

---

# An agent-based layered middleware as tool integration

---

Flavio Corradini  
University of L'Aquila  
ITALY

Leonardo Mariani  
University of Milano  
ITALY

**Emanuela Merelli**  
University of Camerino  
ITALY

Helsinki  
FSE/ESEC 2003  
Tool Integration Workshop

---

# Outline

- The Tool Integration problem in the Bioinformatics Domain
- The Workflow-based Task Coordination (High Level Tool Integration)
- The Wrapper-based Data Integration (Low level Tool Integration)
- The Proposed Approach:  
An Agent-based Middleware for Tool Integration
- Preliminary Results
- Future Activities and Conclusions

# The Tool Integration problem in Bioinformatics Domain

**Problem:** To find the crystallographic structure of the 10 proteins more similar to a new genetic sequence, e.g. X=MEEP ... DSD,

**Objective:** To use several Bioinformatics Software Tools available on Internet in order to find the wanted result

For **Tool Integration** we mean

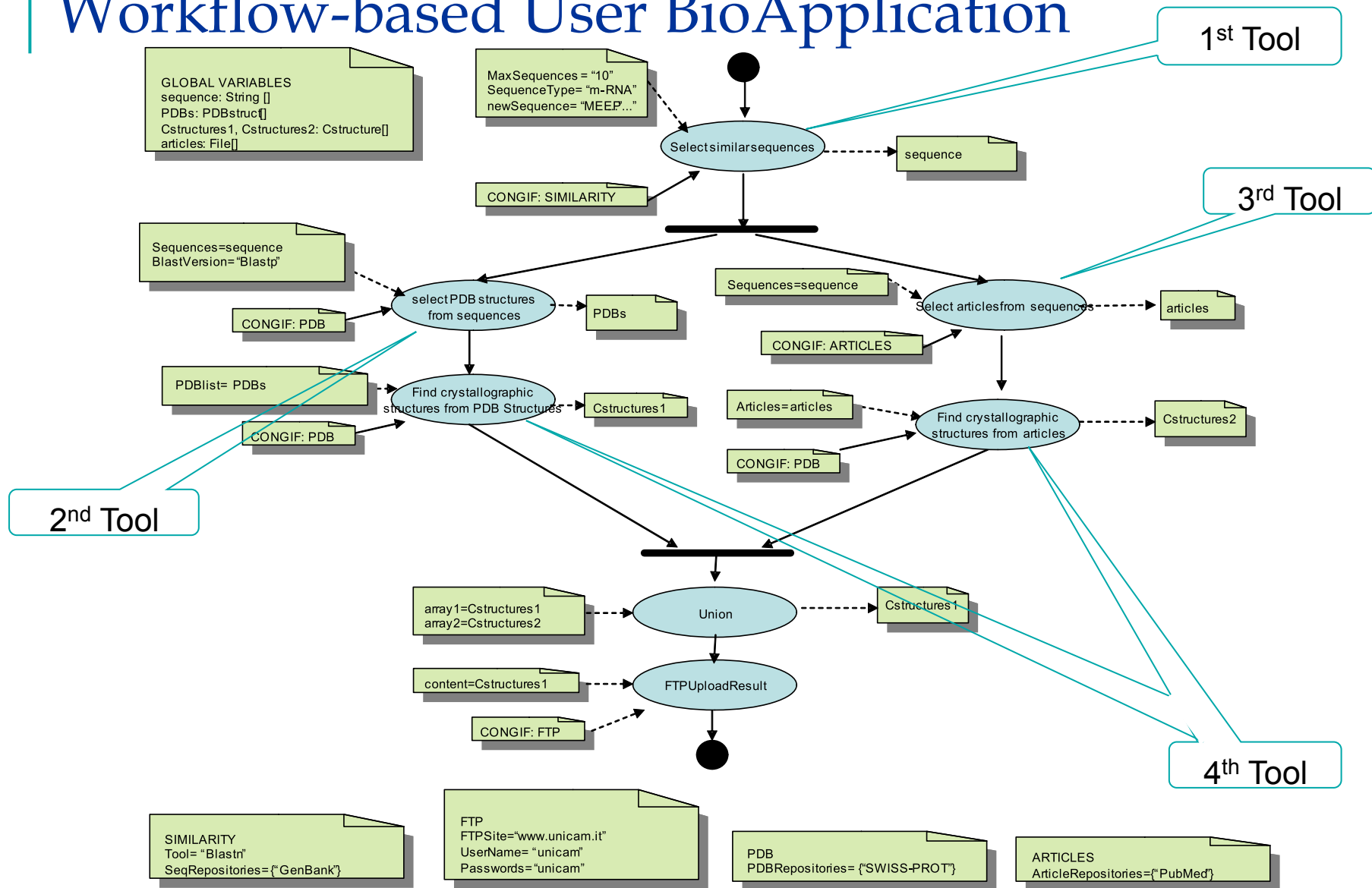
- 1) Supporting Tasks Coordination
- 2) Allowing Data Integration

in order to automatically execute an experiment

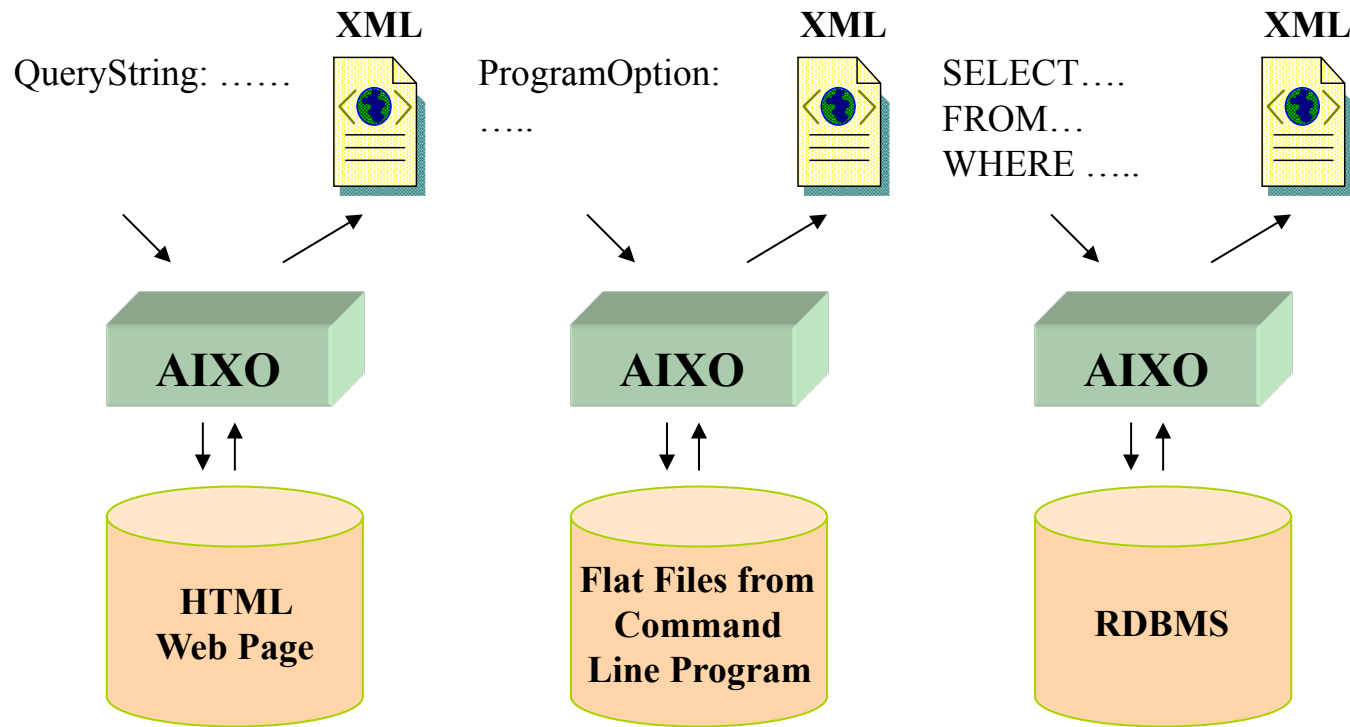
1. *Select the 10 proteins more similar to the X=MEEP ... DSD sequence*
  - by using **BLAST** in **Protein DataBank**
2. *Search for the PDB ID (crystallographic structure identifier) of each selected proteins,*
  - by using **BLAST** in **SWISS-PROT** at **EMBL/EBI**
  - by retrieving from **PubMed** via **Entrez Retrieval System** at **NCBI**, abstracts containing PDB-ID information
3. *Search for the Crystallographic Structure of any selected PDB ID*
  - find 3-D biological macromolecular structure in **Protein DataBank** repository

**Aim:** To **integrate** the four **Bioinformatics tools** freeing the Bioscientist from the need to continuous interact with remote sites.

# Workflow-based User BioApplication



# Wrapper-based System: general scenario



---

# Wrapper-based System: Bioinformatics Tools

## *Tool 1:*

Environment: NCBI (WebSite): [html format](#)  
Data: GenBank (DB): **proprietary format**  
Tool: BLASTn (Algorithm): Takes nucleotides sequences in FASTA format, GenBank Accession numbers or GI numbers and compares them against the NCBI [nucleotide databases](#)  
Output: GenBank Format

## *Tool 2:*

Environment: EMBL-EBI (WebSite): **html format**  
Data: Swiss-Prot (DB): **proprietary format**  
Tool: BLASTp (Algorithm): Takes protein sequences in **FASTA format**, GenBank Accession numbers or GI numbers and compares them against the NCBI [protein databases](#).  
Output: FASTA format

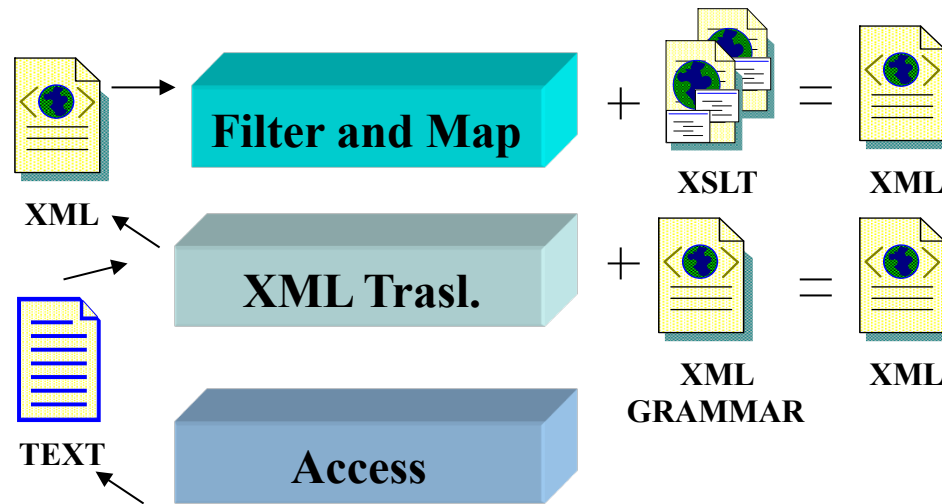
## *Tool 3:*

Environment: NCBI (WebSite): [html format](#)  
Data: PubMed & MEDLINE: **ANS.1 format**  
Tool: [Entrez Retrieval System](#)  
Output: XML

## *Tool 4:*

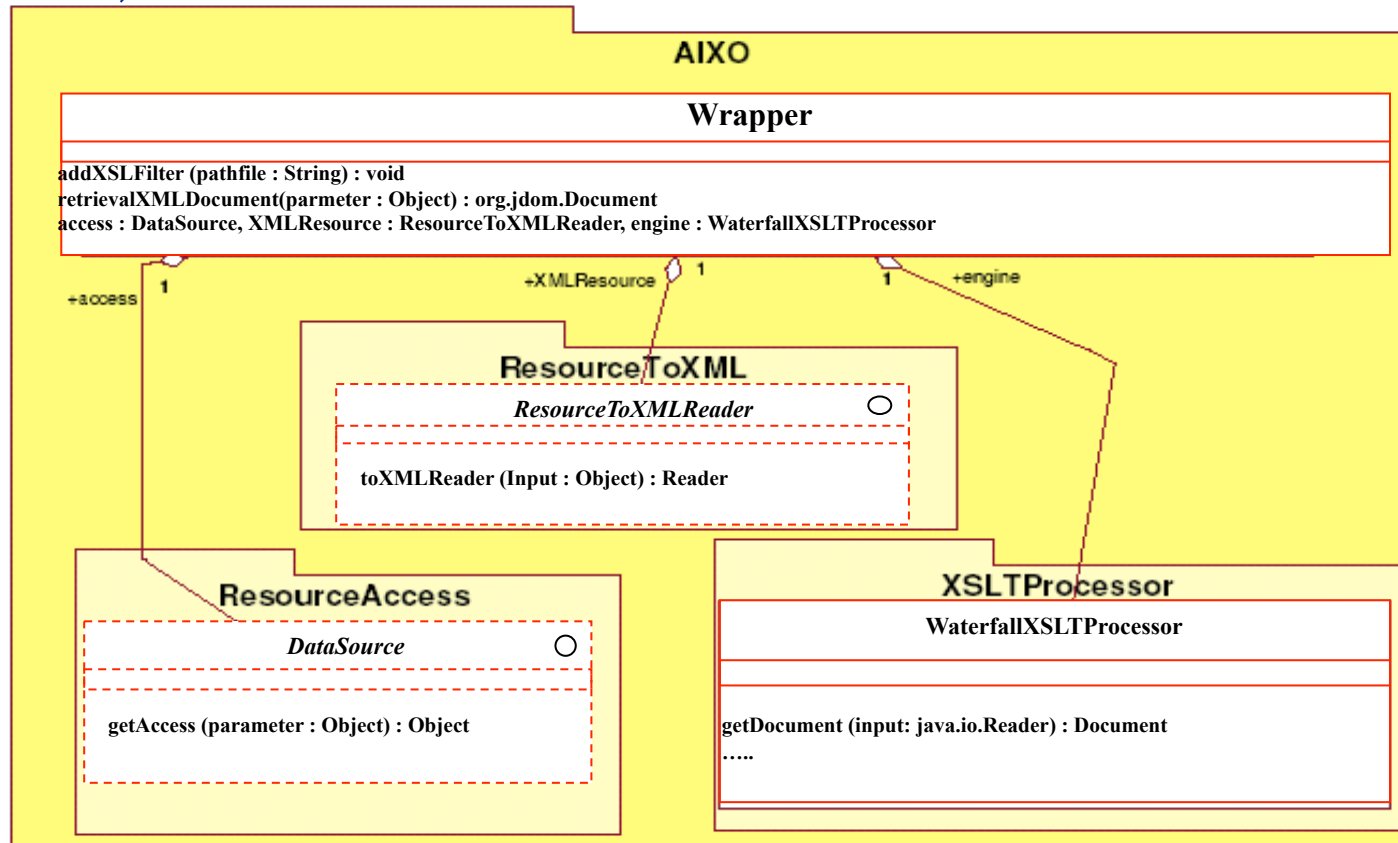
Environment: Protein DataBank web site  
Data: PDB(DB): **proprietary format**  
Tool: FASTA (Algorithm):  
Output: FASTA Format

# Wrapper-based System: Retrieval MedLine articles about P53 proteine



```
<...>
<entry>
  <ID name="P53_HUMAN" type="STANDARD" molecule="PRT" lenght="393"/>
  <AC value="P04637"/> <AC value="Q16848"/> <AC value="Q9UB12"/>
  <DT day="13" month="AUG" year="1987" rel="05"/>
</entry>
```

# Wrapper-based System: the software architecture (AIXO)



- **DataSource**: HTTP, RDBMS, Command Line program,....
- **ResourceToXMLReader**: HTML, FlatFile, ...



---

# The Tool Integration Problem in Activity-Based Applications

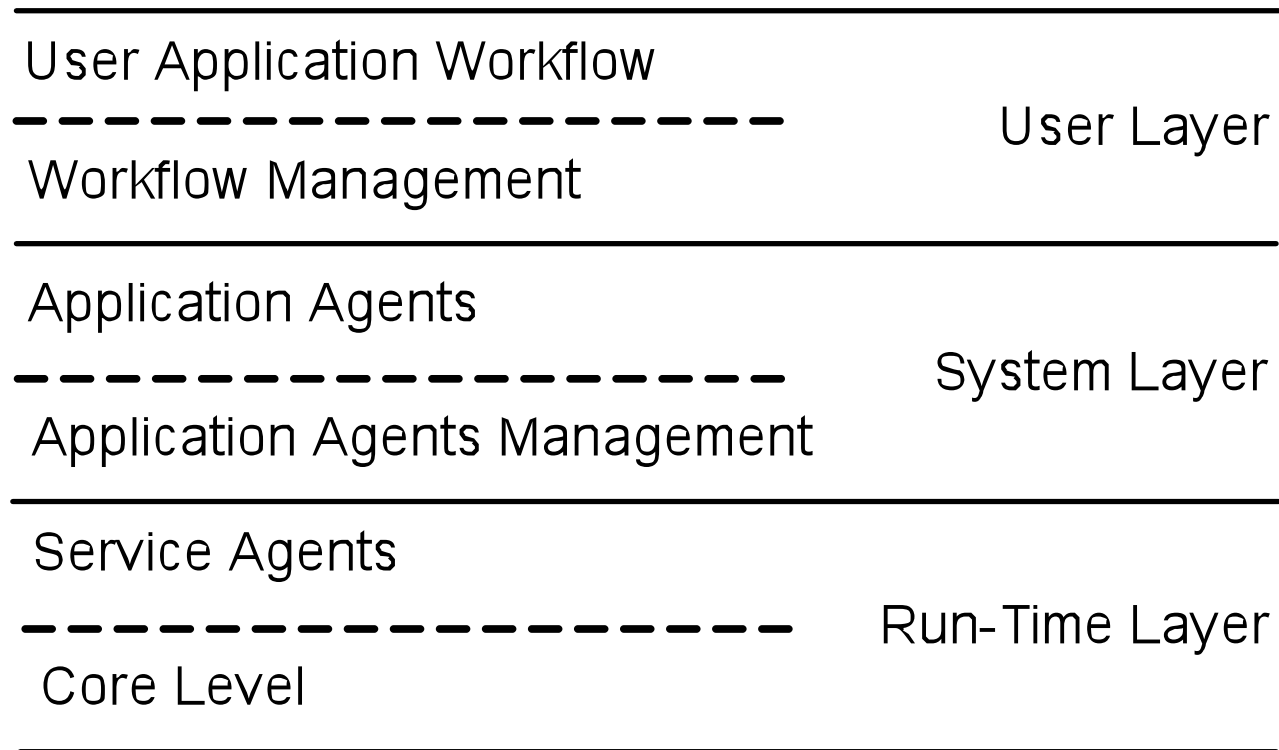
**Problem:** To *integrate and coordinate multiple software tools* for retrieving and integrating heterogeneous, distributed and frequently redundant data

**Objective:** To *integrate and coordinate several software tools* in order to provide a uniform way and an high level of abstraction for users

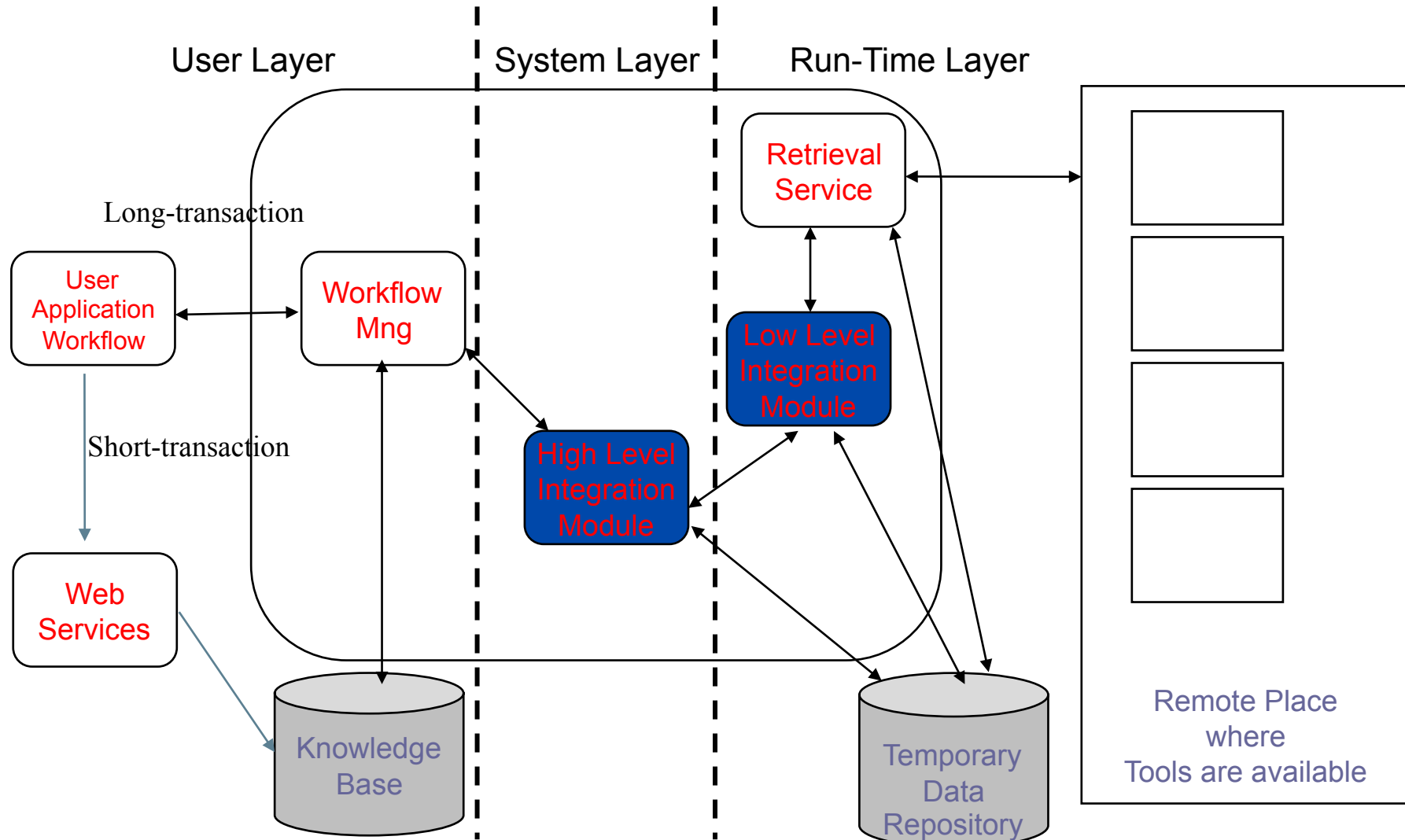
**Aim:** To define an *integrated environment* freeing the user from the need to know details on data repository and to coordinate the intermediate steps of an experiment (tasks)

**Proposed Approach:** To define an application as a workflow of tasks; to coordinate the execution of cooperative tasks by using software agent tools

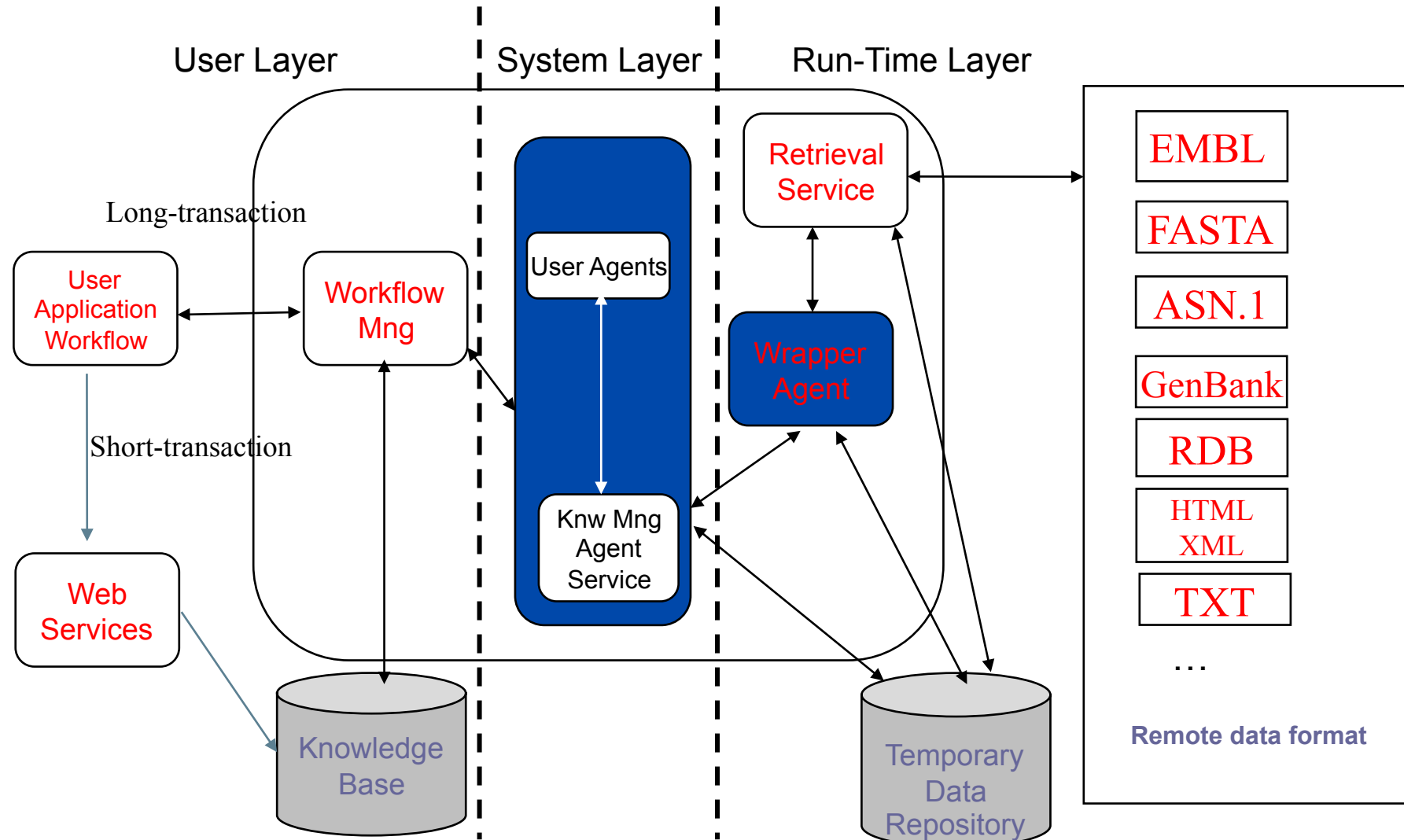
# System's software architecture



# A general system's architecture

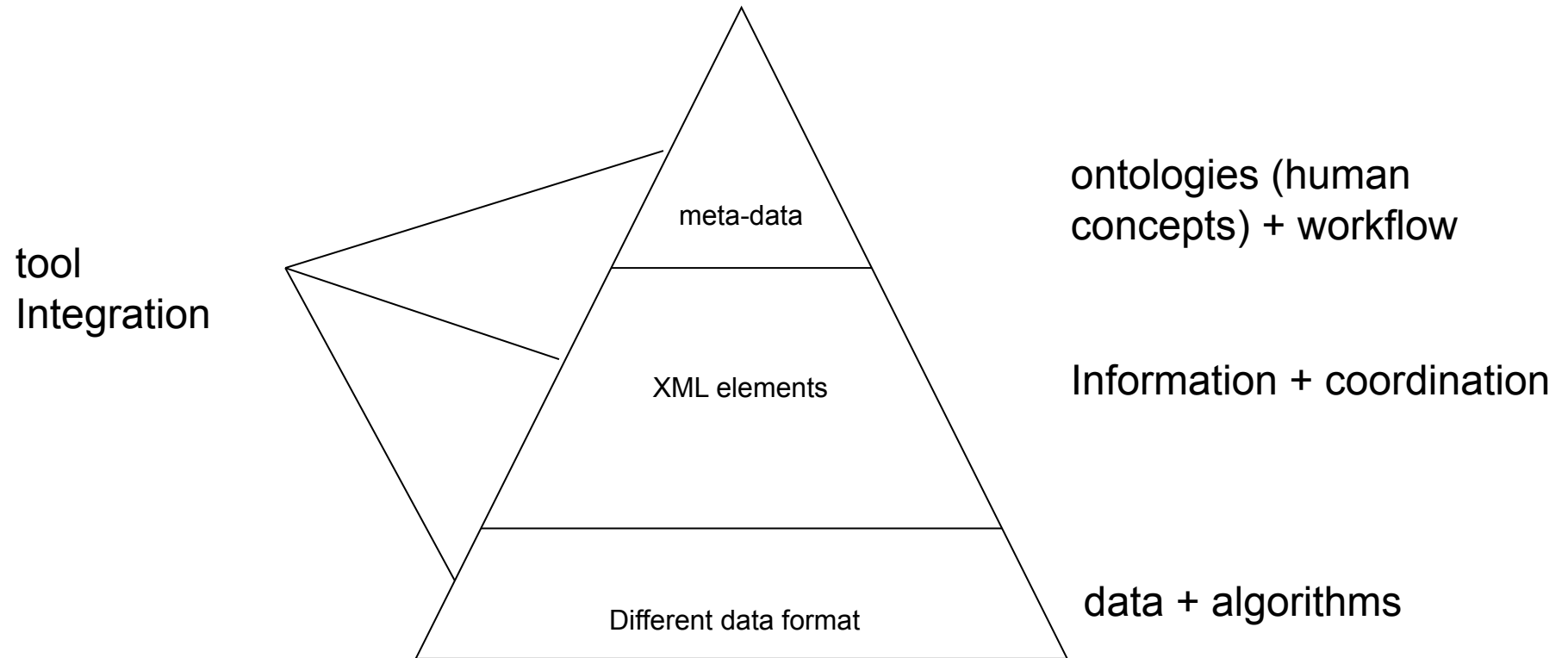


# Agent-based System Architecture

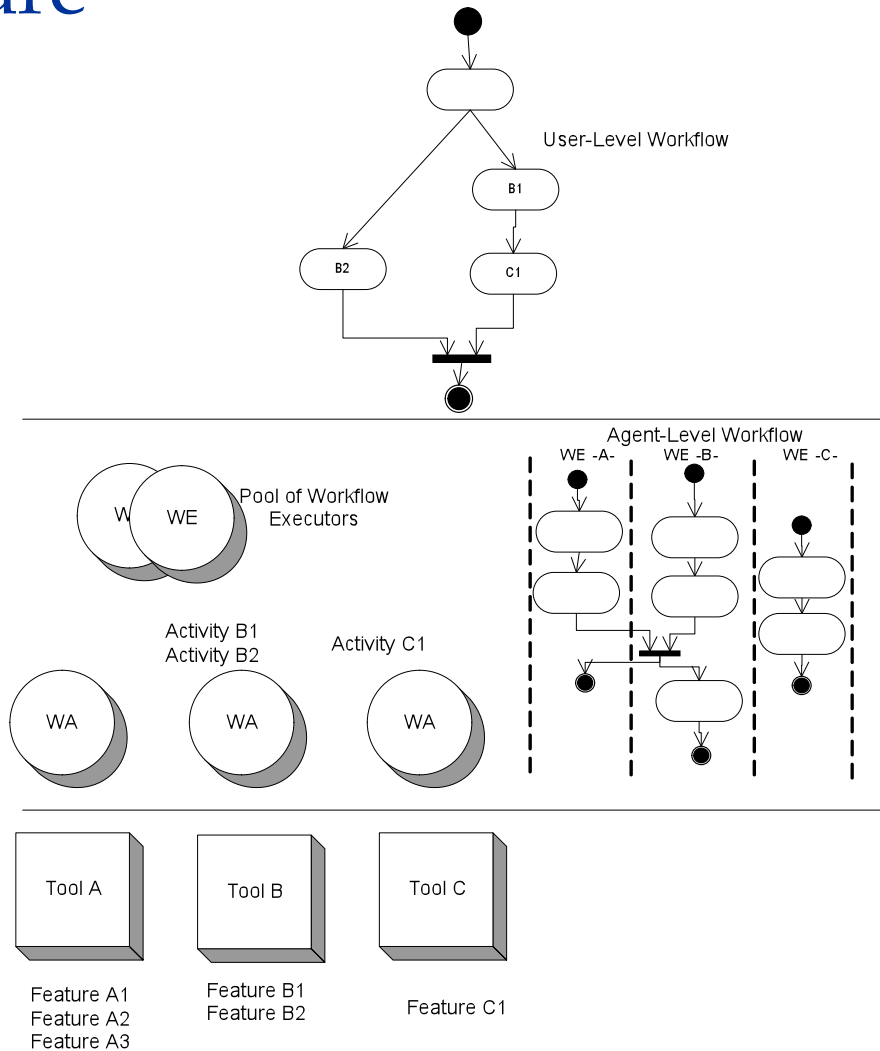


---

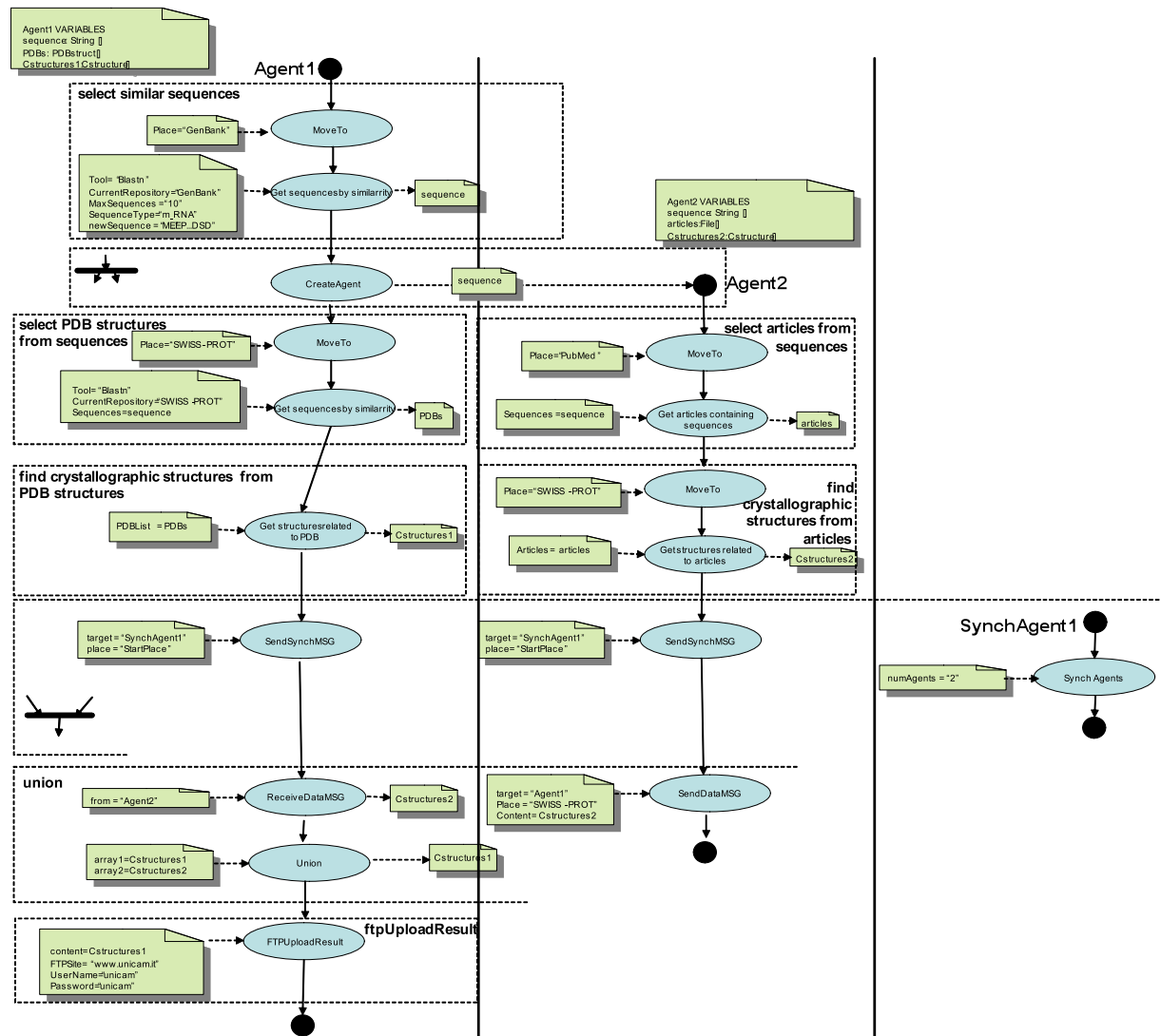
# From Data to Knowledge and vice versa



# The Proposed Approach: an Agent-based Middleware



# Preliminary Results: User-agent as high level *Tool Integration*

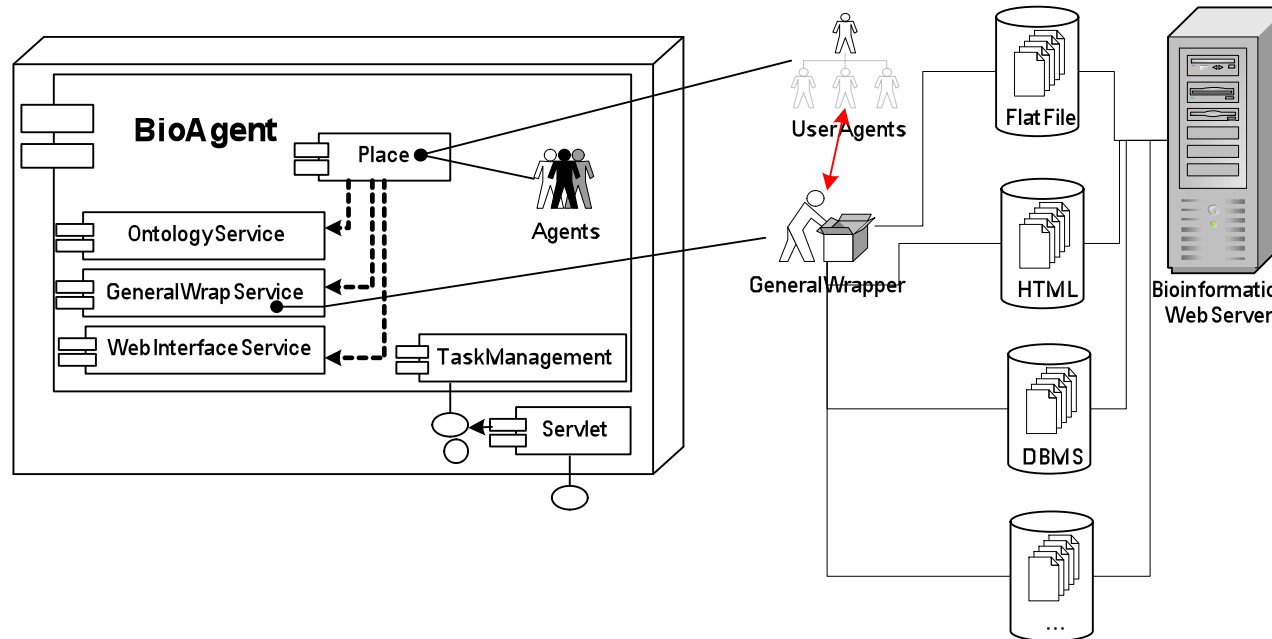


---

# Preliminary Results: Wrapper-agent as low level Tool Integration



# BioAgent



---

## Future Activities and Conclusions

For different application domains (i.e supply chain, components traceability for testing...) we plan to:

- Develop wrapper agents
- Design and develop the knowledge database to manage software tools
- Develop the compiler to allow the automatic generation of user-agents
- Evaluate the possibility to include mobility to user-agents in order to minimize the data transfer during tasks execution.

We conclude saying that software tool integration for real applications, as those in Bioinformatics domain, is a very difficult task due to both heterogeneity of data format and wide variety of tools which continuously evolve.

# NCBI - Home page

NCBI  
National Center for Biotechnology Information  
National Library of Medicine National Institutes of Health

PubMed Entrez BLAST OMIM Books TaxBrowser Structure

Search Nucleotide for MEEEP...DSD Go

**SITE MAP**  
Guide to NCBI resources

**About NCBI**  
The science behind our resources. An introduction for researchers, educators and the public.

**GenBank**  
Sequence submission support and software

**Literature databases**  
PubMed, OMIM, Books and PubMed Central

**Molecular databases**  
Sequences, structures, and taxonomy

**Genomic biology**  
The human genome, whole genomes and related resources

**Tools**  
Data mining

**What does NCBI do?**

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

**PubMed Central**  
An archive of life sciences journals

- Free fulltext
- 80,000 articles from over 100 journals
- Linked to PubMed and fully searchable

Use of PubMed Central requires no registration or fee. Access it from any computer with an Internet connection.

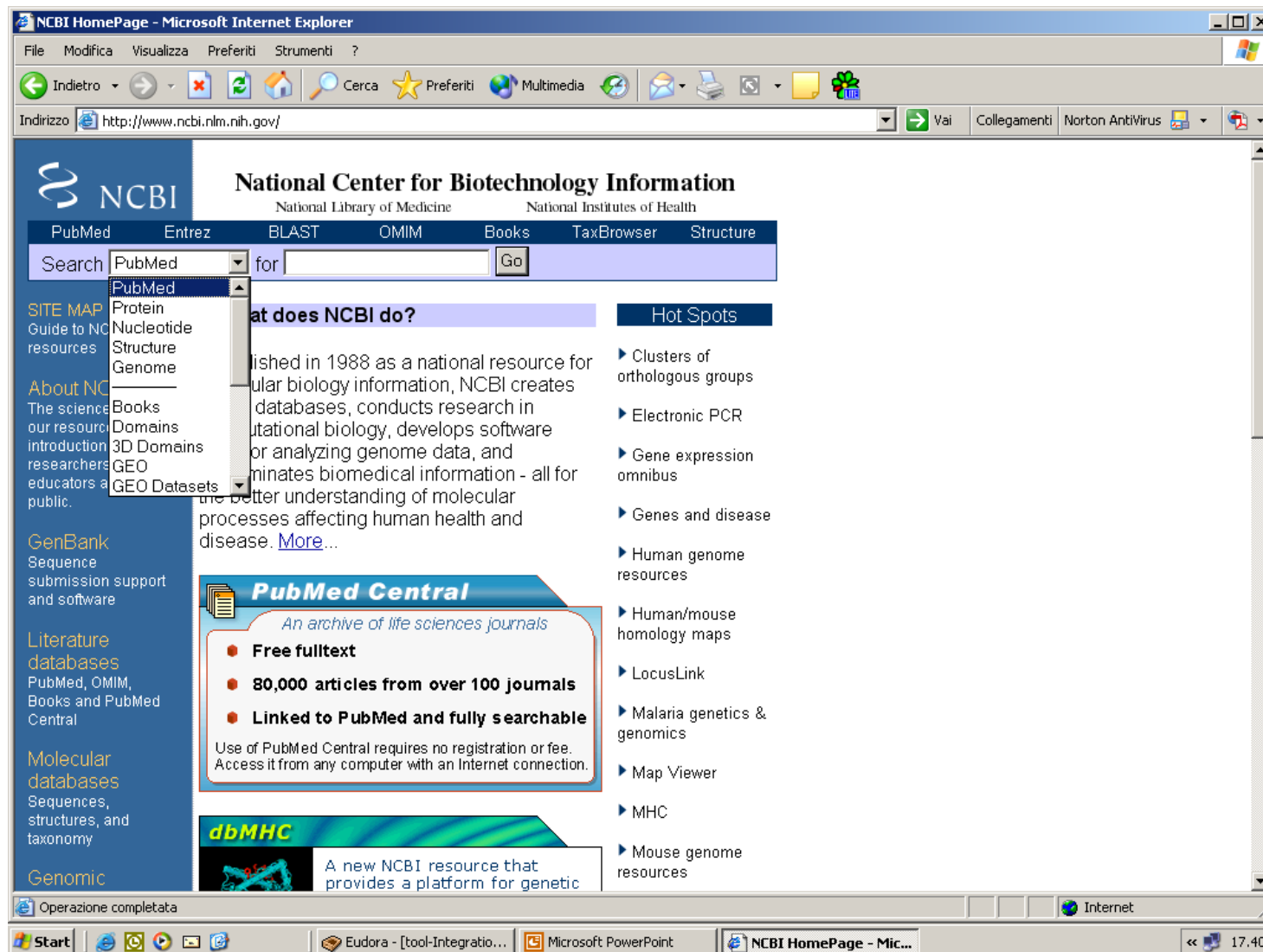
**NCBI Web Site Search**

A function in Entrez is now available allowing one to search the NCBI web site and ftp site. Choose 'NCBI site search' from the Entrez pulldown menu to find information from any area of our web site.

**Hot Spots**

- Clusters of orthologous groups
- Electronic PCR
- Gene expression omnibus
- Genes and disease
- Human genome resources
- Human/mouse homology maps
- LocusLink
- Malaria genetics & genomics
- Map Viewer
- MHC
- Mouse genome resources
- NCBI Handbook
- ORF finder
- Reference sequence project
- Retrovirus

# NCBI - main databses



# NCBI - site map

This site map is also a **guide to NCBI resources**. Each link leads to a **brief description of the resource** on this page, then to the resource itself. A **Quick Links** table is also available. It provides only an alphabetical list of the major resources with **direct links** to those resources, bypassing the descriptions.

## RESOURCE CATEGORIES

### About NCBI

[programs and services](#), [NCBI handbook](#), [what's new](#), [NCBI News](#), [postdoctoral fellowships](#), [organizational structure](#), [contact information](#), [e-mail lists](#), [site search](#)

### GenBank

[overview](#), [submit sequences](#), [submit genomes](#), [sample record](#), [GenBank divisions](#), [statistics](#), [release notes](#), [international collaboration](#), [FTP GenBank](#)

### Molecular Databases

[nucleotides](#), [proteins](#), [structures](#), [taxonomy](#)

### Literature Databases

[PubMed](#), [PubMedCentral](#), [OMIM](#), [Books](#), [Citation Matcher](#)

### Genomes and Maps

[organism collections](#) (including [Entrez Genomes](#), [Map Viewer](#), and [UniGene](#)), [human](#), [mouse](#), [rat](#), [cow](#), [zebrafish](#), [Drosophila](#), [nematode](#), [plant genomes](#), [yeast](#), [malaria](#), [microbial genomes](#), [viruses](#), [viroids](#), [plasmids](#), [eukaryotic organelles](#)

### Tools

[Entrez](#), [LinkOut](#), [Cubby](#), [BLAST](#), [nucleotide sequence analysis](#), [protein sequence analysis](#), [3-D structure display and similarity searching](#)

### Research at NCBI

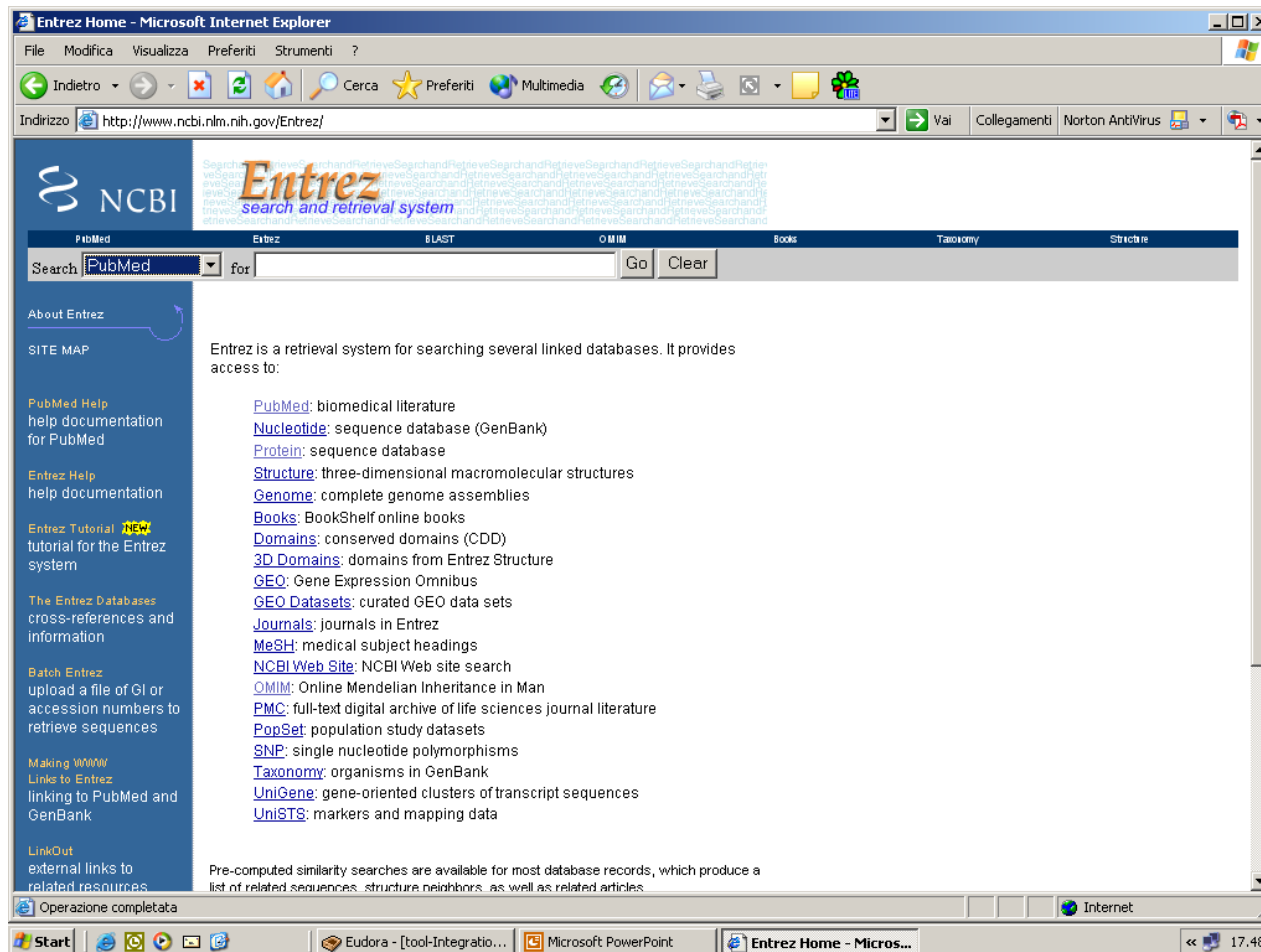
[Computational Biology Branch \(CBB\)](#), [senior](#)

## ALPHABETICAL INDEX

WITH LINKS TO RESOURCE DESCRIPTIONS  
(To bypass descriptions, use the [Quick Links](#) table.)

<a href="#">About NCBI</a>	<a href="#">Genomes and Maps</a>	<a href="#">PubMed</a>
<a href="#">ASN.1</a>	<a href="#">GEO (Expression)</a>	<a href="#">PubMed Central</a>
<a href="#">BankIt</a>	<a href="#">Glossaries</a>	<a href="#">RefSeq</a>
<a href="#">BLAST</a>	<a href="#">HTGs</a>	<a href="#">Research at NCBI</a>
<a href="#">Books</a>	<a href="#">HomoloGene</a>	<a href="#">Retroviruses</a>
<a href="#">CDART</a>	<a href="#">Human Genome Resources</a>	<a href="#">SAGEmap</a>
<a href="#">CDD</a>	<a href="#">Human-Mouse Homology Maps</a>	<a href="#">Science Primer</a>
<a href="#">CGAP</a>	<a href="#">LinkOut</a>	<a href="#">Seminars</a>
<a href="#">Clones</a>	<a href="#">LocusLink</a>	<a href="#">Sequin</a>
<a href="#">Cn3D</a>	<a href="#">Malaria</a>	<a href="#">Site Search</a> <b>NEW</b>
<a href="#">Coffee Break</a>	<a href="#">Map Viewer</a>	<a href="#">SKY/CGH</a>
<a href="#">COGs</a>	<a href="#">MGC</a>	<a href="#">Software Engineering</a>
<a href="#">Computational Biology Branch</a>	<a href="#">Microbial Genomes</a>	<a href="#">Spidey</a>
<a href="#">dbEST</a>	<a href="#">MMDB</a>	<a href="#">Structures</a>
<a href="#">dbGSS</a>	<a href="#">Model Maker</a> <b>NEW</b>	<a href="#">Submit Data</a>
<a href="#">dbSNP</a>	<a href="#">Mutation Databases</a>	<a href="#">Taxonomy</a>

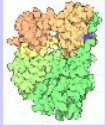
# NCBI - Entrez



# PDB -Home page

**DEPOSIT data**  
**DOWNLOAD files**  
**browse LINKS**  
**BETA TEST new features**  
**BETA mmCIF and XML files**

**Current Holdings**  
**21838 Structures**  
**Last Update: 22-Jul-2003**  
**PDB Statistics**



**Molecule of the Month:**  
**Src Tyrosine Kinase**

The Protein Data Bank (PDB) is operated by Rutgers, The State University of New Jersey; the San Diego Supercomputer Center at the University of California, San Diego; and the Center for Advanced Research in Biotechnology of the National Institute of Standards and Technology -- three members of the [Research Collaboratory for Structural Bioinformatics \(RCSB\)](#). The PDB is supported by funds from the [National Science Foundation](#), the [Department of Energy](#), and two units of the National Institutes of Health: the [National Institute of General Medical Sciences](#) and the [National Library of Medicine](#).


**PROTEIN DATA BANK**

Welcome to the PDB, the single worldwide repository for the processing and distribution of 3-D biological macromolecular structure data.

[RCSB Home](#) [Contact Us](#) [Help](#)

**Did you find what you wanted?**

[ABOUT PDB](#) | [DATA UNIFORMITY](#) | [RECENT FEATURES](#) | [USER GUIDES](#) | [FILE FORMATS](#) | [EDUCATION](#) | [STRUCTURAL GENOMICS](#) | [PUBLICATIONS](#) | [SOFTWARE](#)

**Search the Archive** 

Enter a **PDB ID** or keyword [Query Tutorial](#)

1T5R

query by PDB id only  match exact word  
 [remove sequence homologs](#)

[SearchLite](#) keyword search form with examples  
[SearchFields](#) customizable search form  
[Status Search](#) find entries awaiting release

**News** [Complete News Newsletter](#) [pdb-I Archive Subscribe](#)

**22-Jul-2003**  
**Demonstrations, Posters, and More: PDB at the ACA Annual Meeting and the 17th Symposium of the Protein Society**  
The PDB would like to thank those ISMB 2003 attendees who provided valuable feedback at our demonstration session and exhibit during this worthwhile event. PDB staff members will also participate in several other meetings in the near future, including the annual meetings of the Protein Society and the American Crystallographic Association... [\[MORE...\]](#)


**PDB Mirrors**

*\*\*Please bookmark a mirror site\*\**  
[San Diego Supercomputer Center\\*](#)  
[Rutgers University\\*](#)  
[National Institute of Standards and Technology\\*](#)  
[Cambridge Crystallographic Data Centre, UK](#)  
[National University of Singapore](#)  
[Osaka University, Japan](#)  
[Universidade Federal de Minas Gerais, Brazil](#)  
[Max Delbrück Center for Molecular Medicine, Germany](#)

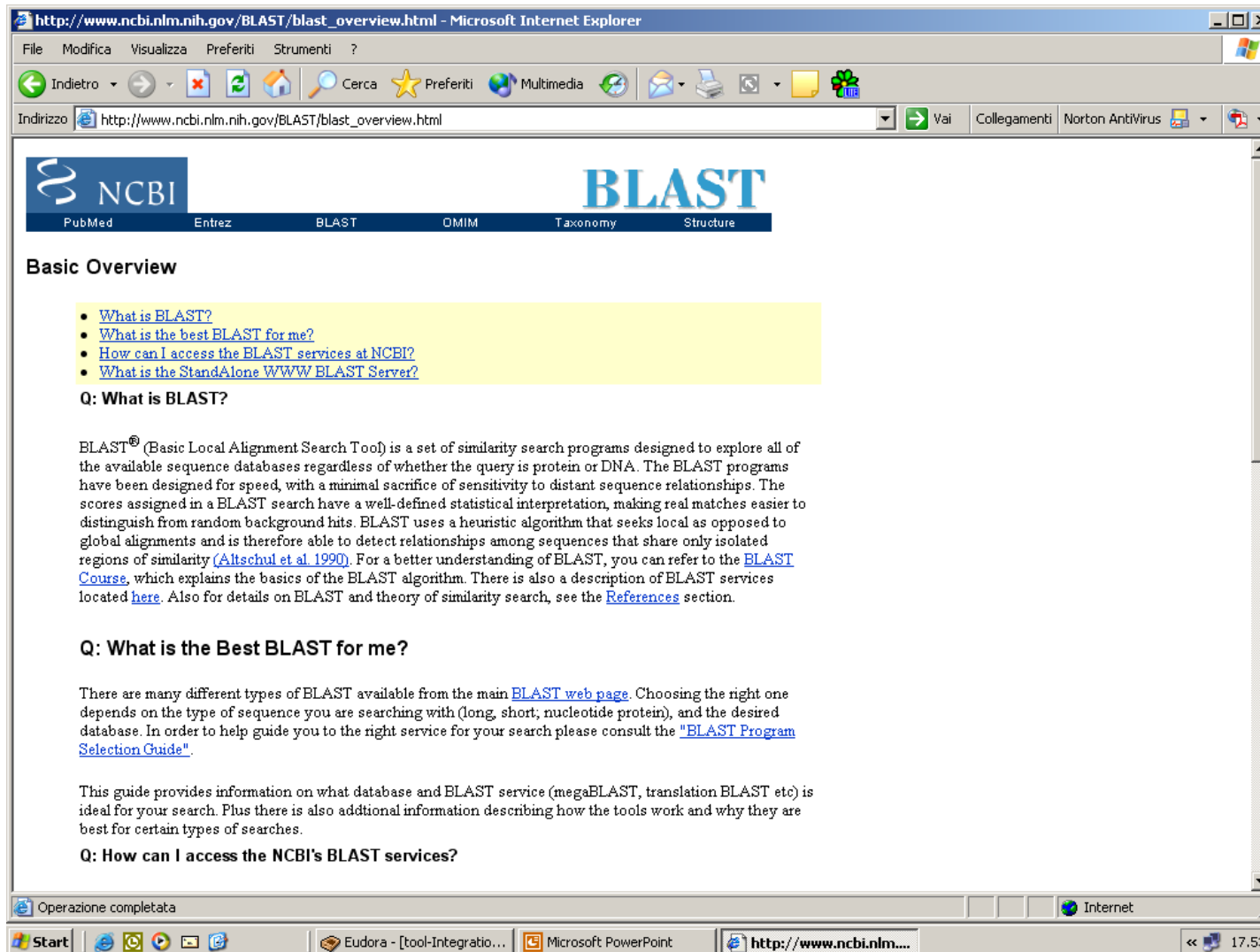
**OTHER SITES** \*RCSB partner

In citing the PDB please refer to:  
H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N. Shindyalov, P.E. Bourne: [The Protein Data Bank](#), *Nucleic Acids Research*, **28** pp. 235-242 (2000)

[ABOUT PDB](#) | [DATA UNIFORMITY](#) | [RECENT FEATURES](#) | [USER GUIDES](#) | [FILE FORMATS](#) | [EDUCATION](#) | [STRUCTURAL GENOMICS](#) | [PUBLICATIONS](#) | [SOFTWARE](#)

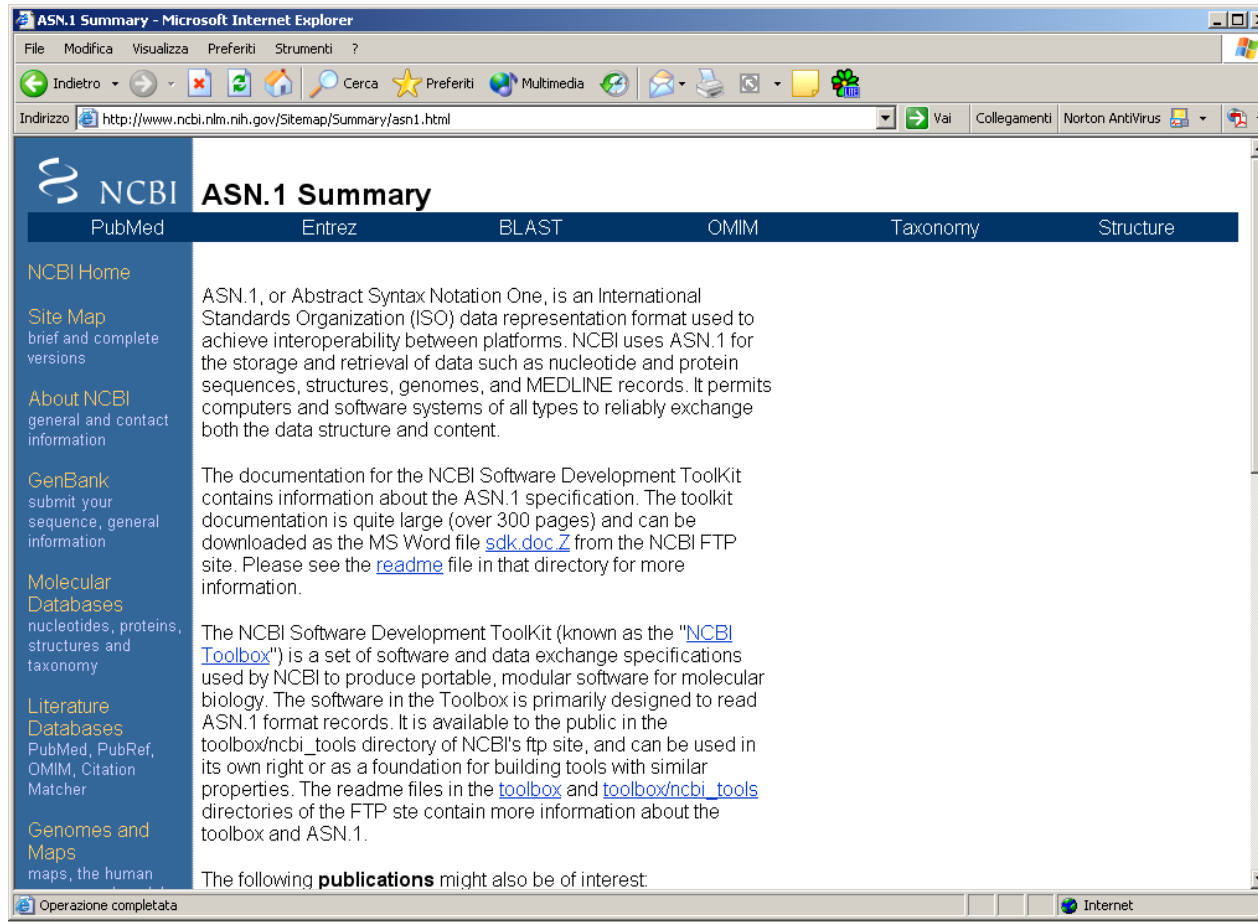


# NCBI - BLAST





# NCBI - ASN.1



# NCBI - fomats

NCBI Sequence Viewer - Microsoft Internet Explorer

File Modifica Visualizza Preferiti Strumenti ?

Indirizzo <http://www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?val=4680228&db=Nucleotide&dopt=GenBank>

NCBI

Search Nucleotide for [ ] Go Clear

Limits Preview/Index History Clipboard Details

Display default Show: 1 Send to File Get Subsequence Features

1: AF1182

Summary

ASN.1

FASTA

TinySeq XML

GenBank 1748 bp mRNA linear PRI 17-APR-2000

DEFINITION Nb-5 mRNA, partial cds.

ACCESSION

GI List

VERSION 4680228

KEYWORDS Graphics

SOURCE XML (human)

ORGANISM default

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 1748)

AUTHORS Amler, L.C., Bauer, A., Corvi, R., Dihlmann, S., Praml, C., Cavenee, W.K., Schwab, M. and Hampton, G.M.

TITLE Identification and characterization of novel genes located at the t(1;15)(p36.2;q24) translocation breakpoint in the neuroblastoma cell line NGP

JOURNAL Genomics 64 (2), 195-202 (2000)

MEDLINE 20195630

PUBMED 10729226

REFERENCE 2 (bases 1 to 1748)

AUTHORS Amler, L.C. and Hampton, G.M.

TITLE Direct Submission

JOURNAL Submitted (06-JAN-1999) Genos Biosciences, 11099 North Torrey Pines Road, La Jolla, CA 92037, USA

FEATURES

Location/Qualifiers

source 1..1748

/organism="Homo sapiens"

/mol\_type="mRNA"

Operazione completata

Start Eudora - [tool-Integratio... Microsoft PowerPoint NCBI Sequence Viewe... 17.39

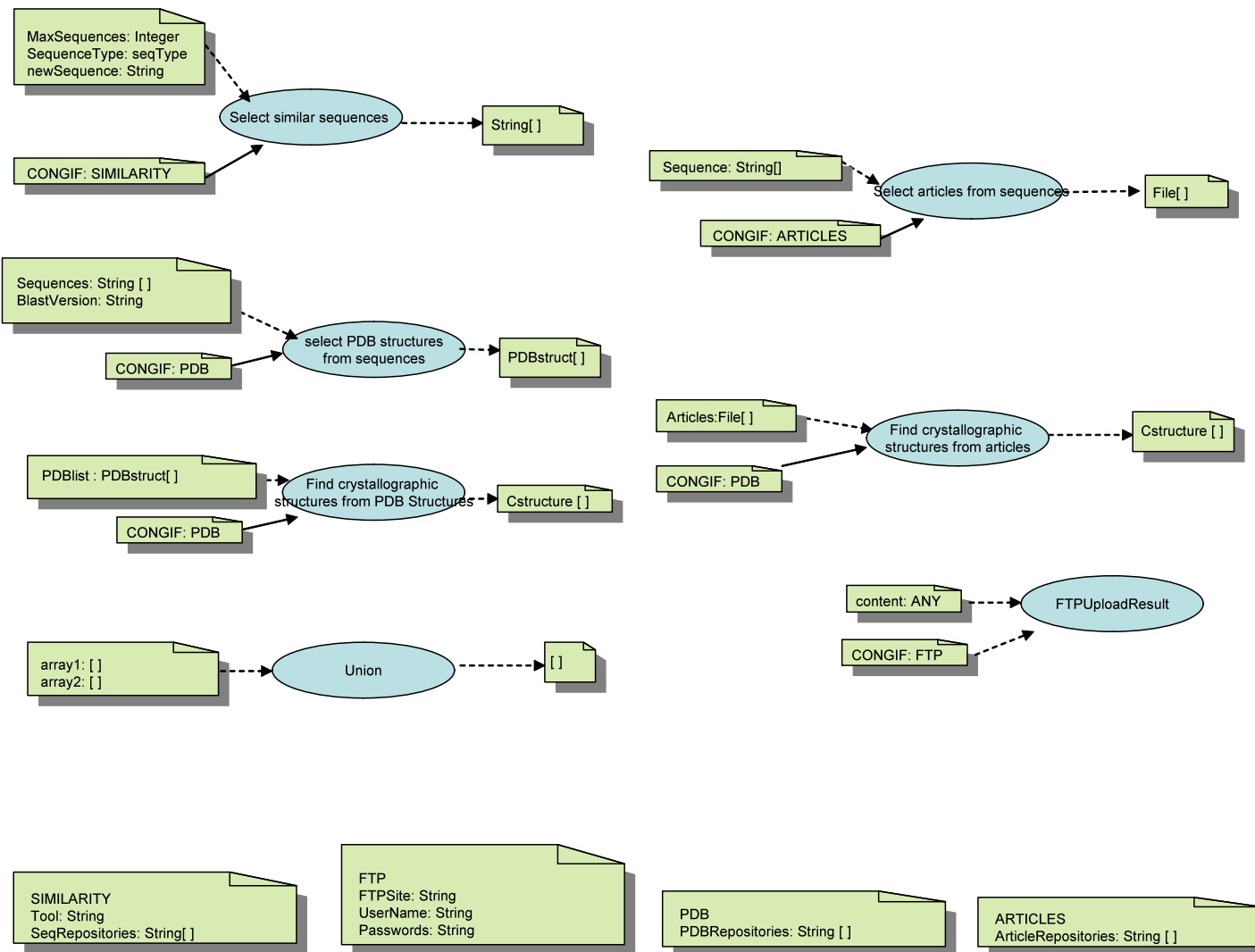
- 
- PDB-ID (P53) = 1TSR
  - [www.rcsb.org](http://www.rcsb.org) (

# DNA and nucleotide sequence

atggaggagccg cagtcagatcctagcgtcgagccccctctgagtcaggaaacattttca  
AtGgAggAgcSg cagtcagatcctagcgtcgagccccctctgagtcaggaaacattttca  
yaEclA P g a s a d t a c s c v t g a a a c a c a c g t e t r t e s s c c t t g c c g t c c c a a g c a a t g  
DalcMalgalaPcFAnNcYgaaAaPaaPcStQAdMcccccttgccgtcccaagcaatg  
DalgaAttgattgPgtEcncyyaccSattgaaScAgyMtcactgaagaccaggtcca  
D.D L M L S P D D I E Q W F T E D P G P  
gatgaagctcccagaatgccagaggctgctccccgcgtggccccctggaccagcagctcct  
DcAggPabDacPaaAgagDadVgAcPaaPaaAcPcagctcctctccccagccaaagaag  
PccSgtgBgBaclPcNcTgScScStEdgPckdktgtcatcttctgtcccttcccag  
aPaaAdtgatgAgaaWattcascStcVgAtScQtggggcgtgagcgcttcgagatg  
kaacDtaCcaYgJdagcaBgGtEcPgtPgJcMttgcattctgggacagccaag  
KctYagGcSgYagFgBdctGgFalcAagJagKccaggctgggaaggagccaggg  
fctgtactNcaqtEcLccDgCkaGaagagatgtttgccaaactggccaagacc  
SgyAgCaYgStPaAtcNgcMctGaaJtaKaAaagggtcagctacctcccgccat  
tGcPgaGdagSgtgggttsaktkaCaacScDcSrdggcaccgcgtccgcgccatg  
GaaAaaclatWtdaagTcaPaaggtRgaRcaGActga  
KckatDtaFakgtAgGacPdgSaDdatgacggagggtgtgaggcgctgccccaccatgag

ttgcgtgtggaglatllggatgacagaaacacilllcgacatagltgtggltgggtccctat  
L R V E Y I D D R N T E R H S V V V P Y

# ULAD



---

# Significant References

Y. Papakonstantinou, H. Garcia-Molina & J. Widom '95

OEM: Object Exchange Across Heterogeneous Information Sources

S. Bergamaschi ... '00

Momis: Mediator envirOnment for Multiple Information Sources

G. Cabri, L. Leonardi & F. Zambonelli '00

MARS: A Programmable coordination Architecture for Mobile Agents

...

E. Bartocci, L. Mariani & E. Merelli '03

...  
AIXO: Any Input XML Output, a generalized wrapper

F. Corradini, L. Mariani & E. Merelli '03

PEGAA: A Programming Environment for Global Activity-based Applications