



Giovanni Boniolo

FIRC Institute of Molecular Oncology (IFOM)- Milano
European School of Molecular Medicine (SEMM) - Milano
Faculty of Medicine - University of Milano

Could We Formalize Biochemical Processes?
And why?

- The theoretical framework
- The proposal
- The possible research programme

mathematics	geometry	logic	dynamics (differential equations)
-------------	----------	-------	--------------------------------------

Classical mechanics	Vectorial calculus	Euclidean geometry	Classical logic	Newton law
----------------------------	--------------------	--------------------	-----------------	------------

General relativity	Differential topology	Riemannian geometry	Classical logic	Einstein equations
---------------------------	-----------------------	---------------------	-----------------	--------------------

Quantum mechanics	Hilbert space	Euclidean geometry	Quantum logic	Schrödinger equation
--------------------------	---------------	--------------------	---------------	----------------------

Molecular biology	NO	NO	YES	YES
--------------------------	----	----	------------	------------

From statements to molecules



From classical logic to resource logic

NON CLASSICAL COMPUTATIONAL LOGIC

ZSYNTAX

	Linguistic case	Biological case
Premises (Hypothesis): Γ ↓ Inferential process ↓ Conclusion (Thesis): C	Conjunction of statements ↓ Classical logic ↓ Statement	Aggregate of molecules ↓ Biochemical reactions (non classical logic) ↓ Aggregate of molecules

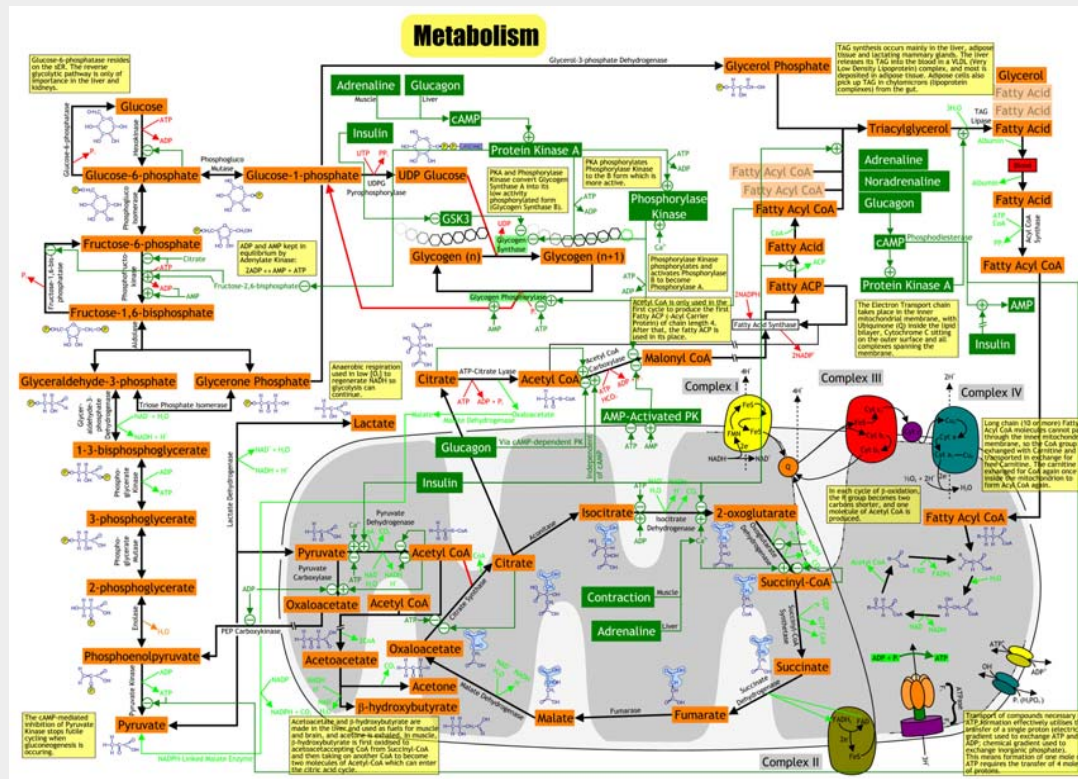
$$\Gamma \quad | \text{---} \quad C$$

2 EXAMPLES

$\Gamma = (\text{D-glucose, hexokinase, phosphoglucoisomerase, phosphofructokinase, ATP, ATP})$

Theorem 1:

$\Gamma \mid - (\text{fructose-1,6-phosphate})$



Theorem 2:

$(MDM2, p53) \mid - MDM2$

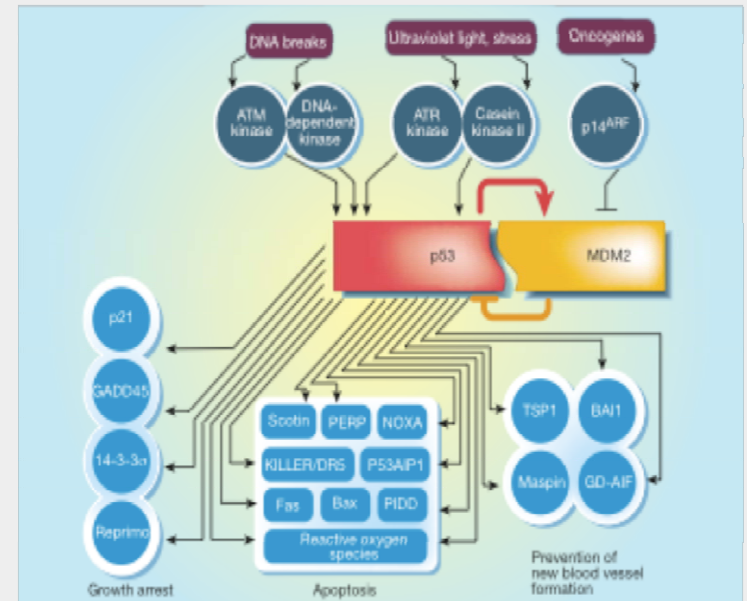
Nature 408 307 (2000)

news and views feature

Surfing the p53 network

Bert Vogelstein, David Lane and Arnold J. Levine

The p53 tumour-suppressor gene integrates numerous signals that control cell life and death. As when a highly connected node in the Internet breaks down, the disruption of p53 has severe consequences.



The three operators of our ZSYNTAX

1) Z-interaction (indicated by \odot)

we denote by $A\odot B$ the type of molecules which results from the interaction of two types of molecules A and B.



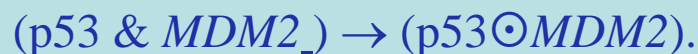
2) Z-conjunction (indicated by $\&$)

we denote by $A\&B\&C\&D$ the type of aggregate constitutes by the types of molecules A, B, C, D.

(D-glucose $\&$ hexokinase $\&$ phosphoglucoisomerase $\&$ phosphofructokinase)

3) Z-conditional (indicated by \rightarrow)

we denote by $A\&C \rightarrow B$ the fact that there is a transition path from an aggregate of type $A\&C$ to an aggregate of type B.



The two kinds of valid formulas for ZSYNTAX

1) The **empirically valid formulas** (EVF). They represent singular reactions and their validity depends on the fact that the processes they describe are empirically corroborated. We can have EVFs representing

(i) that two molecules interact, e.g.



(ii) that the interaction delivers certain products in a biochemical sense, e.g. (D-glucose-6-phosphate \odot phosphoglucosomerase \rightarrow D-fructose-6-phosphate)

(iii) that the interaction allows the gene expression, e.g.,



2) The **logically valid formulas**. Their validity rests only on the definitions of the logical operators inside the language we have constructed. They give the rules to move from one EVF to another EVF

Example of a logically valid formula:

$$(A \rightarrow B \ \& \ A) \rightarrow B$$

This can be expressed as the *modus ponens* rule:

1. If A then B
2. A

B

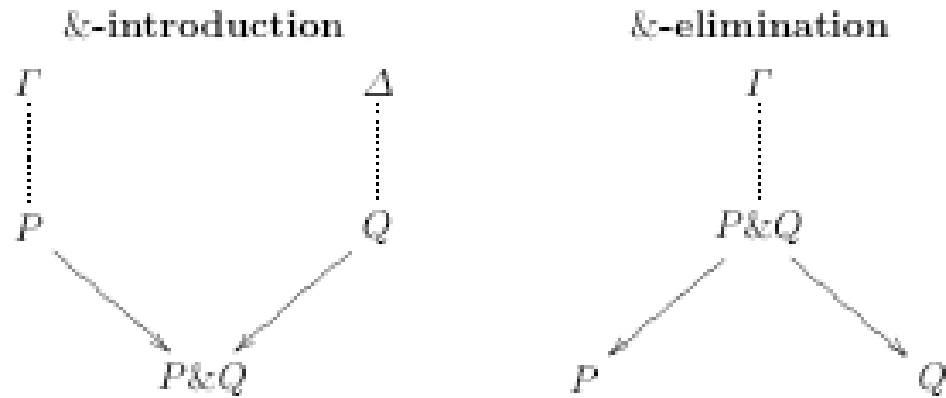
1. $A \rightarrow B$
2. A

B

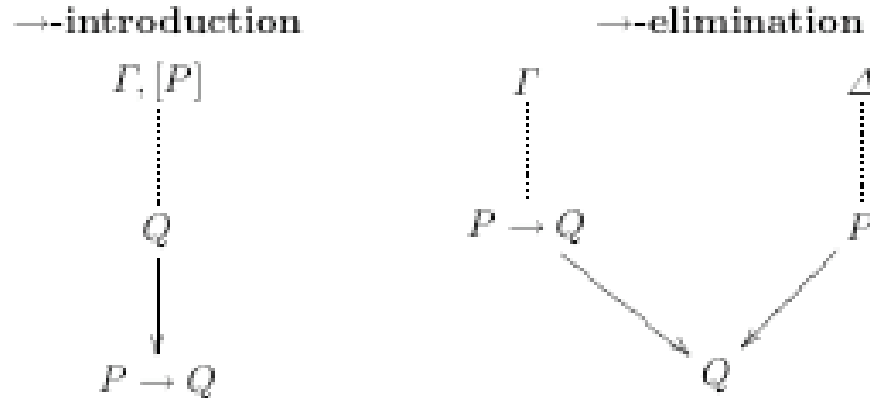
1. $(p53 \ \& \ MDM2) \rightarrow (p53 \odot MDM2)$
2. $(p53 \ \& \ MDM2)$

 $(p53 \odot MDM2)$

Our primitive logical rules for $\&$ are the following:



The rules for \rightarrow are the following:

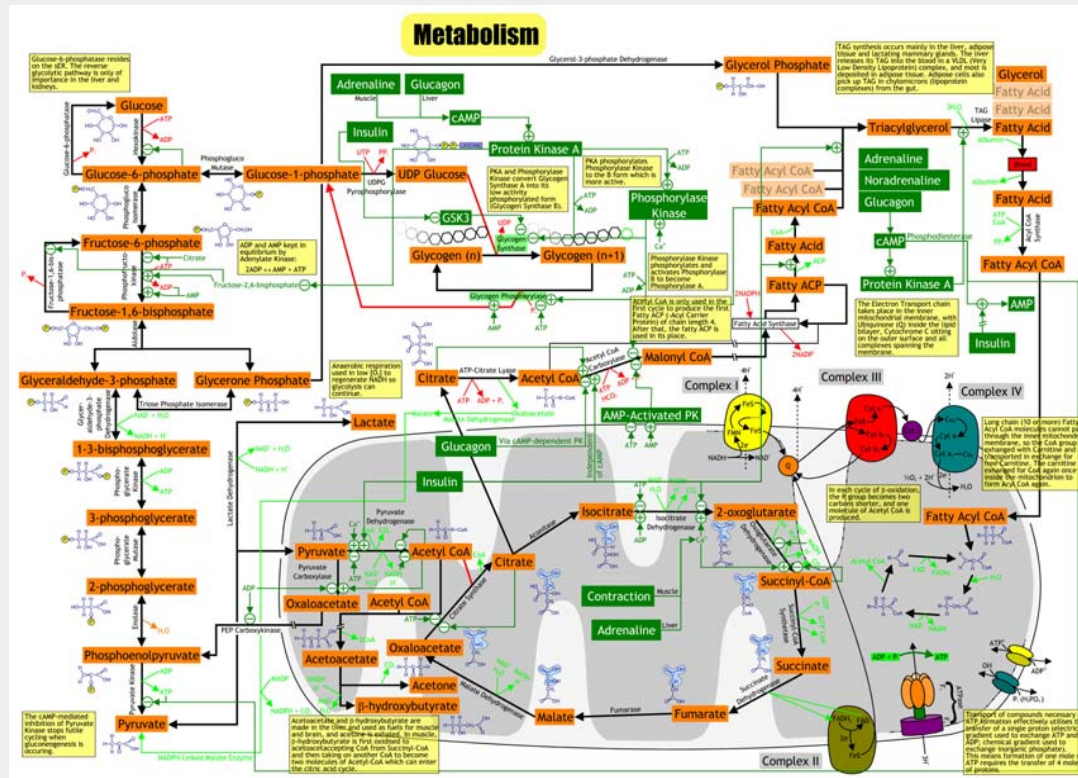


2 EXAMPLES
OF
DEMONSTRATION

$\Gamma =$ (D-glucose, hexokinase, phosphoglucoisomerase, phosphofructokinase, ATP, ATP)

Theorem 1:

$\Gamma \mid -$ (fructose-1,6-phosphate)



1. $\text{Glc} \ \& \ \text{HK} \ \& \ \text{GPI} \ \& \ \text{PFK} \ \& \ \text{ATP} \ \& \ \text{ATP}$	IA
2. $\text{Glc} \ \& \ \text{HK}$	From 1 by $\&E$
3. GPI	From 1 by $\&E$
4. PFK	From 1 by $\&E$
5. ATP	From 1 by $\&E$
6. ATP	From 1 by $\&E$
7. $\text{Glc} \ \& \ \text{HK} \ \rightarrow \ \text{Glc} \ \odot \ \text{HK}$	EVF
8. $\text{Glc} \ \odot \ \text{HK}$	From 2,7 by $\rightarrow E$
9. $(\text{Glc} \ \odot \ \text{HK}) \ \& \ \text{ATP}$	From 5,8 by $\&I$
10. $(\text{Glc} \ \odot \ \text{HK}) \ \& \ \text{ATP} \ \rightarrow \ (\text{Glc} \ \odot \ \text{HK}) \ \odot \ \text{ATP}$	EVF
11. $(\text{Glc} \ \odot \ \text{HK}) \ \odot \ \text{ATP}$	From 9,10 by $\rightarrow E$
12. $(\text{Glc} \ \odot \ \text{HK}) \ \odot \ \text{ATP} \ \rightarrow \ \text{G6P} \ \& \ \text{HK} \ \& \ \text{ADP}$	EVF
13. $\text{G6P} \ \& \ \text{HK} \ \& \ \text{ADP}$	From 11,12 by $\rightarrow E$
14. G6P	From 13 by $\&E$
15. HK	From 13 by $\&E$
16. ADP	From 13 by $\&E$
17. $\text{G6P} \ \& \ \text{GPI}$	From 3,14 by $\&I$
18. $\text{G6P} \ \& \ \text{GPI} \ \rightarrow \ \text{G6P} \ \odot \ \text{GPI}$	EVF
19. $\text{G6P} \ \odot \ \text{GPI}$	From 17,18 by $\rightarrow E$
20. $\text{G6P} \ \odot \ \text{GPI} \ \rightarrow \ \text{F6P} \ \& \ \text{GPI}$	EVF
21. $\text{F6P} \ \& \ \text{GPI}$	From 19,20 by $\rightarrow E$
22. F6P	From 21 by $\&E$
23. GPI	From 21 by $\&E$
24. $\text{F6P} \ \& \ \text{PFK}$	From 4,22 by $\&I$
25. $\text{F6P} \ \& \ \text{PFK} \ \rightarrow \ \text{F6P} \ \odot \ \text{PFK}$	EVF
26. $\text{F6P} \ \odot \ \text{PFK}$	From 24,25 by $\rightarrow E$
27. $(\text{F6P} \ \odot \ \text{PFK}) \ \& \ \text{ATP}$	From 6,26 by $\&I$
28. $(\text{F6P} \ \odot \ \text{PFK}) \ \& \ \text{ATP} \ \rightarrow \ (\text{F6P} \ \odot \ \text{PFK}) \ \odot \ \text{ATP}$	EVF
29. $(\text{F6P} \ \odot \ \text{PFK}) \ \odot \ \text{ATP}$	From 27,28 by $\rightarrow E$
30. $(\text{F6P} \ \odot \ \text{PFK}) \ \odot \ \text{ATP} \ \rightarrow \ \text{F1,6P} \ \& \ \text{PFK} \ \& \ \text{ADP}$	EVF
31. $\text{F1,6P} \ \& \ \text{PFK} \ \& \ \text{ADP}$	From 29,30 by $\rightarrow E$
32. F1,6P	From 31 by $\&E$
33. $(\text{Glc} \ \& \ \text{HK} \ \& \ \text{GPI} \ \& \ \text{PFK} \ \& \ \text{ATP} \ \& \ \text{ATP}) \ \rightarrow \ \text{F1,6P}$	From 1-32 by $\rightarrow I$

Theorem

$$\text{Glc} \ \& \ \text{HK} \ \& \ \text{GPI} \ \& \ \text{PFK} \ \& \ \text{ATP} \ \& \ \text{ATP} \ \vdash \ \text{F1,6P}$$

Demonstration

1. $\text{Glc} \ \& \ \text{HK} \ \rightarrow \ \text{Glc} \ \odot \ \text{HK}$
2. $(\text{Glc} \ \odot \ \text{HK}) \ \& \ \text{ATP} \ \rightarrow \ (\text{Glc} \ \odot \ \text{HK}) \ \odot \ \text{ATP}$
3. $(\text{Glc} \ \odot \ \text{HK}) \ \odot \ \text{ATP} \ \rightarrow \ \text{G6P} \ \& \ \text{HK} \ \& \ \text{ADP}$
4. $\text{G6P} \ \& \ \text{GPI} \ \rightarrow \ \text{G6P} \ \odot \ \text{GPI}$
5. $\text{G6P} \ \odot \ \text{GPI} \ \rightarrow \ \text{F6P} \ \& \ \text{GPI}$
6. $\text{F6P} \ \& \ \text{PFK} \ \rightarrow \ \text{F6P} \ \odot \ \text{PFK}$
7. $(\text{F6P} \ \odot \ \text{PFK}) \ \& \ \text{ATP} \ \rightarrow \ (\text{F6P} \ \odot \ \text{PFK}) \ \odot \ \text{ATP}$
8. $(\text{F6P} \ \odot \ \text{PFK}) \ \odot \ \text{ATP} \ \rightarrow \ \text{F1,6P} \ \& \ \text{PFK} \ \& \ \text{ADP}$

Theorem 2:

$(MDM2, p53) \mid - MDM2$

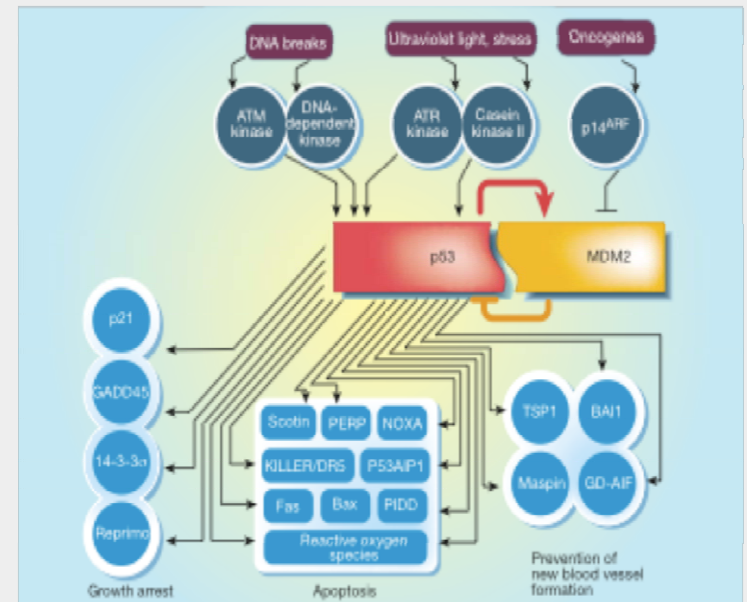
Nature 408 307 (2000)

news and views feature

Surfing the p53 network

Bert Vogelstein, David Lane and Arnold J. Levine

The p53 tumour-suppressor gene integrates numerous signals that control cell life and death. As when a highly connected node in the Internet breaks down, the disruption of p53 has severe consequences.



1.	$TP53 \ \& \ TP53 \ \& \ MDM2 \ \& \ U \ \& \ P$	IA
2.	$TP53$	From 1 by $\&E$
3.	$TP53$	From 1 by $\&E$
4.	$MDM2$	From 1 by $\&E$
5.	U	From 1 by $\&E$
6.	P	From 1 by $\&E$
7.	$TP53 \ \& \ MDM2$	From 2,4 by $\&I$
8.	$TP53 \ \& \ MDM2 \ \rightarrow \ TP53 \ \odot \ MDM2$	EVF
9.	$TP53 \ \odot \ MDM2$	From 7,8 by $\rightarrow E$
10.	$TP53 \ \odot \ MDM2 \ \rightarrow \ MDM2$	EVF
11.	$MDM2$	From 9,10 by $\rightarrow E$
12.	$MDM2 \ \& \ TP53$	From 3,11 by $\&I$
13.	$MDM2 \ \& \ TP53 \ \rightarrow \ MDM2 \ \odot \ TP53$	EVF
14.	$MDM2 \ \odot \ TP53$	From 12,13 by $\rightarrow E$
15.	$(MDM2 \ \odot \ TP53) \ \& \ U$	From 5,14 by $\&I$
16.	$(MDM2 \ \odot \ TP53) \ \& \ U \ \rightarrow \ (MDM2 \ \odot \ TP53) \ \odot \ U$	EVF
17.	$(MDM2 \ \odot \ TP53) \ \odot \ U$	From 15,16 by $\rightarrow E$
18.	$(MDM2 \ \odot \ TP53) \ \odot \ U \ \rightarrow \ MDM2 \ \& \ (TP53 \ \odot \ U)$	EVF
19.	$MDM2 \ \& \ (TP53 \ \odot \ U)$	From 17,18 by $\rightarrow E$
20.	$TP53 \ \odot \ U$	From 19 by $\&E$
21.	$(TP53 \ \odot \ U) \ \& \ P$	From 6,20 by $\&I$
22.	$(TP53 \ \odot \ U) \ \& \ P \ \rightarrow \ (TP53 \ \odot \ U) \ \odot \ P$	EVF
23.	$(TP53 \ \odot \ U) \ \odot \ P$	From 21,22 by $\rightarrow E$
24.	$(TP53 \ \odot \ U) \ \odot \ P \ \rightarrow \ \underline{d}(TP53) \ \& \ U \ \& \ P$	EVF
25.	$\underline{d}(TP53) \ \& \ U \ \& \ P$	From 23,24 by $\rightarrow E$
26.	$\underline{d}(TP53)$	From 25 by $\&E$
	$(TP53 \ \& \ TP53 \ \& \ MDM2 \ \& \ U \ \& \ P) \ \rightarrow \ \underline{d}(TP53)$	From 1-26 by $\rightarrow I$

Theorem

$TP53 \ \& \ TP53 \ \& \ MDM2 \ \& \ U \ \& \ P \ \vdash \ \underline{d}(TP53)$

Demonstration

1. $TP53 \ \& \ MDM2 \ \rightarrow \ TP53 \ \odot \ MDM2$
2. $TP53 \ \odot \ MDM2 \ \rightarrow \ MDM2$
3. $MDM2 \ \& \ TP53 \ \rightarrow \ MDM2 \ \odot \ TP53$
4. $(MDM2 \ \odot \ TP53) \ \& \ U \ \rightarrow \ (MDM2 \ \odot \ TP53) \ \odot \ U$
5. $(MDM2 \ \odot \ TP53) \ \odot \ U \ \rightarrow \ MDM2 \ \& \ (TP53 \ \odot \ U)$
6. $(TP53 \ \odot \ U) \ \& \ P \ \rightarrow \ (TP53 \ \odot \ U) \ \odot \ P$
7. $(TP53 \ \odot \ U) \ \odot \ P \ \rightarrow \ \underline{d}(TP53) \ \& \ U \ \& \ P$

Why should we use it?

- 1) To formalize molecular biology
- 2) To perform text mining
- 3) To predict biological reactions and biological products

TEXT MINING

Seeking a New Biology through Text Mining

Andrey Rzhetsky,^{1*} Michael Seringhaus,² and Mark Gerstein²

¹University of Chicago, Chicago, IL 60637, USA

²Yale University, New Haven, CT 06510, USA

*Correspondence: arzhetsky@uchicago.edu

DOI 10.1016/j.cell.2008.06.029

Cell 134, July 11, 2008 ©2008 Elsevier Inc.

Tens of thousands of biomedical journals exist, and the deluge of new articles in the biomedical sciences is leading to information overload. Hence, there is much interest in text mining, the use of computational tools to enhance the human ability to parse and understand complex text.

Imagine that a graduate student enters the U.S. Library of Congress with the goal of retrieving all texts relevant to protein glycosylation. Her problem is straightforward, known among text miners as *information retrieval* (IR). If the student must not only find the books but also flag the most important concepts she encounters in each, she is performing *named entity recognition* (NER). Undaunted by her workload, imagine she decides to identify relations between concepts, such as “protein BAD binds to protein BAX” (called *information extraction* or IE). Then she takes on additional tasks such as question/answer (QA) and text summarization (TS). Computational IR, NER, IE, QA, and TS are all part of text mining and belong to the larger field of *natural language processing* (NLP), which itself is a part of *artificial intelligence* (AI) that aims to recreate or surpass the computational ability of the human brain. Although multiple definitions exist, text mining is typically associated with information retrieval, extraction, and synthesis, with a special emphasis on gaining new knowledge (Table 1).

available online). This is because biology and medicine are unusually rich in terminology: the collective vocabulary used by biomedicine incorporates many millions of terms. The exact number is unknown and constantly in flux. Because this vocabulary is large and dynamic, new terms emerge rapidly and erratically. As a result, the same real-world object may have numerous names (synonyms), whereas distinct objects can be identified with the same name (homonyms). The terms that most notoriously suffer from synonym and homonym abundance are gene and protein names (Hirschman et al., 2002; Wilbur et al., 1999). A given gene may be denoted by several dozen synonymous names—for example, the *Drosophila* genes *br* and *mod/modg4* have 82 and 64 aliases, respectively. Even worse, vastly different and incompatible naming systems are used in different species, and gene aliases themselves merge into intricate semantic networks that connect gene names, pop-culture catch phrases, idioms, and common everyday utterances borrowed from multiple languages.

PREDICTION

Abduction

Let us suppose that we know the final molecule of a biological process (call it **B**), but we do not know exactly either **1**) which is the initial aggregate (call it Γ), or **2**) which are the right biological reactions leading from that initial aggregate to the final molecule.

In this case the automated prover may indicate that something is missing:

maybe **1**) an additional molecule in the initial aggregate (that is, **we should find a new molecule**),

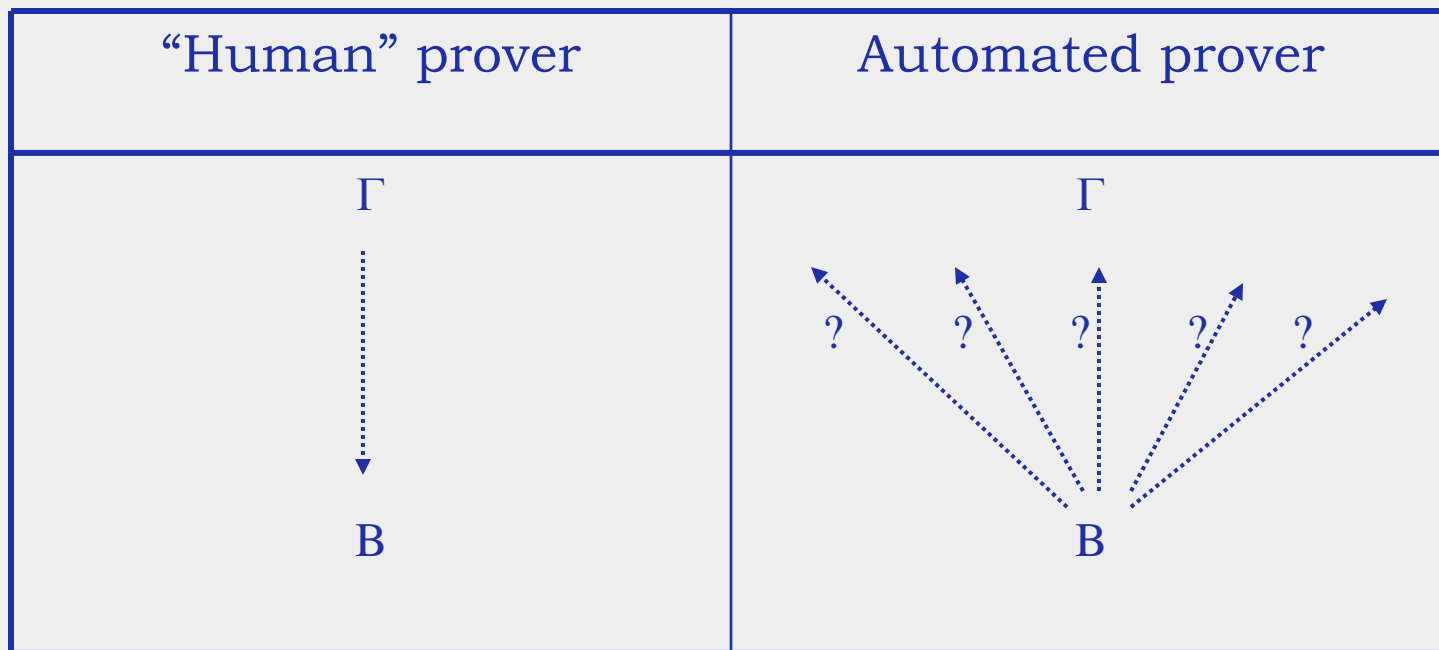
maybe **2**) an additional reaction (that is, **we should find a new reaction**).

But not only A well-constructed automated prover should be able to extract a range of alternative hypotheses concerning how to adjust the data and/or the reactions involved and, thus, to indicate us what should miss.

At this point, we could go back to the lab and investigate the hypotheses which have been automatically generated by the computer.

Typically, in automated theorem proving, one starts from the conclusion B and looks back for the possible premises Γ .

Reiterating this inverse process the automated prover constructs several possible paths. Going through all possible paths (and using suitable heuristics to prune the wrong ones), it eventually arrives at the right initial aggregate of molecules Γ .



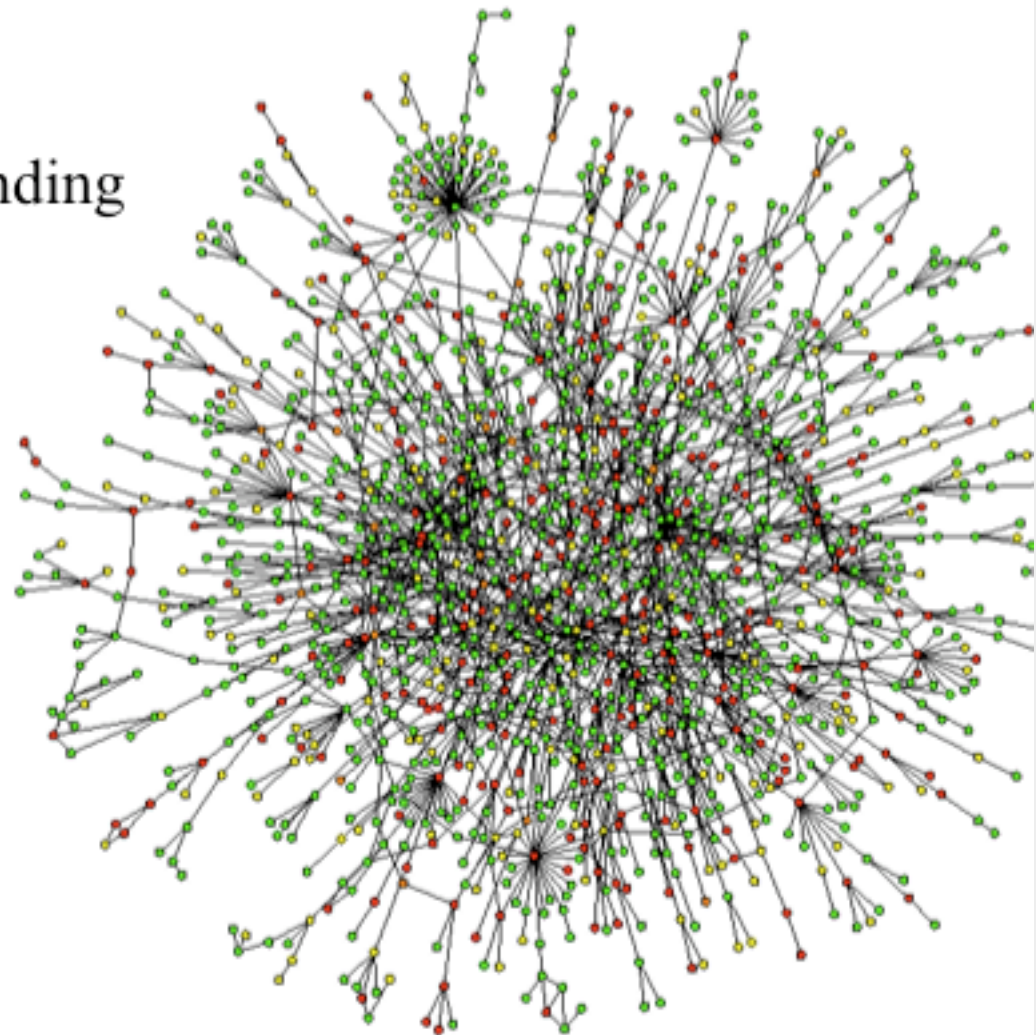
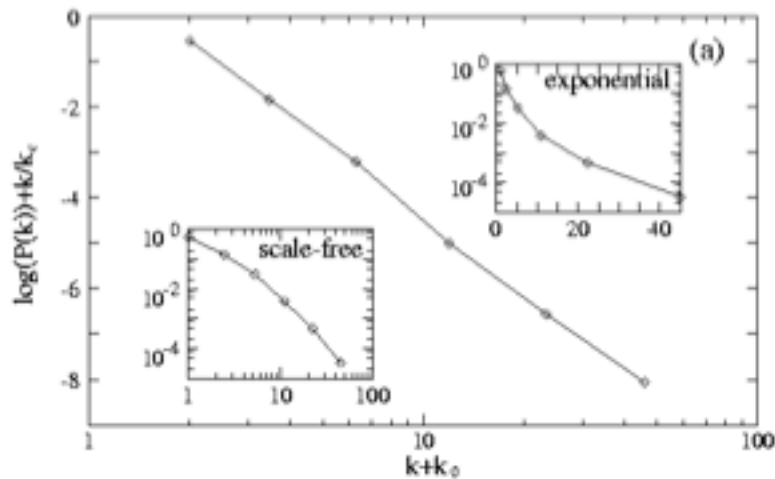
Are there competitors to ZSYNTAX?

Zsyntax and biological networks

Topology of the protein network

Nodes: proteins

Links: physical interactions-binding



$$P(k) \sim (k + k_0)^{-\gamma} \exp\left(-\frac{k + k_0}{k_\tau}\right)$$

to elucidate the age of retrogene movements, we dated the human duplications involving X-linked parents or retrogenes both by comparison to the mouse genome sequence and by sequence divergence analysis (16). Most copies that escape X linkage (12/15) as well as most copies that obtain X linkage (10/13) originated before the human-mouse split (Fig. 2, tables S7 and S8). Duplicates in the mouse genome show the same pattern, consistent with this notion. Thus, both patterns result from ancient evolutionary forces common to eutherian mammals. However, this process appears to be an ongoing characteristic of eutherian X evolution, because 6/28 events have occurred subsequent to the human-mouse split in the human lineage, 6/33 retropositions have occurred within the past ~80 million years in the mouse lineage, and some of these retroduplicate pairs have high sequence similarity (>95%) at synonymous sites. This chromosome-biased gene origination appears to be an important process actively driving the differentiation of the X chromosome in mammals and suggests that this differentiation is still in progress.

A Map of the Interactome Network of the Metazoan *C. elegans*

Siming Li,^{1*} Christopher M. Armstrong,^{1*} Nicolas Bertin,^{1,4} Hui Ge,^{1*} Stuart Milstein,^{1,4} Mike Boxem,^{1*} Pierre-Olivier Vidalain,^{1*} Jing-Dong J. Han,^{1,4} Alban Chesneau,^{1,2*} Tong Hao,¹ Debra S. Goldberg,³ Ning Li,¹ Monica Martinez,¹ Jean-François Rual,^{1,4} Philippe Lamesch,^{1,4} Lai Xu,^{5†} Muneesh Tewari,¹ Sharyl L. Wong,³ Lan V. Zhang,³ Gabriel F. Berriz,³ Laurent Jacotot,^{1‡} Philippe Vaglio,^{1‡} Jérôme Reboul,^{1§} Tomoko Hirozane-Kishikawa,¹ Qianru Li,¹ Harrison W. Gabel,¹ Ahmed Elewa,^{6||} Bridget Baumgartner,⁵ Debra J. Rose,⁶ Haiyuan Yu,⁷ Stephanie Bosak,⁸ Reynaldo Sequerra,⁸ Andrew Fraser,⁹ Susan E. Mango,¹⁰ William M. Saxton,⁶ Susan Strome,⁶ Sander van den Heuvel,¹¹ Fabio Piano,¹² Jean Vandenhaute,⁴ Claude Sardet,² Mark Gerstein,⁷ Lynn Doucette-Stamm,⁸ Kristin C. Gunsalus,¹² J. Wade Harper,^{5†} Michael E. Cusick,¹ Frederick P. Roth,³ David E. Hill,^{1¶} Marc Vidal^{1¶#}

References and Notes

1. B. T. Lahn, N. M. Pearson, K. Jegalian, *Nature Rev. Genet.* **2**, 207 (2001).
2. H. Skalařsky et al., *Nature* **423**, 825 (2003).
3. J. A. Marshall Graves et al., *Cytogenet. Genome Res.* **96**, 161 (2002).
4. J. R. McCarrey, *BioScience* **44**, 20 (1994).
5. M. J. Lercher, A. O. Urrutia, L. D. Hurst, *Mol. Biol. Evol.* **20**, 1113 (2003).
6. P. J. Wang, J. R. McCarrey, F. Yang, D. C. Page, *Nature Genet.* **27**, 422 (2001).
7. J. C. Venter et al., *Science* **291**, 1304 (2001).
8. B. Lewin, *Genes VII* (Oxford University Press, New York, 2000).
9. E. Betran, K. Thornton, M. Long, *Genome Res.* **12**, 1854 (2002).
10. P. M. Harrison et al., *Genome Res.* **12**, 272 (2002).
11. L. Z. Strickman-Almashanu, M. Bustin, D. Lindman, *Genome Res.* **13**, 800 (2003).
12. E. Betran, W. Wang, L. Jin, M. Long, *Mol. Biol. Evol.* **19**, 654 (2002).
13. J. Brosius, *Science* **251**, 753 (1991).
14. E. S. Lander et al., *Nature* **409**, 860 (2001).
15. T. Hubbard et al., *Nucleic Acids Res.* **30**, 38 (2002).
16. Materials and methods are available as supporting material on Science Online.
17. R. H. Waterston et al., *Nature* **420**, 520 (2002).
18. Z. Zhang, P. Harrison, M. Gerstein, *Genome Res.* **12**, 1466 (2002).
19. C.-I. Wu, E. Y. Xu, *Trends Genet.* **19**, 243 (2003).

To further understand biological processes, it is important to consider protein functions in the context of complex molecular networks. The study of such networks requires the availability of proteome-wide protein-protein interaction, or "interactome," maps. The yeast *Saccharomyces cerevisiae* has been used to develop a eukaryotic unicellular interactome map (1-6). *Caenorhabditis elegans* is an ideal model for studying how protein networks relate to multicellularity. Here we investigate its interactome network with HIT-Y2H.

As Y2H baits, we selected a set of 3024 worm predicted proteins that relate directly or indirectly to multicellular functions (7). Gateway-cloned open reading frames (ORFs) were available in the *C. elegans* ORFeome 1.1 (8) for 1978 of these selected proteins. Of these, 81 autoactivated the Y2H *GAL1::HIS3* reporter gene as Gal4 DNA binding domain fusions (DB-X), and 24 others conferred toxicity to yeast cells. The remaining 1873 baits were screened against two different Gal4 activation domain libraries (AD-wrmcDNA and

To initiate studies on how protein-protein interaction (or "interactome") networks relate to multicellular functions, we have mapped a large fraction of the *Caenorhabditis elegans* interactome network. Starting with a subset of metazoan-specific proteins, more than 4000 interactions were identified from high-throughput, yeast two-hybrid (HT-Y2H) screens. Independent coaffinity purification assays experimentally validated the overall quality of this Y2H data set. Together with already described Y2H interactions and interologs predicted *in silico*, the current version of the Worm Interactome (WIS) map contains ~5500 interactions. Topological and biological features of this interactome network, as well as its integration with phenome and transcriptome data sets, lead to numerous biological hypotheses.

REPORTS

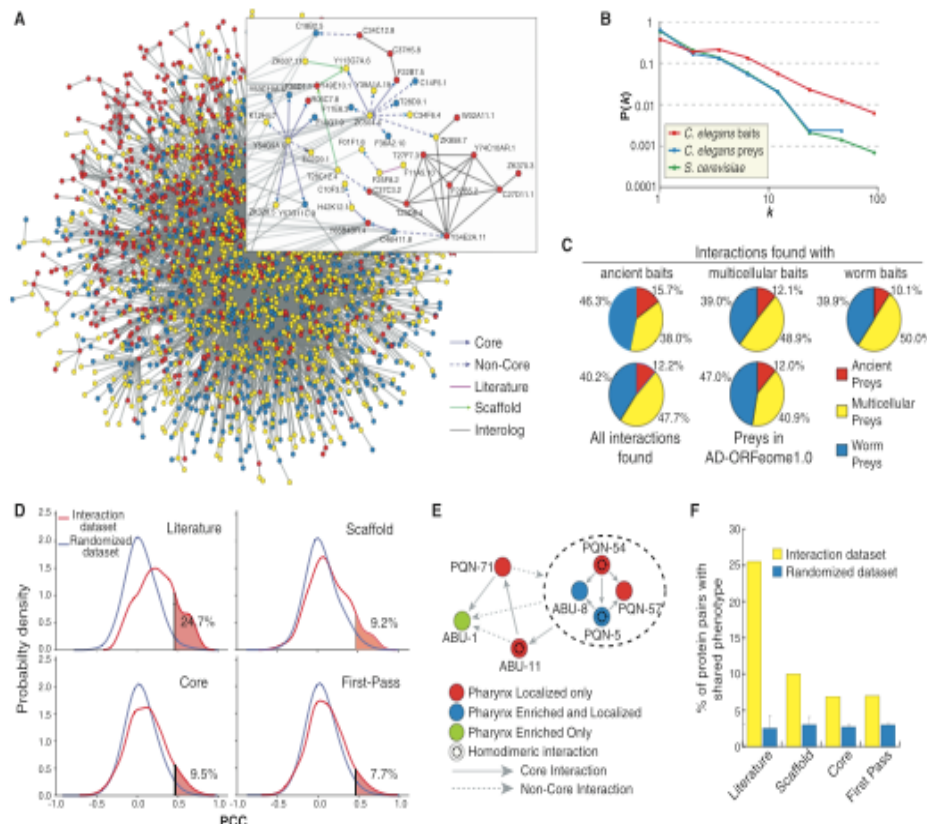


Fig. 2. Analysis of the WIS network. (A) Nodes (representing proteins) are colored according to their phylogenetic class: ancient (red), multicellular (yellow), and worm (blue). Edges represent protein-protein interactions. The inset highlights a small part of the network. (B) The proportion of proteins, $P(k)$, with different numbers of interacting partners, k , is shown for *C. elegans* proteins used as baits or preys and for *S. cerevisiae* proteins. (C) The pie charts show the proportion of interacting preys found in Y2H screens that fall into each phylogenetic class. Also shown is the distribution of all preys found and all preys searched in the AD-ORFeome1.0 library. (D) Overlap with transcriptome (see text) (18), Pearson correlation coefficients (PCCs) were calculated and graphed for each pair of proteins in the interaction data sets and their corresponding randomized data sets. The red area to the right corresponds to interactions that show a significant relationship to expression profiling data ($P < 0.05$). (E) Interactions between proteins in Topomap mountain 29 (18). The dash-circled proteins belong to the same paralogous family (sharing more than 80% homology) and are thus collapsed into one set of interactions. (F) Proportion of interaction pairs where both genes are embryonic lethal ($P < 10^{-2}$).

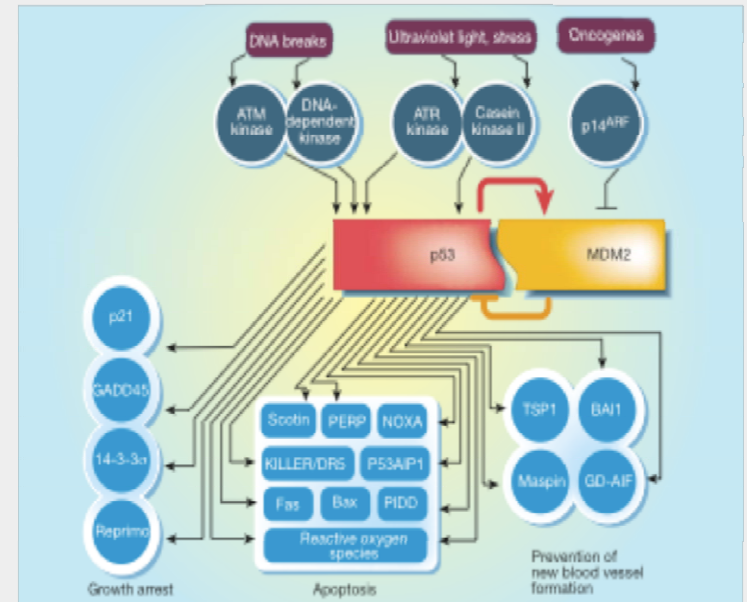
Nature 408 307 (2000)

news and views feature

Surfing the p53 network

Bert Vogelstein, David Lane and Arnold J. Levine

The p53 tumour-suppressor gene integrates numerous signals that control cell life and death. As when a highly connected node in the Internet breaks down, the disruption of p53 has severe consequences.



Molecular Interaction Map of the Mammalian Cell Cycle Control and DNA Repair Systems

Kurt W. Kohn*

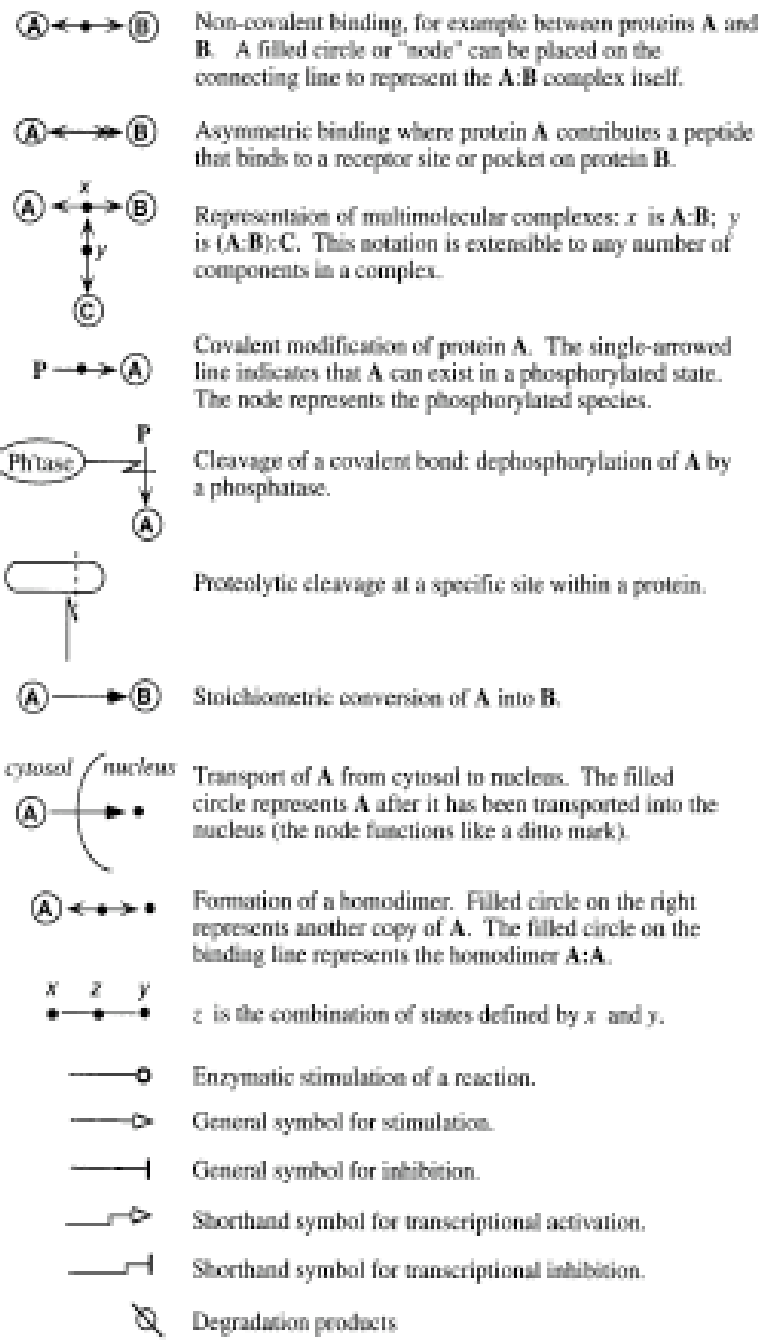
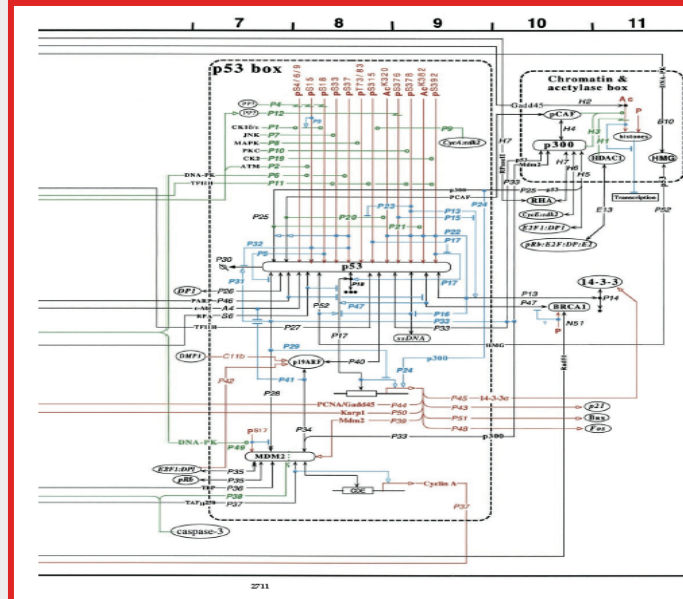
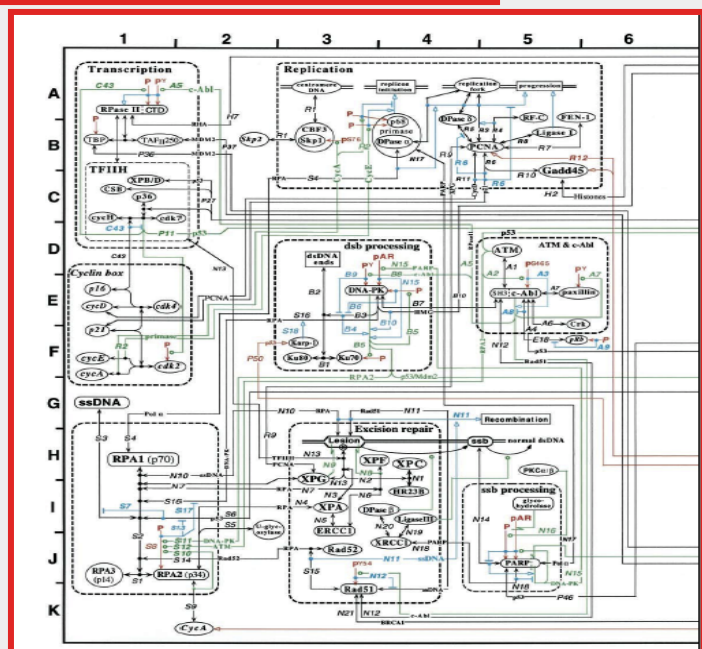
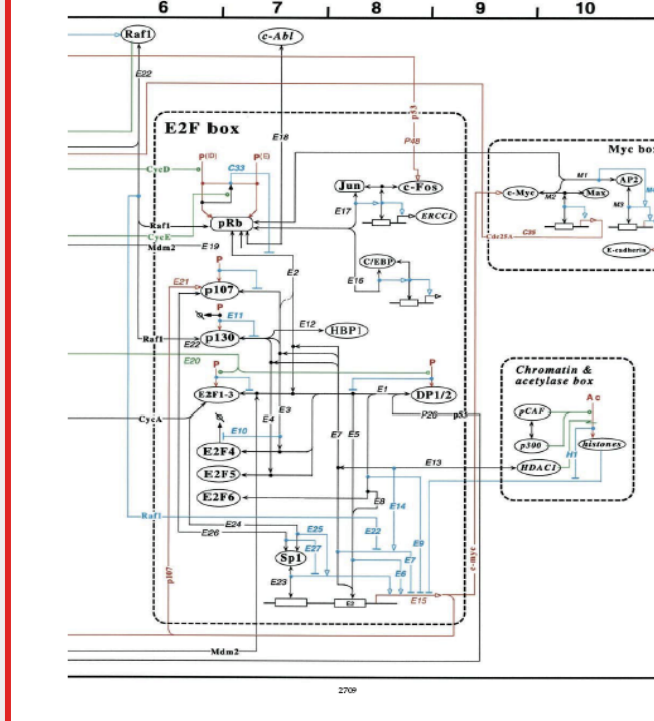
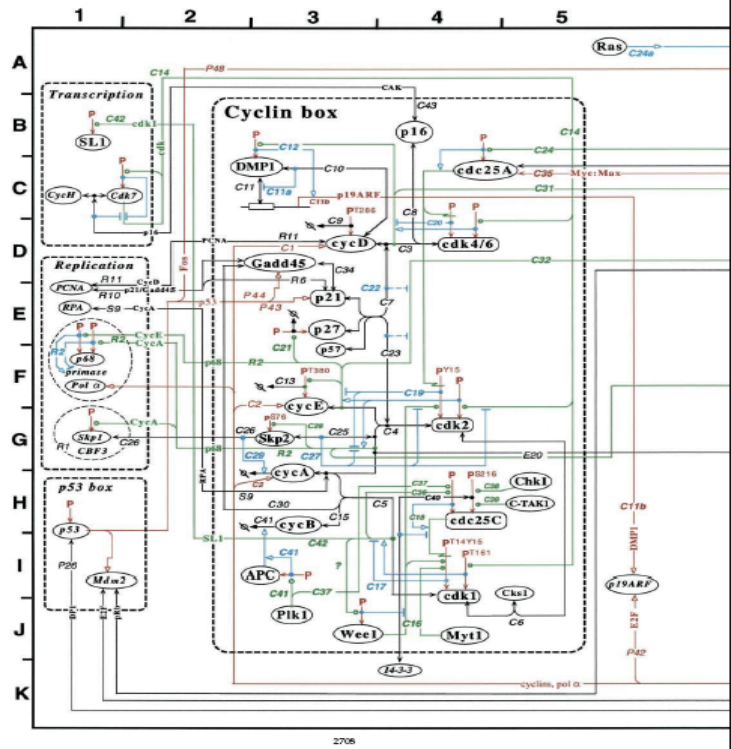


Figure 1. Summary of symbols.



Using process diagrams for the graphical representation of biological networks

Hiroaki Kitano¹⁻⁴, Akira Funahashi^{1,3,4}, Yukiko Matsuoka^{1,3,4} & Kanae Oda^{1,4}

With the increased interest in understanding biological networks, such as protein-protein interaction networks and gene regulatory networks, methods for representing and communicating such networks in both human- and machine-readable form have become increasingly important. Although there has been significant progress in machine-readable representation of networks, as exemplified by the Systems Biology Mark-up Language (SBML) (<http://www.sbml.org>) issues in human-readable representation have been largely ignored. This article discusses human-readable diagrammatic representations and proposes a set of notations that enhances the formality and richness of the information represented. The process diagram is a fully state transition-based diagram that can be translated into machine-readable forms such as SBML in a straightforward way. It is supported by CellDesigner, a diagrammatic network editing software (<http://www.celldesigner.org/>), and has been used to represent a variety of networks of various sizes (from only a few components to several hundred components).

Drawing diagrams with nodes and arrows is the common approach for representing how proteins and genes interact, and papers frequently include such informal node-and-arrow diagrams. Although such diagrams are useful in providing an intuitive idea of how proteins and genes interact, the information contained in such diagrams is not precise because the syntax and semantics of the symbols used tend to be ambiguously defined. Often, arrows adopt multiple different meanings, so that correct interpretation of the diagram depends upon the knowledge of the reader. For example, Figure 1a shows a typical diagram often found in signal transduction papers. In this example, an arrow symbol could be interpreted four different ways: activation, translocation, dissociation of protein complex and residue modification. Correct interpretation of which biological process the arrow refers to depends entirely on the reader's knowledge. In general, such ambiguities and lack of information are not a major problem as long as the diagrams are small and represent genes, proteins and their local interactions. However, problems emerge

when representing interactions within larger networks. Therefore, there is a need for diagrams that contain unambiguous process information in the symbols used and that can be transferred to standard machine-readable codes such as SBML for computational analysis¹.

Circuit schematic diagrams used in electronics are ideal examples of a graphical diagram. Engineers can reproduce the circuits drawn in the schematic diagrams without substantial additional information, because the diagrams are unambiguously defined, contain sufficient information and are based on well-accepted standards.

Kurt Kohn was the first to produce canonical representations for molecular interactions^{2,3}, and other researchers have been working on alternative representations⁴⁻⁸. Unfortunately, none of the proposed schemes has been widely used for a variety of reasons. For example, there is no software tool to create a Kohn Map efficiently, and this type of representation does not explicitly display temporal processes, which makes it difficult for readers to understand the sequence of events. Diagrammatic Cell Language (DCL) modifies Kohn's notation⁹, but suffers from similar problems in that it does not explicitly display a temporal sequence of events and lacks publicly accessible documents and supporting software. Other notations have different shortcomings.

A successful diagram scheme must: (i) allow representation of diverse biological objects and interactions, (ii) be semantically and visually unambiguous, (iii) be able to incorporate notations, (iv) allow software tools to convert a graphically represented model into mathematical formulas for analysis and simulation, (v) have software support to draw the diagrams, and (vi) ensure that the community can freely use the notation scheme.

We have accumulated substantial experience in creating molecular interaction diagrams of various sizes, ranging from several components and interactions to several hundred components and interactions^{10,11}. Whereas associations and combinatorial bindings of molecular species can be compactly described by an entity-relationship diagram (as exemplified by Kohn's diagram), temporal orders of reactions are made implicit so that intuitive understanding of the process of reactions is difficult. The process diagram explicitly represents the temporal order of reactions and states of molecules and complexes at the cost of an increased number of nodes and lines in the diagram. We have previously argued that either approach can be used, depending upon the purpose of the diagram, and both notations can maintain compatible information internally, but differ in visualization⁷. In our experience, however, a process diagram graphically representing state transitions of the molecules involved is more intuitively understandable than an entity-relationship diagram. This article describes in detail how process diagrams can be a vehicle for representing biological networks.

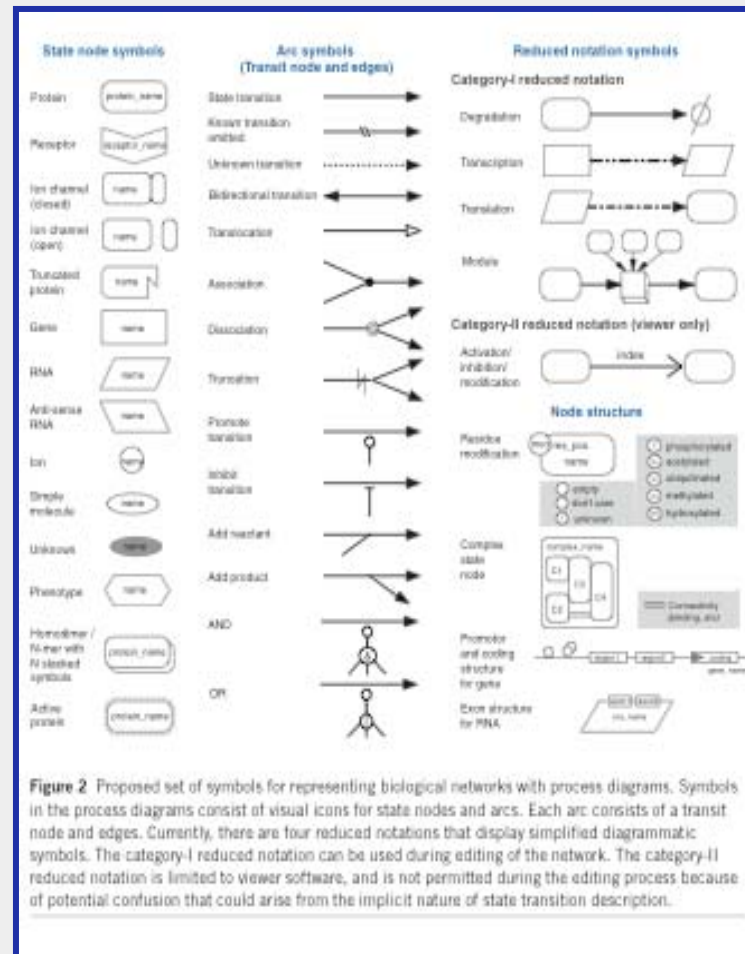


Figure 2 Proposed set of symbols for representing biological networks with process diagrams. Symbols in the process diagrams consist of visual icons for state nodes and arcs. Each arc consists of a transit node and edges. Currently, there are four reduced notations that display simplified diagrammatic symbols. The category-I reduced notation can be used during editing of the network. The category-II reduced notation is limited to viewer software, and is not permitted during the editing process because of potential confusion that could arise from the implicit nature of state transition description.

¹The Systems Biology Institute, Suite 6A, M31 6-31-15 Jingumae, Shibuya, Tokyo, 150-0001 Japan. ²Sony Computer Science Laboratories, Inc., 3-14-13 Higashi-gotanda, Shinagawa, Tokyo, 141-0022 Japan. ³ERATO-SORST Kitano Symbiotic Systems Project, Japan Science and Technology Agency, Suite 6A, M31 6-31-15 Jingumae, Shibuya, Tokyo, 150-0001 Japan. ⁴Department of Fundamental Science and Technology, Keio University, 3-14-1 Hiyoshi, Kohoku, Yokohama 223-8522 Japan. Correspondence should be addressed to H.K. (kitano@symbio.jst.go.jp)

Theorem

Nec. Cond. (the unpacking side)

Given a biological network satisfying the condition C^1, \dots, C^n , it can be rewritten as a conjunction of N Zsyntax theorems

Suff. Cond. (the packing side)

Given a conjunction of N Zsyntax theorems satisfying the condition R^1, \dots, R^m , it can be rewritten as a biological network

**Is ZSYNTAX too simple
to grasp biological complexity?**

- *Informational content and context dependence*
- *Function and context dependance*
- *Post-translational modifications*
- *Allosteric modifications*
- *Different kinds of interactions*
- *Compartmentalisations*
- *Quantitative aspects*
- *Role of time*

Z-conjunction: $\&_i$

Z-interaction: \odot_i

Z-conditional: \rightarrow_i

The shift from plain formulas, such as A, B, C , etc., to *labeled* formulas, such as $A:\mathbf{x}, B:\mathbf{y}, C:\mathbf{z}$. etc. – where $\mathbf{x}, \mathbf{y}, \mathbf{z}$ are labeling strings, consisting of suitable variables, parameters or function symbols – greatly enhances the expressive power of Zsyntax and is well-grounded in contemporary logic.

The basic idea is that of generalizing the inference rules to deal with labeled formulas in which the labels are used to specify any kind of additional information concerning the entities to which the formulas refer.

Then, the purely logical mechanism is integrated with a **labeling algebra** specifying the way in which the values of the parameters should be propagated by each application of the rule.

The labels are strings of suitable variables, parameters or function symbols.

For example, EVF's can be replaced by more precise *empirically valid rules* such as:

$$\frac{A:[x,t_0] \quad B[y,t_0]}{A \odot B:[g(x,y),t]}$$

so that, in this specific case, the corresponding instance of the (logically valid) rule of MP would take the form:

The logical content is expressed by the formulas.

The empirical content is expressed by the labels.

$$\frac{A \rightarrow A \odot B:[y,t_0] \quad A:[x,t_0]}{A \odot B:[g(x,y), t]}$$

Theorem

$TP53[a, t_0] \ \& \ MDM2[b, t_0] \ \& \ U[c, t_0] \ \& \ P[d, t_0] \ \vdash \ d(TP53)[e, t_f]$

Demonstration

1. $MDM2[a, t_0] \ \& \ TP53[b, t_0] \ \rightarrow \ MDM2 \odot TP53[l, t_1]$
2. $(MDM2 \odot TP53) [l, t_1] \ \& \ U[c, t_1] \ \rightarrow \ (MDM2 \odot TP53) \odot U[m, t_2]$
3. $(MDM2 \odot TP53) \odot U[m, t_2] \ \rightarrow \ MDM2[n, t_3] \ \& \ (TP53 \odot U) [o, t_3]$
4. $(TP53 \odot U) [o, t_3] \ \& \ P[p, t_3] \ \rightarrow \ (TP53 \odot U) \odot P[q, t_4]$
5. $(TP53 \odot U) \odot P[q, t_4] \ \rightarrow \ d(TP53) [e, t_f] \ \& \ U[u, t_f] \ \& \ P[v, t_f]$

Zsyntax: A Formal Language for Molecular Biology with Projected Applications in Text Mining and Biological Prediction

Giovanni Boniolo^{1,2}, Marcello D'Agostino³, Pier Paolo Di Fiore^{1,2,4*}

1 IFOM, Istituto FIRC di Oncologia Molecolare, Milano, Italy, 2 Dipartimento di Medicina, Chirurgia ed Odontoiatria, Università di Milano, Milano, Italy, 3 Dipartimento di Scienze Umane, Università di Ferrara, Ferrara, Italy, 4 Istituto Europeo di Oncologia, Milano, Italy

Abstract

We propose a formal language that allows for transposing biological information precisely and rigorously into machine-readable information. This language, which we call Zsyntax (where Z stands for the Greek word ζωή, life), is grounded on a particular type of non-classical logic, and it can be used to write algorithms and computer programs. We present it as a first step towards a comprehensive formal language for molecular biology in which any biological process can be written and analyzed as a sort of logical “deduction”. Moreover, we illustrate the potential value of this language, both in the field of text mining and in that of biological prediction.

The research programme

Future steps:

- 1) elaborating the correct formal framework for Z_{syntax}
- 1) constructing the right algorithm and software to implement our logic in a computer (both for text mining and for prediction)
- 2) showing the foundational relevance of the language
- 3) joining the logical language and the differential equations of molecular biology

