

# Data-Driven Medicine

**Nataša Pržulj, PhD, MAE**

ICREA Research Professor  
Barcelona Supercomputing Center



**Barcelona  
Supercomputing  
Center**

*Centro Nacional de Supercomputación*

# Overview

**Medicine: complex world of inter-connected entities**

## **1. Motivation**

## **2. New Methods – Examples: mine inter-connected data**

- i. **Single type of omics data: Molecular networks** → function, disease
- ii. **Multiple layers of heterogeneous data:**
  - **Patient-centered data integration** → Precision medicine
  - Disease re-classification
  - Gene Ontology reconstruction
  - Network alignment

## **3. Vision**

# Overview

**Medicine: complex world of inter-connected entities**

## **1. Motivation**

## **2. New Methods – Examples: mine inter-connected data**

i. Single type of omics data: Molecular networks → function, disease

ii. Multiple layers of heterogeneous data:

- Patient-centered data integration → Precision medicine
- Disease re-classification
- Gene Ontology reconstruction
- Network alignment

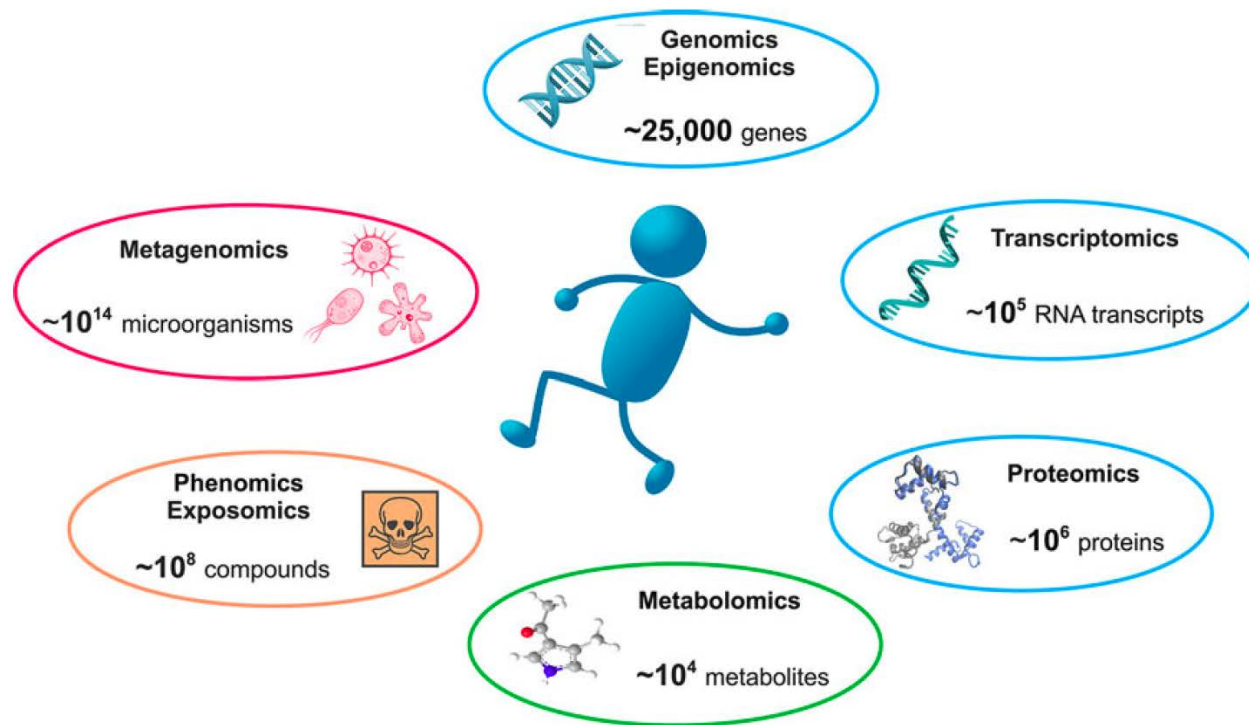
## **3. Vision**

# 1. Motivation

## Medicine: complex world of inter-connected entities

Technological advances →  
astounding harvest of various molecular and clinical data

*Proteomics* 2016, 16, 741–758



REVIEW

## Integrative methods for analyzing big data in precision medicine

Vladimir Gligorijević, Noël Malod-Dognin and Nataša Pržulj



# 1. Motivation

## Medicine: complex world of inter-connected entities

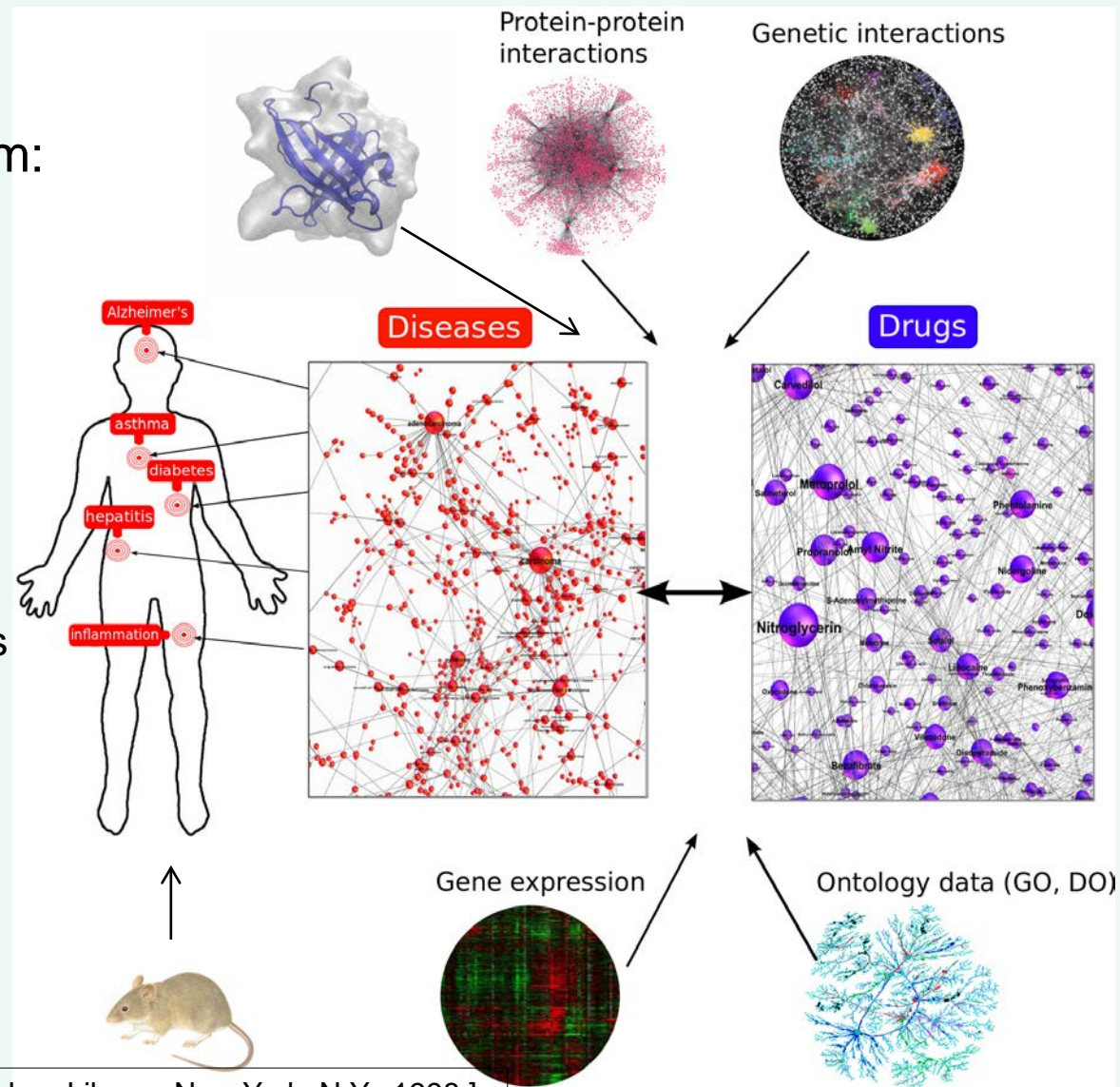
### Data growth:

➤ Guided by empirical reductionism:

- Striving to dissect a biological entity into its constituent parts
- To better understand it

➤ However, knowing parts is not enough:

- 1859 — Darwin<sup>1</sup> saw biology as a “tangled bank” with all its aspects interconnected
- 1855 — Virchow<sup>2</sup>: all diseases involve changes in normal cells



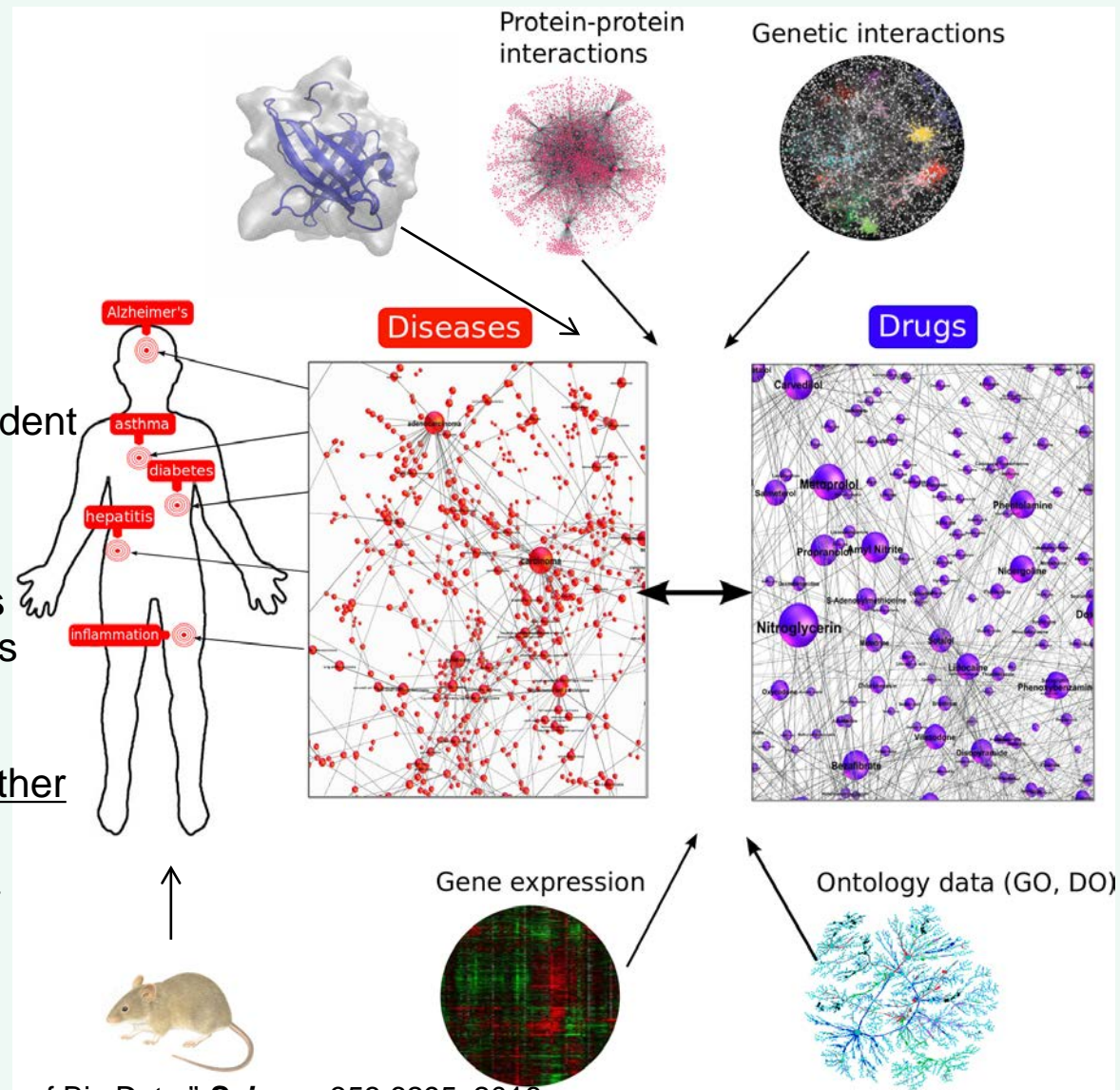
1. Darwin, C., On the origin of species, 1859 [Reprint, Modern Library, New York, N.Y., 1998.]  
2. Virchow R., Arch. Pathol. Anat. u. Physiol. u. klin. Med. 8:3, 1855

# 1. Motivation

## Medicine: complex world of inter-connected entities

### Data growth about a cell:

- Hit the wall of bio-complexity
- Cells:
  - are not just loosely coupled arrangements of quasi-independent molecules
  - highly intricately and precisely integrated **networks** of **entities** and **interactions** within the cells and with the environment
  - Data types complement each other
  - Seek joint modeling and mining



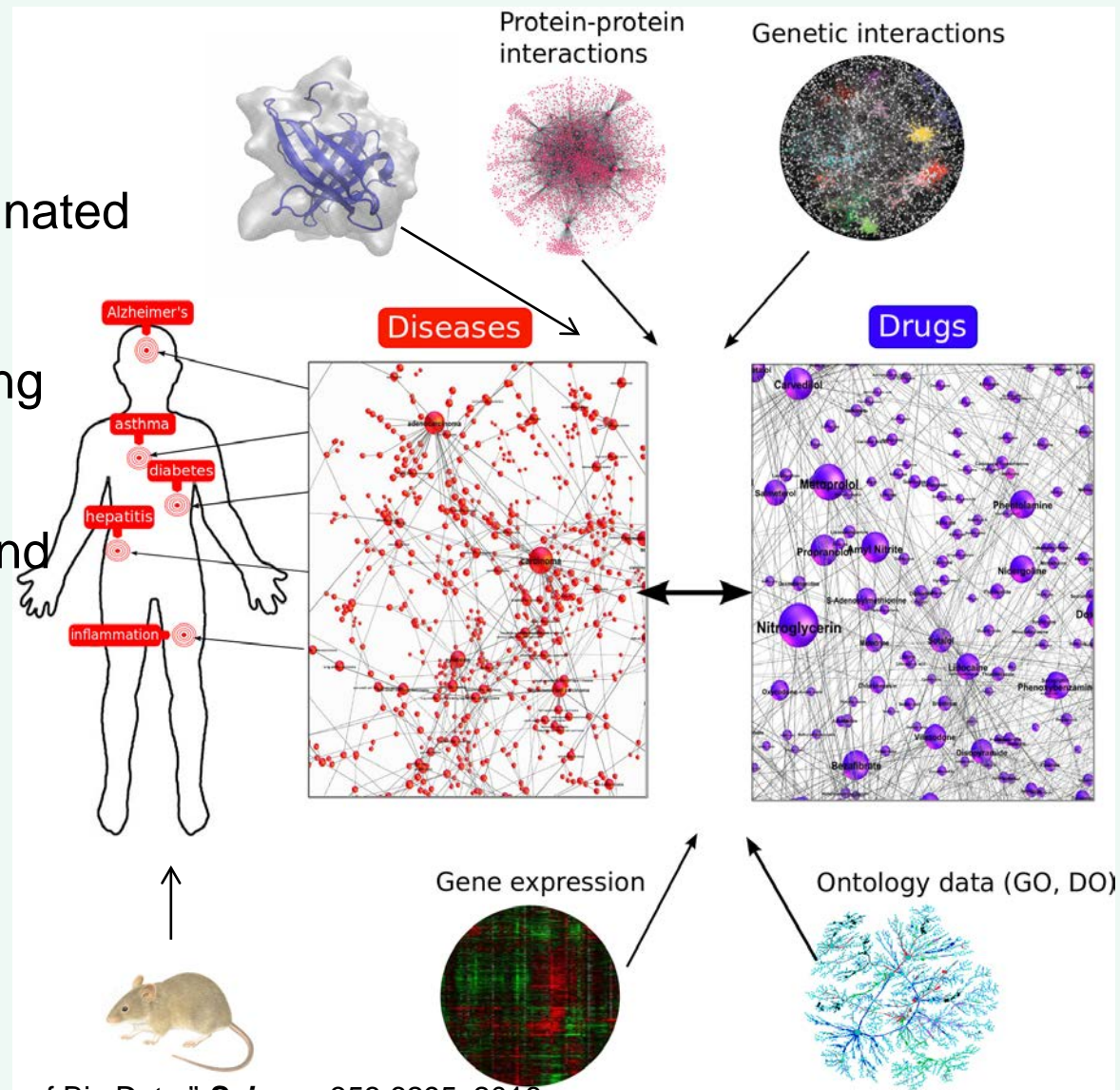


# 1. Motivation

## Medicine: complex world of inter-connected entities

### Time to:

- Replace the mostly reductionist molecular perspective that dominated the 20<sup>th</sup> century
- New and holistic view of the living world
- Required to explain biological and medical phenomena
- Biology's innate complexity



# 1. Motivation

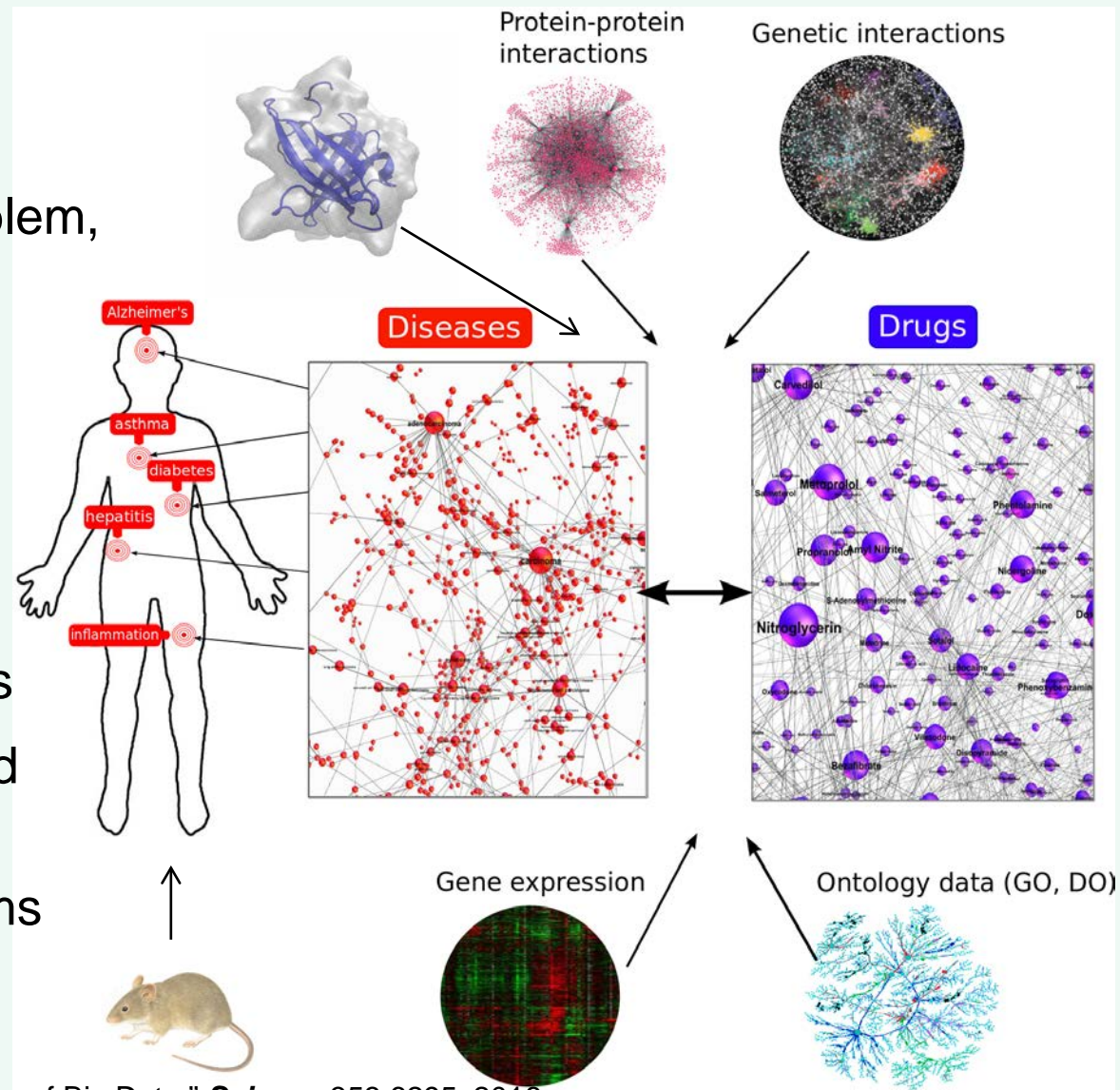
## Medicine: complex world of inter-connected entities

### Requires:

- Establishing a perspective and framework not only for one problem, but for biology and medicine in general

### A foremost challenge:

- How to re-synthesize biology
- Put the elements back into their complex, dynamic environments
- Connect them all within a unified framework
- Reformulate biological paradigms within the non-linear world



# 1. Motivation

## Medicine: complex world of inter-connected entities

### Vision:

- Bridge this gap by developing a mathematically principled framework for integration of networked data
- Marry biomedical problems and data with algorithms from:
  - ML, such as NMTF
  - Mathematical non-linear optimization
  - Network science
  - Algebraic topology...
  - High-performance computing
- Propose modelling & computational advances that will link the medicine's:
  - reductionist past with its holistic future
- Enable
  - displacement of the dominant molecular representation of biology
  - by a new, integrative paradigm that is deeper, more comprehensive and inspiring

**€2M ERC Consolidator Grant for 2018-2023**

**Title: “Integrated Connectedness for a New Representation of Biology”**

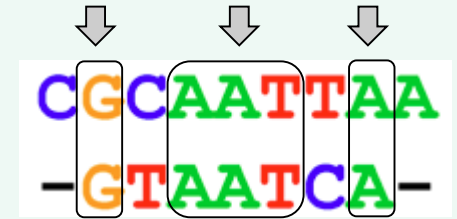


# 1. Motivation

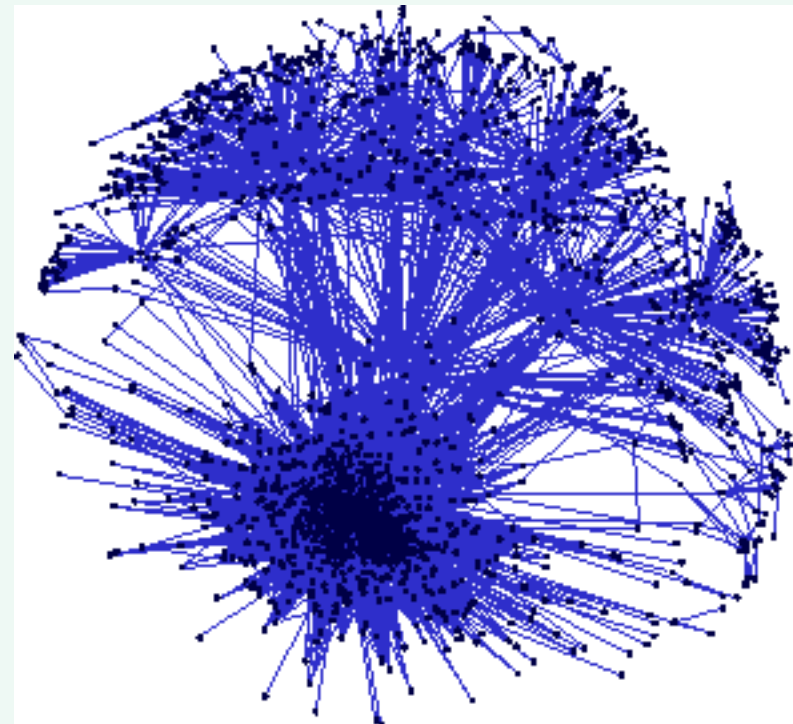
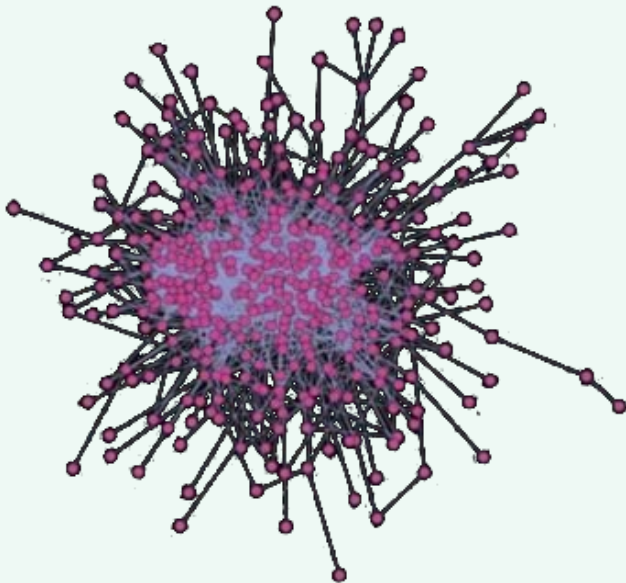
Medicine: complex world of inter-connected entities

## Computational challenges

- Need new tools to mine complex data systems
- Why?
  - Analysing sequences: “computationally easy” → still lacking
  - Analysing interconnected heterogeneous data: “computationally hard”



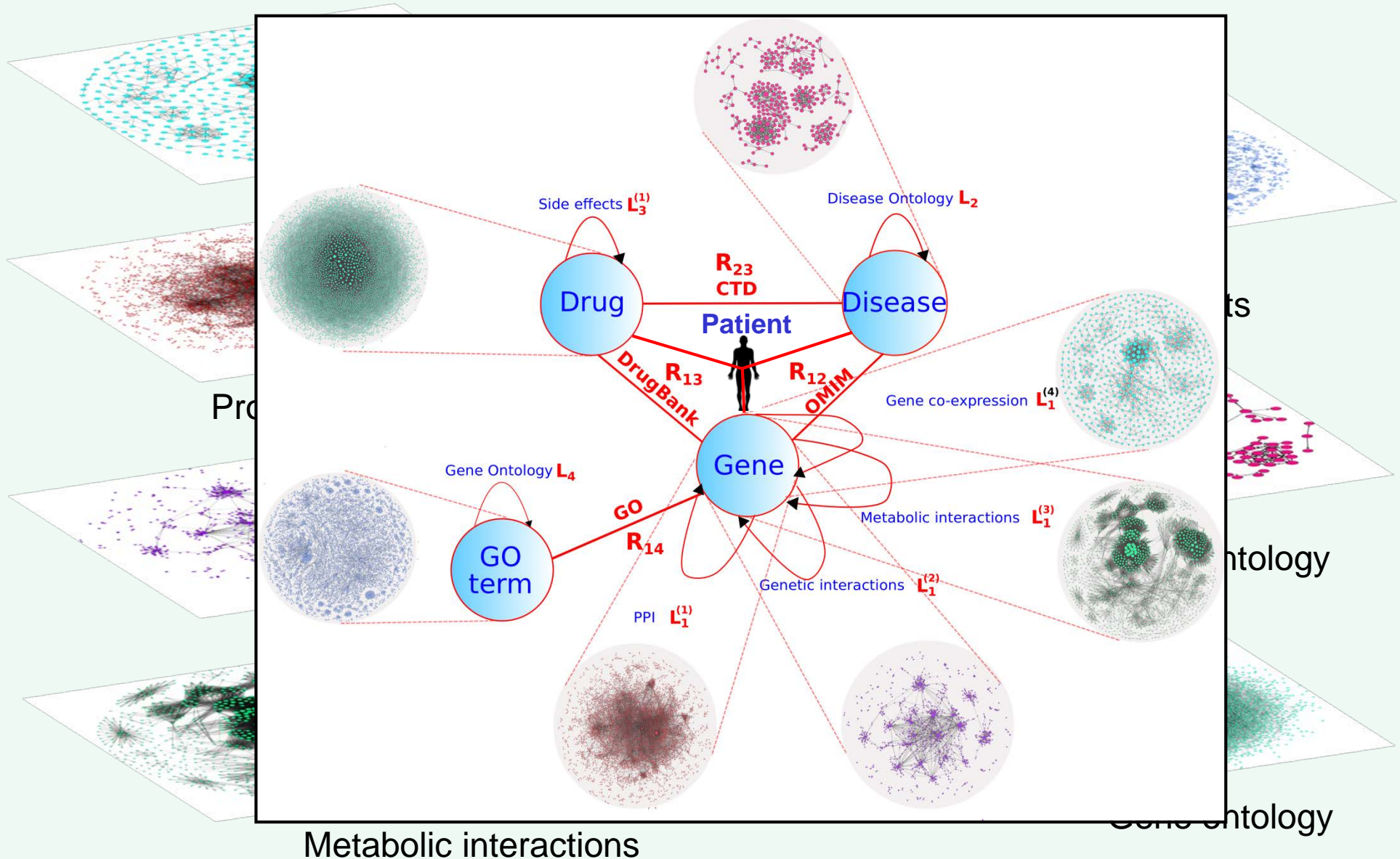
- **Sophisticated** methods **carefully tuned** to extract new knowledge from **particular data**



# 1. Motivation

Medicine: complex world of inter-connected entities

## Computational challenges





# Overview

**Medicine: complex world of inter-connected entities**

## **1. Motivation**

## **2. New Methods – Examples: mine inter-connected data**

i. Single layer of omics data: Molecular networks → function, disease

ii. Multiple layers of heterogeneous data:

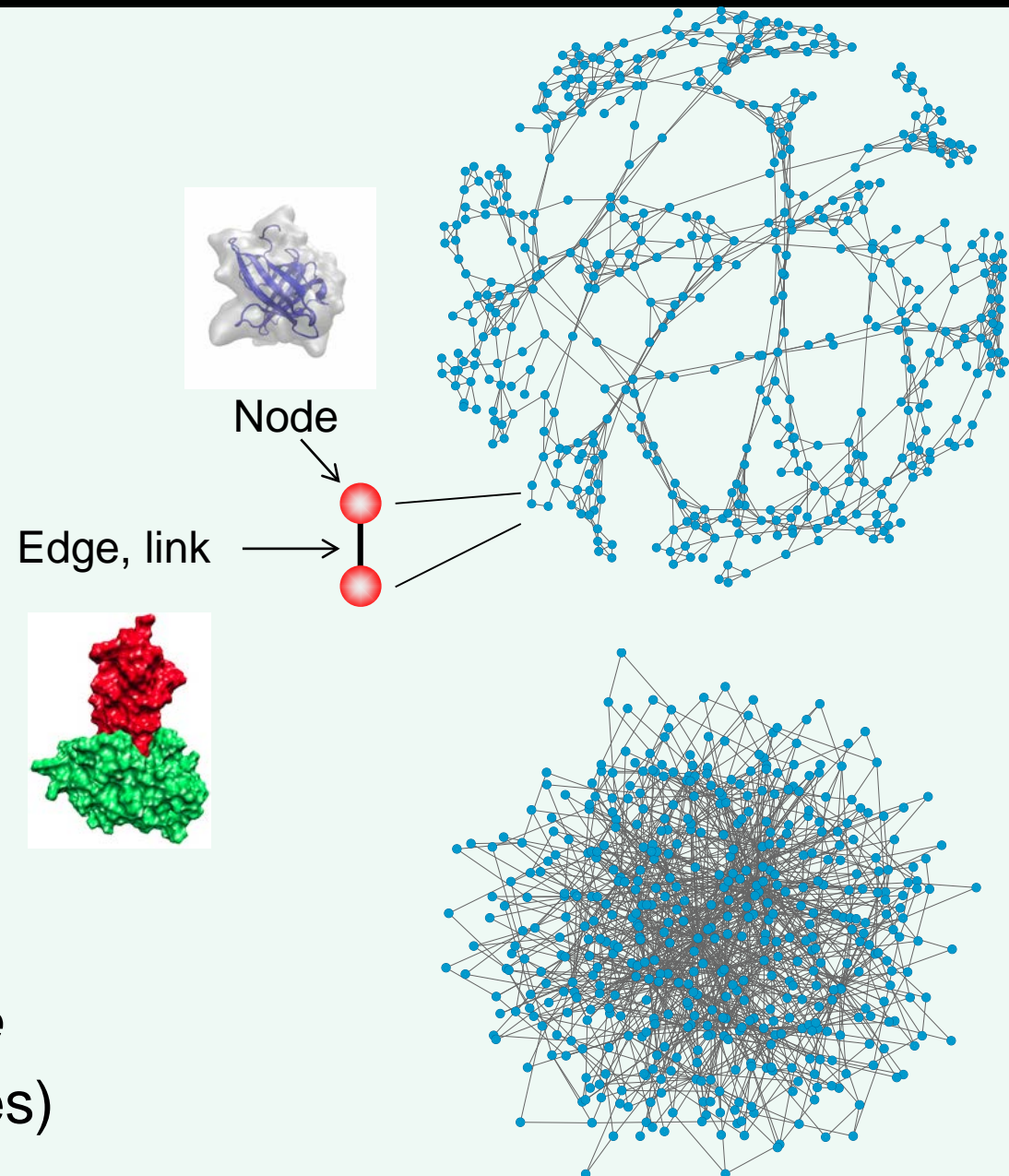
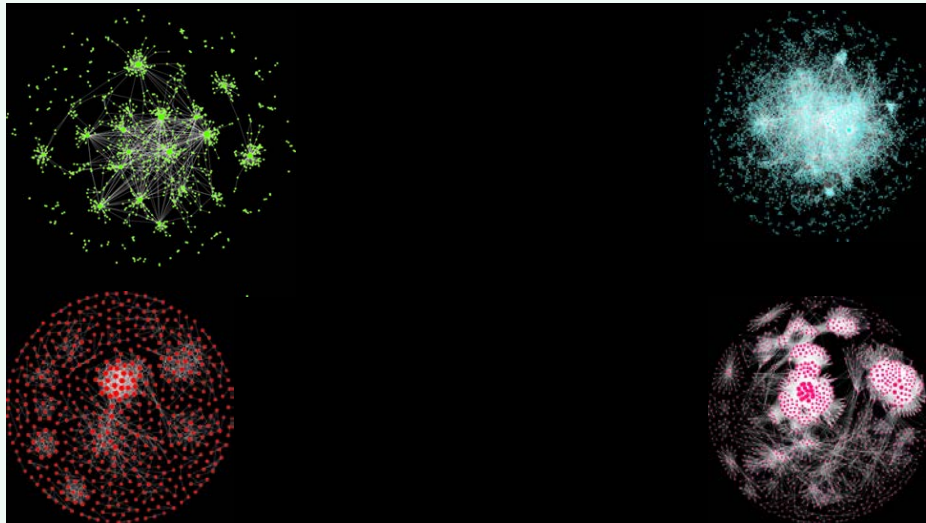
- Patient-centered data integration → Precision medicine
- Disease re-classification
- Gene Ontology reconstruction
- Network alignment

## **3. Vision**

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### i. Molecular Networks

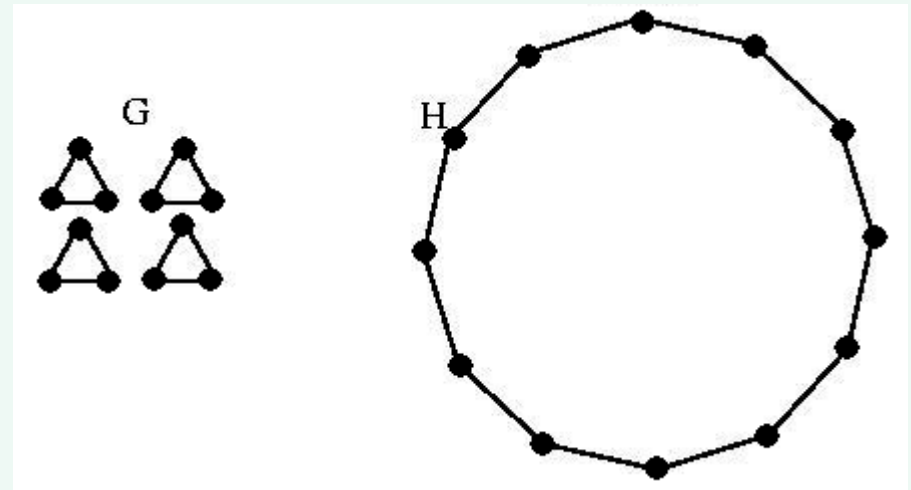
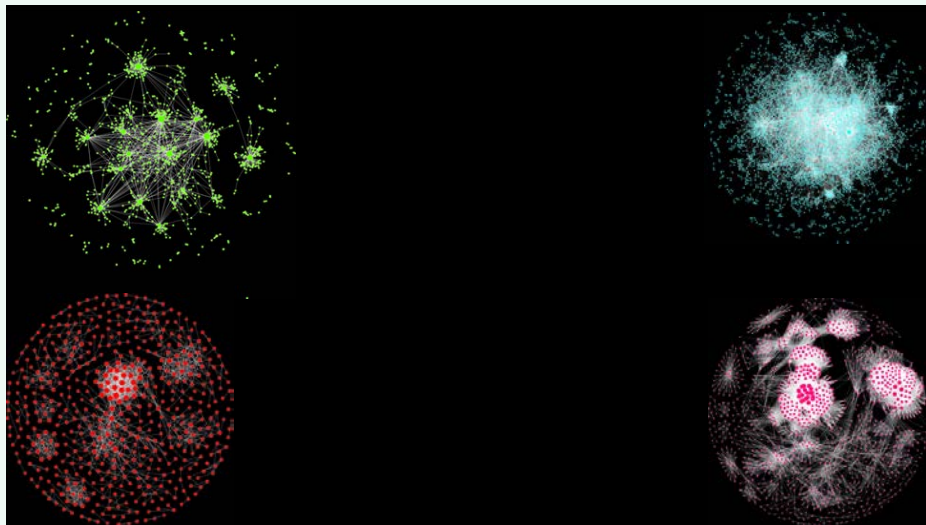


- The number of nodes
- The number of links
- Links of each node: *degree*
- Distribution of links (degrees)

## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

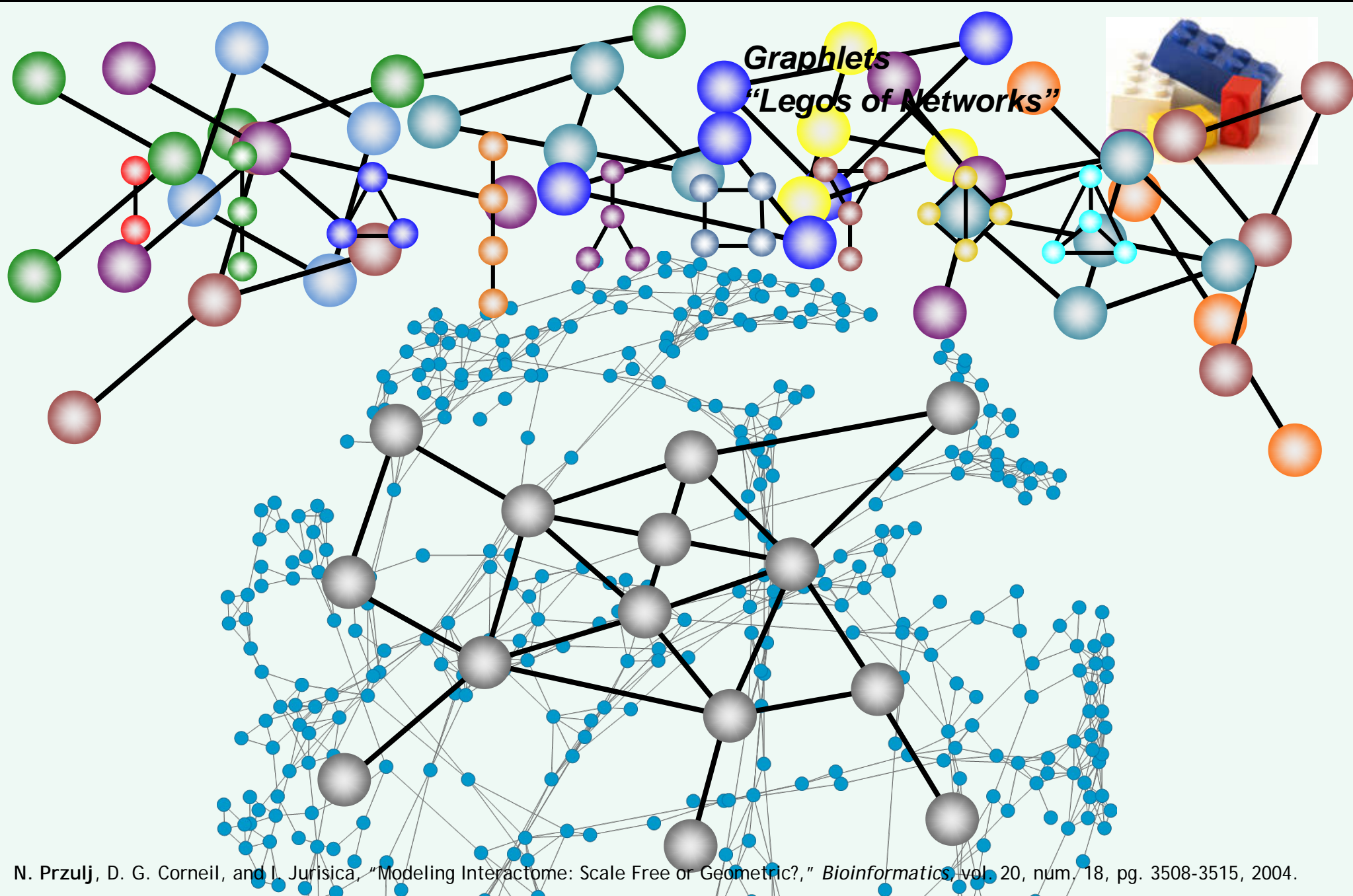
#### i. Molecular Networks



- The number of nodes
- The number of links
- Links of each node: *degree*
- Distribution of links (degrees)

## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities



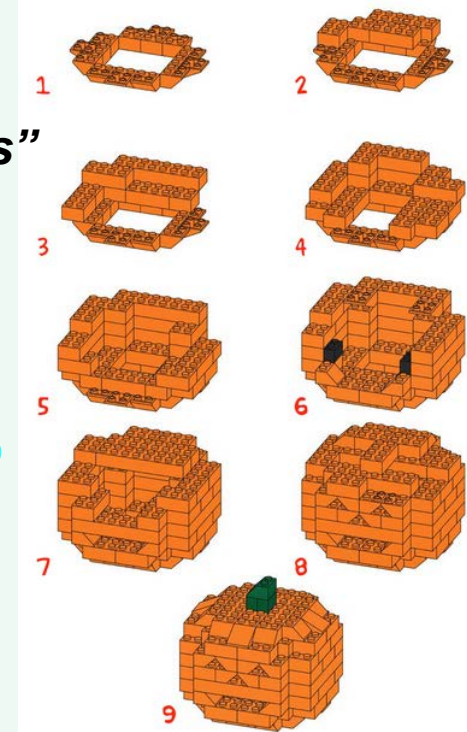
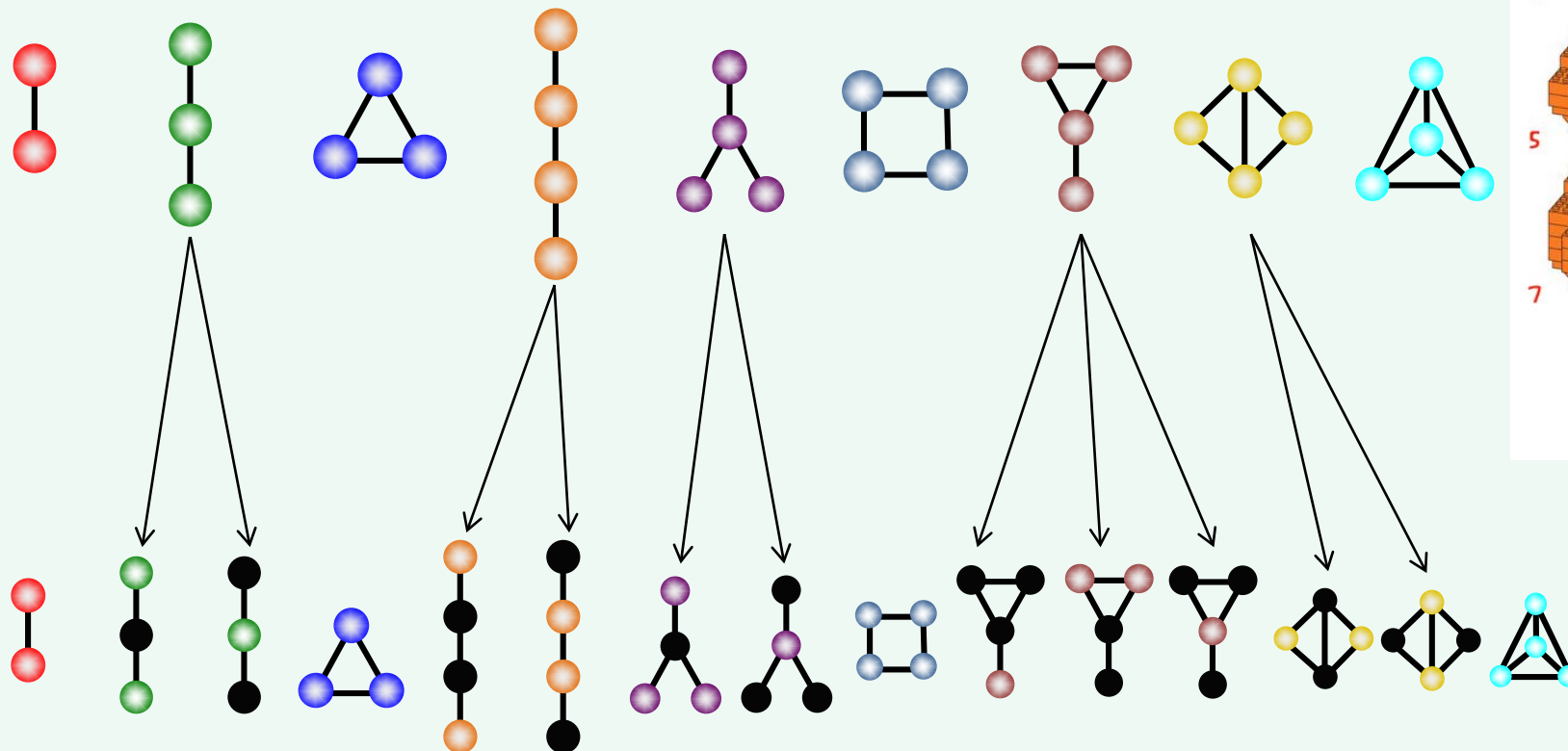


# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

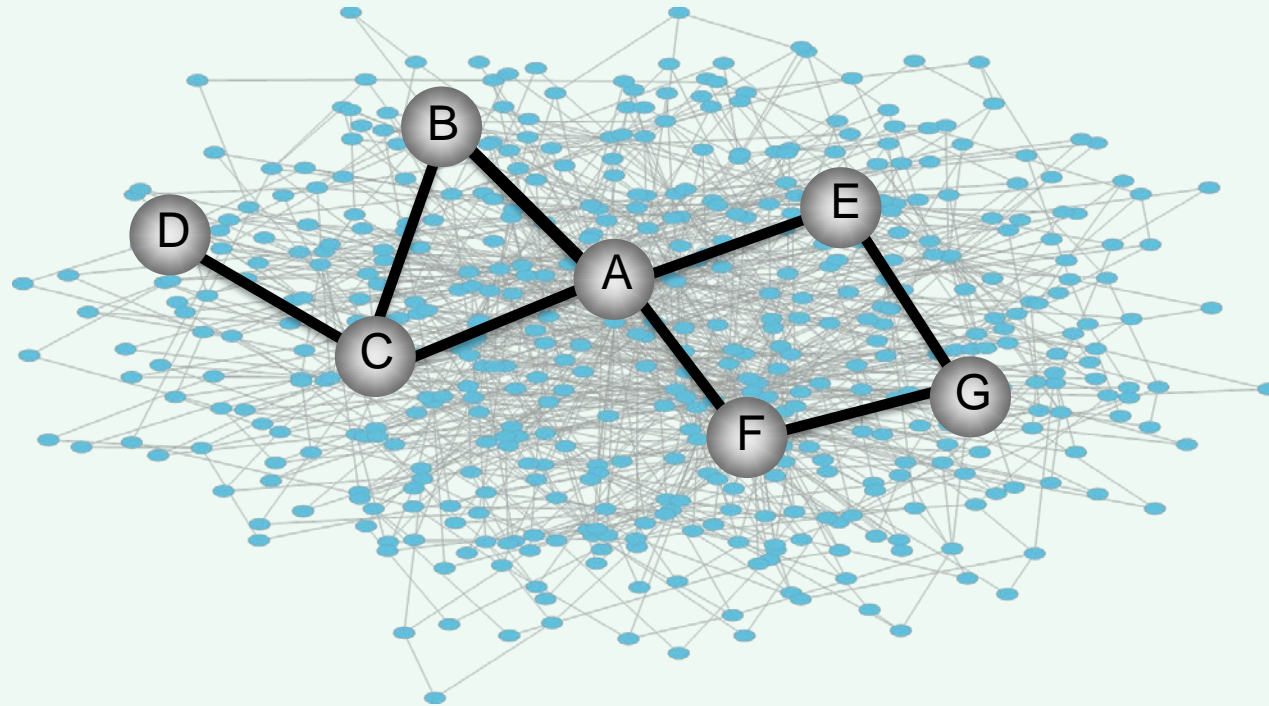
**ERC StG: 278212 (2012-2017): "Biological Network Topology Complements Genome as a Source of Biological Information"**

**Graphlets**  
"Legos of Networks"



## 2. Novel Methods

Mine the Medical World of Inter-Connected Entities

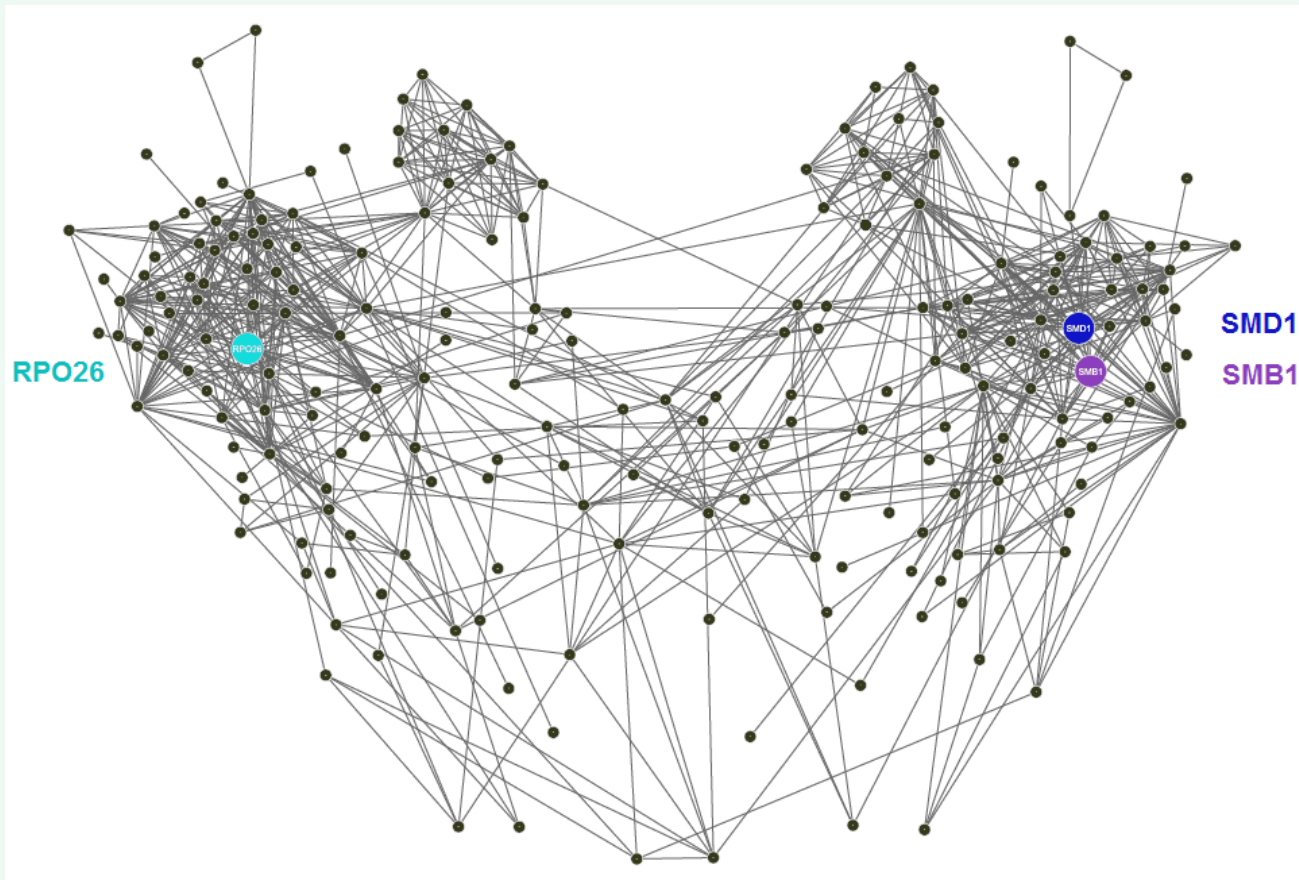


Orbit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
A	4	3	5	1	0	6	0	2	1	0	1	2	0	0	0
B	2	3	0	1	2	0	1	0	0	0	3	0	0	0	0
C	3	2	2	1	2	2	1	0	0	0	2	1	0	0	0
D	1	2	0	0	2	0	0	0	0	1	0	0	0	0	0
E	2	4	1	0	1	2	2	0	1	1	0	0	0	0	0
F	2	4	1	0	1	2	2	0	1	1	0	0	0	0	0
G	2	2	1	0	4	0	0	0	1	0	0	0	0	0	0

## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

90% similar wiring – significantly enriched:

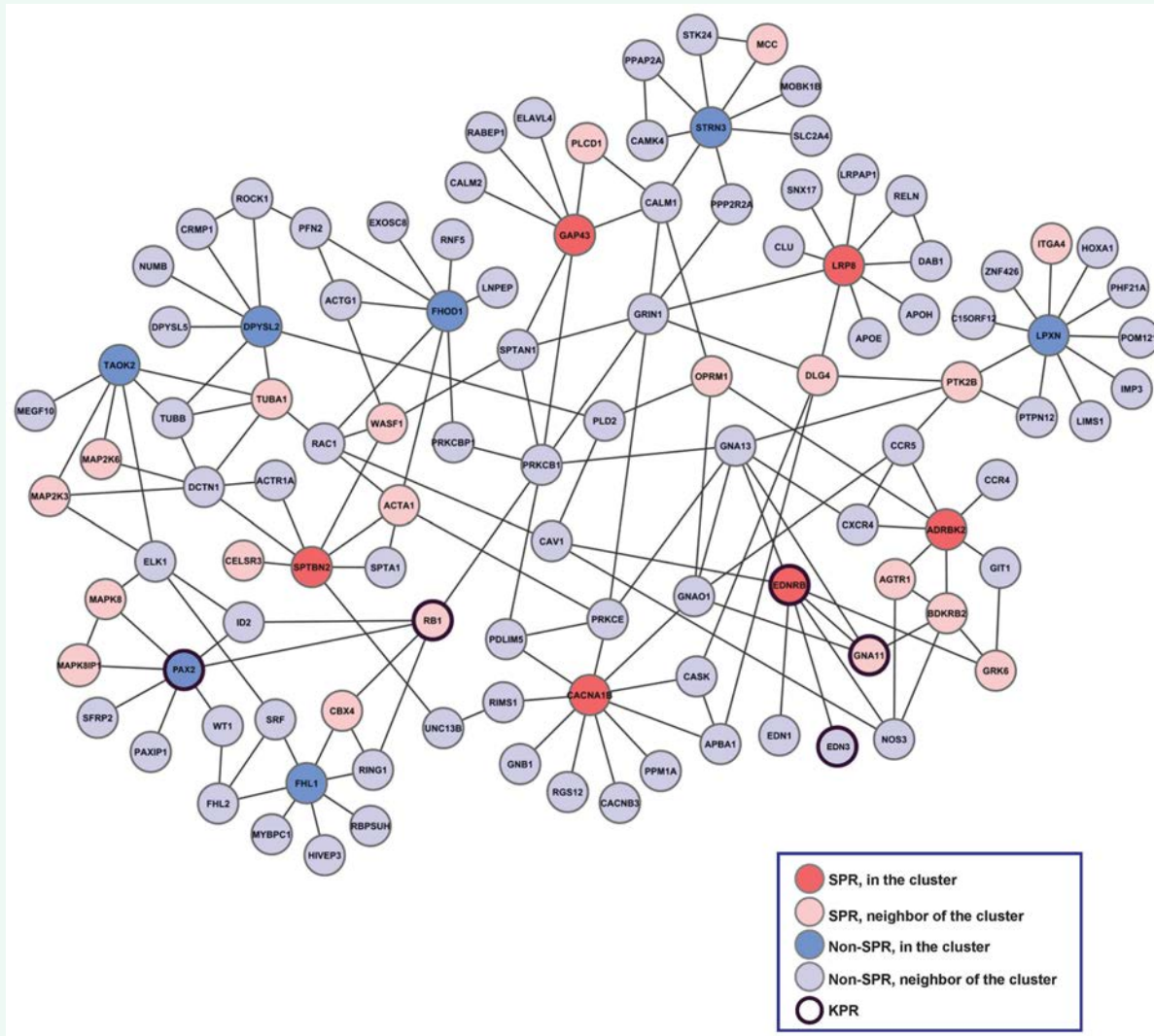


- Biological function
- Protein complexes
- Sub-cellular localization
- Tissue expression
- Disease



# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

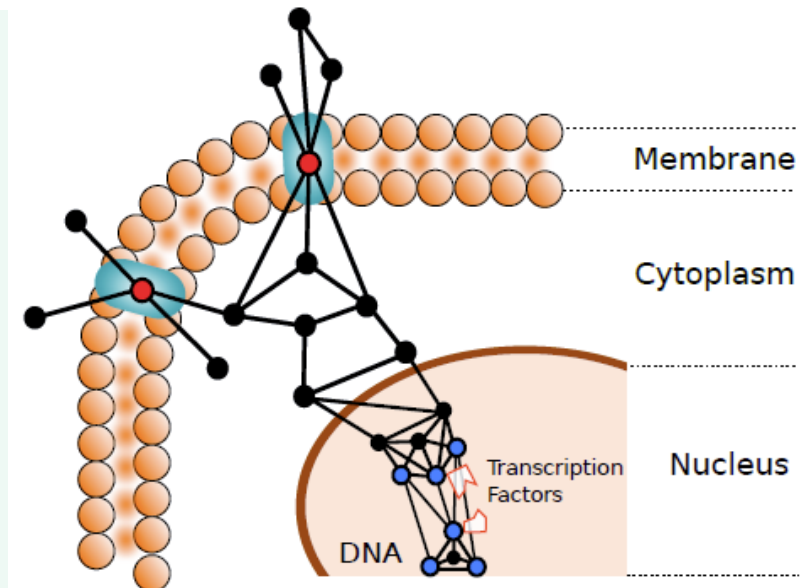
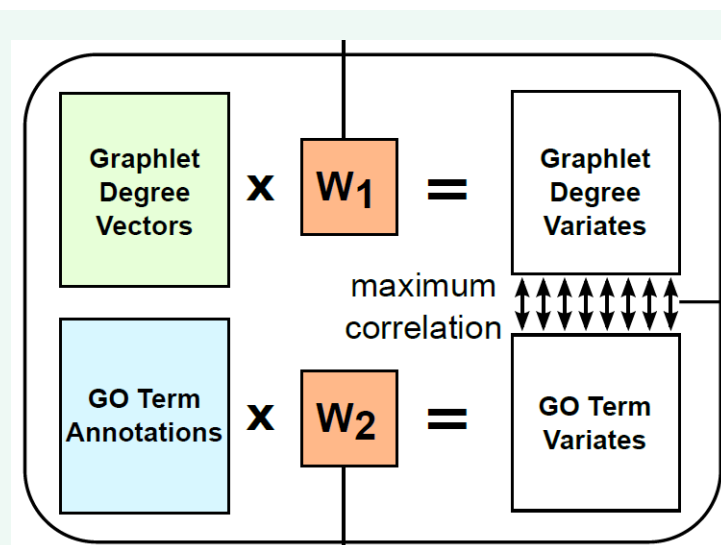
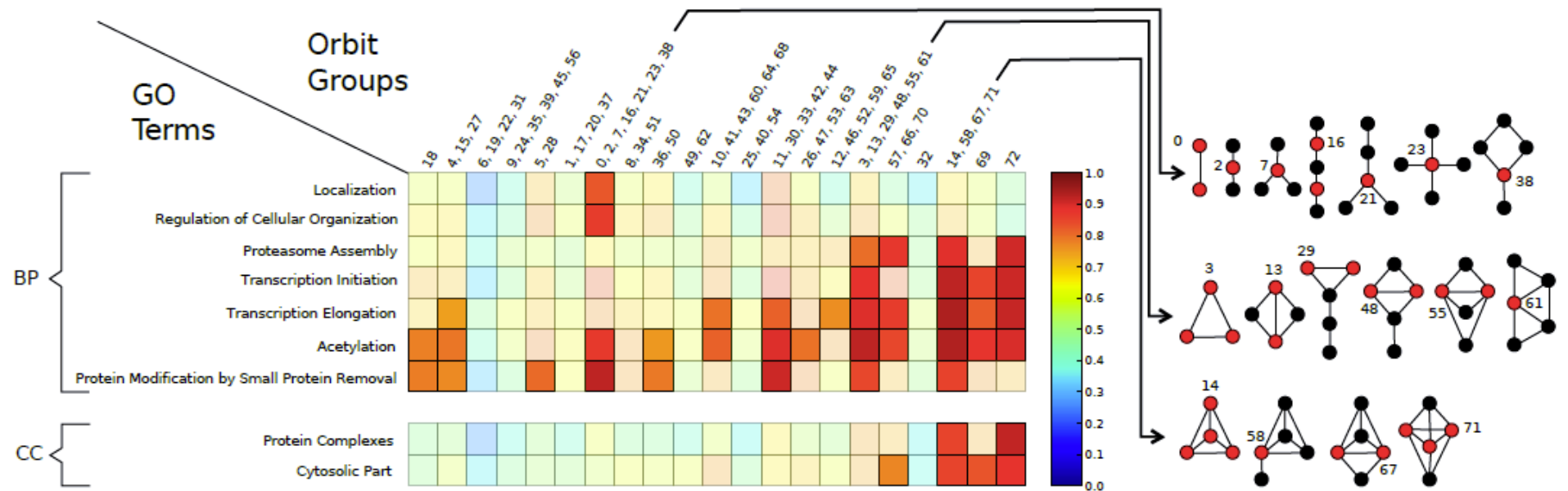


### Cancer research:

- New proteins for melanin production
- Same cancer type: more similar wiring
- Far away in the network

# 2. Novel Methods

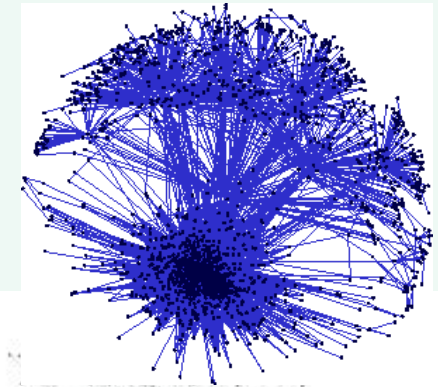
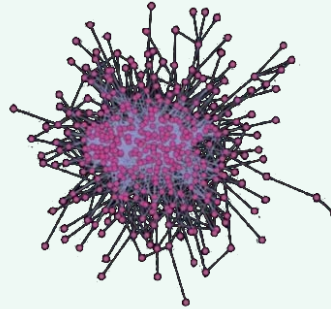
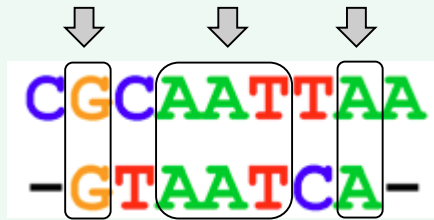
## Mine the Medical World of Inter-Connected Entities



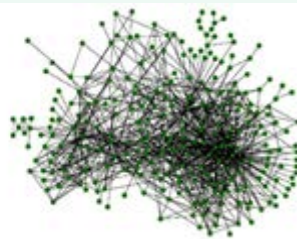
# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

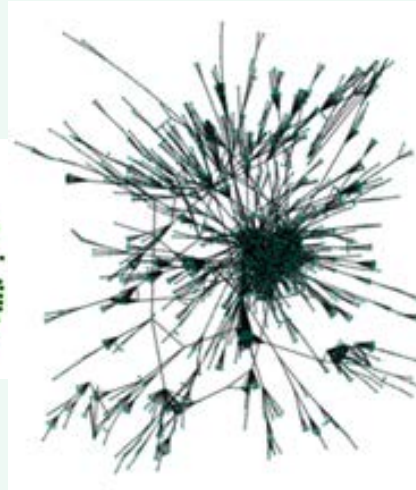
### Network Alignment



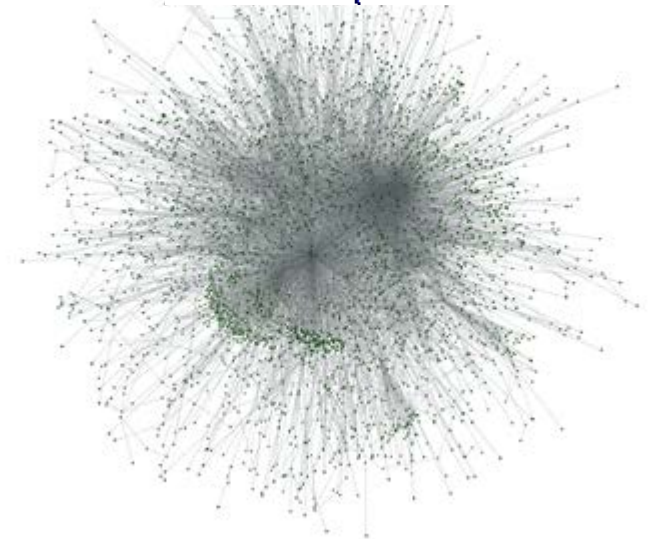
Isorank:  
116 nodes  
261 edges



GRAAL:  
267 nodes  
900 edges



MI-GRAAL:  
1,858 nodes  
3,467 edges



L-GRAAL:  
5,726 nodes  
16,084 edges  
**Yeast: 98% proteins  
21% interactions**

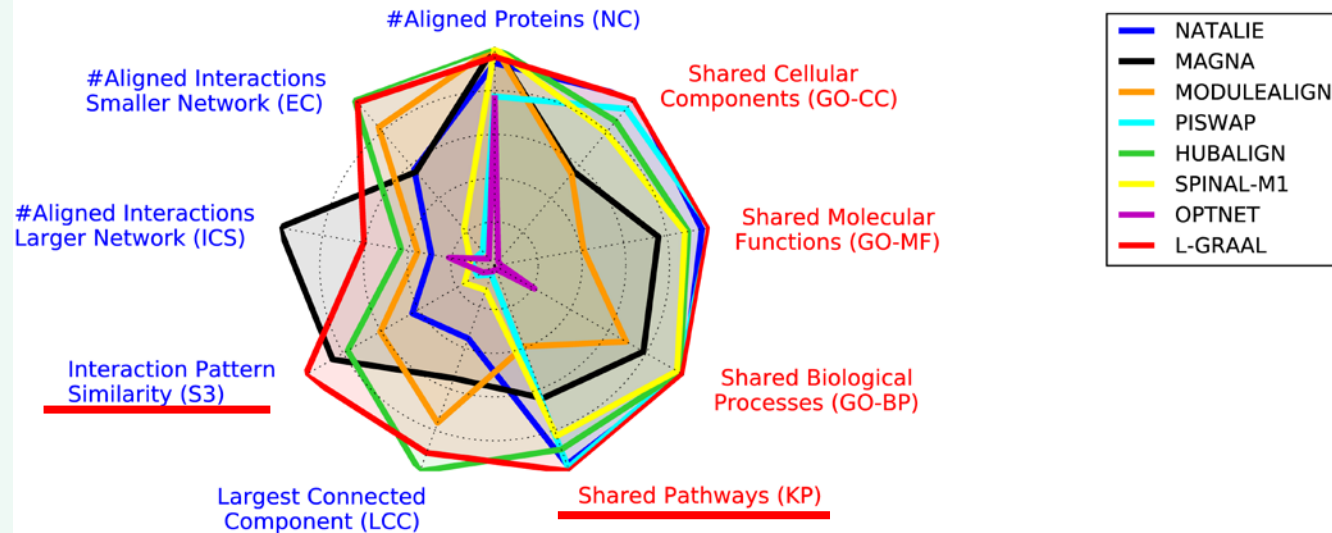
N. Malod-Dognin & N. Pržulj, L-GRAAL, *Bioinformatics*, doi: 10.1093/bioinformatics/btv130, 2015  
N. Malod-Dognin & N. Pržulj, GR-ALIGN, *Bioinformatics*, doi:10.1093/bioinformatics/btu020, 2014  
V. Memisevic & N. Pržulj, C-GRAAL, *Integrative Biology*, doi:10.1039/c2ib00140c, 2012  
O. Kuchaiev & N. Pržulj, MI-GRAAL, *Bioinformatics*, 27(10): 1390-6, 2011  
O. Kuchaiev, T. Milenkovic, V. Memisevic, W. Hayes, & N. Pržulj, *J. Royal Society Interface*, 7:1341-1354, 2010  
T. Milenkovic, W.L. Wong, W. Hayes, & N. Pržulj, *Cancer Informatics*, 9:121-37, June 30, 2010 (Highly visible)

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Alignment of PPI Networks – Ualign

- Many methods
- All heuristic
- No gold standard
- Questions:
  - Which aligner for which data?
  - Which scoring scheme for evaluation?
  - Coverage: biological and topological?
  - Contribution of topology vs sequence?



- Map biologically and topologically **different** network regions
- Each covers only about 50% of the proteins of the larger network
- **Together** – map **entire** networks → Ualign
  - Biologically coherent

- The most topologically coherent – using topology only
- The most biologically coherent – using sequence only

Why?

- Existing annotations ill-suited?
- **Methodological limitations?**

→ **Combine** topology and sequence information



# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

- ✓ The best performing
- ✓ Robust
- ✓ ...

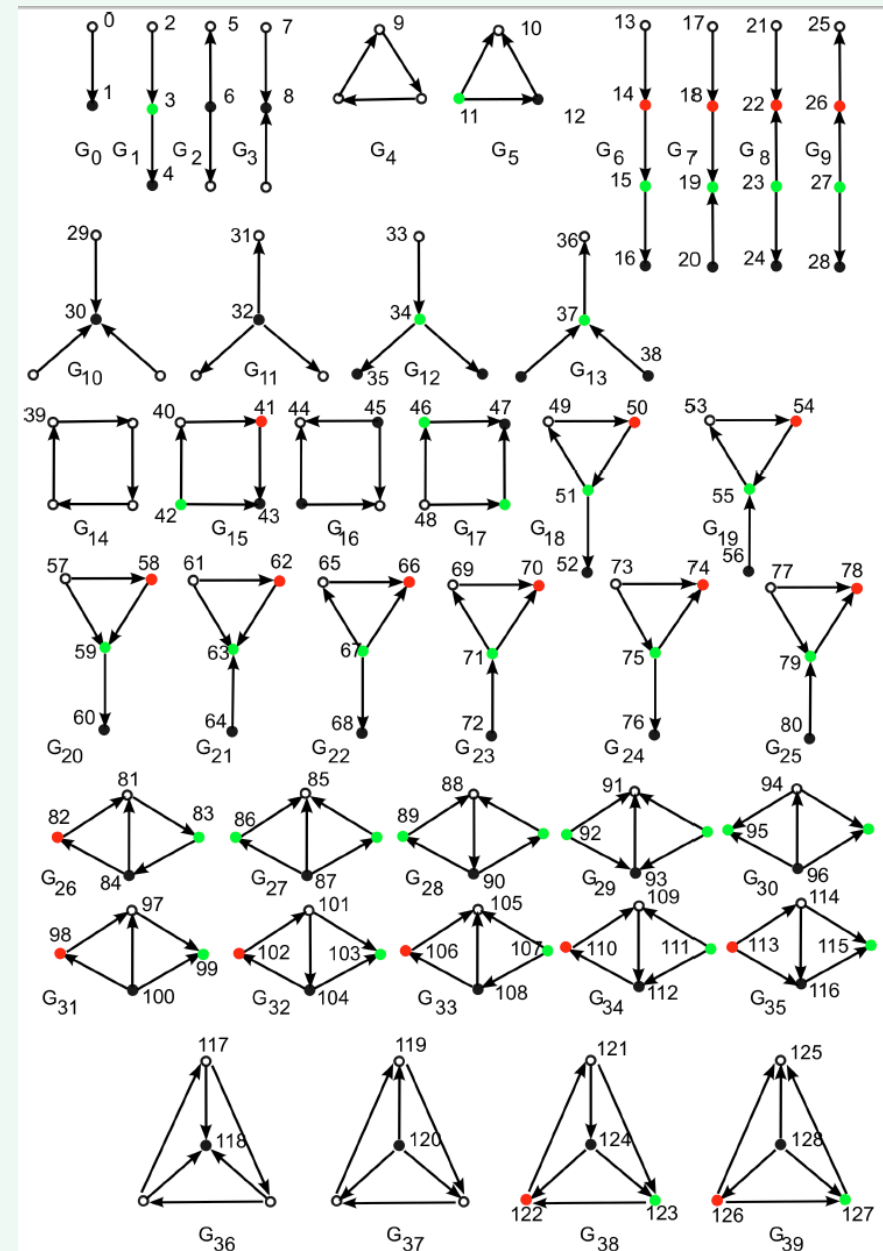
### □ PPI networks are *geometric*

N. Przulj, D. G. Corneil, and I. Jurisica, "Modeling Interactome: Scale Free or Geometric?," *Bioinformatics*, vol. 20, num. 18, pg. 3508-3515, 2004.

N. Przulj, "Biological Network Comparison Using Graphlet Degree Distribution," Proceedings of the 2006 European Conference on Computational Biology, ECCB '06, Eilat, Israel, January 21-24, 2007, acceptance rate 18%. *Bioinformatics*, volume 23, pages e177-e183, 2007

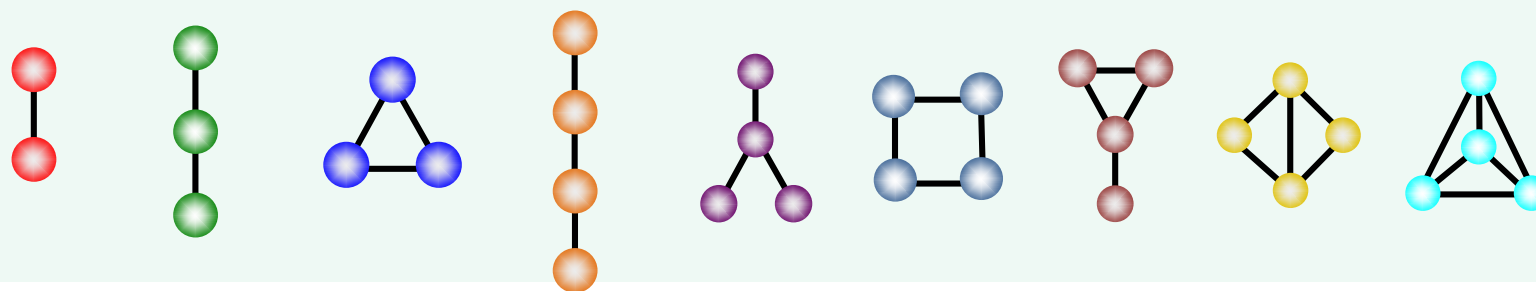
...

- Directed Networks
- Track dynamics

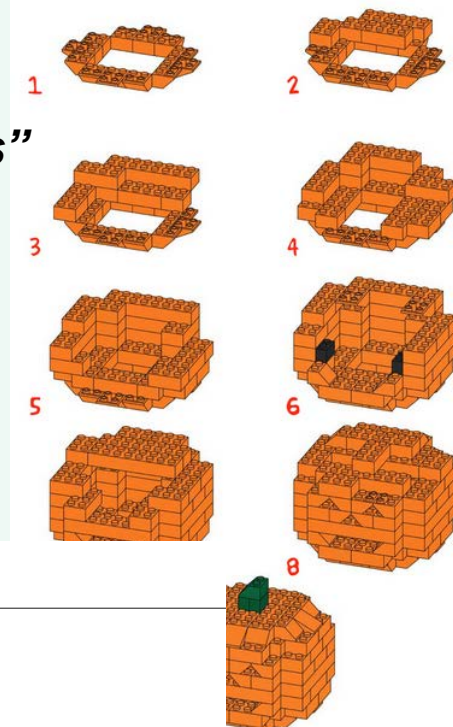


# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities



**Graphlets**  
“Legos of Networks”



INSIGHTS | PERSPECTIVES

## Network analytics in the age of Big Data

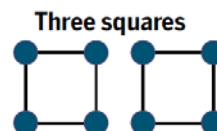
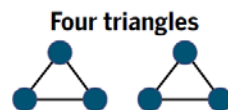
How can we holistically mine big data?

By **Nataša Pržulj** and **Noël Malod-Dognin**

**W**e live in a complex world of inter-connected entities. In all areas of human endeavor, from biology to medicine, economics, and climate science, we are flooded with large-scale data sets. They describe intricate real-world systems from different and complementary viewpoints, with entities being modeled as nodes and their connections as edges, comprising large networks. This is

### Network structures

The four networks shown have exactly the same size (the same number of nodes and edges), and each node within each network has the same degree (the number of interactions with other nodes), but each network canis of very different structure.



into RNAs and translated into proteins, which adopt various three-dimensional structures to carry out particular cellular functions. Molecular interactions are captured by different high-throughput biotechnologies and modeled with different types of networks. Individual analyses of molecular networks have revealed that molecules involved in similar functions tend to group together in a network and are similarly wired (13), leading to better understanding of gene functions (6) and molecular organization of the cell (7) and to im-

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### A global genetic interaction network maps a wiring diagram of cellular function

*Science* 23 Sep 2016:  
Vol. 353, Issue 6306, aaf1420  
DOI: 10.1126/science.aaf1420

Michael Costanzo<sup>1\*</sup>, Benjamin VanderSluis<sup>2,3\*</sup>, Elizabeth N. Koch<sup>2\*</sup>, Anastasia Baryshnikova<sup>4\*</sup>, Carles Pons<sup>2\*</sup>, Guihong Tan<sup>1\*</sup>, Wen Wang<sup>2</sup>, Matej Usaj<sup>1</sup>, Julia Hanchard<sup>1,5</sup>, Susan D. Lee<sup>6</sup>, Vincent Pelechano<sup>7\*</sup>, Erin B. Styles<sup>1,5</sup>, Maximilian Billmann<sup>8</sup>, Jolanda van Leeuwen<sup>1</sup>, Nydia van Dyk<sup>1</sup>, Zhen-Yuan Lin<sup>9</sup>, Elena Kuzmin<sup>1,5</sup>, Justin Nelson<sup>2,10</sup>, Jeff S. Piotrowski<sup>1,11\*</sup>, Tharan Srikumar<sup>12</sup>, Sondra Bahr<sup>1</sup>, Yiqun Chen<sup>1</sup>, Raamesh Deshpande<sup>2</sup>, Christoph F. Kurat<sup>1\*</sup>, Sheena C. Li<sup>1,11</sup>, Zhijian Li<sup>1</sup>, Mojca Mattiazzi Usaj<sup>1</sup>, Hiroki Okada<sup>13</sup>, Natasha Pascoe<sup>1,5</sup>, Bryan-Joseph San Luis<sup>1</sup>, Sara Sharifpoor<sup>1</sup>, Emira Shuteriqi<sup>1</sup>, Scott W. Simpkins<sup>2,10</sup>, Jamie Snider<sup>1</sup>, Harsha Garadi Suresh<sup>1</sup>, Yizhao Tan<sup>1</sup>, Hongwei Zhu<sup>1</sup>, Noel Malod-Dognin<sup>1</sup>, Vuk Janjić<sup>1</sup>, Nataša Pržulj<sup>1</sup>, Olga G. Troyanskaya<sup>3,4</sup>, Igor Stagljar<sup>1,5,16</sup>, Tian Xia<sup>1,5,17</sup>, Yoshikazu Ohya<sup>13</sup>, Anne-Claude Gingras<sup>5,9</sup>, Brian Raught<sup>12</sup>, Michael Boutros<sup>8</sup>, Lars M. Steinmetz<sup>7,18</sup>, Claire L. Moore<sup>6</sup>, Adam P. Rosebrock<sup>1,5</sup>, Amy A. Caudy<sup>1,5</sup>, Chad L. Myers<sup>2,10\*</sup>, Brenda Andrews<sup>1,5\*</sup>, and Charles Boone<sup>1,5\*</sup>

INSIGHTS | PERSPECTIVES

## Network analytics in the age of Big Data

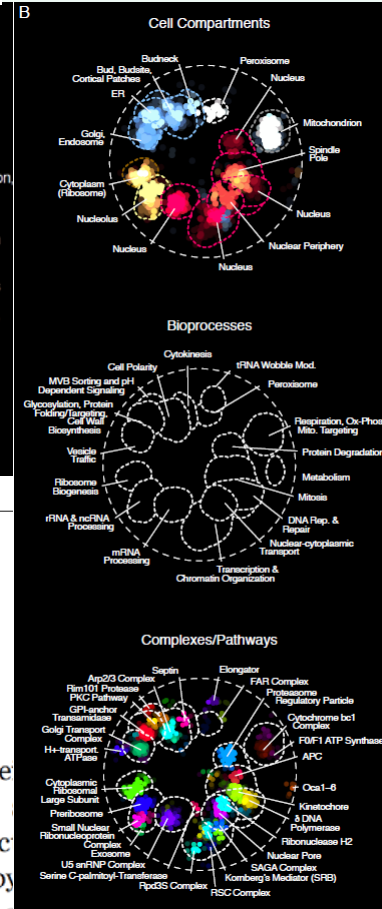
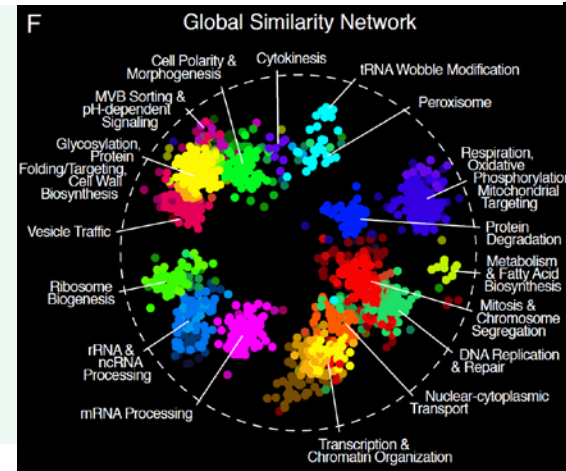
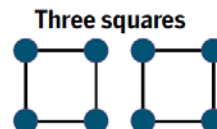
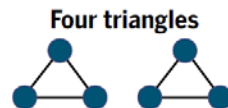
How can we holistically mine big data?

By Nataša Pržulj and Noël Malod-Dognin

We live in a complex world of inter-connected entities. In all areas of human endeavor, from biology to medicine, economics, and climate science, we are flooded with large-scale data sets. They describe intricate real-world systems from different and complementary viewpoints, with entities being modeled as nodes and their connections as edges, comprising large networks. This is

### Network structures

The four networks shown have exactly the same size (the same number of nodes and edges), and each node within each network has the same degree (the number of interactions with other nodes), but each network canis of very different structure.



into RNAs and translated into proteins. adopt various three-dimensional to carry out particular cellular functions. molecular interactions are captured by high-throughput biotechnologies and modeled with different types of networks. Individual analyses of molecular networks have revealed that molecules involved in similar functions tend to group together in a network and are similarly wired (13), leading to better understanding of gene functions (6) and molecular organization of the cell (7) and to im-



# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

Published online: March 15, 2017

Article

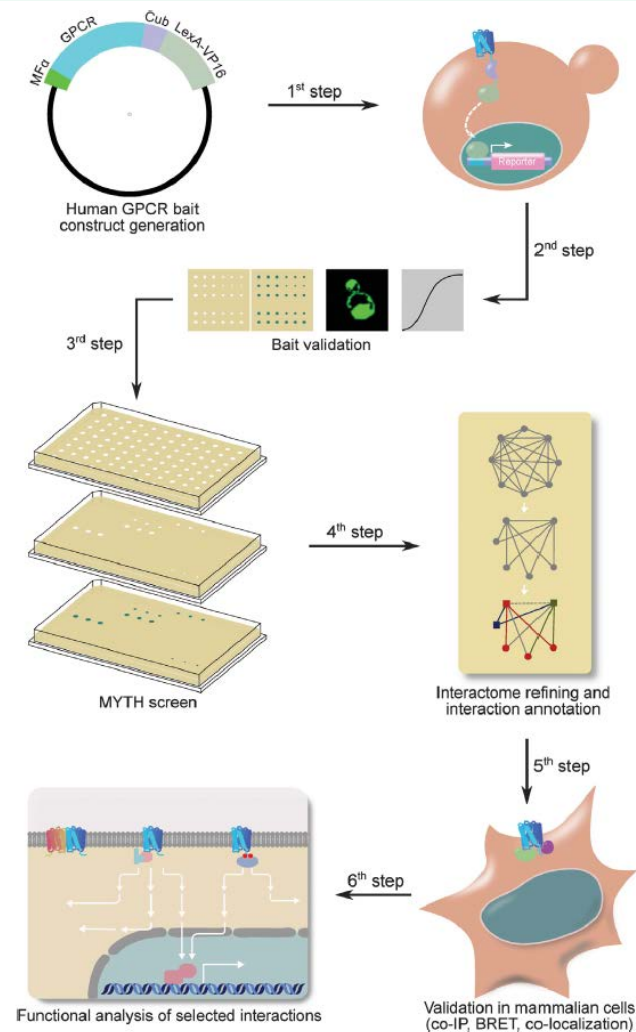
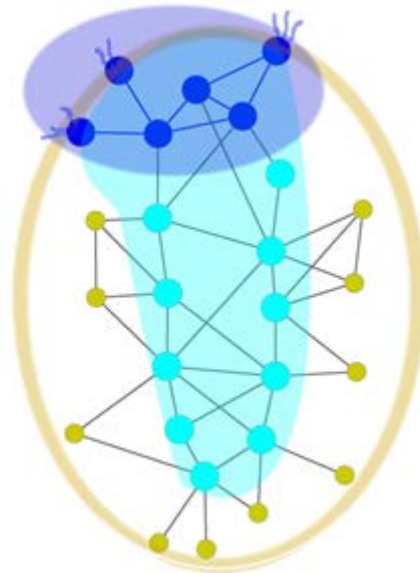


molecular  
systems  
biology

### Systematic protein–protein interaction mapping for clinically relevant human GPCRs

Kate Sokolina<sup>1,†</sup>, Saranya Kittanakom<sup>1,†</sup>, Jamie Snider<sup>1,†</sup>, Max Kotlyar<sup>2</sup>, Pascal Maurice<sup>3,4,5,6</sup> , Jorge Gandía<sup>7,8</sup> , Abba Benleulmi-Chaachoua<sup>3,4,5</sup>, Kenjiro Tadagaki<sup>3,4,5</sup>, Atsuro Oishi<sup>3,4,5</sup>, Victoria Wong<sup>1</sup>, Ramy H Maltz<sup>9</sup>, Viktor Deineko<sup>9</sup>, Hiroyuki Aoki<sup>9</sup>, Shahreen Amin<sup>9</sup>, Zhong Yao<sup>1</sup>, Xavier Morató<sup>7,8</sup>, David Otasek<sup>2</sup>, Hiroyuki Kobayashi<sup>10</sup>, Javier Menendez<sup>1</sup>, Daniel Auerbach<sup>11</sup>, Stephane Angers<sup>12</sup>, **Natasa Pržulj**<sup>13</sup> , Michel Bouvier<sup>10</sup>, Mohan Babu<sup>9</sup>, Francisco Ciruela<sup>7,8</sup>, Ralf Jockers<sup>3,4,5</sup>, Igor Jurisica<sup>2,14,15</sup> & **Igor Stagljarić**<sup>16,17,\*</sup>

- ✓ “Spine” of the network
  - “Dominating set” heuristic
- ✓ Functionally and topologically separates the cell
- ✓ Predict new GPCRs:
  - e.g., chromosome 20 open reading frame 39 (TMEM90B)



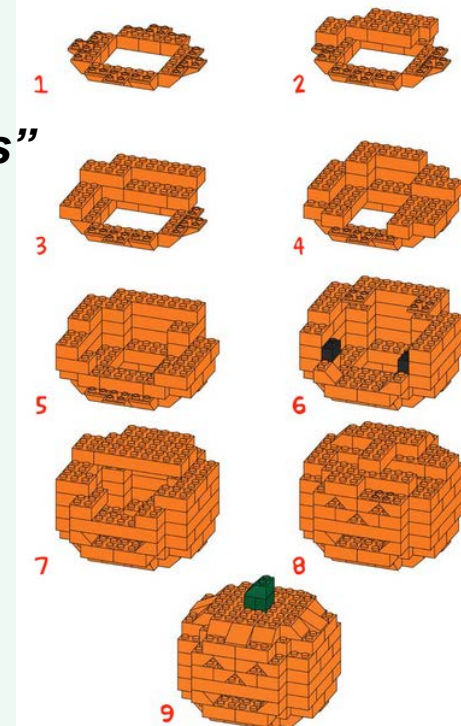
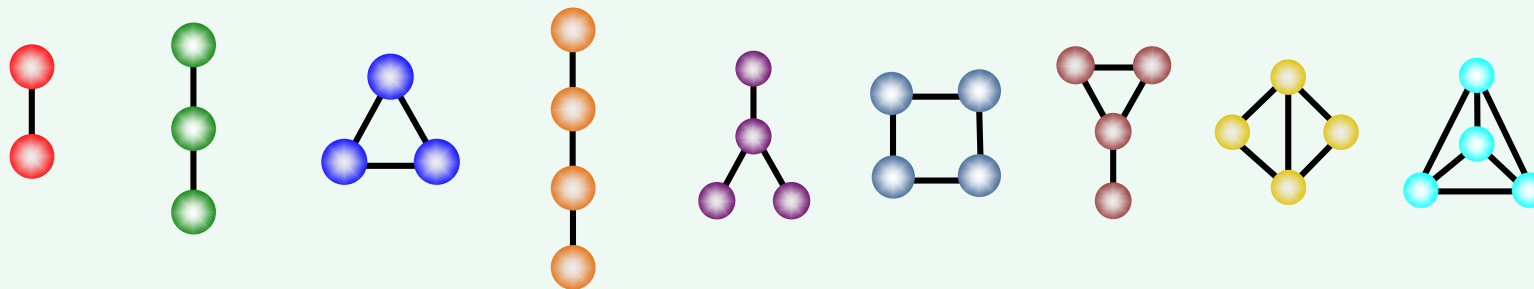
Workflow for generating the human full-length GPCR interactome.

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Graphlets

#### "Legos of Networks"



← → ⓘ https://scholar.google.ca/scholar?hl=en&q=graphlets&btnG=&as\_sdt= 90% Search ☆ 📁 📄 ⬇ 🏠 S ☰

Web Images More... Sign In

Google graphlets 🔍

**Scholar** About 1,490 results (0.05 sec) My Citations ▾

<p><b>Articles</b></p> <p>Case law</p> <p>My library</p> <p><b>Any time</b></p> <p>Since 2017</p> <p>Since 2016</p> <p>Since 2013</p> <p>Custom range...</p> <p><b>Sort by relevance</b></p> <p>Sort by date</p> <p><input checked="" type="checkbox"/> include patents</p> <p><input checked="" type="checkbox"/> include citations</p>	<p><u><b>Discovering discriminative graphlets for aerial image categories recognition</b></u></p> <p>L Zhang, Y Han, Y Yang, M Song... - IEEE Transactions on ..., 2013 - ieeeexplore.ieee.org</p> <p>Abstract: Recognizing aerial image categories is useful for scene annotation and surveillance. Local features have been demonstrated to be robust to image transformations, including occlusions and clutters. However, the geometric property of an aerial image (ie, ...)</p> <p>Cited by 111 Related articles All 10 versions Cite Save</p> <p><u><b>Probabilistic graphlet transfer for photo cropping</b></u></p> <p>L Zhang, M Song, Q Zhao, X Liu, J Bu... - IEEE Transactions on ..., 2013 - ieeeexplore.ieee.org</p> <p>Abstract: As one of the most basic photo manipulation processes, photo cropping is widely used in the printing, graphic design, and photography industries. In this paper, we introduce graphlets (ie, small connected subgraphs) to represent a photo's aesthetic features, and ...</p> <p>Cited by 95 Related articles All 16 versions Cite Save</p> <p><u><b>Model selection for social networks using graphlets</b></u></p> <p>J Janssen, M Hurshman, N Kalyaniwalla - Internet Mathematics, 2012 - Taylor &amp; Francis</p> <p>Several network models have been proposed to explain the link structure observed in online social networks. This paper addresses the problem of choosing the model that best fits a given real-world network. We implement a model-selection method based on unsupervised ...</p> <p>Cited by 25 Related articles All 10 versions Cite Save</p>	<p><u><b>Integrating local features into discriminative graphlets for scene classification</b></u></p> <p>L Zhang, W Bian, M Song, D Tao, X Liu - International Conference on ..., 2011 - Springer</p> <p>Abstract Scene classification plays an important role in multimedia information retrieval. Since local features are robust to image transformation, they have been used extensively for scene classification. However, it is difficult to encode the spatial relations of local features in ...</p> <p>Cited by 9 Related articles All 4 versions Cite Save</p> <p><u><b>Ego-centric graphlets for personality and affective states recognition</b></u></p> <p>S Teso, J Staiano, B Lepri, A Passerini... - Social Computing ( ..., 2013 - ieeeexplore.ieee.org</p> <p>Abstract: Do we tend to perceive ourselves more creative when surrounded by creative people? Or rather the opposite holds? Such information is very valuable to understand how to optimize work processes and boost people's productivity along with their happiness and ...</p> <p>Cited by 5 Related articles All 14 versions Cite Save</p> <p><u><b>From quasirandom graphs to graph limits and graphlets</b></u></p> <p>F Chung - Advances in Applied Mathematics, 2014 - Elsevier</p> <p>Abstract We generalize the notion of quasirandomness which concerns a class of equivalent properties that random graphs satisfy. We show that the convergence of a graph sequence under the spectral distance is equivalent to the convergence using the (normalized) cut ...</p> <p>Cited by 4 Related articles All 10 versions Cite Save</p>
--	--	---

# Overview

**Medicine: complex world of inter-connected entities**

## **1. Motivation**

## **2. New Methods – Examples: mine inter-connected data**

- i. **Single layer of omics data: Molecular networks** → function, disease
- ii. **Multiple layers of heterogeneous data:**
  - **Patient-centered data integration** → Precision medicine
  - Disease re-classification
  - Gene Ontology reconstruction
  - Network alignment

## **3. Vision**

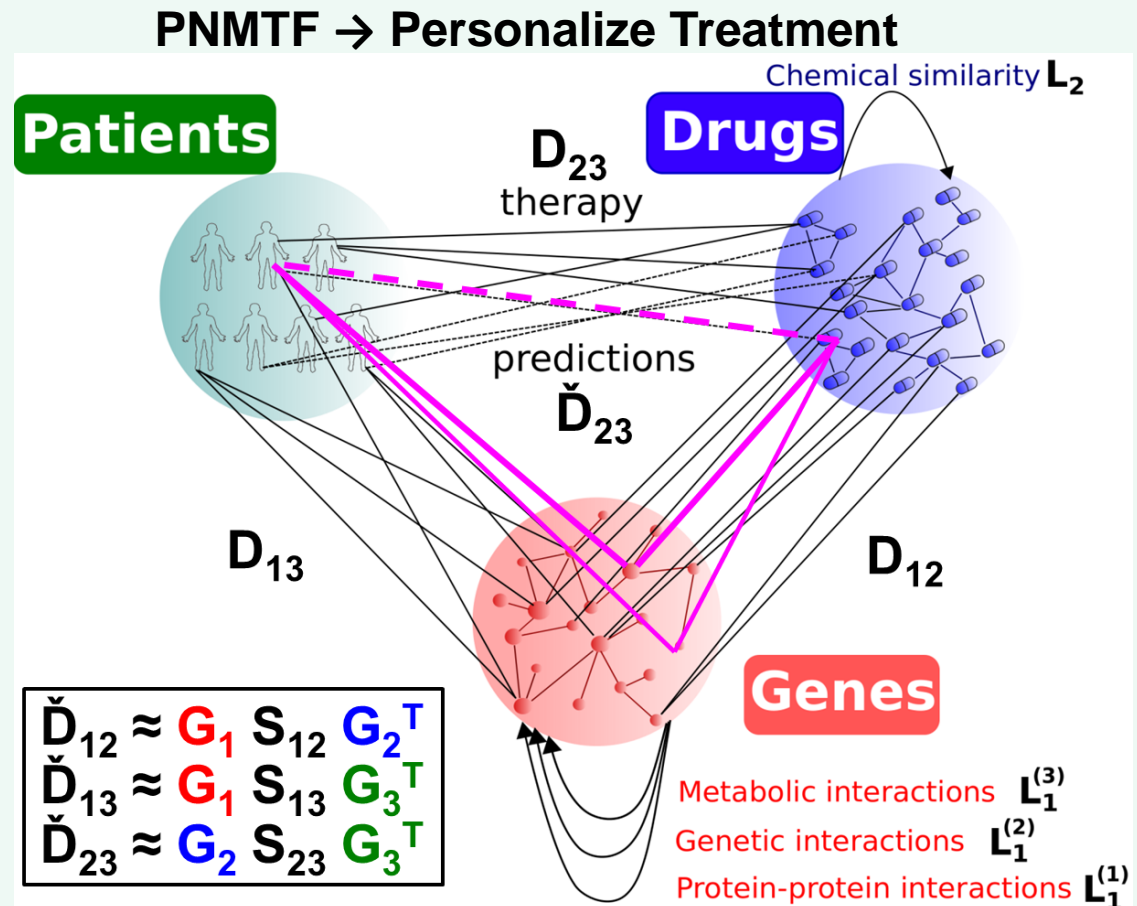
# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Multi-disciplinary, data-fusion methodology

#### Motivation:

- Captures all systems-level
- Captures how data relate
- **Mechanistic explanations**



$$\min\{\sum_{1 \leq i \leq j \leq p} [ ||W_{ij} \circ (D_{ij} - G_i S_{ij} G_j^T) ||^2 + \alpha ||S_{ij}||^2 + \alpha_i \text{tr}(G_i^T L_i G_i) + \alpha_j \text{tr}(G_j^T L_j G_j) ] : G_i, S_{ij} \geq 0\}$$

$\alpha ||S_{ij}||^2$  maintain sparsity of  $S_{ij}$ ,  $\alpha_i \text{tr}(G_i^T L_i G_i)$  and  $\alpha_j \text{tr}(G_j^T L_j G_j)$  adding prior knowledge (penalties),  $G_i, S_{ij} \geq 0$  is needed for cluster interpretation



# 2. Novel Methods

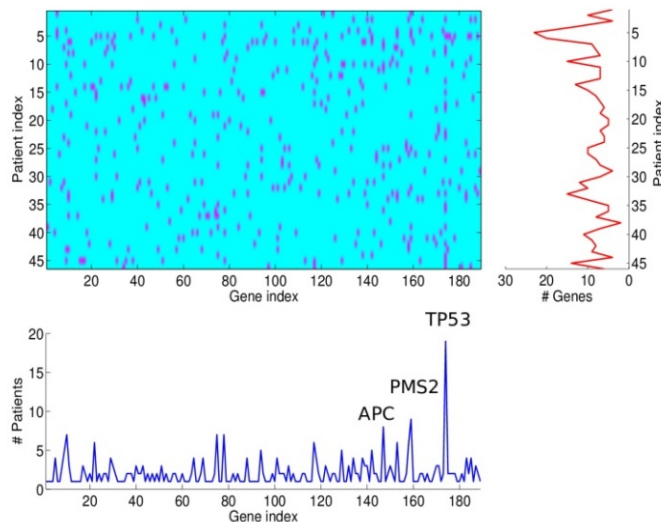
## Mine the Medical World of Inter-Connected Entities

### Multi-disciplinary, data-fusion methodology

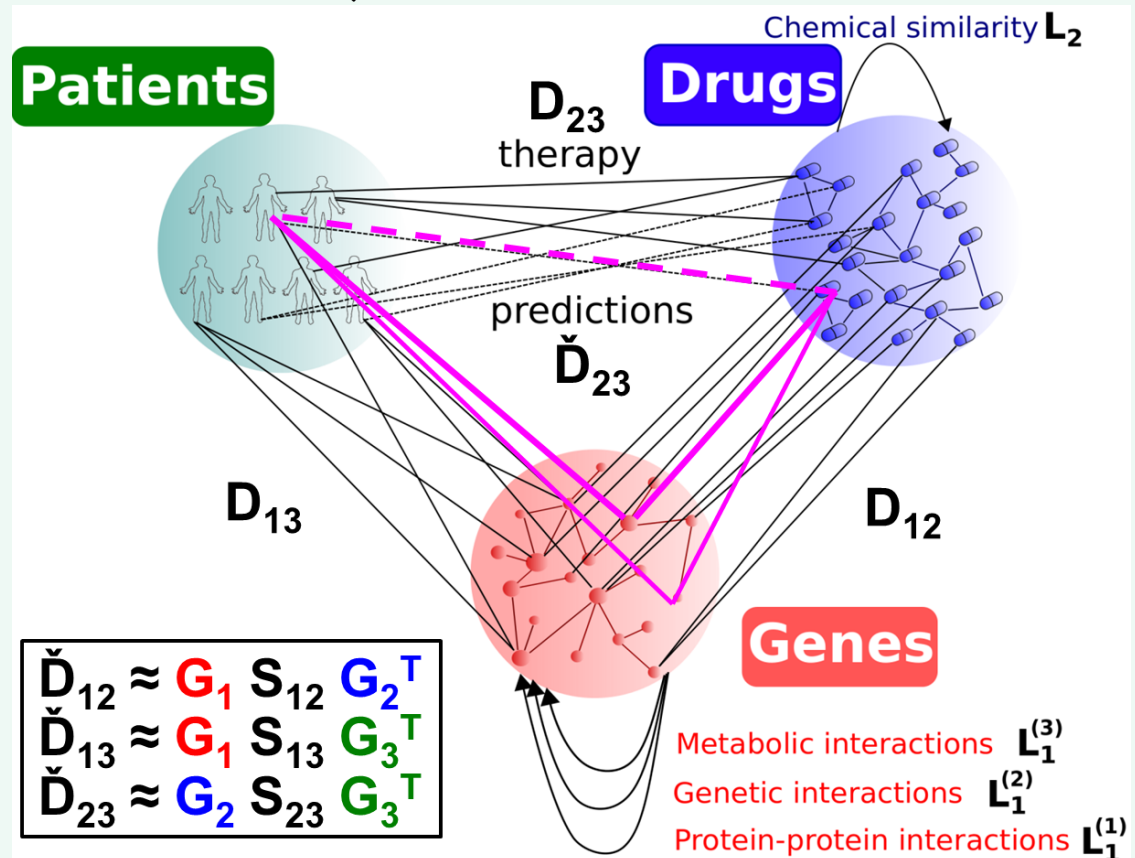
#### Motivation:

- Captures all systems-level
- Captures how data relate
- **Mechanistic explanations**

Next-Gen Sequencing



#### PNMTF → Personalize Treatment



$$\min\{\sum_{1 \leq i \leq j \leq p} [ ||W_{ij} \circ (D_{ij} - G_i S_{ij} G_j^T) ||^2 + \alpha ||S_{ij}||^2 + \alpha_i \text{tr}(G_i^T L_i G_i) + \alpha_j \text{tr}(G_j^T L_j G_j) ] : G_i, S_{ij} \geq 0\}$$

$\alpha ||S_{ij}||^2$  maintain sparsity of  $S_{ij}$ ,  $\alpha_i \text{tr}(G_i^T L_i G_i)$  and  $\alpha_j \text{tr}(G_j^T L_j G_j)$  adding prior knowledge (penalties),  $G_i, S_{ij} \geq 0$  is needed for cluster interpretation

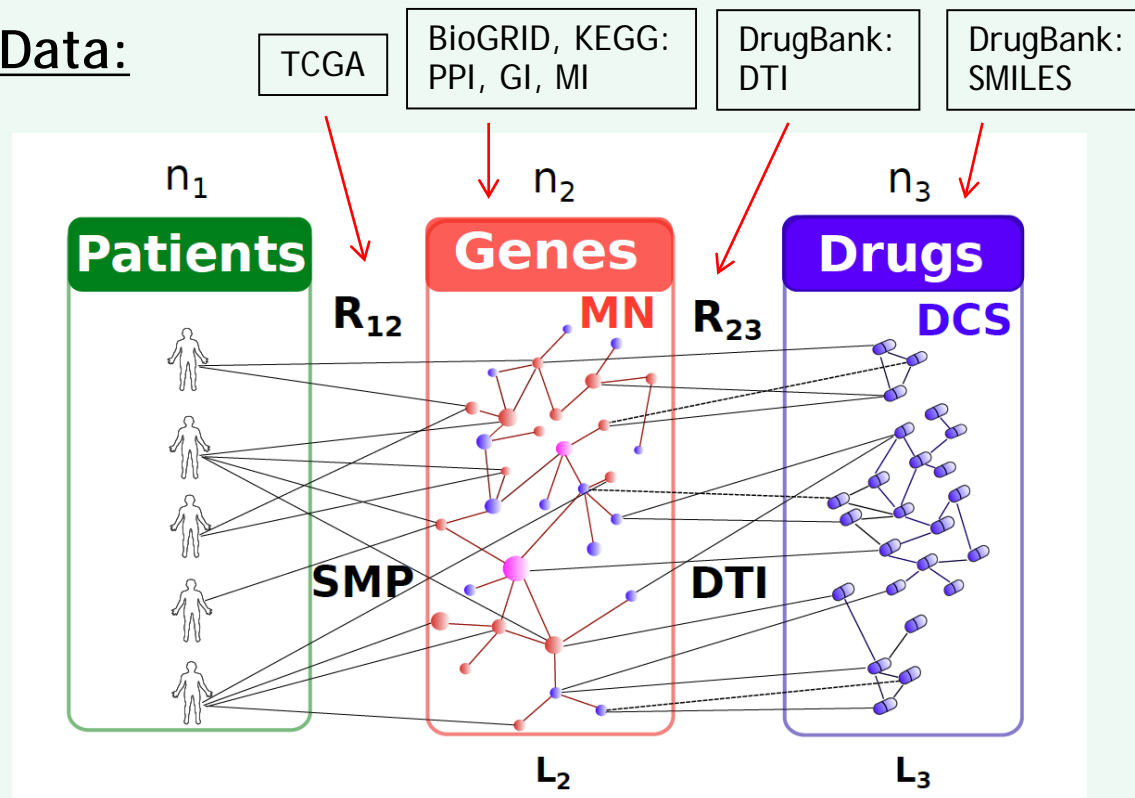
## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

#### Patient-Specific Data Fusion → Personalized Treatment

Co-clustering: **patients**, **genes** and **drugs**

Data:



#### 353 serous ovarian cancer patients from TCGA:

1. Patient stratification
2. Driver gene prediction
3. Drug repurposing

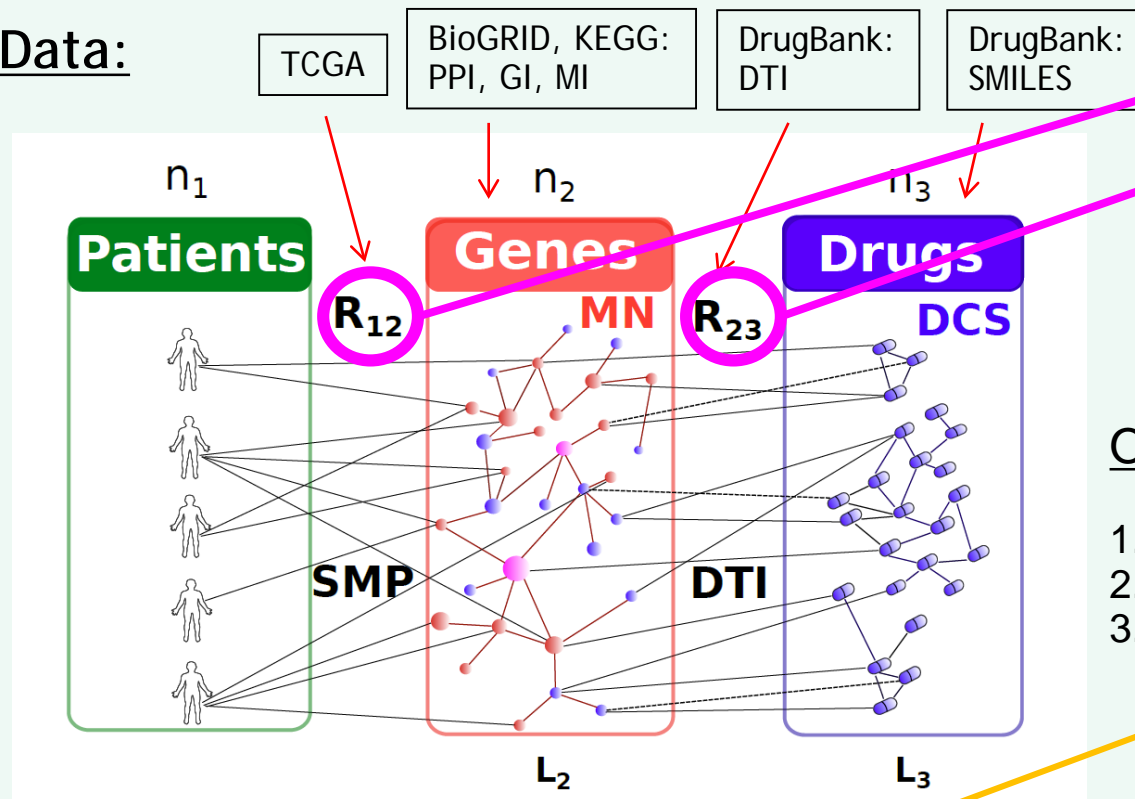
# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Patient-Specific Data Fusion → Personalized Treatment

Co-clustering: **patients**, **genes** and **drugs**

Data:



$$R_{12} \approx G_1 H_{12} G_2^T$$

$$R_{23} \approx G_2 H_{23} G_3^T$$

$k_1 \ll n_1$  – patient clusters  
 $k_2 \ll n_2$  – gene clusters  
 $k_3 \ll n_3$  – drug clusters  
 $G_1, G_2$  and  $G_3$  are cluster indicator matrices

Ovarian cancer patients:

1. Patient stratification →  $\hat{C}_1$
2. Driver gene prediction →  $\hat{C}_2$
3. Drug repurposing →  $\hat{R}_{23}$

$$\min_{G_i \geq 0, 1 \leq i \leq 3} J = \min_{G_i \geq 0, 1 \leq i \leq 3} \left[ \| R_{12} - G_1 H_{12} G_2^T \|_F^2 + \| R_{23} - G_2 H_{23} G_3^T \|_F^2 + \right. \\ \left. tr(G_2^T L_2 G_2) + tr(G_3^T L_3 G_3) \right]$$

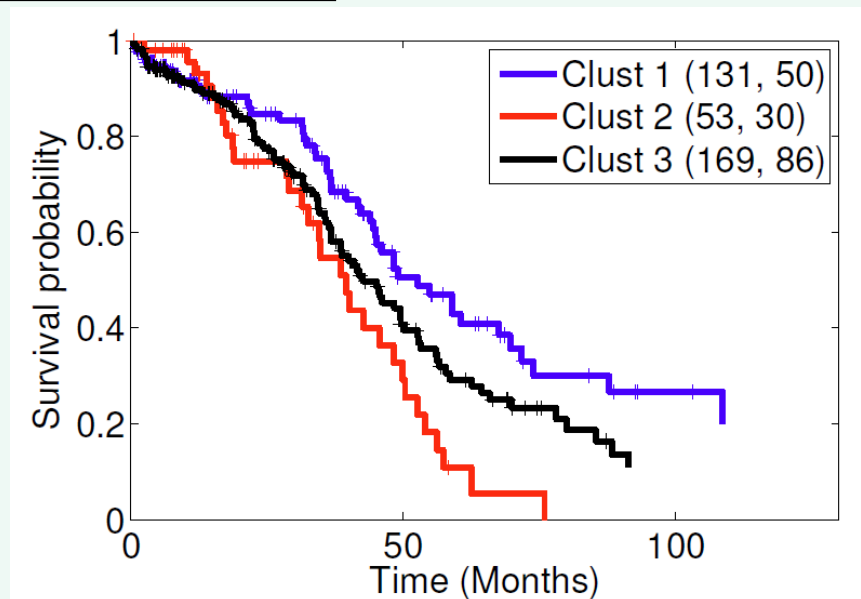


# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Patient-Specific Data Fusion → Personalized Treatment

Some results:



Kaplan-Meier survival curves for 3 patient groups found by GNMTF (log-rank p-val =  $5.3 \times 10^{-3}$ )

$$\mathbf{R}_{12} \approx \mathbf{G}_1 \mathbf{H}_{12} \mathbf{G}_2^T$$

$$\mathbf{R}_{23} \approx \mathbf{G}_2 \mathbf{H}_{23} \mathbf{G}_3^T$$

$k_1 \ll n_1$  – patient clusters

$k_2 \ll n_2$  – gene clusters

$k_3 \ll n_3$  – drug clusters

$\mathbf{G}_1$ ,  $\mathbf{G}_2$  and  $\mathbf{G}_3$  are cluster indicator matrices

Ovarian cancer patients:

1. **Patient stratification** →  $\hat{\mathbf{C}}_1$
2. Driver gene prediction →  $\hat{\mathbf{C}}_2$
3. Drug repurposing →  $\hat{\mathbf{R}}_{23}$

$$\min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} J = \min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} \left[ \|\mathbf{R}_{12} - \mathbf{G}_1 \mathbf{H}_{12} \mathbf{G}_2^T\|_F^2 + \|\mathbf{R}_{23} - \mathbf{G}_2 \mathbf{H}_{23} \mathbf{G}_3^T\|_F^2 + \text{tr}(\mathbf{G}_2^T \mathbf{L}_2 \mathbf{G}_2) + \text{tr}(\mathbf{G}_3^T \mathbf{L}_3 \mathbf{G}_3) \right]$$

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Patient-Specific Data Fusion → Personalized Treatment

Some results: ~40% of our 809 predicted driver genes in CCGD, Census, or IntOGen

New driver	Known drivers	Score	DB
ADAM32	BMPR2	1.000	—
REG1P	CLASP2	1.000	—
PCDHA2	CHD4	1.000	—
NCR1	BMPR2	1.000	—
USPL1	CLASP2	1.000	—
GDPD3	DDX5	1.000	—
LECT1	CLASP2	1.000	CCGD
IL25	CDK12, CCAR1	0.975	—
BAK1	ATRX, TFDP1, NDRG1	0.967	—
MOGAT2	ATRX, TFDP1, NDRG1	0.967	—
CHAF1A	ATRX, TFDP1, NDRG1	0.967	CCGD
PITX2	ATRX, TFDP1, NDRG1	0.967	—
SIN3B	ATRX, TFDP1, NDRG1	0.967	—
RPL30	ATRX, TFDP1, NDRG1	0.967	—
GRWD1	ATRX, TFDP1, NDRG1	0.967	—
SNAI1	ATRX, TFDP1, NDRG1	0.967	CCGD
RBMXP4	ATRX, TFDP1, NDRG1	0.967	—
CPNE7	ATRX, TFDP1, NDRG1	0.967	—
HIPK3	ATRX, TFDP1, NDRG1	0.967	CCGD
EPOR	ATRX, TFDP1, NDRG1	0.967	CCGD

↔ TGFs, cell proliferation & progression

↔ proliferation, migration, anti-apoptosis; prognosis markers

#### Ovarian cancer patients:

1. Patient stratification →  $\hat{C}_1$
2. **Driver gene prediction** →  $\hat{C}_2$
3. Drug repurposing →  $\hat{R}_{23}$

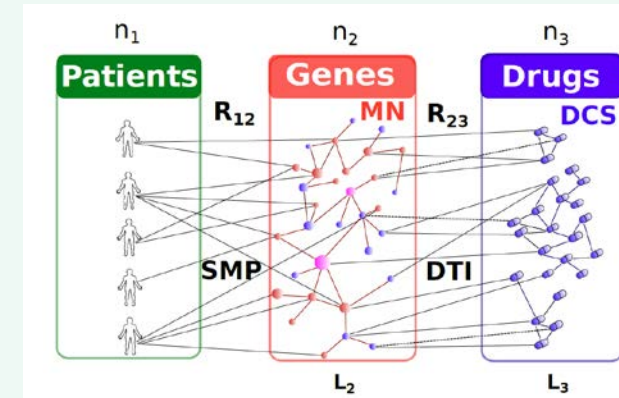
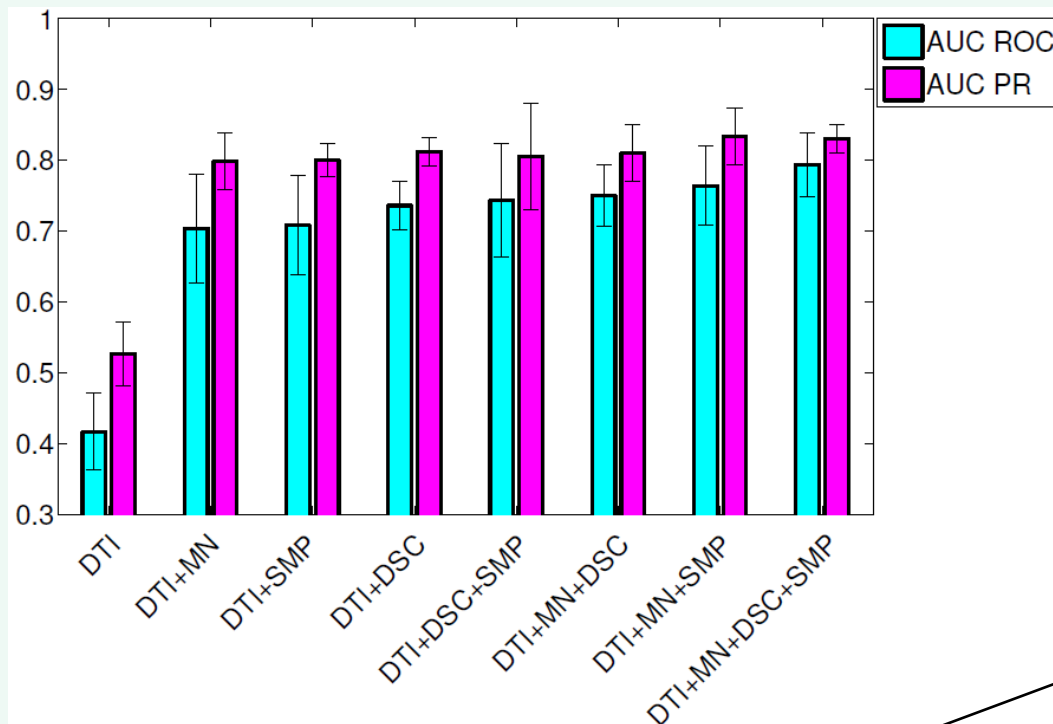
$$\min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} J = \min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} \left[ \|\mathbf{R}_{12} - \mathbf{G}_1 \mathbf{H}_{12} \mathbf{G}_2^T\|_F^2 + \|\mathbf{R}_{23} - \mathbf{G}_2 \mathbf{H}_{23} \mathbf{G}_3^T\|_F^2 + \right. \\ \left. tr(\mathbf{G}_2^T \mathbf{L}_2 \mathbf{G}_2) + tr(\mathbf{G}_3^T \mathbf{L}_3 \mathbf{G}_3) \right]$$

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Patient-Specific Data Fusion → Personalized Treatment

Some results: 5-fold cross validation, average AUC: ROC and PR



Ovarian cancer patients:

1. Patient stratification →  $\hat{C}_1$
2. Driver gene prediction →  $\hat{C}_2$
3. Drug repurposing →  $\hat{R}_{23}$

$$\min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} J = \min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} \left[ \left\| \mathbf{R}_{12} - \mathbf{G}_1 \mathbf{H}_{12} \mathbf{G}_2^T \right\|_F^2 + \left\| \mathbf{R}_{23} - \mathbf{G}_2 \mathbf{H}_{23} \mathbf{G}_3^T \right\|_F^2 + \right. \\ \left. tr(\mathbf{G}_2^T \mathbf{L}_2 \mathbf{G}_2) + tr(\mathbf{G}_3^T \mathbf{L}_3 \mathbf{G}_3) \right]$$

## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

#### Patient-Specific Data Fusion → Personalized Treatment

Some results: 37% of our ~225K predicted DTIs confirmed in MATADOR or CTD

Gene	Drug	Score	Clusters	DB
KIT	ATP	0.873	1, 2, 3	–
GABRQ	Adinazolam	0.808	1	M
GABRQ	Fludiazepam	0.808	1	M
GABRQ	Cinolazepam	0.809	1	M
GABRQ	Clotiazepam	0.809	1	M
HTR2A	Dopamine	0.809	1, 3	C, M
<b>GRIN3A</b>	<b>Pethidine</b>	0.801	1, 2	–
CACNA2D1	Verapamil	0.761	1, 3	M
PDGFRB	ATP	0.724	1, 2	–
KDR	ATP	0.724	1, 3	C
HTR1A	Mirtazapine	0.720	1, 2	C, M
GABRA6	Adinazolam	0.688	1	M
GABRA6	Fludiazepam	0.688	1	M
GABRA6	Cinolazepam	0.688	1	M
GABRA6	Clotiazepam	0.688	1	M
GABRA4	Adinazolam	0.687	1, 3	M
GABRA4	Fludiazepam	0.687	1, 3	M
GABRA4	Cinolazepam	0.687	1, 3	M
GABRA4	Clotiazepam	0.687	1, 3	M
CACNA1D	Magnesium Sulfate	0.676	1, 2, 3	M

Ovarian cancer patients:

1. Patient stratification →  $\hat{C}_1$
2. Driver gene prediction →  $\hat{C}_2$
3. **Drug repurposing** →  $\hat{R}_{23}$

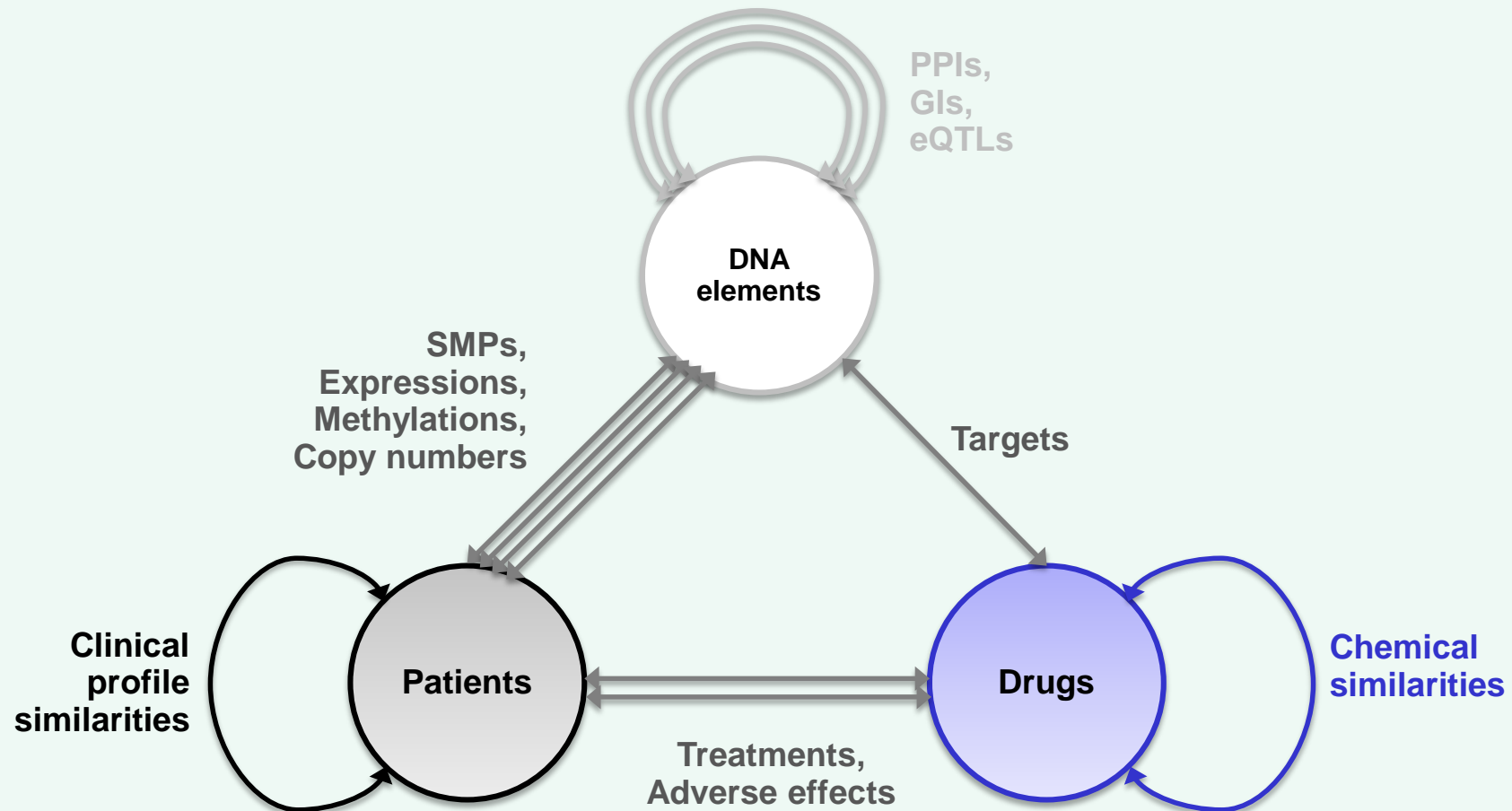
$$\min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} J = \min_{\mathbf{G}_i \geq 0, 1 \leq i \leq 3} \left[ \left\| \mathbf{R}_{12} - \mathbf{G}_1 \mathbf{H}_{12} \mathbf{G}_2^T \right\|_F^2 + \left\| \mathbf{R}_{23} - \mathbf{G}_2 \mathbf{H}_{23} \mathbf{G}_3^T \right\|_F^2 + \right. \\ \left. tr(\mathbf{G}_2^T \mathbf{L}_2 \mathbf{G}_2) + tr(\mathbf{G}_3^T \mathbf{L}_3 \mathbf{G}_3) \right]$$



## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

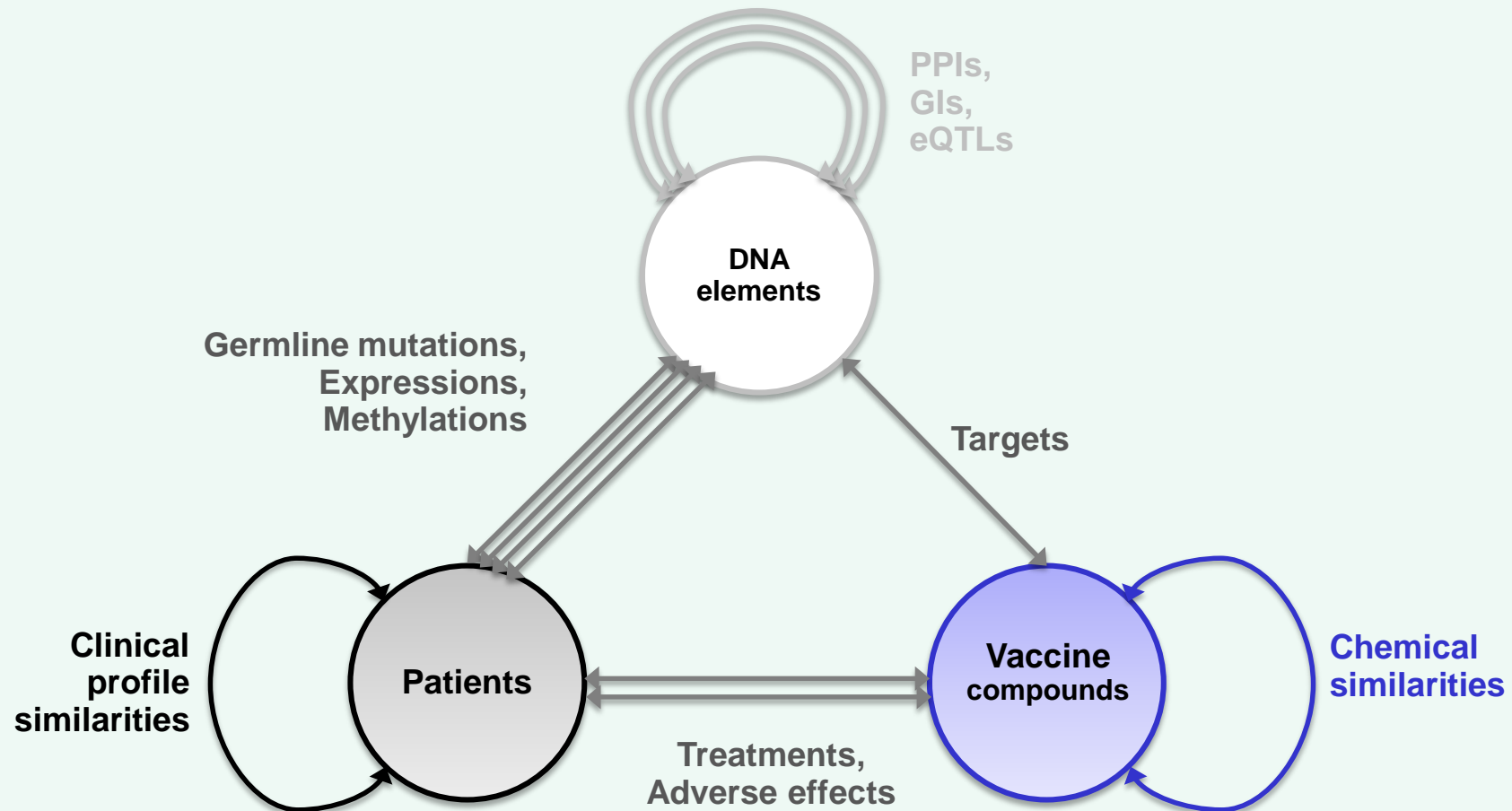
Patient-Specific Data Fusion → Personalized Treatment



## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

Patient-Specific Data Fusion → Personalized Treatment



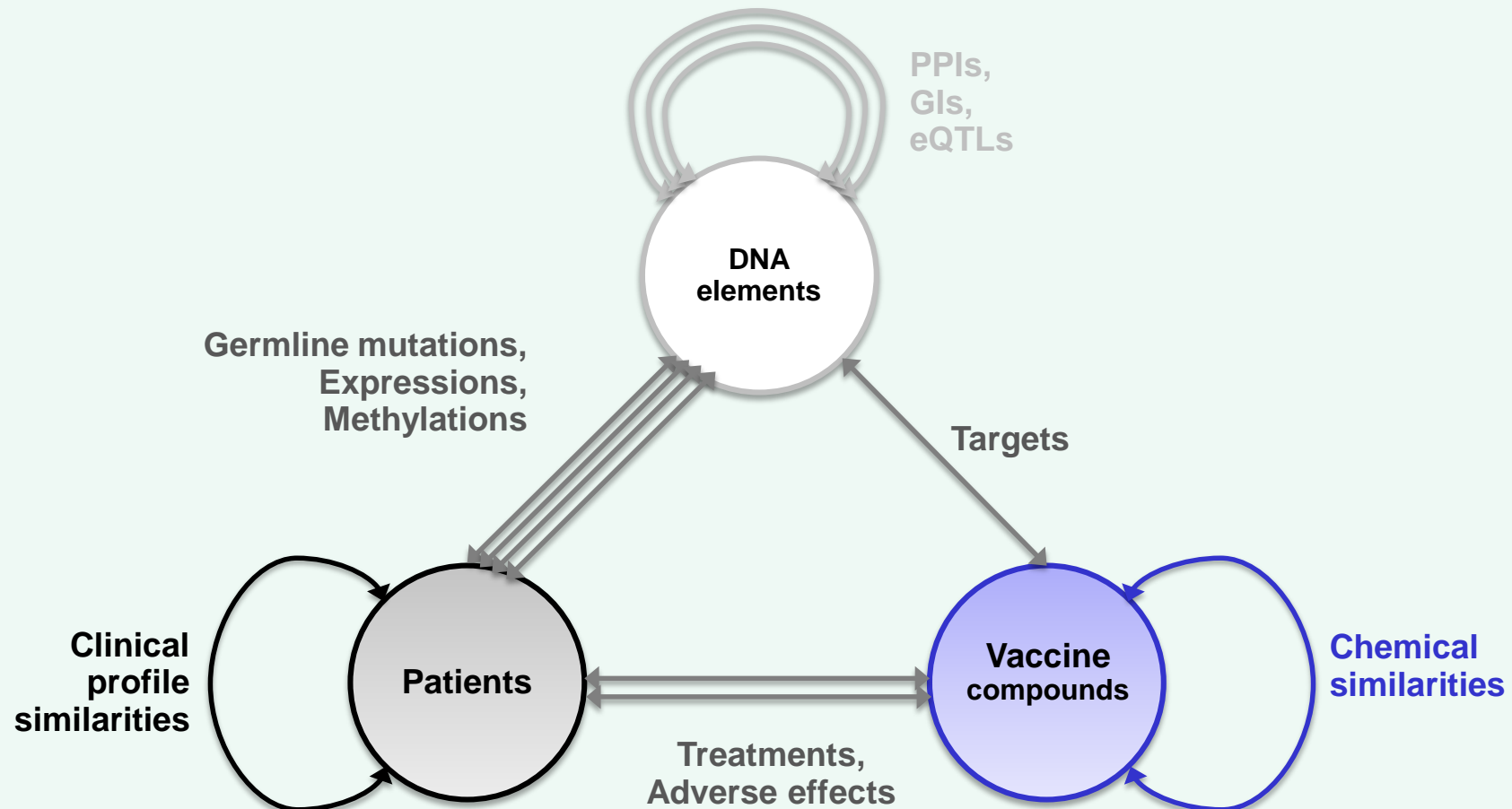
#### ❑ Systems vaccinology

- With Dr. Nuria Izquierdo, IGTP IrsiCaixa, Badalona
- Scientific Advisory Board of the Helmholtz Centre for Infection Research (HZI / Braunschweig, Germany)

## 2. Novel Methods

### Mine the Medical World of Inter-Connected Entities

Patient-Specific Data Fusion → Personalized Treatment



#### Obstacles:

##### 1. Different NP-hard continuous optimization problem:

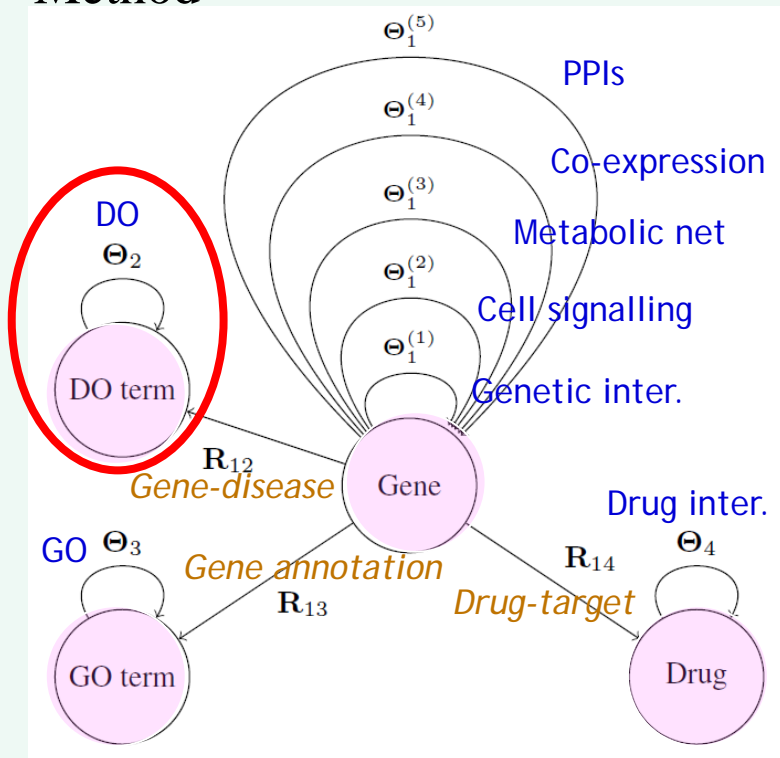
- propose objective function,
- optimization solver — prove convergence and correctness

##### 2. Optimization is slow → HPC

# Mine the Medical World of Inter-Connected Entities

# Disease Classification from Systems-Level Molecular Data

## Method



## 4 Objects: Genes, GO terms, DO terms, Drugs

Constraints:  $\Theta_i$  (*network topology, ontology relations*)

Relation matrices:  $R_{ij}$

## Some Results:

→ 14 disease-disease associations currently not present in DO:

- evidence for their relationships through *comorbidity* data and *literature curation*

→ GI the most important predictor  
of a link between diseases, despite small

→ Omission of any one of the included data sources reduces prediction quality

- Importance of systems-level data fusion

→ DO  $\cap$  disease class → 80% DO from *only* network data



# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Disease Classification from Systems-Level Molecular Data

- Co-clustering GO terms, DO terms, Genes and Drugs under pairwise constraints:

$$\text{Input: } \mathbf{R} = \begin{bmatrix} 0 & \mathbf{R}_{12} & \mathbf{R}_{13} & \mathbf{R}_{14} \\ \mathbf{R}_{12}^T & 0 & 0 & 0 \\ \mathbf{R}_{13}^T & 0 & 0 & 0 \\ \mathbf{R}_{14}^T & 0 & 0 & 0 \end{bmatrix} \quad \Theta = \begin{bmatrix} \Theta_1^{(t)} & 0 & 0 & 0 \\ 0 & \Theta_2 & 0 & 0 \\ 0 & 0 & \Theta_3 & 0 \\ 0 & 0 & 0 & \Theta_4 \end{bmatrix}$$

$$\text{Output: } \mathbf{S} = \begin{bmatrix} 0 & \mathbf{S}_{12} & \mathbf{S}_{13} & \mathbf{S}_{14} \\ \mathbf{S}_{21} & 0 & 0 & 0 \\ \mathbf{S}_{31} & 0 & 0 & 0 \\ \mathbf{S}_{41} & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{G} = \begin{bmatrix} \mathbf{G}_1 & 0 & 0 & 0 \\ 0 & \mathbf{G}_2 & 0 & 0 \\ 0 & 0 & \mathbf{G}_3 & 0 \\ 0 & 0 & 0 & \mathbf{G}_4 \end{bmatrix}$$

- Minimizing Frobenious distance between  $R_{ij}$  and  $G_i S_{ij} G_j^T$ , for all relation matrices:

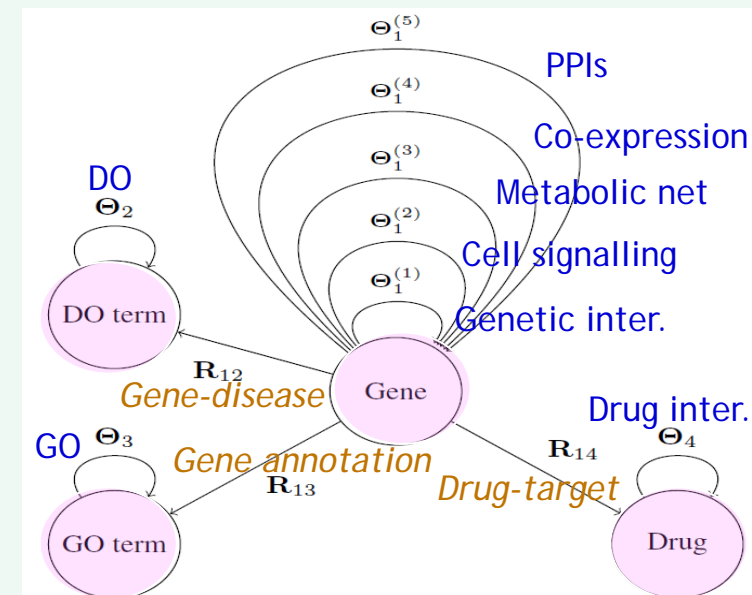
- $i = \{\text{Genes}\}, j = \{\text{DO terms}, \text{GO terms}, \text{Drugs}\}$
- $G_i$  is a cluster indicator matrix for data type  $i$  (genes, DO terms, GO terms and Drugs)

with additional penalty terms:

$$\min_{\mathbf{G} \geq 0} J = \min_{\mathbf{G} \geq 0} \left[ \|\mathbf{R} - \mathbf{G} \mathbf{S} \mathbf{G}^T\|_F^2 + \left( \sum_{t=1}^5 \text{tr}(\mathbf{G}^T \Theta^{(t)} \mathbf{G}) \right) \right]$$

- Interested in  $\mathbf{G}_2$  (DO terms)

- used for cluster assignment and inferring new disease associations from clusters

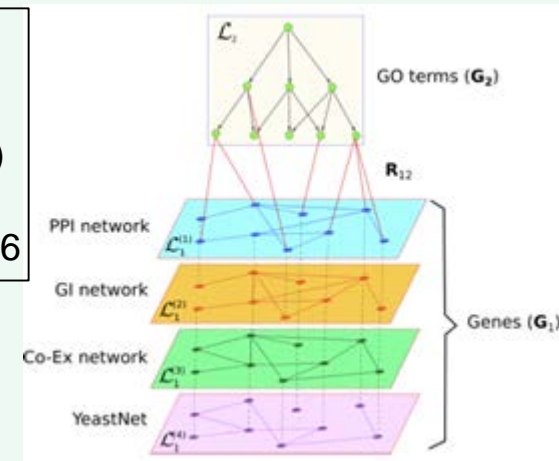


# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Gene Ontology from Systems-Level Molecular Data

- Outperform Dutkowski *et al.* [2013]
- 96% of GO reconstructed!
- Correct assignment of GO terms to genes (3-fold cross-validation,  $AUC=0.874 \pm 0.002$ )
- Graphlets improve results
- **Validated biologically** by Bonne's yeast Genetic Interaction profile data, *Science*, 2016



→ Optimization problem which minimizes  $\|R_{12} - G_1 S_{12} G_2^T\|_F^2$   
under the guidance of *pairwise constraints*  
(*connectivity* and *GDV similarity*) between genes in networks:

$$\min_{G_1 \geq 0, G_2 \geq 0} J = \min_{G_1 \geq 0, G_2 \geq 0} \left[ \|R_{12} - G_1 S_{12} G_2^T\|_F^2 + \left( \sum_{i=1}^4 \text{tr}(G_1^T L_1^{(i)} G_1) \right) + \left( \sum_{i=1}^4 \text{tr}(G_1^T \Lambda_1^{(i)} G_1) \right) + \text{tr}(G_2^T L_2 G_2) \right]$$

→ using topology of molecular networks as constraints (penalty terms) in this optimization problem:

→  $L_1^{(i)}$  is Laplacian of *adjacency matrix* of a molecular network  $i=1,2,3,4$ :

$L_1^{(i)} = D^i - A^i$ ,  $D^i$  is diagonal matrix of degrees (row summation of  $A^i$ ),  $A^i$  is adjacency matrix

→  $\Lambda_1^{(i)}$  are Laplacians of *GDV similarity matrices* over all genes for each molecular network  $i$ :

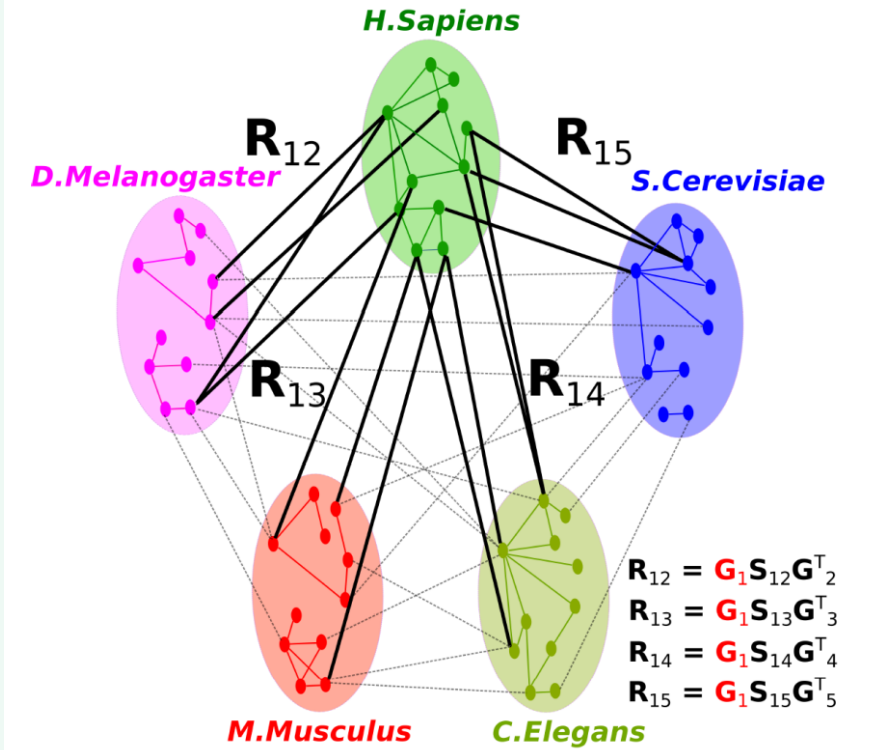
$\Lambda_1^{(i)} = D^i - \sigma^{(i)}$ ,  $D^i$  is diagonal matrix of row summation of  $\sigma^{(i)}$ ,  $\sigma^{(i)}$  is binary GDV similarity matrix (containing only significantly similar gene/protein pairs)

→  $L_2$  is Laplacian of *Gene Ontology graph*

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Multiple Network Alignment: Fuse



**Algorithm 1** Approximate maximum weight  $k$ -partite matching.

**Input**  $G = (\bigcup_{i=1}^k V_i, E, W)$   
**for**  $i = \{2, \dots, k\}$  **do**  
    Find maximum weight bipartite matching  $F_{1,i}$  of  $G[V_1, V_i]$   
    Construct  $G_{1i}$ , the merge of  $V_1$  and  $V_i$  from  $G$  along  $F_{1,i}$   
    Set  $G = G_{1i}$ , and relabel  $V_{1i}$  as  $V_1$   
 $C = \{\emptyset\}$   
**for** each merged node  $u$  in  $V_1$  **do**  
    Cluster  $C_u$  is the set of nodes that are merged into  $u$   
    Add  $C_u$  to  $C$   
**Output**  $C$

We use a block-based representation of relation ( $\mathbf{R}$ ) and Laplacian ( $\mathbf{L}$ ) matrices and matrix factors ( $\mathbf{S}$  and  $\mathbf{G}$ ) for our 5 PPI networks as follows:

$$\mathbf{R} = \begin{bmatrix} 0 & \mathbf{R}_{12} & \dots & \mathbf{R}_{15} \\ \mathbf{R}_{12}^T & 0 & \dots & \mathbf{R}_{25} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{R}_{15}^T & \mathbf{R}_{25}^T & \dots & 0 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} \mathbf{L}_1 & 0 & \dots & 0 \\ 0 & \mathbf{L}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{L}_5 \end{bmatrix};$$

$$\mathbf{S} = \begin{bmatrix} 0 & \mathbf{S}_{12} & \dots & \mathbf{S}_{15} \\ \mathbf{S}_{12}^T & 0 & \dots & \mathbf{S}_{25} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{S}_{15}^T & \mathbf{S}_{25}^T & \dots & 0 \end{bmatrix}, \quad \mathbf{G} = \begin{bmatrix} \mathbf{G}_1 & 0 & \dots & 0 \\ 0 & \mathbf{G}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{G}_5 \end{bmatrix}$$

To simultaneously factorize all relation matrices,  $\mathbf{R}_{ij} \approx \mathbf{G}_i \mathbf{S}_{ij} \mathbf{G}_j^T$ ,  $0 \leq i, j \leq 5$ , under the constraints of PPI networks, we minimize the following objective function:

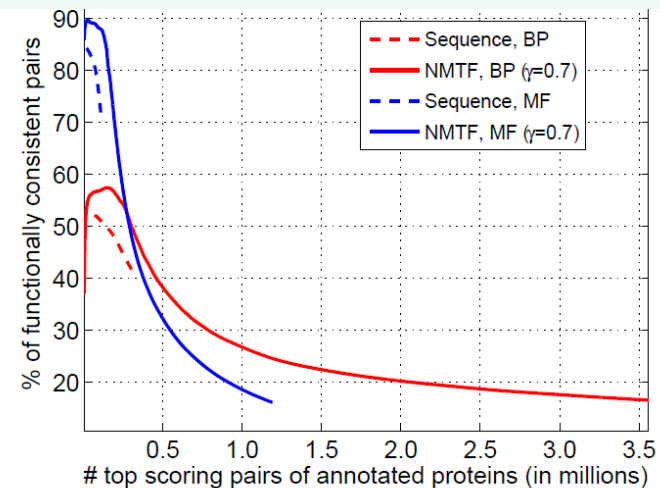
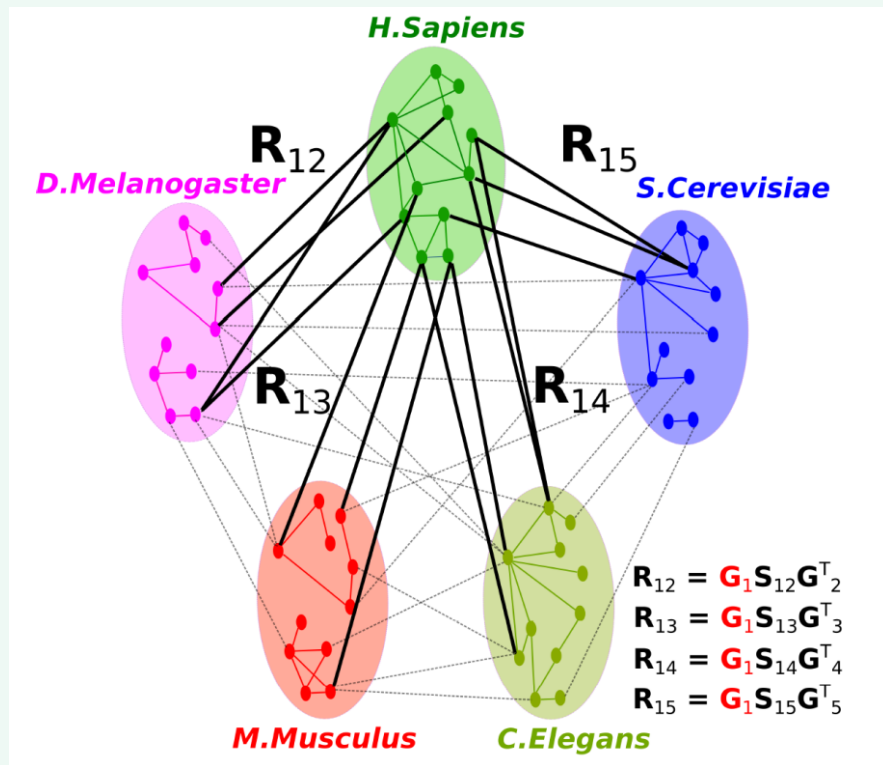
$$\min_{\mathbf{G} \geq 0} J = [\|\mathbf{R} - \mathbf{G} \mathbf{S} \mathbf{G}^T\|_F^2 + \gamma \text{Tr}(\mathbf{G}^T \mathbf{L} \mathbf{G})] \quad (2)$$

where  $\text{Tr}$  denotes the trace of a matrix and  $\gamma$  is a regularization parameter which balances the influence of network topologies in reconstruction of the relation matrix. The second term of equation 2 is the penalization term.

# 2. Novel Methods

## Mine the Medical World of Inter-Connected Entities

### Multiple Network Alignment: Fuse



**Fig. 2. Functional consistency of NMTF associations.** For both NMTF associations and sequence similarity of protein pairs, we plot the cumulative number of protein pairs with both proteins annotated (x-axis) against the percentages of them sharing GO terms (y-axis). Biological process (BP) and molecular function (MF) annotations are considered separately.



# Overview

**Medicine: complex world of inter-connected entities**

## **1. Motivation**

## **2. New Methods – Examples: mine inter-connected data**

i. Single layer of omics data: Molecular networks → function, disease

ii. Multiple layers of heterogeneous data:

- Patient-centered data integration → Precision medicine
- Disease re-classification
- Gene Ontology reconstruction
- Network alignment

## **3. Vision**

# 3. Vision

## **Biomedical Data: complex system of heterogeneous interacting entities**

- Large
- Heterogeneous
- Highly dimensional
- Growing Complexity
- Noisy
- Dynamic
- Different time and space scales

- World Economic Forum in Davos 2016:
    - “Big data” potential to transform medicine
    - Make it more effective due to increased life expectancy and exposure to environmental risks
  - *Nature* Insight and Outlook of 2015 and 2016
- 
- I was awarded 2014 BCS Roger Needham Award in recognition of “the potential my work and research have to revolutionize health and pharmaceuticals”

# 3. Vision

## Biomedical Data: complex system of heterogeneous interacting entities

- Large
  - Heterogeneous
  - Highly dimensional
  - Growing Complexity
  - Noisy
  - Dynamic
  - Different time and space scales
- Each type: *limited*, but *complementary* information
  - **Seek principled, joint organization and mining within the same framework**

- World Economic Forum in Davos 2016:
    - “Big data” potential to transform medicine
    - Make it more effective due to increased life expectancy and exposure to environmental risks
  - *Nature* Insight and Outlook of 2015 and 2016
- I was awarded 2014 BCS Roger Needham Award in recognition of “the potential my work and research have to revolutionize health and pharmaceuticals”

**€2M ERC Consolidator Grant for 2018-2023**  
**Title: “Integrated Connectedness for a New Representation of Biology”**

# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

1. Conceptual

2. Methodological



# 3. Vision

## Holistically Mine All Available Data

### → Paradigm shifts

#### 1. Conceptual

Do not analyze single data type in isolation of others (e.g., sequence align.)

- **Analyze all types of data within a single framework**
- **New, bottom-up, data-driven biological concepts**
  - Elucidate that a cell may be governed by yet undiscovered principles of life
  - Point to ways to re-think biology and approaches to medicine

# 3. Vision

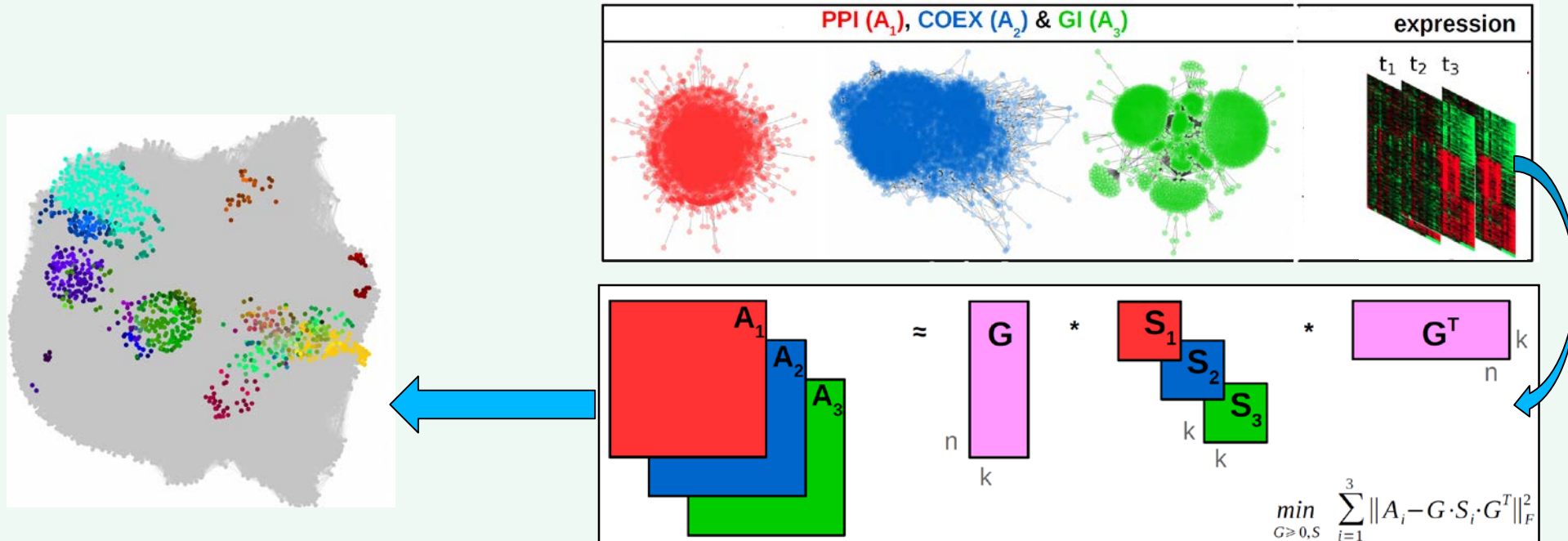
## Holistically Mine All Available Data

### → Paradigm shifts

#### 1. Conceptual

Do not analyze single data type in isolation of others (e.g., sequence align.)

- Introduce a concept of an **“Integrated Cell (iCell)”**

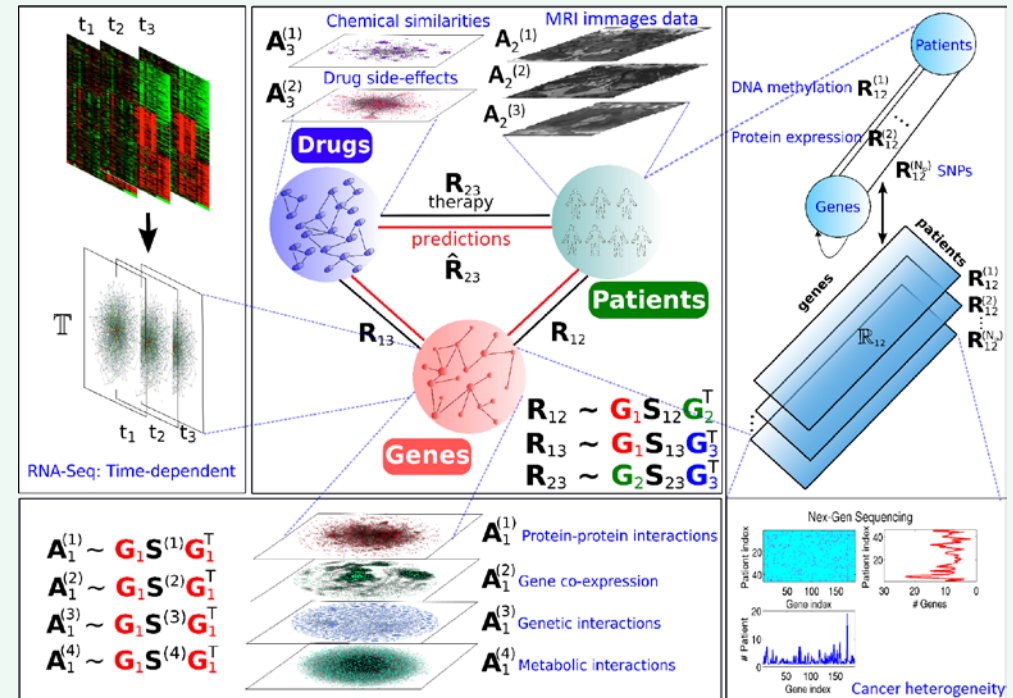


# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

## 2. Methodological



# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

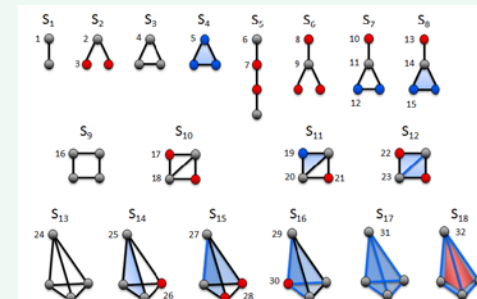
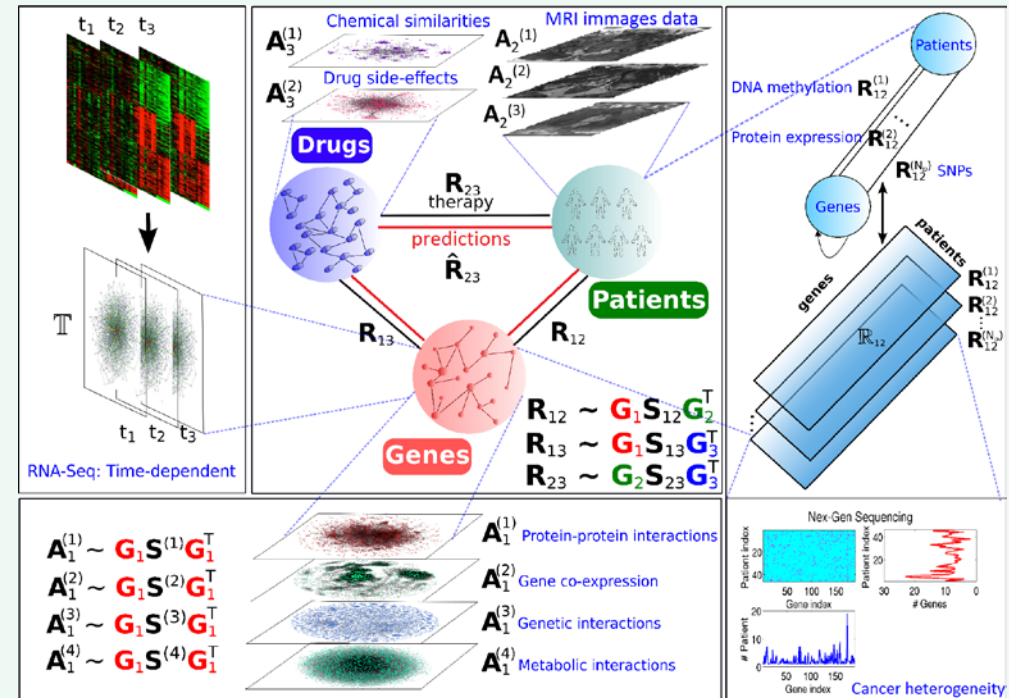
### 2. Methodological

- **Mathematical formalisms**

- Capture multi-scale organization
- Dynamics, stochasticity of the data,...

- E.g., multiplex networks, hypergraphs, simplicial complexes ...

- **Algorithms to compute and extract information from those formalisms**



T. Gaudelet, N. Malod-Dognin and **N. Przulj**, "Higher order molecular organisation as a source of biological function," *Bioinformatics*, ECCB'18

N. Malod-Dognin and **N. Przulj**, "Functional geometry of protein-protein interaction networks," arXiv:1804.04428, 2018

Noël Malod-Dognin, Julia Petschnigg, Sam F. L. Windels, Janez Povh, Harry Hemmingway, Robin Ketteler and **Nataša Pržulj**, "iCell: integrated cells uncover new cancer genes," *Nature Communications*, 2019



# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

### 2. Methodological

- **Mathematical formalisms**

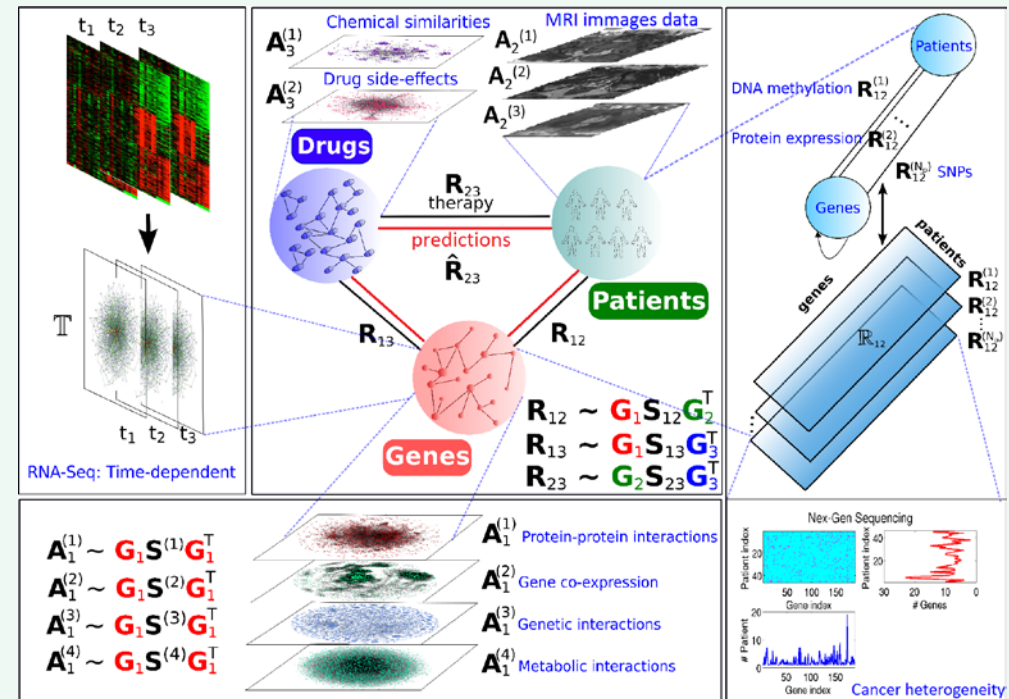
- Capture multi-scale organization
- Dynamics, stochasticity of the data,...

- E.g., multiplex networks, hypergraphs, simplicial complexes ...

- **Algorithms** to compute and **extract information** from those formalisms

How: e.g.

- Utilize dependencies in local network topology (orbits) — data set dependent
- Uncover latent low-dimensional structure of data
- Exploit structure for developing efficient toolsets for particular data



# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

### 2. Methodological

- **Mathematical formalisms**

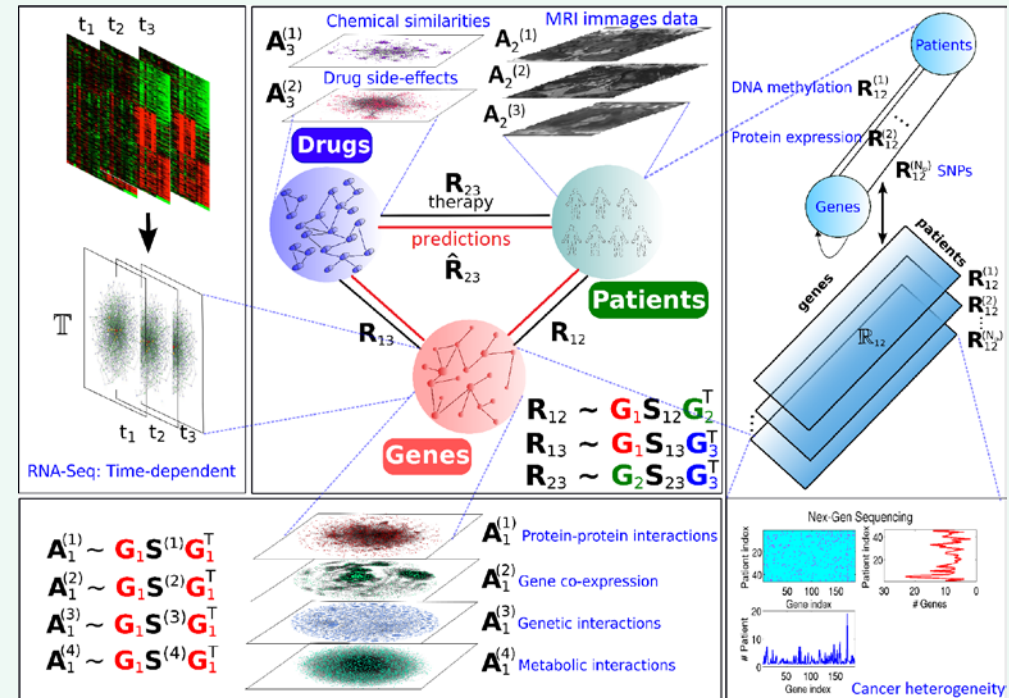
- Capture multi-scale organization
- Dynamics, stochasticity of the data,...
- E.g., multiplex networks, hypergraphs, simplicial complexes ...

- **Algorithms** to compute and **extract information** from those formalisms

**Computational issues** remain to be addressed, arising from intractability:

- large sizes, complexity, heterogeneity, noisiness, and
- different time and space scales of the data

**“Embedded” data scientists:** problem-specific heuristic methods, HPC



# 3. Vision

## Holistically Mine All Available Data

→ Paradigm shifts

**Guided by Needs of Biomedical Collaborators and Industry**

**E.g.:**

- Cancer
- Rare genetic diseases
- Viral medicines
- JnJ
- GSK
- Medium, start-ups, ...

# Acknowledgements

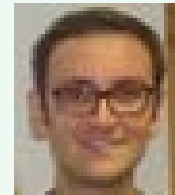
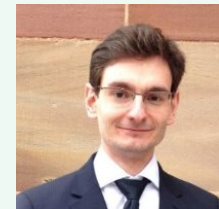
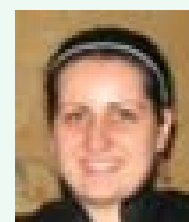
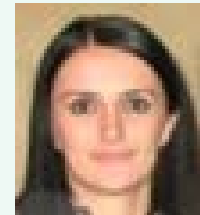
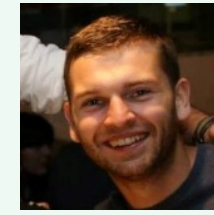


## ➤ Funding:



## ➤ Group members (present and past):

1. Dr. Noel Malod-Dogning
2. Dr. Julia Petschnigg
3. Dr. Chhedi Gupta
4. Sam Windels
5. Thomas Gaudet
6. Dr. Omer Yaveroglu
7. Prof. Tijana Milenković
8. Dr. Oleksii Kuchaiev
9. Dr. Vesna Memišević
10. Dr. Vladimir Gligorijevic



## ➤ Collaborators:

Robin Ketteler, Harry Hemmingway,  
Igor Stagljjar, Charles Boone, ...

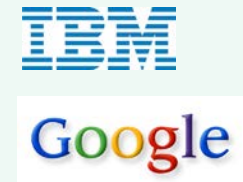




# Acknowledgements



## ➤ Funding:

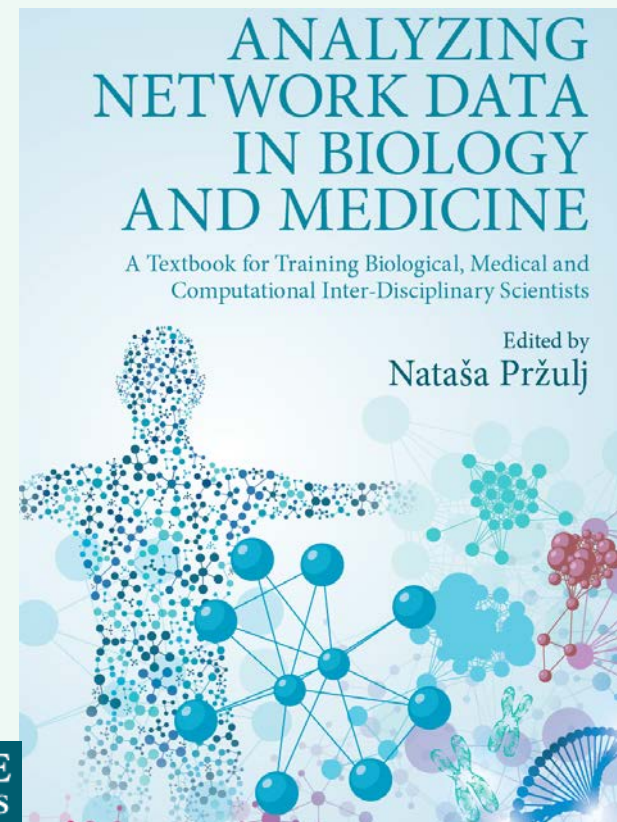


### ERC Consolidator Grant:

- Post-Doc positions
- PhD student positions

### JnJ:

- Post-Doc position



CAMBRIDGE  
UNIVERSITY PRESS

# Thank you



## Comments and Questions