



HyBrow: a prototype system for computer-aided hypothesis evaluation

S. A. Racunas[†], N. H. Shah^{*,†}, I. Albert and N. V. Fedoroff

The Huck Institute for Life Sciences, Penn State University, University park,
PA 16801, USA

Received on ...; revised on ...; accepted on ...

ABSTRACT

Motivation: Experimental design, hypothesis-testing and model-building in the current data-rich environment require the biologists' to collect, evaluate and integrate large amounts of information of many disparate kinds. Developing a unified framework for the representation and conceptual integration of biological data and processes is a major challenge in bioinformatics because of the variety of available data and the different levels of detail at which biological processes can be considered.

Results: We have developed the HyBrow (Hypothesis Browser) system as a prototype bioinformatics tool for designing hypotheses and evaluating them for consistency with existing knowledge. HyBrow consists of a modeling framework with the ability to accommodate diverse biological information sources, an event-based ontology for representing biological processes at different levels of detail, a database to query information in the ontology and programs to perform hypothesis design and evaluation. We demonstrate the HyBrow prototype using the galactose gene network in *Saccharomyces cerevisiae* as our test system, and evaluate alternative hypotheses for consistency with stored information.

Availability: www.hybrow.org

Contact: nigam@psu.edu

1 INTRODUCTION

To expand understanding of a biological system, an experimentalist (1) formulates hypotheses about relationships that exist within that system, (2) gathers information from various repositories about the components of the system, (3) evaluates the hypotheses to assess whether they are supported or contradicted by this information, (4) revises hypotheses as needed, (5) perturbs the system in informative ways and (6) integrates all available information to deepen the understanding of how the system works. Understanding grows as hypotheses are accumulated.

Current investigations of signal transduction and gene regulation generate large volumes of data, making it increasingly difficult to assemble and organize all the information needed to test hypotheses. To complicate matters, the various kinds of data reside in a wide array of repositories and are stored in different formats. To help biologists make effective use of increasing amounts of diverse data, we have developed the HyBrow (Hypothesis Browser) system to aid in the hypothesis formulation and evaluation cycle.

Most bioinformatics tools are designed to perform specific analytical functions. These tools carry out tasks, such as identifying patterns, categorizing information and probing data sources for similarities. However, information synthesis has remained solely the purview of the biologist (Kuchinsky *et al.*, 2002). To design a system that will support the tasks of formulating and evaluating hypotheses for consistency with prior knowledge, we must address several issues. We must specify a representation for hypotheses that is both machine-understandable and accessible to the experimental biologist. We must choose a conceptual model and methods for storing existing information about the biological system in that model. Finally, we must develop a framework that supports the evaluation of hypotheses with respect to the stored information.

Organisms are (and contain) complex systems which are incompletely understood. Such systems are more readily represented by specifying the events that occur in them than by writing differential equations for the system's constituent reactions because sufficient detailed kinetic information is generally not available (Ho, 1989). Biological events are changes in a biological system for which we can obtain experimental evidence. In a previous publication, we described the development of a framework for conceptualizing biological processes in terms of events, formulating hypotheses about them and evaluating the hypotheses (Racunas *et al.*, 2003). In order to represent hypotheses about a biological process in a machine-understandable format, it is necessary to create a vocabulary of objects (agents) and processes, and define the relationships in which these entities can participate. We refer to this vocabulary as the hypothesis ontology and in our current work we construct and populate such an ontology

*To whom correspondence should be addressed.

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

for a simple test system. We describe biological events by naming the agents from the ontology (such as proteins and nucleic acids) and the processes (such as 'binds') that connect them. We use the term hypothesis event to represent an abstract biological event. Thus, an hypothesis event consists of an acting agent (a 'subject', such as a protein), a relationship (a 'verb', such as induce, repress, . . .), a target agent (an 'object', a gene, protein, . . .), the experimental and cellular contexts in which the event takes place; and a set of associated conditions (such as the presence or absence of other agents), which can accompany the event. This event-based framework, together with our hypothesis ontology, allows us to represent hypotheses in a formal language that specifies the time and context-dependent relationships among the system's objects and processes (Sudkamp, 1988; Racunas *et al.*, 2003).

HyBrow's event-based framework includes methods to evaluate such formal language hypotheses for internal consistency and agreement with existing knowledge (Racunas *et al.*, 2003). Consistency of an hypothesis with observed data and prior knowledge is evaluated by applying constraints and rules. Constraints specify classes of forbidden events. Rules are the operations performed upon available information in order to enforce the constraints. Rules generate judgments of support or conflict, depending upon whether or not an assertion is supported by existing knowledge. The framework also includes neighborhood functions to establish similarity between hypotheses. These facilitate hypothesis revision through the automatic generation of 'neighboring' hypotheses that are variants of an original hypothesis. Neighborhood functions use biologically acceptable notions to generate sets of variant events for events that conflict with existing data or prior knowledge. We examine these variants to find more fitting events and replace conflicted events with superior variants to produce hypotheses that better fit the stored information.

HyBrow's event-based framework makes it possible for the biologists' to deal strictly with experimental evidence and to avoid the unintended assertions that are common artifacts of statistical and equation-based approaches. In HyBrow's framework, hypotheses and evaluation methods are directly compatible with the way information is conceptualized by biologists, making it easier to tap the expertise of experienced biologists. Finally, and most importantly, HyBrow's framework makes it possible to bring together many kinds of data and information in a unified formal language. The inability to combine information sources has been a stumbling block for computational models of biological systems, leading current information integration efforts to focus on only one or two categories of information (Hartemink *et al.*, 2001; Segal *et al.*, 2003).

In this paper, we demonstrate HyBrow's information synthesis capability using the galactose metabolic and regulatory network, which we chose because abundant data and information of many different types are publicly available for this system (Ideker *et al.*, 2001). We designed a small

hypothesis ontology appropriate for the GAL system. We specified the formal grammar that describes how to combine terms from the ontology into hypotheses. We designed a database to store yeast GAL data structured in the ontology and developed hypothesis composition, visualization and evaluation software.

To test the HyBrow prototype, we evaluated and ranked hypotheses about the GAL system. During evaluation, HyBrow assayed all stored data for conflicting or supporting evidence for each statement in each hypothesis. HyBrow modified hypotheses that contained errors to generate variants with fewer flaws. Finally, HyBrow combined the resulting determinations of conflict and support to generate evaluations and rankings for all the original and variant hypotheses.

2 RESULTS

2.1 Hypothesis ontology

Common ontologies for biological objects and processes (Schulze-Kremer, 1998; Ashburner *et al.*, 2000) are being developed to support the intercommunication of diverse databases as well as enable automated annotation and extraction of information from the literature (Fleischmann *et al.*, 1999; Novichkova *et al.*, 2003). Ontologies also provide a foundation for the construction of higher level models of biological systems (Rzhetsky *et al.*, 2000; Peleg *et al.*, 2002). Models vary from abstract Boolean (Akutsu *et al.*, 2000) and Bayesian networks (Hartemink *et al.*, 2001) to highly specific (McAdams and Arkin, 1998) and quantitative models (Sveiczzer *et al.*, 2000). Currently, most databases do not store information in an explicit ontology that facilitates modeling, and groups that design ontologies (Rzhetsky *et al.*, 2000; Peleg *et al.*, 2002) do not store all relevant data structured in these ontologies. Hence, efforts aimed at integrating diverse information sources need to choose or design a representation scheme and convert existing data into that representation. Although specialized ontologies exist (Karp, 2000; Rzhetsky *et al.*, 2000), there is a need for an ontology that allows users to represent biological processes in an event-based manner. Such an ontology should be compatible with existing ones so that hierarchical relationships can be made between terms in existing ontologies and the 'hypothesis ontology' (Karp, 2000; Rzhetsky *et al.*, 2000).

We used Protégé (Crubézy *et al.*, 2003) to design a small hypothesis ontology for representing GAL system information in HyBrow's event-based conceptual framework (Racunas *et al.*, 2003). For guidance, we relied upon the principles used to design the Rzhetsky and the Bioprocess ontologies (Rzhetsky *et al.*, 2000; Peleg *et al.*, 2002). Our ontology (Fig. 1) accommodates currently available literature data, extracted primarily from Yeast Proteome Database (YPD) (Proteome, 2001, <http://www.proteome.com/YPDhome.html>) at a coarse level of resolution. An event consists of an acting agent (the 'subject', such as gene, RNA,

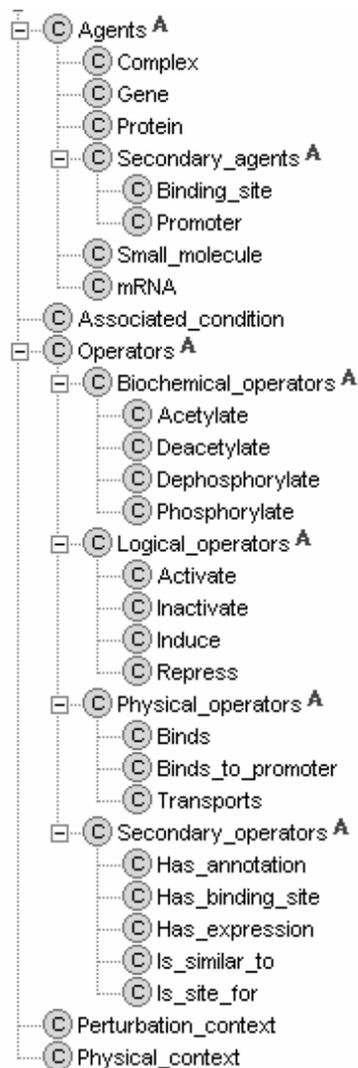


Fig. 1. An overview of the ontology used to represent data in an event-centered way. ‘Operators’ are the relationships that can exist between agents.

protein), a target agent (the ‘object’, such as a gene, protein, complex), a relationship (the ‘verb’, such as induce, repress, bind), a context in which the event takes place and an optional set of associated conditions (such as the presence or absence of other agents) that accompany the event. The construction of events from elements of the ontology, event sets from events and hypotheses from event sets is governed by a context-free grammar. Events that occur in the same context are combined to form event sets and an hypothesis consists of event sets linked by logical and temporal operators. An hypothesis must contain at least one event set and each event set must contain at least one event. Please refer to www.hybrow.org for a formal specification of this grammar.

Contexts specify where events occur in the cell and under what genetic conditions they occur. Our contexts are derived

from established ontologies. For example, terms for specifying physical locations in the cell come from the cellular component division of the Gene Ontology (GO). We currently support genes, proteins, mRNA, small molecules, and complexes of proteins, small molecules and mRNA as agents in our prototype. We define three main categories of relationships: logical (e.g. induce), biochemical (e.g. phosphorylate) and physical (e.g. bind). The key design principle is that the ontology describes a regulatory system in an event-based way consistent with our evaluation framework. Our current hypothesis ontology allows representation of events such as ‘Gal4p binds to the promoter of the *gal1* gene in the presence of galactose in wild type *Saccharomyces cerevisiae*’. Depending on the resolution of the ontology, this approach can represent anything from simple protein phosphorylation to the entire cell cycle (GKB, 2003, <http://www.genomeknowledge.org/>). Formal presentation of the complete ontology is available at www.hybrow.org.

2.2 Inference rules and constraints

We have defined constraints and the rules that determine whether or not a constraint is satisfied for each relationship expressible with terms from our ontology. We define several categories of constraints. Ontology constraints determine what agents can participate in which types of biological relationships. For example, a gene cannot transport a gene, but a protein can transport a small molecule or another protein. Data constraints determine what data values are valid for a particular relationship. For example, for the relationship ‘protein A binds to the promoter of gene B’, it is acceptable for protein A to be annotated as localized in the nucleus or cytoplasm but not on the cell membrane. Existence constraints require an agent’s presence before it can enter a relationship. For example, a protein cannot perform its function when its gene has been deleted. Temporal constraints govern the transmission of modifications made to an agent by previous events. For example, event ‘X phosphorylates Y’ implies that in all subsequent events Y is phosphorylated (unless a dephosphorylation event occurs).

Each rule has sections that correspond to the different constraints that exist in HyBrow. The first section deals with ontology constraints, the second, with constraints on annotation data in GO (Ashburner *et al.*, 2000) format, the third deals with literature-extracted information structured in the ontology and the fourth with constraints on the specific data type(s) for a relationship, such as promoter sequence in the case of the binds to promoter relationship. For each constraint that is violated in any section, the event is assigned a ‘conflict’. For each constraint that is satisfied, the event is assigned a ‘support’. If a constraint is neither violated nor supported, a ‘cannot comment’ is assigned. Sections 1–3 can be generalized because they have a common structure for different relationships, and the operations to be performed on the data are very similar. Section 4 is very specific because of the different ways in

which different data types for each relationship must be used. There are additional general sections that enforce existence and temporal constraints. For example, the rule for protein A binds to promoter of gene B has the following sections: (1) check if A is a protein or a protein-complex and if B is a gene; (2) check whether protein A is annotated (a) to have the molecular function of a transcriptional activator or repressor, (b) to be involved in the biological process of transcriptional regulation and (c) to have a nuclear localization; (3) determine whether the literature reports the postulated event; (4) search the promoter of gene B for a binding site for protein A; (5) ensure that the event is not postulated in a genetic context where the gene for protein A is knocked out.

Rules are coded in Perl as hierarchical function libraries to keep the rule set extensible and flexible. Most of the constraints enforced by the generalized sections are stored in database tables, which are queried at run time, allowing flexibility for changing the stringency of the constraints. A more detailed description of the rule library is provided at www.hybrow.org.

2.3 Database and information gathering

At the heart of HyBrow is the idea that disparate kinds of information can be represented in a unified formal language. Biological information residing in the published literature and electronic databases is expanding at an accelerating rate. Retrieving information and translating it into our ontology presents several problems because the information is in different repositories and storage formats. Further, only a fraction of the published literature is available electronically. The problem of automating extraction of information from the literature is being addressed by a number of research groups, but is far from solved (Andrade and Bork, 2000; Rzhetsky *et al.*, 2000; IBM, 2002). The most promising approach appears to be MedScan (Novichkova *et al.*, 2003), which can parse literature abstracts to identify ‘biological events’. But information extraction is still largely manual, practiced by annotators who read papers for relevant concepts and information.

In this work, we adopted different approaches to gather and structure data in our ontology. For data with standardized representation formats, we designed user agents to access the existing public repositories and retrieve desired information. For example, we designed a user-agent to retrieve promoter sequences from the *S.cerevisiae* Promoter Database (Zhu and Zhang, 1999). In most cases, we were able to access well-annotated information from the *Saccharomyces* Genome Database (SGD, <http://genome-www.stanford.edu/Saccharomyces/>) (Cherry *et al.*, 1998) directly. We used YPD (Proteome, 2001) to access curated literature information about *S.cerevisiae* genes and proteins. We designed a form-based layout for gathering biological information from YPD reports and filled in predefined table fields compatible with the ontology from specific fields of the YPD report. This process is easily automated if direct access

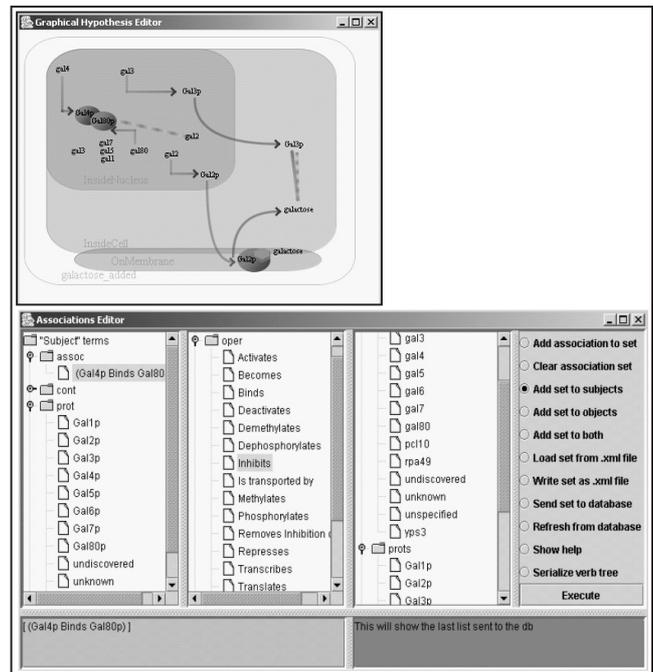


Fig. 2. Screen shots of the visual and widget interface for constructing hypotheses.

to database tables is obtained and can be extended to frame-based ‘loading forms’ like the EcoCyc database (Karp and Riley, 1999). For quantitative data, such as that from microarray expression profiling experiments, we converted the data into our table format using custom Perl scripts. If microarray data are structured in the MAGE object model, this task is more straightforward (Spellman *et al.*, 2002).

We designed a MySQL database and mapped our ontology onto the database for easy extension as our ontology evolved. We created a table in the database for each class in the ontology, at the finest level of resolution. The table has fields for properties, called ‘slots’ and ‘facets’ in the ontology, of the relevant class. This creates more tables than would be present in a well-normalized relational database. However, a prototyping effort requires the backend to be easily modified in response to changes in the ontology. The backend also contains tables to store constraints used during evaluation. A detailed database schema is found at www.hybrow.org.

2.4 User interfaces

An important feature of HyBrow is that it is easy for the user to construct a machine-readable formal language hypothesis. We have created two interfaces for this purpose, a visual interface (Fig. 2, upper panel) and a widget interface (Fig. 2, lower panel). Our visual interface allows users to construct hypotheses using a visual notation designed in accordance with the proposed conventions (Cook *et al.*, 2001). This interface allows users to draw diagrams that are

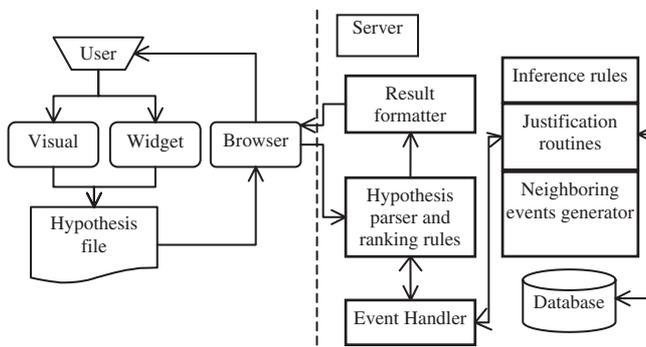


Fig. 3. The visual or the widget interface is used to design hypotheses, which are sent to the server via a browser. The hypothesis parser is the entry point for the system and it uses the event handler, which manages the event library, ranking, justification and event neighborhood generation. The database stores the different data structured into ‘events’.

then automatically translated into hypotheses. The widget interface allows the user to write hypotheses in English-like ‘sentences’ constructed using subject/verb/object selection menus. A user can construct portions of an hypothesis using different interfaces and then combine them. Details on how to use the tools are provided at www.hybrow.org. Hypotheses are saved to local files and then submitted for evaluation via the Web.

2.5 The hypothesis evaluation process

The hypothesis evaluation process is illustrated diagrammatically in Figure 3. When HyBrow receives an hypothesis, it checks the connections between events and event sets for conformity with the hypothesis grammar. If the hypothesis passes these tests for syntax, each event is then checked for validity using the appropriate rule for the relationship proposed in the event. For each event, a support, conflict or cannot comment result corresponding to each of the four sections of the inference rules is returned. Finally, the support and conflict calls are tallied based upon the logical structure of the hypothesis. Each ‘and’ between event sets leads to the inclusion of results from both sets in the final tally. For each ‘or’ connection, the ‘better’ set is chosen using a hierarchical set of rules. [Sample rules: (1) an event set with conflicts it is better than an event set with more conflicts and worse than one with fewer conflicts. (2) An event set for which all events have at least some support is better than an event set for which at least one event is not supported. (3) If one event set’s support is a strict superset of another event set’s support, the superset is superior...]. We apply these rules sequentially until one of the rules returns a clear decision.

For each event, the hypothesis evaluation process finds all conflicts with existing knowledge and indexes them, along with their sources. These are reported to the user to allow them to identify specific problems with the hypothesis and

the conflicting data source. For each event that has a conflict, a set of variant events is generated using biologically motivated heuristics, such as replacing the acting agent with agents that share a sequence similarity or share a similar cellular localization and sequence similarity with the original agent. Neighboring hypotheses that share the logical structure of the original are generated by replacing conflicting events with the best variant events. These neighboring hypotheses are then evaluated, and if a better (more supported, less conflicted) hypothesis is found, it is presented to the user.

After evaluation, the user is shown (1) the support and conflict totals; (2) the least conflicted, most supported event sets that fit the logical structure of the hypothesis; (3) a support-conflict scatter plot of neighboring hypotheses automatically generated from the user submitted hypothesis; and (4) a list of all events that had conflicts, the data that triggered the conflicts, an explanation of why the rules interpret that data as a conflict, and a reference to the original article or data source. The results pages (Fig. 4) allow a user to gauge the ‘fitness’ of his/her hypothesis in the light of all stored data. Iterative refinement of the hypothesis allows the user to reconcile all stored data into a single coherent representation whose level of detail depends on the resolution of the ontology used for constructing hypotheses.

2.6 Test runs with sample hypotheses

In order to test the prototype, we comprised hypotheses about the GAL system and ranked them. The GAL system consists of genes that transport and metabolize galactose and the regulatory network that controls whether the pathway is on or off (Lohr *et al.*, 1995). The process involves three types of proteins as follows: (a) a permease (Gal2p) that transports galactose (encoded by the *gal2* gene); (b) proteins that utilize intracellular galactose, galactokinase (encoded by *gal1*), uridylyltransferase (encoded by *gal7*), epimerase (encoded by *gal10*) and phosphoglucosyltransferase (encoded by *gal5*); and (c) the regulatory proteins Gal3p, Gal4p and Gal80p, which exert transcriptional control over the genes encoding the transporter, the enzymes and to some extent, their own genes (Ideker *et al.*, 2001). HyBrow successfully identified the hypothesis that best explained the current understanding about GAL system regulation (Lohr *et al.*, 1995; Ideker *et al.*, 2001). For six of the seven events that had conflicts, HyBrow was also able to suggest corrections successfully that increased agreement with stored information. All hypotheses used and explanations of their evaluations can be found at www.hybrow.org.

Here, we describe the evaluation of a simple illustrative hypothesis as follows: ‘Gal2p transports galactose into the cell at the cell membrane. In the cytoplasm, galactose activates Gal3p. Gal3p binds to the promoter of the *gal1* gene and induces its transcription in the presence of galactose’. This hypothesis was decomposed into events as shown in Figure 4A. On evaluation, HyBrow reported support from

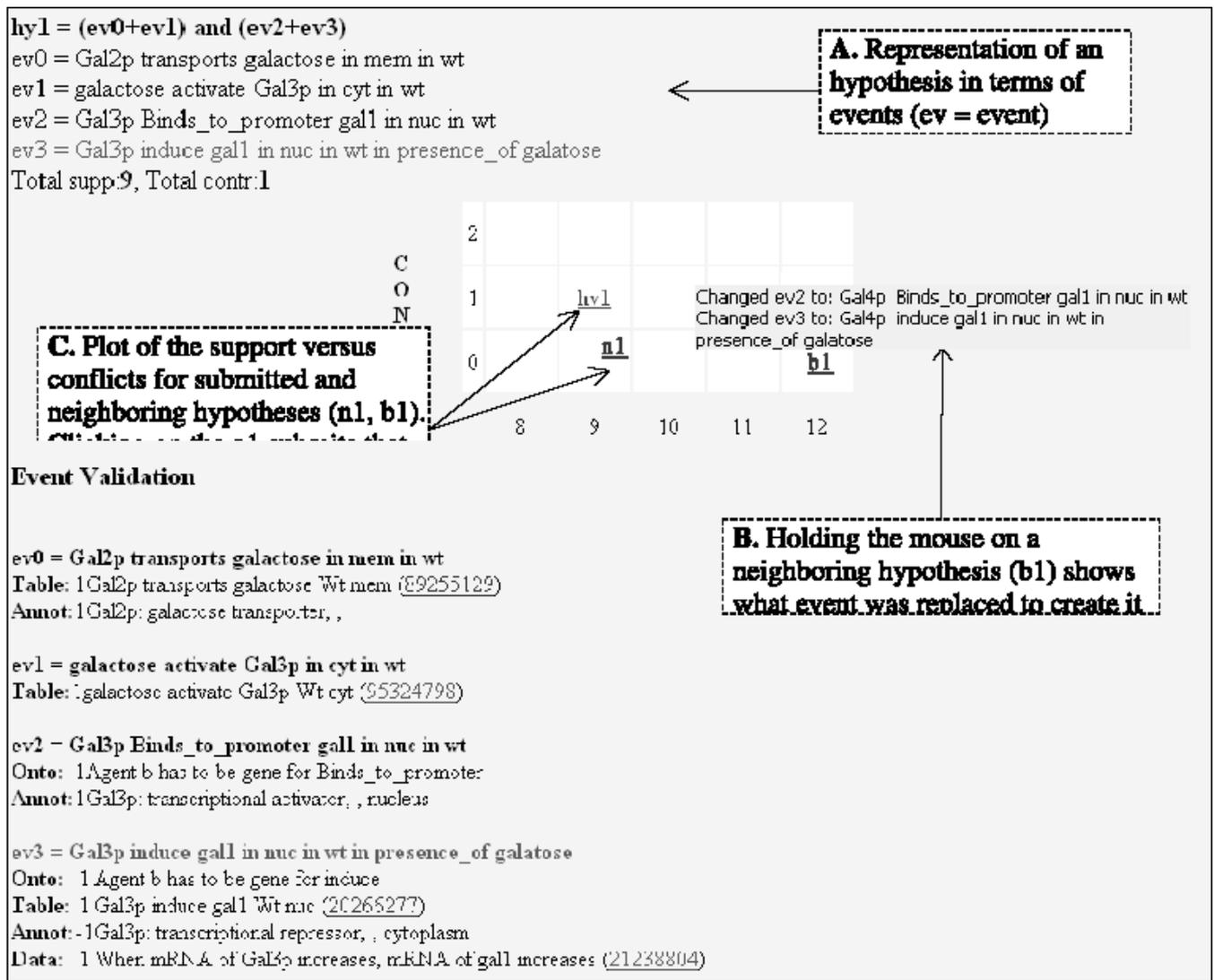


Fig. 4. Screen shot of the results page: see text for a more detailed description.

literature and GO annotation for event number 0 (ev0), support from literature for ev1, support from ontology constraints and annotation for ev2 and support from the ontology, literature and data sections for ev3. It reported a conflict for ev3 (which is marked in red) from the annotation rule section because Gal3p is annotated to be primarily in the cytoplasm (SGD, 2003). HyBrow then searched for variant events. For ev3, it found an event (Gal4p binds to promoter of gal1) with higher support and for ev4 it found the more meaningful event [Gal4p induces gal1 in nucleus in wild-type (wt) in presence of galactose] with the same support but no conflict. These events were inserted in place of the original events to create a neighboring hypothesis that is better than the original hypothesis (Fig. 4B and C).

When a submitted event contains a perturbation, such as the deletion of a gene, HyBrow identifies the agents disabled

because of the perturbation and infers a conflict with events that depend on these agents. For example, if the submitted event is ‘Gal3p induce gal1 in nucleus in gal3-K/O’ HyBrow reports a conflict. (The event ‘Gal3p not induces gal1 in nucleus in gal3-K/O’ gets support.) Some of the inferences suggested by HyBrow are obvious for the small GAL system, but HyBrow’s ability to automate the process offers a substantial advantage for systems containing large numbers of genes and proteins.

3 DISCUSSION AND FUTURE WORK

HyBrow supports the construction, ‘proofreading’ and evaluation of hypotheses expressed in familiar diagram or intuitive text-based formats to aid in synthesizing data into working models. HyBrow’s methodology is evaluation-based. Thus,

Please check the insertion of wild-type in the sentence is ok.

unlike systems that construct statistical or equation-based models, HyBrow is able to provide explicit reasons (and references) for its output. However, HyBrow neither force the user to accept all the stored data nor judge the validity of stored information. Rather, it gives the user links to the exact source of each conflict, leaving it upto the user to judge the relative merits of information sources. The user can choose to ignore conflicts from data sources deemed unreliable.

HyBrow differs fundamentally from existing efforts such as EcoCyc, modeling biological processes as ‘workflows’ and Genome Knowledgebase (GKB). EcoCyc is designed using an explicit ontology for biological function and facilitates functional querying. However, it lacks the notion of an hypothesis or a formal framework to evaluate and rank alternative statements about a biological process (Karp *et al.*, 2002). Modeling of biological processes as ‘workflows’ by Altman’s group includes some of HyBrow’s features, but the underlying conceptual model (which uses hybrid Petri nets) does not support hypothesis neighborhoods and the analysis of models has to be done manually (Peleg *et al.*, 2002). GKB (2003) is an effort to structure biological knowledge in an event-centered data model. It is not a modeling framework by itself, but serves as a public source of structured data which efforts like ours can use.

In our test runs, HyBrow identified the least conflicted hypothesis accurately and suggested valid ‘corrections’ for events with conflicts. We believe that we can build upon this success and plan to extend and strengthen HyBrow in several ways. Currently, improvements to hypotheses are suggested using neighboring events generated using simple heuristics, while our conceptual framework supports neighborhood functions that create similar event sets from a given event set (Racunas *et al.*, 2003). Extending HyBrow to use neighborhoods of event sets as well as of events requires new evaluation routines to track all the biochemical and other modifications that an event set generates and to ensure that the neighboring event sets satisfy them. In future work, we will explore biological notions of similarity between event sets and modify our neighborhood functions and evaluation routines accordingly. The current rule library contains rules for ‘extrapolations’ in the presence of perturbations such as gene knock-outs and constitutive over-expression. In future, we would be able to include extrapolations for more categories of perturbations. Our current implementation can only propagate temporal constraints about the presence or absence of biochemical modifications. We intend to propagate constraints about activation/inhibition and induction/repression in an attempt to model how an event affects downstream agents. To this end, we will extend the current ontology to include ‘modification state’ and ‘activation state’ descriptors for agents; events will then be able to modify these state descriptors. During the ontology extension, we will also define more operators (such as ubiquitinate and methylate) and introduce temporal operators, a new category of relationship that

allows for stating relations such as ‘precedes’, ‘after’ and ‘until’ between events.

As our ontology becomes more complex and refined, reorganization and reloading of the underlying database is unavoidable. The manual gathering and loading of information can become a major bottleneck and hence we need automatic ways to access and structure information. This is possible when the exchange format for the data type under question is standardized or special semantic ontologies like TAMBIS are used (Stevens *et al.*, 2000). Three alternative approaches have been proposed in the community to address the task of integrating frequently changing databases independent of the data source’s storage format. The first is the BioMoby project, whose objective is to develop a web service infrastructure for biological data sources (Wilkinson and Links, 2002). The second approach is to design database ‘transformers’ or ‘wrappers’ and a set of ‘drivers’ for extracting a particular data type from a database for ‘view integration’ or constructing a data warehouse. The third alternative uses Grid architecture to avoid local storage and provide distributed access to biological data. However, the myGrid system is at a conceptual state at present and is itself under development (Stevens *et al.*, 2003). Solutions such as K2, a view integration system, and federated information integration schemes such as Genomics Unified Schema are needed to automate data access to maintain and periodically update all the information structured in the ontology (Davidson *et al.*, 2001). Finally, HyBrow will be extended to identify events that are frequently specified, but for which evaluation was not possible. Identifying such events will allow HyBrow to aid experiment design. For instance, if many users include an event in their hypotheses and there is no experimental evidence for it, HyBrow can indicate a need to obtain such data.

4 CONCLUSIONS AND SUMMARY

Our implementation of HyBrow for the GAL system demonstrates that ontology driven, event-based modeling of biological processes is feasible and that structuring data in HyBrow’s event-based framework facilitates computer-aided hypothesis evaluation. HyBrow can accommodate both more data and types of data as they become available. Moreover, its constraints can be elaborated as understanding about the biological system grows. We believe that the approach we have developed for this HyBrow prototype can significantly inform experimentation by integrating large amounts of information for the evaluation of hypotheses.

REFERENCES

- Akutsu, T., Miyano, S. and Kuhara, S. (2000) Algorithms for identifying Boolean networks and related biological networks based on matrix multiplication and fingerprint function. *J. Comput. Biol.*, **7**, 331–343.

- Andrade, M.A. and Bork, P. (2000) Automated extraction of information in molecular biology. *FEBS Lett.*, **476**, 12–17.
- Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K., Dwight, S.S., Eppig, J.T. et al. (2000) Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.*, **25**, 25–29.
- Cherry, J.M., Adler, C., Ball, C., Chervitz, S.A., Dwight, S.S., Hester, E.T., Jia, Y., Juvik, G., Roe, T., Schroeder, M., Weng, S. and Botstein, D. (1998) SGD: Saccharomyces Genome Database. *Nucleic Acids Res.*, **26**, 73–79.
- Cook, D.L., Farley, J.F. and Tapscott, S.J. (2001) A basis for a visual language for describing, archiving and analyzing functional models of complex biological systems. *Genome Biol.*, **2**, RESEARCH0012.
- Crubézy, M., Ferguson, R., Knublauch, H., Musen, M., Noy, N., Tu, S. and Vendetti, J. (2003) *Protege 2000*. Stanford University, Palo Alto, CA.
- Davidson, S.B., Crabtree, J., Brunk, B.P., Schug, J., Tannen, V., Overton, G.C. and Stoekert, C. (2001) K2/Kleisli and GUS: experiments in integrated access to genomic data sources. *IBM Syst. J.*, **40**, 512–531.
- Fleischmann, W., Moller, S., Gateau, A. and Apweiler, R. (1999) A novel method for automatic functional annotation of proteins. *Bioinformatics*, **15**, 228–233.
- GKB (2003) Genome Knowledge Base.
- Hartemink, A.J., Gifford, D.K., Jaakkola, T.S. and Young, R.A. (2001) Using graphical models and genomic expression data to statistically validate models of genetic regulatory networks. *Pac. Symp. Biocomput.*, 422–433.
- Ho, Y.C. (1989) Special issue on discrete event dynamical systems: editorial. *Proc. IEEE*, **77**, 24–38.
- IBM (2002) Discovery Link, IBM Corporation.
- Ideker, T., Thorsson, V., Ranish, J.A., Christmas, R., Buhler, J., Eng, J.K., Bumgarner, R., Goodlett, D.R., Aebersold, R. and Hood, L. (2001) Integrated genomic and proteomic analyses of a systematically perturbed metabolic network. *Science*, **292**, 929–934.
- Karp, P.D. (2000) An ontology for biological function based on molecular interactions. *Bioinformatics*, **16**, 269–285.
- Karp, P.D. and Riley, M. (1999) EcoCyc the resource and the lessons learned. In Letovsky, S. (ed.), *Bioinformatics Databases and Systems*. Kluwer Academic Publishers, New York, pp. 47–62.
- Karp, P.D., Riley, M., Saier, M., Paulsen, I.T., Collado-Vides, J., Paley, S.M., Pellegrini-Toole, A., Bonavides, C. and Gama-Castro, S. (2002) The EcoCyc Database. *Nucleic Acids Res.*, **30**, 56–58.
- Kuchinsky, A., Graham, K., Moh, D., Creech, M., Babaria, K. and Adler, A. (2002) Biological Storytelling: a software tool for biological information organization based upon narrative structure. *6th International Working Conference on Advanced Visual Interfaces (AVI'02)*, Trento, Italy, 22–24 May.
- Lohr, D., Venkov, P. and Zlatanova, J. (1995) Transcriptional regulation in the yeast GAL gene family: a complex genetic network. *FASEB J.*, **9**, 777–787.
- McAdams, H.H. and Arkin, A. (1998) Simulation of prokaryotic genetic circuits. *Annu. Rev. Biophys. Biomol. Struct.*, **27**, 199–224.
- Novichkova, S., Egorov, S. and Daraselia, N. (2003) MedScan, a natural language processing engine for MEDLINE abstracts. *Bioinformatics*, **19**, 1699–1706.
- Peleg, M., Yeh, I. and Altman, R.B. (2002) Modelling biological processes using workflow and Petri Net models. *Bioinformatics*, **18**, 825–837.
- Proteome (2001) Yeast Proteome Database.
- Racunas, S.A., Shah, N. and Fedoroff, N.V. (2003) A contradiction-based framework for testing gene regulation hypotheses. *IEEE Computer Society CSB Conference*, Stanford University, Palo Alto, CA. IEEE Computer Society.
- Rzhetsky, A., Koike, T., Kalachikov, S., Gomez, S.M., Krauthammer, M., Kaplan, S.H., Kra, P., Russo, J.J. and Friedman, C. (2000) A knowledge model for analysis and simulation of regulatory networks. *Bioinformatics*, **16**, 1120–1128.
- Schulze-Kremer, S. (1998) Ontologies for molecular biology. *Pac. Symp. Biocomput.*, 695–706.
- Segal, E., Wang, H. and Koller, D. (2003) Discovering molecular pathways from protein interaction and gene expression data. *Bioinformatics*, **19**, I264–I272.
- SGD (2003) Saccharomyces Genome Database.
- Spellman, P.T., Miller, M., Stewart, J., Troup, C., Sarkans, U., Chervitz, S., Bernhart, D., Sherlock, G., Ball, C., Lepage, M. et al. (2002) Design and implementation of microarray gene expression markup language (MAGE-ML). *Genome Biol.*, **3**, RESEARCH0046.
- Stevens, R.D., Robinson, A.J. and Goble, C.A. (2003) myGrid: personalised bioinformatics on the information grid. *Bioinformatics*, **19**, I302–I304.
- Stevens, R., Baker, P., Bechhofer, S., Ng, G., Jacoby, A., Paton, N.W., Goble, C.A. and Brass, A. (2000) TAMBIS: transparent access to multiple bioinformatics information sources. *Bioinformatics*, **16**, 184–185.
- Sudkamp, T.A. (1988) *Languages and Machines*. Addison-Wesley, Reading.
- Sveiczner, A., Csikasz-Nagy, A., Gyorfy, B., Tyson, J.J. and Novak, B. (2000) Modeling the fission yeast cell cycle: quantized cycle times in *wee1-cdc25Delta* mutant cells. *Proc. Natl Acad. Sci., USA*, **97**, 7865–7870.
- Wilkinson, M.D. and Links, M. (2002) BioMOBY: an open source biological web services proposal. *Brief Bioinform.*, **3**, 331–341.
- Zhu, J. and Zhang, M.Q. (1999) SCPD: a promoter database of the yeast *Saccharomyces cerevisiae*. *Bioinformatics*, **15**, 607–611.



Please provide publisher details if possible and page range