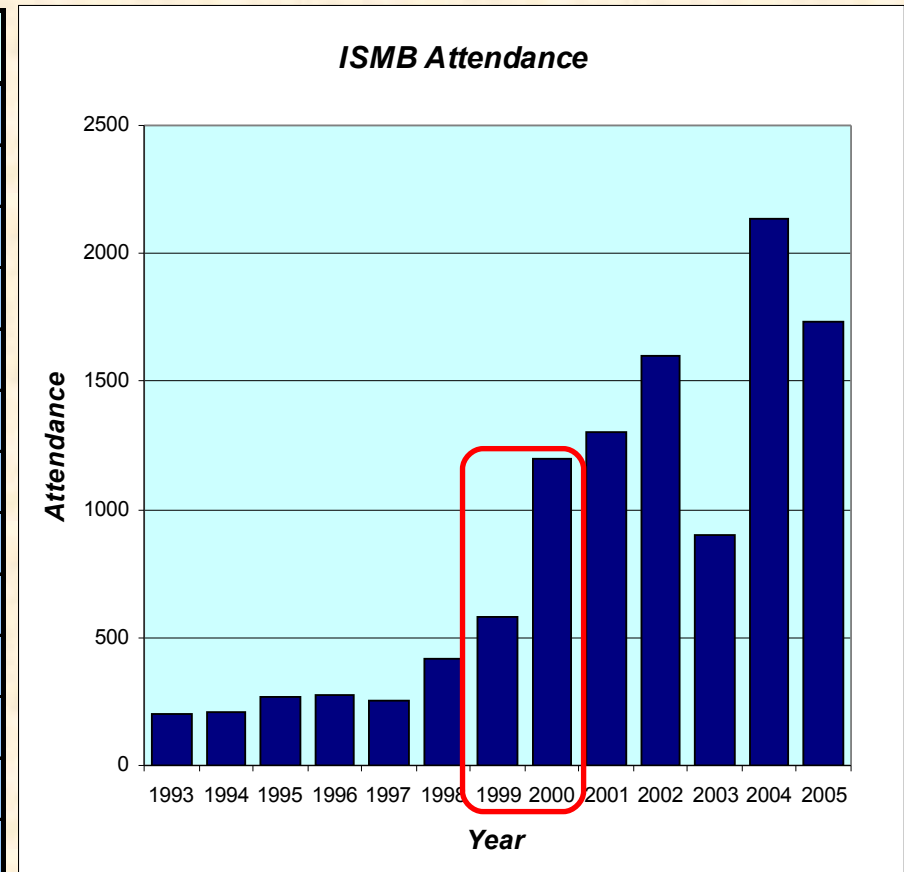


Growth of the International Nucleotide Sequence Database Collaboration



The number of Bioinformaticians is growing...

Year	Place	Attendance
1993	Bethesda, MD, USA	201
1994	Stanford, California, USA	205
1995	Cambridge, United Kingdom	270
1996	St. Louis, MO, USA	279
1997	Halkidiki, Greece	254
1998	Montréal, Québec, Canada	413
1999	Heidelberg, Germany	580
2000	La Jolla/ San Diego, CA, USA	1200
2001	Copenhagen, Denmark	1300
2002	Edmonton, Alberta, Canada	1600
2003	Brisbane, Australia	900
2004	Glasgow, Scotland, UK (With EOCB)	2136
2005	Detroit, Michigan, USA	1731



ISMB conferences attendance

...and is growing the number of their “products”!

Year	Articles	DB listed
2006	164	858
2005	137	719
2004	142	548
2003	95	386
2002	94	335
2001	73	281
2000	95	226
1999	86	201
1998	77	
1997	64	
1996	51	

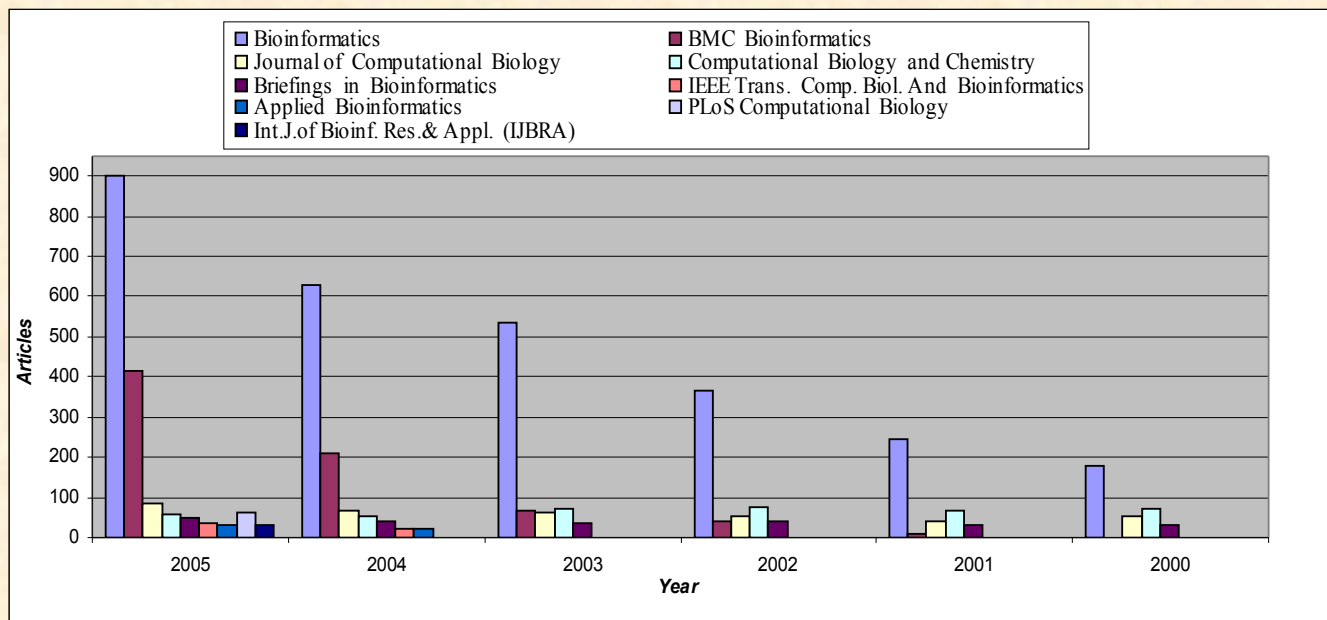
Database Issue

Year	Articles	WS presented
2005	159	166
2004	137	137
2003	106	131

Web server Issue

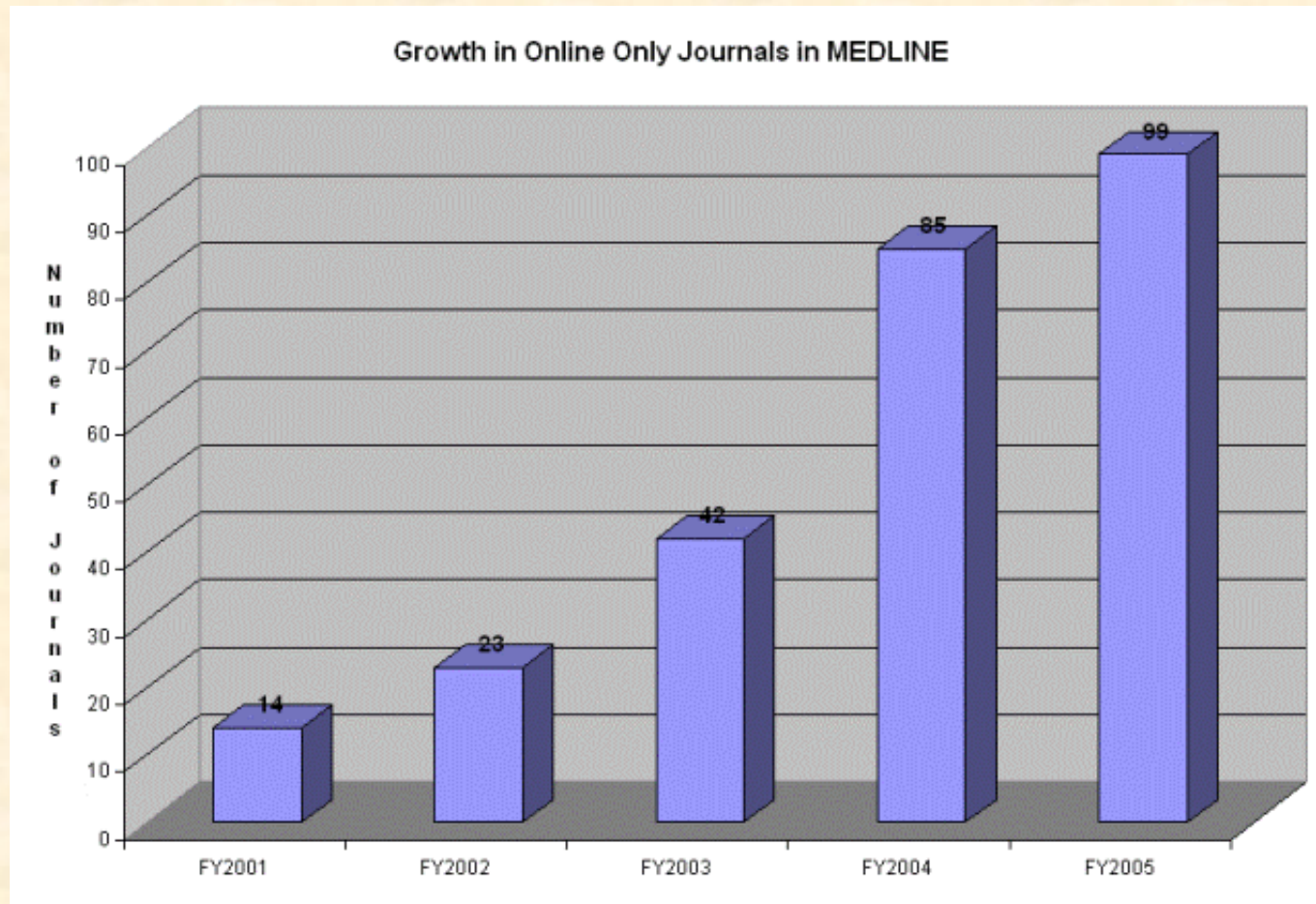
Nucleic Acid Research special issues

Year	Bioinformatics	BMC Bioinformatics	Journal of Computational Biology	PLOS Computational Biology	Computational Biology and Chemistry	Briefings in Bioinformatics	IEEE Trans. Comp. Biol. And Bioinformatics	Applied Bioinformatics	Int.J.of Bioinf. Res.&Appl. (IJBRA)
2005	90	44	8	6	5	4	3	2	2
2004	67	20	6		5	4	2	2	
2003	54	6	6		7	3			
2002	35	4	5		7	4			
2001	25	9	3		6	3			
2000	18	1	5		2	3			



Articles published “ONLY” from main bioinformatics journals

Another impressive growth



It has become really hard to find resources!

- It is difficult to track the evolution of a research area
- It is not existing any classification schema for bioinformatics resources
- For the bioinformatics domain is not existing any classification schema defined from an authority (like e.g. ACM, AMS)
- Keyword search vs Semantic search
- Take into account semantic relationships between resources
- Resources disappears (Many 404s)
- Multidisciplinarity is hard
- Easy to re-invent of the wheel

Perspective

Time to Organize the Bioinformatics Resourceome

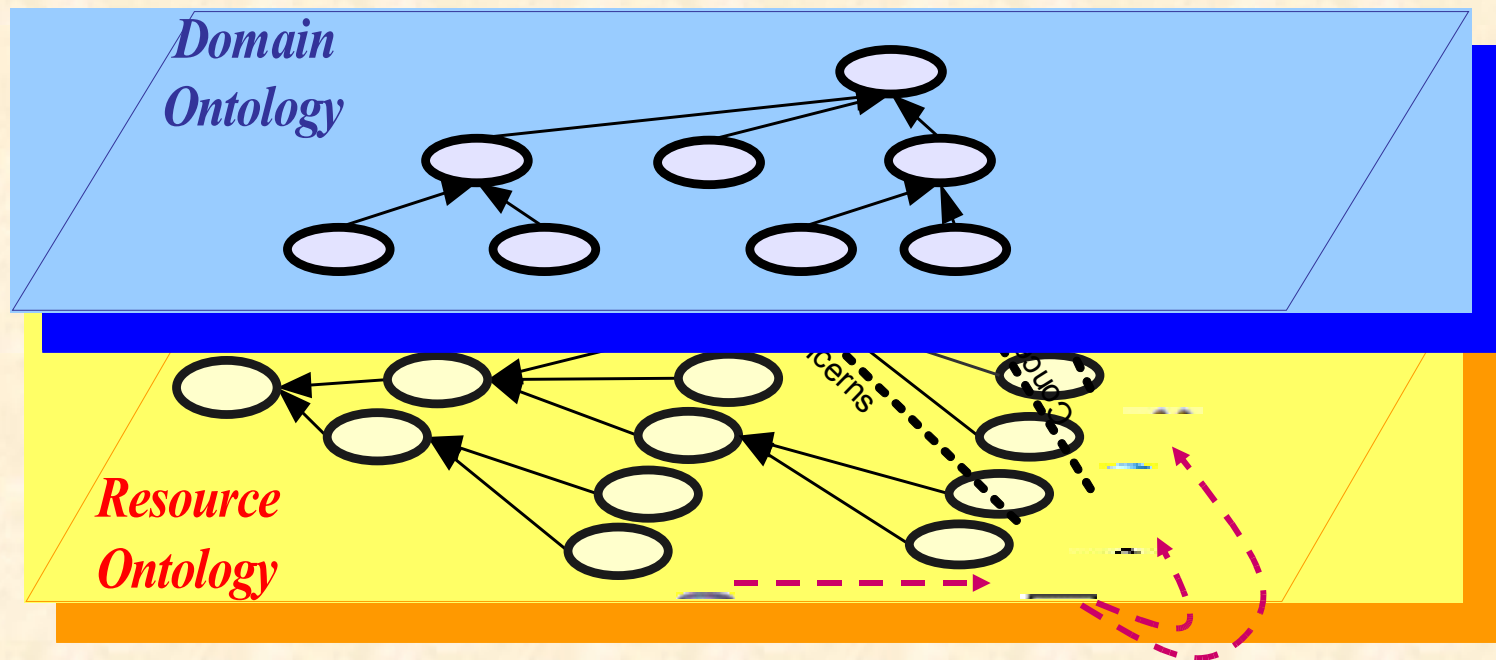
Nicola Cannata, Emanuela Merelli, Russ B. Altman*

We will be witnessing the birth of the artificial, or in-silico, scientist. — J. D. Wren [1]

The field of bioinformatics has blossomed in the last ten years, and as a result, there is a large and increasing number of researchers generating computational tools for solving problems relevant to biology. Because the number of artifacts has increased greatly, it is impossible for many bioinformatics researchers to track tools, databases, and methods in the field—or even perhaps within their own specialty area. More critically, however, biologist users and scientists approaching the field do not have a comprehensive index of bioinformatics algorithms, databases, and literature annotated with information about their context and appropriate use. We suggest that the full set of bioinformatics resources—the “resourceome”—should be explicitly characterized and organized. A hierarchical and machine-understandable organization of the field, along with rich cross-links (an ontology!) would be a useful start. It is likely that a distributed development approach would be required so that those with focused expertise can classify resources in their area, while providing the metadata that would allow easier access to useful existing resources.

keyword searching [5]. However, the lack of standard terms makes sensitive and specific searches difficult. In addition, most search hits confound papers, Web sites, tools, departments, and people in a manner that makes extracting useful information very difficult.

Recognizing this limitation, there have been some grassroots attempts to organize the bioinformatics resourceome. Among the most famous are the “archaeological” Pedro’s List—a list of computer tools for molecular biologists (http://www.public.iastate.edu/~pedro/research_tools.html)—and the Expasy Life Sciences Directory, formerly known as the Amos’s WWW links page (<http://www.expasy.org/links.html>). The Bioinformatics Links Directory (http://www.bioinformatics.ubc.ca/resources/links_directory/) today contains more than 700 curated links to bioinformatics resources, organized into eleven main categories, including all the databases and Web servers yearly listed in the dedicated *Nucleic Acids Research* special issues [6]. The National Center for Biotechnology Institute has tried to make access to its suite of tools transparent, with moderate success. Many Web sites can be found listing “useful sites,” especially concerning special interest or limited topics (e.g., microarrays, text mining, and gene regulation). But all of these efforts are limited by the difficulty in maintaining



Proposal of an architecture for Resourceomes (Cannata et al. Submitted)

OWL : Thing

isa

isa

Auxiliary Concept

Resource

isa isa isa

isa

isa

Data Feature Affiliation Resource Access Type

Actor

isa

isa

Institution

Person

isa isa

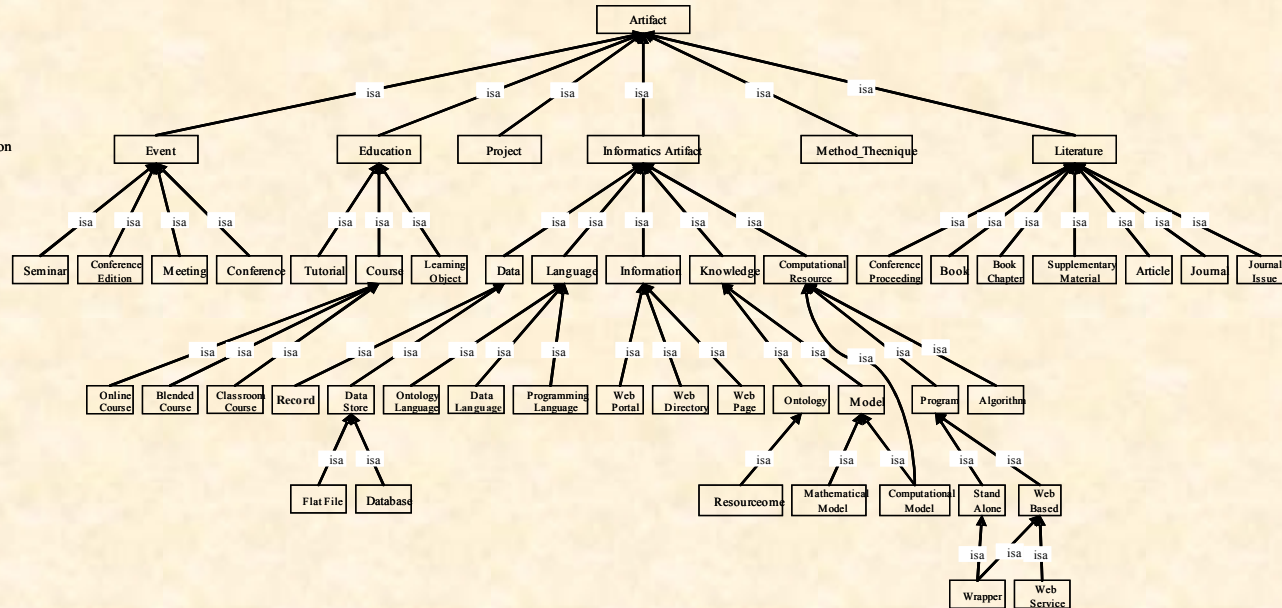
isa isa isa isa isa isa

Data Format Abstract Data Type

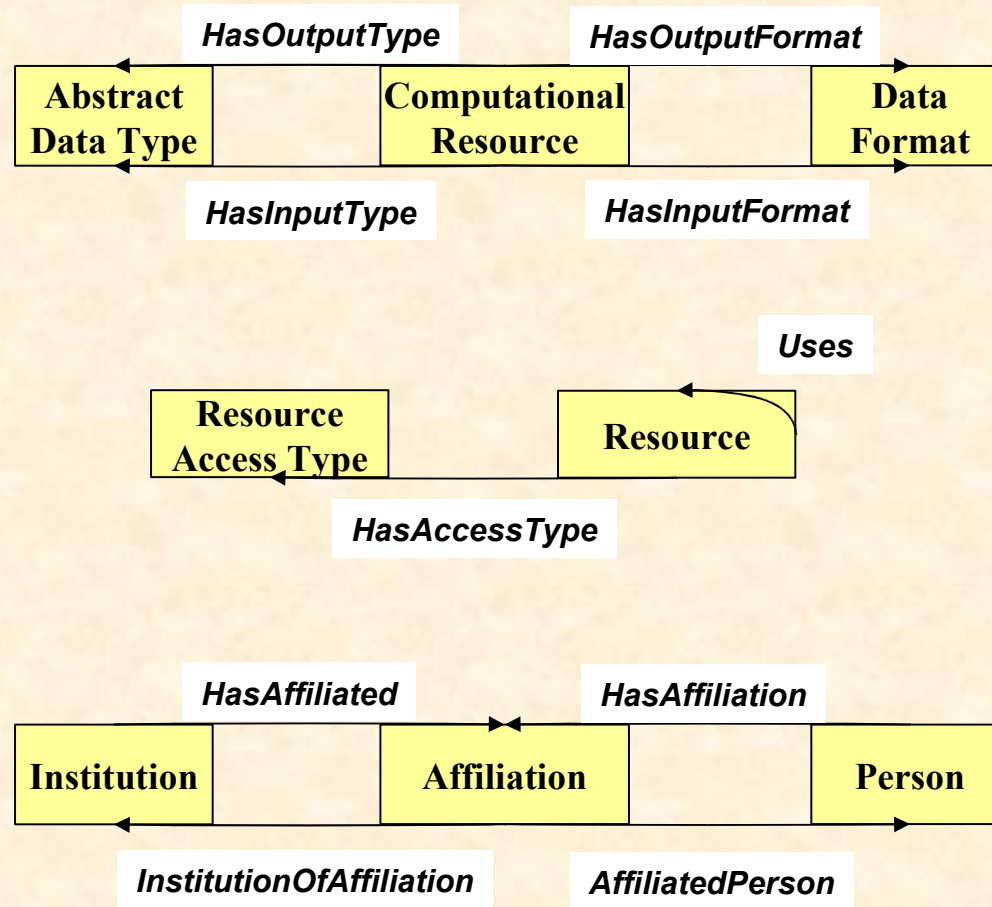
Domain Data Type Basic Data Type

isa isa Institute/ Department/ Division Research Group Firm Organization University

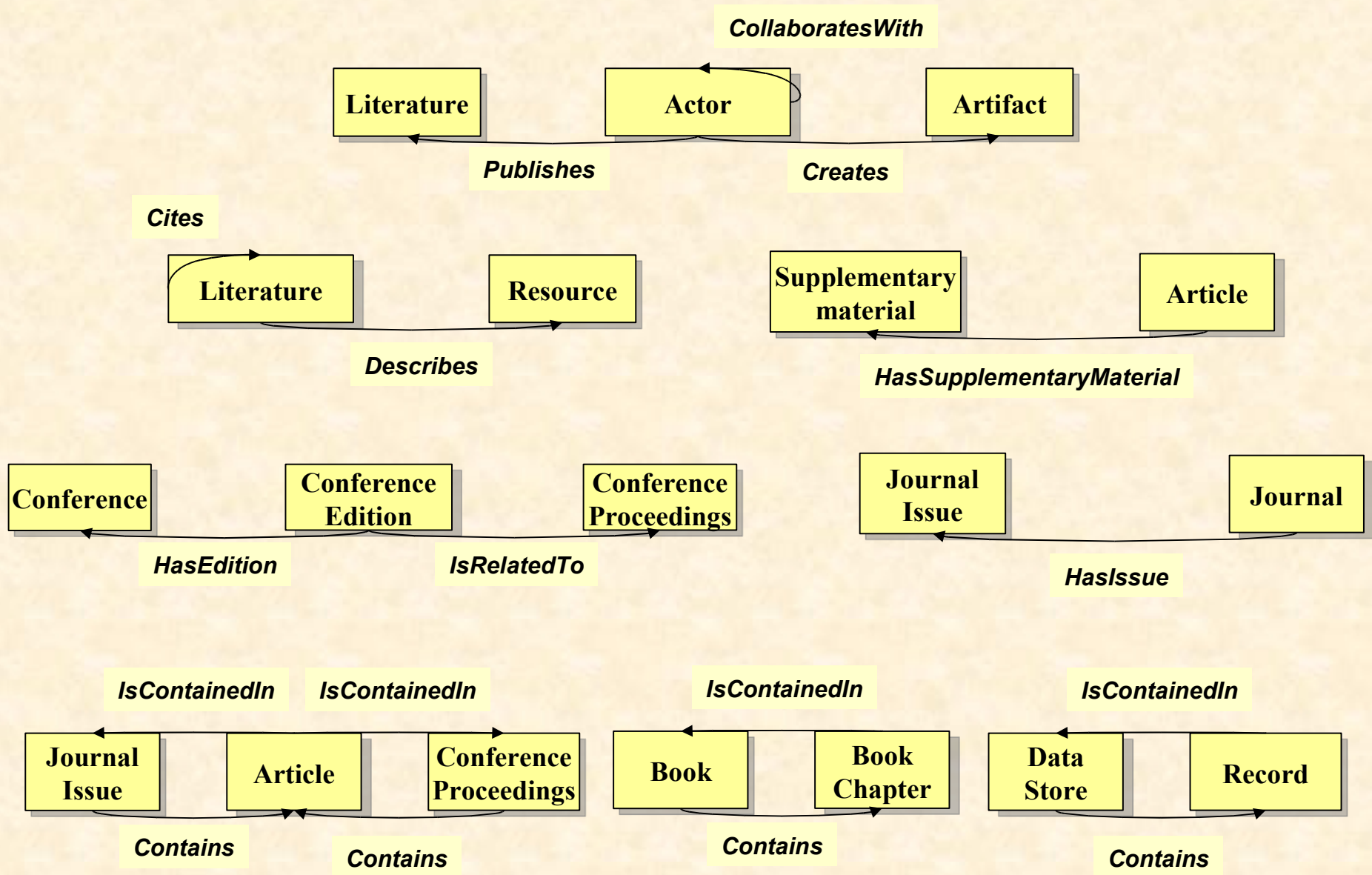
Special Interest Group



A proposal for a resource ontology

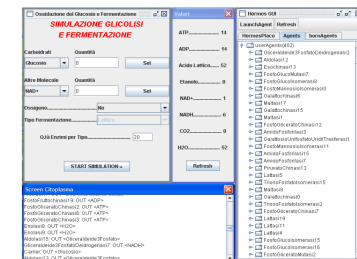
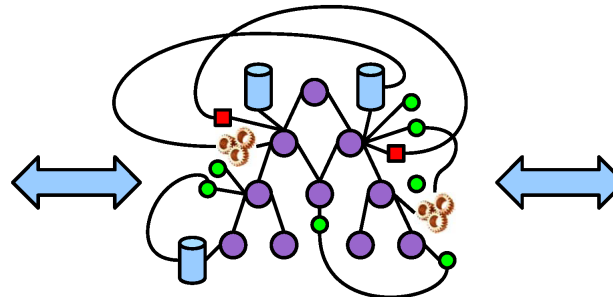
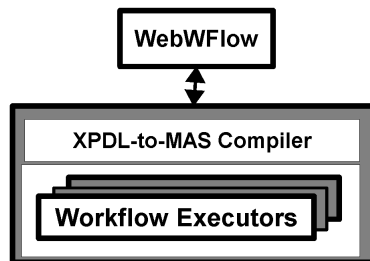
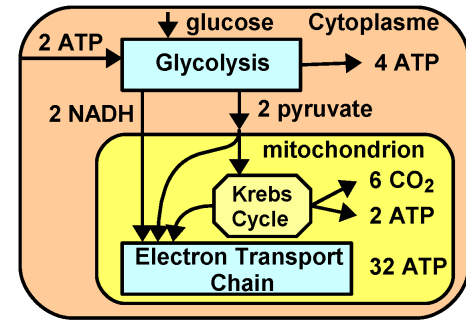
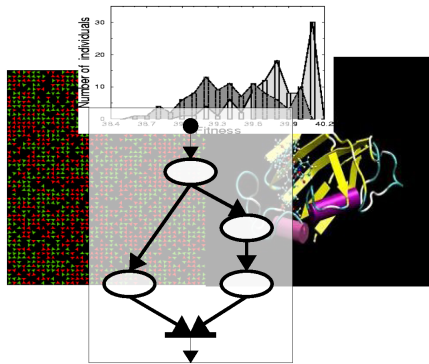


Semantic relationships between resources



Semantic relationships between resources

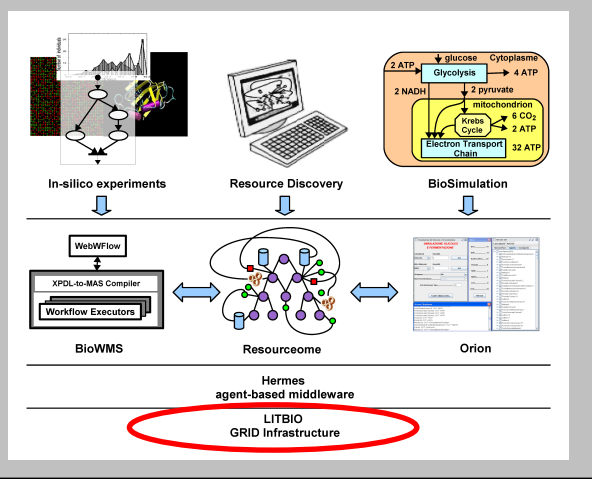
A GRID-based multilayer architecture for bioinformatics



Hermes agent-based middleware

LITBIO
GRID Infrastructure

A GRID-based multilayer architecture for bioinformatics



The Grid

- The concept of Grid is evolving (or better specializing)
- In our case it will be a “classical” Grid infrastructure provided in the LITBIO project (nodes with HPCs)

BIOINFORMATICS



myGrid: personalised bioinformatics information grid

Robert D. Stevens¹, Alan J. Robinson

The Anatomy of the Grid

Enabling Scalable Virtual Organizations *

Ian Foster^{*†} Carl Kesselman[§] Steven Tuecke^{*}

{foster, tuecke}@mcs.anl.gov, carl@isi.edu

Abstract

“Grid” computing has emerged as an important new field, distinguished from conventional distributed computing by its focus on large-scale resource sharing, innovative applications, and, in some cases, high-performance orientation. In this article, we define this new field. First, we review the “Grid problem,” which we define as flexible, secure, coordinated resource sharing among dynamic collections of individuals, institutions, and resources—what we refer to as *virtual organizations*. In such settings, we encounter unique authentication, authorization, resource access, resource discovery, and other challenges. It is this class of problem that is addressed by Grid technologies. Next, we present an extensible and open *Grid architecture*, in which

The Semantic Grid: A Future e-Science Infrastructure

David De Roure, Nicholas R. Jennings and Nigel R. Shadbolt¹

Dept of Electronics and Computer Science,
University of Southampton,
Southampton SO17 1BJ, UK

{dder,nrj,nrs}@ecs.soton.ac.uk

Available at
www.ComputerScienceWeb.com
POWERED BY SCIENCE @ DIRECT®

Future Generation Computer Systems 20 (2004) 101–111

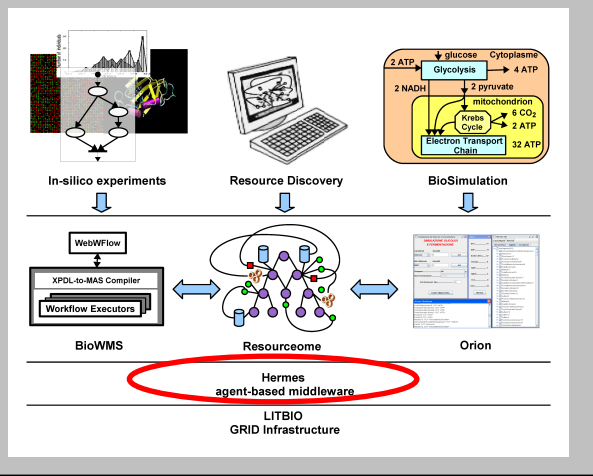


www.elsevier.com/locate/future

Comparative assessment approach for Knowledge Grid

Hai Zhuge*, Jie Liu

Key Group, Key Lab of Intelligent Information Processing, Institute of Computing Technology,
Chinese Academy of Sciences, P.O. Box 2704-28, 100080 Beijing, China



Hermes

- An agent-based middleware developed at Camerino
- Provide a common runtime support to agent based applications

What is an agent?

A computer system capable of **flexible**, **autonomous** (problem-solving) actions, situated in dynamic, unpredictable and typically multi-agent environment.

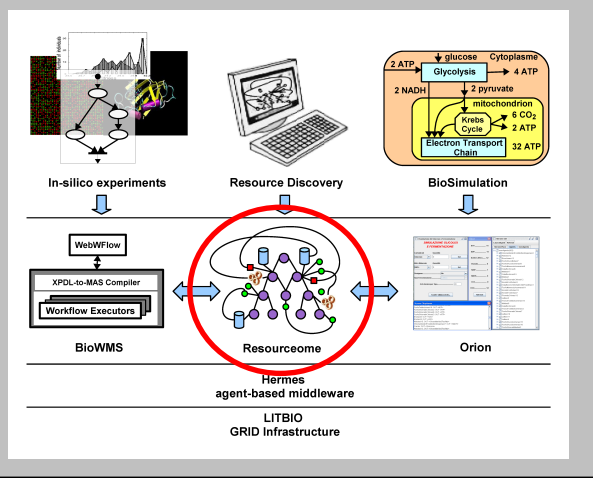
Active and reactive: responds in timely fashion to environmental change

proactive: acts in anticipation of future goals

cooperative: communicates with other agents to reach its goal

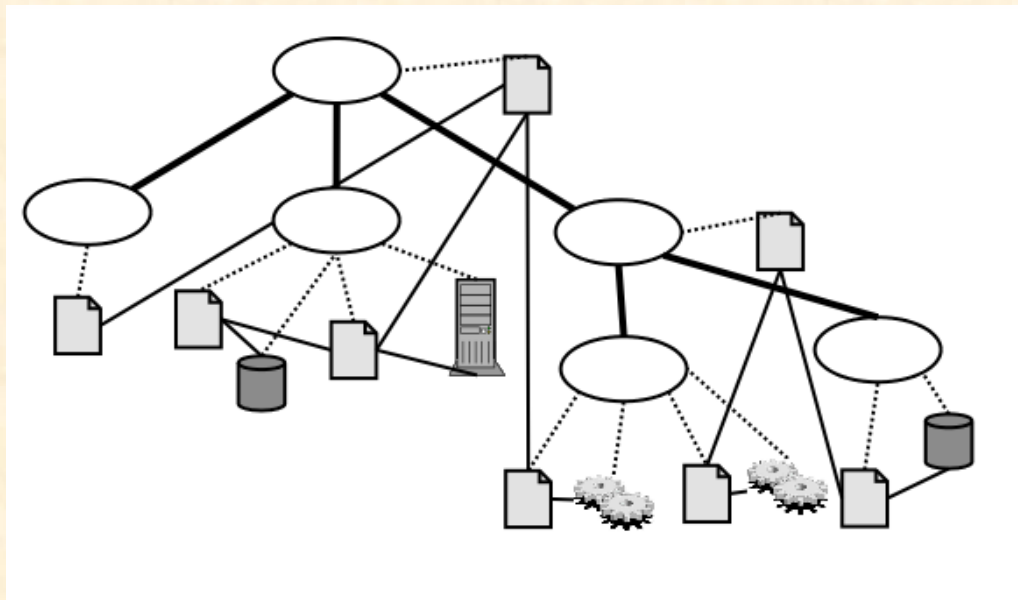
mobile: moves across distributed environments (execution platforms)

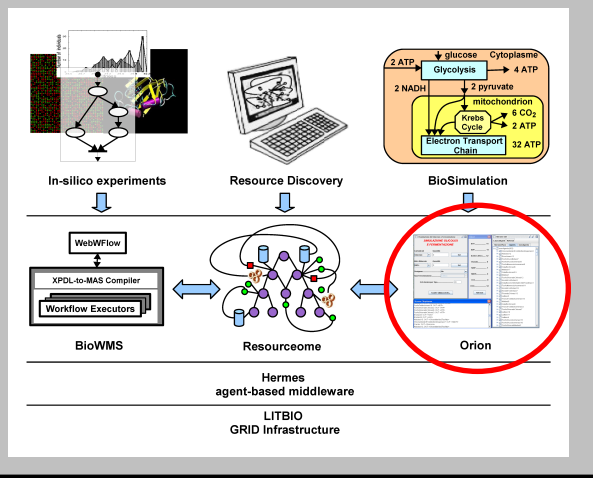
intelligent: reasons over its knowledge base and by managing ontologies



Resourceome

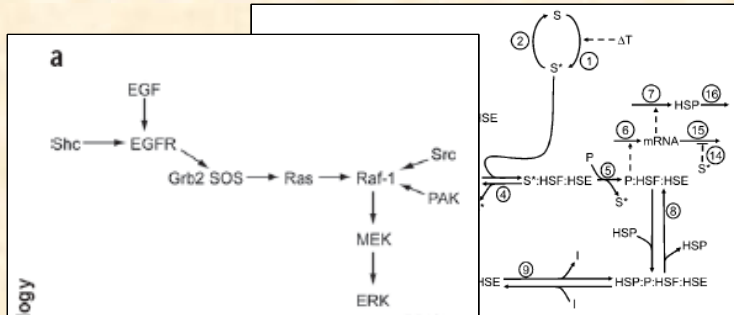
- A conceptual map of a domain with related resources
- It is the Pivot of the architecture
- Will permit to semantically organize workflows, activities, services
- Will permit to semantically organize systems biology resources and knowledge
- Will permit to “reason” over resources
- Agents will build, maintain and keep it “alive”





Orion

- An agent-based framework for systems biology
- Behavioural modeling of molecular entities
- We are starting with metabolic reactions...



<http://www.nature.com/naturebiotechnology>

$$\frac{dx_3}{d\tau} = \kappa_1(\Gamma_1 \times x_1 \times x_2 - x_3) + \kappa_2(x_4 - \Gamma_2 \times x_{12} \times x_3) + \kappa_5(x_8 - \Gamma_5 \times x_{11} \times x_3) \quad (9)$$

$$\frac{dx_4}{d\tau} = \kappa_2(\Gamma_2 \times x_{12} \times x_3 - x_4) - \kappa_6 x_4 \quad (10)$$

$$\frac{dx_5}{d\tau} = \kappa_5 x_4 + \kappa_3(x_6 - \Gamma_3 \times x_{11} \times x_3) \quad (11)$$

$$\frac{dx_6}{d\tau} = \kappa_3(\Gamma_3 \times x_{11} \times x_5 - x_6) + \kappa_4(x_7 - \Gamma_4 \times x_{13} \times x_6) \quad (12)$$

$$\frac{dx_7}{d\tau} = \kappa_4(\Gamma_4 \times x_{13} \times x_6 - x_7) - \kappa_1 x_7 \quad (13)$$

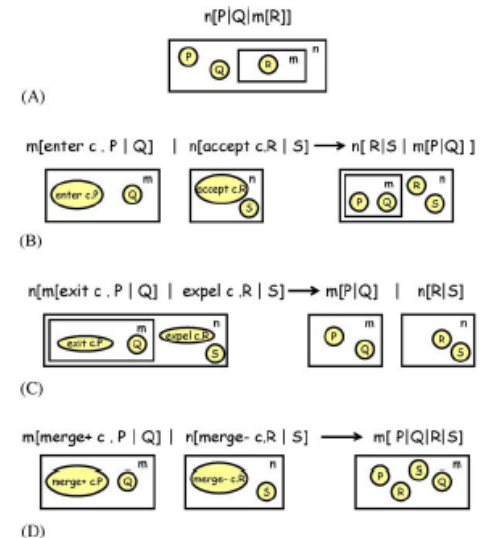
$$\frac{dx_8}{d\tau} = \kappa_1 x_7 + \kappa_5(\Gamma_5 \times x_{11} \times x_3 - x_8) + \kappa_6(\Gamma_6 \times x_2 \times x_9 - x_8) \quad (14)$$



```

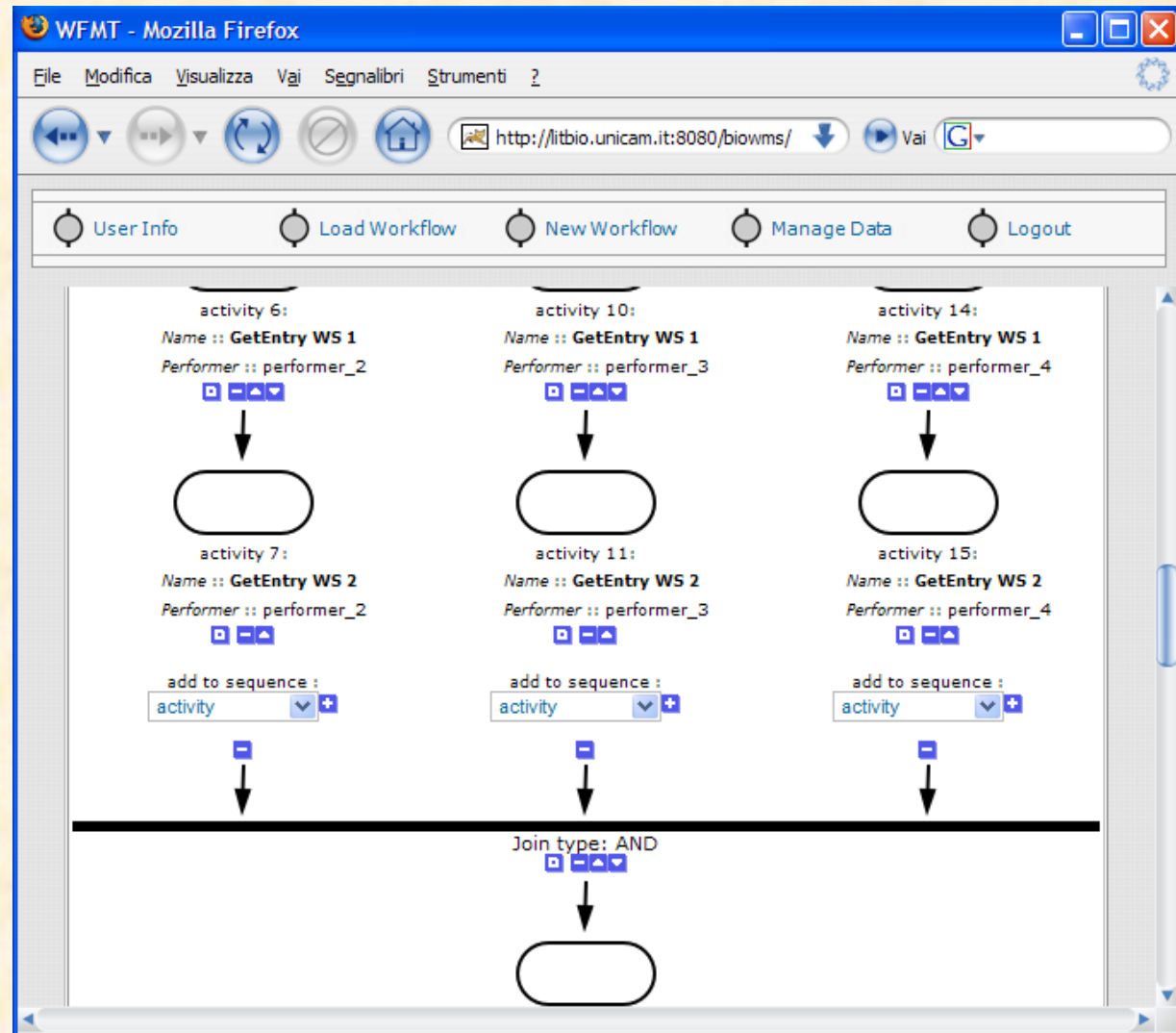
SYSTEM ::= CYCLIN|CDK|CDH1|CDC14|CKI|CLOCK
CYCLIN ::= (ν bb) BINDING_SITE
BINDING_SITE ::= (lb(bb), R4).CYCLIN_BOUND
CYCLIN_BOUND ::= DEGCYC + DEGCKI + CYC_CDK_CKI
DEGCYC ::= (degp, R1).degc.0
DEGCKI ::= (degd, R3).CYCLIN_BOUND
CYC_CDK_CKI ::= (bind(bb), R11).bb.TRIM
TRIM ::= DIM + NOTHING
DIM ::= (removed, R10).DIM
NOTHING ::= (d, R9).NOTHING
CDK ::= (lb(cbb), R2).CDK_CATALYTIC
INACTCDH1 ::= (inact, R7).INACTCDH1
INACTCAT ::= (inact, R8).INACTCAT
  
```

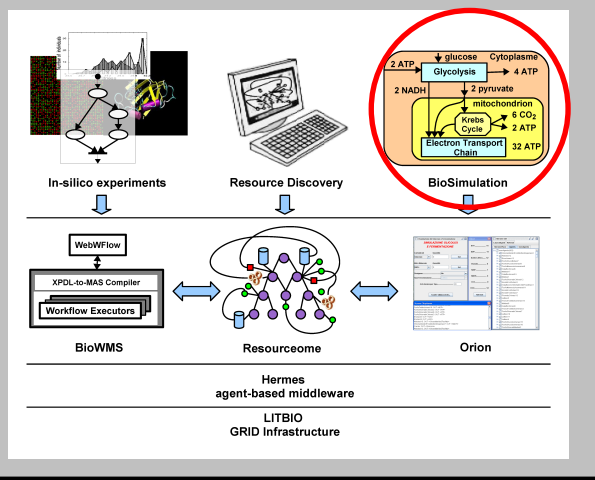
A. Regev et al. / Theoretical Computer Science 325 (2004) 141–167



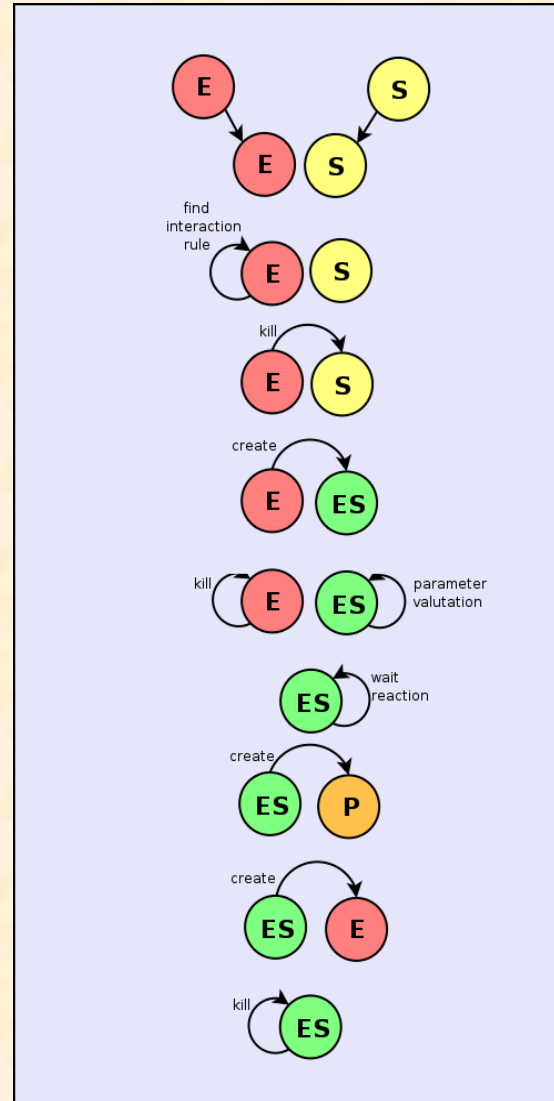
Will agents permit to “bridge” the “two” systems biology?

The diagram illustrates a bioinformatics workflow and its underlying infrastructure. At the top left, a red circle highlights the **In-silico experiments** stage, which involves data analysis (histogram and scatter plot) and a network diagram. This stage leads to **Resource Discovery**, represented by a laptop. The next step is **BioSimulation**, which is detailed with a metabolic pathway diagram. This diagram shows **Glycolysis** in the **Cytoplasm** converting glucose to pyruvate, producing 2 ATP and 2 NADH. Pyruvate then enters the **mitochondrion**, where it undergoes the **Krebs Cycle** (producing 8 CO₂ and 2 ATP) and the **Electron Transport Chain** (producing 32 ATP). The **BioSimulation** stage leads to the **Orion** web interface, which displays resource management information. The **Orion** interface is supported by the **BioWMS** (BioWorkflow Management System), which includes a **WebWorkflow** component, an **XPDL-to-MAS Compiler**, and **Workflow Executors**. The **BioWMS** is connected to a central **Resource** component, which is represented by a network diagram of nodes and edges. This resource component is supported by the **Hermes agent-based middleware** and the **LITBIO GRID Infrastructure**.





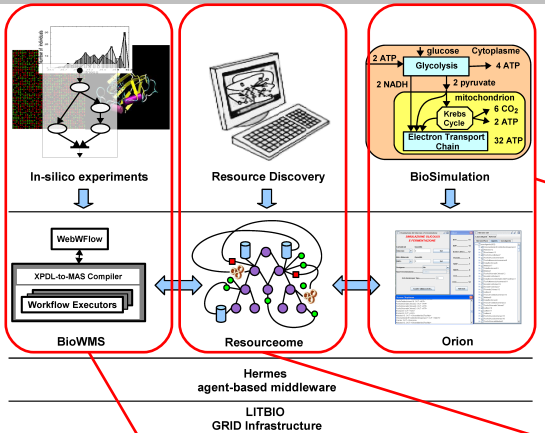
Simulation of biological process



Conclusions

- LITBIO aims to offer a Laboratory for Interdisciplinary Technologies in Bioinformatics
- Bioinformatics is fast changing and growing
- Resourceome can make interdisciplinarity easier
- We propose an agent-based multi-layer architecture which will lay on the LITBIO Grid
- User-friendly and user-assisting
- We hope it will make easier the life of molecular biologists, bioinformaticians and systems biologists of today and tomorrow

Acknowledgment



BIOWMS
Lorenzo Scortichini

RESOURCEOME
Sergio Gabrielli
Luana Leoni
Francesca Piersigilli
Leonardo Vito

Rosario Culmone

ORION
Arianna Baldoncini
Claudio Forcato
Michele Mattioni

Mauro Angeletti
Riccardo Piergallini

