

# Linking publication, gene and protein data

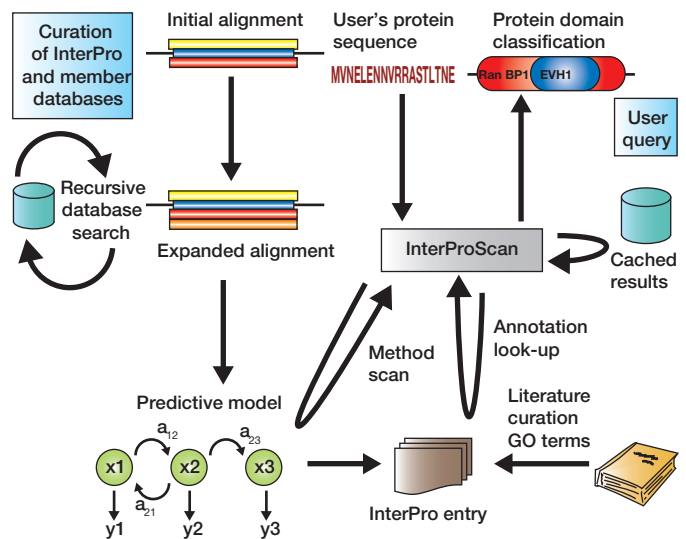
Paul Kersey and Rolf Apweiler

The computational reconstruction of biological systems, 'systems biology', is necessarily dependent on the existence of well-annotated data sets defining and describing the components of these systems, especially genes and the proteins they encode. Information about these components can be accessed either through structured bioinformatics databases, which store basic chemical and functional information abstracted from (or supplementing) the scientific literature, or through the literature itself, which is richer in content but essentially unstructured.

The systematic application of automated high-throughput molecular biology techniques has led to the generation of an immense quantity of data compared with that produced by the hypothesis-driven application of earlier technologies. For example, advances in DNA sequencing technology have made the study of complete genomes routine, and have been followed by similar progress in the study of transcripts, proteins and metabolites, although these generate far more complex datasets. In fact, the need for databases that store molecular biological data, and that allow analysis through computational algorithms, was apparent long before experimental techniques were as powerful — and as data generative — as they are today. The first attempts to catalogue information about proteins began in 1965, and by 1984 these had evolved into the Protein Information Resource<sup>1</sup>, the Protein Data Bank<sup>2</sup> (which stores information about protein structures and was founded in 1971) and the EMBL data library<sup>3</sup> (the first database to store information about nucleic-acid sequences, founded in 1981). Today, these databases and their successors have been joined by numerous other resources, which store information on chemical entities, gene expression, molecular interactions and biochemical pathways. These databases often contain more raw data than could be practically published in a conventional hard-copy journal. However, the interpretation of these data is still dependent on inference drawn from hypothesis-driven experimentation, the details of which reside in free-text articles. The value of bioinformatics data is thus utterly dependent on the ability to make the correct links from the sequences to the scientific literature, and to extract the information that literature contains.

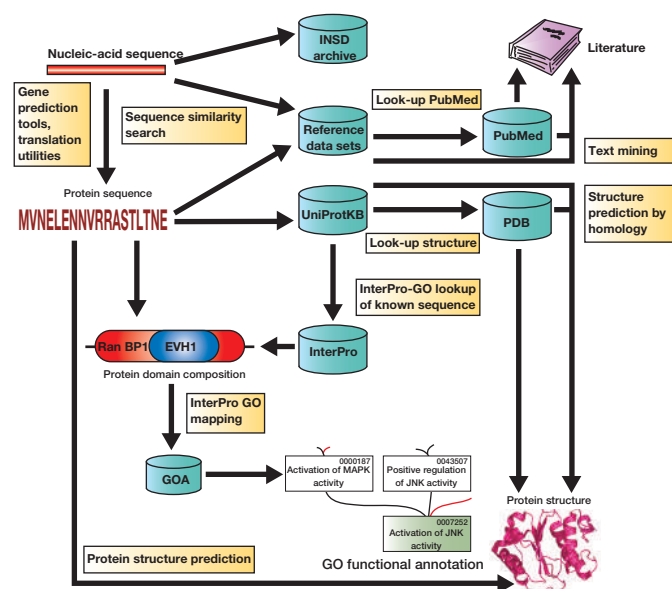
This article is a user's guide to linking gene, protein and publication data. It considers the principal primary resources in which such data can be found, as well as the difficulties and potential pitfalls when linking them, and some of the general approaches and specific tools that can be used to overcome these problems. We also discuss the semantic web, a recent technological framework proposed for the integration of distributed data that is attracting particular interest within the bioinformatics community. The need for new solutions is acute, as the quantity and

Paul Kersey and Rolf Apweiler are in the EMBL-European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK. e-mail: pkersey@ebi.ac.uk



**Figure 1** Protein analysis using InterProScan. Curators of member databases of InterPro identify proteins that share a known functional domain and manually identified sequences can be supplemented through iterative searching of the public databases until a comprehensive set is found. The sequences are aligned and a classifying function is built from the alignment, which is designed to recognise other proteins that possess the same domain. Many of these classifiers are built using hidden Markov models. Alternative classifying functions judged (by curators) to identify the same domain are grouped into a single entry in the InterPro database, and annotated with relevant information (expressed both in free text and in the terms of the GO controlled vocabulary). InterProScan is a programme that allows users to characterize a sequence by applying these classifying functions. Performance can be improved by precomputing the results for known sequences.

variety of available data, including transcriptome, proteome, metabolome, molecular and genetic interaction, and image data continues to rise. However, this deluge of data is also an opportunity for development and has led to the growth of the discipline of systems biology, which aims to reconstruct entire biological systems through the modelling of their components. This requires that all components are clearly identified and well-described, which for a large system means, in practice, that descriptive annotation must be inferred, often automatically, from related



**Figure 2** A schematic representation of a typical workflow for bioinformatics analysis. The figure shows a typical workflow for bioinformatics analysis, progressing from sequence to functional annotation, protein structure and literature. The complete analysis may be carried out entirely within a bioinformatics warehousing system (such as SRS), or as a sequence of separate operations performed in different environments. Starting with the sequence of a gene or protein, identical and/or similar sequences are identified in the public databases. The database records describing these sequences also contain general information about the sequence, curated links to structural information and relevant scientific literature. Protein sequence itself can also be directly analysed to predict its domain composition and structure. These may provide a more reliable indication of protein function than overall sequence similarity. It is often useful to express the results of these analyses in standard controlled vocabularies (such as GO), to allow the comparison and correspondence of information derived from different resources.

components in other systems. Therefore, the establishment and exploitation of reliable connections between the underlying source data is not only necessary for primary sequence analysis, but is also an essential part of the platform needed to support the computational, synthetic approaches to biology that are emerging in the post-genomic era.

## PRIMARY BIOINFORMATICS DATABASES

Here, we discuss some of the main resources that function as primary repositories of bioinformatics data and the types of data that they contain, and introduce MEDLINE and PubMed, two databases that offer computerised access to scientific literature.

### BOX 1 DATA WAREHOUSES

A data warehouse is a database constructed to support efficient querying of the data it contains (in contrast to normalized databases designed to support data integrity, which are widely used to maintain primary resources). A feature of many data warehouses used in bioinformatics is that they provide generic query interfaces (for example, computer languages and graphical user interfaces) applicable to all the data they contain, thus enabling the addition of new data without the need for interface redesign. A single warehouse may be built from several different resources, but to allow the construction of queries that filter and/or extract information derived from more than one of these, the data must be fitted into a single model that captures the relationships between the different sources. This is often done by exploiting the cross-references that many of these sources contain. The centralization of data in a data warehouse is an alternative to the use of technologies designed to support distributed queries. Distributed approaches can avoid certain problems associated with data warehousing, such as the synchronization of updates, but queries accessing large quantities of data spread over many locations can be difficult to optimise.

## Nucleotide databases

The EMBL<sup>3</sup>–Genbank<sup>4</sup>–DDBJ<sup>5</sup> International Nucleotide Sequence Database (INSD) is the product of a collaborative effort to maintain a single global archive for nucleotide sequences. As most other bioinformatic data is dependent on inference from nucleotide sequence, it is arguably the single most important bioinformatics resource. Its comprehensive coverage is maintained through the editorial policy of most scientific journals — submitting authors have to deposit the sequence information that is associated with articles they publish. The joint database currently contains over 80 million sequences.

The information contained within the archive, however, is not uniformly well annotated. The contents of existing databases are used in most methods for gene prediction, and thus, once these databases have been updated, earlier predictions may differ from those that could now be made. In addition, individual submitters are free to choose what annotation they consider pertinent to add to a sequence. As a result, nothing can be inferred from the absence of information in database records — this particular information may simply have been unknown to, or ignored by, the submitter.

With the advent of whole genome sequencing, a new class of database has emerged for data from higher eukaryotic species, in which sequence and automatically generated annotation is made available in advance of not only experimental confirmation, but even the existence of definitive methods for gene prediction. The Ensembl<sup>6</sup> database of metazoan genome annotation is a good example of such a resource. The contents of this type of database are subject to frequent revision, especially when the sequence has only recently been deciphered, but they often provide more complete and up-to-date annotation than the traditional archives.

## Protein databases

Most efforts at recording functional information have occurred in protein databases. Swiss-Prot, a subsection of the UniProt Knowledgebase<sup>7</sup> (UniProtKB) that is manually curated from the scientific literature, is widely considered to be the richest and most reliable store of functional information, and is frequently used to infer information about proteins not yet subjected to specific experimental study. The world-wide Protein Data Bank<sup>2</sup> functions as an archive for protein structures and information about posttranslational modifications can be found in the RESID database<sup>8</sup>

## Other bioinformatics databases

Several hundred bioinformatics databases exist in total, with many resources focusing on particular taxonomic or methodological domains. For example, there are specific databases for many model organisms

**BOX 2 PROBLEMS IN LINKING BIOINFORMATICS DATA TO LITERATURE**

In theory, data warehousing technologies would seem to support the automatic generation of a single integrated data resource from independent, but cross-referenced, databases. However, there are certain inherent limitations to this approach. Accurate cross-references are vital for good results, but the volume of data now extant exceeds the capacity for manual curation of almost all bioinformatics databases — for example, about 220,000 records (7% of the total) in the UniProtKB are manually curated. Automatic methods for generating cross-references can be used instead, through the tracking of identifiers (as data is transferred between resources), or by a comparison of the properties (such as the sequence) of entities to establish their equivalence. However, these methods may not always produce the correct answers. Different databases may maintain different identifiers, and use different names, for the same biological object (synonymy and homonymy are particularly problematic when searching literature, in which references to genes and proteins occur in free text and do not necessarily follow well-defined standards). And the 'same' object may be assigned different properties in different resources.

Other problems may be encountered even when cross-references are well-maintained: first, the records in different databases may not be conceptually compatible. Numerous definitions exist for words like protein and gene, and their use is not necessarily standardised between resources. Second, different resources may differ in content, that is, the data they choose to include may not be the same. Third, the huge quantity of experimentally generated data, and the generally limited frequency at which submissions are updated, means that archive databases may contain much information that can be considered redundant or out-of-date. Fourth, the use of different standards and syntax for annotation is a further barrier to integration.

(such as bakers' yeast<sup>9</sup>, worm<sup>10</sup>, fruit fly<sup>11</sup> and mouse<sup>12</sup>), which are often the most detailed source of information about gene and protein function within their taxonomic scope. Other resources concentrate on particular families of (protein and functional RNA-encoding) genes.

**Literature databases**

The scientific literature itself can be accessed through computerized databases. The largest primary database containing biological literature is MEDLINE, which contains over 10 million citations. MEDLINE is frequently accessed as part of PubMed ([www.ncbi.nih.gov/entrez/query/static/overview.html](http://www.ncbi.nih.gov/entrez/query/static/overview.html)), a larger collection of literature databases. Many publishers are now making the full text of articles publicly available from their own archives over the World Wide Web, frequently using digital object identifiers (DOIs; [www.doi.com](http://www.doi.com)), which provide a common mechanism for data access and many full-text articles are also freely available in the PubMed Central archive<sup>13</sup>.

The simplest way that bioinformatics and literature databases can be linked is through the use of the cross-references, which many of these resources maintain to each other. The UniProtKB, for example, maintains cross-references to over 100 other resources (including MEDLINE and PubMed) in which more detailed information on particular aspects of a protein's function can be found. These cross-references are widely exploited by data warehouses (see Box 1).

**PUBLIC RESOURCES FOR INTEGRATED BIOINFORMATICS**

Many of the main bioinformatics databases are made available by a small number of institutions whose mission is to provide services to the scientific community (such as the European Bioinformatics Institute<sup>14</sup> (EBI) in the UK, and the National Centre for Biotechnology Information<sup>15</sup> (NCBI) in the USA). These resources provide an obvious first point of call for anyone attempting to integrate bioinformatics data.

The EBI maintains a large data warehouse that contains over 100 bioinformatics databases<sup>16</sup> using the SRS data warehousing system<sup>17</sup>. The NCBI Entrez server<sup>18</sup> offers similar functionality (see Box 2). However, there is a difference in the focus of the two systems: the interface to SRS is powerful but complex, allowing the construction

of inter-database queries in a generic manner (for example, "get me the DNA sequences encoding proteins of known structure in these species"). Entrez offers a simpler interface, with less support for structured queries, but provides rapid retrieval of data of all types linked to a given search term (for example, "get me all the genes, proteins, genomes and literature that relate to ras1"). Both warehouses contain not only gene and protein databases, but also literature databases (for example, PubMed is incorporated into Entrez), allowing direct inter-linking between these resources.

Ultimately, a warehouse's usefulness is dependent on the coherence of the data it contains (see Box 2). Approaches to improve the coherence among biological databases include the establishment of international collaborations to share data, the use of common controlled vocabularies (such as the Gene Ontology<sup>19</sup>, GO, a hierarchical controlled vocabulary for describing the function of gene products) by different resources, and the wider development of standards for data representation (particularly in new areas of research, such as transcriptomics<sup>20</sup> and proteomics<sup>21</sup>), where the data complexity can be orders of magnitude greater than that of sequence data. Another approach is the production of well-annotated, non-redundant subsets of the complete data and the development of services based on these. The EBI's Integr8 project<sup>22</sup> provides reference data sets and analysis tools for species with completely deciphered genomes from which redundancy has been removed. Annotation in Integr8 is enhanced, updated and corrected (with respect to the original data submissions) through the transfer of data from actively curated resources and the performance of new computational analysis. The re-annotated genomes are available in EMBL-like format as Genome Reviews. The NCBI's RefSeq database<sup>23</sup> similarly provides comprehensive, non-redundant sequences for genomes, genes and proteins. The related resource Gene<sup>24</sup> provides a framework (a defined gene set) against which sequences from RefSeq and other resources are cross-referenced, and includes two-way links to literature in PubMed. Searching against resources such as these may provide clearer answers than those obtained when using the primary repositories, in which data redundancy, and non-standard or out-of-date annotation, may obscure the correct results.



```
atgcagggaataaatacaactataagagagataaagatagtagtggctgggggaggtggc
gttggtaaatctgccttaacaattcaattcaatcaatcactcttgggaagaaatagac
cctactatcgaaagattcttaacagaaaacagttgbcactgatgacaaagatccatctt
gacatcttagatactgctggacaagaaagatctctgagtagagagaacagatcaatgag
actggggaaggttctactggtctatctcgtcaacctatagaaatctcttggatgagtta
ctgtcttattatcagcaaatcaaaagatcaaaagattctgactacatctctgtagctgt
gtaggttaacaaattggacottgaaaatgaaagcaagctctctatgaagacgggttao
ctggccaagcagttgaatgcaacctctttagaaacgtctgogaacaagocataacgta
gacgagccttttatagccttattcgtttggttaagggaogacoggtgggaaatcaaatg
atgaaatcgtcaactggatcaatcgaatgaaatgaaagattcggagctaaacctcactgca
acagcggatatagaaaaaagaaacacgggtcttattgactcgaatctcttgaccaat
gctggcactggctccagtcaaaagtcagcogttaaocataaaggtgaaactactaaaag
actgataaaaagaattacgcttaacaaacaaatacaaaagaaagaaatacaagtaactcc
agtacggcaacggaaatcgaagtgatattagctcgtgtaatacaaaataatgcttaaat
tcgagaagtcaacagctctgctgagccacaataaaatcaaacgocaaacogctagaaagtt
ctactggtggtggtgataaattggtgaa
```

```
MQGNKSTIREYKIVVVGGGVGSALTIQFIQSYFVDEYDPTIEDSYRKQVVIDDKVSI
LDLDTAGQEEYSAMREQYMRIGEFLLVYSVTSRNSFDLLESLYQQIQRVKSDSYIPVVV
VGNKLDLENERQVSYEDGLRLAKQLNPFLETSAKQAINVDEAFYSLIRLVRDDGGKYN
MNRQLDNTNIRDSLETSSATADTEKKNNSYVLDNLSLTNAGTGSSSKSAVNHNGETTKR
TDEKNYVNVQNNNEGNKYVSNGNGNRSDISRGNQNNALNSRSKQSAEPQKNSANARKV
SSGGCCIIIC
```

Sequence	1	309
1 UNIPROT <a href="#">RAS1 YEAST</a>	1:309	1:309
2 UNIPROT <a href="#">RAS2 YEAST</a>	4:308	3:320
3 UNIPROT <a href="#">Q6FKM8 CANGA</a>	4:303	2:305
4 UNIPROT <a href="#">Q75B39 ASHGO</a>	5:308	5:269
5 UNIPROT <a href="#">Q9UWU4 CANAL</a>	8:308	2:288
6 UNIPROT <a href="#">RAS1 CANAL</a>	8:308	2:289
7 UNIPROT <a href="#">Q5XU5 CANAL</a>	8:308	2:290
8 UNIPROT <a href="#">Q6K00 DEBHA</a>	8:263	2:250

```
MQGNKSTIREYKIVVVGGGVGSALTIQFIQSYFVDEYDPTIEDSYRKQVVIDDKVSI
LDLDTAGQEEYSAMREQYMRIGEFLLVYSVTSRNSFDLLESLYQQIQRVKSDSYIPVVV
VGNKLDLENERQVSYEDGLRLAKQLNPFLETSAKQAINVDEAFYSLIRLVRDDGGKYN
MNRQLDNTNIRDSLETSSATADTEKKNNSYVLDNLSLTNAGTGSSSKSAVNHNGETTKR
TDEKNYVNVQNNNEGNKYVSNGNGNRSDISRGNQNNALNSRSKQSAEPQKNSANARKV
SSGGCCIIIC
```

SEQUENCE: Sequence\_1 CRC64: 17169CD5DBC8F4E0 LENGTH: 309 aa

InterPro IPR01804 Family	Ras GTPase PR00449	RASTRNSFRMNG
InterPro IPR02577 Family	Ras small GTPase, Ras type SM00173	no description
InterPro IPR02525 Domain	Small GTP-binding protein domain TIGR0221	small_GTP: small GTP-binding protein domain
InterPro IPR013753 Family	Ras PF00071	Ras
noIPR unintegrated	G3D:3.40.50.300 PTHR11708	no description RAS-RELATED GTPASE

**References**

- NUCLEOTIDE SEQUENCE [GENOMIC DNA]. DOI=10.1016/0092-8674(84)90340-4; MEDLINE=84130171; PubMed=6365329; [NCBI, ExPASy, EBI, Israel, Japan] Powers S., Kataoka T., Fasano O., Goldfarb M., Strathern J., "Genes in *S. cerevisiae* encoding proteins with domains homologous to the mammalian ras proteins."; Cell 36:607-612(1984).
- NUCLEOTIDE SEQUENCE [GENOMIC DNA]. MEDLINE=84221383; PubMed=6328429; [NCBI, ExPASy, EBI, Israel, Japan] Dhar R., Nieto A., Koller R., DeFeo-Jones D., Scolnick E.M., "Nucleotide sequence of two rasH-related genes isolated from the yeast *Saccharomyces cerevisiae*."; Nucleic Acids Res. 12:3611-3618(1984).
- NUCLEOTIDE SEQUENCE [GENOMIC DNA]. DOI=10.1002/(SICI)1097-0061(19970615)13:7<655::AID-HEA120>3.0.CO;2-I

Characterization of *Saccharomyces cerevisiae* **Ras1** and chimeric constructs of **Ras** proteins reveals the hypermutable region and farnesylation as critical elements in the **adenylyl cyclase** signaling pathway. This was found to result from the higher expression of **Ras1** and **Ras1-Ras2**, which compensate for their lower efficacy in activating **adenylyl cyclase**. We have attempted to identify amino acid residues of the yeast **adenylyl cyclase** that are involved in the regulation of its activity, by isolating adenylyl cyclase-limited spontaneous mutations capable of suppressing the temperature-sensitive phenotype of **ras1**-**ras2**-**ras1** strains. In this pathway **Ras1** interacts with the **protein kinase** **By2** and leads to its activation in conjunction with a signal from the receptor-coupled, heterotrimeric G protein. The single **Ras** homologue (**Ras1**) of *S. pombe* regulates two distinct processes: (1) Signal transduction through a MAP kinase-like **protein kinase** cascade in response to mating pheromones. The **ras1-15** mutation does not alter the level of **RAS1** mRNA in cells grown on glucose. **Carbon** source regulation of **RAS1** expression in *Saccharomyces cerevisiae* and the phenotypes of **ras2** cells. The amount of **RAS1** mRNA is significantly repressed in cultures grown on the nonfermentable carbon sources ethanol and acetate. The **ras1** mutant strain was also assayed in an animal model of cryptococcal meningitis. The phosphorylation of **RAS1** protein is demonstrated by treating with alkaline phosphatases as well as by labeling with [32P]orthophosphate. To test whether **Ras** function depends on a cytosolic factor such as GAP, we microinjected into *Xenopus oocytes* a form of *Saccharomyces cerevisiae* **RAS1** (Leu59**RAS1** terminated at residue 182, called [Leu59**RAS1**term]) that lacks the consensus membrane localization site, does not respond to GAP in a GTPase assay, but binds to GAP 100-fold more tightly than [Val12**Ras**: [Leu48**RAS1**term] alone did not stimulate inositol geminal-oxide breakdown. **RAS1** mRNA levels and protein synthesis are very low at all stages of growth when ethanol rather than glucose is provided as the sole carbon source. The genes studied include the normal cellular and **bladder cancer** ras genes, recombinant viral/cellular ras genes, recombinant yeast/mammalian ras genes, and a constructed gene with yeast **RAS1** sequences significantly modified by deletions and an oncogenic mutation.

**Figure 3** From gene to protein to literature: a sample analysis. Beginning with a gene sequence (top left), a protein sequence (below gene sequence) is generated and can be analysed with InterProScan to detect sequences indicative of the presence of known protein domains and families within the sequence, generating a graphical overview (top right). The sequence is also compared with known sequences in the UniProtKB database using the

BLAST algorithm and the output displays an alignment (middle) between the query sequence and its best matches. The best match for the sample query data is the **RAS1** protein from *Saccharomyces cerevisiae*. The corresponding UniProtKB entry features curated links for this protein to known publications (lower left). Additional publications for this protein, retrieved using text-mining methods, can be found by querying iHOP (lower right).

**SMALL-SCALE DATA MINING**

Data mining refers to two phenomena: the analysis of large-scale data sets for the purpose of general inference (discussed in the next section) and the extraction of specific information that relates to some initial data of interest. If that initial data is a gene or protein name or ID, querying resources such as UniProtKB (for protein-centric queries), Integr8 (to obtain information in a genomic context) or SRS (for complex filtering over many resources) will usually identify a database record that contains much of the information that is already known about that entity. However, if the data consists of nucleotide or protein sequence, a comparison needs to be made to known sequences to identify similar or identical molecules that have already been annotated. Common algorithms used for this purpose include BLAST<sup>25</sup>, FASTA<sup>26</sup> and Smith-Waterman<sup>27</sup>, and are available through the websites and web services of the EBI, NCBI and other bioinformatics service providers.

**Analysing sequence**

High overall sequence similarity may be a good indicator of functional equivalence, but aspects of a protein's function may be inferred from

the presence of particular domains, even if a protein's overall architecture has not been reported or the complete function of closely matching sequences is not known. Several methods exist for domain identification in sequence, mostly using hidden Markov models<sup>28</sup>. InterPro<sup>29</sup> is a curated, integrative resource that combines methods for domain identification from 15 different member databases, in which redundant methods are merged and common annotation attached. InterProScan<sup>30</sup>, the programme which applies these methods (online or as a local installation), is among the most powerful tools for characterizing unknown protein sequence. How InterPro entries are assembled and how a scan is applied to unknown sequence is illustrated in Fig. 1.

Predicting protein structure is more problematic, as the structures of most proteins remain unknown. Sequence similarity can be used to infer structure, providing the contribution of individual residues within the sequence to the structure is taken into account (and not just the overall degree of similarity). Tools such as SWISS-MODEL<sup>31</sup> generate a probable structure for a query sequence by identifying a similar sequence of known structures that can serve as a model. Certain structural features can also be predicted by direct sequence analysis. Programmes

that perform this type of analysis include TMHMM<sup>32</sup>, which predicts trans-membrane helices, and GenThreader<sup>33</sup> (available through the PSIPRED server<sup>34</sup>), a tool for predicting protein folds. However, prediction of higher-order structural features from sequence is more difficult than predicting secondary structure or domain composition, especially for proteins not related to those with known structures.

### Literature mining

In Swiss-Prot, Entrez Gene and other well-curated databases, direct links exist from individual records to relevant publications. However, references to many papers relevant to a gene or protein are not likely to have been directly curated. Further literature can be found by directly searching literature databases with the name(s) associated with the protein, although care must be taken to avoid confusion caused by the existence of synonyms and homonyms. Abstracts in MEDLINE are labelled with terms from a controlled vocabulary (medical subject headings; MeSH<sup>35</sup>), and searching with relevant MeSH terms can be used to identify papers of interest. Unfortunately, most bioinformatics database records are not directly annotated with MeSH terms; instead, they are often annotated with terms from the GO vocabulary. GO is more tightly focused than MeSH — each contains approximately 20,000 terms, but MeSH is split into 16 principal subsections and covers geography and sociology in addition to biology, whereas GO only covers three specific aspects of biology. A resource to translate GO terms automatically into their equivalent MeSH terms is still in development.

Increasingly, natural language programming (NLP) techniques are being used to mine literature databases automatically, thus providing a more sophisticated way of identifying relevant publications. Programs using NLP can typically be rerun with greater frequency than a record can be revisited by a curator. These techniques are limited in their ability to accurately summarise the complete information contained in a single paper, but can effectively identify papers worthy of subsequent manual inspection. iHOP<sup>36</sup> is a web interface that provides access to papers that have been identified using natural language programming as relevant to genes from a number of sample genomes. This interface displays the specific text that was used to identify each paper as relevant, thereby allowing false positives to be quickly discarded.

An outline workflow that might be used in characterizing an unknown sequence is shown in Fig. 2 and a specific example is given in Fig. 3.

### INFORMATION EXTRACTION FROM LARGE DATA SETS

Data mining for large data sets can be quite different to extracting data about individual sequences. The types of analysis performed are generally similar, but the development of automated procedures is usually essential as the data volume is increased. Many public bioinformatics servers place limits on the amount of data that can be requested by users in a single query, although the underlying software is often also available for local installation. Another route by which service providers are increasingly attempting to meet user demands for bulk data analysis is through the provision of programmatic access over the internet, with the widespread implementation of services meeting the Web Service standards proposed by the World Wide Web consortium (<http://www.w3.org/2002/ws>).

For scientists who have to analyse large quantities of data, but who lack programming expertise, the use of a workflow management tool may provide a solution. Taverna<sup>37</sup> is designed explicitly for use in the context of bioinformatics data, providing a graphical user interface for the assembly of multiple services, potentially running at diverse locations, into a single data processing pipeline.

Large data sets enable 'knowledge discovery' through the identification of patterns within the data. For example, individual pieces of protein interaction data may be combined to produce a local or global data interaction network. One such tool, the IntAct Hierarch view application<sup>38</sup>, also allows such information to be overlaid with patterns of GO annotation, which in turn can enable the validation and interpretation of the results. Another example of large-scale data mining can be found on the UniProt website<sup>7</sup>, where statistical patterns in curated data sets have been used to apply annotation to non-curated sequences. Statistical predictions can also be directly tested against the actual annotation in curated records, allowing each type of data to be validated against the other.

### THE FUTURE OF BIOINFORMATICS DATA AND BIOINFORMATICS DATA SERVICES

The importance of automatic methods for linking gene, protein and literature data has grown as the number of known sequences has increased and the proportion of manually curated database records has fallen. As new, increasingly large-scale, experimental techniques continue to be developed, and the number of specialised databases holding data from their application grows in parallel, the ability to perform distributed

#### BOX 3 TOWARDS SYSTEMATIC RECONSTRUCTION OF WHOLE ORGANISM BIOLOGY

In addition to generic data warehousing models such as SRS and Entrez, an increasing range of more specialist data warehousing tools have been developed by groups specifically aiming to integrate data for particular species. The proliferation of different data models and user interfaces can be a problem for users. Therefore, an important development is the emergence of reusable frameworks (applicable to multiple species) such as GMOD39 (the generic model organism database project), GUS40 (the genomics unified schema), and BioCyc41 (a collection of 205 species-specific metabolic pathway databases). BioCyc is particularly interesting — for its content (a partial reconstruction of cellular metabolomes), but also for its mode of generation, whereby databases for most species are inferred computationally from a combination of genome sequence and a curated database of pathway information largely derived from well-studied species (such as *Escherichia coli*). As the number of biological systems, and the attributes of these systems, continues to increase, it is likely that only a decreasing portion of the generated data will be subject to individual experimental verification and subsequent curation into databases. Therefore, the approximate reconstruction of less-characterized systems by inference from well-characterized models will become an increasingly important factor in the interpretation of large-scale experiments.

queries over resources with separate physical locations is also growing in importance, as there is a limit to the size and complexity at which integrated data warehouses can be efficiently maintained. This ability is also necessary if local, and potentially private, data is to be integrated with information available from public resources. New bioinformatics warehousing systems, such as BioMart<sup>42</sup>, have been designed with the explicit aim of supporting queries that can be executed across differently located resources, but ensuring the efficient execution of such queries is inherently difficult.

An alternative to accessing a single data warehouse is to retrieve data from multiple sources and merge this as needed. Projects such as MyGrid<sup>43</sup> are developing software to provide solutions to generic problems in analysing bioinformatics data distributed across many resources. However, it remains difficult to integrate data from different resources. For example, to write a programme to combine data from Entrez and the EBI SRS server requires domain-specific expertise to retrieve and parse the data from each source, and additional expertise

to merge the data correctly. These are common problems for programmers attempting to mine information available on the internet, not just those working with biological data, and there is growing interest among bioinformaticians in adopting the emerging 'semantic web' technologies<sup>44</sup> designed to support distributed computing over a wider range of domains. The key concept of the semantic web is the publication of self-describing data, that is, data published together with the metadata that describes it. If such publication accords to standardised protocols, and if the descriptions themselves are machine-readable, standardised and shared across resources, then a programmer should be able to retrieve and integrate distributed data by specifying a logical request to a semantically aware search engine without needing specific knowledge about the peculiarities of individual sources. This model is particularly attractive, not only due to the highly dispersed nature of much bioinformatics data, but also because the descriptive standards necessary to make it work have already been developed for particular sub-domains, such as microarray and molecular interaction data.

**Table 1** Glossary

<b>Controlled vocabulary</b>	A limited range of terms used in bioinformatics databases to ensure standardised annotation (as opposed to free text). If terms are linked by defined relationships, a controlled vocabulary may constitute an ontology.
<b>Data warehouse</b>	A database constructed to support efficient querying of the data it contains, possibly constructed from many primary data sources.
<b>Distributed queries</b>	A query addressed, not to a single database, but to multiple data resources at potentially disparate locations on a network. Distributed queries work on distributed data, and are one application in the wider field of distributed computing.
<b>Digital object identifier</b>	A standard way for representing a piece of intellectual property expressed in a digital environment (for example, an online journal article).
<b>GO</b>	Gene ontology; a structured controlled vocabulary widely used for the annotation of gene products.
<b>MeSH</b>	Medical subject headings; a structured controlled vocabulary used in MEDLINE to summarise the contents of scientific articles.
<b>Hidden markov models</b>	HMMs; a statistical modelling technique that can be applied to alignments of sequences and used to identify other sequences with similar properties to those in the alignment.
<b>Homonymy</b>	The condition where two different entities share a common name (often true of genes).
<b>Natural language programming</b>	NLP; the discipline of computer science that attempts to understand the meaning of natural language. Widely used to automate processing of scientific literature, sometimes accurately.
<b>-ome, -omics</b>	Suffixes indicating the complete data (and the study of that data) of a particular type (indicated by the root) from particular biological systems (for example, genome/genomics, proteome/proteomics).
<b>Parsing</b>	The process of analysing the contents of a document (such as a flat-file database entry) in accordance with a set of rules defining its structure.
<b>Ontology</b>	In computer science, a model of a particular domain, amounting to a logical description of a set of related concepts.
<b>Semantic web</b>	A network of interlinked data available on the internet, capable of being understood by computer software; also used to refer to technologies associated with the development of such a network.
<b>SRS</b>	Sequence retrieval system; a data warehousing system widely used in bioinformatics, originally developed at the EBI and now maintained by BioWisdom Ltd.
<b>Synonymy</b>	The condition where a single entity has multiple names (often true of genes)
<b>Web services</b>	Programmatic interfaces available for remote computation, using the protocols associated with the world wide web; one potential component of the semantic web.



A significant challenge for the bioinformatics community is to continue to maintain and promote these standards, and to develop additional standards for data types such as small RNAs, where suitable representations do not yet exist, at a sufficient pace to match the ever-growing diversity and quantity of data. □

Note: Supplementary Information is available on the Nature Cell Biology website.

#### COMPETING FINANCIAL INTERESTS

The authors declare that they have no competing financial interests.

- Apweiler, R., Bairoch, A. & Wu, C. H. Protein sequence databases. *Curr. Opin. Chem. Biol.* **8**, 76–80 (2004).
- Berman, H., Henrick, K. & Nakamura, H. Announcing the worldwide Protein Data Bank. *Nature Struct. Biol.* **10**, 980 (2003).
- Cochrane, G. *et al.* EMBL Nucleotide Sequence Database: developments in 2005. *Nucleic Acids Res.* **34**, D10–D15 (2006).
- Benson, D. A., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J. & Wheeler, D. L. GenBank. *Nucleic Acids Res.* **34**, D16–D20 (2006).
- Okubo, K., Sugawara, H., Gojobori, T. & Tateno, Y. DDBJ in preparation for overview of research activities behind data submissions. *Nucleic Acids Res.* **34**, D6–9 (2006).
- Birney, E. *et al.* Ensembl 2006. *Nucleic Acids Res.* **34**, D556–D561 (2006).
- Wu, C. H. *et al.* The Universal Protein Resource (UniProt): an expanding universe of protein information. *Nucleic Acids Res.* **34**, D187–D191 (2006).
- Garavelli, J. S. The RESID Database of Protein Modifications as a resource and annotation tool. *Proteomics* **4**, 1527–1533 (2004).
- Christie, K. R. *et al.* Saccharomyces Genome Database (SGD) provides tools to identify and analyze sequences from *Saccharomyces cerevisiae* and related sequences from other organisms. *Nucleic Acids Res.* **32**, D311–D314 (2004).
- Schwarz, E. M. *et al.* WormBase: better software, richer content. *Nucleic Acids Res.* **34**, D475–D478 (2006).
- Grumblin, G. & Strelets, V. FlyBase: anatomical data, images and queries. *Nucleic Acids Res.* **34**, D484–D488 (2006).
- Blake, J. A., Eppig, J. T., Bult, C. J., Kadin, J. A. & Richardson, J. E. The Mouse Genome Database (MGD): updates and enhancements. *Nucleic Acids Res.* **34**, D562–D567 (2006).
- Sequeira, E., McEntyre, J. & Lipman, D. PubMed Central decentralized. *Nature* **410**, 740 (2001).
- Lopez, R., Duggan, K., Harte, N. & Kibria, A. Public services from the European Bioinformatics Institute. *Brief Bioinform.* **4**, 332–340 (2003).
- Jenuth, J. P. The NCBI. Publicly available tools and resources on the Web. *Methods Mol. Biol.* **132**, 301–312 (2000).
- Zdobnov, E. M., Lopez, R., Apweiler, R. & Etzold, T. The EBI SRS server-new features. *Bioinformatics* **18**, 1149–1150 (2002).
- Etzold, T., Ulyanov, A. & Argos, P. SRS: information retrieval system for molecular biology data banks. *Methods Enzymol.* **266**, 114–128 (1996).
- Geer, R. C. & Sayers, E. W. Entrez: making use of its power. *Brief Bioinform.* **4**, 179–184 (2003).
- Gene Ontology Consortium. The Gene Ontology (GO) project in 2006. *Nucleic Acids Res.* **34**, D322–D326 (2006).
- Whetzel, P. L. *et al.* The MGED Ontology: a resource for semantics-based description of microarray experiments. *Bioinformatics* **22**, 866–873 (2006).
- Orchard, S. *et al.* Autumn 2005 Workshop of the Human Proteome Organisation Proteomics Standards Initiative (HUPO-PSI) Geneva, September, 4–6, 2005. *Proteomics* **6**, 738–741 (2006).
- Kersey, P. *et al.* Integr8 and Genome Reviews: integrated views of complete genomes and proteomes. *Nucleic Acids Res.* **33**, D297–D302 (2005).
- Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI Reference Sequence (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* **33**, D501–D504 (2005).
- Maglott, D., Ostell, J., Pruitt, K. D. & Tatusova, T. Entrez Gene: gene-centered information at NCBI. *Nucleic Acids Res.* **33**, D54–D58 (2005).
- McGinnis, S. & Madden, T. L. BLAST: at the core of a powerful and diverse set of sequence analysis tools. *Nucleic Acids Res.* **32**, W20–W25 (2004).
- Pearson, W. R. Using the FASTA program to search protein and DNA sequence databases. *Methods Mol. Biol.* **24**, 307–331 (1994).
- Smith, T. F. & Waterman, M. S. Identification of common molecular subsequences. *J. Mol. Biol.* **147**, 195–197 (1981).
- Eddy, S. R. What is a hidden Markov model? *Nature Biotechnol.* **22**, 1315–1316 (2004).
- Mulder, N. J. *et al.* InterPro, progress and status in 2005. *Nucleic Acids Res.* **33**, D201–D205 (2005).
- Quevillon, E. *et al.* InterProScan: protein domains identifier. *Nucleic Acids Res.* **33**, W116–W120 (2005).
- Kopp, J. & Schwede, T. The SWISS-MODEL Repository: new features and functionalities. *Nucleic Acids Res.* **34**, D315–D318 (2006).
- Sonnhammer, E. L., von Heijne, G. & Krogh, A. A hidden Markov model for predicting transmembrane helices in protein sequences. *Proc. Int. Conf. Intell. Syst. Mol. Biol.* **6**, 175–182 (1998).
- McGuffin, L. J. & Jones, D. T. Improvement of the GenTHREADER method for genomic fold recognition. *Bioinformatics* **19**, 874–881 (2003).
- McGuffin, L. J., Bryson, K. & Jones, D. T. The PSIPRED protein structure prediction server. *Bioinformatics* **16**, 404–405 (2000).
- Nelson, S. J., Schopen, M., Savage, A. G., Schulman, J. L. & Arluk, N. The MeSH translation maintenance system: structure, interface design, and implementation. *Medinfo* **11**, 67–69 (2004).
- Hoffmann, R. & Valencia, A. Implementing the iHOP concept for navigation of biomedical literature. *Bioinformatics* **21**, ii252–ii258 (2005).
- Oinn, T. *et al.* Taverna: a tool for the composition and enactment of bioinformatics workflows. *Bioinformatics* **20**, 3045–3054 (2004).
- Hermjakob, H. *et al.* IntAct: an open source molecular interaction database. *Nucleic Acids Res.* **32**, D452–D455 (2004).
- Stein, L. D. *et al.* The generic genome browser: a building block for a model organism system database. *Genome Res.* **12**, 1599–1610 (2002).
- Davidson, S. B. *et al.* K2/Kleisli and GUS: Experiments in integrated access to genomic data sources. *IBM Systems Journal* **40**, 512–531 (2001).
- Karp, P. D. *et al.* Expansion of the BioCyc collection of pathway/genome databases to 160 genomes. *Nucleic Acids Res.* **33**, 6083–6089 (2005).
- Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics* **21**, 3439–3440 (2005).
- Stevens, R. D., Robinson, A. J. & Goble, C. A. myGrid: personalised bioinformatics on the information grid. *Bioinformatics* **19**, i302–i304 (2003).
- Berners-Lee, T. & Hendler, J. Publishing on the semantic web. *Nature* **410**, 1023–1024 (2001).

**Table S1** Some important bioinformatics resources and their applications

<b>Institutions</b>		
<b>Name</b>	<b>URL</b>	<b>Description</b>
European Bioinformatics Institute (EBI)	<a href="http://www.ebi.ac.uk">www.ebi.ac.uk</a>	European Centre for research and services in bioinformatics, part of EMBL (European Molecular Biology Laboratory)
National Center for Biotechnology Information	<a href="http://www.ncbi.nih.gov">www.ncbi.nih.gov</a>	Centre for bioinformatics research and services in the United States, part of the National Institutes of Health.
<b>Primary Databases</b>		
<b>Name</b>	<b>URL</b>	<b>Description</b>
International Nucleotide Sequence Database Collaboration (INSDC)	<a href="http://www.insdc.org">www.insdc.org</a>	Collaboration that maintains archival repository of nucleotide sequence data
EMBL Nucleotide Sequence Database	<a href="http://www.ebi.ac.uk/embl">www.ebi.ac.uk/embl</a>	European partner in the INSDC
Genbank	<a href="http://www.ncbi.nih.gov/Genbank">www.ncbi.nih.gov/Genbank</a>	American partner in the INSDC
DNA Database of Japan (DDBJ)		Japanese partner in the INSDC
UniProt Knowledgebase (UniProtKB)	<a href="http://www.uniprot.org">www.uniprot.org</a>	Central access point for protein-related information maintained by the UniProt consortium
Swiss-Prot		Manually curated portion of the UniProtKB, commonly regarded as a gold standard for functional annotation
Gene	<a href="http://www.ncbi.nih.gov/query.fcgi?db=gene">www.ncbi.nih.gov/query.fcgi?db=gene</a>	Database of known genes maintained by NCBI. Provides reference system for NCBI annotation efforts.
RefSeq	<a href="http://www.ncbi.nih.gov/RefSeq">www.ncbi.nih.gov/RefSeq</a>	Database of reference DNA and protein sequences maintained by NCBI. Closely linked to Gene.
Ensembl	<a href="http://www.ensembl.org">www.ensembl.org</a>	Metazoan genome annotation database
Saccharomyces Genome Database (SGD)	<a href="http://www.yeastgenome.org">www.yeastgenome.org</a>	<i>Saccharomyces</i> genome database
Mouse Genome Informatics (MGI)	<a href="http://www.informatics.jax.org">www.informatics.jax.org</a>	Mouse genome database
FlyBase	<a href="http://flybase.bio.indiana.edu">flybase.bio.indiana.edu</a>	<i>Drosophila</i> genome database
WormBase	<a href="http://www.wormbase.org">www.wormbase.org</a>	<i>Caenorhabditis</i> genome database
GOA	<a href="http://www.ebi.ac.uk/goa">www.ebi.ac.uk/goa</a>	Database of functional annotation, part-curated, part derived from InterPro analysis
IntAct	<a href="http://www.ebi.ac.uk">www.ebi.ac.uk</a>	Database of protein interactions
BioCyc	<a href="http://www.biocyc.org">www.biocyc.org</a>	Collection of pathway and genome databases
World-wide Protein Data Bank (wwPDB)	<a href="http://www.wwpdb.org">www.wwpdb.org</a>	Repository for protein structure information
ArrayExpress	<a href="http://www.ebi.ac.uk/arrayexpress">www.ebi.ac.uk/arrayexpress</a>	European repository for microarray data
Gene Expression Omnibus (GEO)	<a href="http://www.ncbi.nlm.nih.gov/geo">www.ncbi.nlm.nih.gov/geo</a>	American repository for microarray data
MEDLINE	<a href="http://www.nlm.nih.gov/pubs/factsheets/jsel.html">www.nlm.nih.gov/pubs/factsheets/jsel.html</a>	Database of information on biomedical literature



PubMed	<a href="http://www.ncbi.nih.gov/entrez/query/static/overview.html">www.ncbi.nih.gov/entrez/query/static/overview.html</a>	Expanded version of MEDLINE, searchable through Entrez at NCBI
<b>Sequence Analysis</b>		
<b>Name</b>	<b>URL</b>	<b>Description</b>
Smith-Waterman (SW)	Programs running most sequence analysis algorithms are available from many bioinformatics service providers.	Dynamic programming algorithm for sequence analysis, widely considered good, but slow. Various hardware-optimised approximations to SW also exist.
Basic Local Alignment Search Tool (BLAST)		High speed algorithm for sequence comparison.
FAST-AII (FASTA)		Another widely used sequence comparison algorithm
InterProScan	<a href="http://www.ebi.ac.uk/InterProScan">www.ebi.ac.uk/InterProScan</a>	Domain analysis tool incorporating algorithms and annotations from 16 consortium members. Results for known sequences stored in InterPro database.
TMHMM	<a href="http://www.cbs.dtu.dk/services/TMHMM-2.0">www.cbs.dtu.dk/services/TMHMM-2.0</a>	Prediction of transmembrane helices in proteins
SWISS-MODEL	<a href="http://swissmodel.expasy.org">swissmodel.expasy.org</a>	Automated protein modelling server
PSIPRED	<a href="http://bioinf.cs.ucl.ac.uk/psipred">bioinf.cs.ucl.ac.uk/psipred</a>	Protein structure prediction server, offers access to the GenTHREADER program
<b>Warehouses and web portals</b>		
EBI SRS server	<a href="http://srs.ebi.ac.uk">srs.ebi.ac.uk</a>	Primary data warehouse/query engine maintained by EBI
NCBI Entrez server	<a href="http://www.ncbi.nih.gov/database/index.html">www.ncbi.nih.gov/database/index.html</a>	Primary data warehouse/query engine maintained by NCBI
Integr8	<a href="http://www.ebi.ac.uk/integr8">www.ebi.ac.uk/integr8</a>	Provides structured overview of genes and genomes
iHOP	<a href="http://www.ihop-net.org/UniPub/iHOP">www.ihop-net.org/UniPub/iHOP</a>	Portal for assessing automatically retrieved literature linked to genes and proteins
<b>Generic database schemas</b>		
GMOD	<a href="http://www.gmod.org">www.gmod.org</a>	Schema maintained by several model organism databases
GUS	<a href="http://www.gusdb.org">www.gusdb.org</a>	Schema for functional genomics data
<b>Integration tools</b>		
BioMart	<a href="http://www.biomart.org">www.biomart.org</a>	Data warehousing system designed for easy local install, interface generation and distributed querying
Taverna	<a href="http://taverna.sourceforge.net">taverna.sourceforge.net</a>	Workflow manager for bioinformatics. Part of the MyGRID project.
MyGRID	<a href="http://mygrid.org.uk">mygrid.org.uk</a>	Project developing a suite of tools for distributed computing in bioinformatics