# Changiz Eslahchi

AmirKabir
University of Technology

Computational Biology
Research Center
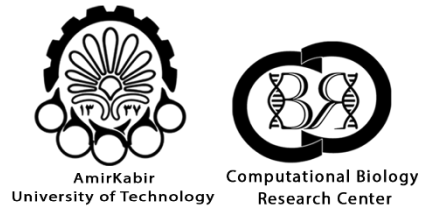
Ph.D. in Mathematics, Sharif University of Technology, 1998

Professor in Department of Computer Sciences, Shahid Beheshti University, Tehran

Senior Researcher in School of Biological Sciences, Institute for Research in Fundamental Sciences (IPM)

Honors…

# Comparison of Different Approaches for Identifying Subnetworks in Metabolic Networks

Changiz Eslahchi

Department of Computer Sciences, Shahid Beheshti University

School of Biological Sciences,  IPM Institute for Research in Fundamental Sciences
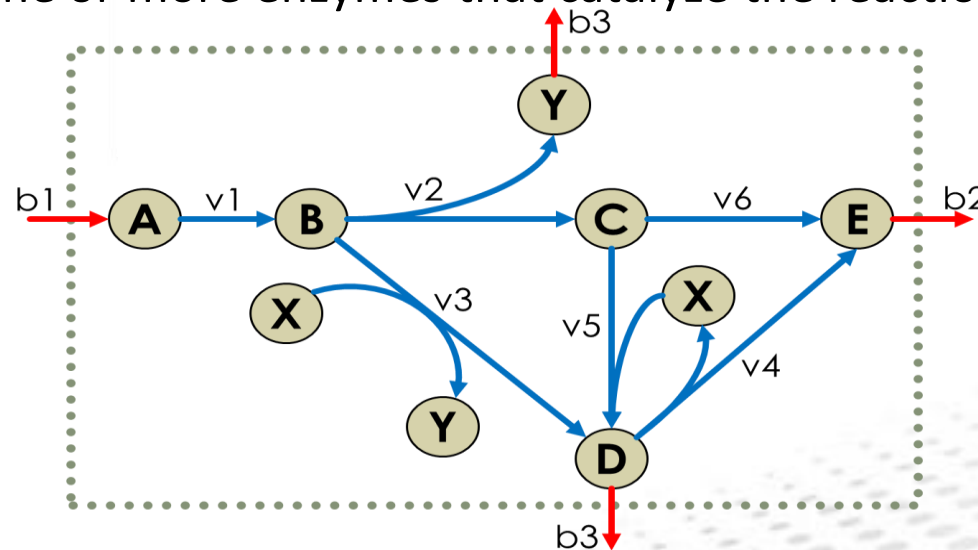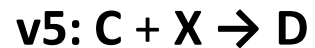
# Outline

- <span style="color:red">Structure of Metabolic of Networks</span>
- Decomposing Metabolic Network Models
- Comparison Framework
  - Definition of Criteria
- Comparison Results
- Discussion
  - Verifying Ranking Stability
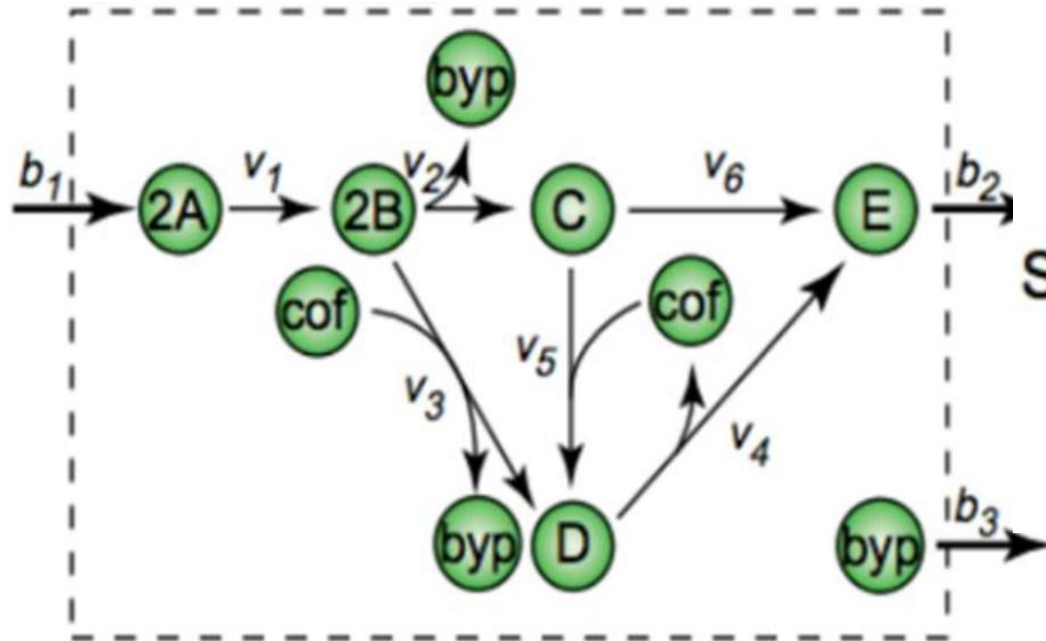  - Evaluation of Subnetworks
- Future Work

1st International Computational Biology workshop

# What is a Metabolic Network?

- The biochemical "engine" of the cell
  - Converts raw materials into energy and polymer building blocks
  - Makes survival, growth, and reproduction feasible

- Consists of metabolites (bio-molecules) and reactions (that converts metabolites)
  - Reactions may be reversible or irreversible (thermodynamic constraints)
  - May be associated with one or more enzymes that catalyze the reaction

**v1: A → B**
**v2: B → C + Y**
**v3: B + X → Y + D**
**v4: D → X + E**
**v5: C + X → D**
**v6: C → E**

# Metabolic Network Modeling

**Hypergraph**

**Stoichiometric Matrix**

# Definition of Flux and Flux Distribution

- Flux of a reaction: the rate at which the reaction works

- Flux distribution: for a network with N reactions, any N-tuple which specifies the flux of each reaction

E.g. **V=(3, 2, 0, 2, 2, 0)** is a flux distribution which means:
- v1 works with rate 3
- v2 works with rate 2
- …

With such a flux distribution, B is gradually increased over time, but abundance of C does not change over time



**v1: A → B**
**v2: B → C + Y**
**v3: B + X → Y + D**
**v4: D → X + E**
**v5: C + X → D**
**v6: C → E**

# Steady State Analysis

■ **Steady-state:**

- No changes in metabolite concentrations

- Metabolite production and consumption rates are equal

- It is shown that cell is in steady state in normal condition

$$\frac{d\overline{m}}{dt} = S \cdot \overline{v} = 0$$

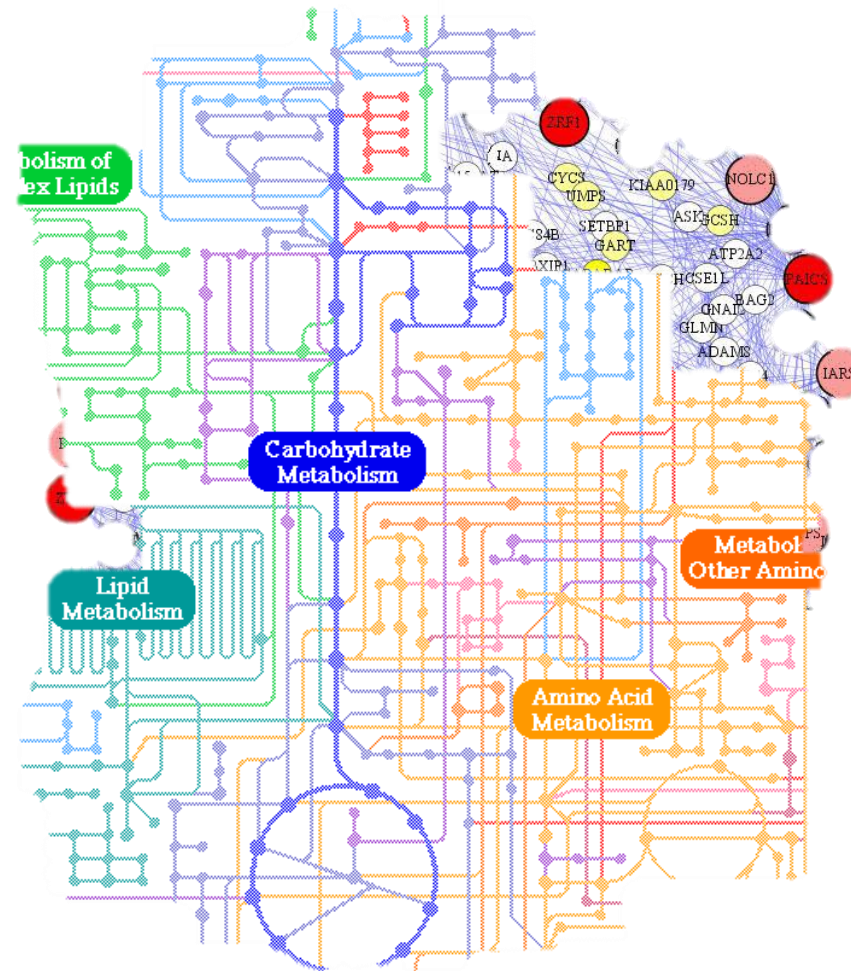|  | $R_1$ | $R_2$ | $R_3$ | $R_4$ | $R_5$ | $R_6$ | $R_7$ | $R_8$ | $R_9$ | $R_{10}$ | $V_m$ | $V_{growth}$ | $A_{up}$ | $D_{up}$ | $F_{up}$ | $H_{up}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **A** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **B** | 1 | −1 | 0 | 0 | −1 | 0 | 0 | −1 | 0 | 0 | 0 | −1 | 0 | 0 | 0 | 0 |
| **C** | 0 | 2 | −1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **D** | 0 | 0 | 1 | −1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **E** | 0 | 0 | 0 | 0 | 1 | −1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **F** | 0 | 0 | 0 | 0 | 0 | 1 | −1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **G** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | −1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **H** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | −1 | 0 | −2 | 0 | 0 | 0 | 0 |
| **I** | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | −1 | 0 | 0 | 0 | 0 | 0 |
| **A**$_{external}$ | −1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| **D**$_{external}$ | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| **F**$_{external}$ | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| **H**$_{external}$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

$$
\begin{bmatrix} R_1 \\ R_2 \\ R_3 \\ R_4 \\ R_5 \\ R_6 \\ R_7 \\ R_8 \\ R_9 \\ R_{10} \\ V_m \\ V_{growth} \\ A_{up} \\ D_{up} \\ F_{up} \\ H_{up} \end{bmatrix}
=
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}
$$

■**m**: *metabolite concentrations vector (mol/mg)*
■**S**: *stoichiometric matrix*
■**v**: *reaction rates vector*

# Outline

- Structure of Metabolic of Networks

- <span style="color:red">Decomposing Metabolic Network Models</span>

- Comparison Framework
  - Definition of Criteria

- Comparison Results

- Discussion
  - Verifying Ranking Stability
  - Comparing Subnetworks with KEGG

- Future Work

1st International Computational Biology workshop

# Decomposition Facilitates Analysis

# Metabolite-based decomposition     vs.     Reaction-based decomposition

# Metabolic Network Decomposition History

- Jeong (2000): 43 metabolic networks are analyzed and suggested that these networks have small-world structure properties
  - Power-law distribution
  - High cluster coefficient
  - Short network diameter

- Schilling and Palsson (2000): Defined several **manual** instructions for properly decomposing networks
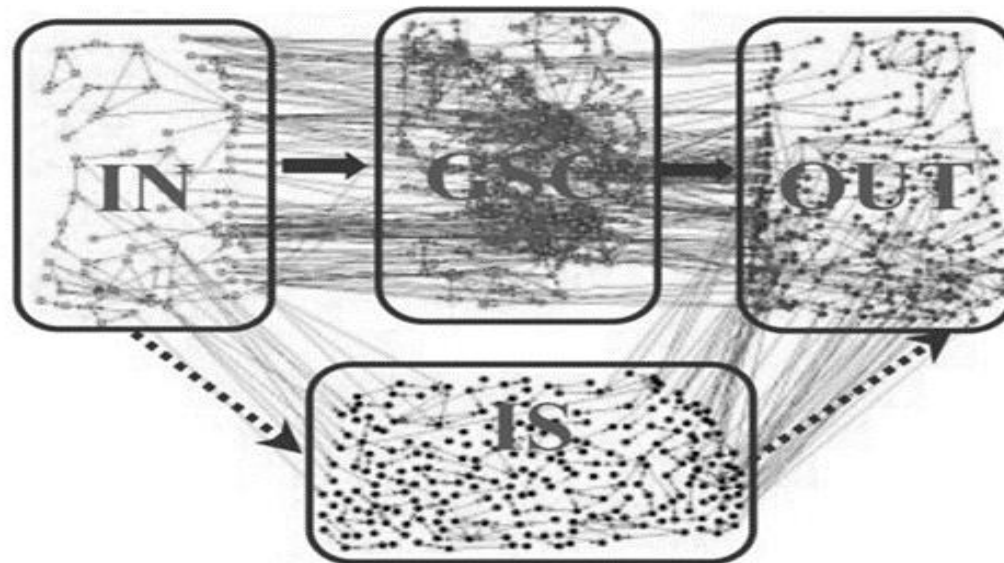
# Metabolic Network Decomposition Methods

- Schuster (2002): Partitioning by removing the "hub" metabolites of the network Internal metabolite would be subnetworks
  - Hubs: high connectivity degree metabolites

- Holme (2003): Partitioning by removing "central" metabolites
  - Based on betweenness centrality
  - Iterative removal produces a hierarchical decomposition of the network

# Metabolic Network Decomposition Methods

- Ma (2004): Decomposition into predefined "**bow-tie**" structure
  - IN (input), GSC (core), OUT (output), and IS (isolated) components
  - Simple hierarchical clustering of **reactions** (instead of metabolites) in GSC component based on shortest-path distance

# Metabolic Network Decomposition Methods

- Guimera (2005): Finding modules by maximizing "**modularity**" (community detection)
  - Uses Simulated Annealing to maximize modularity
  - The main goal of the method is to assign biological roles to each metabolite based on its position in its subnetwork
- Newman (2006): Finding modules by maximizing "modularity"
  - Using spectral graph partitioning
- Yoon (2007): Adding edge (reaction) weights to hypergraph representation and then removing central metabolites
  - Define edge weights based on **reaction flux data**
  - Suggests that functional organization of a metabolic network differs in different physiological conditions

# Metabolic Network Decomposition Methods

- Poolman (2007): Defining distance between reaction based on "**correlation between reaction flux values**" in "**steady-state**"
  - Defines "reaction correlation coefficient" which is a measure of "correlation between reaction flux values"
  - Reaction correlation coefficient is computed directly using stoichiometric matrix representation of the network
- Verwoerd (2011): Extending Schuster method by redefining "hub" metabolites
  - Defined a **global connection degree** based on random walks on the network (similar to MCL inflammation step)
  - A method similar to Schuster method is applied based on this global connection degree
  - Interactive software which allows complete user adjustments in the process of decomposition

# Metabolic Network Decomposition Methods

- Sridharan (2011): Finding communities based on maximizing "retroactive interactions" (cycle) inside subnetworks
  - "Modularity" is redefined so that the number of cycles is maximized instead of number of edges
  - Recursively divides network into two subnetworks which produces a hierarchical decomposition of the network
- Muller (2014): Finding modules based on linear algebra
  - "Module-finding" rather than "decomposition" method

# Summary of Implemented Methods

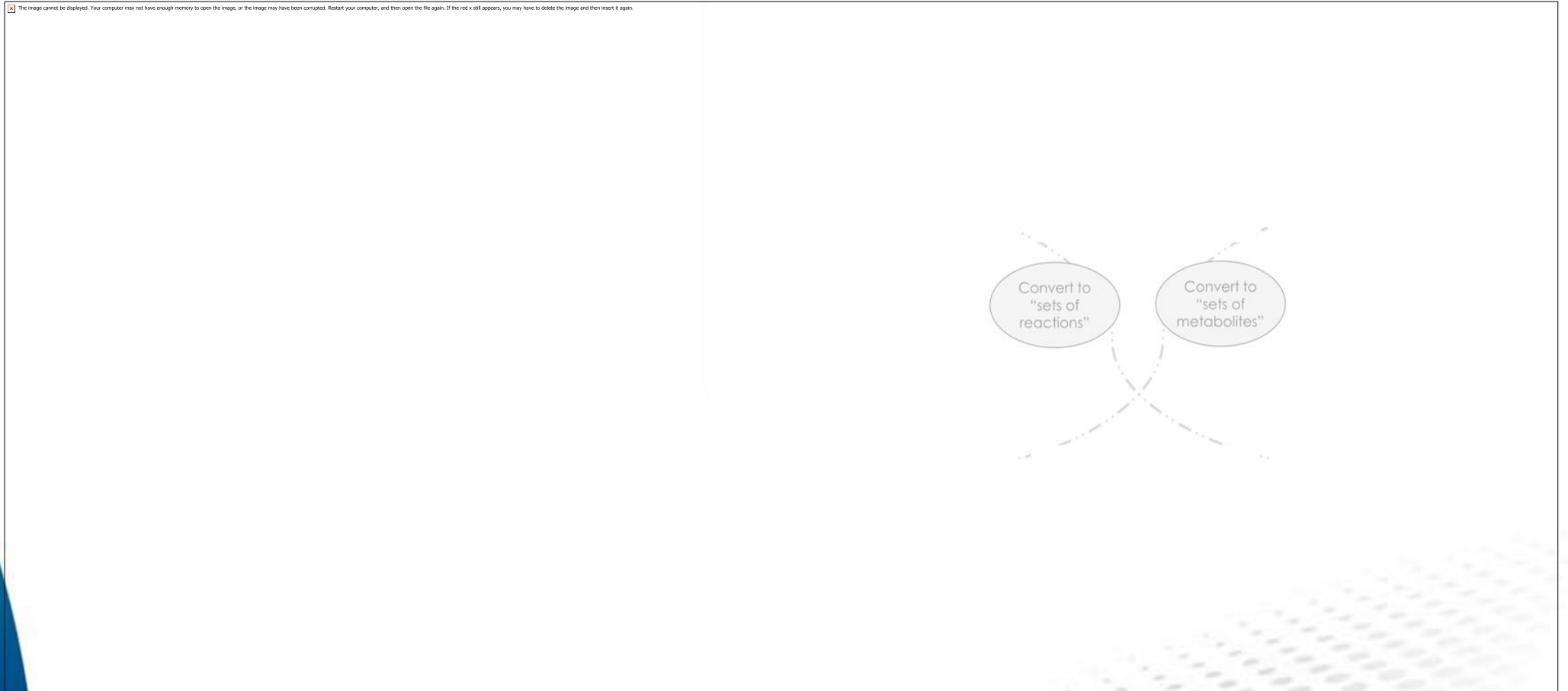| Method | Output Subnetwork | Module Finding vs. Decomposition | Hierarchical Output |
|---|---|---|---|
| Schuster et al. (2002) | Sets of metabolites | Decomposition | No |
| Newman (2006) | Sets of metabolites | Decomposition | No |
| Guimera and Amaral (2005) | Sets of metabolites | Decomposition | No |
| Holme et al. (2003) | Sets of metabolites | Decomposition | Yes |
| Verwoerd (2011) | Sets of metabolites | Decomposition | Yes |
| Poolman et al. (2007) | Sets of reactions | Decomposition | Yes |
| Sridharan et al. (2011) | Sets of reactions | Decomposition | Yes |
| Muller (2014) | Sets of reactions | Module finding | No |

# Outline

- Structure of Metabolic of Networks

- Decomposing Metabolic Network Models

- <span style="color:red">Comparison Framework</span>
  - Definition of Criteria

- Comparison Results

- Discussion
  - Verifying Ranking Stability
  - Comparing Subnetworks with KEGG

- Future Work

# The Comparison Framework

The image cannot be displayed. Your computer may not have enough memory to open the image, or the image may have been corrupted. Restart your computer, and then open the file again. If the red x still appears, you may have to delete the image and then insert it again.

Convert to "sets of reactions"

Convert to "sets of metabolites"

# Outline

- Structure of Metabolic of Networks

- Decomposing Metabolic Network Models

- Comparison Framework
  - <span style="color:red">Definition of Criteria</span>

- Comparison Results

- Discussion
  - Verifying Ranking Stability
  - Comparing Subnetworks with KEGG

- Future Work

# Criteria: Modularity

$$M = \sum_{i=1}^{K} \left[ \frac{l_i}{L} - \left( \frac{d_i}{2L} \right)^2 \right]$$

- Decomposition of a network into $K$ subnetworks
- $L$ is the total number of edges in the network
- $l_i$ is the number of edges connecting nodes in subnetwork $i$
- $d_i$ is the sum of degrees of the nodes in subnetwork $i$
- Proposed by Newman (2006)
- Can be applied to both metabolite-based and reaction-based methods

# Criteria: Modularity

- Zero expected value for both:
  - Random decompositions
  - Trivial decomposition: the whole network as the only subnetwork

- Approximates the following value:

$$\#\left[\begin{array}{c}\textbf{Edges within}\\\textbf{subnetworks}\end{array}\right] - E\left[\begin{array}{c}\textbf{Number of such edges in}\\\textbf{randomized decomposition}\\\textbf{of the network}\end{array}\right]$$

# Criteria: GO Similarity (for reaction-based methods)

- Gene Ontology is a valuable source of information about:
  - Functions of gene products (molecular function)
  - Locations and sublocations of gene products (cellular compartment)
  - Processes which gene products involve (biological process)

- We define three different scores based on Resnick "semantic similarity" between genes in Gene Ontology
  - GO molecular function
  - GO cellular compartment
  - GO biological process

# Criteria: GO Similarity (for reaction-based methods)

- For a given decomposition **D**, "GO similarity score" is defined as:

$$GoScore(D) = \sum_i \left( \frac{ModSim(m_i, m_i)}{|m_i|^2} - \sum_k \frac{ModSim(m_i, m_k)}{|m_i||m_k|} ; \quad k \neq i \right)$$

- $GoScore$ is a measure of relatedness of reactions in each subnetwork and their distance from reactions in other subnetworks

- $ModSim$ denotes the similarity between modules. For a pair of modules $m_u$ and $m_v$, it is defined as:

$$ModSim(m_u, m_v) = \sum_{r \in m_u} \sum_{s \in m_v} RxnSim(r, s); \quad r \neq s$$

# Criteria: GO Similarity (for reaction-based methods)

- *RxnSim* denotes the similarity between reactions. For a given pair of reactions $r_i$ and $r_j$, it is defined as:

$$RxnSim(r_i, r_j) = \frac{\sum_{e \in G_i} \sum_{f \in G_j} SS(e,f)}{|G_i| \times |G_j|}$$

  - $G_i$ is the set of all genes associated with enzymes that catalyze reaction $r_i$
  - $SS(e,f)$ is the Resnik similarity of genes associated with genes $\boldsymbol{e}$ and $\boldsymbol{f}$

helico_iit341 rmods go_distance_mf_F

# Criteria: Module Coupling (for reaction-based methods)

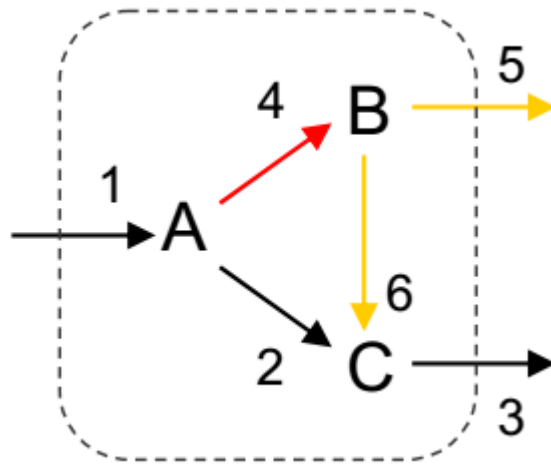- What is a Flux coupling relation!?
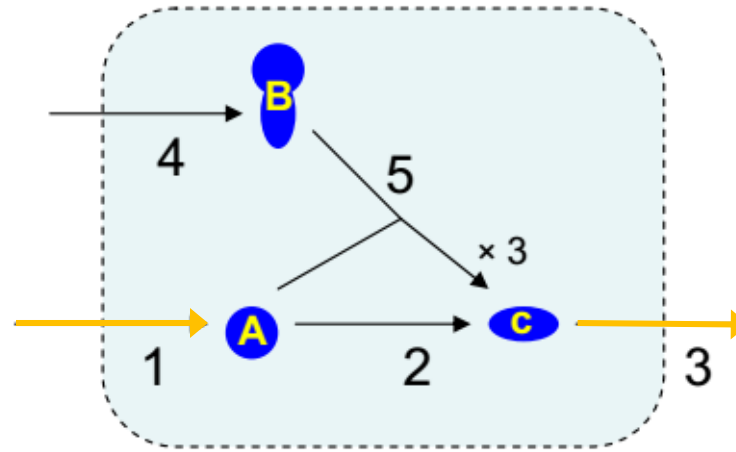
# Flux Coupling Relation

- Flux coupling represent how metabolic reactions cooperate

- Formal definition ($V_i$ denotes flux of reaction $r_i$ )
  - Fully coupled
    - $V_1 = c \, V_2 \; (c > 0)$
  - Partially coupled
    - $V_1 \neq 0 \leftrightarrow V_2 \neq 0$
  - Directionally coupled
    - $V_1 \neq 0 \rightarrow V_2 \neq 0$
  - Uncoupled

- Computing the set of flux coupling relations in a whole-genome network is fast (minutes)
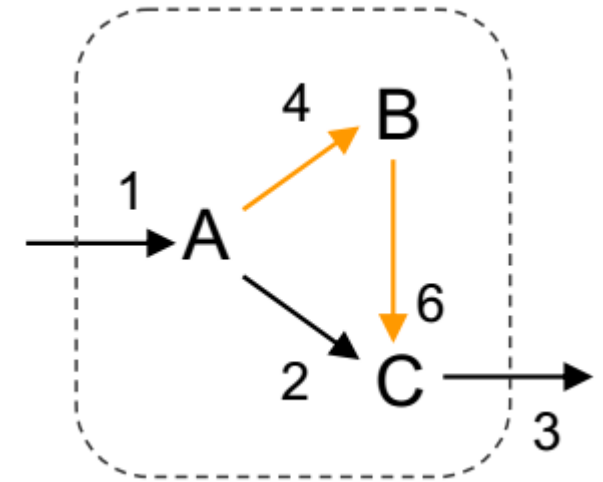
# Flux Coupling Relation Examples



Reactions 5 and 4
are directionally coupled
(Also 6 and 4)

Reactions 1 and 3
are partially coupled

Reactions 4 and 6
are fully coupled

# Criteria: Module Coupling (for reaction-based methods)

- $CM = [cm_{ij}]$: Reaction coupling matrix where $cm_{ij}$ denotes the type of coupling between reaction pair $r_i$ and $r_j$

- $SCM = [scm_{ij}]$: Simple coupling matrix where $scm_{ij} = \begin{cases} 1, & cm_{ij} \text{ is coupled} \\ 0, & otherwise \end{cases}$

- Based on simple reaction coupling matrix, we define "Module coupling score"

# Criteria: Module Coupling (for reaction-based methods)

- For a given decomposition $D$, Module coupling score, $McScore(D)$, is defined as:

$$\sum_i \left( Couplings(m_i, m_i) - Uncouplings(m_i, m_i) + \sum_j Uncouplings(m_i, m_j); j \neq i \right)$$

- For a pair of subnetworks $m_u$ and $m_v$

  - Number of coupling between two subnetworks:

  $$Coupulings(m_u, m_v) = \sum_{r \in m_u} \sum_{s \in m_v} scm_{rs}$$

  - Number of uncoupling between two subnetworks:

  $$Uncouplings(m_u, m_v) = \sum_{r \in m_u} \sum_{s \in m_v} (1 - scm_{rs})$$

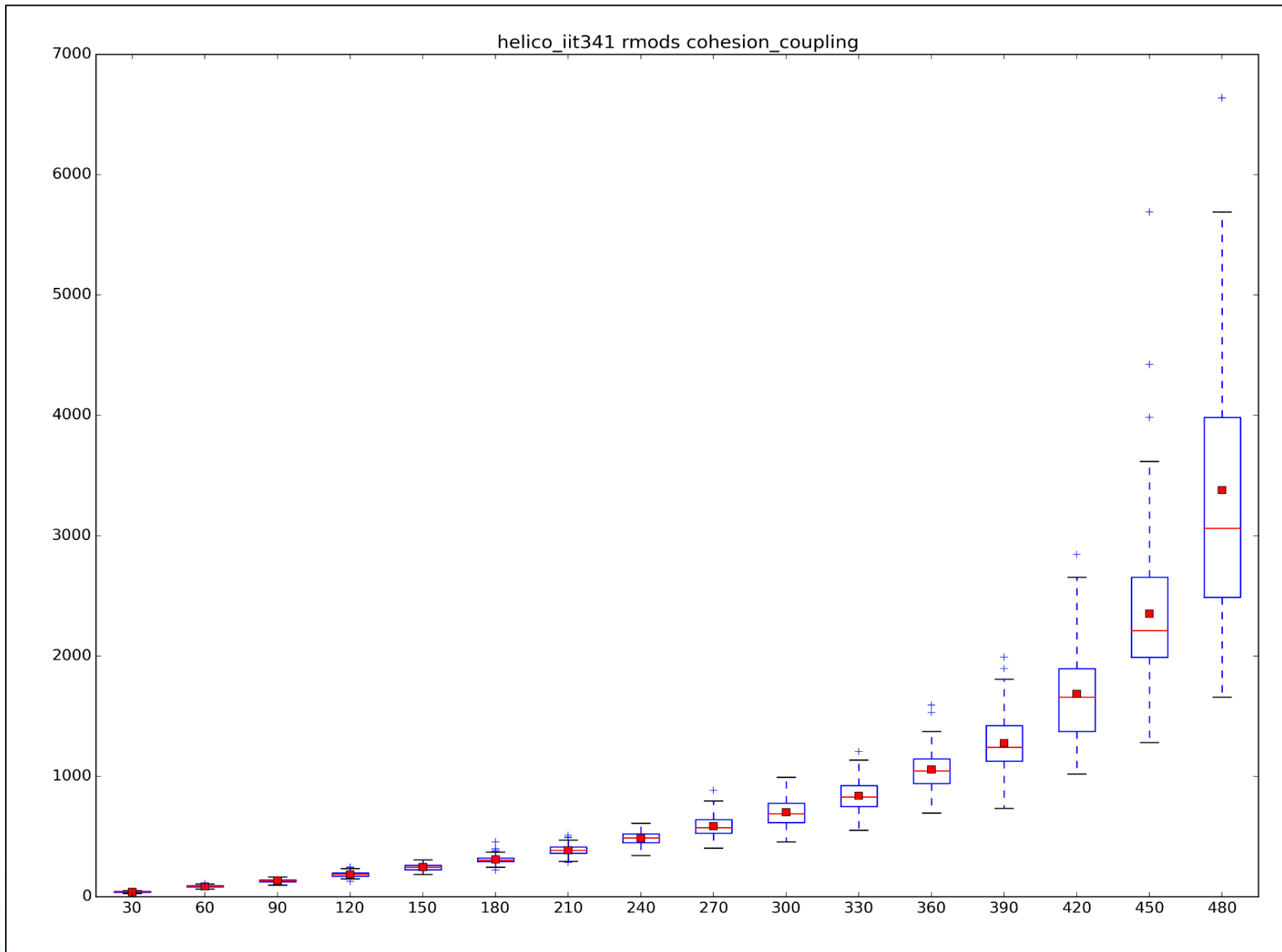1st International Computational Biology workshop

# Criteria: Module Coupling (for reaction-based methods)

- As with GO Similarity, $McScore(D)$ depends on the number of subnetworks

- The same procedure as GO similarity score is executed and "module coupling score" will be:

   the **p-value** of $McScore(D)$ against $McScore$ values for random samples

# Criteria: Efficacy (for metabolite- and reaction- based methods)

- Proposed by Verwoerd (2011)

- It is a measure of how much:
  - Sizes of subnetworks are balanced
  - The number of subnetworks is far from trivial (1 or N)

- Evaluates to zero (or small negative values) for trivial decomposition

# Criteria: Efficacy (for metabolite- and reaction- based methods)

- Assumes $f(n)$ as "the effort needed to analyze a network" of size **n**
- Efficacy

$$E = 100 \frac{Log[f(N)] - Log[f(k) + 1/k \sum_{i=1}^{k} f(n_i)]}{Log[f(N)] - Log[2f(\sqrt{N})]}$$

- $E_{max}$: for decompositions with $\sqrt{N}$ subnetworks of size $\sqrt{N}$
- The general behavior does not change dramatically with the choice of $f(N)$
  - A suggested choice for metabolic networks: $f(N) = \alpha N^p$ with $p = 0.25\sqrt{N}$

# Outline

- Structure of Metabolic of Networks

- Decomposing Metabolic Network Models

- Comparison Framework
    - Definition of Criteria

- Comparison Results

- Discussion
    - Verifying Ranking Stability
    - Comparing Subnetworks with KEGG

- Future Work

# Evaluated Datasets

- Model organisms from different domains of life
    - *Methanosarcina barkeri (Archaea, 628 mets, 690 rxns)*
    - *Helicobacter pylori (Small bacteria, 485 mets, 554 rxns)*
    - *Escherichia coli* (Bacteria, 1668 mets, 2382 rxns)
    - *Arabidopsis thaliana (Plant, 1913 mets, 1576 rxns)*
    - *Saccharomyces cerevisiae (Yeast, 1059 mets, 1266 rxns)*
    - *Mus musculus (Eukaryote, 2775 mets, 3726 rxns)*

# High Ranking Metabolites-based Methods

|  | **Modularity** | **Efficacy** |
| --- | --- | --- |
| *H. pylori* | Guimera & Amaral | Verwoerd |
| *M. barkeri* | Guimera & Amaral | Verwoerd |
| *S. cerevisiae* | Guimera & Amaral | Verwoerd |
| *A. Thaliana* | Guimera & Amaral | Verwoerd |
| *E. coli* | Guimera & Amaral | Verwoerd |
| *M. musculus* | Verwoerd | Verwoerd |

# High Ranking Reaction-based Methods

| | Efficacy | Module coupling | GO similarity molecular function | GO similarity biological process | GO similarity cell compartment |
|---|---|---|---|---|---|
| *H. pylori* | Poolman *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* | - |
| *M. barkeri* | Muller & Bockmayr Poolman *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* | - |
| *S. cerevisiae* | Poolman *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* |
| *A. thaliana* | Sridharan *et al.* | Sridharan *et al.* | Poolman *et al.* | - | - |
| *E. coli* | Poolman *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* | - |
| *M. musculus* | Poolman *et al.* | Poolman *et al.* | Sridharan *et al.* | Sridharan *et al.* | Sridharan *et al.* |

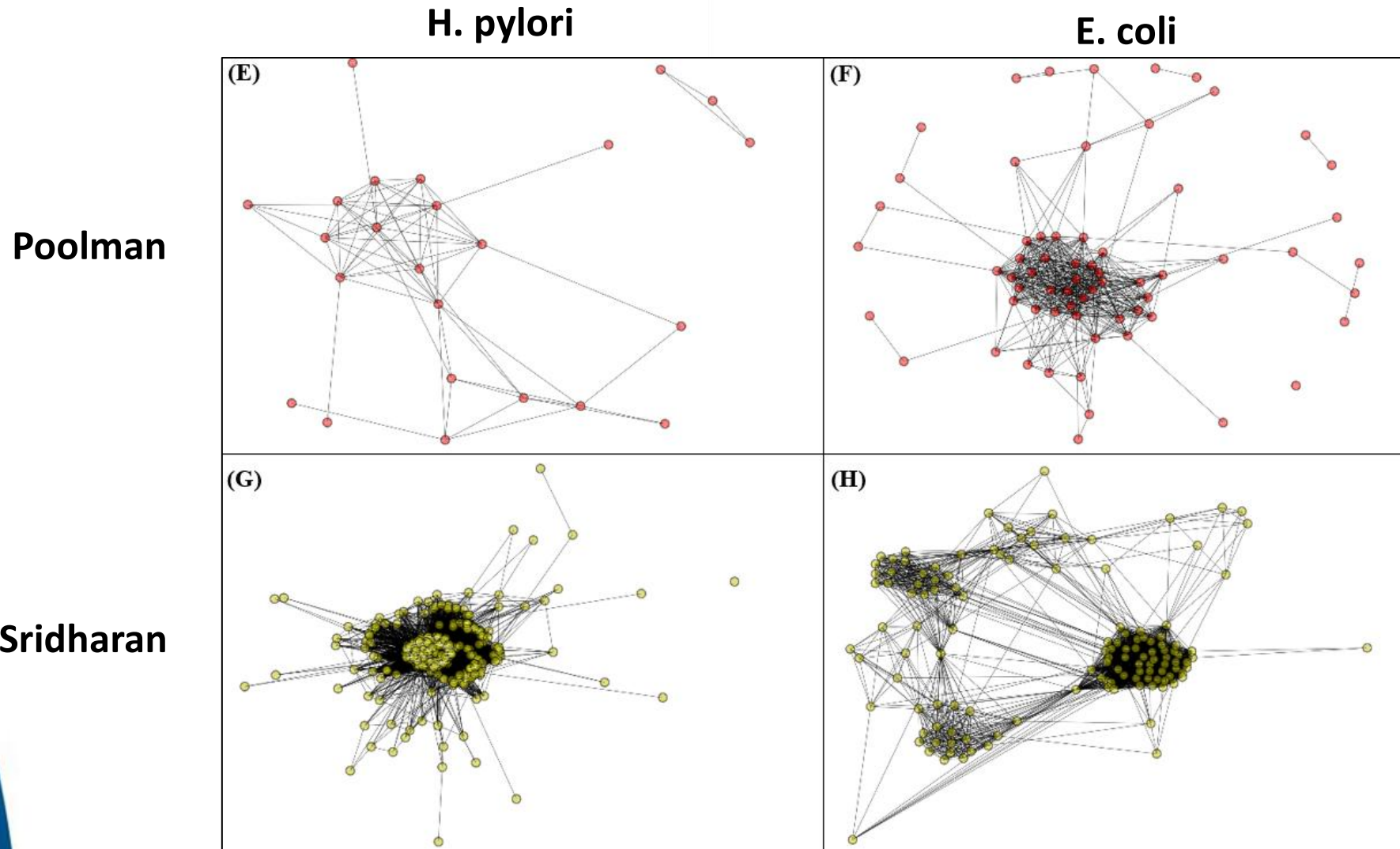# Sample Subnetworks (metabolite-based methods)

**H. pylori**

**E. coli**

**Guimera**

**Verwoerd**

# Sample Subnetworks (reaction-based methods)

**H. pylori**

**E. coli**

**Poolman**

**Sridharan**

# Outline

- Structure of Metabolic of Networks

- Decomposing Metabolic Network Models

- Comparison Framework
  - Definition of Criteria

- Comparison Results

- Discussion
  - Verifying Ranking Stability
  - Comparing Subnetworks with KEGG
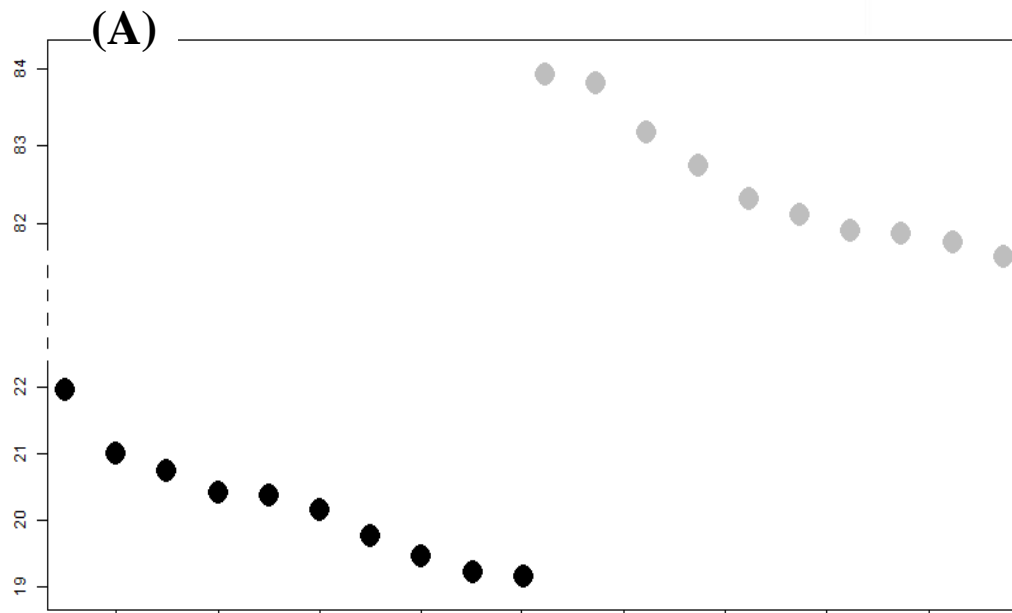
- Future Work

# Verifying Ranking Stability

- GO similarity and module coupling scores are based on *p-value*s

- **Question:** May a different set of random samples (as null distribution) affect the ranking of the methods?

- **To Answer:** An approach similar to *k*-fold cross-validation
  - Randomly divide the set of random samples into 10 equally-sized parts
  - Remove one part at a time ➔ a new set of random samples is generated
  - A p-value score is computed based on this new set
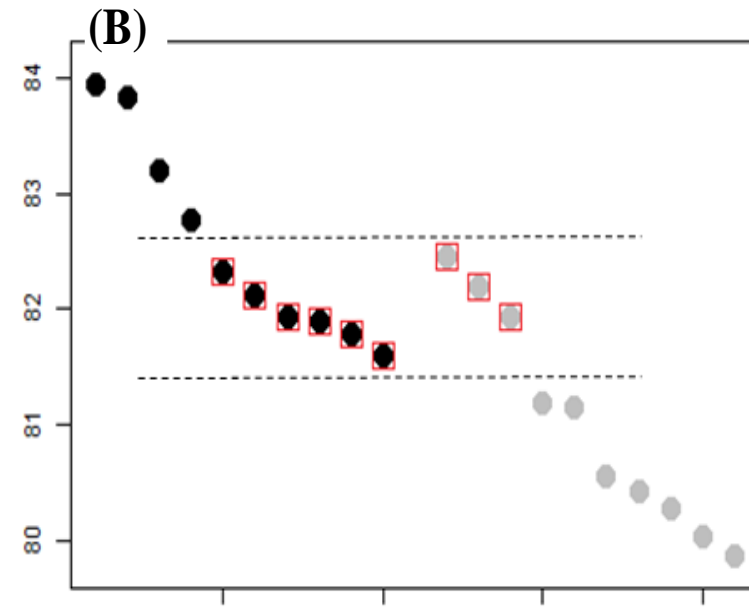  - 10 different p-value scores are computed for each original score

# Verifying Ranking Stability

- A given ranked pair may be either stable or unstable



p-value scores for a stable ranked pair

An unstable ranked pair

$$\text{Unstable Percent} = \frac{\#[\ \blacksquare\ \square\ ]}{\#[\text{all points}]}\ \%$$

# Verifying Ranking Stability

- Stability of all pairs in all ranking are checked
- List of all found unstable ranked pairs:

|  | Criterion | Unstable pairs | Unstable Percents |
|---|---|---|---|
| *E. coli* | Module Coupling | Poolman > Sridharan | 35% |
| *M. Musculus* | GO (biological process) | Poolman > Sridharan | 25% |

# Outline

- Structure of Metabolic of Networks

- Decomposing Metabolic Network Models

- Comparison Framework
  - Definition of Criteria

- Comparison Results

- Discussion
  - Verifying Ranking Stability
  - Comparing Subnetworks with KEGG

- Future Work

# Comparing Subnetworks and KEGG Pathways

- KEGG categorizes its metabolic pathways in 11 different major pathways

- We merge several random metabolic pathways (2 to 5 pathways) to create artificial networks with known modules (each pathways as one module)

- Apply methods to the aritifical networks and check how successful are they in detecting original metabolic pathways
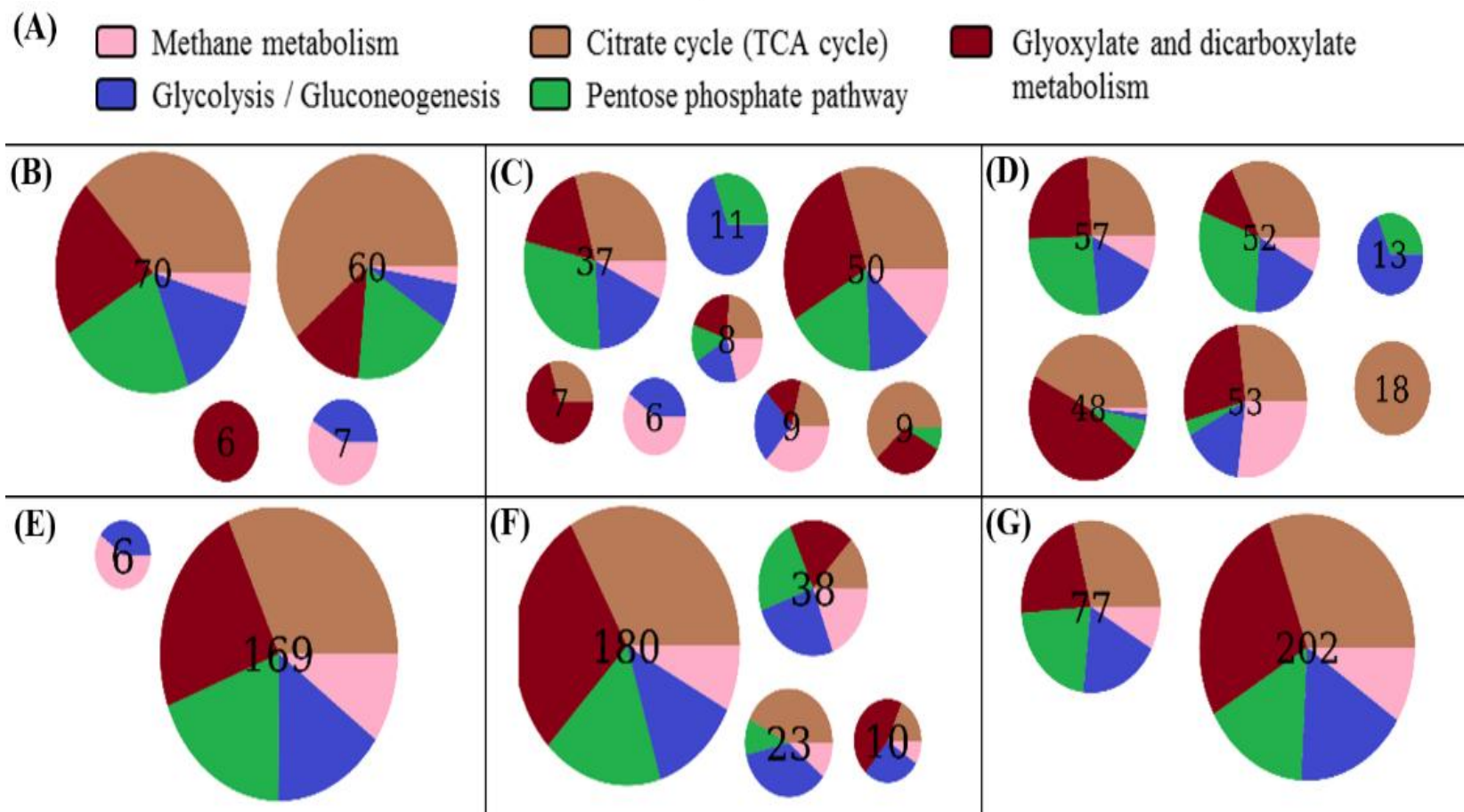
$$AS(C, C') = \frac{\sum_{i \in \mathcal{M}(C)} \mathcal{S}(c_i, \ C') \ - \ \mathcal{S}(\mathcal{P}(C), C')}{|\mathcal{M}(C)| + 1}$$

$$\mathcal{S}(c_i, C') = max_j \frac{|c_i \cap c'_j|}{|c_i|}$$

|  | Agreement Score | | | | Number of Networks |
|---|---|---|---|---|---|
| | **2** | **3** | **4** | **5** | |
| **Schuster *et al.*** | 0.35 | 0.46 | 0.55 | 0.59 | 94 |
| **Newman** | 0.57 | 0.56 | 0.60 | 0.62 | 100 |
| **Guimera & Amaral** | 0.81 | 0.76 | 0.69 | 0.64 | 100 |
| **Holme *et al.*** | 0.23 | 0.25 | 0.13 | 0.22 | 100 |
| **Verwoerd** | 0.50 | 0.51 | 0.48 | 0.45 | 90 |
| **Poolman *et al.*** | 0.40 | 0.40 | 0.39 | 0.42 | 100 |
| **Sridharan *et al.*** | 0.62 | 0.48 | 0.41 | 0.37 | 97 |

# Future Work

- Publicly available software package!

- More rigorous checking against KEGG

- Adding new criteria
  - Agreement of subnetworks with KEGG pathways
  - Co-expression of enzymes related to reactions in each subnetworks
  - Semantic similarity for metabolite-subnetworks based on ChEBI ontology

- Thorough investigation on the types of modules created by each method
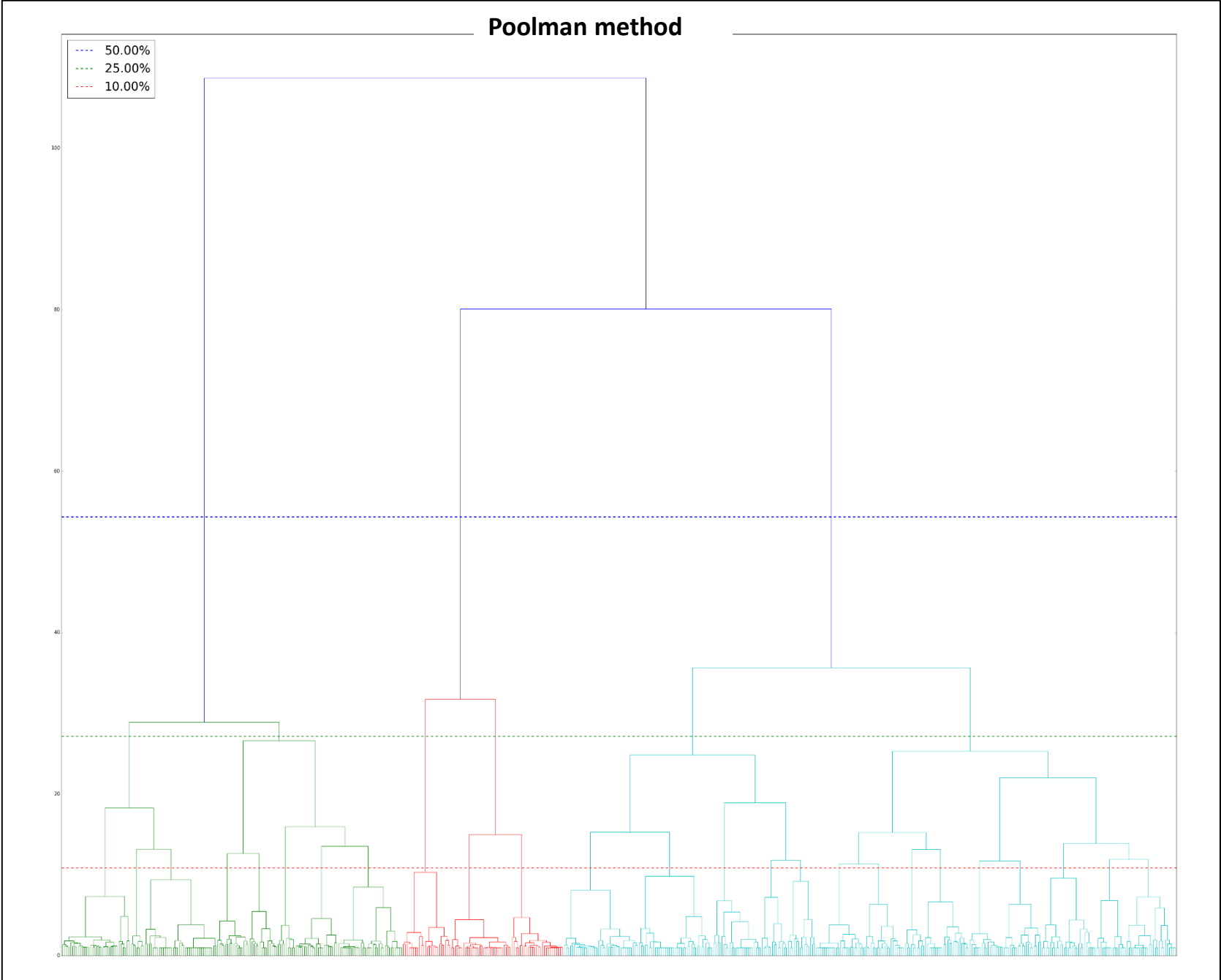
# Thanks and Questions

# Dealing with Methods with Hierarchical Output

- Holme, Poolman, and Sridharan methods produce hierarchical decompositions
- Cutting dendrograms at different levels produce different decompositions
- We have chosen several cut-thresholds for each hierarchical method manually
  - At the top, middle, and bottom of the dendrogram
  - E.g. Poolman dendrogram cut at 10%, 25%, and 50% height of dendrogram

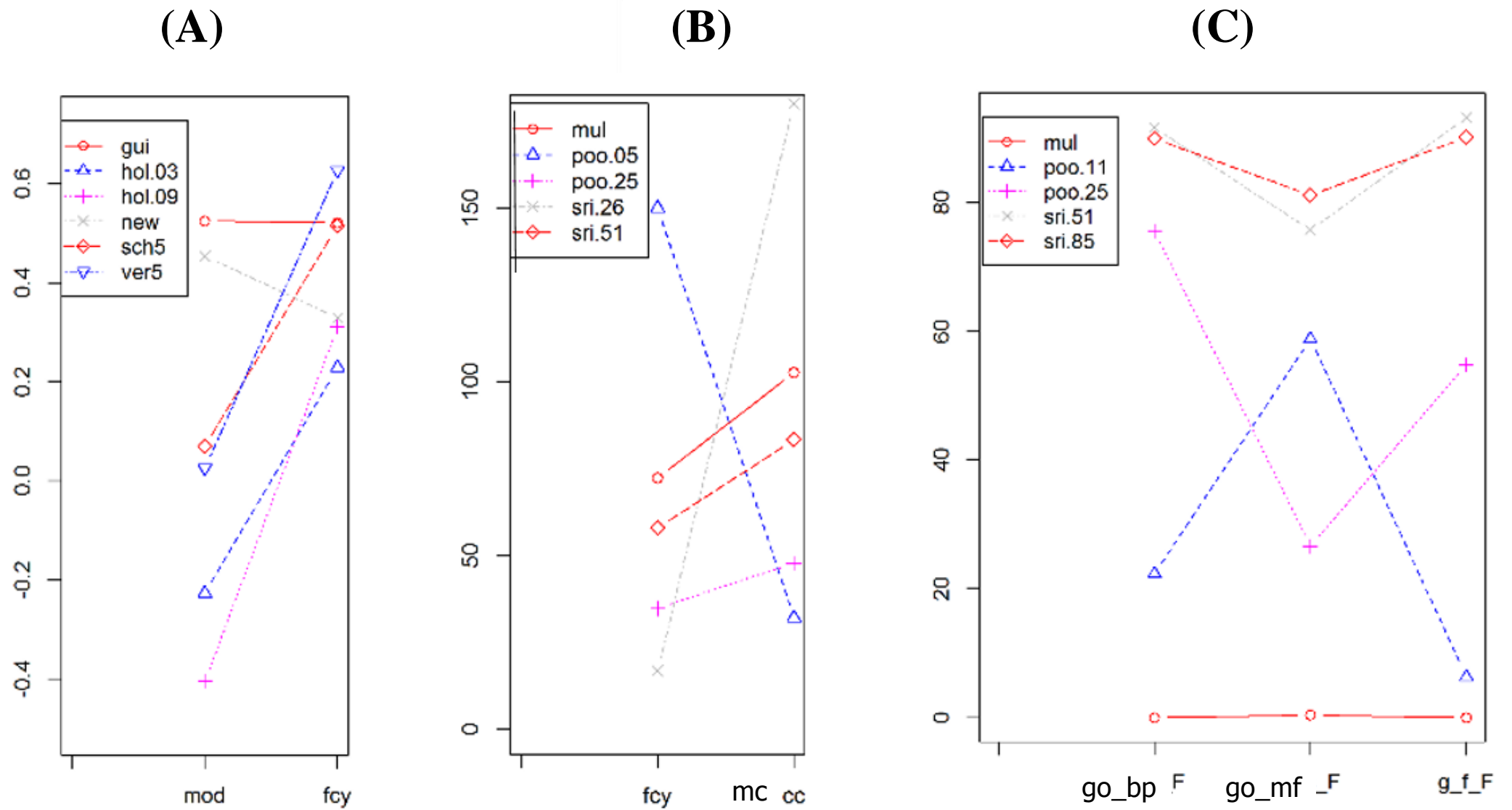# Dealing with Methods with Hierarchical Output

**Fig. 5.** Scores of methods in different criteria for *S. cerevisiae*. **(A)** metabolite-based methods;

# Verifying Ranking Stability

- Stability of all pairs in all ranking are checked
- List of all found unstable ranked pairs:

| | Criterion | Unstable pairs | UP | | Criterion | Unstable pairs | UP |
|---|---|---|---|---|---|---|---|
| *E.Coli* | Module Coupling | poo.5 > sri.51 | 35% | *S.cerevisiae* | GO (mf) | sri.51 > sri.85 | 10% |
| *A.Thaliana* | GO (mf) | poo.25 > poo.68 | 10% | | GO (bp) | sri.51 > sri.85 | 45% |
| | | poo.65 > poo.68 | 10% | *M.barkeri* | Module Coupling | sri.53 > sri.23 | 55% |
| | | poo.25 > poo.34 | 15% | | GO (mf) | sri.23 > sri.53 | 70% |
| | | poo.65 > poo.34 | 30% | *M.musculus* | GO (bp) | poo.15 > sri.85 | 25% |
| | | poo.25 > poo.65 | 90% | | GO (cc) | sri.26 > sri.51 | 40% |
| | | poo.34 > poo.68 | 95% | | GO (mf) | sri.26 > sri.51 | 15% |
| | | | | | | poo.36 > poo.15 | 95% |