

ONTOLOGIES FOR MOLECULAR BIOLOGY

S. Schulze-Kremer

*Max-Planck-Institute for Molecular Genetics, Ihnestraße 73, Dept. Lehrach,
D-14195 Berlin, Germany, email steffen@chemie.fu-berlin.de,
<http://mycroft.rz-berlin.mpg.de/~steffen>*

Molecular biology has a communication problem. There are many databases using their own labels and categories for storing data objects and some using identical labels and categories but with a different meaning. A prominent example is the concept “gene” which is used with different semantics by major international genomic databases. Ontologies are one means to provide a semantic repository to systematically order relevant concepts in molecular biology and to bridge the different notions in various databases by explicitly specifying the meaning of and relation between the fundamental concepts in an application domain. Here, the upper level and a database branch of a prospective ontology for molecular biology (OMB) is presented and compared to other ontologies with respect to suitability for molecular biology (<http://igd.rz-berlin.mpg.de/~www/oe/mbo.html>).

1 Introduction

There are a multitude of databases accessible over the Internet that cover genomic¹, cellular², structure³, phenotype⁴ and other types of biologically relevant information⁵. Even for one type of information, e.g. DNA sequence data, there exist several databases of different scope and organisation^{1,6,7}.

Unfortunately, naming conventions of data objects, object identifier codes and record labels differ between databases and do not follow a unified scheme. But worse, even the meaning of important high level concepts that are fundamental to many molecular biology databases is ambiguous.

One prominent example is the concept *gene*. For GDB¹, a *gene* is a DNA fragment that can be transcribed and translated into a protein; for Genbank⁷ and GSDB⁶, however, a *gene* is a “DNA region of biological interest with a name and that carries a genetic trait or phenotype” which includes non-structural coding DNA regions like intron, promoter and enhancer. There is a clear semantic difference between those two notions of *gene* but both continue to be used interchangeably causing misunderstanding and making the integration of databases non-trivial.

To eliminate semantic confusion in molecular biology, it will be therefore necessary to have a list of the most important and frequently used concepts coherently defined so that database managers could use such set of definitions either to create new database schemata or to provide an exact, semantic specification of the concepts used in an existing schema.

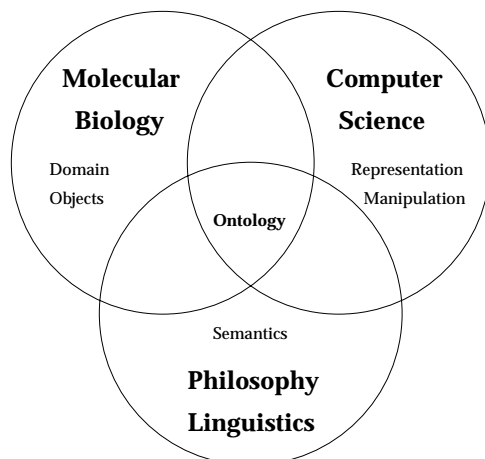


Figure 1: Molecular biologists discover facts that need to be organised and stored in databases. Computer scientists provide techniques for data representation and manipulation. Philosophers and linguists help organise the *meaning* behind database labels.

To become generally acceptable in the molecular biological community such semantic compendium (also often called controlled vocabulary) must be accessible electronically and without licensing charges, preferably using a world wide web browser; be intuitively comprehensible without special computer programming background; be able to cope with natural language features as e.g. homonyms; be capable of performing logical inference over the set of concepts to provide for generalisation and explanation facilities; exhaustively cover the application domain; and be coordinated but open for input from the community. Also, software to manage a semantic repository must be created.

One way to consistently and transparently create such set of definitions for molecular biology is by using an ontology. By adhering to a commonly agreeable ontology, uncertainty and misunderstanding about the semantic relations between database entries from different databases can be eliminated. When all relevant concepts of an application domain will have been specified in an ontology, a computer program can search for *concepts* instead of words in a set of heterogeneous, autonomous databases⁸; carry out semantic consistency checks; and detect ill-formed statements and interpret well-formed ones⁹.

In this report, two well-known ontologies, Cyc and μ Kosmos are examined with respect to their applicability in molecular biology. Since both are found not to fully satisfy that purpose the foundations of a new, prospective ontology for molecular biology are laid out.

2 What is an Ontology (and what isn't)?

Ontology was originally perceived by ancient philosophers as the study of *being*. They asked “What does the statement ‘X *is*’ mean?” and “Which things *are*?”¹⁰. In modern times, computer science uses ontology in a narrower sense as a “specification of a conceptualization”¹¹ or, in other words, as a concise and unambiguous description of what principal entities are relevant in an application domain and how they can relate to each other. The entities can be objects, processes, functions, predicates, or of other type depending on the selected representation formalism (Figure 1). The formal definition of the components of an ontology and heuristics for ontology construction are given elsewhere¹².

An ontology is *not* a collection of facts that arise from an actual, specific situation but it defines and provides all semantic entities and their potential interactions that would be needed to completely describe such situation. *Neither* is an ontology a model for an application domain (which would be a theory), but a compendium that holds all necessary “building blocks” with rules of how and which entities can relate to each other and which ones are semantically incompatible. For example, “transcription of a gene” is an ontologically valid expression whereas a “transcription of a cell” is not. *Nor* is an ontology a database schema which defines categories and their data types in a database but which need not represent ontological relations between entities in the real world.

The graphical representation of an ontology in general is not a tree but a semantic net or conceptual graph¹³ because there are two or more types of links (“is a member of” and “is a subset of” plus additional domain specific relations) which can give rise to circular loops in the graph. However, if only the “is a subset of” or “is a member of” relation is displayed the ontological graph becomes a tree.

3 Ontologies for Molecular Biology

Two ontologies are discussed with respect to their applicability in molecular biology. Several other approaches exist but they are either sparsely populated, specialised to other domains or in general do not easily connect to molecular biology. Then, parts of a new, prospective ontology for molecular biology (OMB) are presented. An interactive, graphical representation of all public classes and instances of μ Kosmos, Cyc and the prospective ontology for molecular biology was prepared with a new ontology editor¹⁴ and is accessible for browsing on the world wide web at <http://igd.rz-berlin.mpg.de/~www/oe/mbo.html>.

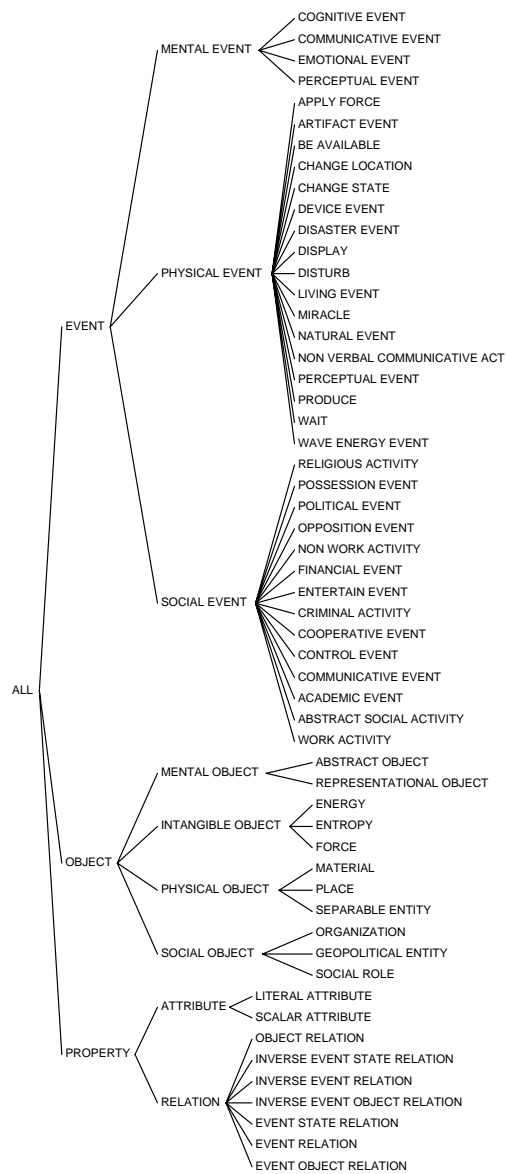


Figure 2: Upper Level of μ Kosmos Ontology. Links represent the “is a subclass of” relation. Hyphenation was removed from original concept names for better reading.

3.1 μ Kosmos

μ Kosmos¹⁵ is an ontology developed for machine translation of business and financial texts between English, Spanish and other languages. The ontology contains 4790 annotated concepts which are publicly accessible as HTML hypertext and Lisp code. Ontological entries in μ Kosmos are linked to dictionaries in several languages. Thus, the ontology can serve as an inter-lingual framework to map the meaning of words from different languages. The links to the language dictionaries are not publicly available.

The upper level of the μ Kosmos ontology is graphically depicted in Figure 2. There are several features of μ Kosmos that limit its suitability for use in molecular biology, some of which are discussed here.

In general, the criterion used to subclassify a concept in μ Kosmos is not made explicit. This makes it difficult to locate the exact position of a given concept. For example, the definition of OBJECT is “ontological concepts that are not actions, or properties; the static things that exist in the physical, mental, and social world” and for EVENT it is “any activity, action, happening, or situation”. However, it can be argued that a situation is a static, mental object describing the relations between actors and components within a context at a given time. Because the criterion to discriminate between OBJECT and EVENT is not stated explicitly an ambiguity about the exact classification of “situation” in μ Kosmos remains.

The definitions of concepts are deliberately kept vague in many places which suits the task of linking different natural languages because of the imprecise nature of non-scientific concepts. For example, the definition for PHYSICAL-OBJECT is an “object which is observable, has position, and has physical dimensions” whereas INTANGIBLE-OBJECT is defined as an “object that cannot be seen or touched but is evident in its influence on the physical world, such as momentum, energy, entropy, etc”. These definitions are not precise enough for natural sciences. Energy in the form of visible light has a physical dimension, can be located and seen and thus qualifies also for a PHYSICAL-OBJECT.

Entropy, because it can only be indirectly accessed, could also be classified a SOCIAL-OBJECT according to μ Kosmos since it is an “object which exists only by the agreement of some people”, in that case, theoretical physicists. The discriminating criterion between PHYSICAL-OBJECT, INTANGIBLE-OBJECT and SOCIAL-OBJECT remains unclear.

μ Kosmos does not include a lot of concepts from molecular biology because it was developed mainly for translation of financial and corporate texts. μ Kosmos cannot store homonyms, i.e. one word with different (and possibly

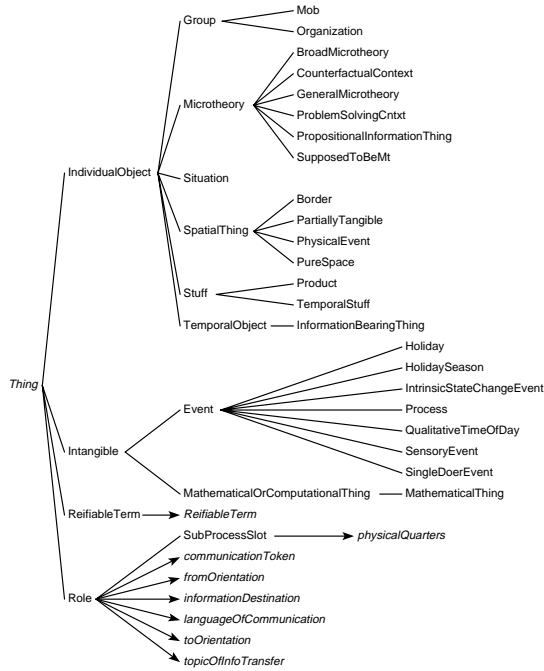


Figure 3: Upper Level of Cyc Ontology. Straight lines indicate “is a subclass of” relation, arrows and italics denote “is a member of” relation (instances). The leading special characters “#” in concept names have been removed for clarity.

disjunct) meanings at different places in the ontological tree. There is no ontology editor available with μ Kosmos and no interactive graphical browsing facilities as at <http://igd.rz-berlin.mpg.de/~www/oe/mbo.html>.

3.2 Cyc

Cyc¹⁶ is an ontology originally developed to cover everyday common-sense knowledge. Of the reported “tens of thousands” ontological entries a subset of about 2200 are publicly available as HTML hypertext with ample documentation. Cyc was not built to support a specific application but with the intention to cover even subtle semantic distinctions that a person has to consider when communicating in daily life. There is no ontology editor available with Cyc and no interactive graphical browsing facilities are provided except static HTML hypertext. Cyc has an inference engine and a natural language interface with

the ontological entries linked to an English dictionary but all of which are not available to the public. The complete version of Cyc is commercially available.

Although Cyc contains a large and detailed collection of well documented concepts it is of limited use for molecular biology for several reasons.

Cyc does not include a significant portion of concepts relevant to molecular biology since it was designed to be a universal ontology. Only very basic knowledge about chemistry and biology has been added.

Although the authors of Cyc state that they “generally only list a non-redundant series of supersets” or “the incommensurably most specific (i.e., smallest) supersets of each collection” this rule is violated on several occasions. For example, `#$Collection` has listed the supersets `#$Intangible`, `#$Thing` and `#$Set` of which `#$Thing` is a superset of `#$Intangible` which in turn is a superset of `#$Set`. There are also several cases where two concepts are listed to be the superset of each other, e.g. `#$Stuff` and `#$IndividualObject`.

`#$Thing`, the “universal set of everything”, has as its immediate subclasses `#$IndividualObject`, `#$Intangible`, and `#$Role` of which all three are overlapping because there exist intangible `#$IndividualObject`(s) and a `#$Role` is something both individual and intangible (Figure 3). The definition of `#$Thing` as the set of everything also faces Russel’s set dilemma.

Though most definitions in Cyc seem philosophically well established, what is visible to the public is counterintuitive in some places. For example, `#$Situation` is defined to be “a state of affairs” with superclass `#$IndividualObject` which is a “discrete, not abstract entity that can have parts but not elements or subsets”, suggesting that not only objects involved in a `#$Situation` but also `#$Situation` itself is a tangible entity since no link to `#$Intangible` exists.

The concept `#$Stuff`, defined as a discrete object that “when divided into pieces remains of the same type” (e.g. water) includes “physical entities like wood”, “temporal entities like the event of a person running” and abstract things like “a piece of English text”. One problem with the definition of `#$Stuff` is its granularity: on a molecular scale wood can well be divided into components that no longer are wood. Similarly, English text can be divided into letters which are neither distinctively English nor text anymore.

The criterion used to subclassify a concept in Cyc is not always stated explicitly. In many cases, subclasses in one class overlap semantically or are created using different subclassifying criteria. No homonyms are found in the public parts of Cyc. Naming of concepts is sometimes confusing, e.g. `#$Thing` vs. `#$SomethingExisting`; `#$PartiallyTangible` vs. `#$PartiallyIntangible`; `#$IntangibleObject` vs. `#$IntangibleStuff`.

Cyc contains a hierarchy of classes containing only classes that in some cases mirrors a similar hierarchy of classes containing instances but which does

not convey any new information. This adds to the confusion when searching for a concept. All these properties of the Cyc ontology make it difficult to locate the appropriate position for an existing concept or for a new one to be added.

3.3 Molecular Biology Ontology

An ontology for molecular biology should become a repository for all relevant concepts that are required to describe biological objects, experimental procedures and computational aspects of molecular biology. Although this looks like an impossible task at first sight it does not mean to compile all knowledge about molecular biology nor does it imply being able to explain every biological phenomena. It just means collecting all types of entities that molecular biologists include in their professional thinking and placing those concepts appropriately in a “is a subset of” and “is a member of” hierarchy plus annotating them with additional properties.

By doing so in a consistent manner, where the discriminating criterion for subclassifying each concept is made explicit, the definition of a concept becomes the path from its own node to the root node of the ontology. As an example, it can be read as “*Name* is an \Rightarrow *Identifier* is an \Rightarrow *Attribute* is a \Rightarrow *Property* is a \Rightarrow ...”, for the case of only “is a subclass of” relationships (\Rightarrow). Relations of type “is a member of” (\rightarrow) inherit only from their direct parent node(s) because being an instance of a concept is different from being a subclass of a class. For example, “Aristotle is a \rightarrow *Name* is an \Rightarrow *Identifier* is an \Rightarrow *Attribute* is a ...” means that Aristotle is one *Name* but not an *Attribute*, whereas *Name* is an *Attribute*.

One difficulty when compiling an ontology for molecular biology is identifying the subtle connotations that are hidden in everyday language. For example, the chromosome of an *E. coli* bacterium is a DNA molecule which is a physical object. The sequence of that DNA, however, is an abstract object that is not contained in *E. coli* at all but that can be subjected to mathematical analysis. Therefore the concepts DNA and DNA sequence will reside on quite remote branches in an ontological graph.

The upper level of a prospective Ontology for Molecular Biology (OMB) is shown in Figure 4. Starting from the root node *Being* which includes anything that *is*, the classes *Object* and *Event* are disjoint and discriminated based on their temporal extent. An *Object* remains an *Object* even in a single moment whereas an *Event* when dissected into single moments loses its identity. This holds also for all subclasses of *Object* and *Event*. The class *Object* is further subclassified into *Individual Object* and *Property*. Both can be thought of as in-

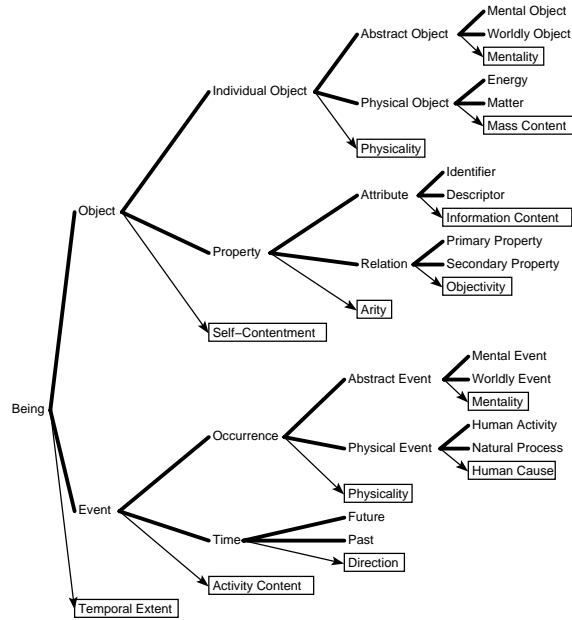


Figure 4: Upper Level of a prospective Molecular Biology Ontology. Links represent the “is a subclass of” relation. No instances are present; discriminating criteria have arrows and boxes; thick lines denote disjunct subclasses.

stantaneous, i.e. they keep their identity even if looked at only for one moment. The two are discriminated based on self-contentment. An *Individual Object* can stand alone whereas a *Property* always needs another *Object* or *Event* to refer to. A *Property* is further subclassified based on arity into *Attribute*, a property with only one argument, and *Relation*, a property relating two or more *Beings*. Hereby, the logical grammar of words, not their surface structure must be considered. For example, in the statement “Paris is beautiful”, beautiful is not a logical attribute to Paris because this statement necessarily involves a second entity, the speaker, and thus becomes one binary and one unary relation: “She thinks, Paris is beautiful” or in Prolog syntax: think (she, beautiful (Paris)).

Attribute can be subclassified into *Identifier* and *Descriptor* based on whether it just labels an entity or whether it carries additional information about it. *Relation* can be subclassified analogous to Locke¹⁷ into *Secondary Property* relations that involve personal judgement and *Primary Property* facts describing intersubjective measurable relations.

Concept	Discriminator	Value
<i>Being</i>	Temporal Extent	instantaneous
Object	Self-Contentment	yes
Individual Object	Physicality	no
Abstract Object	Mentality	no
Worldly Object	Domain Specific Usage	yes
Domain Specific Wordly Object	Subject Domain	mathematics
Mathematical Object	Complexity	high
Composed Mathematical Object	Application Specificity	yes
Applied Mathematical Object	Application Domain	computer science
Computer Science Object	Subject Matter	theory of data
Theory Of Data Object	Subject Matter	data structure
Data Structure Object	Representation Formalism	object-oriented
Object-Oriented Data Structure	Implementation	OPM
OPM Database Object	Database	GDB
GDB Database Object	Subject Matter	genomic
GDB Genome Object	Object Class	DBObject
DBObject	Subclass	MappingObject
MappingObject	Subclass	Map
Map	Subclass	<i>LinkageMap</i>

Figure 5: Semantic Hierarchy for the GDB database category *LinkageMap*. A *LinkageMap* is a *Being* with instantaneous temporal extent, an *Object* which is self-contenting, etc. until the concept *LinkageMap* is reached right at the bottom.

Individual Object is subclassified based on physicality into *Abstract Object*, which has no physical equivalent *per se* (except capable of being represented neurologically or in writing, etc.) and *Physical Object*, which must have a defined spatial extension and/or energy content and is similar to Popper’s “World 1”¹⁸. *Abstract Object* is further subclassified based on mentality, i.e. whether it refers to an object within the mind or to an object in the outside world, into *Mental Object* (similar to Popper’s “World 2”) and *Worldly Object* (similar to Popper’s “World 3”). Although energy and matter are equivalent in nuclear physics a given object can be only of one type at a time. Hence, *Physical Object* has been subclassified based on mass content into *Energy* and *Matter*.

On the other branch of the ontology *Event* is subclassified based on activity into *Occurrence*, where at least one object participates and (pure) *Time*, where nothing happens. This is the notion of absolute time which is no longer valid in relativistic physics and astronomy. The reason for nevertheless holding on to the belief of absolute time here is justified by the intended scope of the ontology for molecular biology: physical processes in living organisms have so far never been known to reach the realm of relativistic physics. *Time* is further subclassified according direction into *Past* and *Future*. Because presence strictly lasts one moment only, it does not appear in this branch.

Analogous to abstract and physical objects, *Occurrence* is subclassified based on physicality into *Abstract Event* and *Physical Event*, and further *Abstract Event* based on mentality into *Mental Event* (similar to Popper’s “World 2”) and *Worldly Event* (similar to Popper’s “World 3”). *Physical Event* is sim-

ilar to Popper’s “World 1” and subclassified based on whether it is done or initiated by human intention into *Human Activity* and *Natural Process*.

Similar reasoning was applied to ontologically capture the meaning of the GDB database category *LinkageMap*, which is defined as a “database object class used to store maps based upon frequency of recombination between genomic segments, resulting in the ordering of markers along a chromosome backbone, usually measured in centiMorgans”. Note that this is ontologically not the same as a linkage map itself which is an abstract concept with certain mathematical properties, nor is it an actual linkage map which is a concrete instance of the class of linkage maps for a particular organism and chromosome.

The complete path from root node *Being* to *LinkageMap* is summarised in Figure 5. The meaning of the database category *LinkageMap* is captured by a series of ontological specifications. This example shows how a semantic definition of a molecular biological concept can be extracted from its mere position in an ontological graph. Similarly, semantic differences and the least general common concept of a pair of concepts can be found by following the graph upwards along “is a subclass of” and “is a member of” links until both paths meet in one concept.

4 Discussion

In this work, the communication problem in molecular biology has been looked at from the viewpoint of semantic integration. To improve the current situation of non-unified and ambiguous vocabulary the only solution is to develop a core of commonly agreeable definitions and using those to implement user interfaces to and between databases. Those definitions must be connected to each other so that no ambiguities over relations between concepts remain and that a computer may infer from specialised to generalised concepts and vice versa.

An ontology as an explicit and hierarchical specification of the relevant concepts in an application domain is one means to develop such semantic repository. Here, two ontologies from the literature, μ Kosmos and Cyc, have been reviewed with respect to suitability in molecular biology. The μ Kosmos ontology is found not to be transparent and precise enough to collect and sort scientific concepts concerning molecular biology. Cyc contains a lot of knowledge about semantic distinctions in daily life but seems to be too complicated and overloaded with concepts not relevant for molecular biology (and probably too expensive) to be of use here.

Work has begun on a prospective ontology for molecular biology. The upper level of the ontology is lean and in agreement with traditional arguments from philosophy of science. In contrast to other ontologies, the criterion used

for subclassifying a concept is explicitly stated and therefore essential decisions and assumptions behind the ontology are made transparent. Homonyms can be handled by explicit reference to their superclasses. Currently, the ontology has about 1300 concepts with emphasis on molecular biology database entities and should be regarded as work in progress.

References

1. K. H. Fasman, S. I. Letovsky, R. W. Cottingham, and D. T. Kingsbury, *Nucleic Acids Research*, vol. 24, no. 1, pp. 57–63, 1996.
2. D. Jacobson and A. Anagnostopoulos, *Trends in Genetics*, vol. 12, pp. 117–118, Mar. 1996.
3. F. C. Bernstein, T. F. Koetzle, G. J. B. Williams, E. F. Meyer, M. D. Brice, J. R. Rodgers, O. K. T. Shimanouchi, and M. Tasumi, *Journal of Molecular Biology*, vol. 112, pp. 535–542, 1977.
4. V. A. McKusick, *Mendelian Inheritance in Man. Catalogs of Human Genes and Genetic Disorders*. Baltimore, MD: Johns Hopkins University Press, 11 ed., 1994.
5. A. Bairoch, *Nucleic Acids Research*, vol. 21, pp. 3155–3156, July 1993.
6. G. Keen, C. Fields, *et al.*, *Nucleic Acids Research*, vol. 24, no. 1, pp. 13–16, 1996.
7. D. A. Benson, M. S. Boguski, D. J. Lipman, and J. Ostell, *Nucleic Acids Research*, vol. 25, no. 1, pp. 1–6, 1997.
8. B. R. Schatz, *Science*, vol. 275, pp. 327–334, Jan. 1997.
9. S. Schulze-Kremer, *Molecular Bioinformatics - Algorithms and Applications*, pp. 13–108. Walter de Gruyter, Berlin, 1996.
10. Aristotle, Translated by E. M. Edghill, 350 BC.
11. T. R. Gruber, *Knowledge Acquisition*, vol. 5, no. 2, pp. 199–220, 1993.
12. S. Schulze-Kremer, in *Proceedings of the Fifth International Conference on Intelligent Systems for Molecular Biology*, pp. 272–275, 1997.
13. J. F. Sowa, *Conceptual Structures: Information Processing in Mind and Machine*. Reading, MA: Addison-Wesley, 1984.
14. S. Schulze-Kremer, in *Molecular Bioinformatics, Sequence Analysis - The Human Genome Project*, pp. 43–56, Shaker Verlag, Aachen, 1997.
15. K. Mahesh and S. Nirenburg, in *Proceedings of the FLAIRS-96 Track on Information Interchange, Florida AI Research Symposium*, May 1996.
16. D. B. Lenat, *Communications of the ACM*, vol. 38, Nov. 1995.
17. J. Locke, *An Essay Concerning Human Understanding*, vol. 2, 4, ch. 8, 3, pp. 23–26, 11–16. Oxford University Press, 1975, 3 ed., 1690.
18. K. R. Popper and J. C. Eccles, *The Self and Its Brain*. Springer, 3 ed., 1985.