# Data Mining for the Americas: Biowebs

January 29, 2003
AMPATH: Pathway of the Americas

Robert Grossman
University of Illinois at Chicago (grossman at uic.edu)
Open Data (rlg at opendata.biz)

# Part 1.  Three Trends

Three trends driving the
emergence of biowebs.

# Trend 1. Proliferation of Biological Databases

NDB

swissprot

FlyBase

ExPASy

PDB
PROTEIN DATA BANK

e!

from dozens to
hundreds of databases

SGD

UCSC

WormBase

NCBI

# …Usually in the Wrong Format

# Trend 2: More and More Discoveries will be Across Databases



❑ Pearson's Law: The usefulness of a column of data varies as the square of the number of columns it is compared to.

# Example:
# Microarray Data & Clinical Data



\<publication\>
!Citation=Alizadeh AA et al.(2000) Nature 403:503-11
!Title=Distinct types of diffuse large B-cell lymphoma (DLBCL) identified by gene expression profiling.
!PubMedID=10676951

# Trend 3. Near a Trifurcation Point



Biological Databases

2003-2008

Biowebs – remote data analysis and distributed mining

Biogrids – transparent high end computing

Biological semantic webs – sem. webs for biological knowledge

# Data Grids vs. Data Webs

What is more valuable: other peoples' cycles or data?

Browsing & Casual Exploration

Data Webs

- Searching
- Exploration
- Casual correlation

- Security
- Authorization
- Scheduling

Collaborations

Data Grids

Distributed Computer

Web Based Computing

# Part 2.  Examples

Biogrids, biowebs, and all that.

# Example 1. BIRN



❑ NIH Sponsored project developing collaborative infrastructure for studying brains in humans and animals

# Example 2. OptIPuter



Photonic DataSpace

- ❑ Data intensive computing over photonic networks
- ❑ Replication of the protein data bank (PDB).
- ❑ Linked with a chemical library of small organics molecules.
- ❑ Distributed docking algorithms

# Example 3. Molecular DataSpace



- ❑ How do you interactively explore other people's data?
- ❑ How do you overlay other peoples data on your own?
- ❑ How do you do distributed data mining?

# Simplify the Integration of Two or More Distributed Data Sets

About

Select one or more structures and press **Submit** to view, download, visualize and perform a Docking Algorithm.

| Protein Structure Data | Small Organic Compound - Drug Data |
|---|---|

1 – 500 ▭ Go Prev | Next

1 – 100 ▭ Go Prev | Next

```
101m – SPERM WHALE MYOGLOBIN F46...
102l – LYSOZYME INSERTION MUTANT...
102m – SPERM WHALE MYOGLOBIN H64...
103l – PHAGE T4 LYSOZYME INSERTI...
103m – SPERM WHALE MYOGLOBIN H64...
104l – LYSOZYME INSERTION MUTANT...
104m – SPERM WHALE MYOGLOBIN N–B...
```

```
d101d – NETROPSIN...
d102d – PROPAMIDINE...
d107d – DUOCARMYCIN...
d108d – THE BISINTERCALATING DYE...
d11gs – ETHACRYNIC ACID–GLUTATHIO...
d12gs – S–NONYL–GLUTATHIONE...
d13gs – SULFASALAZINE...
```

Submit Click here to View Structure, View or download O... files

Clear Clear Selection

**Database 1**

**Database 2**

Molecular Home

# Easy to Overlay External, Third Party Data with Local Data

| NSCID | Weight | H-Donors | H-Acceptor | Formula | LogP | SmileString | RecpS | POSorNEG | Hf | NCSP3-R22 | NS-R13 | NOH-R3 | Activity |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| d99 | 364.415 | 2 | 7 | C17H20N2O5S | 1.807 | CC(C)C(=O)c2c(Sc1ccccc1)n(COCCO)c(=O)[nH]c2=O | 1.6776 | -1.1048 | -142.56 | 9 | 0 | 1 | 4.92 |
| d98 | 323.366 | 4 | 7 | C14H17N3O4S | 0.732 | Cc2c(Sc1cccc(N)c1)n(COCCO)c(=O)[nH]c2=O | 1.5957 | -1.1131 | -105.18 | 1 | 27 | 1 | 3.6 |
| d97 | 394.444 | 2 | 6 | C21H18N2O4S | 2.224 | O=c2[nH]c(=O)n(COCCO)c(Sc1ccccc1)c2C#Cc3ccccc3 | 1.8155 | -1.0701 | -23.39 | 0 | 0 | 1 | 5.47 |
| d96 | 318.347 | 2 | 6 | C15H14N2O4S | 0.923 | C#Cc2c(Sc1ccccc1)n(COCCO)c(=O)[nH]c2=O | 1.8947 | -1.0856 | -44.31 | 0 | 0 | 1 | 4.74 |
| d95 | 402.482 | 2 | 6 | C19H18N2O4S2 | 2.891 | O=c2[nH]c(=O)n(COCCO)c(Sc1ccccc1)c2Sc3ccccc3 | 1.5547 | -0.9164 | -66.97 | 0 | 0 | 1 | 4.68 |
| d94 | 413.447 | 3 | 8 | C20H19N3O5S | 2.004 | O=C(Nc1ccccc1)c3c(Sc2ccccc2)n(COCCO)c(=O)[nH]c3=O | 1.8096 | -1.084 | -99.11 | 0 | 0 | 1 | 4.74 |
| d93 | 352.361 | 2 | 8 | C15H16N2O6S | 1.276 | COC(=O)c2c(Sc1ccccc1)n(COCCO)c(=O)[nH]c2=O | 1.6067 | -1.0866 | -177.38 | 1 | 0 | 1 | 5.18 |
| d92 | 351.376 | 4 | 8 | C15H17N3O5S | 0.087 | Cc2c(Sc1cccc(C(N)=O)c1)n(COCCO)c(=O)[nH]c2=O | 1.6345 | -0.9326 | -142.88 | 1 | 27 | 1 | 3.51 |
| d91 | 352.361 | 3 | 8 | C15H16N2O6S | 1.05 | Cc2c(Sc1ccc(C(=O)O)cc1)n(COCCO)c(=O)[nH]c2=O | 1.7677 | -1.1281 | -193.97 | 1 | 27 | 1 | 3.45 |
| d90 | 350.389 | 2 | 7 | C16H18N2O5S | 1.282 | CC(=O)c2ccc(Sc1c(C)c(=O)[nH]c(=O)n1COCCO)cc2 | 1.6622 | -1.1213 | -140.23 | 1 | 64 | 1 | 3.96 |
| d9 | 342.797 | 2 | 6 | C14ClH15N2O4S | 2.059 | Cc2c(Sc1cccc(Cl)c1)n(COCCO)c(=O)[nH]c2=O | 1.6239 | -1.0905 | -111.42 | 1 | 27 | 1 | 4.89 |
| d89 | 338.378 | 2 | 7 | C15H18N2O5S | 1.377 | COc2ccc(Sc1c(C)c(=O)[nH]c(=O)n1COCCO)cc2 | 1.6824 | -1.1077 | -143.98 | 1 | 64 | 1 | 3.6 |
| d88 | 324.351 | 3 | 7 | C14H16N2O5S | 1.048 | Cc2c(Sc1ccc(O)cc1)n(COCCO)c(=O)[nH]c2=O | 1.6807 | -1.1031 | -148.42 | 1 | 64 | 1 | 3.56 |
| d87 | 333.361 | 2 | 7 | C15H15N3O4S | 0.918 | Cc2c(Sc1ccc(C#N)cc1)n(COCCO)c(=O)[nH]c2=O | 1.7283 | -1.08 | -76.67 | 1 | 64 | 1 | 3.6 |
| d86 | 353.349 | 2 | 9 | C14H15N3O6S | 1.028 | Cc2c(Sc1ccc(N(=O)=O)cc1)n(COCCO)c(=O)[nH]c2=O | 1.7135 | -1.5648 | -104.8 | 1 | 64 | 1 | 3.72 |
| d85 | 359.804 | 3 | 7 | C14ClH15N2O4S | 1.256 | Cc2c(Sc1ccc(Cl)cc1)n(COCCO)c(=O)[nH]c2=OO | 1.6179 | -1.0692 | -115.7 | 1 | 64 | 1 | 3.6 |
| d84 | 359.349 | 3 | 8 | C14FH15N2O4S | 0.306 | Cc2c(Sc1ccc(F)cc1)n(COCCO)c(=O)[nH]c2=OOO | 1.5366 | -1.0717 | -154.3 | 1 | 64 | 1 | 3.6 |

Local Database 1

External Database 2

# Data and Metadata Separated

# Data Can be Streamed…

# Support PMML Based Analytics



Activity vs. Predicted Activity



Decision Tree generated by WEKA

```
Options: -B 10 -W weka.classifiers.j48.J48 -- -C 0.25 -M 2
Regression by discretization
Class attribute discretized into 10 values
Subclassifier: weka.classifiers.j48.J48
J48 pruned tree
------------------
NCSP3-R22 <= 1
|   NS-R13 <= 27
|   |   H-Donors <= 2
|   |   |   H-Acceptor <= 5
|   |   |   |   Weight <= 248.299: '(-inf-3.962]' (2.0)
|   |   |   |   Weight > 248.299
|   |   |   |   |   LogP <= 2.687: '(5.498-6.01]' (6.0/1.0)
|   |   |   |   |   LogP > 2.687
|   |   |   |   |   |   Hf <= -51.57: '(4.986-5.498]' (2.0)
|   |   |   |   |   |   Hf > -51.57: '(5.498-6.01]' (2.0/1.0)
|   |   |   H-Acceptor > 5
|   |   |   |   H-Acceptor <= 6
|   |   |   |   |   Weight <= 384.51
|   |   |   |   |   |   LogP <= 1.734
|   |   |   |   |   |   |   NCSP3-R22 <= 0: '(4.474-4.986]' (3.0/1.0)
|   |   |   |   |   |   |   NCSP3-R22 > 0: '(4.986-5.498]' (3.0/1.0)
|   |   |   |   |   |   LogP > 1.734
|   |   |   |   |   |   |   Weight <= 338.378: '(5.498-6.01]' (6.0/2.0)
|   |   |   |   |   |   |   Weight > 338.378
|   |   |   |   |   |   |   |   Weight <= 374.497: '(4.474-4.986]' (3.0)
|   |   |   |   |   |   |   |   Weight > 374.497: '(3.962-4.474]' (3.0)
|   |   |   |   |   Weight > 384.51
|   |   |   |   |   |   LogP <= 2.461: '(4.986-5.498]' (4.0)
|   |   |   |   |   |   LogP > 2.461
|   |   |   |   |   |   |   LogP <= 3.781: '(4.474-4.986]' (2.0)
|   |   |   |   |   |   |   LogP > 3.781: '(4.986-5.498]' (2.0/1.0)
```
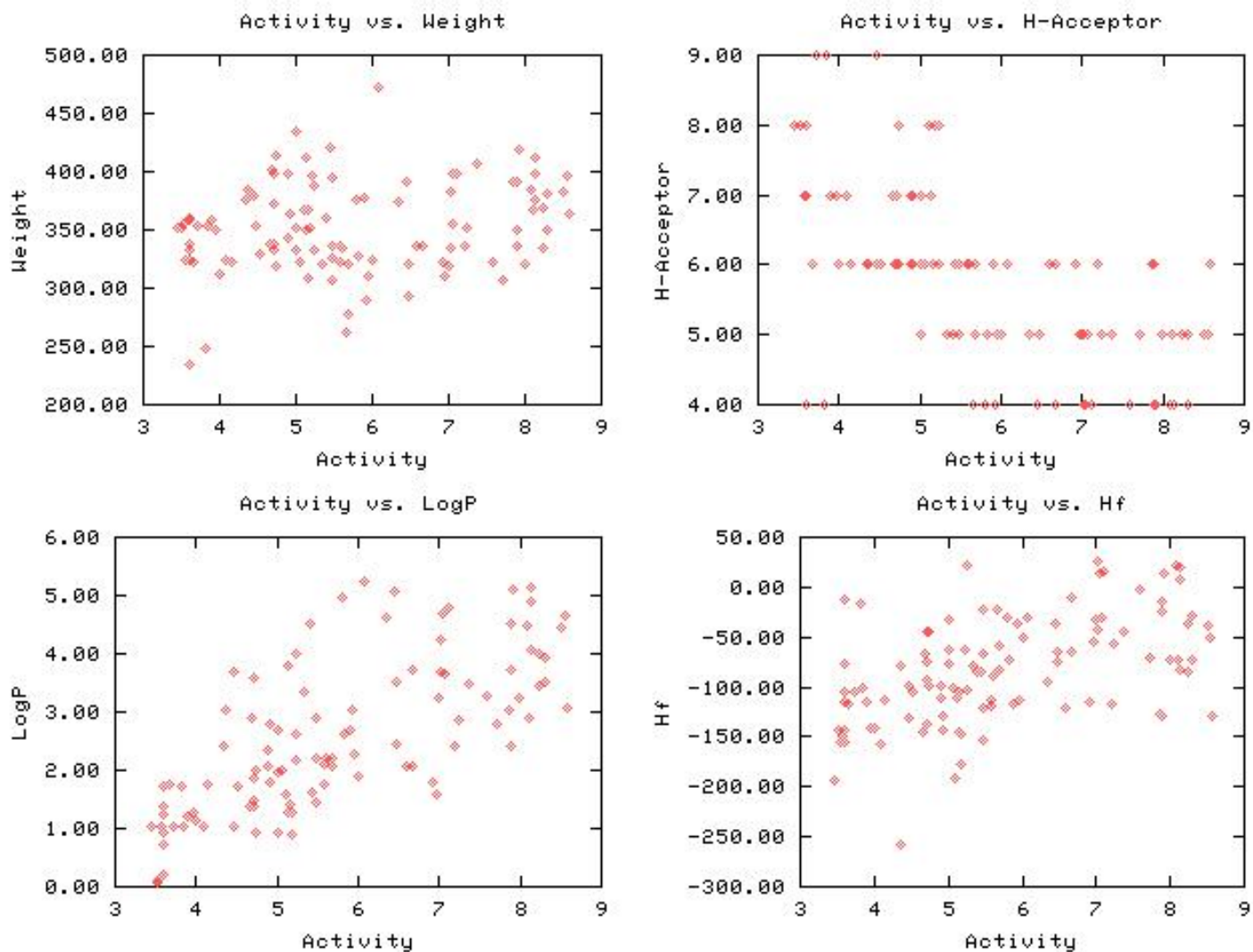
# Integrate with Open Source Analytics

# Part 3.  Strategy for the Americas

It's the data.

# 1. Think Small, Medium and Large

❑ Large Science: high energy physics, astronomy, …

❑ Medium science: sequencing an organism, biodiversity surveys, …

❑ Small science: creating interesting bioinformatics databases and resources, overlaying external data over your data to do new science.

# 2. Open Data:
# Free & Controlled Biodata

❑ Lawrence Lessig of Stanford Law School has highlighted the battle between "free" and "controlled" web resources

❑ Genbank created a culture of free sequence data

❑ There is also a culture of proprietary data

❑ Consider part of your mission to create *open data repositories*

# Open Source Projects

❑ Open Source libraries
  – Bioperl, Biojava, Biopython

❑ Open Source protocols
  – DWTP, DAS, MOBY, OmniGene, G2G,…

❑ Open Source end-user applications
  – Genquire, Generic Genome Browser, PyMol, Molecular DataSpace…

# For More Information

- ❑ Data webs – www.dataspaceweb.net
- ❑ Data web servers – www.sourceforge.net/projects/dataspace
- ❑ Robert Grossman – grossman @ uic.edu or rlg @ opendata.biz