

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	

en	Assump tions	Comme nt

SUPPLEMENTARY INFORMATION

In format provided by Mitra *et al.* (O

Supplementary information S1 (table). Sources of molecular interaction and ‘omics’ profiling data

	Interaction type(s)	Detection methodologies	Databases
Physical	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 , S ¹² , MINT ¹³ , DIP ¹⁴ , HPRD ¹⁵ , MIPS 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 , D ²⁴ , BioGRID ⁹ , TcoF-DB ²⁴ , BIPA ²⁵ , EDGEDb ²⁷ , NPIDB ^{28,29}
	Protein-RNA	RNA electro-mobility shift (RNA-EMSA) ³⁰ , RNA-pull down ³¹	PRID ³² , BIPA ²⁵ NPInter ³³ , RBP 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 ⁴² , HMDB ^{43,44} , EcoCyc ⁴⁵ , HumanC 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 , Drug ⁵⁵ , ChemProt ⁵⁶ , STITCH ⁵⁷ , AffinD ⁵⁸ , MatrixDB ⁵⁹ , PSMDB ⁶⁰ , PDB-L ⁶¹ , ChEMBL ⁶² , ConsensusPathDB ⁴⁷
Functional	Genetic (gene-gene)	Synthetic genetic array (SGA) ⁶³ , Epistatic Miniarray Profiling (E-MAP) ⁶⁴ , co-expression profiling ^{65,66}	BioGRID ⁹ , DRYGIN ⁶⁷ , CYGD ⁶⁸ , ConsensusPathDB ⁴⁷
	Gene-Disease	Literature curation, clinical and sequence information	OMIM ⁶⁹ , HuDiNe ⁷⁰ , Diseases 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 , si ⁸³

OCTOBER 2013)

TRING ¹² , -MPPI ¹⁶ ,
biD ¹⁸ , hPDI ²⁶ ,
DB ³⁴ , PRIDB ³⁵ ,
and MetaCyc ⁴² , Cyc ⁴⁶ ,
Bank ⁵⁵ B ⁵⁸ , igand ⁶¹ , ⁴⁷
⁶⁸ , DroID ¹⁸ ,
e ⁷¹
⁷⁷ , , Records 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 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

- GOAL: develop a complete map of biological modules underlying cellular architecture and
Integrative approaches for finding modular structure in biological networks.

		Functional relations

Epigenomics	Methylation profiling ⁸⁴	DAnCER ⁸⁵ , DiseaseMeth ⁸⁶ , PubMedMethDB ⁸⁸ , MethCancerDB ⁸⁹ , N
Mutation / SNP	SNP Array ⁹¹ , genome sequencing ⁹²	TCGA ⁷⁶ , dbSNP ⁹³ , dbQSNP ⁹⁴ , OMIM ⁶⁹
Proteomics	CHIP ^{11,2,96} , Mass Spectrometry ^{5,6}	PDB ⁹⁷ , ExPASy ⁹⁸ , InterPro ⁹⁹ , WJASPAR ¹⁰¹
Phosphorylation profile	Mass Spectrometry ⁵ , literature curation	PhosphoGRID ¹⁰² , PhosphoELM
For more information and databases, visit www.pathguide.org		

Newman networks an

)

d function ->

bmeth ⁸⁷ , MethyCancer ⁹⁰
GWAS Central
World-2DPage ¹⁰⁰ ,
¹⁰³ , PHOSIDA ¹⁰⁴

- Levels of biological - molecular networks -> **Genetic interaction networks/ better understanding**
 - For a computer, the lower level of abstraction would contain details on the hard-ware. The first level will represent the logic of the program. In agreement with this approach, a system can be considered a biological system with all its complexity and identify, from the genomic data, several networks representing the physical structure and organization of the genome. Nodes could be genes/coding sequences, single-nucleotide polymorphisms (SNPs) 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 partitioning of the genome into physical components: proteins and RNA. Edges between these elements would represent if they are co-expressed in different contexts or that their expression profiles through different experimental conditions are highly correlated (**Figure 2**, level II; Ge et al., 2003; Vidal et al., 2002). 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). The fourth level of abstraction will allow the visualization of the functional relationships between physical elements. This level would contain GI networks, signaling and metabolic pathways (level IV). The fifth level would represent biological processes. This level would contain genes/proteins implicated in the same biological process would be linked by an edge (**Figure 2**, level V) and last level of abstraction would represent phenotypes and show the relationships between genes/proteins associated with similar phenotypes and diseases (**Figure 2**, level VI).
- Metabolic networks convert raw materials from the environment into value-added products and dispose of intracellular materials. Regulatory networks alter the output of the genomic program. A long-term regulatory network only encompasses transcriptional and translational regulation; however, a short-term regulatory network encompasses all regulatory features associated with all cellular processes. Signalling networks are the components that connect cellular components and coordinate the phenotypic response to perturbations. Sensing networks simultaneously monitor a wide variety of external and internal parameters, including nutrient availability and DNA damage. The observations are continuously processed in an integrated fashion, and metabolic outputs are modulated to deliver the appropriate response. For simplicity's sake, we focus on the molecular components of a cell and their interactions — intercellular interactions with 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]. To better understand gene regulation, a large and growing body of transcriptomic data can be used to infer relationships between genes, and has prompted the development of a number of gene regulatory network reconstruction algorithms, including Graphical Gaussian models, Bayesian networks and Markov Blanket models. **Information theory and signal transduction systems/From molecular information processing to network inference**
 - the aims of large-scale network inference, which are typically to highlight and explore key features of the data sets—and thereby help to generate new hypotheses worthy of further investigation—rather than to develop a single grand unifying model of the system. -> **Systems biology (uncertainties)**
 - Models are simplified (but not simplistic) representations of real systems, and this is 8pre

stand to better predict

ware while the higher
systems biologist will
c sequence to the
ructure, we would find
me. In these networks,
or coding sequences
mosomes, their
ne expression of that
nts would indicate that
hout multiple
la et al., 2011). The
elements – protein–
l; Vidal et al., 2011).
ships linking these
pathways (**Figure 2**,
ain networks where
ure 2, level V). The sixth
os between elements

ucts and recycle or
programme. Here, the
however, there are
ommunication networks
ons. Signalling
cluding nutrient levels
and the genomic and
ke, this Review focuses
are treated as

his regulation occurs at
. In order better to
d to infer interactions
work (GRN)
relevance networks ##
essing to network

pendencies in data
cher than to uncover a

ecisely the property

that makes them attractive to explore the consequences of our assumptions, and to identify the principles governing a biological system. Models are tools to uncover what cannot be directly observed, akin to microscopes or nuclear magnetic resonance machines. Interpreted appropriately, with due attention paid to inherent uncertainties, the mathematical modeling of biological systems allows the exploration of hypotheses. But the success of models depends on the ability to assess, communicate, and, ultimately, understand the

Systems biology (uncertainties)

- For some model organisms, however, protein interaction data covers 20% of the proteins 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.
scale-free networks are not scale-free: Sampling properties of networks
- We emphasize that the degree distribution alone does not suffice to characterize a network. Networks, e.g., some with many cross-connections (loops) and others with “tree-like” form, can 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 essentiality in functional genomics experiments**
- However, negative feedback plays an important role in biological robustness and evolution. It 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: negative feedback, which enables growth, with negative feedback, which contributes to robustness. <- Prebiotic network evolution/ six key parameters**
- This reverse engineering is extremely difficult. Although an electrical engineer could design circuits 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 a high-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 data. Finding general principles that govern the structure and function of modules ## **From molecular biology 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 biology 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 Stumm and Mark Girolami (editors)-Handbook of Statistical Systems Biology ## Chapter 12**
- main obstacles is the difficulty of choosing the appropriate experimental data, together with the 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

tity where we lack
er mechanisms that
es (15). Used and
natical and
the relevance of these
ir uncertainties. ->

s known to exist in that
oses the interesting
work -> **Subnets of**

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

k analysis of gene

on52 and may have
x key parameters
n is to balance positive
ness and stability (Fig.

gn many different
of an unknown
he circuitry or a higher-
on and physical
f information and
lecular to modular cell

ful as the tools
ein sequence data has
odular cell biology
de mechanisms for
iation, post-
mpf, David J. Balding,

with the appropriate

ng regulatory
ook of Statistical

is possible only using

- From the biological side, it is not that clear that the identification of regulatory networks concentrations of mRNA due to the role of additional sources of regulation such as microRNAs. **Biological Networks 2007 ## chapter 3**
- What are the limits of the definition of a regulatory protein? We are interested here in a small 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 proteins, like RNA polymerase itself, are essential for the transcription of all genes that encode proteins. However, because 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(!)** Transcription factors bind to form protein complexes and then they interact with other proteins and/or catalyze reactions. *are not suitable when it comes to expressing interactions that involve more than two participants in a single fashion. The formalism of hypergraphs allows to express non-binary interactions and couplings in a more general context.* -> **Biological Networks 2007 ## chapter 4**
- because high-throughput datasets can be filled with technical and biological noise and have biases and coverage, improved statistics are needed to distinguish signal from noise, as well as the need 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 generate output from the signal input or system perturbation; because no analysis and modeling needed, 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 pathways and genes and pathways. This information forms the skeleton of the entire complete network. However, the annotation standards. For now, such information covers only a tiny fraction of the full network. Progress is made 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 conservative estimate 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 networks. **Understanding biological functions through molecular networks ## Problems of the single gene approach!!**
- However, the change to network biology does not simply entail handing over biology to mathematicians. A good understanding of biology is needed to ask the right questions, to use network analysis tools, and to confirm analysis results by solid experimentation. After all, network biology is a part of biology. The fundamental goal of network biology is the same as molecular biology: to understand biological processes and the mechanisms of human diseases. -> **Understanding biological functions through molecular networks**
- Furthermore, yeast two-hybrid assays produce many false-positive outcomes, and the current network maps may be heavily biased towards connection to functionally important genes. However, these genes have been popular targets for research -> **Kitano, Hiroaki Computational systems biology**
- despite the known poor overall correlation between mRNAs and their protein products, the common assumption that differentially expressed mRNAs impact their respective experimental conditions.

is possible only using
RNA. -> **Biological**

proteins that regulate a
general role in initiating
se generalists, like RNA
however, since their
y for all transcription

!!), but some proteins
ze reactions!!! *Graphs*
tners in an obligate
ld be used in the above

ave different technical
vell as better data
retation of the whole
eded to predict the
method is perfect, more
nesis can be used to
molecular networks ##

relationships between
k and can serve as
twork and is biased
urated by HPRD
ervative estimate of 150
formation, table S1).
tic networks ->
gle gene - process

physicists and
o choose proper
network biology is
nderstand basic
al functions through

rrent hand-crafted
s simply because these
y.pdf
here is an implicit
nditions via differences

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 biological system

- Although clustering analysis provides insight into the “correlation” among genes and biological processes, it 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 (7–10). These methods use information derived from mRNA abundance, so there is limited scope to infer relationships involving transcriptional regulation. Posttranscriptional and posttranslational mechanisms of regulation can 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 components. 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 graph theory. Further integration of other information-flow properties such as edge-direction will enable more 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 regulatory relationships. **Direct regulation relationships must include both phenotype and 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 influence a given functional module to influence and even coordinate multiple cellular processes. It is the high interactional complexity, that single genes rarely specify a phenotype in its entirety. The complexity of genetic networks are now apparent, but a compelling case can be made for a deeper and broader exploration of this model system as an exemplar for more complex eukaryotes. -> **Exploring genetic interactions and networks with yeast**
- The physical-interaction map, generated by large-scale two-hybrid^{51,52} or affinity purification-mass spectrometry identification^{26,43,53,54}, provides a view of the gene products that act as protein complexes and function together as biochemical machines. Rather than physical interactions, a **genetic-interaction map provides functional information**, largely identifying gene products that are 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 approximately 1/2). -> **Exploring genetic interactions and networks with yeast**
- Large-scale genetic analyses reveal that mutations in most eukaryotic genes have little or no effect on the phenotype.

relations via differences

explicitly tested. ->

a xenograft model

biological phenomena, it

proposed to

(9). At present, such

infer causality based on

relation must be

measured with sufficient

regulating proteins, by which

e.g. phosphorylation,

networks with

structure balance

would enable

networks with

reactions

phenotype?

at a network based on

physical interactions and

by the context of the

and the general

s extend beyond a

s no wonder, given this

outlines of a yeast

d more complete

ring genetic

cation followed by

assemble into soluble

information, the

ts that operate in

n interactions more

f less than 1% (REF.

discernable effect. For

example, systematic gene deletion in *S. cerevisiae*, discussed in detail below, produced a set of essential genes. Only ~20% of yeast genes are essential for viability when deleted individually in haploids under laboratory conditions

- Comparisons between aggravating and alleviating effects revealed that, for most functions, effects were either largely aggravating or largely alleviating, but not mixed, an asymmetry termed 'monochromatic'.
- The set of interactions that are observed for a particular query gene can be suggestive of the position of a gene in a genetic-interaction network being highly predictive of its molecular function.

Genetic interactions and networks with yeast

- **Because most genetic interactions do not overlap with physical interactions, the two types of interactions are said to be largely orthogonal**
- Integration attempts in yeast, combining physical protein-protein interaction maps with gene expression data, 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 experimental 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 Interactome
- Experimental artefacts, variability in coverage across data sets, sampling bias towards well-studied 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 protein-protein 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 analysis of protein-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, 20% of proteins in the CYC2008 hand-curated yeast complex data set participate in more than one complex, indicating that many overlapping protein complexes exist in protein-protein interaction networks
- The existence of a GI between two genes does not necessarily imply that these two genes encode proteins that share a functional relationship 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 compartments: Endoplasmic, Cytoskeleton, Golgi, Cytosol, Lysosome(or Vacuole), Mitochondrion, Plasma, Nucleus, Peroxisome and Extracellular, where Lysosome only exists in animal cells. For each compartment, a PSLIN of this compartment is constituted by the proteins localized in this compartment and their interactions. With the subcellular localization information of proteins, the PSLINs can be derived by mapping the global PIN to each compartment separately. -> **Rechecking the Centrality-Likelihood of Protein Subcellular Localization Interaction Networks**

remarkable result:
growing in standard

nal groups, interac-
cal feature that they

f its function, with the
r role. -> **Exploring**

types of interaction are

coexpression profiles,
pression profiles than
et al., 2002). These
eyond the funda-
networks and
sed to estimate the
an Disease

well-studied processes,
posi- tives and false
biological networks.

nected clusters rather
integration of protein-

functions, and
07 of 1,628 proteins in
x. -> detecting

s code for interacting
ies that the two genes
process or pathway; or
> **Genetic interaction**

d into 11
nondrion, Endosome,
s. For each
s compartment and
n be generated by
lethality Rule in the

Network construction methods and evaluation

- Information theoretic approaches to GRN reconstruction have two major strengths. The first is that mutual information is able to capture nonlinear associations between variables, a feature seen in many biological systems. Second, it has been demonstrated that the use of the data processing inequality (DPI), which states that in a noisy system $X \rightarrow Y \rightarrow Z$, knowledge of Z cannot give more information about X than Y can, is useful 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 other techniques in the precision and recall of direct regulatory links **## Information theory and its applications to biological systems/ From molecular information processing to network inference**
- Several alternative network inference approaches exist (4, 5) that provide better, more robust networks, and can incorporate expert or domain knowledge. 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 targets. 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 essential for their 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; network inference performance depends largely on the specific implementation of each individual method. **for robust gene network inference**
- We next analyzed how method-specific biases influenced the recovery of different network motifs (network motifs), and we observed characteristic trends for different method categories. 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 regulators.

n

first is that mutual
n expression data.
which states that for a
n give about X [19,65],
onstructing simulated
and relevance networks
d signal transduction

robust candidate
d to investigate gene-
between disease cases

advent of high-
shots of the
ge is to
d its target genes; the
m of crowds for

critical for their
network inference
superior and that
-> **Wisdom of crowds**

ectivity patterns
(**Fig. 2c**). For example,
based methods,
. The reason for this is
to the expression of
ow mutually dependent

target genes. However, the assumption of independence is violated for genes regulated by transcription factors, as in the case of feed-forward loops. Indeed, linear cascades were not predicted by regression and Bayesian-network methods. This shows that current methods trade off between performance on cascades and performance on feed-forward loops. -> **Wisdom of crowds for robust gene network inference**

robust gene network inference

- Network inference methods have complementary advantages and limitations under different conditions. This suggests that combining the results of multiple inference methods could be a good strategy for network 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. The modules perform unique functions to most of the identified modules in both networks
- When constructing a compendium of microarrays for global network inference, one should avoid moving 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 due to heuristic searches, which are often too costly for systematic data resampling and may be biased on smaller networks. Information theoretic methods performed better than correlation-based methods. Both two approaches had similar biases in predicting regulatory relationships. They also performed poorly on regression and Bayesian- network methods on feed-forward loops, fan-ins and fan-outs (i.e., disconnected parts of the network), but they had an increased rate of false positives for cascades. -> **Wisdom of crowds for robust gene network inference**

A fundamental assumption of network inference algorithms is that mRNA levels of transcription factors and their targets tend to be correlated; we found that this is true for *E. coli*, but not for *S. cerevisiae*. 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, such as promoter sequences and data sets for transcription-factor binding and chromatin modifications. -> **Wisdom of crowds for robust gene network inference**

Wisdom of crowds for robust gene network inference

- We will show in this chapter that there are many ways to express the reverse-modeling problem 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 optimization problem is NP-hard ->
- Transcription of DNA into RNA is the first — and often most regulated — step in gene expression. Many genes are regulated at the transcriptional level, and 5-10% of protein-encoding genes encode regulatory proteins, the interactions between regulated genes and regulatory proteins constitute a transcriptional network or genetic network. -> Biological Networks 2007 Chapter 4
- It is therefore likely that comprehensive mapping of the quantitative genetic interaction network requires the integration of a number of datasets from different screening approaches, similar to the recent success in mapping the physical protein-protein interaction (PPI) networks in yeast and human -> **Quantitative genetic interactions in yeast - Comparative evaluation and integrative analysis**
- Compared to PPI networks, an additional challenge originates from the quantitative nature of gene expression data.

er, making dependent
more accurately
is experience a trade-
om of crowds for

erent contexts, which
egy for improving

ms (Supplementary
is allowed us to assign

uld thus avoid any bias
robust gene network

y because they use
better suited for
ed methods, but the
rmed better than
the more densely
cades. -> **Wisdom of**

cription factors and
evisiae. Although the
nown regulatory
visiae is likely due to
eukaryotes, which
onal inputs, such as
cation7. -> **Wisdom of**

f regulatory networks
st point to solve is the
one is known to be

pression. As most
code regulatory
complex web called

networks will require
ent efforts to complete
ve maps of genetic

re of the genetic

interaction datasets; instead of comparing the overlap in binary terms, such as presence of physical interaction, here we should take into account the full spectrum of genetic interactions from extreme cases of negative interactions (i.e., synthetic sick and lethality) to the positive class of pairs (e.g., masking and suppression subcategories) [2,3,17].

1. 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 suggests 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 suggests that the gene products operate in concert or in series within the same pathway [4].

2. Genome-scale, quantitative analysis of genetic interactions. We consider a digenic interaction where a double mutant that shows a significant deviation in fitness compared with the expected multiplicative effect of combining two single mutants (6). Negative interactions refer to a more severe fitness defect than expected, with the extreme case being synthetic lethality; positive interactions refer to double mutants with a less severe fitness defect than expected. To quantitatively score genetic interactions in large-scale SGA, we have developed a model to estimate fitness defects directly from double-mutant colony sizes.

Landscape of a Cell

3. There have been some attempts to investigate the temporal properties for individual protein-protein interactions by integrating PPI data with time-course gene expression data [21-29]. <- Deconvolution of 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. -> A node's centrality is a measure of its influence on the network.
- degree of a vertex alone, as a specific centrality measure, is not sufficient to distinguish between essential and non-essential genes (Wuchty (2002)), that in protein networks there is no relation between degree and robustness against amino-acid substitutions (Hahn et al. (2004)), and that for biological networks 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 networks**
- Topological features of the protein networks have been demonstrated to reflect the functional properties of the interacting genes. For example, essential genes in yeast tend to be well connected and highly central in the protein network (Jeong et al. 2001; Wuchty and Almaer 2005). Furthermore, globally central nodes are often found to be essential (Jeong et al. 2001; Wuchty and Almaer 2005).

or absence of a
ctions, ranging from
asses of interacting

enotype that is more
interaction indicates
robustness of the

notype that is less
raction indicates that

ction as a double
cative effect of
effect than expected,
ants with a less severe
GA screens, we
s -> **The Genetic**

teins and protein
tecting temporal

Axioms of centrality
ethal proteins clearly
network connectivity
ical network analysis
ki and Schreiber
Gene regulatory

ctionality of the
obally centered in the
ntered interactions are

protein network (Jeong et al., 2001, Wuchty and Almaas, 2005). Furthermore, globally conserved 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 gene expression profiling study by using interactome-transcriptome analysis. We find high centrality of 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 generally conserved. Non-essential genes were not (Wuchty and Almaas, 2005). The study also indicates that topological features are conserved throughout different species. -> **Interactome-transcriptome analysis reveals the high centrality of genes differentially expressed in lung cancer tissues**
- The problem of choosing initial nodes as source spreaders to achieve maximum scale of spreading is a *source maximization problem*¹². Our research focuses on the strategy of choosing a set of 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, closeness centrality^{4,5}, betweenness centrality⁸ counting the number of shortest paths through a certain node, eigenvector centrality⁹ based on the idea that relations with more influential neighbors confer greater influence, and shell decomposition¹⁰ that correlates with the outcome of supercritical spreading original nodes^{11–13}—**rely only on topology**, even if an underlying process can be indirectly associated with the network. 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 centrality (0.16), and betweenness centrality (0.001) of essential genes in yeast are about twice those in nonessential genes -> **Predicting essential genes based on network and sequence data**
- many bottlenecks also tend to be hubs). Therefore, we further investigate which one of the two is a better predictor of protein essentiality in both regulatory and interaction networks. -> **Bottlenecks in Protein Networks/ Correlation with Gene Essentiality and Expression Dynamics**
- we observed that bottlenecks (both nonhub–bottlenecks and hub–bottlenecks) have a statistically higher number of products of essential genes, whereas hub– nonbottlenecks are surprisingly not essential. This suggests that it is the betweenness that is a stronger determinant of the essentiality of a protein in the 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 topology implies biological importance. One of the first connections between the two in the context of the protein interaction network, the so-called centrality-lethality rule, was observed by Jeong and co-workers (Jeong et al., 2001). **Hubs in the Yeast Protein Interaction Network Tend To Be Essential/ Reexamining the Centrality-Lethality Rule in the Network Topology and Essentiality**
- Since then the correlation between degree and essentiality was confirmed by other studies

interacted interactions are
(Wuchty and Almaas,
essentially expressed in

analysis by considering
significance of
recent microarray
with centrality in these
proteome-transcriptome
ties
global hubs, whereas
these global hubs are
the high centrality of

leading is defined as *in*
of critical nodes as
networks

ors a node interacts
in node, eigenvector
er importance, or the k-
ating in specific
ciated in some cases. In
tably depends on the

gree (11.55), clustering
ut twice as large as
ce analysis #

these two quantities is
The Importance of
namics

rong tendency to be
Thus, we determined
n the regulatory
ation with Gene

eristics of a protein,
logical prominence
xt of a protein
lleagues ->**Why Do**
Connection between

ies [4–7]. but until

since then the correlation between degree and essentiality, was confirmed by other studies. recently there was no systematic attempt to examine the reasons for this correlation ->W

Yeast Protein Interaction Network Tend To Be Essential/ Reexamining the Connection between 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 essential in yeast but nonessential in mouse (E2N), and 1,149 genes are nonessential in both yeast and mouse (N2N). 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: gene knockout (Giaever et al., 2002; Chen et al., 2015), gene knockdown (Harborth et al., 2001; Roemer et al., 2003) and transposon mutagenesis (Gallagher et al., 2007; Langridge et al., 2007). 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, especially for humans. -> **Predicting Essential Genes and Proteins Based on Machine Learning and Network Topological Features/ A Comprehensive Review**

- 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 methods are those based on machine learning, a method in which computers make and improve predictions based on data through learning algorithms (see next sections for details). On the other hand, of the 16 articles reporting the utilization of machine learning for the prediction of essential genes and proteins, 12 (75%) utilize the utilization of network topological features as learning attributes (features that describe a node or an edge in details in the next section). Therefore, this brief analysis of these 34 papers suggests a strong correlation between machine learning and network topological features regarding the prediction of essential genes and proteins.

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 (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 essential in one certain growth condition might not play as critical role in another condition (Tong et al., 2002; Tong et al., 2011). These are the so-called conditionally essential genes, i.e., nonessential genes that become essential depending on the environment conditions. The development of machine learning approaches to predict 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 unclear. 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 individually deleted cause lethality when they are combined as a double mutant. -> **Exploring genetic interactions and synthetic lethality in yeast**
- we estimate that a global network will contain ~200,000 synthetic-lethal interactions. To

... [1, 2], ...

Why Do Hubs in the between the Network

their essentialities
se genes into four
essential in both yeast
genes are essential in
and mouse (N2N) ->

essential genes: gene
L; Ji et al., 2001;
, 2009). These

ne-consuming and
organisms, especially

Network Topological

utilization of network
nal methods used were
ictions based on some
e 14 articles reporting
86%) report the
certain instance; see
rong link between
genes and proteins. ->

logical Features/ A

of essential genes
eld (Chen and Xu, 2005;
ntial genes, that is,
ed as essential under
2004; Nichols et al.,
become essential
ches for the prediction
n challenge in the near
explored. ## global

on Machine Learning

ually mutated cause
and networks with

put this number in

context, there are $\sim 1,000$ essential genes in yeast, for which a single mutation leads to a phenotype. However, 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, explaining why the analytical power of single-gene effects on many phenotypes has been so limited.

interactions and networks with yeast

- Moreover, our approach of focusing on nonhub–bottle-necks is useful for finding proteins involved in different processes and are involved in cross-talk. -> The Importance of Bottlenecks in Protein-Protein Correlation with Gene Essentiality and Expression Dynamics
- This indicates the existence of a large number of nodes with high betweenness but low connectivity (hubs and nodes). Importantly, such nodes are absent in computer-generated, random scale-free networks. mention of bottlenecks: **High-Betweenness Proteins in the Yeast Protein Interaction Network**

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 network
 - Even a comparatively small subnetwork of this size is still a challenge to visualize in a network. The authors assess its quality in two ways. *First, they determine whether the genes of the subnetwork are enriched for specific cellular process categories in the Gene Ontology database, which is a standard method for assessing the biological relevance of a set of genes, and this is a wonderful strength of the paper—*the authors experimentally validate the interactions of 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 connectivity 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 direct binding. -> **Wisdom of crowds for robust gene network inference**
- Gene essentiality data of yeast were manually compiled from the Comprehensive Yeast Gene Deletion Project (<http://mips.helmholtz-muenchen.de/genre/proj/yeast/>) and large-scale experiments [26]. 1,178 essential and 4,904 nonessential yeast genes. -> Network rewiring is an important mechanism for essentiality change
- [currently, the available essential genes and protein databases are DEG (Zhang and Lin, 2013), 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, identify features associated with essentiality and, finally, develop computational methods for predicting essential genes and proteins. -> **Predicting Essential Genes and Proteins Based on Machine Learning / Network Topological Features / A Comprehensive Review**

lethal phenotype, but
ic-lethal interaction.
s, and begins to explain
-> **Exploring genetic**

ns that mediate
rotein Networks/
onnectivity (HBLC
etworks ->First
network 2005

etworks!! e.g
sightfully, so the
the subnetwork are
which they are. Second—
some of the first

the 'expression space'
onal algorithms ->
g of regulatory

on factor on the target
- mine physical

Genome Database
The dataset contained
mechanism of gene

009), CEG (Ye et al.,
data have enabled
loration, reveal which
pposed to identify
ine Learning and

Network Topological Features: A Comprehensive Review

- There exist many different databases of a certain type of interaction (e.g., DIP, BioGRID and usually, these databases are regularly updated. Different databases or newer versions of have different sets of interactions that, in turn, will give rise to new networks with distinct consequently, different values of network topological features. As an example, we can cite Hwang et al. (2009) and Acencio and Lemke (2009). In both studies, PINs of *Saccharomyces* were created; however, the interactions of the PIN constructed in the study by Hwang et al. (2009) were 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 results must be the same kind and the same data! That means not only the same database because it is updated regularly, 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 have more studies and therefore tend to have higher values of network topological)** present in the data 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 the prediction performance.
- ROC curves in R <http://blog.revolutionanalytics.com/2016/11/calculating-auc.html>
- Dygraphs for R Interactive plots for HTML files

and IntAct for PPIs) and, a given database will extract structures and, re the studies by *Saccharomyces cerevisiae* were (2009) were collected from the study by Acencio. Therefore, the prediction of the reference networks because they are updated

essential from non-essential **essentiality and proteins are the focus of** in the networks mainly

essential genes and this moderate