

CHAPTER 17

TECHNOLOGIES FOR BIOLOGICAL DATA INTEGRATION

Limsoon Wong

Institute for Infocomm Research
limsoon@i2r.a-star.edu.sg

The process of building a new database relevant to some field of study in biomedicine involves transforming, integrating, and cleansing multiple data sources, as well as adding new material and annotations. We review in this chapter some of the requirements and technologies relevant to this data integration problem.

ORGANIZATION.

Section 1. We begin with a detailed account of the motivations and requirements for a general integration system for biological data. We also discuss the motivations and requirements for locally warehousing such data.

Section 2. Then we review some representative technologies for data integration in biological and medical research. The technologies surveyed include EnsEMBL, GenoMax, SRS, DiscoveryLink, OPM, Kleisli, and XML.

Section 3. Following that, we highlight some of the features that distinguish the more general data integration technologies from the more specialized ones. The features considered include data model, data exchange format, query capability, warehousing capability, and application programming interface.

Section 4. Lastly, we compare the surveyed technologies and comment on selecting such technologies. We also briefly discuss TAMBIS and the semantics aspect of data integration.

1. Requirements of Integration Systems for Biological Data

In a dynamic heterogeneous environment such as that of bioinformatics, many different databases and software systems are used. A large proportion of these databases were designed and built by biologists. When these databases were first created, the amount of data was small and it was important that the database entries were human readable. Database entries were therefore often created as flat

files. As new types of data were captured, new databases were created using a variety of flat file formats. We ended up with a large number of different databases in different formats, typically using non-standard query softwares.⁵² Many questions that a biologist is interested in could not be answered using any single data source. However, some of these queries can be satisfactorily solved by using information from several sources. Unfortunately, this has proved to be quite difficult in practice.

These databases and systems often do not have anything that can be thought of as an explicit database schema, which is a formalized queryable catalogue of all the tables in the database, the attributes of each of these tables, and the meaning of and indices on each of these attributes. Further compounding the problem is that research biologists demand flexible access and queries in ad-hoc combinations. Simple retrieval of data is not sufficient for modern bioinformatics. The challenge is how to manipulate the retrieved data derived from various databases and restructure the data in such a way to investigate specific biomedical problems.¹

As observed by Baker and Brass,⁵² many existing biology data retrieval systems are not fully up to the demand of painless and flexible data integration. These systems rely on low-level direct manipulation by the user, where she uses a keyword to extract summary records, then clicks on each resulting record to view its contents or to perform operations. This works well for simple actions. However, as the number of actions or records increases, such direct manipulations rapidly becomes a repetitive drudgery. Also when the questions become more complex and involve many databases, assembly of the data needed is likely to exceed the skill and patience of the biologist. Merely providing a library package that interfaces to a large number of databases and analysis softwares is also not useful if it requires long-winded and tedious programming to make use of and adding to the package.

The systems provided by bioinformaticians in answer to the challenge above can roughly be divided into “point” and “general” solutions. A point solution is a highly specialized system: the data sources to be considered are small and fixed; the biomedical research questions to be addressed are small and fixed; and the point solution is a specific software that provides the expected answers and nothing else.¹ Hence, there is little database design and consideration for extensibility nor for flexibility. In contrast, a general solution is not designed with a specific set of biomedical research questions in mind nor with a specific set of data sources in mind. It must be designed with extensibility and flexibility in mind. A general solution can serve as the platform upon which to shorten the time needed for constructing various point solutions, just as a relational database management system can serve as the platform upon which to build specific accounting systems.

A system that aims to be a general integration mechanism in the bioinformatics environment described earlier must satisfy at least the following four conditions, which were identified previously by Wong.⁸⁹⁴

- (1) It must not count on the availability of schemas. It must be able to compile any query submitted based solely on the structure of that query. If it needs a schema before it can compile a query, then it would be hard to use for our purpose because biomedical databases often do not have usable schemas.
- (2) It must have a data model that the external database and software systems can easily translate to, without doing a lot of type declarations. If it does not have such a data model, then there would be a significant impedance in moving external data into the system, in moving internal data into external databases, and in manipulating the data when they are brought into the system.
- (3) It must shield existing queries from evolution of the external sources as much as possible. For example, an extra field appearing in an external database table must not necessitate the recompilation or rewriting of existing queries over that data source. The external data sources used by a bioinformatician are typically owned by different organizations who have autonomous right to evolve their databases. It is therefore important for a general data integration solution to be robust when the data sources evolve.
- (4) It must have a data exchange format that is straightforward to use, so that it does not demand too much programming effort or contortion to capture the variety of structures of output from external databases and softwares. The data exchange format is the standard by which the system exchange data with the external data sources. If it is not straightforward to use, then great effort would be needed for connecting the system to the external data sources.

Besides the ability to query, assemble, and transform data from remote heterogeneous sources, it is also important to be able to conveniently warehouse the data locally. The reasons to create local warehouses are given below, some of which were identified previously by Davidson *et al.*¹⁸⁸

- (1) It increases efficiency. It is clear that we do not want to be choked by the slowest external data source nor by communication latency in the execution of our queries, especially if we own a fast computer. Warehousing gives us as much efficiency as we can afford to pay for.
- (2) It increases availability. It is clear that we do not want to be unable to run our queries at a time we wish because a needed external source is unavailable. Warehousing guarantees that the data we need for our queries are always available whenever we need them.

- (3) It reduces risk of unintended “denial of service” attacks on the original sources. Some data sources, such as the Entrez website at the National Center for Bioinformatic Information, impose a strict limit on the number of times or the amount of data that we can access within a single day. If we exceed that limit, then we risk being banned from the site. Unfortunately, some of our queries may require very intensive access to data held in such sites. Warehousing protects us from this risk by rendering it unnecessary for us to access the remote site.
- (4) It allows more careful data cleansing that cannot be done on the fly. It is widely acknowledged that many of the biomedical sources contain a large number of errors.¹⁰⁴ For example, Schönbach *et al.*⁷⁵⁶ reported that up to 30% of the database records that they accessed when constructing their warehouse on swine major histocompatibility complexes contained errors. Some of these errors can be detected and corrected on the fly, some of these errors cannot. It therefore makes sense that if our queries are sensitive to certain errors that cannot be detected nor corrected on-the-fly, then we should warehouse the data after careful cleansing.

Creating warehouses leads to other requirements on a general data integration solution. Specifically the general data integration solution must provide for the construction of warehouses that have the following properties.

- (1) The warehouse should be efficient to query.
- (2) The warehouse should be easy to update. There are two aspects to this issue of ease of update. The first aspect is of making an individual change to the warehouse, such as modifying an existing record, deleting an existing record, or adding a new record. This aspect is a fundamental characteristic of the data integration tools that are used for maintaining the warehouse. The second aspect is that of the number of such individual changes that need to be made to bring the warehouse up to date. The second aspect is more a consideration for the strategy for maintaining the warehouse and is dictated by the interval between updates to the warehouse and the amount of changes that the underlying data sources can accumulate during the interval. A data integration tool that offers greater ease on the first aspect obviously also allows a greater range of strategies on the second aspect.
- (3) Equally important in the biology arena is that the warehouse should model the data in a conceptually natural form. Although a relational database system is efficient for querying and easy to update, its native data model of flat tables forces us to unnaturally and unnecessarily fragment our data in order to fit our data into the third normal form.⁸⁵⁰ For example, a record in the popular

SWISS-PROT database⁴¹ would be fragmented into almost 30 tables in order to be stored in accordance to the third normal form. This unnatural fragmentation brings forth two problems. Firstly, it increases the mental load of the programmer and the possibility of programming errors in answering a queries for several reasons: (i) the implementer of a query at a later date may not be the same person who did the third normal form conversion and (ii) the implementer of a query may not be the biologist who asks the query. Secondly, it increases the cost of certain queries significantly. For example, if the query needs to reconstruct a large portion of a SWISS-PROT record, we would be required to perform 10-20 joins on the tables.

It is also important to realise that no single system is complete for all possible uses. A data integration system is rightly focused on

- (1) reading data from multiple sources for integration,
- (2) simple database-style transformation of data to facilitate data being passed from one application to the next, and
- (3) writing data to warehouses.

There are certain types of analysis and manipulations of data that a data integration system is not expected to perform but is merely expected to facilitate. These analysis and manipulations include bioinformatics specific operations such as multiple alignment and visualization-specific operations such as display data in a graphical user interface. These operations are best implemented either in a specialized scripting language designed for those purpose or in a full strength common programming language. In order to facilitate the programming of these operations, the general data integration system must provide a means for these scripting and programming languages to interface to it, via a language embedding or via an application programming interface for these languages.

Lastly, the semantics issue may also be important.¹⁸⁸ This issue concerns the equivalence and consistency between parts of records in different data sources, as well as the mappings between these parts. A data integration technology that understands which parts of two data sources have the same meanings and should be consistent with each other is desirable. However, it must be recognized that the same record in a database can sometimes be interpreted in different ways depending upon the purpose and requirement of the user. Consequently, this issue is sometimes considered as a part of building a specific application or integrated database, as opposed to as a part of the tools used for building that integrated database or application.

2. Some Data Integration Solutions

We survey here a few alternative solutions to the data integration and warehousing problem in biomedicine. The surveyed solutions include EnsEMBL,³⁷² GenoMax, SRS,²³⁵ DiscoveryLink,³¹¹ OPM,¹⁴⁵ Kleisli,⁸⁹⁴ and XML.^{2, 5} These examples are chosen to span specialized point solutions to increasingly general solutions. For each of these systems, we provide an overview and a discussion of their strong and weak points.

2.1. *EnsEMBL*

EnsEMBL is a software system jointly developed by the European Bioinformatics Institute and the Sanger Institute.³⁷² It provides easy access to eukaryotic genomic sequence data. It also performs automatic prediction of genes in these sequence data and assembles supporting annotations for these predictions. It is not so much an integration technology. However, it is an excellent example of a very successful integration of data and tools for the highly specific purpose of genome browsing.

EnsEMBL organizes raw sequence data from public databases into its internal database. It then assembles these sequences into their proper place in the genome. After that, it runs GenScan¹²³ to predict the location of genes and applies various analysis programs to annotate these predicted genes. Finally, the results of the process described above are presented for public access.

The main “entry points” to these results on the EnsEMBL Genome Browser are by

- (1) searching by sequence similarity via the built-in BLAST component of the EnsEMBL Genome Browser;
- (2) browsing from the chromosome level all the way down to the DNA sequence level;
- (3) searching using special EnsEMBL identifiers; and
- (4) free-text matching using annotation of databases linked to EnsEMBL, including OMIM,³¹⁸ SWISS-PROT,⁴¹ and InterPro.²⁹

It can also dump its data into Excel spreadsheets for use by external datamining softwares. Alternatively, the EnsMart data retrieval tool can also be used to access these results. EnsMart has a good query builder interface that allows a user to conveniently specify certain types of genomic regions and filters on these results. As a last resort, EnsEMBL provides a Perl-based programmatic interface for the most flexible access to its stored results.

Its strengths lie in its highly tailored functionalities for genome browsing. Once the sequences are imported into the system, assembly and annotation are

automatically performed, the results are automatically prepared for browsing in a nice graphical user interface.

Its weaknesses lie also in its highly tailored point solution nature. It is not possible to ask EnsEMBL to perform an *ad hoc* query in general, unless that particular type of query has been anticipated by the designer of the EnsEMBL and its associated access tools. For example, while it is possible to ask a query such as “extract 500 bases flanking the translation initiation site of each confirmed gene in the database” using EnsMart, it does not seem possible to ask a query such as “extract the first exon of each confirmed gene in the database” using EnsMart at this moment. For the latter query, the user can resort to accessing EnsEMBL and extracting the required information by Perl programming. EnsEMBL also does not have a flexible data model nor exchange format, other than the structure of its highly specialized internal database. Thus, it is not straightforward to add new kinds of data sources, and it is also not straightforward to output or export data from EnsEMBL other than in the fixed export formats.

The weaknesses mentioned above are viewed from the perspective of the requirements of a general data integration system. However, one has to remember that EnsEMBL is intended as a point solution for the specific purpose of genome browsing. Within the context of this specific purpose, EnsEMBL works much better than virtually any other alternatives, as its design has anticipated the common queries a biologist may want to ask and makes it possible for her to ask them without requiring the help of a programmer.

2.2. *GenoMax*

GenoMax is an enterprise-level integration of bioinformatics tools and data sources developed by InforMax; see <http://www.informaxinc.com/solutions/genomax>. It is a good illustration of an amalgamation of a few point-solutions, including a sequence analysis module and a gene expression module, developed on top of a data warehouse of fixed design.¹ The warehouse is an ORACLE database designed to hold sequence data, gene expression data, 3D protein structures, and protein-protein interaction information. Load routines are built in for standard data sources such as GenBank and SWISS-PROT. The specialized point-solution modules provide capabilities such as performing BLAST²⁴ and GenScan¹²³ runs on sequences and computing differentially expressed genes from microarray experiments. A special scripting language of limited expressive power is also supported for building analytical pipelines.

Its strengths are twofold. Firstly, each of GenoMax’s component point-solution modules is a very well designed application for a specific purpose. For

example, its gene expression module provides self-organizing map clustering, principal component analysis, and so forth on microarray data via simple-to-use graphical user interfaces. Secondly, these components are integrated in a tight way via the specially designed data warehouse.

Its weakness is its tight point-solution-like application integration. While GenoMax has a broader scope than EnsEMBL, it does cover less data types and products than products such as SRS, DiscoveryLink, and Kleisli. For example, these latter systems can easily incorporate chemical assay data which are beyond the current data warehouse design of GenoMax. In addition, GenoMax's scripting language is not designed for large-scale database style manipulations and hence this type of *ad hoc* queries are not always straightforward nor optimized in GenoMax. There are also difficulties in adding new kinds of data sources and analysis tools. For example, it is probably impossible to express in the GenoMax scripting language the "rosetta stone" method for extracting protein interactions.⁵⁴¹

2.3. SRS

SRS²³⁵ is marketed by LION Bioscience and is arguably the most widely used database query and navigation system for the Life Science community. It provides easy-to-use graphical user interface access to a broad range of scientific databases, including biological sequences, metabolic pathways, and literature abstracts. SRS provides some functionalities to search across public, in-house and in-licensed databases.

In order to add a new data source into SRS, this data source is generally required to be available as a flat file and a description of the schema or structure of the data source must be available as an Icarus script, which is the special built-in wrapper programming language of SRS. The notable exception to this flat file requirement on the data source is when the data source is a relational database. SRS then indexes this data source on various fields parsed and described by the Icarus script. A biologist then accesses the data by supplying some keywords and constraints on them in the SRS Query Language. Then all records matching those keywords and constraints are returned.

The SRS Query language is primarily a navigational language. This query language has limited data joining capabilities based on indexed fields and has limited data restructuring capabilities. The results are returned as a simple aggregation of records that matched the search constraints. In short, in terms of querying power, SRS is essentially an information retrieval system. It brings back records matching specified keywords and constraints. These records can contain embedded links that a user can follow individually to obtain deeper information. However, it does

not offer much help in organizing or transforming the retrieved results in a way that might be needed for setting up an analytical pipeline.

There is also a browser-based interface for formulating SRS queries and viewing results. In fact, this interface of SRS is often used by biologists as a unified front end to independently access multiple data sources, rather than learning the idiosyncrasies of the original search interfaces of these data sources. For this reason, SRS is sometimes considered¹ to serve “more of a user interface integration role rather than as a true data integration tool.”

In summary, SRS has two main strengths. Firstly, it is very straightforward to add new data sources into the system, due to the use of the Icarus scripting language and due to the simplicity of flat file indexing. In fact, several hundred data sources have been incorporated into SRS to date. Secondly, it has a nice user interface that greatly simplifies query formulation, making the system usable by a biologist without the assistance of a programmer. In addition, SRS has an extension known as PRISMA that is designed for automating the process of maintaining a SRS warehouse. PRISMA integrates the tasks of monitoring remote data sources for new data sets, and downloading and indexing such data sets.

On the other hand, SRS also has some weaknesses. Firstly, it is basically a retrieval system that simply returns entries in a simple aggregation. If the biologist wishes to perform further operations or transformations on the results, she has to do that by hand or writes a separate postprocessing program using some external scripting languages like C or Perl, which is cumbersome. Secondly, its principally flat-file based indexing mechanism rules out the use of certain remote data sources—in particular, those that are not relational databases—and does not provide for straightforward integration with dynamic analysis tools. However, this latter shortcoming is mitigated by the SCOUT suite of applications marketed by LION Bioscience that are specifically designed to interact with SRS.

2.4. *DiscoveryLink*

DiscoveryLink³¹¹ is an IBM product and, in principle, it goes one step beyond SRS as a general data integration system for biomedical data. The first thing that stands out—when DiscoveryLink is compared to SRS, EnsEMBL, and GenoMax—is the presence of an explicit data model. This data model dictates the way a DiscoveryLink user views the underlying data, the way she views results, as well as the way she queries the data.

The data model is the relational data model.¹⁷⁰ The relational data model is the de facto data model of most commercial database management systems, including the IBM's DB2 database management system upon which DiscoveryLink is based.

As a result, DiscoveryLink comes with a high-level query language, SQL, that is a standard feature of all such database management systems.

This gives DiscoveryLink several advantages over SRS. Firstly, not only can a user easily express SQL queries that go across multiple data sources—which a SRS user is able to do, but she can also perform further manipulations on the results—which a SRS user is unable to do. Secondly, not only are the SQL queries more powerful and expressive than those of SRS, the SQL queries are also automatically optimized by DB2. The use of query optimization allows a user to concentrate on getting her query right without worrying about getting it fast.

However, DiscoveryLink still has a some way to go in practice. The reason is twofold. The first reason is that DiscoveryLink is tied to the relational data model. This implies every piece of data that it handles must be a table of atomic objects like strings and numbers. Unfortunately, most of the data sources in biology are not that simple and are deeply nested. Therefore, there is severe impedance mismatch between these sources and DiscoveryLink. Consequently, it is not straightforward to add new data sources or analysis tools into the system. For example, to put the SWISS-PROT database into a relational database in the third normal form would require us to break every SWISS-PROT record into nearly 30 pieces in a normalization process! Such a normalization process requires a certain amount of skill. Similarly, to query the normalized data in DiscoveryLink requires some mental and performance overhead, as we need to figure out which part of SWISS-PROT has gone to which of the 30 pieces and we need to join some of the pieces back again.

The second reason is that DiscoveryLink supports only wrappers written in C++, which is not the most suitable programming language for writing wrappers. In short, it is difficult to extend DiscoveryLink with new sources. In addition, DiscoveryLink does not store nested objects in a natural way and is limited in its capability for handling long documents. It also has limitations as a tool for creating and managing data warehouses for biology.

In spite of these weaknesses, in theory, DiscoveryLink has greater generality than point solutions like EnsEMBL, specialized application integration like GenoMax, and user interface integration solutions like SRS. Unfortunately, this greater generality is achieved at the price of requiring that SQL be used for expressing queries. While writing queries in SQL is generally simpler than writing in Perl, it is probably still beyond the skill of an average biologist. This is a disadvantage in comparison to EnsEMBL, GenoMax, and SRS, which have good user interfaces for a biologist to build the simpler queries.

2.5. OPM

OPM¹⁴⁵ was developed at Lawrence-Berkeley National Labs. OPM is a general data integration system. OPM was marketed by Gene Logic, but its sales were discontinued some time ago. It goes one step beyond DiscoveryLink in the sense that it has a more powerful data model, which is an enriched form of the entity-relationship data model.¹⁴⁹

This data model can deal with the deeply nested structure of biomedical data in a natural way. Thus it removes the impedance mismatch. This data model is also supported by a SQL-like query language that allows data to be seen in terms of entities and relationships. Queries across multiple data sources, as well as transformation of results, can be easily and naturally expressed in this query language. Queries are also optimized. Furthermore, OPM comes with a number of data management tools that are useful for designing an integrated data warehouse on top of OPM.

However, OPM has several weaknesses. Firstly, OPM requires the use of a global integrated schema. It requires significant skill and effort to design a global integrated schema well. If a new data source needs to be added, the effort needed to re-design the global integrated schema potentially goes up quadratically with respect to the number of data sources already integrated. If an underlying source evolves, the global integrated schema tends to be affected and significant re-design effort is potentially needed. Therefore, it may be costly to extend OPM with new sources.

Secondly, OPM stores entities and relationships internally using a relational database management system. It achieves this by automatically converting the entities and relationships into a set of relational tables in the third normal form. This conversion process leads to an entity being broken up into many pieces when stored. This process is transparent to the OPM user. So she can continue to think and query in terms of entities and relationships. Nevertheless, the underlying fragmentation often causes performance problems, as many queries that required no join—when viewed at the conceptual level of entities and relations—are mapped to queries that required many joins on the physical pieces that entities are broken into.

Thirdly, OPM does not have a simple format to exchange data with external systems. At one stage, it interfaces to external sources using CORBA. The effort required for developing CORBA-compliance wrappers is generally significant.⁷⁶⁴ Furthermore, CORBA is not designed for data intensive applications.

Although OPM's query language is at a higher level and is simpler to use than the SQL of DiscoveryLink, it shares the same disadvantage as DiscoveryLink

from the perspective of an average biologist. The programming of queries other than the simplest kind is probably still beyond her expertise.

2.6. Kleisli

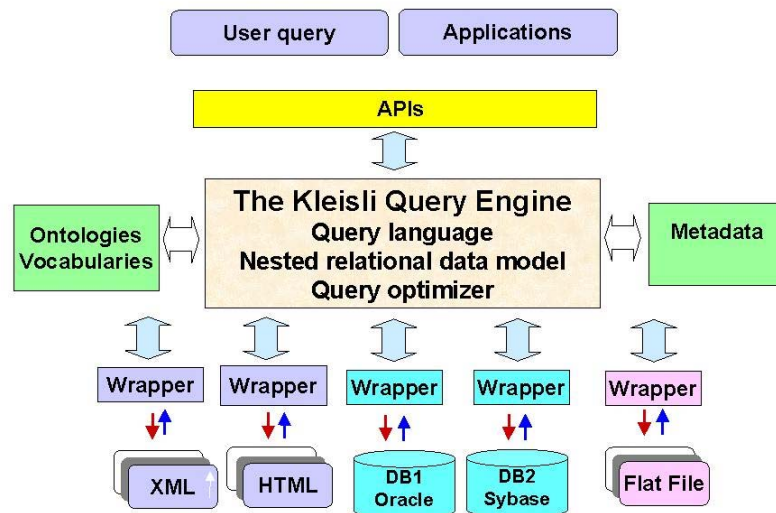


Fig. 1. Kleisli, positioned as a mediator

Kleisli^{162, 189, 894} is marketed by geneticXchange Inc of Menlo Park. It is one of the earliest systems that have been successfully applied to some of the earliest data integration problem in the human genome project, including the so-called US Department of Energy's "impossible" queries in early 1994.

The approach taken by the Kleisli system is illustrated by the diagram in Figure 1. It is positioned as a mediator system encompassing a nested relational data model, a high-level query language, and a powerful query optimizer. It runs on top of a large number of light-weight wrappers for accessing various data sources. There are also a number of application programming interfaces that allow Kleisli to be accessed in a ODBC- or JDBC-like fashion in various programming languages for a various applications.

The Kleisli system is highly extensible. It can be used to support several different high-level query languages by replacing its high-level query language module. Currently, Kleisli supports a “comprehension syntax”-based language called CPL⁸⁹⁴ and a “nested relationalized” version of SQL called sSQL. The Kleisli system can also be used to support many different types of external data sources by adding new wrappers, which forward Kleisli’s requests to these sources and translate their replies into Kleisli’s exchange format. These wrappers are light weight and new wrappers are generally easy to develop and insert into the Kleisli system. The optimizer of the Kleisli system can also be customized by different rules and strategies.⁸⁹⁴

Besides the ability to query, assemble, and transform data from remote heterogeneous sources, it is also important to be able to conveniently warehouse the data locally. Kleisli does not have its own native database management system. Instead, Kleisli has the ability to turn many kinds of database systems into an updatable store conforming to its nested relational data model. In particular, Kleisli can use flat relational database management systems such as Sybase, Oracle, MySQL, *etc.* to be its updatable store. It can even use all of these systems simultaneously. It is also worth noting that Kleisli stores nested relations into flat relational database management systems using an encoding scheme that does not require these nested relations to be fragmented over several tables.

Kleisli possesses the following strengths.⁸⁹⁴ It does not require data schemas to be available. It has a nested relational data model and a data exchange format that external databases and software systems can easily translate into. It shields existing queries, via a type inference mechanism, from certain kinds of structural changes in the external data sources. Kleisli also has the ability to store, update, and manage complex nested data. It has a good query optimizer. Finally, Kleisli is also equipped with two application programming interfaces so that it can be accessed in a JDBC-like manner from Perl and Java.⁸⁹⁵

However, Kleisli shares a common weakness with DiscoveryLink and OPM. Even though CPL and sSQL are both high-level query languages and protect the user from many low level details—such as communication protocols, memory management, thread scheduling, and so on—the programming of queries using CPL or sSQL other than the simplest kind is probably still beyond the expertise of an average biologist.

2.7. XML

XML is a standard for formatting document. As such, XML is not a data integration system by itself. However, there is a growing suite of tools based on XML

that, taken as a whole, can be used as a data integration system. We therefore believe it is pertinent to include a discussion on XML and its associated tools in the context of this paper.

XML allows for a hierarchical nesting of tags and the set of tags can be defined flexibly. Thus XML can be viewed as a powerful data model² and a useful data exchange format, providing directly for two of the important ingredients of a general data integration solution for biomedicine. As a result, an increasing number of tools and sources in biomedicine such as PIR,⁹⁰³ Entrez, and so on are becoming XML compatible.⁵

The intense interest in the development of query languages for semi-structured data² in the database community has also resulted in a number of powerful XML query languages such as XQL⁷¹⁶ and XQuery,²⁰² which provide the means for querying across multiple data sources and for transforming the results into more suitable form for subsequent analysis steps. Research and development works are also in progress on XML query optimization²⁴⁶ and on XML data stores.²⁵²

A robust and stable XML-based general data integrating and warehousing system does not yet exist for biomedicine. However, once high-performance XML data stores become available, we can also expect the database research community to begin more research and development on data warehousing using these stores.

Consequently, we believe that given sufficient time, XML and the growing suite of XML-based tools can mature into an alternative data integration system in biomedicine that is comparable to Kleisli in generality and sophistication.

3. Highlight of Selected Features

This section highlights some features that distinguish the more general data integration technologies from the more specialized data integration solutions surveyed earlier.

3.1. Data Model and Data Exchange Format

A key feature that separates the more general data integration technologies—DiscoveryLink, OPM, Kleisli—from the more specialized technologies—Ensembl, GenoMax, SRS—is the explicit presence of a data model. From the point of view of the traditional database world,³ a data model provides the means for specifying particular data structures, for constraining the data associated with these structures, and for manipulating the data *within* a database system. In order to handle data outside of the database system, this traditional concept of a data model is extended to include a data exchange format, which is a means for bring-

ing data outside the database system into it and also for bringing data inside the database system outside. We use Kleisli's data model to illustrate this concept.

The data model underlying the Kleisli system is a complex object type system that goes beyond the "sets of records" or "flat relations" type system of relational databases.¹⁷⁰ It allows arbitrarily nested records, sets, and a few other data types.⁸⁹⁴ Having such a "nested relational" data model is useful and matches the structure of biomedical data sources well. For example, if we are restricted to the flat relational data model, the GenPept report in Example 1 must necessarily be split into many separate tables in order to be losslessly stored in a relational database. The resulting multi-table representation of the GenPept report is conceptually unnatural and operationally inefficient.

Example 1: The GenPept report is the format chosen by the US National Center for Biotechnology Information to present amino acid sequence information. The feature table is the part of the GenPept report that documents the positions and annotations of regions of special biological interest. The following type represents the feature table of a GenPept report from Entrez.⁷⁵⁸ Here we use {– and }–brackets for sets, (– and)–brackets for records, [– and]–brackets for lists, and #*l* : to label the field *l* of a record. In fact, the same bracketing scheme is used as the data exchange format of Kleisli.⁸⁹⁵

```
(#uid:num, #title:string, #accession:string, #feature:{(
  #name:string, #start:num, #end:num, #anno:[(
    #anno_name:string, #descr:string)])})
```

The feature table of GenPept report 131470, a tyrosine phosphatase 1C sequence, is shown partially below. The particular feature displayed goes from amino acid 0 to amino acid 594, which is actually the entire sequence, and has two annotations: The first annotation indicates that this amino acid sequence is derived from mouse DNA sequence. The second is a cross reference to the US National Center for Biotechnology Information taxonomy database.

```
(#uid:131470, #accession:"131470", #title:"... (PTP-1C)...",
 #feature:{(
  #name:"source", #start:0, #end:594, #anno:[
    (#anno_name:"organism", #descr:"Mus musculus"),
    (#anno_name:"db_xref", #descr:"taxon:10090")]),
 ...})
```

□

It is generally easy to develop a wrapper for a new data source, or modifying an existing one, and insert it into Kleisli. The main reason is that there is no impedance mismatch between the data model supported by Kleisli and the data model that is necessary to capture the data source. The wrapper is often a very

light-weight parser that simply parses records in the data source and prints it out in Kleisli's very simple data exchange format.

Example 2: Suppose we want to implement a function `webomim-get-detail` that uses an OMIM³¹⁸ identifier to access the OMIM database and returns a set of objects matching the identifier. Suppose the output is of type

```
{(#uid: num, #title: string, #gene_map_locus: {string},
  #alternative_titles: {string}, #allelic_variants: {string})}
```

Note that is this a nested relation: it is a set of records, and each record has three fields that are also of set types, viz. `#gene_map_locus`, `alternative_titles`, and `allelic_variants`. This type of output would definitely present a problem if we had to give it to a system based on the flat relational model, as we would need to arrange for the information in these three fields to be sent into separate tables. Fortunately, such a nested structure can be mapped directly into Kleisli's exchange format. So the wrapper implementor would only need to parse each matching OMIM records and to write it out in a format like this:

```
{(#uid: 189965,
  #title: "CCAAT/ENHANCER-BINDING PROTEIN, BETA; CEBPB",
  #gene_map_locus: "20q13.1",
  #alternative_titles: {
    "C/EBP-BETA",
    "INTERLEUKIN 6-DEPENDENT DNA-BINDING PROTEIN; IL6DBP",
    "LIVER ACTIVATOR PROTEIN; LAP",
    "LIVER-ENRICHED TRANSCRIPTIONAL ACTIVATOR PROTEIN",
    "TRANSCRIPTION FACTOR 5; TCF5"},
  #allelic_variants: {}}}
```

Here, instead of needing to create separate tables to keep the sets nested inside each record, the wrapper would simply print the appropriate set brackets `{` and `}` to enclose these sets. Kleisli would automatically deal with them as they were handed over by the wrapper. This kind of parsing and printing is extremely easy to implement. \square

OPM shares with Kleisli a nested relational data model, except that the former lacks a data exchange format. Hence the mapping of the examples to OPM's data model is conceptually just as straightforward, but the practical implementation in OPM demands considerably more effort. It is worth pointing out that while SRS does not have an explicit data model, it does have an implicit one supported by its Icarus language for scripting parsers. In the case of SRS, this implicit data model in Icarus also greatly facilitates the rapid scripting of wrappers.

3.2. Query Capability

Another feature that separates the more general data integration technologies from the more specialized ones is the presence of a flexible high level query language for manipulating data conforming to the data model. We use sSQL, the primary query language of Kleisli to illustrate this feature. sSQL is based on the *de facto* commercial database query language SQL, except for extensions made to cater for the nested relational model and for the federated heterogeneous data sources.

Example 3: The feature table of a GenBank report has the type below. The field #position of a feature entry is a list indicating the start and stop positions of that feature. If the feature entry is a CDS, this list corresponds to the list of exons of the CDS. The field #anno is a list of annotations associated with the feature entry.

```
(#uid: num, #title: string, #accession: string, #seq: string,
 #feature: { (
   #name: string,
   #position: [(#start:num,#end:num,#negative:bool, ...)],
   #anno: [(#anno_name:string,#descr:string)], ... )
```

Given a set DB of feature tables of GenBank chromosome sequences, we can extract the 500 bases up stream of the translation initiation sites of all disease genes—in the sense that these genes have a cross reference to OMIM—on the positive strand in DB as below. Here, `l2s` is a function that converts a list into a set.

```
select
  uid: x.uid,
  protein: r.descr,
  flank: string-span(x.seq, p.start - 500, p.start)
from
  DB x, x.feature f,
  {f.position.list-head} p,
  f.anno.l2s a, f.anno.l2s r
where not (p.#negative)
and   a.descr like "MIM:%" and a.anno_name = "db_xref"
and   r.anno_name = "protein_id"
```

Similarly, we can extract the first exons of these same genes as follows:

```
select
  uid: x.uid,
  protein: r.descr,
  exon1: string-span(x.seq, p.start, p.end)
from
  DB x, x.feature f,
  {f.position.list-head} p,
```

```

f.anno.l2s a, f.anno.l2s r
where not (p.#negative)
and a.descr like "MIM:%" and a.anno_name = "db_xref"
and r.anno_name = "protein_id"

```

These two example queries illustrate the how a high level query language makes it possible to extract very specific output in a relatively straightforward manner. \square

We illustrate how to combine multiple sources using high level query languages. An *in silico* discovery kit (ISDK) prescribes experimental steps carried out in computers very much like the experimental protocol carried out in wet laboratories for specific scientific investigation. From the perspective of Kleisli, an *in silico* discovery kit is just a script written in sSQL and performs a defined information integration task. It takes an input data set and parameters from the user, executes and integrates the necessary computational steps of database queries and applications of analysis programs or algorithms, and outputs a set of results for specific scientific inquiry.

Example 4: The simple *in silico* discovery kit in Figure 2 demonstrates how to use an available ontology data source to get around the problem of inconsistent naming in genes and proteins, and to integrate information across multiple data sources. It is implemented in the sSQL script below.

```

create function get-info-by-genename (G) as
Select
  hugo: w, omim: y, pmid1-abstract: z,
  num-medline-entries: list-sum(
    lselect ml-get-count-general(n)
    from x.Aliases.s2l n)
from
  hugo-get-by-symbol(G) w,
  webomim-get-id(
    searchtime:0, maxhits:0,
    searchfields:{}, searchterms:G) x,
  webomim-get-detail(x.uid) y,
  ml-get-abstract-by-uid(w.PMID1) z
where
  x.title like ("% ^ G ^ %");

```

With the user input of a gene name G , the ISDK performs the following task: First, it retrieves a list of aliases for G from the Gene Nomenclature database provided by the Human Genome Organization (HUGO). Then it retrieves information for diseases associated with this particular protein in the Online Mendelian Inher-

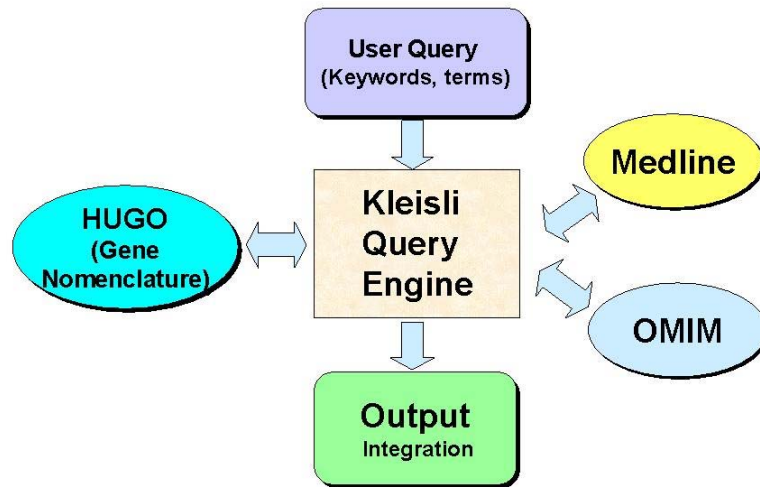


Fig. 2. An “*in silico* discovery kit” that uses an available ontology data source to get around the problem of inconsistent naming in genes and proteins, and integrates information across multiple data sources.

itance of Man Database (OMIM),³¹⁸ and finally it retrieves all relevant references from MEDLINE.

Here, `s2l` is a function that converts a set into a list; `list-sum` is a function to sum a list of numbers; `ml-get-count-general` is a function that accesses the MEDLINE database in Bethesda and computes the number of MEDLINE reports matching a given keyword; `ml-get-abstract-by-uid` is a function that accesses MEDLINE for report given a unique identifier; `webomim-get-id` is a function that accesses the OMIM database in Bethesda to obtain unique identifiers of OMIM reports matching a keyword; `webomim-get-detail` is a function that accesses OMIM for report given a unique identifier; and `hugo-get-by-symbol` is a function that accesses the HUGO database and return HUGO reports matching a given gene name.

For instance, this query `get-info-by-genename` can be invoked with the transcription factor `CEBPB` as input to obtain the following result.

```
{(#hugo: (#HGNC: "1834",
  #Symbol: "CEBPB", #PMID1: "1535333", ...
```

```
#Name: "CCAAT/enhancer binding protein (C/EBP), beta",
#Aliases: {"LAP", "CRP2", "NFIL6", "IL6DBP", "TCF5"},
#omim: (#uid: 189965, #gene_map_locus: "20q13.1",
#allelic_variants: {}, ...),
#pmid1-abstract: (#muid: 1535333,
#authors:"Szipirer C...", #address:"Departement ...",
#title: "Chromosomal localization in man and rat ...",
#abstract: "By means of somatic cell hybrids ...",
#journal: "Genomics 1992 Jun; 13(2):292-300",
#num-medline-entries: 1936)}
```

Such queries fulfill many of the requirements for efficient *in silico* discovery processes:

- (1) their modular nature gives scientists the flexibility to select and combine specific queries for specific research project;
- (2) they can be executed automatically by Kleisli in batch mode and can handle large data volume;
- (3) their scripts are reusable to perform repetitive tasks and can be shared among scientific collaborators;
- (4) they form a base set of templates that can be readily modified and refined to meet different specifications and to make new queries; and
- (5) new databases and new computational tools can be readily incorporated to existing scripts. □

The flexibility and power shown in these sSQL examples can also be experienced in OPM, and to a lesser extent in DiscoveryLink. With good planning, a specialised data integration system can also achieve great flexibility and power within a narrower context. For example, the EnsMart tool of EnSEMBL is a very well designed interface that helps a non-programmer build complex queries in a simple way. In fact, an equivalent query to the first sSQL query in Example 3 can be also be specified using EnsMart with a few clicks of the mouse. Nevertheless, there are some unanticipated cases that cannot be expressed, such as the second sSQL query in Example 3.

3.3. Warehousing Capability

Besides the ability to query, assemble, and transform data from remote heterogeneous sources, a general data integration technology is also expected to be able to conveniently store data locally. Kleisli does not have its own native database management system. Instead, Kleisli has the ability to turn many kinds of database systems into an updatable store conforming to its complex object data model. In

particular, Kleisli can use flat relational database management systems such as Sybase, Oracle, MySQL, *etc.* to be its updatable complex object store. It can even use all of these systems simultaneously. We illustrate this feature using the example of GenPept reports.

Example 5: Create a warehouse of GenPept reports. Initialize it to reports on protein tyrosine phosphatases. Note that Kleisli provides several functions to access GenPept reports remotely from Entrez.⁷⁵⁸ One of them is the function `aa-get-seqfeat-general` used below, which retrieves GenPept reports matching a search string.

```
! connect to our Oracle database system
oracle-cplobj-add (name: "db", ...);
! create a table to store GenPept reports
create table genpept(uid:"NUMBER", detail: "LONG") using db;
! initialize it with PTP data
select (uid: x.uid, detail: x) into genpept
from aa-get-seqfeat-general("PTP") x using db;
! index the uid field for fast access
db-mkindex(table:"genpept",index:"genpeptindex",schema:"uid");
! let's use it now to see the title of report 131470
create view GenPept from genpept using db;
select x.detail.title from GenPept x where x.uid = 131470;
```

In this example, a table `genpept` is created in our Oracle database system. This table has two columns, `uid` for recording the unique identifier and `detail` for recording the GenPept report. A LONG data type is used for the `detail` column of this table. However, recall from Example 1 that each GenPept report is a highly nested complex objects. There is therefore a “mismatch” between LONG and the complex structure of a GenPept report. This mismatch is resolved by the Kleisli system which automatically performs the appropriate encoding and decoding.

Thus, as far as the Kleisli user is concerned, `x.detail` has the type of GenPept report as given in Example 1. So he can ask for the title of a report as straightforwardly as `x.detail.title`. Note that encoding and decoding are performed to map the complex object transparently into the space provided in the `detail` column; that is, the Kleisli system does not fragment the complex object to force it into the third normal form. □

What distinguishes the ability to store or warehouse data in a system like Kleisli is the ease with which the store or warehouse can be specified and populated in theory. However, in practice, one does need to anticipate failures in a large update and to provide recovery mechanisms. OPM also enjoys this capability if

one is not concerned with how the store or warehouse is supposed to be related to the other existing stores or warehouses already in the system. DiscoveryLink enjoys this capability to a lesser extent, as its data model does not permit nesting of data. The more specialized systems can also have warehousing capability, but with greater limitations on what can be warehoused by them.

3.4. *Application Programming Interfaces*

The high-level query languages of the more general data integration systems surveyed are all SQL-like and are thus designed to express traditional (nested relational) database-style queries. Not every query in bioinformatics falls into this class. For these non-database-style queries, some other programming languages can some time be a more convenient or more efficient means of implementation. Therefore, it is useful to develop some application programming interfaces to these more general data integration systems for various popular programming languages.

In the case of Kleisli, there is the Pizzkell suite⁸⁹⁵ of interfaces to the Kleisli Exchange Format for various popular programming languages. Each of these interfaces in the Pizzkell suite is a library package for parsing data in Kleisli's exchange format into an internal object of the corresponding programming language. It also serves as a means for embedding the Kleisli system into that programming language, so that the full power of Kleisli is available within that programming language in a manner similar to that achieved by JDBC and ODBC for relational databases. The Pizzkell suite currently include CPL2Perl and CPL2Java, for Perl and Java.

The presence of such application programming interfaces may be even more crucial for the more specialized integration solutions. While a point solution like EnsEMBL is typically designed with a specific aim in mind, it is not unusual to subsequently discover that a user wants to use the integrated data in an unanticipated way. In such a situation, it would be convenient if an application programming interface is available on the integrated data. For example, in the case of EnsEMBL, as EnsEMBL is implemented in Perl using Bioperl⁷⁹⁷ as the backbone, the same library of routines that have been accumulated in the course of implementing EnsEMBL would be the perfect application programming interface to EnsEMBL.

Bioperl⁷⁹⁷ itself can also be thought of as a low level integration toolkit for biological data. Such a toolkit typically contains a set of library routines, for accessing some commonly used bioinformatics data sources and tools, that can be invoked as procedure or function calls from the underlying programming lan-

guage. No high-level query support and no optimization are provided. There are also similar toolkits in other popular “open source” programming languages, such as BioPython and Biojava.⁵³⁸

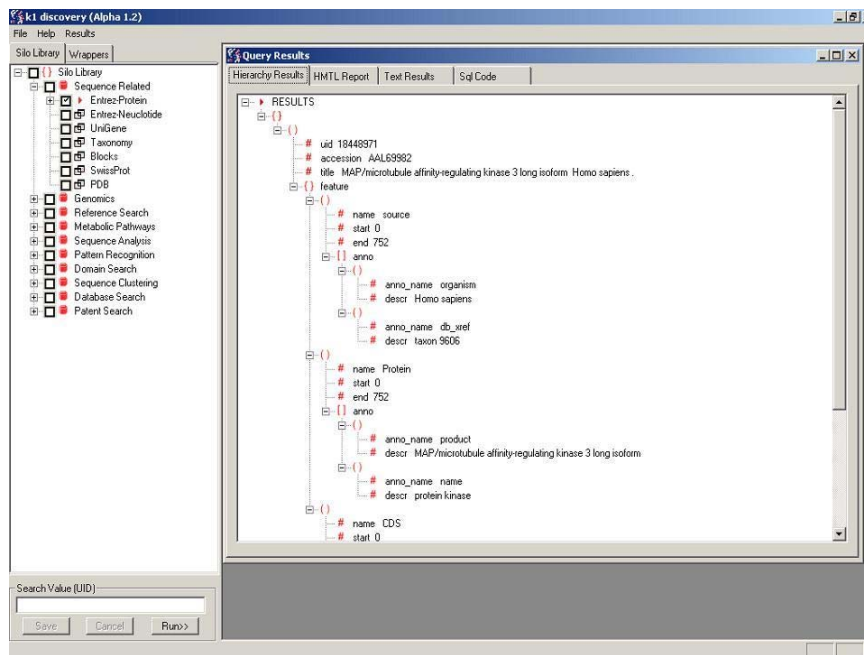


Fig. 3. A screenshot of the Discovery Builder, a graphical interface to Kleisli.

There is also a graphical interface to the Kleisli system that is designed for non-programmers. It is called the Discovery Builder and is developed by folks at geneticXchange Inc. This graphical interface makes it far easier to visualize the source data required to formulate the queries and generates the necessary sSQL codes. It allows a user to see all available data sources and their associated metadata and assists the user to navigate and to specify his query on these sources. A screenshot of Discovery Builder is presented in Figure 3. It provides all of the following key functions:

- (1) a graphical interface that can “see” all of the relevant biological data sources, including metadata—tables, columns, descriptions, *etc.*—and then construct a query “as if” the data were local;
- (2) add new wrappers for any public or proprietary data sources, typically within

hours, and then have them enjoined in any series of ad-hoc queries that can be created;

- (3) execute the queries, which may join many data sources that can be scattered all over the globe, and get fresh result data quickly.

4. Concluding Remarks

Let us first summarize our opinion on how well each of the surveyed systems satisfy the requirements of a general data integration system for biomedicine. EnSEMBL and GenoMax are point solutions and thus naturally do not satisfy the requirements of a general data integration system well.

SRS and DiscoveryLink are claimed by their inventors as general data integration systems for biomedicine. However, in reality, SRS is a form of user interface integration and hence it does not satisfy the requirements well. On the other hand, while DiscoveryLink has most of the components required, these components are probably in the wrong flavour—the adoption of the flat relational model causes it to be less potent in the biomedical data integration arena.

OPM is a well designed system for the purpose of biomedical data integration, except for (i) a problem in performance due to data fragmentation as it unwisely maps all data to the third normal form, and (ii) the lack of a simple data exchange format, and (iii) the need of a global schema.

XML and Kleisli have all the qualities required for a good general data integration. However, compared to Kleisli, XML still need more time to mature especially in terms of query optimization and data warehousing capabilities.

Let us next look at these surveyed systems from the perspective of an average biologist. While general data integration systems such as DiscoveryLink, OPM, and Kleisli simplify the programming of *ad hoc* queries, it must also be acknowledged that the programming skills required are still significant. In contrast, data integration systems that are nearer to the point-solution end of the spectrum—such as EnSEMBL, GenoMax, and SRS—have considerably better user interfaces that help a biologist to build the simpler type of queries.

Of course, a biologist may find it frustrating that the graphical user interfaces of EnSEMBL, GenoMax, and SRS cannot let her express a particular *ad hoc* query such as the one that asks for the sequence of the first exon of all genes in a database. However, it is very likely that the same biologist may also find it equally frustrating that she does not know how to express that query in DiscoveryLink, OPM, and Kleisli, even though she knows that that query is expressible in these systems. In other words, the more general data integration systems can directly increase the productivity of a bioinformatics programmer, but they probably

cannot directly increase the productivity of an average biologist.

Drawing from the remarks above, we see a dichotomy between expressiveness and simplicity. Therefore, which type of data integration system is preferred necessarily depends on the trade-off between these two factors. Many problems in biomedical research on drug targets and candidates require access to many data sources that are voluminous, heterogeneous, complex, and geographically dispersed. If these data sources are successfully integrated into a new database, researchers can then uncover relationships that enable them to make better decisions on understanding and selecting targets and leads. Therefore, a successful integration of data is crucial to improving productivity in this research.

It is important to stress that a successful data integration must be in support of a specific research problem, and different research problems are likely to need different ways of integrating and analysing data. Even though point solutions such as EnsEMBL does not fair well as a general data integration system, it works much better than any general data integration system in the specific context of genome browsing.

If one's data integration needs are of a more *ad hoc* nature, a general data integration system can often ease the implementation significantly as such a system provides greater adaptability. It is also worth remarking that the more specialized solutions may themselves be implemented on top of a more general data integration solution. One such example is TAMBIS,²⁹¹ which is built on top of Kleisli.

The systems surveyed so far generally do not consider the semantics aspect of the underlying data sources. Let us end this chapter with a brief mention of TAMBIS. TAMBIS²⁹¹ is a data integration solution that specifically addresses the semantics aspect. The central distinguishing feature of TAMBIS is the presence of an ontology and a reasoning system over this ontology. The TAMBIS ontology contains nearly 2000 concepts that describe both molecular biology and bioinformatics tasks. TAMBIS provides a user interface for browsing the ontology and for constructing queries.

A TAMBIS query is formulated by starting from one concept, browsing the connected concepts and applicable bioinformatics operations in the ontology, selecting one such connected concepts or applicable bioinformatics operation, and browsing and selecting for further connected concepts and applicable bioinformatics operations. The ontology and the associated reasoning component thereby guide the formulation of the query, ensuring that only a query that is logically meaningful can be formulated. The query is then translated by TAMBIS and passed to an underlying Kleisli system for execution.

From the point of view of TAMBIS, Kleisli significantly simplifies the task of implementing TAMBIS, as the TAMBIS implementors can concentrate on the

ontology and reasoning components and leave the details of handling the underlying data sources to Kleisli. From the point of view of Kleisli, TAMBIS makes it possible for a biologist to ask more complicated *ad hoc* queries on the data sources integrated by Kleisli without the assistance of a programmer. The ontology of TAMBIS is currently being enriched by its inventors in the University of Manchester to allow an even larger range of complicated queries to be expressed.