

# DATABASE INTEGRATION

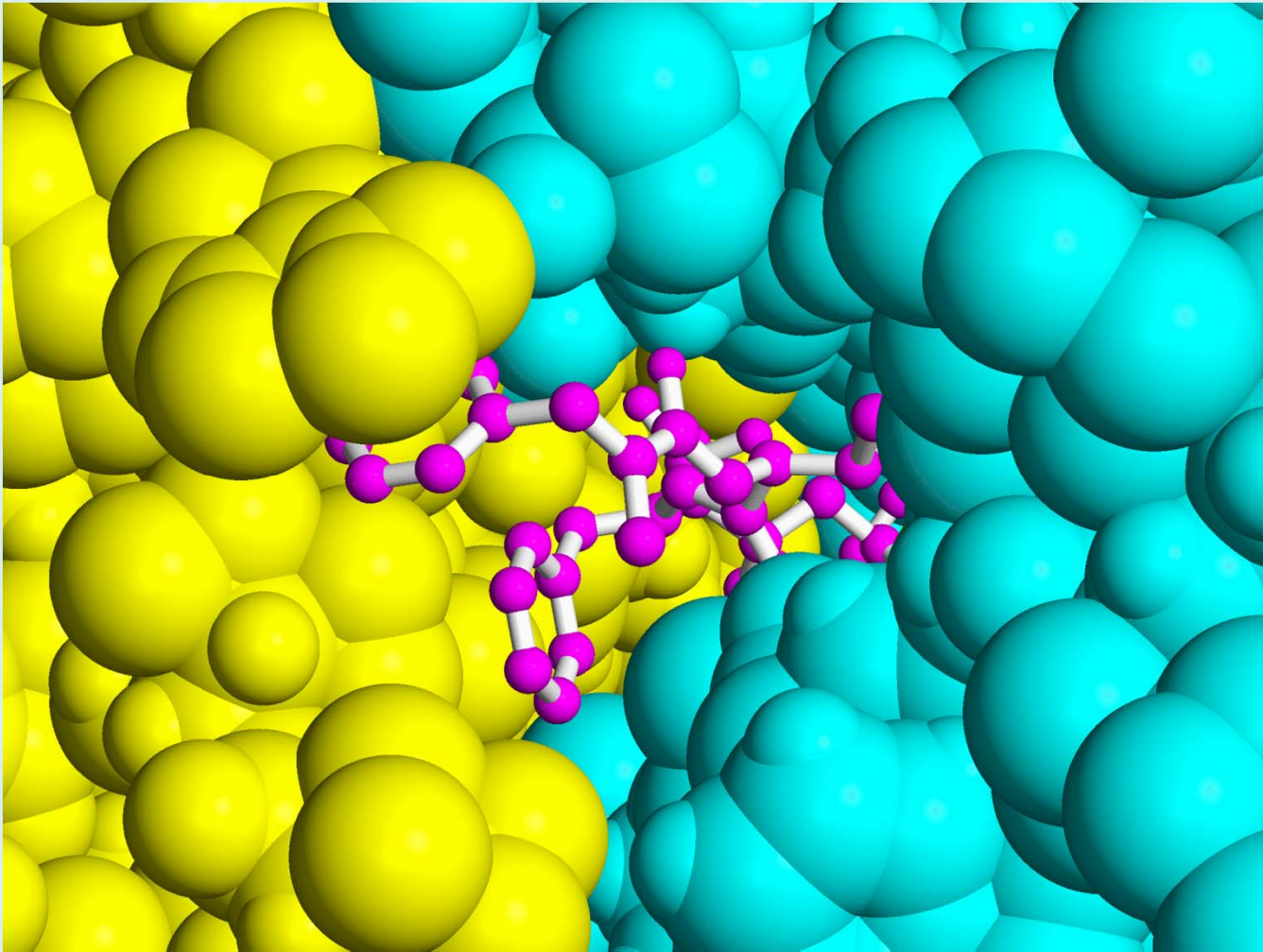
# Genes.....

cagaaagaaatcagaa  
aggggcttggaag  
cttttgctataag  
cgggtggcagtg  
ccaacgtagta  
nratnratnratct<sup>2</sup>

.....make proteins

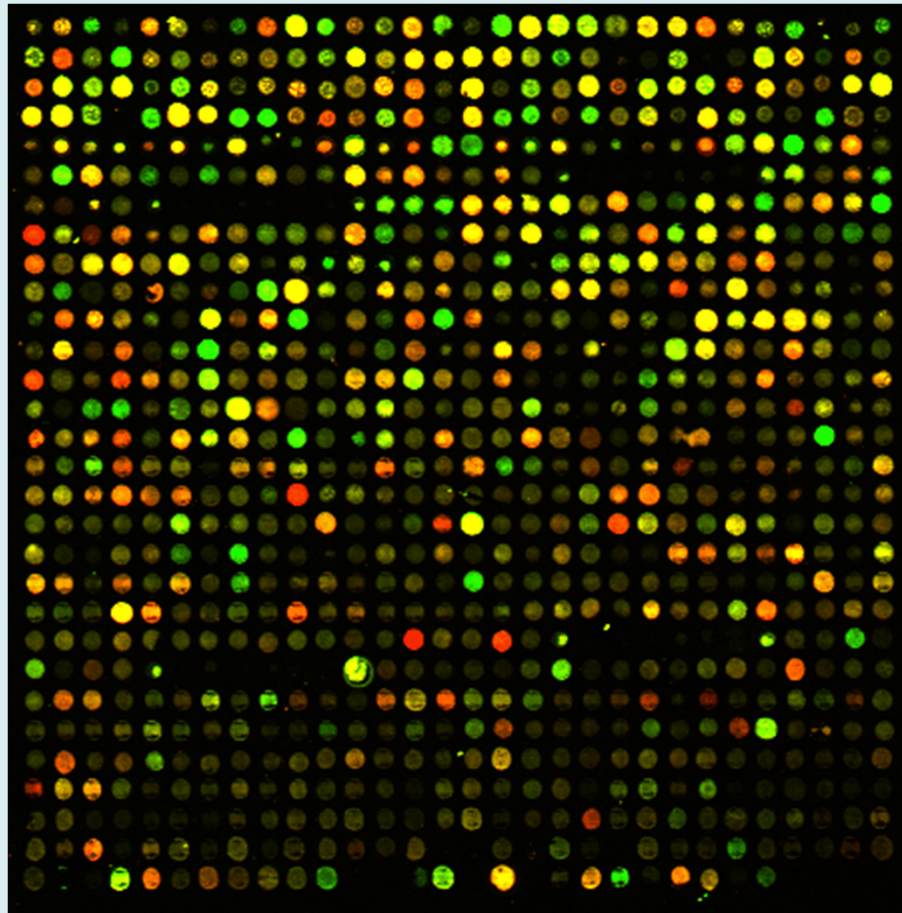
FREDLALLOGR  
FSSSEQTRANSH  
RELQVWGGENNS  
AGADRQGTVSH  
QITLWQRPLVT  
GGQLKEALLDT  
DVT FEMNT D

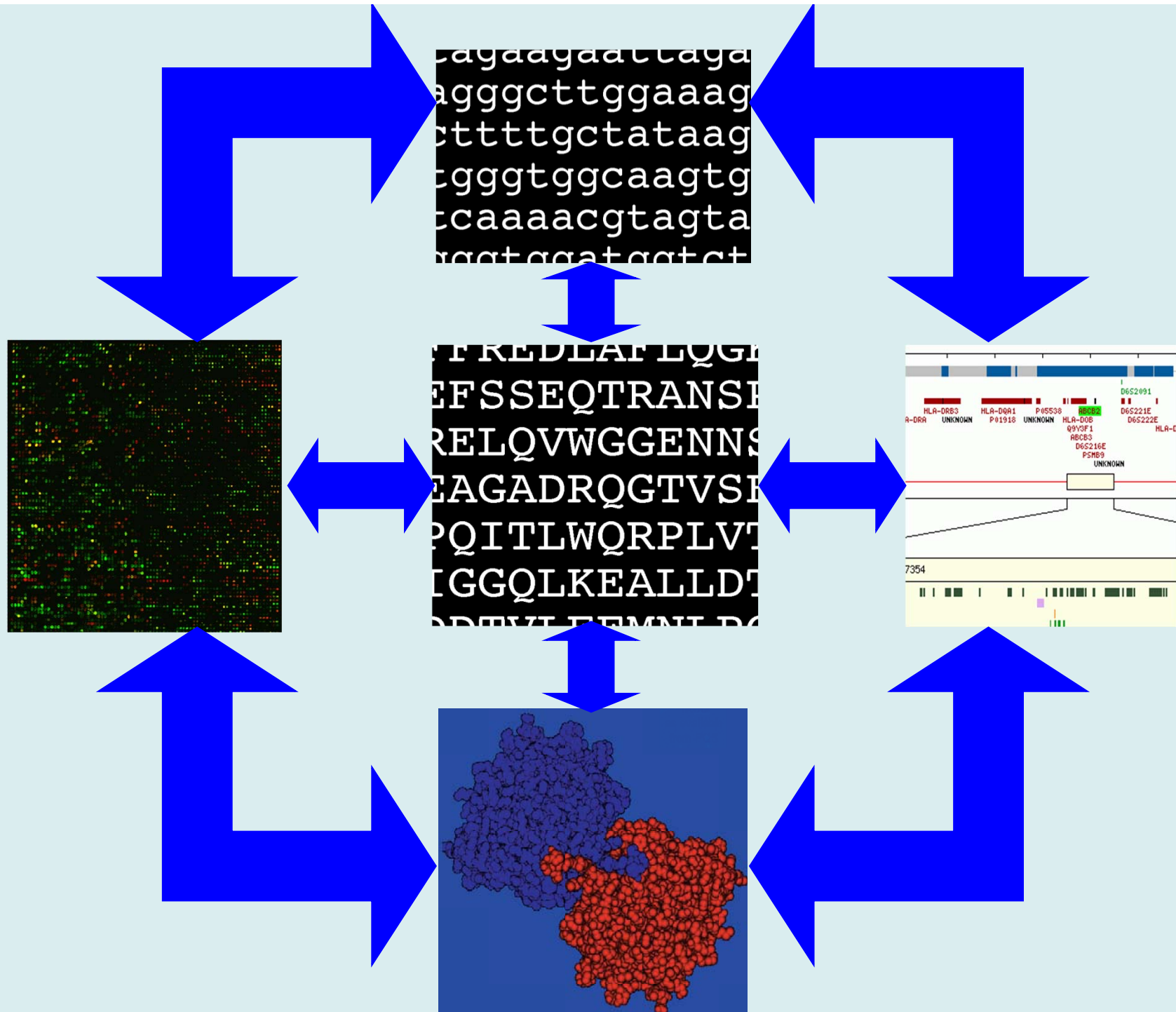
Proteins form complex 3D structures...



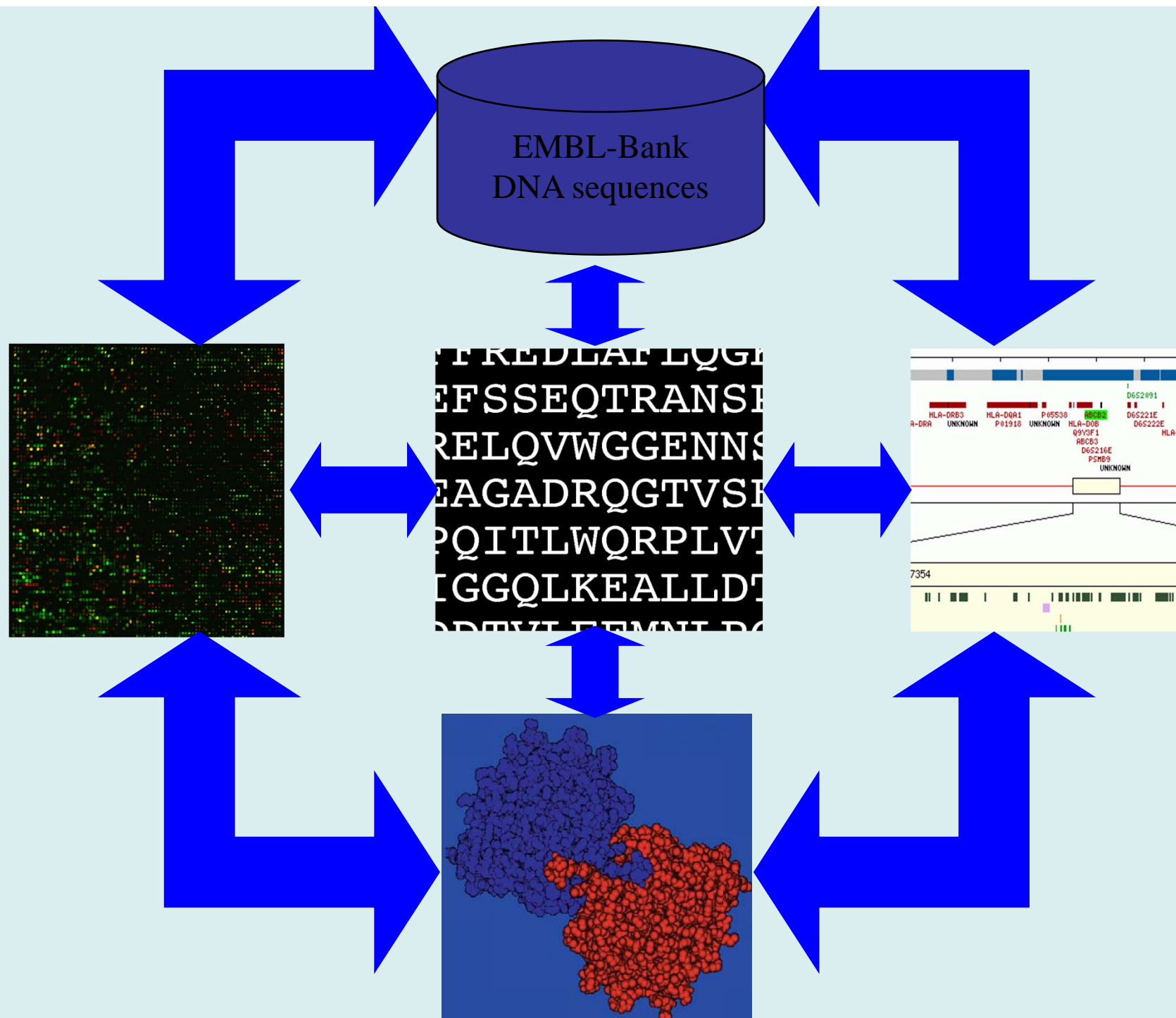
**... and molecules interact**

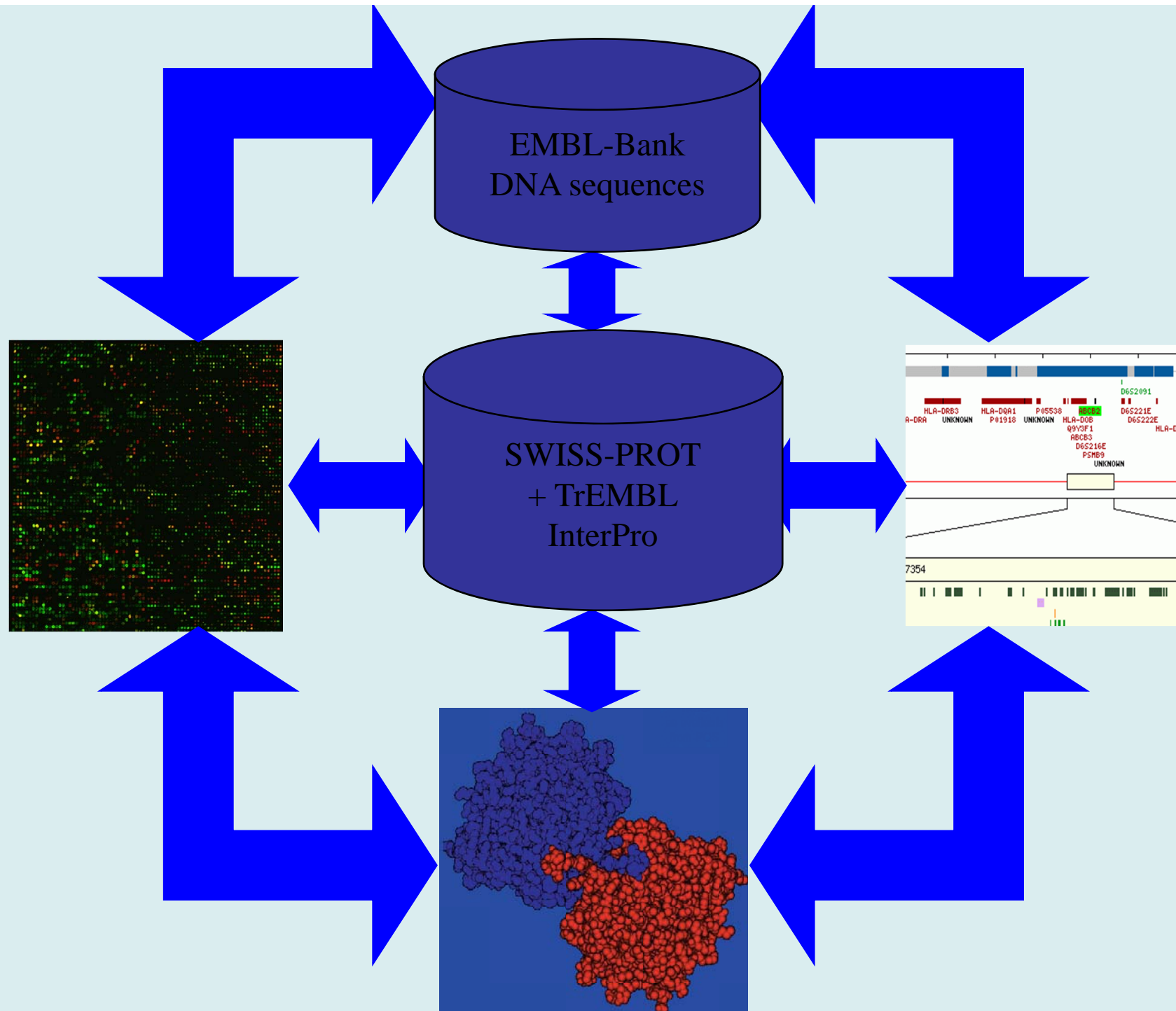
the right molecules need to be  
present at the right time



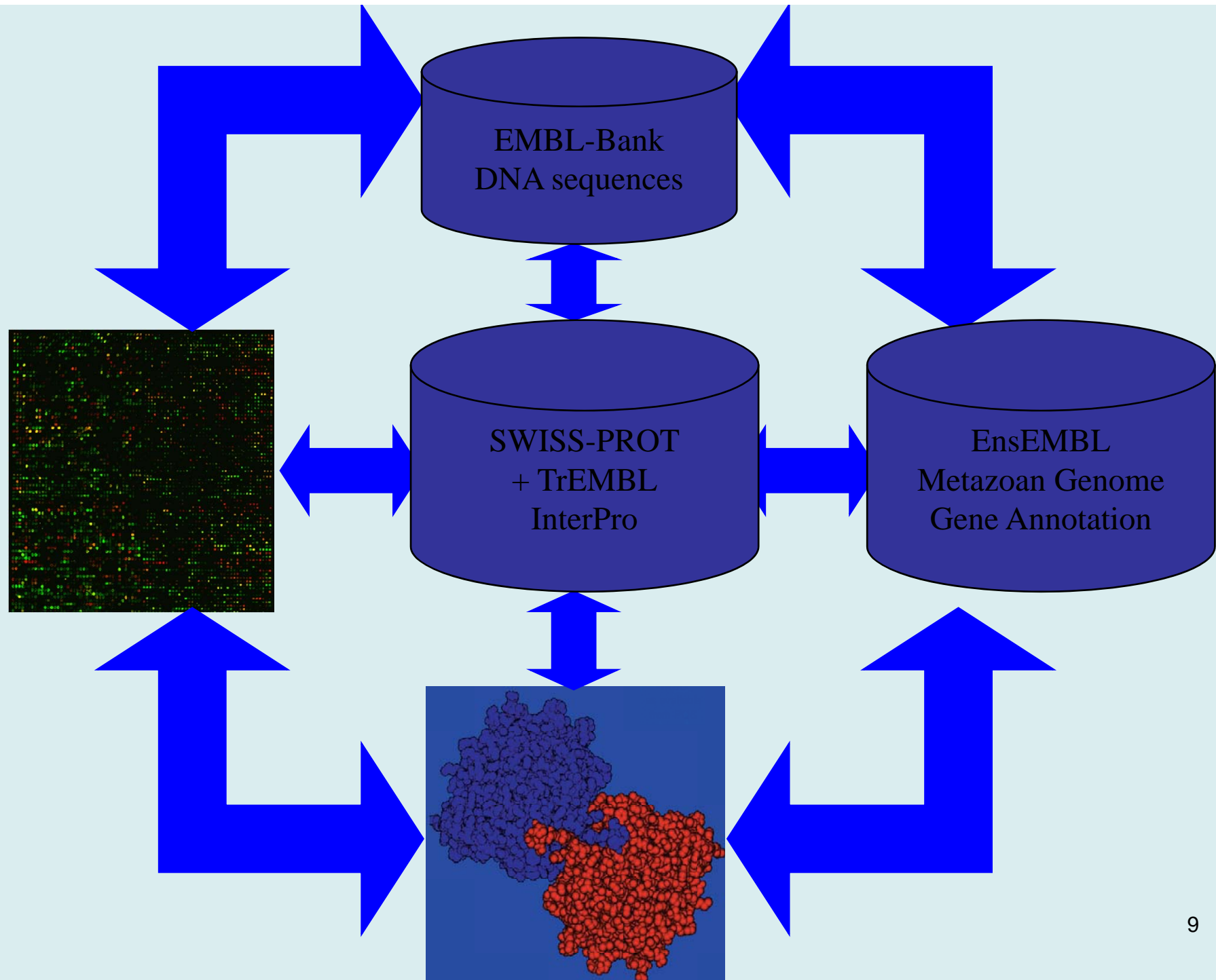


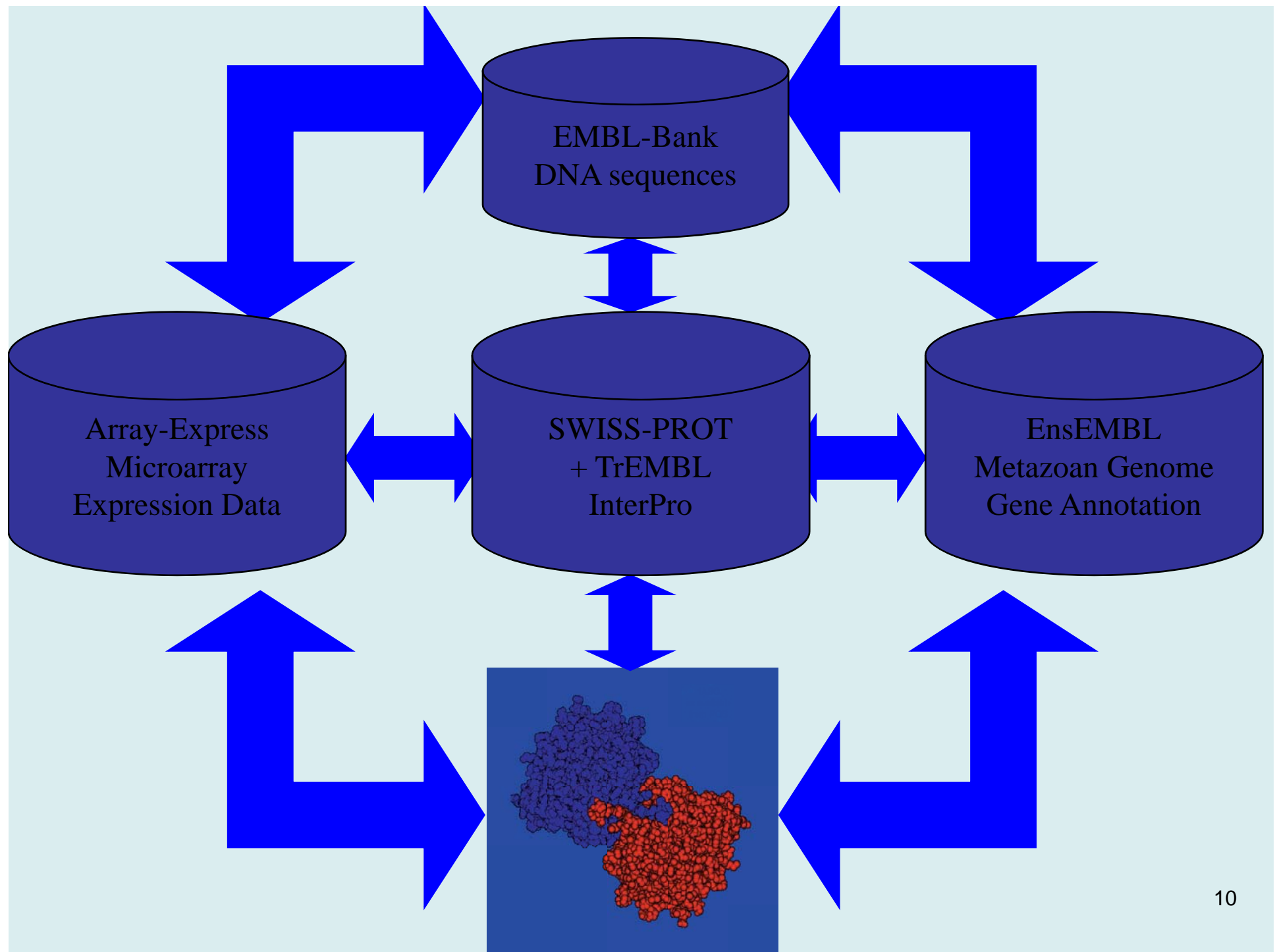


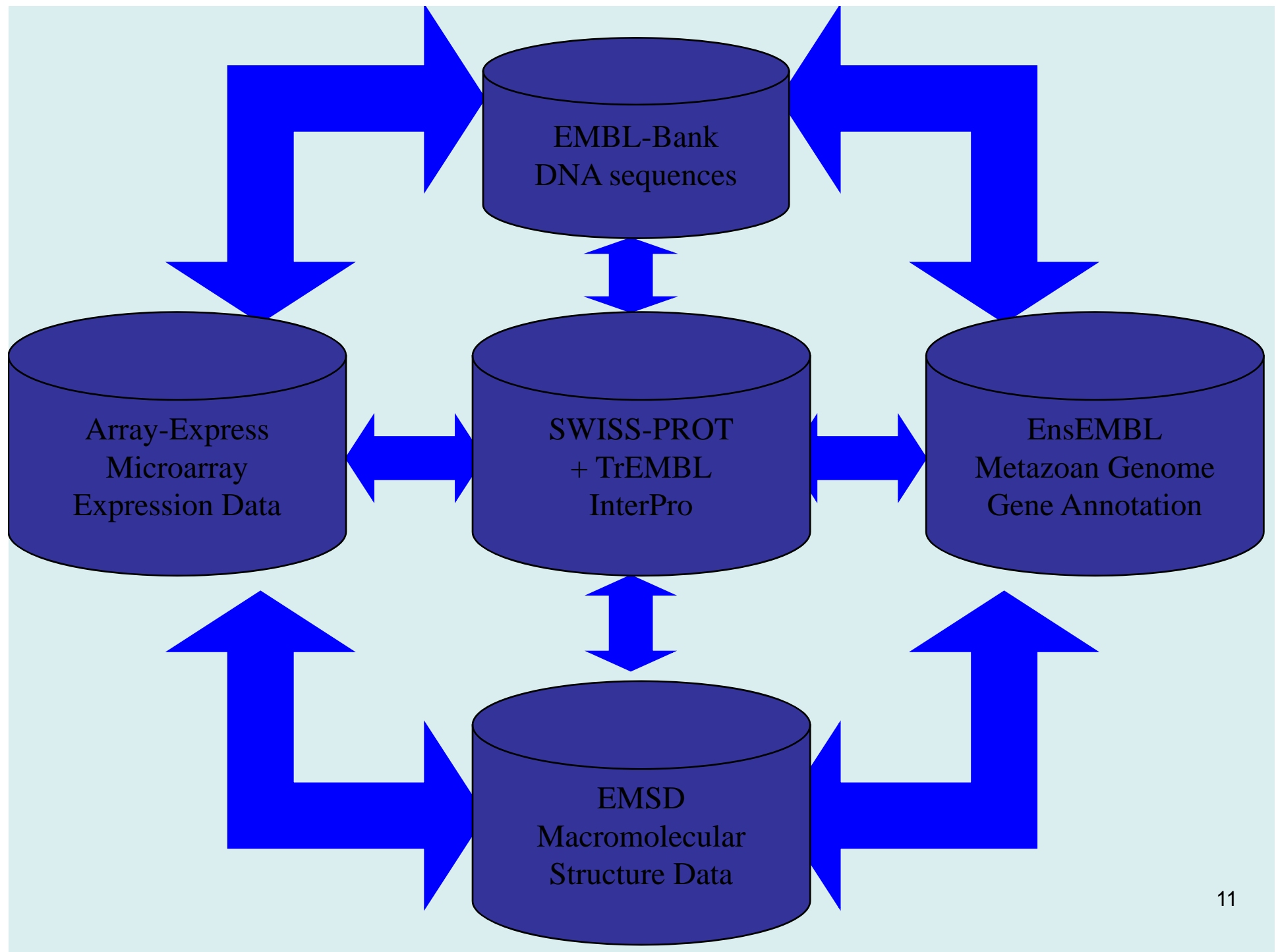


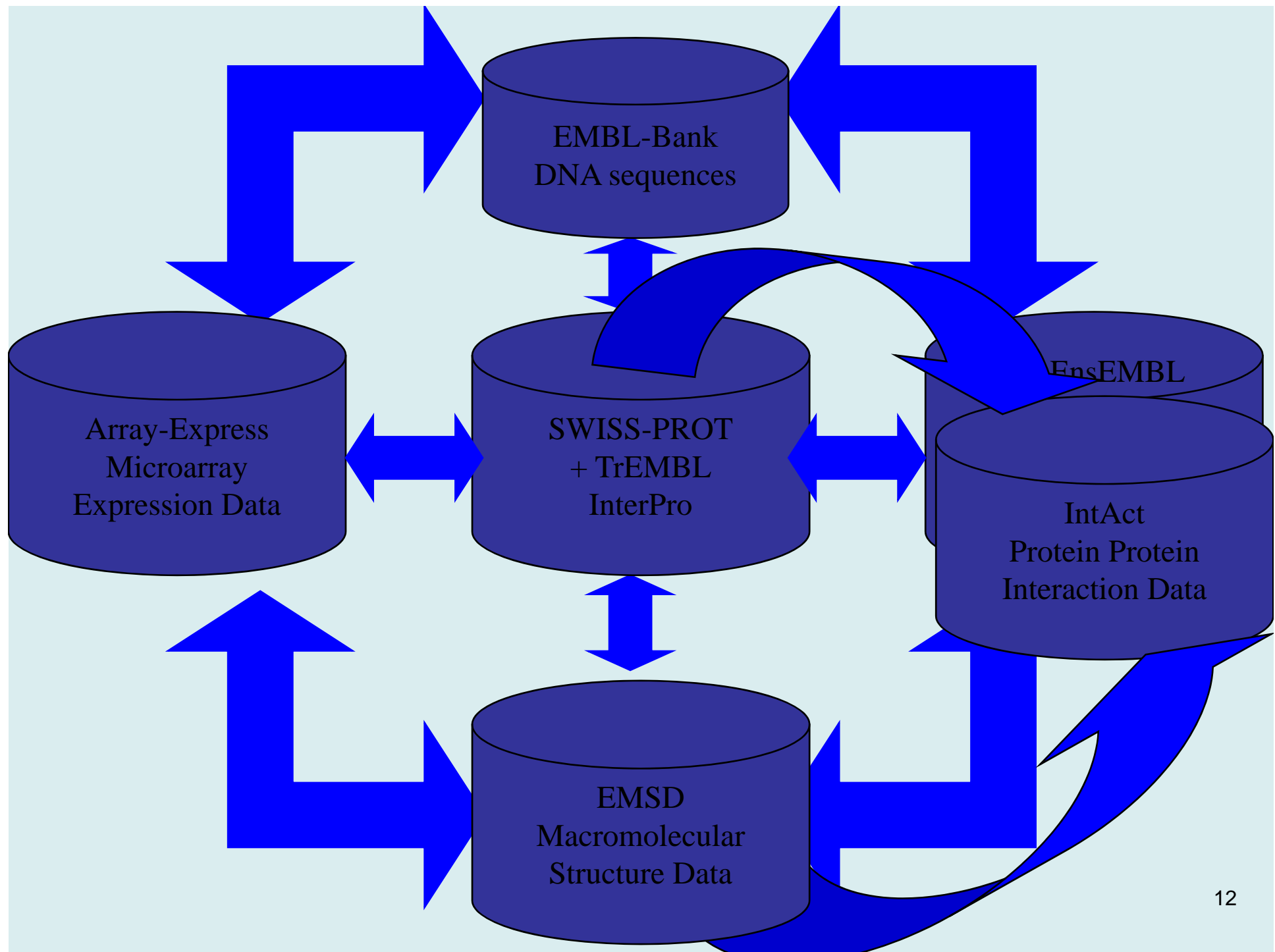












## EMBnet

Developed in the 80s to link European laboratories which used bioinformatics as a research tool in Biology.

Takes advantage of the Internet potential for global communication and resources centralization.



# EMBnet Members

Argentina  
Australia  
Austria  
Belgium  
Brasil  
Canada  
Chile  
China  
Colombia  
Cuba  
Denmark  
Finland  
France  
Germany

Greece  
Hungary  
India  
Ireland  
Israel  
Italy  
Mexico  
Netherlands  
Norway  
Poland  
Portugal  
Russia  
Slovakia  
South Africa

Spain  
Sweden  
Switzerland  
UK  
  
EBI  
ETI  
ICGEB  
MIPS  
UMBER  
Hoffman-La Roche  
LION Bioscience




## The National Nodes

<a href="#"><u>Argentina</u></a>	<a href="#"><u>Node info</u></a>		
<a href="#"><u>Australia</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Israel</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Austria</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Italy</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Belgium</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Mexico</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Brazil</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Netherlands (The)</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Canada</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Norway</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Chile</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Poland</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>China</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Portugal</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Colombia</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Russia</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Cuba</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Slovakia</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Finland</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>South Africa</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>France</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Spain</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Germany</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Sweden</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>Hungary</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>Switzerland</u></a>	<a href="#"><u>Node info</u></a>
<a href="#"><u>India</u></a>	<a href="#"><u>Node info</u></a>	<a href="#"><u>United Kingdom</u></a>	<a href="#"><u>Node info</u></a>

# EBI


Database [Search for](#)  in [Nucleotide sequences](#) [Go](#) [Search EBI Website](#)  [Go](#)




European Bioinformatics Institute

Groups at the European Bioinformatics Institute

Bioinformatics Products and Services

 [FLASH Intro](#)



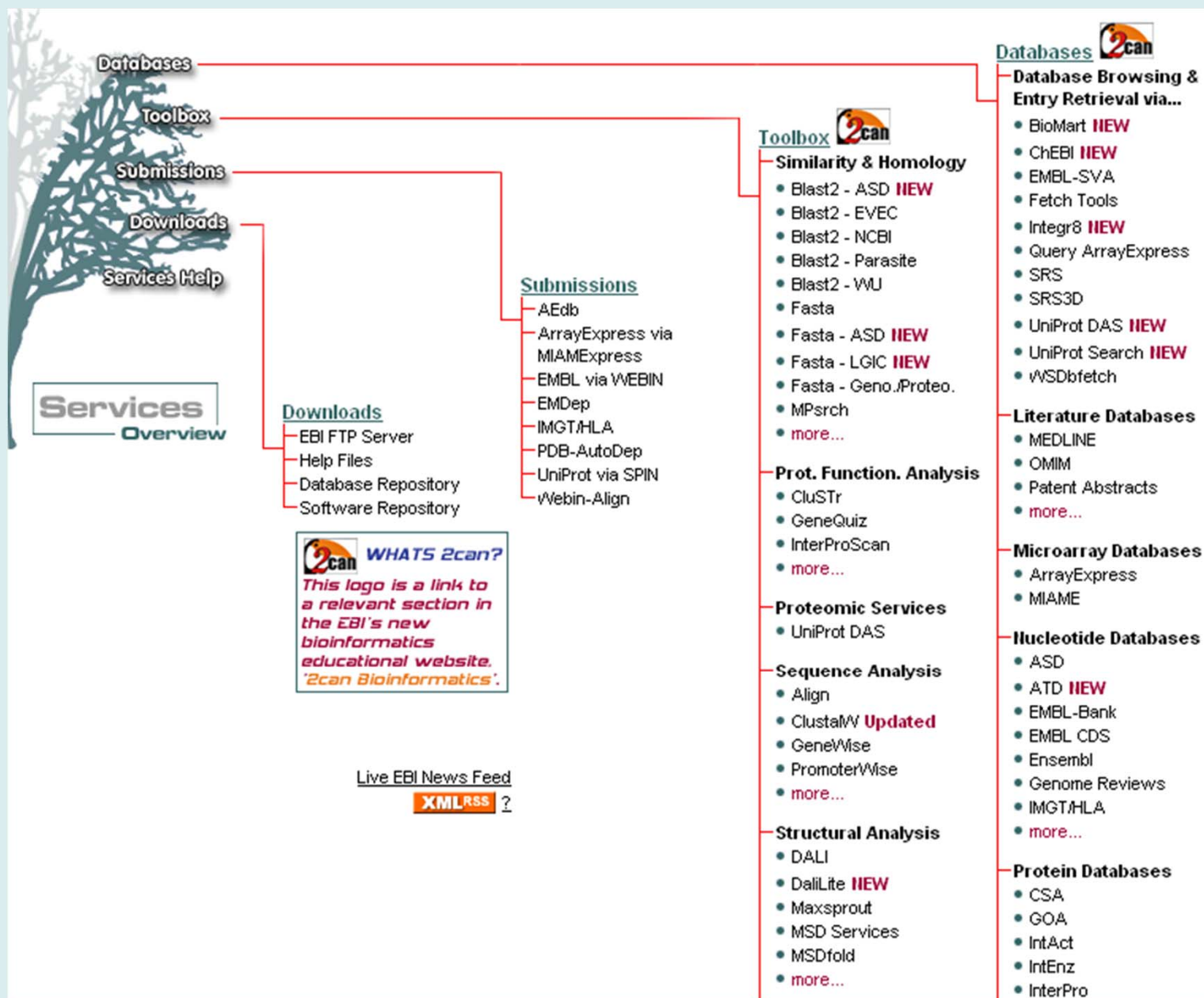
European Bioinformatics Institute  
a part of the European Molecular Biology Laboratory

## the path to knowledge

[About EBI](#) | [Funding](#) | [Whats New](#) | [Groups](#) | [Services](#) | [Toolbox](#) | [EBI Databases](#) | [Downloads](#) | [Submissions](#)  
[BioMart Database Queries](#) | [SRS Database Queries](#) | [Site Search](#) | [Site Map](#) | [Services Map](#) | [Contact Us](#) | [Terms of Use](#)

 [2can](#) [2can Bioinformatics - Training and Education at the EBI](#) [XMLRSS](#) ?

© Copyright European Bioinformatics Institute 2002-2004. All Rights and Trademarks Reserved.





## National Center for Biotechnology Information

National Library of Medicine      National Institutes of Health

[PubMed](#)
[All Databases](#)
[BLAST](#)
[OMIM](#)
[Books](#)
[TaxBrowser](#)
[Structure](#)

Search  for

### SITE MAP

Alphabetical List  
Resource Guide

### About NCBI

An introduction to NCBI

### 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

### Research at NCBI

### 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...](#)



The new My NCBI has replaced the Cubby and includes automatic e-mailing of search updates and filtering search results. A tab format is used for features such as Limits and displaying filtered search results.

### Entrez Gene

You can now use Entrez to search for information centered on the concept of a gene, and connect to many sources of related information both within and outside NCBI.



### PubMed Central

*An archive of life sciences journals*

- Free fulltext
- Over 300,000 articles from over 150 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.

### Hot Spots

- Assembly Archive
- Clusters of orthologous groups
- Coffee Break, Genes & Disease, NCBI Handbook
- Electronic PCR
- Entrez Home
- Entrez Tools
- Gene expression omnibus (GEO)
- Human genome resources
- Malaria genetics & genomics
- Map Viewer
- dbMHC
- Mouse genome resources
- My NCBI
- ORF finder

All Databases
PubMed
Nucleotide
Protein
Genome

Search PubMed for

PubMed
Protein
Nucleotide
Structure
Genome
Books
CancerChromosomes
Conserved Domains
3D Domains
Gene
Genome Project
GENSAT
GEO Profiles
GEO DataSets
HomoloGene
Journals
MeSH
NCBI Web Site
NLM Catalog
OMIM
PMC
PopSet
Probe
PubChem BioAssay
PubChem Compound
PubChem Substance
SNP
Taxonomy
UniGene
UniSTS

Preview/Index
History
Clipboard
Details

Enter one or more search terms, or click [Preview/Index](#) for advanced searching.

Enter [author names](#) as smith jc. Initials are optional.

Enter [journal titles](#) in full or as MEDLINE abbreviations. Use the [Journals Database](#) to find journal titles.

PubMed is a service of the National Library of Medicine, includes over 15 million citations for biomedical articles back to the 1950's. These citations include MEDLINE and additional life science journals. PubMed includes links to many sites providing full text articles and other related resources.

### Bookshelf Additions

*Molecular Biology of the Cell, 4th Ed.* and *The Genetic Landscape of Diabetes* are now available for interactive searching on the [Bookshelf](#).

### PubMed Enhancements!

[Full author](#) searching is now available and the Single Citation Matcher has been enhanced to include first author searching and an autocomplete feature for journal titles.

### Global NCBI Search Engine

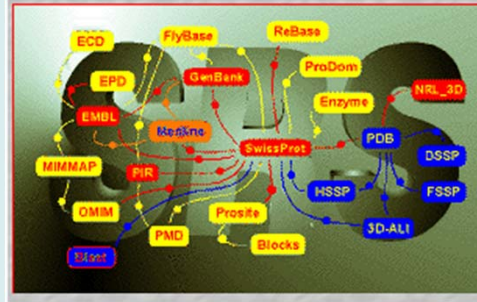
A growing number of Entrez databases can now be searched at once! [Go](#)



## How to get to data

There are several systems to **retrieve resources**:

SRS



ENTREZ



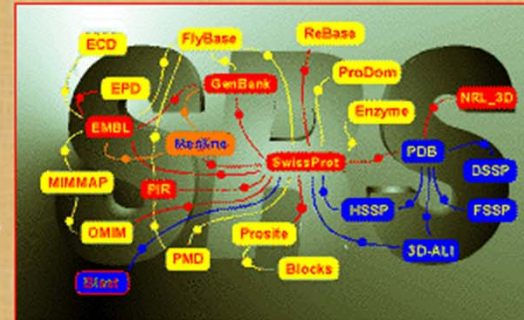
They are **integrated systems to retrieve sequences**, based on text search, to be used with the main databases, including PubMed, Nucleotide and Protein Sequences, Protein Structures, Whole Genomes, Taxonomy, etc.



SRS is an **integrated system of sequence retrieval**, created by EBI, to be used with the main databases.

## Sequence Retrieval System

*Network Browser for Databanks in Molecular Biology*



Start a new SRS session



The SRS Manual



List of all SRS5 servers



The SRS newsgroup



SRS Maintained



Back to the EBI home page

# SRSWWW


Easy to use

Allows recording of “sessions”

[Top Page](#) [Query Form](#) [Query Manager](#) [View Manager](#) [Databanks](#) [Help](#)

---

Select one or more databanks and continue (explode ☐ or collapse ☐ all groups)

 [Continue](#) [Reset](#)

☐ **Sequence** ☐ all

<input type="checkbox"/> <a href="#">SWISSPROT</a>	<input type="checkbox"/> <a href="#">SWISSNEW</a>	<input type="checkbox"/> <a href="#">NRDB</a>	<input type="checkbox"/> <a href="#">SWALL</a>
<input type="checkbox"/> <a href="#">UNIPROT SPROT</a>	<input type="checkbox"/> <a href="#">UNIPROT TREMBL</a>	<input type="checkbox"/> <a href="#">TREMBLNEW</a>	<input type="checkbox"/> <a href="#">TREMBL</a>
<input type="checkbox"/> <a href="#">SPTREMBL</a>	<input type="checkbox"/> <a href="#">SPTREMBLNEW</a>	<input type="checkbox"/> <a href="#">REMTREMBL</a>	<input type="checkbox"/> <a href="#">PIR</a>
<input type="checkbox"/> <a href="#">WORMPEP</a>	<input type="checkbox"/> <a href="#">DROSOPHILA</a>	<input type="checkbox"/> <a href="#">EMBLNEW</a>	<input type="checkbox"/> <a href="#">EMBL</a>
<input type="checkbox"/> <a href="#">EMBLEST</a>	<input type="checkbox"/> <a href="#">EMBLWGS</a>	<input type="checkbox"/> <a href="#">GENBANK</a>	<input type="checkbox"/> <a href="#">GENBANKEST</a>
<input type="checkbox"/> <a href="#">REFSEQP</a>	<input type="checkbox"/> <a href="#">SUBTILIST</a>		

☐ *SeqRelated*

☐ *TransFac*

☐ *Protein3DStruct*




☐ *Genome*

☐ *Mutations*

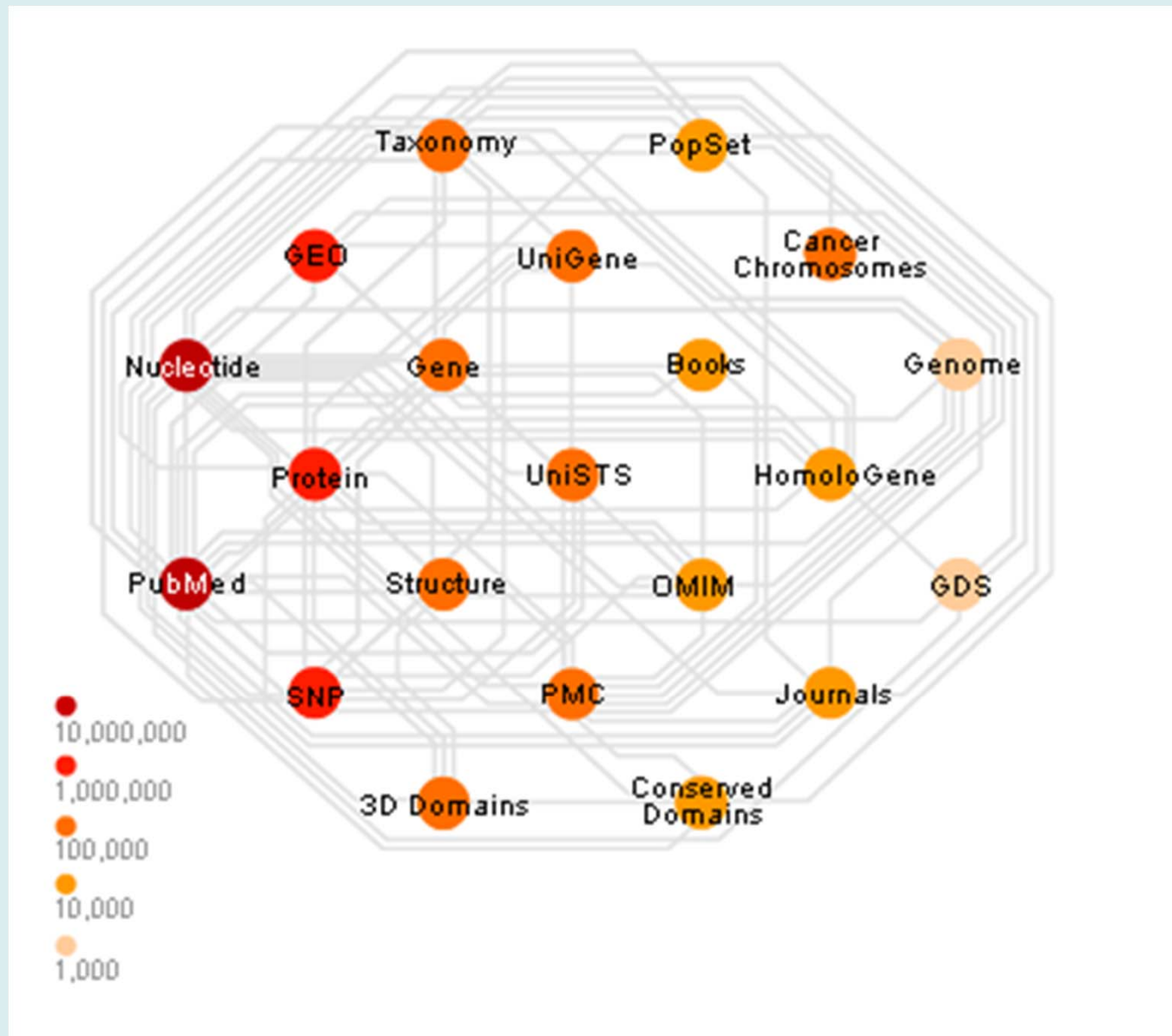
☐ *Others*

# ENTREZ

ENTREZ is an **integrated system of sequence retrieval**, created by NCBI, to be used with the main databases.

Welcome to the Entrez cross-database search page			
	<b>PubMed:</b> biomedical literature citations and abstracts		
	<b>PubMed Central:</b> free, full text journal articles		
	<b>Books:</b> online books		
	<b>OMIM:</b> online Mendelian Inheritance in Man		
	<b>Site Search:</b> NCBI web and FTP sites		
	<b>Nucleotide:</b> sequence database (GenBank)		
	<b>Protein:</b> sequence database		
	<b>Genome:</b> whole genome sequences		
	<b>Structure:</b> three-dimensional macromolecular structures		
	<b>Taxonomy:</b> organisms in GenBank		
	<b>SNP:</b> single nucleotide polymorphism		
	<b>Gene:</b> gene-centered information		
	<b>HomoloGene:</b> eukaryotic homology groups		
	<b>PubChem Compound:</b> small molecule chemical structures		
	<b>PubChem Substance:</b> chemical substances screened for bioactivity		
	<b>Genome Project:</b> genome project information		
	<b>UniGene:</b> gene-oriented clusters of transcript sequences		
	<b>CDD:</b> conserved protein domain database		
	<b>3D Domains:</b> domains from Entrez Structure		
	<b>UniSTS:</b> markers and mapping data		
	<b>PopSet:</b> population study data sets		
	<b>GEO Profiles:</b> expression and molecular abundance profiles		
	<b>GEO DataSets:</b> experimental sets of GEO data		
	<b>Cancer Chromosomes:</b> cytogenetic databases		
	<b>PubChem BioAssay:</b> bioactivity screens of chemical substances		
	<b>GENSAT:</b> gene expression atlas of mouse central nervous system		
	<b>Journals:</b> detailed information <i>about</i> the journals indexed in PubMed and other Entrez databases		
	<b>NLM Catalog:</b> catalog of books, journals, and audiovisuals in the NLM collections		
	<b>MeSH:</b> detailed information about NLM's controlled vocabulary		

## DATABASE INTEGRATION WITH ENTREZ



# Software Tools For Database Analysis

# Software Tools for Sequence Analysis

## General Packages:

Packages that offer a comprehensive range of bioinformatics tools for sequence analysis.

Most researchers would expect to use such packages at some time.

## Specialised Packages

Packages that offer tools for a particular type of analysis.

Used intensely by researchers in the relevant area, not at all by everyone else.

## WWW Resources

Tools whose nature inclines them to be primarily accessed over the network.

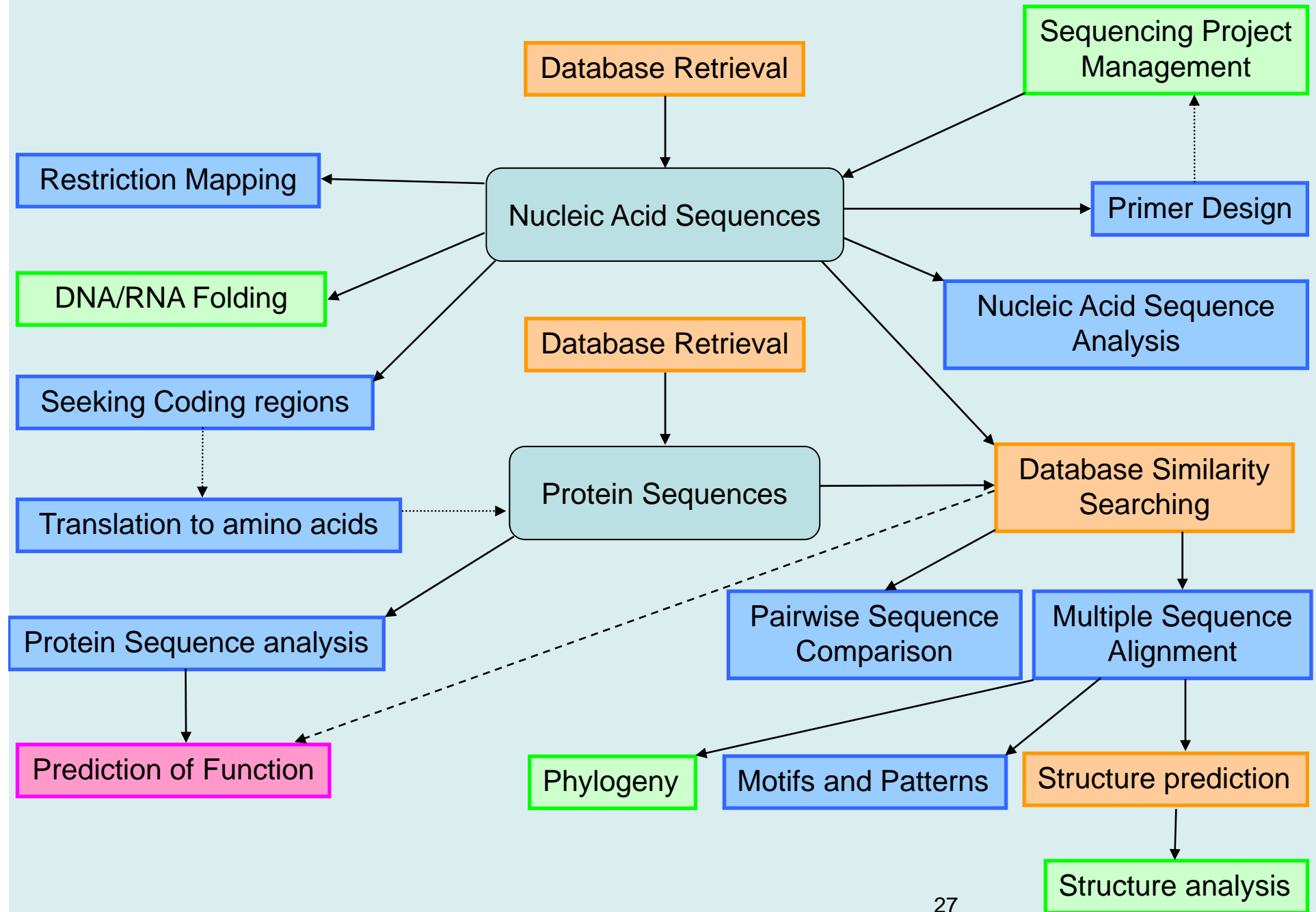
## These categorisations are very general

Many specialist programs are incorporated into the general packages.

Most things can be done at a web site somewhere.



# Sequence Analysis – an Overview



# Software Tools for Sequence Analysis

## General Packages:



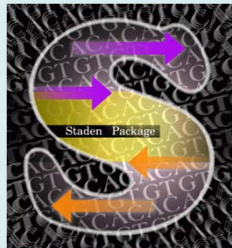
**Commercial**  
**WWW and X GUIs**  
**Widely available**

**UNIX only**  
**Comprehensive**



**Open source**  
**Several GUIs (java, WWW, X)**  
**Similar structure to the GCG package**

**UNIX only**  
**Comprehensive**



**Open source**  
**Excellent GUI including interactive graphical output**  
**Not comprehensive but allows access to EMBOSS**

**Windows, MacOS X, UNIX**

# Software Tools for Sequence Analysis

## General Packages:

Commercial

Expensive

## Other options

Windows PCs or Macintoshes

Good GUIs



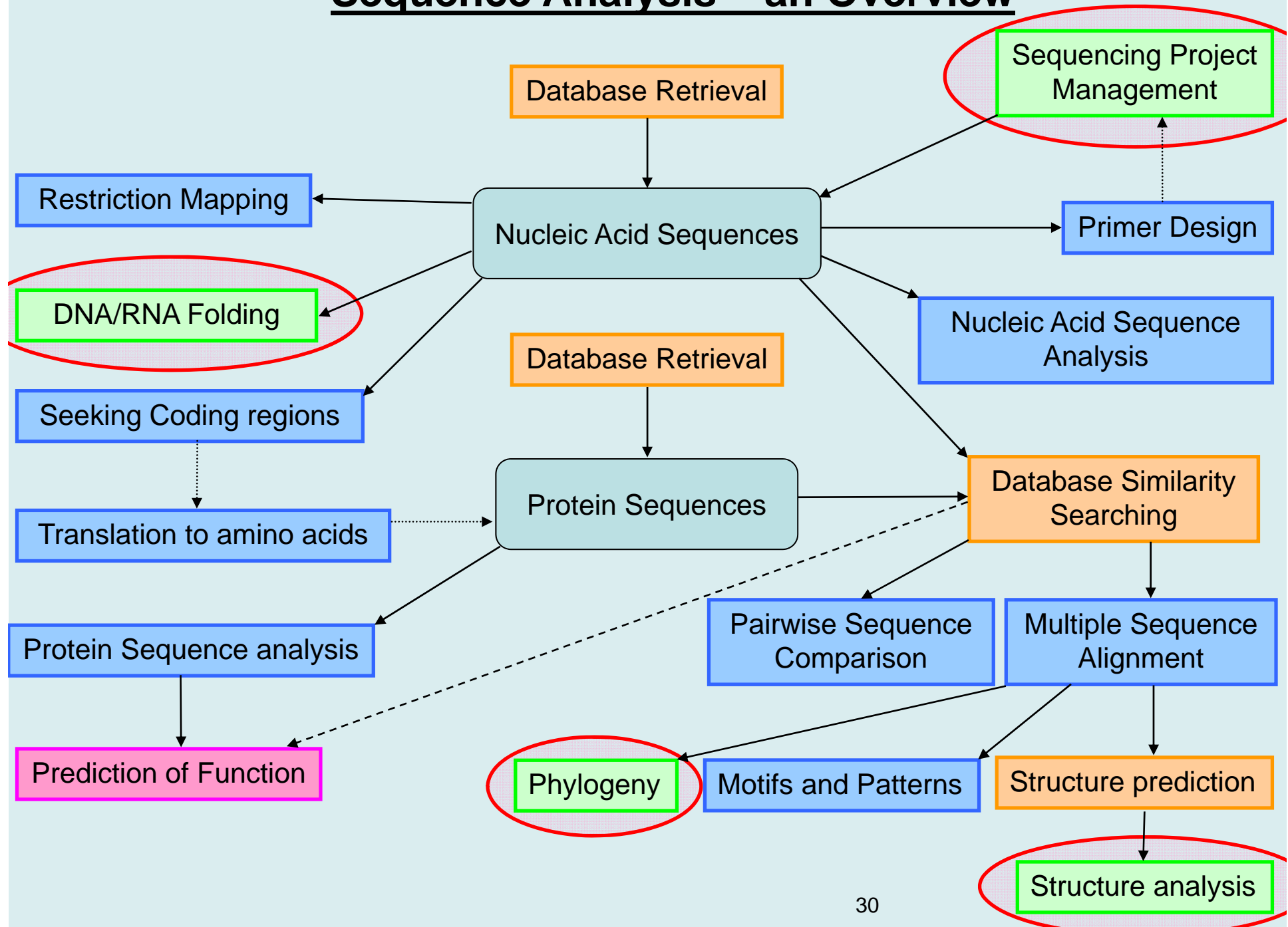
Public Domain

Windows, Macintosh, UNIX

Modern intuitive GUI

Access remote databases

# Sequence Analysis – an Overview

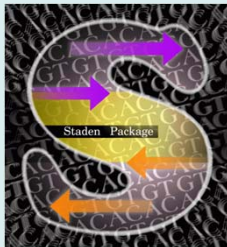


# Software Tools for Sequence Analysis

## Specialised Packages

### Sequencing Project Management

**“The Phred - Phrap Package”**  
By Phil Green et al



Free academic licence

Excellent base call confidence estimation (phred)

Excellent large scale contig assembler (phrap)

Available by anonymous ftp

Excellent GUI

Excellent contig editor

Excellent finishing tools

Simple confidence estimation

Contig assembler – not good for big projects

**BUT**

phred and phrap can be accessed from Staden GUI

# Software Tools for Sequence Analysis

## Specialised Packages

### DNA/RNA Folding



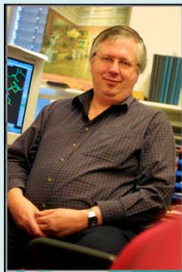
Michael Zuker's Programs

Free for academic use

Can be installed locally or run via a WWW page

Incorporated into the GCG general package

### Protein Structure Analysis



Whatif by Gert Vriend

Nominal fee for academic use

LINUX, IRIX, Windows



# Software Tools for Sequence Analysis

## Specialised Packages

### Protein Structure Analysis – for very rich people

SYBYL



IRIX, HP-UX, LINUX



IRIX, AIX, LINUX

Both systems are **very** impressive @ **very** expensive

# Software Tools for Sequence Analysis

## Specialised Packages

### Phylogeny



Available by anonymous ftp

Windows, Macintosh, UNIX

Incorporated into the EMBOSS general package



Commercial, but reasonable

UNIX, VMS, DOS and windows

Incorporated into the GCG general package



# Software Tools for Sequence Analysis

## WWW Resources

### Database Retrieval

#### Sequence Retrieval System

Retrieves MUCH more than sequences

The logo for the Sequence Retrieval System (SRS), consisting of the letters "SRS" in a bold, blue, sans-serif font.

Core elements free to academic sites



Bioscience AG

Implemented in many places

It is possible to integrate analysis tools

Elements of SRS are incorporated into EMBOSS

# Software Tools for Sequence Analysis

## WWW Resources

### Database Retrieval

Retrieves MUCH more than sequences



Access to NCBI databases only



Entrez, The Life Sciences Search Engine.

Entrez client software available by anonymous ftp

Most general packages include tools to access local sequence databases

**EMBOSS** programs can access sequences from remote **SRS** servers

# Software Tools for Sequence Analysis

## Database Similarity Searching



**BOTH** blast & fasta

## WWW Resources

Very popular, very widely available

Not sensitive – But extremely fast

Popular, widely available

Not sensitive – much slower than blast

Can be installed locally or run via a WWW page

Available by anonymous ftp ([blast](#), [fasta](#))

DNA/Protein query V DNA/Protein database

Incorporated into the GCG general package



# Software Tools for Sequence Analysis

## WWW Resources

### Structure prediction



Was consensus service now **JNet** only

**JNet** available by anonymous ftp

Older service, similar approach to **JNet**

Main element is called PHD



**Burkhard Rost**

Both **JPred** and **PHD** work best from aligned protein families

Simpler methods predicting from single sequences in most general packages

# Software Tools for Sequence Analysis

## WWW Resources

### Other WWW services

#### General Services:

EBI

Pasteur Institute

And many more

#### Protein sequence analysis

Expasy

#### Gene finding

**genscan** at the MIT (Free academic license)

Simple gene finding in most general packages

#### Primer design

**primer3** at the MIT (Available by anonymous ftp)

Primer design in most general packages

Primer design in EMBOSS is primer3