



# Unifying Bio-information on the Grid

Hideo Matsuda  
Department of Bioinformatic Engineering  
Osaka University

E-mail: matsuda@ist.osaka-u.ac.jp



2

## Data Grid for Bio-information

- Bio-information has tremendous amount data and it's still rapidly increase.
- Distributed into many databases (500+ DBs)
  - ✦ Very heterogeneous (different contents, different formats)
  - ✦ Many similar and related data (need approximate matching, difficult to be indexed)
  - ✦ Physical integration is not easy
- ➔ Coordinating DBs by the Grid

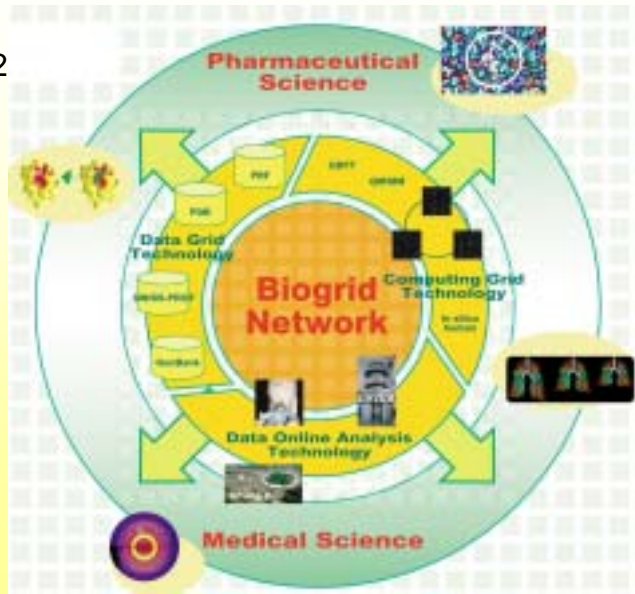


# Our BioGrid

<http://www.biogrid.jp>

3

- Started from 2002
- Goals
  - Technology Development for: Pharmaceutical (drug discovery) and Medical Sciences.
- Leader: Shinji Shimojo (CMC, Osaka Univ.)
- Government Support (MEXT): 5years (3~4M\$/year)



## Heterogeneity of Bio-Databases (Flat Files)

4

DNA sequence DB(DDBJ)	Protein sequence DB (SWISS-PROT)	Protein 3D structure DB (PDB)
<b>LOCUS</b> ECORBS 6197 bp DNA 30-OCT-1994 <b>DEFINITION</b> E.coli ...rbsB ... transport system, ... <b>ACCESSION</b> M13169 M13517 <b>NID</b> g147511 <b>KEYWORDS</b> high affinity ribosome <b>SOURCE</b> E.coli K12 DNA. <b>ORGANISM</b> Escherichia coli <b>REFERENCE</b> 1 (sites) <b>AUTHORS</b> ... <b>TITLE</b> ... <b>JOURNAL</b> ... <b>MEDLINE</b> 84032513 <b>MEDLINE</b> 84032513 <b>FEATURES</b> <b>CDS</b> 3112..4002 /gene="rbsB" /translation="MKKGTVLNS... <b>ORIGIN</b> 94 bp upstream of BclI 1 ctcaggttcg aaatcctaac	<b>ID</b> RBSB_ECOLI ... <b>AC</b> P02925; <b>DT</b> 21-JUL-1986 (...CREATED) <b>DT</b> 21-JUL-1986 (... UPDATE) <b>DT</b> 01-NOV-1995 (... UPDATE) <b>DE</b> D-RIBOSE-BINDING ... <b>GN</b> RBSB OR RBSP OR PRLB. <b>OS</b> ESCHERICHIA COLI. <b>OC</b> PROKARYOTA; GRACILICUTES; <b>OC</b> ENTEROBACTERIACEAE. <b>RN</b> [1]RP SEQUENCE FROM N.A... <b>RX</b> MEDLINE; 84032513. <b>RA</b> GROARKE J.M., MAHONEY W.C., <b>RA</b> ZALKIN H., HERMODSON M.A.; <b>RL</b> J. BIOL. CHEM. 258:... <b>CC</b> -I- FUNCTION: INVOLVED IN ... <b>KW</b> TRANSPORT; SUGAR TRANSPORT; <b>KW</b> 3D-STRUCTURE. <b>SQ</b> SEQUENCE 296 AA; ... MNMKKLATLV SAVALSATVS ANA... ...	<b>HEADER</b> SUGAR TRANSPORT 23-SEP-9 <b>COMPND</b> D-RIBOSE-BINDING PROTEIN <b>COMPND</b> 2 (G134R) COMPLEXED WITH <b>SOURCE</b> (ESCHERICHIA COLI)... <b>SOURCE</b> 2 EXPRESSION PLASMID) <b>AUTHOR</b> S.L.MOWBRAY,A.J.BJORKMAN <b>REVDAT</b> 1 26-JAN-95 1DRJ 0 <b>JRNL</b> AUTH A.J.BJORKMAN <b>JRNL</b> TITL PROBING PR... <b>JRNL</b> TITL 2 RIBOSE-BINDING <b>JRNL</b> REF TO BE PUBLISHED <b>JRNL</b> REFN ASTM <b>SEQRES</b> 1 271 LYS ASP THR ... <b>SEQRES</b> 2 271 PRO PHE PHE ... ... <b>HELIX</b> 1 A PRO 14 LEU ... <b>HELIX</b> 2 B PRO 43 LEU ... ... <b>ATOM</b> 1 N LYS 1 x1 y1 z1 <b>ATOM</b> 2 CA LYS 1 x2 y2 z2 ...



# BioDBs now provide XML formats

## But they are still heterogeneous!

5

```

LOCUS      ECORBS 6197 bp DNA ..
           30-OCT-1994
DEFINITION E.coli ...rbsB ....
           transport system, ...
ACCESSION M13169
KEYWORDS  high affinity ribose
SOURCE    E.coli K12 DNA.
ORGANISM  Escherichia coli
REFERENCE 1 (sites)
AUTHORS   ...
TITLE     ...
JOURNAL   ...
MEDLINE   84032513
FEATURES
CDS       3112..4002
           /gene="rbsB"
           /translation="MKKGTVLNS...
           ...
ORIGIN    94 bp upstream of BclI
           1 ctcaggttcg aaatctaac. ....
    
```

DDBJ FF  
to  
DDBJ XML

```

<?xml version="1.0" standalone="no"?>
<!DOCTYPE DDBJXML SYSTEM "DDBJXML.dtd">
<DDBJXML>
<LOCUS> ECORBS </LOCUS>
  <LAST_UPDATE> 30-OCT-1994 </LAST_UPDATE>
<DEFINITION> E.coli ...rbsB ....
  transport system, ...</DEFINITION>
<ACCESSION> M13169 </ACCESSION>
<KEYWORDS> high affinity ribose ..</KEYWORDS>
<SOURCE> Escherichia coli </SOURCE>
<ORGANISM> Escherichia coli </ORGANISM>
<REFERENCE ID="1">
  <AUTHORS> ... </AUTHORS>
  <TITLE> ... </TITLE>
  <JOURNAL> ... </JOURNAL>
<FEATURES>
<CDS> <location>3112..4002</location>
  <qualifiers name="gene"> rbsB </gene>
  <qualifiers name="translation"> MKKGTVLNS...
  </qualifiers>
  ...
<SEQUENCE> ctcaggttcgaaatctaac. ...</SEQUENCE>
    
```



# Common Data Translation Table

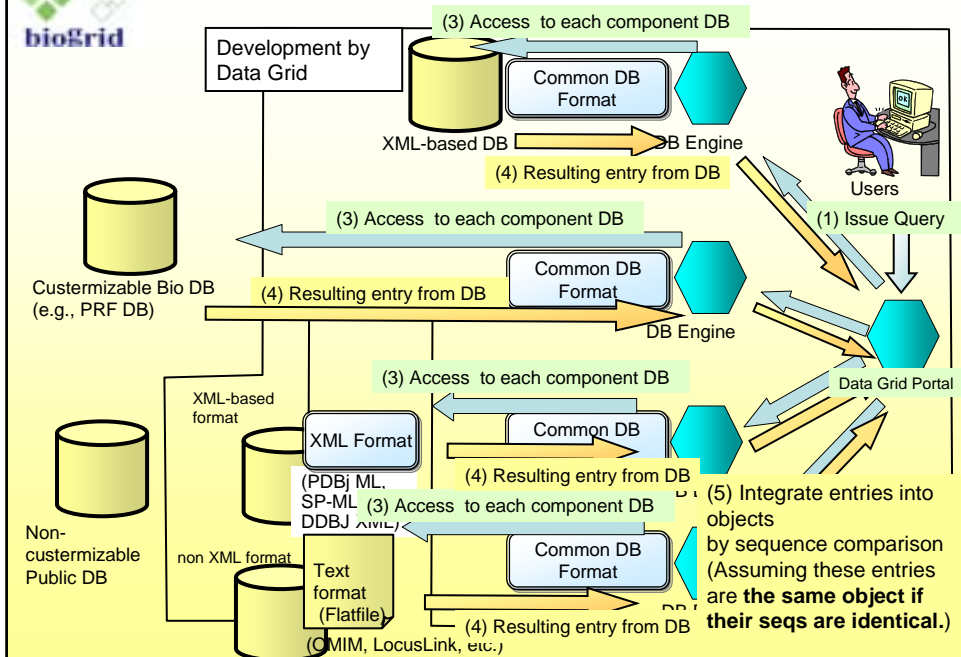
6

	Common Format	SWISS-PROT	PIR	PDB
Name space prefix	---	sp	pir	pdb
Entry (1 protein)	entry	entry	ProteinEntry	PDBj
Entry ID	entry/id	entry/@name	ProteinEntry/@id	PDBj/@entry_ID
Accession	entry/accession	entry/@primaryACC	ProteinEntry/header/accession	PDBj/@entry_ID
Protein name	entry/name	entry/proteinName	ProteinEntry/protein/alt-name	PDBj/main/entity/entity_item/description
Protein alt name	entry/alt-name	entry/proteinName/synonyms/synonym	ProteinEntry/protein/name	PDBj/main/entity/entity_item/common name
Gene name	entry/gene	entry/genes/gene/@name	ProteinEntry/genetics/gene/uid	PDBj/main/entity/entity_item/src_en/gene
Species (formal)	entry/organism/scientific	entry/organisms/organism/name	ProteinEntry/organism/formal	PDBj/main/entity/entity_item/src_en/scientific_name
Species (common)	entry/organism/common	entry/organisms/organism/name	ProteinEntry/organism/common	PDBj/main/entity/entity_item/src_en/common_name
NCBI Taxonomy ID	entry/organism/taxid	entry/organisms/organism/taxid	---	---
Reference	entry/reference	entry/references/reference	ProteinEntry/reference/refinfo	PDBj/main/citation/citation_item
Reference ID (MEDLINE)	entry/reference/medlineID	entry/references/reference/pubMed/@medlineID	ProteinEntry/reference/refinfo/xrefs/xref/uid[db="MUID"]	---
Reference ID (PubMed)	entry/reference/pubmedID	entry/references/reference/pubMed/@value	ProteinEntry/reference/refinfo/xrefs/xref/uid [db="PMID"]	---
Function	entry/function	entry/comments	---	---
Enzyme EC#	entry/EC_num	entry/proteinName	---	PDBj/main/entity/entity_item/EC
GeneOntology	entry/GO	entry/goTerms	---	---
GO term	entry/GO/term	---	---	---
GO description	entry/GO/description	---	---	---
Keyword	entry/keyword	entry/keywords/keyword	ProteinEntry/keywords/keyword	PDBj/main/struct/keywords
Structure	entry/structure	---	---	PDBj/main/struct/struct
Sequence (amino acid)	entry/sequence	entry/sequence	ProteinEntry/sequence	PDBj/main/entity/entity_item/sequence_letter_code
Feature	entry/feature	entry/features/feature	ProteinEntry/feature	PDBj/struct/site
feature type	entry/feature/type	entry/features/feature/@key	ProteinEntry/feature/feature-type	PDBj/struct/site/@id
feature description	entry/feature/description	entry/features/feature/@description	ProteinEntry/feature/description	PDBj/struct/site/site_gen/details



# Unified Access to DBs

7



# An Example of Common Database Format

8

```

<?xml version="1.0"?>
<entry xmlns="bio data grid" xmlns:sp="swiss-prot" xmlns:pir="pir" xmlns:pdb="pdb" >
  <id>..</id> <acc>..</acc> <date> .. </date>
  <name db="pir" acc=".." >protein</name> <alt-name db="pir" acc="..">..</alt-name>
  <gene db="sp" acc="..">gene name</gene>
  <organism db="pir" acc="..">
    <scientific db="pir" acc="..">Homo sapiens</scientific>
    <common db="pir" acc="..">human</common> <taxId db="pir" acc="..">..</taxId>
  </organism>
  <reference db="pir" acc="..">
    <author>..</author> <citation>Journal of ..</citation> <volume>1</volume>
    <year>2002</year>
    <title>paper title</title> <first_page>10</first_page> <last_page>20</last_page>
    <medlineID>..</medlineID> <pubmedID>..</pubmedID>
  </reference>
  <function db="sp" acc=".." > <desc>..</desc> </function>
  <keyword db="sp" acc="..">..</keyword>
  <GO db="sp" acc=".." > <term>0006118</term><desc>electron trans..</desc> </GO>
  <feature db="pir" acc=".." > <type>active site</type> <desc>..</desc>
    <start>..</start> <end>..</end> </feature>
  <sequence db="sp" acc=".." length=".."> MDPVVVLVLGL.. </sequence>
  <sp:data>SWISS-PROT data</sp:data>
  <pir:data>PIR data</pir:data>
  <pdb:data>PDB data</pdb:data>
</entry>

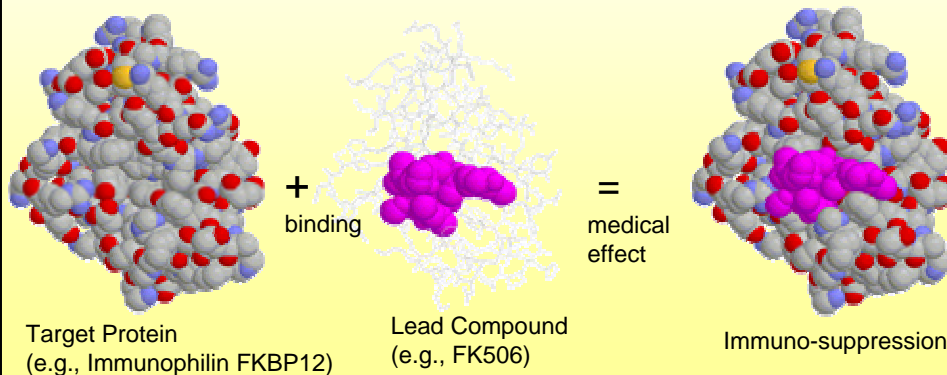
```



# Application of Data Grid to Drug Discovery

9

- Drug Related Data:
  - ➔ Drug Target Proteins (mainly receptors)
  - ➔ Lead Compounds (drug candidates)
  - ➔ Medical Effect
- Correlate all information by unified DB access (active site locations, charges, known binding partners, etc.)



## Summary

10

- Biology is now going to be an *Information Science*.
- Bio DBs are very heterogeneous. We need some unified access with some *common* data format.
- Coordinating DBs are useful for biological applications: Drug Discovery, etc.