# Bioinformatics 3 V 3 – Data for Building Networks

Fri, Oct 31, 2014

Bioinformatics 3 – WS 14/15

# Graph Layout 1

### Requirements:

- fast and stable
- nice graphs
- visualize relations
- symmetry
- interactive exploration



### Force-directed Layout:

based on energy minimization

- $\rightarrow$  runtime
- $\rightarrow$  mapping into 2D



- H3: for hierarchic graphs
- $\rightarrow$  MST-based cone layout
- $\rightarrow$  hyperbolic space





 $\rightarrow$  efficient layout for **biological data???** 



<u>Aim</u>: analyze and visualize **homologies** within the **protein universe** 50 genomes, 145579 proteins,  $21 \times 10^9$  BLASTP pairwise sequence comparisons

Expectations:

- homologs will be close together
- fusion proteins ("Rosetta Stone proteins") will link proteins of related function.

 $\rightarrow$  need to visualize an extremely large network!

 $\rightarrow$  develop a stepwise scheme

### LGL: stepwise scheme

### (0) create network from BLAST E-score

145'579 proteins  $E < 10^{-12} \rightarrow 1'912'684$  links , 30737 proteins in the largest cluster

### (1) separate original network into connected sets

11517 connected components, 33975 proteins w/out links

(2) apply force directed **layout** to each **component independently**, based on a MST

# (3) integrate connected sets into one coordinate system via a **funnel process**, starting from the largest set

The first connected set is placed at the bottom of a potential funnel. Other sets are placed one at a time on the rim of the potential funnel and allowed to fall towards the bottom where they are frozen in space upon collision with the previous sets.

### Component layout I

For each component:

 $\rightarrow$  start from the **root node** of the MST



**Centrality**: minimize total distance to all other nodes in the component

$$v_{root} = \min\left(\sum_{(v,u)\in V} d(v,u)\right)$$

Level *n*-nodes: nodes that are *n* links away from the root in the MST

#### Layout $\rightarrow$ place **root** at the **center**

## Component Layout II

- start with root node of the MST
- place level-1 nodes on circle (sphere) around root, add all links, relax springs (+ short-range repulsion)
- place level-2 nodes on circles (sphere) outside their level-1 descendants, add all links, relax springs
- place level-3 nodes on circles (sphere) outside their level-2 descendants,



### **Combining the Components**

When the components are finished  $\rightarrow$  **assemble** using energy **funnel** 

- place largest component at bottom
- place next smaller one somewhere on the rim, let it slide down
- $\rightarrow$  freeze upon contact



### No information in the relative positions of the components!!!







Adai et al. J. Mol. Biol. 340, 179 (2004)

8

### Annotations in the Largest Cluster



Related functions in the same regions of the cluster  $\rightarrow$  predictions

## **Clustering of Functional Classes**



## **Fusion Proteins**



Fusion proteins connect two protein homology families

A, A', A", AB and B, B', AB

 $\rightarrow$  historic genetic **events**: fusion, fission, duplications, ...

#### Also in the network:

homologies <=> edges

remote homologies <=> in the same cluster

non-homologous functional relations <=> adjacent, linked clusters

### **Functional Relations between Gene Families**

Examples of spatial localization of protein function in the map

**A**: the linkage of the tryptophan synthase  $\alpha$  family to the functionally coupled but non-homologous  $\beta$  family by the yeast tryptophan synthase  $\alpha\beta$  fusion protein,

**B**: protein subunits of the pyruvate synthase and alpha-ketoglutarate ferredexin oxidoreductase complexes

**C**: metabolic enzymes, particularly those of acetyl CoA and amino acid metabolism

 $\rightarrow$  DUF213 likely has metabolic function!



Bioinformatics 3 – WS 14/15

### And the Winner iiiis...



Compare the layouts from

 A: LGL – hierarchic force-directed layout according to MST
→ structure from homology

**B**: **global force**-directed layout without MST  $\rightarrow$  no structure, no components visible

**C**: InterViewer – collapses similar nodes  $\rightarrow$  reduced complexity

## Graph Layout: Summary

Approach	Idea		
Force-directed spring model	relax energy, springs		
Force-directed spring-electric model	relax energy, springs for links, Coulomb repulsion between all nodes		
H3	spanning tree in hyperbolic space		
LGL	hierarchic, force-directed algorithm for modules		

### A "Network"



Bioinformatics 3 – WS 14/15

### **Protein Complexes**



Complex formation may lead to increased diversity

Cooperation and allostery

# **Gel Electrophoresis**

Electrophoresis: directed diffusion of charged particles in an electric field





Put proteins in a spot on a gel-like matrix, apply electric field

- $\rightarrow$  separation according to size (mass) and charge
- $\rightarrow$  identify constituents of a complex

Nasty details: protein charge vs. pH, cloud of counter ions, protein shape, denaturation, ...



### **SDS-PAGE**

For better control: denature proteins with detergent

Often used: sodium dodecyl sulfate (SDS)

 $\rightarrow$  denatures and coats the proteins with a negative charge

- $\rightarrow$  charge proportional to mass
  - $\rightarrow$  traveled distance per time

$$x \propto rac{1}{\log(M)}$$

 $\rightarrow$  SDS-polyacrylamide gel electrophoresis

After the run: **staining** to make proteins visible

For "quantitative" analysis: compare to **marker** (set of proteins with known masses)



Image from Wikipedia, marker on the left lane

## **Protein Charge?**

Protein charge at pH=7 $\cong \sum Lys + \sum Arg - \sum Asp - \sum Glu + \sum co - factors$ 

Main source for charge differences: pH-dependent protonation states

<=> Equilibrium between

- density (pH) dependent H<sup>+</sup>-binding and
- density independent H<sup>+</sup>-dissociation

Probability to have a proton:

$$P = \frac{1}{1 + 10^{pH - pK}}$$

pKa = pH value for 50% protonation

Asp 3.7–4.0 ... His 6.7–7.1 ... Lys 9.3-9.5

### Each H<sup>+</sup> has a +1e charge

 $\rightarrow$  **Isoelectric point**: pH at which the protein is **uncharged** 

 $\rightarrow$  protonation state cancels permanent charges



### 2D Gel Electrophoresis

# **Two steps**: i) separation **by isoelectric** point via pH-gradient ii) separation **by mass** with SDS-PAGE



 $\rightarrow$  Most proteins differ in mass and isoelectric point (pl)

### Mass Spectrometry

Identify constituents of a (fragmented) complex via their mass/charge patterns, detect by pattern recognition with machine learning techniques.





http://gene-exp.ipkgatersleben.de/body\_methods.html

# Tandem affinity purification

Yeast 2-Hybrid-method can only identify binary complexes.

In **affinity purification**, a protein of interest (bait) is tagged with a molecular label (dark route in the middle of the figure) to allow easy purification.

The tagged protein is then co-purified together with its interacting partners (W–Z). This strategy can also be applied on

a genome scale





Gavin *et al. Nature* 415, 141 (2002) V 3 - 22

## TAP analysis of yeast PP complexes

Identify proteins by scanning yeast protein database for protein composed of fragments of suitable mass.

Here, the identified proteins are listed according to their localization (a). (b) lists the number of proteins per complex.





Gavin et al. Nature 415, 141 (2002)

# Validation of TAP methodology

![](_page_23_Figure_1.jpeg)

Check of the method:

can the same complex be obtained for different choices of the attachment point

(tag protein is attached to different components of complex)?

Yes, more or less (see gel in (a)).

Gavin et al. Nature 415, 141 (2002)

# Bioinformatics 3 – WS 14/15

### Pros and Cons of TAP-MS

### Advantages:

- quantitative determination of complex partners *in vivo* without prior knowledge
- simple method, high yield, high throughput

![](_page_24_Picture_5.jpeg)

### **Difficulties:**

- tag may prevent binding of the interaction partners
- tag may change (relative) expression levels
- tag may be **buried** between interaction partners
  → no binding to beads

![](_page_24_Picture_10.jpeg)

![](_page_24_Picture_11.jpeg)

### Yeast Two-Hybrid Screening

Discover binary protein-protein interactions via physical interaction

![](_page_25_Figure_2.jpeg)

complex of binding domain (BD) + activator domain (AD)

Disrupt BD-AD protein; fuse bait to BD, prey to AD

→ expression only when bait:prey-complex formed

Reporter gene may be fused to green fluorescent protein.

### Pros and Cons of Y2H

### Advantages:

- *in vivo* test for interactions
- cheap + robust  $\rightarrow$  large scale tests possible

#### **Problems:**

investigates the interaction between `
(i) overexpressed
(ii) fusion proteins in the
(iii) yeast
(iv) nucleus

 $\rightarrow$  many false positives (up to 50% errors)

spurious interactions via third protein /

# Synthetic Lethality

Apply two mutations that are viable on their own, but lethal when combined.

In cancer therapy, this effect implies that inhibiting one of these genes in a context where the other is defective should be selectively lethal to the tumor cells but not toxic to the normal cells, potentially leading to a large therapeutic window.

Gene X	Gene Y	
+	+	No effect
_	+	No effect
÷	_	No effect
_	_	Death

http://jco.ascopubs.org/

Synthetic lethality may point to:

- physical interaction (building blocks of a complex)
- both proteins belong to the same pathway
- both proteins have the same function (redundancy)

## **Gene Coexpression**

All constituents of a complex should be present at the same point in the cell cycle  $\rightarrow$  test for correlated expression

This is not a direct indication for complexes (there are too many co-regulated genes), but useful "filter"-criterion Standard tools: DNA micro arrays / WGS

DeRisi, Iyer, Brown, *Science* **278** (1997) 680:

Diauxic shift from fermentation (growth on sugar) to respiration (growth on ethanol) in *S. cerevisiae* 

→ Identify groups of genes with similar expression profiles

![](_page_28_Figure_6.jpeg)

# **DNA Microarrays**

Fluorescence labeled DNA (cDNA) applied to micro arrays

- → hybridization with complementary library strand
- → fluorescence indicates relative cDNA amounts

![](_page_29_Picture_4.jpeg)

A. Butte, Nature Reviews Drug Discovery 1, 951-960, 2002

![](_page_29_Figure_6.jpeg)

http://intmedweb.wfubmc.edu/

two labels (red + green) for experiment and control Usually: red = signal green = control  $\rightarrow$  yellow = "no change"

![](_page_30_Figure_0.jpeg)

Identify groups of genes with similar time courses = expression profiles  $\rightarrow$  "cause or correlation"? — biological significance?

DeRisi, Iyer, Brown, Science 278 (1997) 680

### **Interaction Databases**

#### Bioinformatics: make use of existing databases

3.2 Experimental High-Throughput Methods for Detecting Protein-Protein Interactions

Table 3.1 Some public databases compiling data related to protein interactions: (P) and (D) stand for proteins and domains (the number of interactions reflects the status of June 2007).

	URL	Number of interactions	Туре	Proteins /domains
MIPS	mips.gsf.de/genre/proj/mpact	4300	curated	
BIND	bond.unleashedinformatics.com	200000	curated	Р
MINT	160.80.34.4/mint/	103800	curated	Р
DIP	dip.doe-mbi.ucla.edu	56000	curated	Р
PDB	www.rcsb.org/pdb	800 complexes	curated	
HPRD	www.hprd.org	37500	curated	P, D
Scoppi	www.scoppi.org	102000	automatic	D
UniHI	theoderich.fb3.mdc-berlin. de:8080/unihi/home	209000	integrated data	Р
STRING	string.embl.de	interactions of 1500000 proteins	integrated data from genomic context, high-throughput experiments, coexpression, previous knowledge	Р
iPfam	www.sanger.ac.uk/Software/ Pfam/iPfam	3019	data extracted from PDB	D
YEAST protein complex database	yeast.cellzome.com	232 complexes	experimental	Р
ABC	service.bioinformatik. uni-saarland.de/abc	13000 complexes	semiautomatic	Р

# (low) Overlap of Results

For **yeast**: ~ 6000 proteins = ~18 million potential interactions rough estimates:  $\leq$  100000 interactions occur

- $\rightarrow$  1 true positive for 200 potential candidates = **0.5%**
- $\rightarrow$  decisive experiment must have accuracy << 0.5% false positives

### Different experiments detect different interactions

For yeast: 80000 interactions known, 2400 found in > 1 experiment

Problems with experiments:

- i) incomplete coverage
- ii) (many) false positives
- iii) selective to type of interaction and/or compartment

![](_page_32_Figure_10.jpeg)

## **Criteria for Reliability**

Guiding principles (incomplete list!):

#### 1) mRNA abundance:

most experimental techniques are biased towards high-abundance proteins

#### 2) compartments:

- most methods have their "preferred compartment"
- proteins from same compartment => more reliable

### 3) co-functionality

complexes have a functional reason (assumption !?)

### **In-Silico Prediction Methods**

### Sequence-based:

- gene clustering
- gene neighborhood
- Rosetta stone
- phylogenetic profiling
- coevolution

- "Work on the parts list"
- $\rightarrow$  fast
- $\rightarrow$  unspecific
- → high-throughput methods for pre-sorting

### Structure-based:

- interface propensities
- protein-protein docking
- spatial simulations

![](_page_34_Picture_15.jpeg)

"Work on the parts"

- $\rightarrow$  specific, detailed
- $\rightarrow$  expensive
- $\rightarrow$  accurate

# **Gene Clustering**

# Idea: functionally related proteins or parts of a complex are expressed simultaneously

![](_page_35_Figure_2.jpeg)

Search for genes with a **common promoter** 

 $\rightarrow$  when activated, all are transcribed together as one operon

### Example:

bioluminescence in *V. fischeri*, regulated via quorum sensing  $\rightarrow$  three proteins: I, AB, CDE

![](_page_35_Figure_7.jpeg)

# Gene Neighborhood

Hypothesis again: functionally related genes are expressed together

"functionally" = same {complex | pathway | function | ...}

![](_page_36_Figure_3.jpeg)

 $\rightarrow$  Search for **similar sequences** of genes in **different organisms** 

(<=> Gene clustering: one species, promoters)

### **Rosetta Stone Method**

![](_page_37_Picture_1.jpeg)

**Idea**: find homologous genes ("words") in genomes of different organisms ("texts")

- check if fused gene pair exists in one organism
- $\rightarrow$  May indicate that these 2 proteins form a complex

![](_page_37_Figure_5.jpeg)

Multi-lingual stele from 196 BC, found by the French in 1799  $\rightarrow$  key to deciphering hieroglyphs Enright, Ouzounis (2001): 40000 predicted pair-wise interactions from search across 23 species

# **Phylogenetic Profiling**

- Idea: either all or none of the proteins of a complex should be present in an organism
- → compare presence of protein homologs across species (e.g., via sequence alignment)

![](_page_38_Figure_3.jpeg)

### Distances

![](_page_39_Figure_1.jpeg)

Hamming distance between species: number of different protein occurrences

![](_page_39_Figure_3.jpeg)

Two pairs with similar occurrence: P2-P7 and P3-P6

### Coevolution

Idea: not only similar static occurence, but similar dynamic evolution

![](_page_40_Figure_2.jpeg)

Interfaces of complexes are often better conserved than the rest of the protein surfaces.

Also: look for potential substitutes  $\rightarrow$  anti-correlated

- $\rightarrow$  missing components of pathways
  - $\rightarrow$  function prediction across species
    - $\rightarrow$  novel interactions

## i2h method

Schematic representation of the i2h method.

A: Family alignments are collected for two different proteins, 1 and 2, including corresponding sequences from different species (a, b, c, ).

B: A virtual alignment is constructed, concatenating the sequences of the probable orthologous sequences of the two proteins. Correlated mutations are calculated.

![](_page_41_Figure_4.jpeg)

Pazos, Valencia, Proteins 47, 219 (2002)

### Correlated mutations at interface

Correlated mutations evaluate the similarity in variation patterns between positions in a multiple sequence alignment.

Similarity of those variation patterns is thought to be related to compensatory mutations.

Calculate for each positions *i* and *j* in the sequence a rank correlation coefficient  $(r_{ij})$ :  $\sum (S_{ikl} - \overline{S}_i)(S_{ikl} - \overline{S}_i)$ 

$$\Gamma_{ij} = \frac{\sum_{k,l} (\sigma_{ikl} - \sigma_{l}) (\sigma_{jkl} - \sigma_{l})}{\sqrt{\sum_{k,l} (S_{ikl} - \overline{S}_{i})^{2}} \sqrt{\sum_{k,l} (S_{jkl} - \overline{S}_{j})^{2}}}$$

where the summations run over every possible pair of proteins *k* and *l* in the multiple sequence alignment.

 $S_{ikl}$  is the ranked similarity between residue *i* in protein *k* and residue *i* in protein *I*.  $S_{ikl}$  is the same for residue *j*.

 $S_i$  and  $S_j$  are the means of  $S_{ikl}$  and  $S_{jkl}$ .

Pazos, Valencia, Proteins 47, 219 (2002)

## Summary

What you learned today: how to get some data on PP interactions

![](_page_43_Figure_2.jpeg)

type of interaction? — reliability? — sensitivity? — coverage? — ...

Next lecture: Mon, Nov. 3, 2014

- combining weak indicators: Bayesian analysis
- identifying communities in networks