

Paper Notes

Friday 16 September 2016 16:50

Biology

Data



DATABASE of databases!!! <http://www.pathguide.org>

1. Microarray -> **Gene Expression Omnibus (GEO)**
2. CrispR screens
3. Chip-seq
4. RNAi
5. SGA
6. Yeast 2 hybrid
7. Mass spectrometry

Experimental techniques	Network type	Link	Signed	Direction	Computational methods	Database	Reference
Microarray	Relevance - co-expression gene interaction network	Correlation between gene expression levels				Gene Expression Omnibus	
Yeast 2 hybrid	Protein interaction network	Physical interaction of proteins				BioGRID	
Mass spectrometry	Protein interaction network	Physical interaction of proteins				BioGRID	

Assumptions	Comment	SUPPLEMENTARY INFORMATION			In format provided by Mitra <i>et al.</i> (OCTO)
		Supplementary information S1 (table). Sources of molecular interaction and ‘omics’ profiling data			
Physical		Interaction type(s)	Detection methodologies	Databases	
		Protein-protein	Yeast-2-hybrid (Y2H) ¹⁻³ , co-immuno-precipitation (Co-IP) ⁴ , mass spectroscopy ^{5,6} , affinity purification coupled with mass spectroscopy (AP/MS) ^{7,8}	BioGRID ⁹ , IntAct ¹⁰ , APID ¹¹ , STRING ¹² , MINT ¹³ , DIP ¹⁴ , HPRD ¹⁵ , MIPS-MPI ¹⁶ , Netpath ¹⁷ , DroiD ¹⁸	
		Protein-DNA (e.g., regulatory networks)	Yeast-1-hybrid (Y1H) ¹⁹ , chromatin immuno-precipitation based methods (CHIP-CHIP) ²⁰ , DNA-footprinting ²¹	TRANSFAC ²² , UniProbe ²³ , DroiD ¹⁸ , BioGRID ⁹ , TcoF-DB ²⁴ , BIPA ²⁵ , hPDB ²⁶ , EDGEdb ²⁷ , NPIDB ^{28,29}	
		Protein-RNA	RNA electro-mobility shift (RNA-EMSA) ³⁰ , RNA-pull down ³¹	PRID ³² , BIPA ²⁵ , NPInter ³³ , RBPDB ³⁴ , StarBase ³⁶	
		Metabolic (e.g., enzyme-substrate, ligand-receptor)	Mass spectroscopy based selective reaction monitoring (SRM) ^{6,37} , NMR ³⁸ , affinity purification ⁸ , co-IP ³ , fluorescence spectroscopy ³⁹	Reactome ⁴⁰ , KEGG ⁴¹ , BioCyc and LIPIDMAPS ⁴² , HMDB ^{43,44} , EcoCyc ⁴⁵ , HumanCyc ⁴⁶ , ConsensusPathDB ⁴⁷	
		Protein/gene-compound (e.g., drug-target, chemical-protein)	Chemical structure ^{48,49} , forward or reverse chemo-genomic/proteomic profiling ⁵⁰⁻⁵² , in silico predictions ⁵³	SuperTarget ⁵⁴ , Matador ⁵⁴ , DrugBank ⁵⁵ , ChemProt ⁵⁶ , STITCH ⁵⁷ , AffinDB ⁵⁸ , MatrixDB ⁵⁹ , PSMDB ⁶⁰ , PDB-Ligand ⁶¹ , ChEMBL ⁶² , ConsensusPathDB ⁴⁷	
		Genetic (gene-gene)	Synthetic genetic array (SGA) ⁶³ , Epistatic Miniaarray Profiling (E-MAP) ⁶⁴ , co-expression profiling ^{65,66}	BioGRID ⁹ , DRYGIN ⁶⁷ , CYGD ⁶⁸ , DiseaseMiner ⁶⁹ , ConsensusPathDB ⁴⁷	
		Gene-Disease	Literature curation, clinical and sequence information	OMIM ⁶⁹ , HuDiNe ⁷⁰ , Diseaseome ⁷¹	
	Omics data type	Detection methodologies	Databases		
	Transcriptomics	Microarray ⁷² , RNASeq ⁷³	GEO ⁷⁴ , SMD ⁷⁵ , TCGA ⁷⁶ , GXD ⁷⁷ , ONCOMINE ⁷⁸ , ArrayExpress ⁷⁹		
	RNAi (phenomics)	RNAi interference assay ⁸⁰	RNAiDB ⁸¹ , GenomeRNAi ⁸² , siReco ⁸³		

OBER 2013)

IG ¹² ,
PI ¹⁶ ,
,
PI ²⁶ ,
, PRIDB ³⁵ ,
MetaCyc ⁴² ,
,
K ⁵⁵
d ⁶¹ ,
roidID ¹⁸ ,
ords ⁸³

SGA	Gene interaction network based on phenotypes	Phenotype quantification, genetic interaction			
	Signaling network				
	Metabolic network				
Combined					

Assumptions of mRNA levels network inference

1. mRNA levels correlate to the protein levels
2. mRNA levels of transcription factors and their targets tend to be correlated
3. binary interactions
 1. This can be solved by using hypergraphs or the equivalent bipartite networks! See Network introduction
4. Static representation of a nonlinear dynamic process
5. ????

Assumptions and biases of SGA phenotype networks (inference of genetic interactions)

- i. Phenotype and function refers to growth measurements
- ii. One condition, usually the optimal conditions -> Differential network biology
- iii. Biased towards negative interactions

Biological Networks at the molecular level

Functional relations	Epigenomics	Methylation profiling ⁸⁴	DAnCER ⁸⁵ , DiseaseMeth ⁸⁶ , PubMetDB ⁸⁷ , MethDB ⁸⁸ , MethCancerDB ⁸⁹ , MethylDB ⁹⁰
	Mutation / SNP	SNP Array ⁹¹ , genome sequencing ⁹²	TCGA ⁷⁶ , dbSNP ⁹³ , dbQSNP ⁹⁴ , GWAS Catalog ⁹⁵ , OMIM ⁶⁹
	Proteomics	CHIP ^{11,2,96} , Mass Spectrometry ^{5,6}	PDB ⁹⁷ , ExPASy ⁹⁸ , InterPro ⁹⁹ , World Protein Database ¹⁰⁰ , JASPAR ¹⁰¹
	Phosphorylation profile	Mass Spectrometry ⁵ , literature curation	PhosphoGRID ¹⁰² , PhosphoELM ¹⁰³ , ProPhenotype ¹⁰⁴
	For more information and databases, visit www.pathguide.org		

wman networks an

h ⁸⁷ ,
Cancer ⁹⁰
AS Central
-2DPage ¹⁰⁰ ,
HOSIDA ¹⁰⁴

1.

Biological evidence and framework

- GOAL: develop a complete map of biological modules underlying cellular architecture and function
Integrative approaches for finding modular structure in biological networks.
- Levels of biological - molecular networks -> **Genetic interaction networks/ better understanding of biological systems**
 - For a computer, the lower level of abstraction would contain details on the hardware and software. At the higher level will represent the logic of the program. In agreement with this approach, a system needs to consider a biological system with all its complexity and identify, from the genomic sequence, different levels of abstractions. At the lower level of this conceptual structure there would be several networks representing the physical structure and organization of the genome. The nodes could be genes/coding sequences, single-nucleotide polymorphisms (SNPs) or other genomic features linked by edges representing their physical proximity and organization within chromosomes, gene homology etc. (**Figure 2**, level I). The second level of abstraction would represent the organization of the genome into physical components: proteins and RNA. Edges between these elements would indicate that they are co-expressed in different contexts or that their expression profiles throughout different experimental conditions are highly correlated (**Figure 2**, level II; Ge et al., 2003; Vidal et al., 2007). The third level of abstraction would represent physical interactions between different elements: protein-protein (PPI), protein-DNA (PDI) or protein-RNA (PRI) interactions (**Figure 2**, level III; Vidal et al., 2007). The fourth level of abstraction will allow the visualization of the functional relationships between different physical elements. This level would contain GI networks, signalling and metabolic pathways (Figure 2, level IV). The fifth level would represent biological processes. This level would contain sets of proteins implicated in the same biological process would be linked by an edge (**Figure 2**, level V). The last level of abstraction would represent phenotypes and show the relationships between different phenotypes associated with similar phenotypes and diseases (**Figure 2**, level VI).
- Metabolic networks convert raw materials from the environment into value-added products and help the cell to dispose of intracellular materials. Regulatory networks alter the output of the genomic program. A short-term regulatory network only encompasses transcriptional and translational regulation; however, a long-term regulatory network encompasses all cellular processes. Signalling networks are the communication pathways that connect cellular components and coordinate the phenotypic response to perturbations. These networks simultaneously monitor a wide variety of external and internal parameters, including hormones and DNA damage. The observations are continuously processed in an integrated fashion, and the metabolic outputs are modulated to deliver the appropriate response. For simplicity's sake, we will focus on the molecular components of a cell and their interactions — intercellular interactions are influenced by environmental factors. -> **Towards genome-scale signalling-network reconstructions**
- The state of a cell is governed by the complex regulation of the expression of its genes. This regulation occurs at many levels, ranging from chromatin remodelling to post translational modifications [60]. In order to better understand gene regulation, a large and growing body of transcriptomic data can be used to analyse the interactions between genes, and has prompted the development of a number of gene regulatory network reconstruction algorithms, including Graphical Gaussian models, Bayesian networks and related methods.

unction ->

and to better predict

while the higher

systems biologist will

sequence to the

ture, we would find

. In these networks,

coding sequences

somes, their

expression of that

would indicate that

ut multiple

et al., 2011). The

ments – protein–

idal et al., 2011).

ps linking these

pathways (**Figure 2**,

n networks where

Figure 2, level V). The sixth

between elements

s and recycle or

gramme. Here, the

ever, there are

unication networks

s. Signalling

ng nutrient levels

nd the genomic and

this Review focuses

e treated as

regulation occurs at

in order better to

o infer interactions

rk (GRN)

evance networks ##

Information theory and signal transduction systems/ From molecular information processing to inference

- the aims of large-scale network inference, which are typically to highlight and explore dependent sets—and thereby help to generate new hypotheses worthy of further investigation—rather than a single grand unifying model of the system. -> **Systems biology (un)certainties**
- Models are simplified (but not simplistic) representations of real systems, and this is what makes them attractive to explore the consequences of our assumptions, and to identify an understanding of the principles governing a biological system. Models are tools to uncover mechanisms that cannot be directly observed, akin to microscopes or nuclear magnetic resonance machines (NMR). Interpreted appropriately, with due attention paid to inherent uncertainties, the mathematical and computational modeling of biological systems allows the exploration of hypotheses. But the validity of models depends on the ability to assess, communicate, and, ultimately, understand their uncertainty.

Systems biology (un)certainties

- For some model organisms, however, protein interaction data covers 20% of the proteins known in the organism (ignoring multiple isoforms due to alternative splicing, etc.). This observation poses an important question of just how representative a random subnet is for the global network.
- We emphasize that the degree distribution alone does not suffice to characterize a network. There are many different types of networks, e.g., some with many cross-connections (loops) and others with “tree-like” form. Networks with the same degree distribution have the same degree distribution. -> **Subnets of scale-free networks are not scale-free: Sampling properties of networks**
- Essential genes are those genes critical for cell viability under certain contexts -> **Network analysis of essentiality in functional genomics experiments**
- However, negative feedback plays an important role in biological robustness and evolution. It has also played an equally prominent role in the origins of life. <- **Prebiotic network evolution/ six key parameters**
- A possible resolution to the instabilities that arise in network expansion and evolution is negative feedback, which enables growth, with negative feedback, which contributes to robustness** (2c). <- **Prebiotic network evolution/ six key parameters**
- This reverse engineering is extremely difficult. Although an electrical engineer could design a circuit that would amplify signals, he would find it difficult to deduce the circuit diagram of an amplifier by correlating its outputs with its inputs. It is thus unlikely that we can deduce the detailed level description of a module solely from genome-wide information about gene expression and interactions between proteins. Solving this problem is likely to require additional types of information, such as finding general principles that govern the structure and function of modules ## **From molecular to systems biology**
- But techniques for collecting information about the entire genome will be only as powerful as the tools available to analyse it, just as our ability to infer protein structure and function from protein sequences has increased with the sophistication of tools for sequence analysis. ## **From molecular to systems biology**
- Cell functioning requires the tight linking between multiple regulatory systems that include controlling initiation of gene transcription, RNA splicing, mRNA transport, translation initiation, translational protein modifications, or the degradation of mRNA/protein -> **Michael Stumpf**

ing to network

tendencies in data
r than to uncover a

esly the property
y where we lack
mechanisms that
15). Used and
ical and
e relevance of these
uncertainties. ->

nown to exist in that
es the interesting
rk -> **Subnets of**

: very different
(no loops at all), can
Sampling properties

analysis of gene

52 and may have
key parameters
to balance positive
and stability (Fig.

many different
an unknown
circuitry or a higher-
and physical
formation and
cular to modular cell

as the tools
n sequence data has
ular cell biology
mechanisms for
on, post-
f, David J. Balding,

Mark Girolami (editors)-Handbook of Statistical Systems Biology ## Chapter 12

- main obstacles is the difficulty of choosing the appropriate experimental data, together with computational approaches. -> **Handbook of Statistical Systems Biology ## Chapter 12**
- the goal for methodology to infer regulatory networks consists in modelling and recovering interactions such as ‘protein i activates (or inhibits) the transcription of gene j’. -> **Handbook of Statistical Systems Biology ## Chapter 12**
- From the biological side, it is not that clear that the identification of regulatory networks is possible from mRNA concentrations due to the role of additional sources of regulation such as micro RNA. -> **Biological Networks 2007 ## chapter 3**
- What are the limits of the definition of a regulatory protein? We are interested here in proteins that regulate a subset of the genes of an organism. Other proteins also called regulatory play a more general role in transcription (for example the eukaryotic transcription factor of type II). In principle, these genes, like RNA polymerase itself, are essential for the transcription of all genes that encode proteins. However, their action is non-specific, they are not covered in this Section ## RNA polymerase is necessary for transcription so is abundant in the cell!!! -> **Biological Networks 2007 ## chapter 4**
- **Transcriptional Networks represent only binary interactions between mRNA molecules(!!!!)** because proteins bind to form protein complexes and then they interact with other proteins and/or catalyze reactions. *Binary interactions are not suitable when it comes to expressing interactions that involve more than two partners in a complex fashion. The formalism of hypergraphs allows to express non-binary interactions and could be used in this context.* -> **Biological Networks 2007 ## chapter 4**
- because high-throughput datasets can be filled with technical and biological noise and have inherent biases and coverage, improved statistics are needed to distinguish signal from noise, as well as methods for integration to annotate the biologically relevant relationships; because the logical interpretation of a network is not easily comprehensible to the human brain, computational modeling is needed to predict the output from the signal input or system perturbation; because no analysis and modeling methods exist, detailed targeted biological experiments are needed to validate models before a hypothesis can be tested. This approach health and medical problems -> **Understanding biological functions through molecular networks ## Problems of the high throughput data!!**
- Small-scale gene-centric studies have delineated many valuable genetic and biochemical relationships between genes and pathways. This information forms the skeleton of the entire complete network and needs standard annotation standards. For now, such information covers only a tiny fraction of the full network and is biased towards certain biological functions. For example, protein interactions in the literature curated by the Human Protein Reference Database [14] cover only 20 000 interactions out of a conceivable total of 100 000 interactions and are strongly biased for cancer-related processes (Supplementary information). Large-scale experiments and data mining are obviously more effective for mapping genetic and protein interactions. -> **Understanding biological functions through molecular networks ## Problems of the single approach!!**
- However, the change to network biology does not simply entail handing over biology to physicists and mathematicians. A good understanding of biology is needed to ask the right questions, to choose appropriate network analysis tools, and to confirm analysis results by solid experimentation. After all, network biology is a discipline of biology. The fundamental goal of network biology is the same as molecular biology: to understand the function of living systems.

h the appropriate

regulatory **k of Statistical**

possible only using
NA. -> **Biological**

proteins that regulate a
small role in initiating
generalists, like RNA
However, since their
role for all transcription

, but some proteins
in reactions!!! *Graphs*
vers in an obligate
be used in the above

different technical
as better data
ation of the whole
ed to predict the
method is perfect, more
sis can be used to
olecular networks ##

relationships between
and can serve as
ork and is biased
ted by HPRD
tive estimate of 150
rmation, table S1).
networks ->
gene - process

physicists and
choose proper
etwork biology is
erstand basic

biological processes and the mechanisms of human diseases. -> **Understanding biological molecular networks**

- Furthermore, yeast two-hybrid assays produce many false-positive outcomes, and the current pathway maps may be heavily biased towards connection to functionally important genes that have been popular targets for research -> **Kitano, Hiroaki Computational systems biology**.

- despite the known poor overall correlation between mRNAs and their protein products, the assumption that differentially expressed mRNAs impact their respective experimental conditions in proteins. However, to the best of our knowledge, this assumption has never been explicitly tested.

Relationship between differentially expressed mRNA and mRNA- protein correlations in a yeast system

- Although clustering analysis provides insight into the “correlation” among genes and biological pathways does not reveal the “causality” of regulatory relationships. Several methods have been proposed to automatically discover regulatory relationships solely on the basis of microarray data (7–9). These methods use information derived from mRNA abundance, so there is limited scope to infer transcriptional regulation. Posttranscriptional and posttranslational mechanisms of regulation will be incorporated as large-scale data become available, but many properties have yet to be measured with accuracy or in high throughput. -> **Systems Biology: A Brief Overview**

- There is a need for systematic approaches to infer causal relationships between interacting proteins. In this context we refer to the “direction” (edge direction), “sign” (activation/inhibition) and “mode” (e.g. phosphorylation, ubiquitination) of signal flow in PPI networks. -> **Integrating protein-protein interaction networks with phenotypes reveals signs of interactions**

- The signed network opens up new scope for network analysis, such as the application of structural theory. Further integration of other information-flow properties such as edge-direction would enable sophisticated flow-based network analysis. -> **Integrating protein-protein interaction networks with phenotypes reveals signs of interactions**

- **Integrating protein-protein interaction networks with phenotypes reveals signs of interactions**

- Each RNAi screen identifies positive and negative regulators of a particular phenotype
 - How?
- PPI is for the physical interactions, from gene expression one can reconstruct a signed directed graph. Direct regulation relationships must include both physical activation - repression data.

- Biological-network reconstructions provide us with a framework that allows us to identify the whole and, subsequently, understand the relationship between the molecular phenotype and the organism phenotype. -> **Towards genome scale signalling network reconstructions**

- a minority of genes are essential, and these define hubs of activity that can in some cases exert influence over a given functional module to influence and even coordinate multiple cellular processes. It is not surprising, given the interactional complexity, that single genes rarely specify a phenotype in its entirety. The outcome of the genetic network are now apparent, but a compelling case can be made for a deeper and more detailed exploration of this model system as an exemplar for more complex eukaryotes. -> **Exploring interactions and networks with yeast**

- The physical-interaction map, generated by large-scale two-hybrid^{51,52} or affinity purification-mass spectrometry identification^{26,42,53,54}, provides a view of the gene products that are involved in the yeast interactome.

unctions through

nt hand-crafted
imply because these
pdf

re is an implicit
tions via differences
ly tested. ->

xenograft model

ical phenomena, it
posed to

At present, such
causality based on
on must be
sured with sufficient

proteins, by which
phosphorylation,
tworks with

ucture balance
uld enable
woks with

ctions
notype?

network based on
ical interactions and

the context of the
nd the general

xtend beyond a
o wonder, given this
tlines of a yeast
more complete
g genetic

ion followed by
mble into soluble

mass spectrometry identification^{120,45,55,54}, provides a view of the gene products that associate into protein complexes and function together as biochemical machines. Rather than physical interactions, a genetic-interaction map provides functional information, largely identifying gene products that act together in functionally related pathways. Although genetic interactions overlap with protein–protein interactions more often than expected by chance, such overlap is relatively rare, occurring at a frequency of less than 1%.

2). -> Exploring genetic interactions and networks with yeast

- Large-scale genetic analyses reveal that mutations in most eukaryotic genes have little discernible effect on growth. For example, systematic gene deletion in *S. cerevisiae*, discussed in detail below, produced a remarkable finding: only ~20% of yeast genes are essential for viability when deleted individually in haploids grown under standard laboratory conditions
- Comparisons between aggravating and alleviating effects revealed that, for most functional interactions, mutations were either largely aggravating or largely alleviating, but not mixed, an asymmetrical pattern termed ‘monochromatic’.
- the set of interactions that are observed for a particular query gene can be suggestive of its position of a gene in a genetic-interaction network being highly predictive of its molecular requirements
- **genetic interactions and networks with yeast**
- **Because most genetic interactions do not overlap with physical interactions, the two types of interaction are said to be largely orthogonal**
- Integration attempts in yeast, combining physical protein–protein interaction maps with coexpression data, have revealed that interacting proteins are more likely to be encoded by genes with similar expression patterns than noninteracting proteins (Ge et al., 2001; Grigoriev, 2001; Jansen et al., 2002; Kemmeren et al., 2002). These observations were subsequently confirmed in many other organisms (Ge et al., 2003). Beyond the technical aspect of finding significant overlaps between interaction edges in interactome networks and coexpression edges in transcription profiling networks, these observations have been used to assess the overall biological significance of interactome datasets -> Interactome Networks and Human Disease
- Experimental artefacts, variability in coverage across data sets, sampling bias towards well-known proteins, limitations in screening power and inherent sensitivities in various assays can yield false positives and false negatives in interaction data -> **Integrative approaches for finding modular structure in biological interaction networks**
- In addition, as reported by Zotenko et al. [17], essential proteins tend to form highly connected clusters rather than function independently -> **A new essential protein discovery method based on the integration of protein interaction and gene expression data**
- However, standard clustering is not ideal for PPI networks: proteins may have multiple functions, therefore the corresponding nodes may belong to more than one cluster; for example, 207 of the 4300 proteins in the CYC2008 hand-curated yeast complex data set3 participate in more than one complex. -> **Identifying overlapping protein complexes in protein-protein interaction networks**
- The existence of a GI between two genes does not necessarily imply that these two genes code for the same proteins or that the two genes are even expressed in the same cell. In fact, a GI only implies that the two genes share a functional relationship. These two genes may be involved in the same biological process, but they may also be involved in compensatory pathways with unrelated apparent function -> **Genetic interaction networks: better understand to better predict**
- A network in the left represents a global PIN. A eukaryotic cell in the center can be divided into different compartments: Endoplasmic Reticulum, Cytoskeleton, Golgi, Cytosol, Lysosome (or Vacuole), Mitochondria, Nucleus, Peroxisome, and Cell wall.

erns into soluble
ormation, the
that operate in
nteractions more
ss than 1% (REF.

ernable effect. For
markable result:
owing in standard

groups, interac-
feature that they

s function, with the
ole. -> **Exploring**

s of interaction are

expression profiles,
ession profiles than
(al., 2002). These
nd the funda-
works and
l to estimate the
Disease
-studied processes,
osi- tives and false
biological networks.
ected clusters rather
Integration of protein-

nctions, and
of 1,628 proteins in
-> detecting

ode for interacting
that the two genes
cess or pathway; or
Genetic interaction

nto 11
ndrion Endosome

Plasma, Nucleus, Peroxisome and Extracellular, where Lysosome only exists in animal cells. compartment, a PSLIN of this compartment is constituted by the proteins localized in this compartment, a PSLIN of this compartment is constituted by the proteins localized in this compartment, their interactions. With the subcellular localization information of proteins, the PSLINs can be mapped the global PIN to each compartment separately. -> **Rechecking the Centrality-Lethal Scope of Protein Subcellular Localization Interaction Networks**

- In contrast to pathways, protein complexes are the functional units of proteome organization. dynamic assembly is fundamental to induce cellular responses to different internal and external

Protein Complex-Based Analysis Framework for High-Throughput Data Sets

- The vast majority of proteins work as parts of assemblies composed of several elements, thus protein complexes as essential cellular functional units. Given the fundamental importance of interactions, proteome-wide “interactome” maps based on pairwise protein interactions using hybrid (Y2H) system have been determined for several organisms (Giot et al., 2003; Ito et al., 2004; Rual et al., 2005; Stanyon et al., 2004; Stelzl et al., 2005; Uetz et al., 2000). Alternative isolation based on co-affinity purification combined with tandem mass spectrometry (coAP) to generate protein complex maps at proteome-scale for *Saccharomyces cerevisiae* ->**A Protein Network of Drosophila melanogaster**

- However, these abovementioned approaches for extracting dense subgraphs fail to take into account inherent organization. Recent analysis of experimentally detected protein complexes [23] has shown that a complex consists of a core component and attachments. Core proteins are highly co-expressed, functionally similar, and each attachment protein binds to a subset of core proteins to form a complex. Based -> **Identification of protein complexes from multi-relationship protein interaction networks**
- we construct a multi-relationship protein interaction network (MPIN) by integrating PPI network, gene ontology (GO) annotation information.->**Identification of protein complexes from multi-relationship protein interaction networks**

Network construction methods and evaluation

- Information theoretic approaches to GRN reconstruction have two major strengths. The first is that information is able to capture nonlinear associations between variables, a feature seen in ex-

tion, Enzyme,

For each compartment and be generated by **Chaliy Rule in the**

on, and their internal cues. ->

ereby defining of protein ing the yeast two-, 2001; Li et al., eley, protein complex -MS) has been used **Protein Complex**

o account the as revealed that a sed and share high n a biological **Interaction networks** work topology with **Multi-relationship**

it is that mutual expression data.

Second, it has been demonstrated that the use of the data processing inequality (DPI), which noisy system $X \rightarrow Y \rightarrow Z$, knowledge of Z cannot give more information about X than Y can give in distinguishing two genes regulated by a third from a trio of co-regulating genes. In reconstructing GRNs, the combined approach using mutual information and DPI outperforms Bayesian and frequentist techniques in the precision and recall of direct regulatory links.

Information theory and systems/ From molecular information processing to network inference

- Several alternative network inference approaches exist (4, 5) that provide better, more robust and accurate network reconstructions. These methods may be used to infer regulatory relationships, and how interactions between genes change over time or differ between diseased and healthy controls. -> **Systems biology (un)certainties**
- Genome-scale inference of transcriptional gene regulation has become possible with the advent of high-throughput technologies such as microarrays and RNA sequencing, as they provide snapshots of the transcriptome under many tested experimental conditions. From these data, the challenge is to computationally predict direct regulatory interactions between a transcription factor and its target genes. The aggregate of all predicted interactions comprises the gene regulatory network. -> **Wisdom of crowds for robust gene network inference**
- Understanding the advantages and limitations of different network inference methods is critical for effective application in a given biological context -> **Wisdom of crowds for robust gene network inference**
- We conclude that there is no category of network inference methods that is inherently superior. Performance depends largely on the specific implementation of each individual method. -> **Wisdom of crowds for robust gene network inference**
- We next analyzed how method-specific biases influenced the recovery of different connectivity patterns (network motifs), and we observed characteristic trends for different method categories (Figure 6). Feed-forward loops were recovered most reliably by mutual-information and correlation-based methods, whereas sparse-regression and Bayesian-network methods performed worse at this task. This suggests that the latter approaches preferentially select regulators that independently contribute to the expression of target genes. However, the assumption of independence is violated for genes regulated by multiple transcription factors, as in the case of feed-forward loops. Indeed, linear cascades were more often predicted by regression and Bayesian-network methods. This shows that current methods exhibit a trade-off between performance on cascades and performance on feed-forward loops. -> **Wisdom of crowds for robust gene network inference**
- Network inference methods have complementary advantages and limitations under different circumstances. This suggests that combining the results of multiple inference methods could be a good strategy for improving predictions. -> **Wisdom of crowds for robust gene network inference**
- After identifying these modules²⁴, we tested them for enrichment of Gene Ontology terms (see Note 7). Network modules are strongly enriched for very specific biological processes. This analysis reveals unique functions to most of the identified modules in both networks.
- When constructing a compendium of microarrays for global network inference, one should not be tempted toward oversampling a narrow set of experimental conditions. -> **Wisdom of crowds for robust gene network inference** **#### QUESTION!!!!**
- Bayesian- network methods exhibited below-average performance in this challenge, likely because of their heuristic searches, which are often too costly for systematic data resampling and may be biased.

n states that for a
ive about X [19,65],
structing simulated
l relevance networks
signal transduction

ust candidate
o investigate gene-
tween disease cases

event of high-
ts of the
is to
s target genes; the
of crowds for

tical for their
work inference
erior and that

Wisdom of crowds

vity patterns
g. 2c). For example,
sed methods,
ne reason for this is
the expression of
mutually dependent
re accurately
experience a trade-
of crowds for

nt contexts, which
for improving

(Supplementary
allowed us to assign

thus avoid any bias
robust gene network

because they use
atter suited for

~~neuristic strategies, which are often too costly for systematic data resampling and may be better suited for smaller networks. Information theoretic methods performed better than correlation-based methods, but the two approaches had similar biases in predicting regulatory relationships. They also performed similarly to regression and Bayesian- network methods on feed-forward loops, fan-ins and fan-outs (the most common topologically connected parts of the network), but they had an increased rate of false positives for cascades and feedback loops.~~

crowds for robust gene network inference

- ~~A fundamental assumption of network inference algorithms is that mRNA levels of transcripts and their targets tend to be correlated; we found that this is true for *E. coli*, but not for *S. cerevisiae*. This suggests that the lower coverage of *S. cerevisiae* gold standards may also play a role (*E. coli* has the best-known regulatory network of any free-living organism¹⁶), the poor correlation at the mRNA level in *S. cerevisiae* may reflect the increased regulatory complexity and prevalence of post-transcriptional regulation in eukaryotes. This would suggest that accurate inference of eukaryotic regulatory networks requires additional data on transcription-factor binding sites and chromatin modifications.~~
- **crowds for robust gene network inference**
- We will show in this chapter that there are many ways to express the reverse-modeling of regulatory networks as a machine learning problem. However, regardless of the adopted approach, the hardest part is the identification of the network structure. This problem when expressed as a combinatorial one is known to be NP-hard ->
- Transcription of DNA into RNA is the first — and often most regulated — step in gene expression. In eukaryotes, genes are regulated at the transcriptional level, and 5-10% of protein-encoding genes encode proteins that regulate other genes. Between regulated genes and regulatory proteins, the interactions between regulated genes and regulatory proteins constitute a core component of the transcriptional network or genetic network. -> Biological Networks 2007 Chapter 4
- It is therefore likely that comprehensive mapping of the quantitative genetic interaction network will require the integration of a number datasets from different screening approaches, similar to the recent work on the physical protein-protein interaction (PPI) networks in yeast and human -> **Quantitative interactions in yeast - Comparative evaluation and integrative analysis**
- Compared to PPI networks, an additional challenge originates from the quantitative nature of genetic interaction datasets; instead of comparing the overlap in binary terms, such as presence or absence of a physical interaction, here we should take into account the full spectrum of genetic interactions, from the extreme cases of negative interactions (i.e., synthetic sick and lethality) to the positive class of interacting gene pairs (e.g., masking and suppression subcategories) [2,3,17].

2. SGA relationships -> Systematic Mapping of Genetic Interaction Networks

- a. An aggravating (or negative) interaction occurs when a double mutant exhibits a phenotype more severe than expected from the phenotypes of the individual mutants. This type of interaction indicates that the gene products function in redundant parallel pathways, and highlights the robustness of the molecular network in tolerating genetic variations [4].
 - b. An alleviating (or positive) interaction occurs when a double mutant exhibits a phenotype less severe than expected from the phenotypes of individual mutants. This type of interaction indicates that the gene products operate in concert or in series within the same pathway [4].
3. Genome-scale, quantitative analysis of genetic interactions. We consider a digenic interaction between two genes, where each gene is mutated in a single copy. A double mutant that shows a significant deviation in fitness compared with the expected multiplicative effect of the two single mutants is considered to have a significant interaction.

other suited for

methods, but the
ed better than
e more densely
des. -> **Wisdom of**

tion factors and
siae. Although the
wn regulatory
ae is likely due to
karyotes, which
l inputs, such as
ion7. -> **Wisdom of**

regulatory networks
point to solve is the
e is known to be

ssion. As most
de regulatory
mplex web called

tworks will require
efforts to complete
maps of genetic

of the genetic
absence of a
ons, ranging from
es of interacting

hotype that is more
eraction indicates
bustness of the

type that is less
eraction indicates that

on as a double
ive effect of

combining two single mutants (6). Negative interactions refer to a more severe fitness defect with the extreme case being synthetic lethality; positive interactions refer to double mutant fitness defect than expected. To quantitatively score genetic interactions in large-scale SGA developed a model to estimate fitness defects directly from double-mutant colony sizes ->

Landscape of a Cell

4. There have been some attempts to investigate the temporal properties for individual protein interactions by integrating PPI data with time-course gene expression data [21-29]. <- Detected protein complexes from dynamic protein-protein interaction networks

Centralities and node influence

- centrality is a quantitative measure that aims at revealing the importance of a node. -> **Axioms of centrality**
- degree of a vertex alone, as a specific centrality measure, is not sufficient to distinguish lethal from viable ones (Wuchty (2002)), that in protein networks there is no relation between network size and robustness against amino-acid substitutions (Hahn et al. (2004)), and that for biological systems several centrality measures have to be considered (Wuchty and Stadler (2003); Koschützki et al. (2004)) -> **Centrality analysis methods for biological networks and their application to gene regulatory networks**
- Topological features of the protein networks have been demonstrated to reflect the functional roles of interacting genes. For example, essential genes in yeast tend to be well connected and globally centered in the yeast protein network (Jeong et al., 2001; Wuchty and Almaas, 2005). Furthermore, globally centered genes are more likely to be well conserved and serve as an evolutionary backbone for the network (Wuchty and Almaas, 2005). -> **Interactome-transcriptome analysis reveals the high centrality of genes differentially expressed in lung cancer tissues**
- Integration of protein network data may extend the reach of the established method of analyzing the genes in a broader context. Based on this notion, we seek to reveal the biological significance of differentially expressed genes in squamous cell lung cancer that is identified through our recent gene expression profiling study by using interactome-transcriptome analysis. We find high centrality for differentially induced genes, but not for the genes that are suppressed in cancer. -> **Interactome-transcriptome analysis reveals the high centrality of genes differentially expressed in lung cancer tissues**
- k-core analysis has been performed on the yeast essential genes and were shown to be globally centered. Non-essential genes were not (Wuchty and Almaas, 2005). The study also indicates that these genes are well conserved throughout different species. -> **Interactome-transcriptome analysis reveals the high centrality of genes differentially expressed in lung cancer tissues**
- A problem of choosing initial nodes as source spreaders to achieve maximum scale of spreading

ct than expected,
ts with a less severe
creens, we
> **The Genetic**

ins and protein
cting temporal

oms of centrality
al proteins clearly
twork connectivity
l network analysis
and Schreiber
e regulatory

onality of the
ally centered in the
ered interactions are
'uchty and Almaas,
tially expressed in

alysis by considering
cance of
cent microarray
centrality in these
tome-transcriptome

bal hubs, whereas
se global hubs are
high centrality of

ding is de ned as *in*

uence maximization problem. Our research focuses on the strategy of choosing a set of critical source spreaders in this report. -> **Identifying a set of influential spreaders in complex networks**

- Classical centrality measures in complex networks—like the degree or number of neighbors with 4, 5, betweenness centrality⁸ counting the number of shortest paths through a certain node—centrality⁹ based on the idea that relations with more influential neighbors confer greater influence¹⁰ shell decomposition¹⁰ that correlates with the outcome of supercritical spreading originating from nodes^{11–13}—rely only on topology, even if an underlying process can be indirectly associated with it. In contrast, the impact of individual elements in the global performance of the system inevitably depends on the specificities of the dynamics. -> **A measure of individual role in collective dynamics**

- Topological analysis of the yeast protein interaction network shows that the means of degree coefficient (0.16), and betweenness centrality (0.001) of essential genes in yeast are about three times higher than those in nonessential genes -> **Predicting essential genes based on network and sequence information**
- many bottlenecks also tend to be hubs). Therefore, we further investigate which one of these metrics is a better predictor of protein essentiality in both regulatory and interaction networks. -> **The Importance of Bottlenecks in Protein Networks/ Correlation with Gene Essentiality and Expression Dynamics**
- we observed that bottlenecks (both nonhub–bottlenecks and hub–bottlenecks) have a strong correlation with products of essential genes, whereas hub–nonbottlenecks are surprisingly not essential. This finding suggests that it is the betweenness that is a stronger determinant of the essentiality of a protein in the yeast protein interaction network, not the degree -> **The Importance of Bottlenecks in Protein Networks/ Correlation with Gene Essentiality and Expression Dynamics**

- An intriguing question in the analysis of biological networks is whether biological characteristics, such as essentiality, can be explained by its placement in the network, i.e., whether topological position implies biological importance. One of the first connections between the two in the context of the yeast protein interaction network, the so-called centrality-lethality rule, was observed by Jeong and colleagues.

Hubs in the Yeast Protein Interaction Network Tend To Be Essential/ Reexamining the Connection between the Network Topology and Essentiality

- Since then the correlation between degree and essentiality was confirmed by other studies. However, recently there was no systematic attempt to examine the reasons for this correlation ->**Why Hubs in the Yeast Protein Interaction Network Tend To Be Essential/ Reexamining the Connection between the Network Topology and Essentiality**

- We found that a significant portion of 2,144 mouse genes with yeast orthologs changed their essentiality between mouse and yeast (Fig. 1a). We arranged the orthologous pairs of yeast and mouse genes into phenotypic groups based on their changing essentiality patterns. We found 91 genes are essential in yeast and mouse (E2E), 246 genes are nonessential in yeast but essential in mouse (N2E), 659 genes are nonessential in yeast but nonessential in mouse (E2N), and 1,149 genes are nonessential in both yeast and mouse (NN). Network rewiring is an important mechanism of gene essentiality change

- currently three main types of experimental strategies for the genome-wide discovery of essential genes are knockout (Giaever et al., 2002; Chen et al., 2015), gene knockdown (Harborth et al., 2001; Jantsch et al., 2003; Roemer et al., 2003) and transposon mutagenesis (Gallagher et al., 2007; Langridge et al., 2008). These methods can generate accurate collections of essential genes, but they are expensive, time-consuming, and laborious. Furthermore, these experimental methods are not suitable for some complex organisms, such as humans. -> **Predicting Essential Genes and Proteins Based on Machine Learning and Network Theory**

**critical nodes as
works**

a node interacts
node, eigenvector
importance, or the k-
ing in specific
ed in some cases. In
only depends on the

ree (11.55), clustering
twice as large as

analysis #

se two quantities is
**the Importance of
omics**

ng tendency to be
us, we determined
he regulatory
on with Gene

istics of a protein,
logical prominence
of a protein
agues ->**Why Do
nnection between**

[4–7], but until
**Why Do Hubs in the
between the Network**

ir essentialities
genes into four
sential in both yeast
nes are essential in
mouse (N2N) ->

ential genes: gene
iet al., 2001;

009). These

consuming and

organisms, especially

Network Topological

- Another observation in this search of the literature is that of the 28 articles reporting the utilization of network topological features to predict essential genes and proteins, in 12 (43%) the computational models were based on machine learning, a method in which computers make and improve predictions by learning from data through learning algorithms (see next sections for details). On the other hand, of the 16 articles reporting the utilization of network topological features as learning attributes (features that describe a complex system), 14 (86%) utilized machine learning for the prediction of essential genes and proteins (see details in the next section). Therefore, this brief analysis of these 34 papers suggests a strong relationship between the utilization of machine learning and network topological features regarding the prediction of essential genes.

Predicting Essential Genes and Proteins Based on Machine Learning and Network Topological Features/A Comprehensive Review

- despite the great progress that has been made in the machine learning-based prediction of essential genes using network topological features since the publication of the pioneering study in the field (Tong et al., 2001; **Table 1**), it is worth to mention that all these models are predictive of “constitutive” essential genes that are essential regardless the growth condition. However, many genes considered as essential under a certain growth condition might not play as critical role in another condition (Tong et al., 2001; Cagney et al., 2005; Tong et al., 2011). These are the so-called conditionally essential genes, i.e., nonessential genes that become essential under specific conditions. The development of machine learning approaches for the prediction of these conditionally essential genes using network topological features will be the main challenge in the future considering that the network topological characteristics of these genes remain unexplored.
- Essentiality VS conditional essentiality -> Predicting Essential Genes and Proteins Based on Machine Learning and Network Topological Features/A Comprehensive Review**

- Synthetic lethality** : The situation in which two genes that are **non-essential** when individual mutations do not lead to a lethal phenotype, result in **synthetic lethality** when they are combined as a double mutant. -> **Exploring genetic interactions and networks with yeast**
- we estimate that a global network will contain ~200,000 synthetic-lethal interactions. To put this in context, there are ~1,000 essential genes in yeast, for which a single mutation leads to a lethal phenotype. In contrast, there are 200-fold more ways to generate a similar phenotype through a digenic synthetic-lethal interaction. This finding indicates that digenic interactions might underlie many inherited phenotypes, and provides a possible explanation why the analytical power of single-gene effects on many phenotypes has been so limited. -> **Exploring genetic interactions and networks with yeast**
- Moreover, our approach of focusing on nonhub–bottle-neck nodes is useful for finding proteins that are involved in different processes and are involved in cross-talk. -> **The Importance of Bottlenecks in Protein Interaction Networks**
- This indicates the existence of a large number of nodes with high betweenness but low connectivity (bottleneck nodes). Importantly, such nodes are absent in computer-generated, random scale-free networks. -> **Correlation with Gene Essentiality and Expression Dynamics**
- This indicates the existence of a large number of nodes with high betweenness but low connectivity (bottleneck nodes). Importantly, such nodes are absent in computer-generated, random scale-free networks. -> **High-Betweenness Proteins in the Yeast Protein Interaction Network**

ilization of network methods used were ons based on some 4 articles reporting (%) report the certain instance; see g link between genes and proteins. -> **Biological Features/ A**

essential genes (Chen and Xu, 2005; al genes, that is, as essential under 04; Nichols et al., come essential es for the prediction hallenge in the near olored. ## global

Machine Learning

lly mutated cause
I networks with

t this number in al phenotype, but ethal interaction.
and begins to explain -> **Exploring genetic**

that mediate
Protein Networks/

nnectivity (HBL
works ->First
ork 2005

Evaluation of results

- gene enrichment ## **Gene Ontology: tool for the unification of biology**
- New experiments to test new interactions that appeared from reverse engineering of networks
 - Even a comparatively small subnetwork of this size is still a challenge to visualize insights. authors assess its quality in two ways. *First, they determine whether the genes of the network are enriched for specific cellular process categories in the Gene Ontology database, which is a wonderful strength of the paper*—the authors experimentally validate some neighbors of MYC -> **Reverse engineering gene regulatory network**
 - Basso *et al.* demonstrate that as long as the available data explore a wide range in the behavior of the system, biologically meaningful interactions can be recovered by computational methods. **Reverse engineering gene regulatory network** for the paper -> Reverse engineering of gene regulatory networks in human B cells
 - We note that these data support a direct regulatory effect of the tested transcription factor gene, but chromatin immunoprecipitation experiments would be required to determine binding. -> **Wisdom of crowds for robust gene network inference**
- Gene essentiality data of yeast were manually compiled from the Comprehensive Yeast Genome Database (<http://mips.helmholtz-muenchen.de/genre/proj/yeast/>) and large-scale experiments²⁶. There are 1,178 essential and 4,904 nonessential yeast genes. -> Network rewiring is an important mechanism of essentiality change
- [currently, the available essential genes and protein databases are DEG (Zhang and Lin, 2009), OGEE (Chen et al., 2012), and EGGS (<http://www.nmpdr.org/FIG/eggs.cgi>)]. These databases allow researchers to explore the features of essential genes and proteins and, through this exploration, features are associated with essentiality and, finally, develop computational methods for predicting essential genes and proteins. -> **Predicting Essential Genes and Proteins Based on Machine Learning Methods Using Network Topological Features/ A Comprehensive Review**
- There exist many different databases of a certain type of interaction (e.g., DIP, BioGRID and others). usually, these databases are regularly updated. Different databases or newer versions of a given database have different sets of interactions that, in turn, will give rise to new networks with distinct structures and, consequently, different values of network topological features. As an example, we can cite the studies by Hwang et al. (2009) and Acencio and Lemke (2009). In both studies, PINs of *Saccharomyces cerevisiae* were created; however, the interactions of the PIN constructed in the study by Hwang et al. (2009) were taken from the version ScereCR20070107 of DIP and the interactions of the PIN constructed in the study by Acencio and Lemke (2009) were gathered from the version 2.0.42 of the BioGRID database. Therefore, the performances of the models created by these authors cannot be reliably compared. ## The data used must be the same kind and the same data! That means not only the same database because databases are regularly updated, but the exact same datasets.
- Thus, it seems that only network topological features are not enough to distinguish essential genes and proteins. This raises the following question: **is the positive correlation between network topological features only an artifact of a possible bias (essential genes and proteins are more studied and therefore tend to have higher values of network topological features)**

works!! e.g
htfully, so the
subnetwork are
in they are. Second—
me of the first

e ‘expression space’
al algorithms ->
f regulatory

factor on the target
mine physical

ome Database
ne dataset contained
echanism of gene

9), CEG (Ye et al.,
ata have enabled
ation, reveal which
osed to identify
Machine Learning and

IntAct for PPIs) and,
given database will
structures and,
the studies by
cerevisiae were
9) were collected
e study by Acencio
re, the prediction
reference networks
e they are updated

al from non-essential
essentiality and
ins are the focus of

more studies and therefore tend to have higher values of network topological present in derived from small scale experiments?

- Regardless the resolution of this debate, a large-scale study for evaluating how well essential proteins can be predicted solely by network topological features is necessary to confirm this prediction performance.
- ROC curves in R <http://blog.revolutionanalytics.com/2016/11/calculating-auc.html>
- Dygraphs for r Interactive plots for HTML files

the networks mainly

al genes and
s moderate