

# **Bioinformatics as a problem of knowledge representation: applications to some aspects of immunoregulation**

Sándor Pongor<sup>1,2</sup> and András Falus<sup>3</sup>

<sup>1</sup>*Protein Structure and Bioinformatics Group, International Centre for Genetic Engineering and Biotechnology, Area Science Park, I-34012 Trieste, Italy, e-mail: pongor@icgeb.org*

<sup>2</sup>*Bioinformatics Group, Biological Research Center, Hungarian Academy of Sciences, Temesvári krt. 62, H-6726 Szeged, Hungary*

<sup>3</sup>*Department of Genetics, Cell- and Immunobiology, Semmelweis University, Immunogenomics Research Group, Hungarian Academy of Sciences, H-1089 Budapest, Nagyvárad tér 4, E-mail: faland@dgci.sote.hu*

## **Abstract**

Bioinformatics uses a variety of models that fall into three broad categories such as linguistic, 3-D and interaction network models. Though latter allow one to capture interactions among molecules and other cellular components, the underlying representations are predominantly static. The main molecular mechanisms of immunology such as VJD recombination, cellular and molecular networks, somatic hypermutations - cannot be and are not adequately covered in current molecular databases. Other aspects, such as the maturation of single, monospecific immune response or that of immunological memory apparently fall outside the scope of current molecular representations. The complexity of the immunological regulation such as polarized T cell cytokine web, Treg subpopulation, idiotypic networks, etc. calls for a new generation of computational approach leading to a new age of immunoinformatics (“immunomics”).

## *1. Introduction*

The growing network of biomedical databases and analysis programs constitute one of the most sophisticated knowledge representation tools mankind ever built. Bioinformatics differs from other informatics applications not so much by the amount of the data but rather by the complexity and the depth of knowledge it communicates. As an example,

bioinformatics deals with a wealth of molecular representations, such as sequences, 3/D structures, symbolic diagrams (e.g. hydrophobicity plots, helical wheel diagrams), as well as with a variety of group-wise representations such as multiple alignments, metabolic pathways, phylogenetic trees, etc., most of which could not have been conceived without computerized methods. It is customary to define bioinformatics as the informatics of biological data, but in fact it is not, or not exclusively a specialized branch of science: it is rather a general approach to all life sciences that makes it possible to study problems previously inaccessible to systematic research. This aspect – the access to new domains of knowledge – is one of the common themes that link the current age of computerized resources to previous innovations in storing and representing information. And as representations of information are at the very heart of cultural evolution, it is in place to introduce our subject within a historic context.

Complexity leaps in evolution are known to be powered by improvements in the way genetic information is stored and transmitted (Szathmáry and Smith, 1995, Maynard Smith and Szathmáry, 1995). By analogy one can point out that major improvement in scientific knowledge representation are correlated with innovations in the way information is shared within human societies, such as the appearance of writing, printing and the Internet. In a traditional society, knowledge is exchanged mainly by repeated, face-to-face communication and is confirmed and stored by an entire community. Writing not only decoupled knowledge transfer from personal communication but it also created a powerful new medium for the storage and manipulation of complex symbols whose interpretation required, at the same time, an increased intellectual effort from the recipient. Few would doubt that the widespread use of written and especially printed information has been a prerequisite of modern science that characterizes industrial societies. The paradigmatic knowledge source of this period is the *encyclopedia*, an organized, searchable knowledge base that is, in many respects, the predecessor of current electronic databases.

The current age of bioinformatics is characterized by vast amounts of biological data collected by computerized methods and distributed via the Internet and stored in electronic databases (**Table 1**). Knowledge in electronic databases is represented and transferred ways that is radically different from those known before. While readers can directly interpret printed text, electronic databases can only be “read” with the mediation of computer programs. Programs carry a large amount of implicit information in themselves that is not

always transparent to the user. For instance, in order to draw a three-dimensional picture of a protein molecule from an input of atomic coordinates, a program needs to know how the atoms of various types of amino acids are connected with each other. This kind of implicit information represents an intermediate layer between the data and the program and is often organized into *ontologies*, i.e. formal sets of definitions and rules that are valid for the data domain (**Table 2**). Also there are conspicuous changes in the way scientific information is confirmed. In the age of printed information, the quality of scientific discoveries was guaranteed by authoritative scientific societies, by the peer review of scientific journals, and last not least by the personal reputation of individual authors. In contrast, electronic databases are often produced by automated data collection, while textual annotations, such as the description of biological function, etc. are added by anonymous teams of database annotators who often rely on computer-based prediction methods. In other words, the amount of data and the number of databases is growing while data quality is less transparent. However, there are signs of integration as well. Large efforts are devoted to validating and interlinking biological data (sequences, structures), and, what is perhaps more important, highly complex scientific resources have been created wherein diverse data are controlled and accessed by uniform methods. In this setting, data integrity and quality will be increasingly controlled by autonomous agents, which will hopefully decrease the quality gap of current databases.

The first goal of this chapter is to provide of overview on how knowledge is represented in bioinformatics today and to show the cognitive roots of the underlying models. Bioinformatics is primarily concerned with the structure of protein and DNA molecules that fulfill functions in a series of interdependent systems such as pathways, cells, tissues, organs and organisms. This complex scenario can be best described with the concepts of systems theory<sup>1</sup>. Models in molecular biology are simplified systems that can be conveniently be described in terms of entities and relationships (Pongor, 1988). We will deal with three main kinds of representational models: language based models, 3-D models, and networks and will briefly review the development of computational tools that operate

---

<sup>1</sup> According to systems theory, a system is a group of interacting elements functioning as a whole and distinguishable from its environment by recognizable boundaries

CSÁNYI, V. (1989) *Evolutionary Systems and Society*, Durham and London, Duke University Press, KAMPIS, G. (1991) *Self-modifying systems in Biology and Cognitive Science*, Oxford, New York, Pergamon Press. Molecules can be regarded as such systems. Generally speaking, structure is fixed state of a system, and the study of a system usually starts with its characteristic structures that are recurrent in space or time. Function on the other hand is not a property of the system rather a role that the system plays in the context of a higher system.

on these models. The last sections of this chapter describe genomic and immunological resources. The second goal of this chapter is to review the current knowledge on the information processing mechanism of the immune system. The immune system has specific algorithms for handling environmental information, and these mechanisms are among the most intensively researched and perhaps best understood phenomena in the life sciences that calls for specific informatics approaches yet to appear.

## *2. Sequences and Languages*

Language-based descriptions are broadly speaking those that use semantic definitions for entities and relationships. The term “Molecular biology” was independently coined by Warren Weaver and John Astbury in the nineteen thirties, at a time when scientific methodology was dominated by linguistic theories (Wittgenstein, 1922, Carnap, 1939). Subsequent breakthroughs in information theory (Shannon, 1948b, Shannon, 1948a) and formal linguistics (Chomsky, 1957) all pointed towards a broad metaphoric context of language, communication and computation that provided the first framework within which genetics was discussed. Cryptography (Shannon, 1948c) and pattern recognition methods (Ripley, 1999), first developed within intelligence communities of the time, also contributed a great deal to the general view that biological sequences represent a code that carries information in a particular language, a metaphor reflected by such terms as the “genetic code” or “the book of life”.

The analysis of biological sequences first used the statistical tools developed to analyze character strings in linguistics (Konopka, 1994), and many of the first methods such as those concerned with the string complexity, became standard bioinformatics tools in the later years. From the 1980’s as bioinformatics became part of laboratory routine, pattern recognition methodologies that use similarity measures and classification algorithms proved to be of immediate interest, and with the onset of the genomic era, heuristic methods of searching biological databases such as BLAST (Altschul et al., 1990) became the most frequently used algorithms not only within bioinformatics, but allegedly in the entire field of scientific computing.

Margaret Dayhoff and her colleagues at the national Biomedical Research Foundation (NBRF), Washington, DC created first sequence database in the 1960s, an Atlas of protein sequences organized into families and superfamilies, and their collection center eventually became known as the PIR resource. Collections of DNA sequences (**Table1**), started at the European Molecular Biology Laboratory (EMBL, Heidelberg), at Los Alamos National Laboratory (New Mexico) and DDBJ, Japan, gained importance with the spreading of productive DNA sequencing technologies. Initially, sequence records included only the sequence the filename. These were eventually expanded to include annotation information such as references, function, regulatory sites, exons and introns, modified amino acids, protein domains etc. The Swiss-Prot collection of protein sequences is an especially good example of a well-annotated sequence database wherein a uniform syntax was developed for annotations. Very soon, separate so called secondary collections were created for annotated segments, such as the first protein domain sequence database (Pongor et al., 1993), as well as for posttranslational modifications, functional annotations, etc. (Table 4). Development of such secondary databases provided an important entry for specialized information, and current databases such as PFAM are excellent examples of this tendency. An important step was the application of WWW technology for cross-referencing the major databases with each other and subsequently with bibliographic database.

Information contained in current sequence databases(**Tables 3-4**) can be best pictured as an annotated sequence, a linear string of characters to which additional items of informations are linked either as an added field within the same record, or as a cross-reference to another database. Annotations items refer either to the entire sequence (global descriptors, such as protein name) or to a segment of it (local descriptors, such a protein domain or an exon). A database record itself can be an annotation items in a different database: for instance, a record in a protein domain database, or in a bibliographic database can be linked as an annotation items to a sequence database. In principle, annotated sequences should point to a unique protein or gene. In practice, a unique protein can be represented by many sequences, and databases differ in the way redundancy is handled. GenBank contains all published DNA sequences, so there is a considerable redundancy. The same is true for the protein sequence collections prepared by automatic translation, or by automated experimental procedures such as EST sequencing. Maintenance of high quality non-redundant databases require human overhead that is increasingly difficult to provide.

Finally, much of annotation information today is provided by automated procedures, such as similarity searches, HMM methods etc. Even though these methods constantly improve, there is no absolute guarantee behind the information, so much of annotation information today is labeled as “putative” or “by homology”.

Interpretation of annotations and bibliographic records requires a uniform, computer-readable language. This need has fostered intensive research into the natural languages used in science. In early cultures, scientific language describing Nature cultures was based on only four elements (earth, water, fire, wind), and a few, dichotomic relations (hot-cold, dry-wet etc.) between them. Descriptions used in database annotations are based on a large number of models for atoms, molecules and reactions (**Figure 1**) This is a stripped-down language that can be kept uniform by the makers of the databases. On the other hand, scientific publications and abstracts use a “free-style” scientific language that is hard to handle by computers. This problem has emphasized the needs for developing common ontologies for molecular biology (Schulze-Kremer, 1997, Ashburner et al., 2000). An ontology defines a common vocabulary and a shared understanding within a domain of knowledge such as molecular biology. An ontology is an explicit description of the knowledge domain using concepts, properties and attributes of the concepts, and constraints on properties and attributes. The Gene Ontology Annotation (GOA) database (<http://www.ebi.ac.uk/GOA>) is a system of concepts and relations that is designed to convert UniProt annotation into a recognized computational format. GOA provides annotated entries for nearly 60,000 species and is the largest and most comprehensive open-source contributor of annotations to the GO Consortium annotation effort.

### *3. Three-dimensional models*

As sequence databases were born within molecular biology, 3-D databases were brought to life by chemistry and later, structural biology. Same as language provides a conceptual framework for sequences, the metaphor for molecular models are common world objects whose handling and recognition is as at least as deeply rooted in human cognition as is language (Pinker, 2001). The first 3-D model of a molecule was constructed in 1874 by van t'Hoff who recognized that optical isomerism can only be explained by a 3-D arrangement of the chemical bonds. As methods of X-ray crystallography became applicable to organic molecules, Olga Kennard and Desmond. Bernal initiated the collection of 3-D structures what later became known as the Cambridge Crystal Structure Database, and access to 3-D information consolidated the use of molecular geometry in chemistry.

3-D descriptions of macromolecules are based on a series of concepts that resulted from several decades of scientific work. For instance structural descriptions of proteins is based on the stereochemistry of the peptide bond, but elements of secondary structures, supersecondary structures and finally protein folds are the joint results of several scientific disciplines that form what is called structural biology today.

The common ancestor of structural databases is the Protein Data Bank (PDB) (Berman et al., 2000), which was established in 1971 and six years later contained only 77 atomic coordinate entries for 47 macromolecules (Bernstein et al., 1977). Over the years, a conspicuous number of secondary databases have evolved from the PDB (see Table 2). Many of them concentrate on various classes of structural features, such as protein domains (Sander and Schneider, 1991, Orengo et al., 1997, Murzin et al., 1995, Siddiqui et al., 2001), loops (Donate et al., 1996), contact surfaces (Jones and Thornton, 1996, Luscombe et al., 2000), quaternary structure (Henrick and Thornton, 1998), small-molecule ligands (Kleywegt and Jones, 1998), metals (Castagnetto et al., 2002) and disordered regions (Sim et al., 2001). Other databases concentrate on biological themes. The very concept of the “protein fold” owes much of its existence to such protein domain databases as CATH, SCOP and FSSP.

The conceptual structure of current 3-D databases is similar to sequence databases, inasmuch as the records contain both a structural description as well as an annotation part.

The definitions of the individual fields reflect the fact that PDB was originally created as a crystallographic database, and despite the fast growing body of NMR data, this remains its legacy. Current databases, such as ?? are now linked to other molecular and bibliographic databases.

The development of ontologies has begun in structural biology. The STAR/mmCIF ontology (Westbrook and Bourne, 2000) of macromolecular structure is description of structural elements and data items in the framework of X-ray crystallographic experiments but it is extensible other kind of data collection techniques.

#### *4. Genomes, proteomes, networks*

Designing representations for genomes, proteomes and networks is a challenge as we deal with a wide variety of entities and relationships, partly predefined, partly discovered during the project. This class of representations can be called ‘general topological model’, wherein the nature of entities and relationships is not limited either to semantic or to 3-D concepts, as in the previous chapters. The resulting general representation is a graph wherein the physical entities are the nodes and their relations are the edges. The common ancestor of this representation is the structural formula, and graph theory itself owes a great deal to the development of chemistry in the nineteenth century. The representations of genomes as linear array of genes and other DNA segments follow a similar traditions tradition. The entities – genes – are predicted with gene-prediction programs or are determined experimental methods, and this adds a new layer of knowledge to the molecular data. The relationships are manifold but are predominantly binary in nature. Examples of relations include physical vicinity, distance along the chromosome, regulatory links extracted from DNA chip data and so on. The resulting picture is a graph of several ten thousand nodes and relatively few edges per node denoting various relationships. The description of proteomes is only somewhat different. The proteins are described in functional, biochemical and structural terms, and the relationships between proteins include metabolic relationships (sharing substrates in metabolic pathways) as well as structural relationships (sequence and structural similarities). Network models used in biology fall into two large categories. Dynamic models (such as metabolic network of a cell) are based on differential equations of the constituent the enzymatic reactions (for a recent review see: ). Topological



models discussed here deal with the static properties of graphs, that can be undirected, directed and weighted.

From the computational point of view, genomes and proteomes are described as very large graphs in which the nodes (genes, proteins) and the edges (relations) are unknown or unsure. These large and fuzzy descriptions are in sharp contrast with the descriptions developed for well-defined graphs molecular structures, but the methods are not dissimilar to those used in other applications of graph theory. Given the large and varied genome sizes as well as the uncertainties of the data, genomic networks are usually characterized and compared in terms of gross global descriptors such as the degree distribution, composition (e.g. composition expressed in terms of gene- or protein classes). Biological networks are also believed to contain recurrent local patterns (network motifs) that are analogous to sequence motifs found in biological sequences.

Even this sketchy introduction implies that we deal with new a kind of complexity that originates, from the numerous and to a large extent, unknown interactions between the entities. On the other hand, the study of network topology in various fields – such as Internet, social- road- and electric networks, etc. – has provided interesting insights that have been successfully applied to genomes, proteomes and bibliographic networks. However these insight are limited by the fact that static network topology is only a very general description of the underlying biological phenomena, in fact it should rather be considered as a “metamodel” i.e. a “model of models”.

#### 4. Computational tools

Bioinformatics came to life at a time when computer technology reached the daily routine of scientific research. The development of bioinformatics tools (e.g. interfaces, database design, programming methods) is to some extent a mere reflection of the concomitant trends in informatics. On the other hand, bioinformatics software is peculiar because of its impact on how lay users access biological data today. In the 1970's and early 80's, the first published programs were written in a basically sequential style for standalone computers such as campus mainframes. The second stage started by the recognition that the input and output of bioinformatics applications can be standardized, and modular packages based on the *software tools approach* (Knuth, 1998) were developed. The best known example of these, the GCG package of John Devereux (Devereux et al., 1984) developed into a battery of several hundred programs over the years, covering virtually the entire scope of biological sequence analysis. However, such a large body of knowledge is difficult to maintain in a commercial context. EMBOSS, which is developed by a collaboration of academic researchers was designed as an open source alternative commercial programs (Rice et al., 2000). The development of Bioperl (Stajich et al., 2002) – a *Perl* library for bioinformatics applications – and Bioconductor (Gentleman et al., 2004) – a statistical programming package for microarray analysis based on the *R* programming language – are further examples of successful open source collaborations. Web servers developed by academic research group represent a different trend since in such cases, the source code is often not released and users can access the programs only on-line. WWW technology thus allows individual researchers to release their programs before the commercial or open source development stage, and the users interested in the most recent computational tools more and more accept the risk of using non-transparent programs.

The most visible tools of current bioinformatics are the complex knowledge resources composed of databases, analysis tools and WWW interfaces that integrate various kinds of data into a navigable data network (**Figure 2**)

#### 5. Information processing in the immune system

The immune response includes a number of regulating mechanisms that are organized into various networks affecting a wide range of phenomena ranging from uptake, processing and presentation of the antigens, to T and B cell activation and performing the effector functions. During the immune homeostasis, the spatial and temporal pattern of the cellular and soluble interaction networks develop the optimal qualitative and quantitative characteristics in starting, amplifying and finishing the immune response in an optimal way. The real understanding of immune response obviously requires a systems biology approach. There are four highly specific aspects of immune functions that warrant a specific immunoinformatic approach.

1. The immune systems itself functions as a highly regulated information processing system. The major informations are the sequence ( $\alpha\beta$ TCR recognition) and the conformation (BCR/Ig and  $\gamma\delta$ TCR recognition) of the antigen/peptide molecule, the genetic (e.g. MHC) background of the antigen presenting cells, the activation of the innate immune systems (e.g. complement, NK pattern), the actual environmental scenario (e.g. PAMP, local pattern of the inflammatory mediators, etc.).
2. The networking habits of the immune system, both at cellular level (enhancing and inhibitory effects, feed-back regulations) and as Ig networks, such as the idiotypic-antiidiotypic webs.
3. The VJD recombinations and other mechanisms generation diversity of the antigen receptor repertoires are basically different from other , more simple nets. These molecular events result in rapidly evolving gene sets serving a more sophisticated recognition tool during the afferent input of immune response.
4. Somatic hypermutations through *activation induced deamination (AID)* and repair machinery as highly effective way for generation of diversity belong to specific tools .

.In the following issues some relevant representations of immune regulation are mentioned

Information management by lymphocytes

What “program” the given cell has, namely what the “output” signal, the response is, issues from its genetic characteristics and the features it acquired during its ontogenesis

(Figure 3). In other words: after the appropriate co-stimulatory and cytokine effects, a cell-specific pattern of transcription factors develops. Upon their effect, the appropriate cell-response develops, thus, the cell divides or/and differentiates; its fate will be the transformation into a memory cell or apoptosis; it releases antibody, cytokines or other secretion products into the outer world. By this, it frequently participates in the activation or inhibition of another cell or functions as effector. Sometimes it happens that the cell “changes” its signal transduction way during the regulation. For instance, the cAMP-signal transduction related to MHC class 2 on B cells “switches” on to the tyrosine-kinase way.

#### The network of co-stimulatory effects

During the immune response, many cell-cell relations are formed and dissociated. Beyond the physical “approach”, these cell-cell relations may cause intracellularly created signs, which change the functioning of the cell.

Besides the antigen-specific interaction (TCR-MHC/peptide), the T and B cells must receive supplementary signals for carrying out a successful immune response. Without these, the antigen-specific interaction in most case produces exactly produces anergy, incapability to respond. This anergy stands behind the phenomena of immune tolerance. Next to the specific antigen-receptor, there are co-receptors (CD4, CD8), whose binding to the appropriate MHC-molecules orienting T cells regarding which antigen-presenting cell (and exogenous or endogenous antigen) should be connected with.

Our present knowledge lists numerous co-stimulatory (formed with adhesion molecules and other membrane-proteins) interactions of positive and negative effect between the T cell and the antigen-presenting cell. Some interactions increase the activation, e.g. the CD28/B7.1, CD2-LFA-3 and CD40-CD40L (CD154) interactions (Sadra et al., 2004). The CTLA-4/B7.2 linkage, for instance, may have positive or negative effect depending on the circumstances (Zhang et al., 2003).

The possibility in relation with the special characteristics of NK-receptors raises a further regulation opportunity. The NKB1- and p58-receptors of NK-cells recognize MHC, but this recognition hinders the activity of NK-cells (KIR) (Wilson et al., 2000). It has been recently

stated that these *NK-receptors are also found on a subgroup of  $\alpha\beta$ - and  $\gamma\delta$ -cells*, where MHC-recognition is exactly a (positive) condition for the stimulus realized through the  $\alpha\beta$ -TCR-receptors. The point is that, within one cell, two receptor-structures (NK-receptor and TCR) are oppositely regulated by the MHC; thus, the outcome of cell-activation is influenced by the local proportion of the two kinds of receptors.

Presumably, the local, cell-level pattern of costimulating effects organized in time-order represent that fine regulation, which is optimal for starting, properly setting and concluding the immune response in the right time (Frauwirth and Thompson, 2002).

T-cell-dependent stimulation and inhibition, the Th1- and Th2-cytokines

Th-cells are heterogeneous considering their cytokine-production. It is proved that the cytokine-pattern of the Th1- and Th2-cells at the two ends of the polarized T-cell-lineage does not only deviate from one another, but, because of the cross-regulation, certain cytokines inhibit each other and each other's effect crosswise. This, together with the negative costimulation, has a significance in halting the response.

One of the most important effects of the IL-10 of Th2-origin is that it strongly inhibits the cytokine-production of Th1-cells, like, for instance, IL-2-synthesis (and its effects). The influence of IL-2 on B-cells is inhibited the same way by the IL-4, which also has mainly Th2-origin. The antagonism between the effects of IL-4 (Th2) and IFN $\gamma$  (Th1) is extremely sharp and two-directional (e.g. on IgE-production, on DTH). Individual cytokines frequently affect the isotype of antibody production in different way, stimulating one and inhibiting the other (e.g. IgG-IgE). The cytokine-pattern determined by the tissue environment of B cells significantly influences the class-switch, that is, the development of the isotype (e.g. IgG or IgA antibodies are produced).

Usually it can be said (with exceptions) that T cells producing IFN- $\gamma$ , IL-2, and TNF- $\beta$  principally stimulate the cell-mediated immune response, while the T cells producing IL-4, IL-5, IL-9 and IL-6 mainly stimulate the humoral immune response. It is also interesting that chiefly B cells present the antigens to Th2-cells, while macrophages present the antigens to Th1-lymphocytes (Figure 4).

Elegant results have proved that the “Th1-Th2” character - related also to chemokine patterns, (Kim et al., 2001) - is not connected strictly to the CD4- marker, since this double-nature has been detected in some CD8+ cells as well. These cells are called Tc1/Tc2 cells. We also know T $\gamma\delta$ -1 which mainly produce IFN- $\gamma$  and T $\gamma\delta$ -2 primarily secreting IL-4.

Th3 cells represent a separate subset characterized by TGF- $\beta$ -production. They stimulate the functioning of Th1- and inhibit the functioning of Th2-subgroups. They also have important role in the IgA-production of the immune system attached to the gastro-intestinal system. Recently markers characteristic of the Th1- (LAG-3) and the Th2- (CD30) population have been found on the plasma membrane. LAG-3 is a molecule belonging to the immune-globulin supergene-family, while CD30 is a protein belonging to the TNF-receptor family, already known on activated T and B cells as well as on the Sternberg-Reed-cells typical of the Hodgkin-lymphoma.

If it is indeed true that these molecules are markers of the two T cell populations, the cytofluorometry which detects the surface markers, will be a supplementary method for measuring the Th1/Th2 rate besides the rather expensive cytokine-mRNA-measurements. As a result of T cell polarization, an important “division of labor” is formed during overcoming different infections. The role of the Th1-type Th cells is important against the Gram-negative bacteria, while the role of the Th2-type Th cells is essential against parasite infections (Figure 4.).

In these inhibiting and stimulating processes, the non-antigen-specific cells and products of the immune system also play an important role. For instance, if NK-cells are activated in the local immune reaction because of the antigen’s nature (tumor-cell, virus) or other factors, this causes increase in the local IL-12 and IFN $\gamma$ -level, that is, a “Th1”-like (+ and -) effect. According to present views, the IL-12 has a central role in the regulation of cellular (cell-mediated) immune response and promising results have come to light in the treatment of metastatic tumors and diseases caused by hepatitis B and C virus-infections, with the help of IL-12.

If the number of basophilic granulocytes increases locally because of the antigen’s nature (allergen, vermin) or other factors, the result will be exactly opposite: it asserts the increase

of IL-4-level, that is, “Th2”-like influences (+ and -) (Figure 5.). On the other hand, IL-4 plays the “conductor’s” role regarding the humoral immune response. There is a therapeutic possibility in IL-10 (which is also supposed to have Th2-origin) in allergic diseases, because it inhibits the attractive effect of IL-5 on eosinophils.

Similar differences are caused by the locally effective prostaglandins produced by macrophages, fibroblasts and follicular dendritic cells. PGE1 and PGE2 inhibit the cytokine-production of Th1-cells but do not influence Th2-lymphocytes. As a consequence, PGE-s shift the balance locally into the direction of humoral immune response. We have also come to know (e.g. from the AIDS/HIV research) that corticosteroids principally inhibit Th1-cells (apoptosis-induction), while certain androgen steroids (e.g. dehydro-epiandrosteron) inhibit the corticosteroids’ Th1-blocking effect, that is, antagonize it.  $\beta$ -endorphin inhibits the rate of the Th1- and stimulates the rate of Th2-cells.  $\beta$ -antagonists weaken the cellular immune response through IL-12 inhibition.

Generally we can say that, in the organism, the Th1 and Th2 found on the two ends of the polarized T-cell-line, participate in different processes as regulating cells. Nevertheless, we should never explain the effects of the Th1/Th2 cytokines “dogmatically”, since the same cytokine can often act oppositely, depending on the concentration, place and time.

#### The CD4+NK.1.1 and Treg subgroup

Recent results have reported on a T cell type called CD4+NK.1.1+. It can be considered the main source of IL-4 and plays central role in anti-microbial immunity. The NK.1.1+ subgroup is a population which is CD4-CD8-  $\alpha\beta$  in the thymus, CD4+CD8-  $\alpha\beta$  on the periphery, having numerous NK-markers, cytotoxic ability and a relatively homogenous  $\alpha\beta$ -repertoire. Its function is supposedly the regulation of haematopoiesis, the development of T cell tolerance, the cytotoxic removal of virus-infected liver-cells and the stimulation of the Th2 population’s maturation by IL-4. It is presumed that these cells primarily recognize microbial antigens presented by the monomorphic CD1 (Jiang and Chess, 2004).

On the basis of all these, the CD4<sup>+</sup>NK1.1 cells could be considered cellular elements of the non-antigen-specific immunity, a new (regulating?) subclass of the T cells (Godfrey and Kronenberg, 2004).

Most recently a new concept of Treg cells (CD4<sup>+</sup>, CD25<sup>+</sup>) develop (Walsh et al., 2004). These cells representing over 10% of CD4<sup>+</sup> Th cells are mostly silencing cells expressing foxp3 transcription factor acting in various regulatory networks of immune response, producing TGFb, IL-10 and negative co-stimulatory molecules such as CTLA4 (Figure 6.). Treg cells are involved in transplantation tolerance, prevent pathological responses induced by the gut flora or microbial infections, play a role in maternal tolerance, can suppress antitumor immunity and protective immunity to pathogens and can lead to enhanced memory T cell response. Recently immunoregulatory disturbances of many autoimmune diseases (Frey et al., 2005) are coupled with dysfunction of Treg subsets.

Idiotype-regulation, idiotype network (Poljak, 1994)

The potential of antibody and TCR diversity developing in the immune response is extremely high. About  $10^{11}$  different antibody- and  $10^{15}$ - $10^{18}$  different TCR-specificity develop in a healthy adult immune system. Since a significant part of this huge repertoire is not expressed (or not in a significant degree) during maturation in the thymus, an auto-tolerance cannot be formed against them. As a consequence, the segments of the variable region (paratope, complexity-determining region-CDR) appear as antigens (idiotypes) in the organism. In 1973, Niels Jerne announced his attractive theory according to which the organism produces anti-idiotype antibodies. Antibodies (anti-anti-idiotypes) are formed again against the antigens in the variable regions of these antibodies and so on. This way, a network develops (Figure 7), where there is a possibility for every second element to contain similar epitopes (idiotopes) and have similar antibody-specificity. Thus, the first antibody activates a (B and Th cell) self-regulating network and the antibody playing the role of “antigen” activates the next element of “antibody” role. This, again activates the next one which may have “first antibody”-like idiotopes among its idiotopes. The network starts to limit itself; smaller and smaller amount of new antibodies are produced. Theoretically, the situation is similar in the variable regions of the TCR's  $\alpha\beta$ - and  $\gamma\delta$ -chains. The idiotype-anti-idiotype network represents a complex network of interacting T and B cells, which can equally stimulate or inhibit the immunological activation. One of the



most important central principles of the network model is that the anti-idiotypic antibody (second antibody) – which reacts with an idiotope in the CDR of another antibody (first antibody) – presents similarity with the original epitope, being the inner image of the original antigen-determinant. This theory is supported by facts; “second and third”-type anti-idiotypic antibodies against monoclonal paraproteins have already been detected. By the improvement of detecting techniques researchers have recently identified anti-idiotypic antibodies during normal immune response, as well. The idiotype network has a great significance in keeping the memory cells of the immune system in small but lasting excitement after the disappearance of the original antigen, in the presence of the second antigen, through an “internal image”. This way, upon the repeated appearance of the “real” antigen, they can react quickly to it.

The proof showing the presence of the idiotype-network is the following: in a type of autoimmune thyroid disease, the anti-idiotypic antibody acting against the autoantibody against the thyreoid-stimulating hormone (TSH) behaves like the TSH and binds to TSH-receptors. In the future, this phenomenon can be possibly used in vaccination. Anti-idiotypic antibodies will be produced in an experimental animal against the antibody (specific for of using a quite dangerous, living pathogen) also produced in an experimental animal. This second antibody may be similar to original infectious antigen, therefore it can be used in immunization (securely, since it is an immunoglobulin), and thus, immunity develops in the organism.

Today debating opinions are published concerning the central role and size of a given idiotype network, but surely, this system has a considerable role in the regulation of the immune response. There are proofs about the existence of idiotype-specific Th-cells and idiotype (antibody-, Th-) -cross reactions are also assigned importance in certain autoimmune diseases.

## *6. Concluding remarks*

The general approach of bioinformatics emerged from a parallel growth in two major fields, life sciences and information technologies. This concomitant development provided access to several new fields, and it also resulted in new conceptual and technological tools for

representing and manipulating scientific knowledge. Integrated databases, analysis programs and ontologies are typical results of this development.

Current bioinformatics deals with a large number of models that are basically static in nature. Molecular databases represent information in terms of linguistic, 3D and network models. Though latter allow one to capture part of the interactions among the molecules and other cellular components, the descriptions are still predominantly static. Large part of the information is stored in cellular, tissue and systemic models that are not part of the databases but – if they are represented at all - are stored in ontologies or are part of the background knowledge of biologists that is not accessible to computers.

The current databases of immunology are similar in structure and in philosophy to other molecular databases. Nevertheless immunology represents a specific case in many respects, and especially in the latter sense, since many of the models in immunology are different from those in other fields. This is on the one hand valid to the main molecular mechanisms of the immune system (such as VDJ recombination, cellular and molecular networks, somatic hypermutations) which cannot be and are not adequately covered in current databases. On the other hand, major aspects of the immune system are not presently considered in molecular databases. As an example, the maturation of a single, monospecific immune response or that of immunological memory calls for a selectionist description of cell populations. Or, the similar structure-similar function paradigm that works magnificently at the level of most protein classes, is blatantly invalid in discriminating IgG molecules specific for different epitopes.

The overview presented in the present chapter suggests that the needs of immunological research will inevitably create a specific bioinformatics approach that will allow one to represent and access knowledge gained in these highly important fields.

Finally we mention that many biological concepts found their ways back to informatics: artificial neural networks and genetic algorithms are success stories of the computer sciences. It is an intriguing possibility that the information processing methods of immunology, the **immunomics concept** will inspire novel computational approaches (De Groot, 2004, Wang and Falus, 2004, Brusica and Petrovsky, 2005).

Table 1 Main types of bioinformatics databases<sup>1</sup>

**Nucleotide Sequence Databases**

International Nucleotide Sequence Database Collaboration  
Coding and non-coding DNA  
Gene structure, introns and exons, splice sites  
Transcriptional regulator sites and transcription factors  
RNA sequence databases

**Protein sequence databases**

General sequence databases  
Protein properties  
Protein localization and targeting  
Protein sequence motifs and active sites  
Protein domain databases; protein classification  
Databases of individual protein families

**Structure Databases**

Small molecules  
Carbohydrates  
Nucleic acid structure  
Protein structure

**Genomics Databases (non-vertebrate)**

Genome annotation terms, ontologies and nomenclature  
Taxonomy and identification  
General genomics databases  
Viral genome databases  
Prokaryotic genome databases  
Unicellular eukaryotes genome databases  
Fungal genome databases  
Invertebrate genome databases

**Metabolic and Signaling Pathways**

**Enzymes and enzyme nomenclature**

**Intermolecular interactions and signaling pathways**

**Human and other Vertebrate Genomes**

Model organisms, comparative genomics  
Human genome databases, maps and viewers  
Human ORFs  
Human Genes and Diseases  
Model organisms, comparative genomics  
Human genome databases, maps and viewers  
Human ORFs

**Microarray Data and other Gene Expression Databases**

**Proteomics Resources**

**Other Molecular Biology Databases**

Images of biological macromolecules  
Bioremediation database  
Drugs and drug design  
Molecular probes and primers  
Organelle databases

**Plant databases**

General plant databases  
Arabidopsis thaliana  
Rice  
Other plants

**Immunological databases**

<sup>1</sup>Based on the Database issue of *Nucleic Acids Research*, 2005, <http://www3.oup.co.uk/nar/database/cat/12/>

Table 2 Databases for annotation terms, ontologies and nomenclature used in bioinformatics

Genew the Human Gene Nomenclature Database,  
<http://www.gene.ucl.ac.uk/cgi-bin/nomenclature/searchgenes.pl>  
GO - Gene Ontology, <http://www.geneontology.org/>  
GOA - Gene Ontology Annotation <http://www.ebi.ac.uk/GOA>  
IUBMB Nomenclature database for enzymes <http://www.chem.qmul.ac.uk/iubmb/>  
IUPAC Nomenclature database for organic and biochemistry, <http://www.chem.qmul.ac.uk/iupac/>  
IUPHAR-RD Pharmacological nomenclature for receptors and drugs,  
<http://www.iuphar-db.org/iuphar-rd/>  
PANTHER Gene products nomenclature, <http://panther.celera.com/>  
STAR/mmCIF: an ontology for macromolecular structure, <http://ndbserver.rutgers.edu/mmCIF> UMLS -  
Unified Medical Language System (Thesaurus, lexicon and semantic networks)  
<http://umlsks.nlm.nih.gov>

Table 3 Examples of protein sequence databases

<i>Primary protein sequence resources</i>	
Uniprot/ <a href="http://www.expasy.org/sprot/">Swiss-Prot</a> - Annotated protein sequence db (University of Geneva, EBI)	<a href="http://www.expasy.org/sprot/">http://www.expasy.org/sprot/</a>
Uniprot/Trembl - Computer annotated protein sequences (EBI)	<a href="http://www.ebi.ac.uk/trembl/">http://www.ebi.ac.uk/trembl/</a>
Uniprot/PIR – Annotated protein sequences (Georgetown University)	<a href="http://pir.georgetown.edu/">http://pir.georgetown.edu/</a>
Uniprot (Universal Protein database, Swissprot + PIR + TREMBL)	<a href="http://www.expasy.uniprot.org/">http://www.expasy.uniprot.org/</a>
<i>Secondary protein sequence resources</i>	
COG - Clusters of Orthologous Groups of proteins	<a href="http://www.ncbi.nlm.nih.gov/COG">http://www.ncbi.nlm.nih.gov/COG</a>
CDD – Conserved Domain Database	<a href="http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml">http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml</a>
<a href="http://pmd.ddbj.nig.ac.jp/">PMD</a> - Protein Mutant db	<a href="http://pmd.ddbj.nig.ac.jp/">http://pmd.ddbj.nig.ac.jp/</a>
<a href="http://www.ebi.ac.uk/interpro/">InterPro</a> - Integrated Resources of Proteins Domains and Functional Sites	<a href="http://www.ebi.ac.uk/interpro/">http://www.ebi.ac.uk/interpro/</a>
<a href="http://www.expasy.org/prosite/">PROSITE</a> - PROSITE dictionary of protein sites and patterns	<a href="http://www.expasy.org/prosite/">http://www.expasy.org/prosite/</a>
<a href="http://www.blocks.fhcrc.org/">BLOCKS</a> - BLOCKS db	<a href="http://www.blocks.fhcrc.org/">http://www.blocks.fhcrc.org/</a>
<a href="http://www.sanger.ac.uk/Pfam/">Pfam</a> - Protein families db (HMM derived) [Mirror at <a href="http://genome.wustl.edu/Pfam/">St. Louis (USA)</a> ]	<a href="http://www.sanger.ac.uk/Pfam/">http://www.sanger.ac.uk/Pfam/</a> <a href="http://genome.wustl.edu/Pfam/">http://genome.wustl.edu/Pfam/</a>
<a href="http://bioinf.man.ac.uk/dbbrowser/PRINTS/">PRINTS</a> - Protein Motif fingerprint db	<a href="http://bioinf.man.ac.uk/dbbrowser/PRINTS/">http://bioinf.man.ac.uk/dbbrowser/PRINTS/</a>
<a href="http://protein.toulouse.inra.fr/prodom.html">ProDom</a> - Protein domain db (Automatically generated)	<a href="http://protein.toulouse.inra.fr/prodom.html">http://protein.toulouse.inra.fr/prodom.html</a>
<a href="http://protomap.stanford.edu/">PROTOMAP</a> - Hierarchical classification of proteins	<a href="http://protomap.stanford.edu/">http://protomap.stanford.edu/</a>
<a href="http://www3.icgeb.trieste.it/~sbasesrv/">SBASE</a> - SBASE domain db	<a href="http://www3.icgeb.trieste.it/~sbasesrv/">http://www3.icgeb.trieste.it/~sbasesrv/</a>
<a href="http://smart.embl-heidelberg.de/">SMART</a> - Simple Modular Architecture Research Tool	<a href="http://smart.embl-heidelberg.de/">http://smart.embl-heidelberg.de/</a>
<a href="http://www.tigr.org/TIGRFAMs/">TIGRFAMs</a> - TIGR protein families db	<a href="http://www.tigr.org/TIGRFAMs/">http://www.tigr.org/TIGRFAMs/</a>
<a href="http://www.bind.ca/">BIND</a> - Biomolecular Interaction Network db	<a href="http://www.bind.ca/">http://www.bind.ca/</a>
<a href="http://dip.doe-mbi.ucla.edu/">DIP</a> - Db of Interacting Proteins	<a href="http://dip.doe-mbi.ucla.edu/">http://dip.doe-mbi.ucla.edu/</a>
<a href="http://cbm.bio.uniroma2.it/mint/">MINT</a> - Molecular INTeractions	<a href="http://cbm.bio.uniroma2.it/mint/">http://cbm.bio.uniroma2.it/mint/</a>
<a href="http://pronet.doubletivist.com/">ProNet</a> - Protein-Protein interaction db	<a href="http://pronet.doubletivist.com/">http://pronet.doubletivist.com/</a>

Table 4 Examples of DNA sequence databases

<i>Primary DNA sequence resources</i>	
<a href="http://www.ebi.ac.uk/embl/">EMBL</a> - EMBL Nucleotide sequence db (EBI)	<a href="http://www.ebi.ac.uk/embl/">http://www.ebi.ac.uk/embl/</a>
<a href="http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html">Genbank</a> - GenBank Nucleotide Sequence db (NCBI)	<a href="http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html">http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html</a>
<a href="http://www.ddbj.nig.ac.jp/">DDBJ</a> - DNA Data Bank of Japan	<a href="http://www.ddbj.nig.ac.jp/">http://www.ddbj.nig.ac.jp/</a>
<a href="http://www.ncbi.nlm.nih.gov/dbEST/">dbEST</a> - dbEST (Expressed Sequence Tags) db (NCBI)	<a href="http://www.ncbi.nlm.nih.gov/dbEST/">http://www.ncbi.nlm.nih.gov/dbEST/</a>
<a href="http://www.ncbi.nlm.nih.gov/dbSTS/">dbSTS</a> - dbSTS (Sequence Tagged Sites) db (NCBI)	<a href="http://www.ncbi.nlm.nih.gov/dbSTS/">http://www.ncbi.nlm.nih.gov/dbSTS/</a>
<i>Secondary DNA sequence resources</i>	
<a href="http://ndbserver.rutgers.edu/NDB/ndb.html">NDB</a> - Nucleic Acid Databank (3D structures)	<a href="http://ndbserver.rutgers.edu/NDB/ndb.html">http://ndbserver.rutgers.edu/NDB/ndb.html</a>
<a href="http://202.41.70.55/www/net/deva.html">BNASDB</a> - Nucleic acid structure db from University of Pune	<a href="http://202.41.70.55/www/net/deva.html">http://202.41.70.55/www/net/deva.html</a>
<a href="http://www.hgc.ims.u-tokyo.ac.jp/~knakai/asdb.html">AsDb</a> - Aberrant Splicing db	<a href="http://www.hgc.ims.u-tokyo.ac.jp/~knakai/asdb.html">http://www.hgc.ims.u-tokyo.ac.jp/~knakai/asdb.html</a>
<a href="http://pbil.univ-lyon1.fr/acuts/ACUTS.html">ACUTS</a> - Ancient conserved untranslated DNA sequences db	<a href="http://pbil.univ-lyon1.fr/acuts/ACUTS.html">http://pbil.univ-lyon1.fr/acuts/ACUTS.html</a>
<a href="http://www.kazusa.or.jp/codon/">Codon Usage Db</a>	<a href="http://www.kazusa.or.jp/codon/">http://www.kazusa.or.jp/codon/</a>
<a href="http://www.epd.isb-sib.ch/">EPD</a> - Eukaryotic Promoter db	<a href="http://www.epd.isb-sib.ch/">http://www.epd.isb-sib.ch/</a>
<a href="http://pbil.univ-lyon1.fr/databases/hoovergen.html">HOVERGEN</a> - Homologous Vertebrate Genes db	<a href="http://pbil.univ-lyon1.fr/databases/hoovergen.html">http://pbil.univ-lyon1.fr/databases/hoovergen.html</a>
<a href="http://www.introns.com/">ISIS</a> - Intron Sequence and Information System	<a href="http://www.introns.com/">http://www.introns.com/</a>
<a href="http://rdp.cme.msu.edu/html/">RDP</a> - Ribosomal db Project	<a href="http://rdp.cme.msu.edu/html/">http://rdp.cme.msu.edu/html/</a>
<a href="http://biosun.bio.tu-darmstadt.de/goringer/gRNA/gRNA.html">gRNAs db</a> - Guide RNA db	<a href="http://biosun.bio.tu-darmstadt.de/goringer/gRNA/gRNA.html">http://biosun.bio.tu-darmstadt.de/goringer/gRNA/gRNA.html</a>
<a href="http://www.dna.affrc.go.jp/htdocs/PLACE/">PLACE</a> - Plant cis-acting regulatory DNA elements db	<a href="http://www.dna.affrc.go.jp/htdocs/PLACE/">http://www.dna.affrc.go.jp/htdocs/PLACE/</a>
<a href="http://sphinx.rug.ac.be:8080/PlantCARE/">PlantCARE</a> - Plant cis-acting regulatory DNA elements db	<a href="http://sphinx.rug.ac.be:8080/PlantCARE/">http://sphinx.rug.ac.be:8080/PlantCARE/</a>
<a href="http://mber.bcm.tmc.edu/smallRNA/smallrna.html">sRNA db</a> - Small RNA db	<a href="http://mber.bcm.tmc.edu/smallRNA/smallrna.html">http://mber.bcm.tmc.edu/smallRNA/smallrna.html</a>
<a href="http://rrna.uia.ac.be/rrna/ssu/">ssu rRNA</a> - Small ribosomal subunit db	<a href="http://rrna.uia.ac.be/rrna/ssu/">http://rrna.uia.ac.be/rrna/ssu/</a>
<a href="http://rrna.uia.ac.be/rrna/lru/">lsu rRNA</a> - Large ribosomal subunit db	<a href="http://rrna.uia.ac.be/rrna/lru/">http://rrna.uia.ac.be/rrna/lru/</a>
<a href="http://rose.man.poznan.pl/5SData/">5S rRNA</a> - 5S ribosomal RNA db	<a href="http://rose.man.poznan.pl/5SData/">http://rose.man.poznan.pl/5SData/</a>
<a href="http://www.indiana.edu/~tmrna/">tmRNA Website</a>	<a href="http://www.indiana.edu/~tmrna/">http://www.indiana.edu/~tmrna/</a>
<a href="http://psyche.uthct.edu/dbs/tmRDB/tmRDB.html">tmRDB</a> - tmRNA dB	<a href="http://psyche.uthct.edu/dbs/tmRDB/tmRDB.html">http://psyche.uthct.edu/dbs/tmRDB/tmRDB.html</a>
<a href="http://www.uni-bayreuth.de/departments/biochemie/sprinzl/trna/">tRNA</a> - tRNA compilation from the University of Bayreuth	<a href="http://www.uni-bayreuth.de/departments/biochemie/sprinzl/trna/">http://www.uni-bayreuth.de/departments/biochemie/sprinzl/trna/</a>
<a href="http://psyche.uthct.edu/dbs/uRNADB/uRNADB.html">uRNADB</a> - uRNA db	<a href="http://psyche.uthct.edu/dbs/uRNADB/uRNADB.html">http://psyche.uthct.edu/dbs/uRNADB/uRNADB.html</a>
<a href="http://www.lifesci.ucla.edu/RNA">RNA editing</a> - RNA editing site	<a href="http://www.lifesci.ucla.edu/RNA">http://www.lifesci.ucla.edu/RNA</a>
<a href="http://medstat.med.utah.edu/RNAmods/">RNAmdb</a> - RNA modification db	<a href="http://medstat.med.utah.edu/RNAmods/">http://medstat.med.utah.edu/RNAmods/</a>
<a href="http://gifts.univ-mrs.fr/SOS-DGDB/SOS-DGDB_home_page.html">SOS-DGDB</a> - Db of Drosophila DNA annotated with regulatory binding sites	<a href="http://gifts.univ-mrs.fr/SOS-DGDB/SOS-DGDB_home_page.html">http://gifts.univ-mrs.fr/SOS-DGDB/SOS-DGDB_home_page.html</a>
<a href="http://www.genlink.wustl.edu/telldb/index.html">TelDB</a> - Multimedia Telomere Resource	<a href="http://www.genlink.wustl.edu/telldb/index.html">http://www.genlink.wustl.edu/telldb/index.html</a>
TRADAT - TRANscription Databases and Analysis Tools	<a href="http://www.itba.mi.cnr.it/tradat/">http://www.itba.mi.cnr.it/tradat/</a>
Subviral RNA db - Small circular RNAs db (viroid and viroid-like)	<a href="http://nt.ars-grin.gov/subviral/">http://nt.ars-grin.gov/subviral/</a>
<a href="http://www.biotech.ist.unige.it/interlab/mpdb.html">MPDB</a> - Molecular probe db	<a href="http://www.biotech.ist.unige.it/interlab/mpdb.html">http://www.biotech.ist.unige.it/interlab/mpdb.html</a>
<a href="http://www.cme.msu.edu/OPD/">OPD</a> - Oligonucleotide probe db	<a href="http://www.cme.msu.edu/OPD/">http://www.cme.msu.edu/OPD/</a>

Table 5 Examples of three-dimensional databases

<b><i>Primary 3-D resources</i></b>	
Protein Data Bank	<a href="http://www.rcsb.org">http://www.rcsb.org</a>
Macromolecular Structure Database	<a href="http://www.ebi.ac.uk/msd/index.html">http://www.ebi.ac.uk/msd/index.html</a>
Nucleic Acid Database Project	<a href="http://ndbserver.rutgers.edu/NDB/index.html">http://ndbserver.rutgers.edu/NDB/index.html</a>
BioMagResBank	<a href="http://www.bmrwisc.edu/Welcome.html">http://www.bmrwisc.edu/Welcome.html</a>
<b><i>Protein domain/fold databases</i></b>	
3Dee	<a href="http://jura.ebi.ac.uk:8080/3Dee/help/help_intro.html">http://jura.ebi.ac.uk:8080/3Dee/help/help_intro.html</a>
CATH	<a href="http://www.biochem.ucl.ac.uk/bsm/cath">http://www.biochem.ucl.ac.uk/bsm/cath</a>
HSSP	<a href="http://www.sander.ebi.ac.uk/hssp">http://www.sander.ebi.ac.uk/hssp</a>
SCOP	<a href="http://scop.mrc-lmb.cam.ac.uk/scop">http://scop.mrc-lmb.cam.ac.uk/scop</a>
<b><i>Examples of specialized resources</i></b>	
BIND - binding database	<a href="http://www.bind.ca/index.phtml?page=databases">http://www.bind.ca/index.phtml?page=databases</a>
BindingDB – binding database	<a href="http://www.bindingdb.org/bind/index.jsp">http://www.bindingdb.org/bind/index.jsp</a>
Decoys ‘R’ Us	<a href="http://dd.stanford.edu">http://dd.stanford.edu</a>
Disordered structures	<a href="http://bonsai.ims.u-tokyo.ac.jp/~klsim/database.html">http://bonsai.ims.u-tokyo.ac.jp/~klsim/database.html</a>
DNA binding proteins	<a href="http://ndbserver.rutgers.edu/structure-finder/dnabind/">http://ndbserver.rutgers.edu/structure-finder/dnabind/</a>
Intramolecular movements	<a href="http://molmovdb.mbb.yale.edu/MolMovDB/">http://molmovdb.mbb.yale.edu/MolMovDB/</a>
Loops	<a href="http://www-cryst.bioc.cam.ac.uk/~sloop/Info.html">http://www-cryst.bioc.cam.ac.uk/~sloop/Info.html</a>
Membrane protein structures	<a href="http://blanco.biomol.uci.edu/Membrane_Proteins_xtal.html">http://blanco.biomol.uci.edu/Membrane_Proteins_xtal.html</a>
Metal cations	<a href="http://metallo.scripps.edu/">http://metallo.scripps.edu/</a>
P450 containing systems	<a href="http://www.icgeb.trieste.it/~p450srv/">http://www.icgeb.trieste.it/~p450srv/</a>
Predicted protein models	<a href="http://guitar.rockefeller.edu/modbase">http://guitar.rockefeller.edu/modbase</a>
Protein-DNA contacts	<a href="http://www.biochem.ucl.ac.uk/bsm/DNA/server/">http://www.biochem.ucl.ac.uk/bsm/DNA/server/</a>
Protein-protein interfaces	<a href="http://www.biochem.ucl.ac.uk/bsm/PP/server/">http://www.biochem.ucl.ac.uk/bsm/PP/server/</a>
ProTherm	<a href="http://www.rtc.riken.go.jp/jouhou/protherm/protherm.html">http://www.rtc.riken.go.jp/jouhou/protherm/protherm.html</a>
Quaternary structure	<a href="http://pgs.ebi.ac.uk">http://pgs.ebi.ac.uk</a>
Small ligands	<a href="http://alpha2.bmc.uu.se/hicup/">http://alpha2.bmc.uu.se/hicup/</a>
Small ligands	<a href="http://www.ebi.ac.uk/msd-srv/chempdb">http://www.ebi.ac.uk/msd-srv/chempdb</a>
The Protein Kinase Resource	<a href="http://pkr.sdsc.edu/html/index.shtml">http://pkr.sdsc.edu/html/index.shtml</a>
<b><i>Examples of search/retrieval facilities and database interfaces</i></b>	
3DinSight – structure/function dbase	<a href="http://www.rtc.riken.go.jp/jouhou/3dinsight/3DinSight.html">http://www.rtc.riken.go.jp/jouhou/3dinsight/3DinSight.html</a>
BioMolQuest – structure/function dbase	<a href="http://bioinformatics.danforthcenter.org/yury/public/home.html">http://bioinformatics.danforthcenter.org/yury/public/home.html</a>
Entrez	<a href="http://www3.ncbi.nlm.nih.gov/entrez/query.fcgi">http://www3.ncbi.nlm.nih.gov/entrez/query.fcgi</a>
Image Library of Macromolecules	<a href="http://www.imb-jena.de/IMAGE.html">http://www.imb-jena.de/IMAGE.html</a>
OCA	<a href="http://bioinfo.weizmann.ac.il:8500/oca-docs/">http://bioinfo.weizmann.ac.il:8500/oca-docs/</a>
PDBSUM	<a href="http://www.biochem.ucl.ac.uk/bsm/pdbsum/">http://www.biochem.ucl.ac.uk/bsm/pdbsum/</a>
ProNIT – protein/nucleic acid interactions	<a href="http://www.rtc.riken.go.jp/jouhou/pronit/pronit.html">http://www.rtc.riken.go.jp/jouhou/pronit/pronit.html</a>
TargetDB	<a href="http://targetdb.pdb.org/">http://targetdb.pdb.org/</a>
SRS	<a href="http://srs.ebi.ac.uk/">http://srs.ebi.ac.uk/</a>

Table 4 Examples of genomic resources

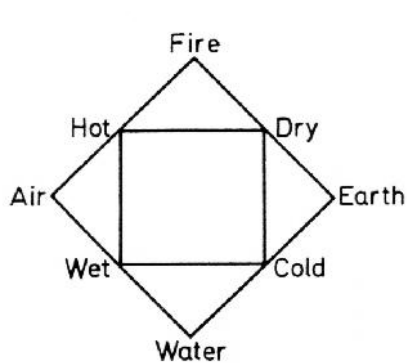
Genomic databases for various organisms <sup>1</sup>	
Flybase - <i>Drosophila melanogaster</i>	<a href="http://flybase.bio.indiana.edu/">http://flybase.bio.indiana.edu/</a>
Subtilist - <i>Bacillus subtilis</i>	<a href="http://genolist.pasteur.fr/SubtiList/">http://genolist.pasteur.fr/SubtiList/</a>
Cyanobase - <i>Synechocystis</i> strain PCC6803	<a href="http://www.kazusa.or.jp/cyano/cyano.html">http://www.kazusa.or.jp/cyano/cyano.html</a>
CYGD - <i>Sacharomyces cerevisiae</i>	<a href="http://mips.gsf.de/genre/proj/yeast/">http://mips.gsf.de/genre/proj/yeast/</a>
ENSEMBLE – human and other invertebrate genomes	<a href="http://www.ensembl.org/">http://www.ensembl.org/</a>
<i>Comparative genomic visualization tools</i>	
VISTA	<a href="http://www-gsd.lbl.gov/vista/">http://www-gsd.lbl.gov/vista/</a>
PipMaker	<a href="http://bio.cse.psu.edu/pipmaker/">http://bio.cse.psu.edu/pipmaker/</a>
Whole-genome annotation browsers	
NCBI Map Viewer	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
UCSC Genome Browser	<a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>
Ensembl	<a href="http://www.ensembl.org/">http://www.ensembl.org/</a>
<i>Whole-genome comparative genomic browsers</i>	
UCSC Genome Browser	<a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>
VISTA Genome Browser	<a href="http://pipeline.lbl.gov/">http://pipeline.lbl.gov/</a>
PipMaker	<a href="http://bio.cse.psu.edu/genome/hummus/">http://bio.cse.psu.edu/genome/hummus/</a>
<i>Custom comparisons to whole genomes</i>	
GenomeVista (AVID)	<a href="http://pipeline.lbl.gov/cgi-bin/GenomeVista">http://pipeline.lbl.gov/cgi-bin/GenomeVista</a>
UCSC Genome Browser (BLAT)	<a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>
ENSEMBL (SSAHA)	<a href="http://www.ensembl.org/">http://www.ensembl.org/</a>
NCBI (BLAST)	<a href="http://www.ncbi.nlm.nih.gov/blast/">http://www.ncbi.nlm.nih.gov/blast/</a>

<sup>1</sup>A more complete list is available at the websites of the EBI, NCBI and DDBJ as well as within the current database issue of Nucleic Acids research, cited in Table 1.



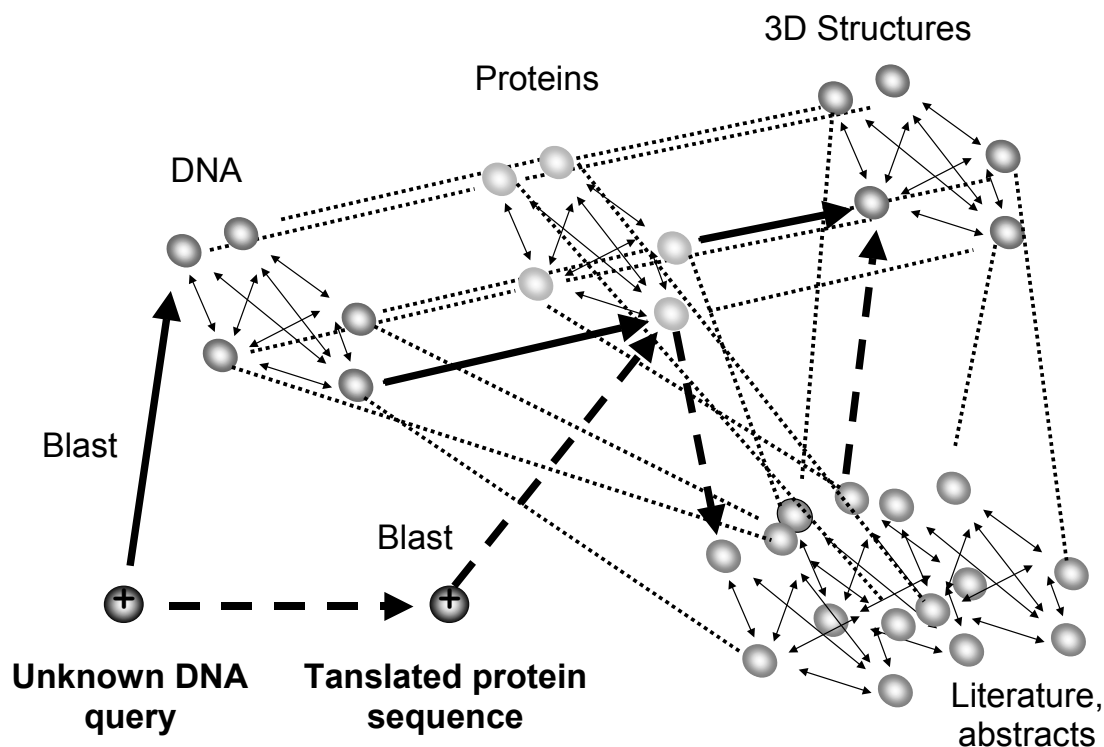
Table 5 Bioinformatics resources for immunology

ALPSbase Autoimmune lymphoproliferative syndrome database	<a href="http://research.nhgri.nih.gov/alps/">http://research.nhgri.nih.gov/alps/</a>
BCIpep Experimentally determined B-cell epitopes of antigenic proteins	<a href="http://bioinformatics.uams.edu/mirror/bcipep/">http://bioinformatics.uams.edu/mirror/bcipep/</a>
dbMHC HLA sequences in human populations	<a href="http://www.ncbi.nlm.nih.gov/mhc/">http://www.ncbi.nlm.nih.gov/mhc/</a>
FIMM Functional molecular immunology, T-cell response to disease-specific antigens.	<a href="http://research.i2r.a-star.edu.sg/fimm/">http://research.i2r.a-star.edu.sg/fimm/</a>
HaptenDB Hapten molecules	<a href="http://www.imtech.res.in/raghava/haptendb/">http://www.imtech.res.in/raghava/haptendb/</a>
HLA Ligand/Motif database A database and search tool for HLA sequences	<a href="http://hlaligand.ouhsc.edu/">http://hlaligand.ouhsc.edu/</a>
IL2Rgbase X-linked severe combined immunodeficiency mutations	<a href="http://research.nhgri.nih.gov/scid/">http://research.nhgri.nih.gov/scid/</a>
IMGT Integrated knowledge resource	<a href="http://imgt.cines.fr">http://imgt.cines.fr</a>
IMGT-GENE-DB Genome database for immunoglobulins (IG) and T cell receptors (TR) genes from human and mouse	<a href="http://imgt.cines.fr/cgi-bin/GENElect.jv">http://imgt.cines.fr/cgi-bin/GENElect.jv</a>
IMGT/HLA HLA sequence database	<a href="http://www.ebi.ac.uk/imgt/hla/">http://www.ebi.ac.uk/imgt/hla/</a>
IMGT/LIGM-DB Immunoglobulin (IG) and T cell receptor (TR) nucleotide sequences, from human and other vertebrate species,	<a href="http://imgt.cines.fr/cgi-bin/IMGTlect.jv">http://imgt.cines.fr/cgi-bin/IMGTlect.jv</a>
Interferon Stimulated Gene Database (by microarray)	<a href="http://www.lerner.ccf.org/labs/williams/xchip-html.cgi">http://www.lerner.ccf.org/labs/williams/xchip-html.cgi</a>
IPD-ESTDAB Polymorphic genes in the immune system	<a href="http://www.ebi.ac.uk/ipd/estdab/">http://www.ebi.ac.uk/ipd/estdab/</a>
IPD-HPA - Human Platelet Antigens	<a href="http://www.ebi.ac.uk/ipd/hpa/">http://www.ebi.ac.uk/ipd/hpa/</a>
IPD-KIR - Killer-cell Immunoglobulin-like Receptors	<a href="http://www.ebi.ac.uk/ipd/kir/">http://www.ebi.ac.uk/ipd/kir/</a>
IPD-MHC Polymorphic genes in the immune system	<a href="http://www.ebi.ac.uk/ipd/mhc">http://www.ebi.ac.uk/ipd/mhc</a>
JenPep Peptide binding to biomacromolecules within immunobiology (epitopes)	<a href="http://www.jenner.ac.uk/Jenpep">http://www.jenner.ac.uk/Jenpep</a>
<a href="#">Kabat</a> - Kabat db of sequences of proteins of immunological interest	<a href="http://immuno.bme.nwu.edu/">http://immuno.bme.nwu.edu/</a>
MHC-Peptide Interaction Database	<a href="http://surya.bic.nus.edu.sg/mpid">http://surya.bic.nus.edu.sg/mpid</a>
MHCBN Peptides binding to MHC or TAP	<a href="http://www.imtech.res.in/raghava/mhcbn/">http://www.imtech.res.in/raghava/mhcbn/</a>
MHCPEP MHC binding peptides	<a href="http://wehih.wehi.edu.au/mhcpep/">http://wehih.wehi.edu.au/mhcpep/</a>
VBASE2 G $\alpha$ erm-line V genes from the immunoglobulin loci of human and mouse	<a href="http://www.vbase2.org">http://www.vbase2.org</a>

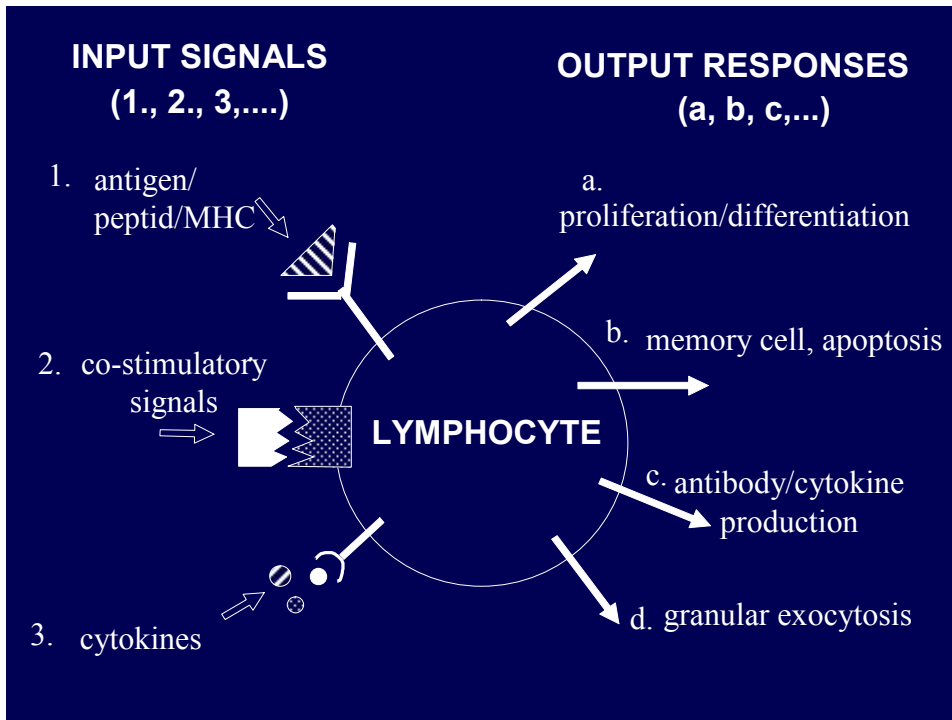


System	Entities	Relationships
Molecules	Atoms	Atomic interactions (chemical bonds)
Assemblies	Proteins, DNA	Molecular contacts
Pathways	Enzymes	Chemical reactions (substrates/products)
Genetic networks	Genes	Co-regulation
Protein structure	Atoms	Chemical bonds
Simplified rotein structure	Secondary structures	Sequential and topological vicinity
Folds	C <sub>α</sub> atoms	3-D vicinity
Protein sequence	Amino acid	Sequential vicinity

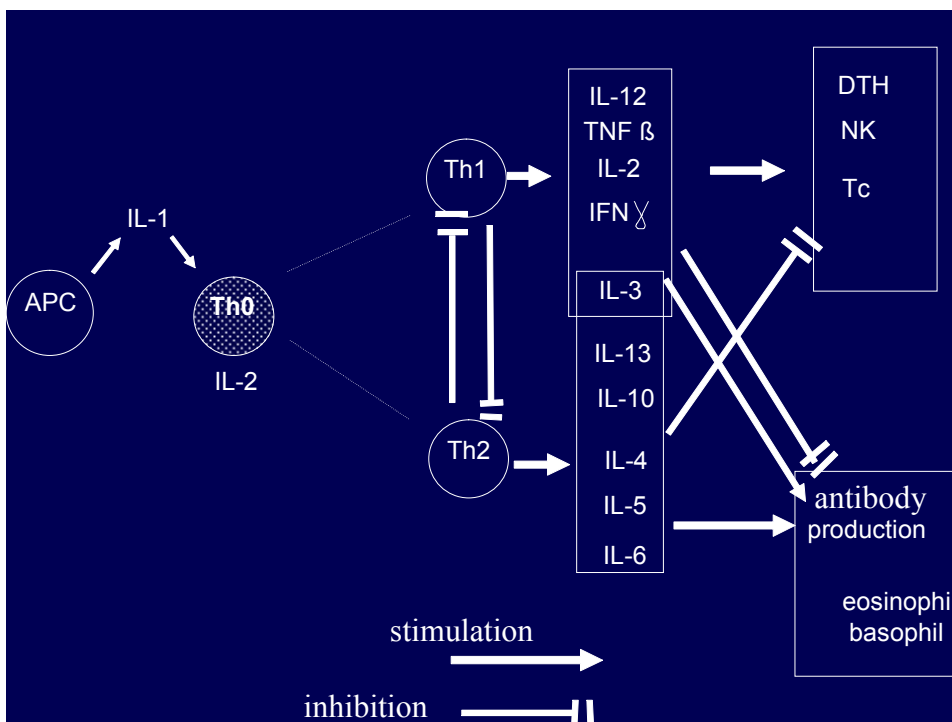
**Figure 1.** Examples of entities and relationships used in molecular models. Left: A simple fourfold semantics (earth, water, air, fire) and binary relations (hot-cold, wet/dry) was used in most early cultures as a conceptual framework to describe Nature. Right: Entities and relationships of modern databases.



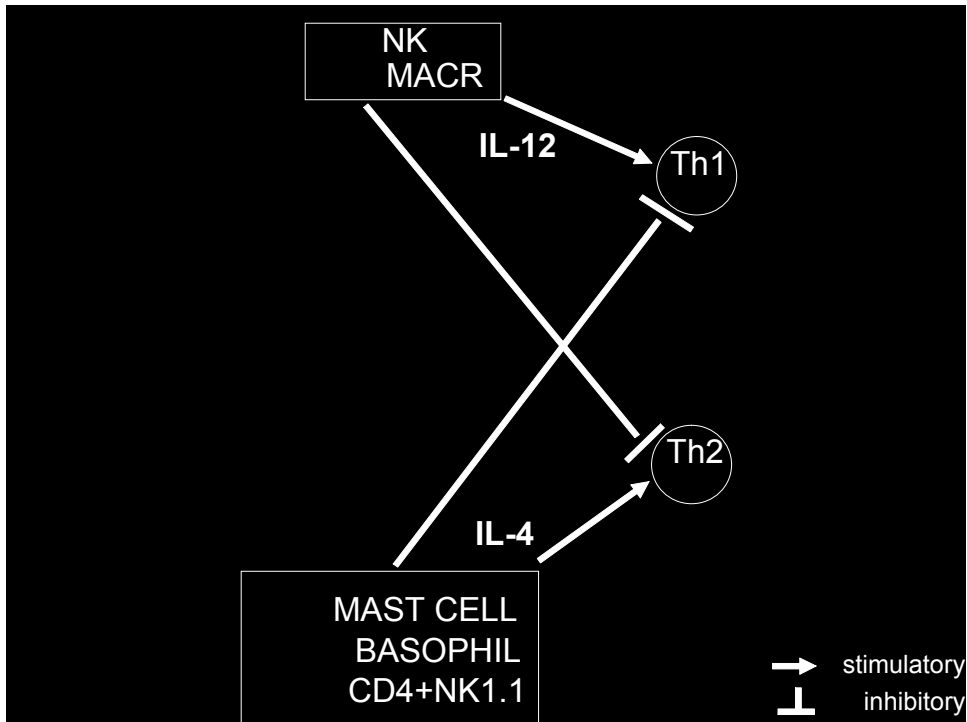
**Figure 2.** Search on an integrated database. Items in the individual databases (DNA sequence, Protein Sequence, 3-D structures, Literature abstracts) are cross-referenced (dotted line) by WWW links. Additional links (thin arrows) connect “neighbourhoods” i.e. similar data-items within each database. Consequently, if similarity search (thick arrows) points e.g. to an unannotated DNA entry, a member of its neighbourhood may help the user to find proteins, or protein structures.



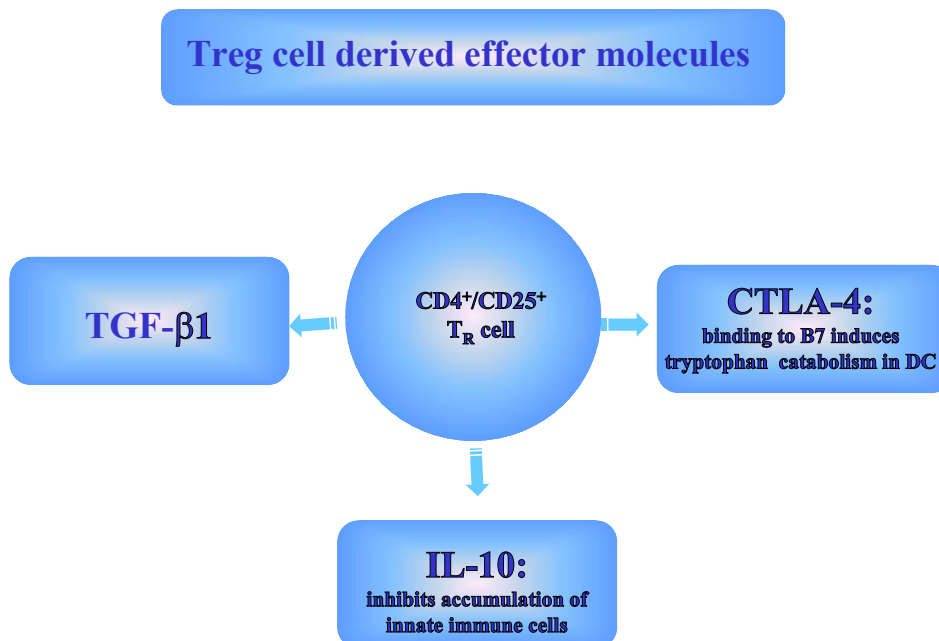
**Figure 3.** Information processing in lymphocytes.



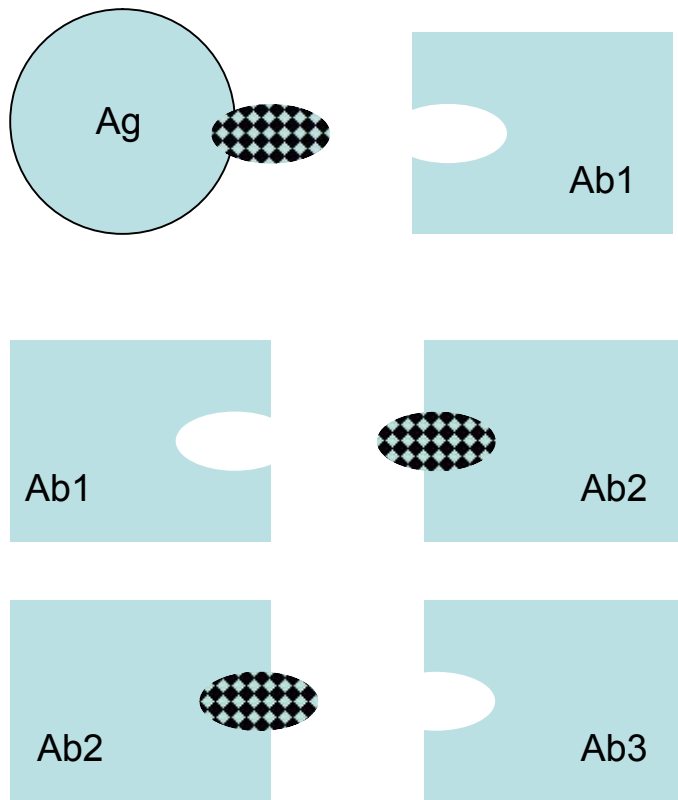
**Figure 4.** Cytokine patterns delivered by T cell polarization. Th1 cells tends to activate cell-mediated and suppress humoral immunity while Th2 cytokines have the opposite effect.



**Figure 5.** Innate immunity and T1/T2-like cytokine polarization. NK cells and macrophages produce preferentially “Th1 –like” cytokines, while mast cells, basophils and NK 1.1 cells secrete “Th2-like” cytokines



**Figure 6.** T regulatory cells expressing the negative costimulatory molecule CTLA-4 and secreting inhibitory cytokine TGF-β1 and IL-10.



**Figure 7.** Idiotypic web. The antigen (Ag) is recognized by a specific antibody (Ab1) carrying idiotope stimulating a second antibody (Ab2). Idiotypic determinant/s/ of Ab2 may further activate Abs, etc. Parts of antigenic structure may be closely similar to that of Ab2.

## References

- ALTSCHUL, S. F., GISH, W., MILLER, W., MYERS, E. W. & LIPMAN, D. J. (1990) Basic local alignment search tool. *J Mol Biol*, 215, 403-10.
- ASHBURNER, M., BALL, C. A., BLAKE, J. A., BOTSTEIN, D., BUTLER, H., CHERRY, J. M., DAVIS, A. P., DOLINSKI, K., DWIGHT, S. S., EPPIG, J. T., HARRIS, M. A., HILL, D. P., ISSELTARVER, L., KASARSKIS, A., LEWIS, S., MATESE, J. C., RICHARDSON, J. E., RINGWALD, M., RUBIN, G. M. & SHERLOCK, G. (2000) Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet*, 25, 25-9.
- BERMAN, H. M., WESTBROOK, J., FENG, Z., GILLILAND, G., BHAT, T. N., WEISSIG, H., SHINDYALOV, I. N. & BOURNE, P. E. (2000) The Protein Data Bank. *Nucleic Acids Res*, 28, 235-42.
- BERNSTEIN, F. C., KOETZLE, T. F., WILLIAMS, G. J., MEYER, E. F., JR., BRICE, M. D., RODGERS, J. R., KENNARD, O., SHIMANOUCI, T. & TASUMI, M. (1977) The Protein Data Bank: a computer-based archival file for macromolecular structures. *J Mol Biol*, 112, 535-42.
- BRUSIC, V. PETROVSKY, N. (2005) Immunoinformatics and its relevance to understanding human immune disease. In press
- CARNAP, R. (1939) *Foundations of Logics and Mathematics*, Chicago, University of Chicago Press.
- CASTAGNETTO, J. M., HENNESSY, S. W., ROBERTS, V. A., GETZOFF, E. D., TAINER, J. A. & PIQUE, M. E. (2002) MDB: the metalloprotein database and browser at the Scripps Research Institute. *Nucleic Acids Res*, 30, 379-382.
- CHOMSKY, N. (1957) *Syntactic Structures*, The Hague, Mouton.
- CSÁNYI, V. (1989) *Evolutionary Systems and Society*, Durham and London, Duke University Press.
- DE GROOT, A. S. (2004) Immunome derived vaccines. *Exp. Opin. Biol. Ther.*, 4, 767-772.
- DEVEREAUX, J., HAEBERLI, P. & O, S. (1984) A comprehensive set of sequence analysis programs for the VAX. *Nucleic Acids Research*, 12, 387-395.
- DONATE, L. E., RUFINO, S. D., CANARD, L. H. & BLUNDELL, T. L. (1996) Conformational analysis and clustering of short and medium size loops connecting regular secondary structures: a database for modeling and prediction. *Protein Sci*, 5, 2600-16.
- FRAUWIRTH, K. A. & THOMPSON, C. B. (2002) Activation and inhibition of lymphocytes by costimulation. *J Clin Invest*, 109, 295-9.
- FREY, O., PETROW, P. K., GAJDA, M., SIEGMUND, K., HUEHN, J., SCHEFFOLD, A., HAMANN, A., RADBRUCH, A. & BRAUER, R. (2005) The role of regulatory T cells in antigen-induced arthritis: aggravation of arthritis after depletion and amelioration after transfer of CD4+CD25+ T cells. *Arthritis Res Ther*, 7, R291-301.
- GENTLEMAN, R. C., CAREY, V. J., BATES, D. M., BOLSTAD, B., DETTLING, M., DUDOIT, S., ELLIS, B., GAUTIER, L., GE, Y., GENTRY, J., HORNIK, K., HOTHORN, T., HUBER, W., IACUS, S., IRIZARRY, R., LEISCH, F., LI, C., MAECHLER, M., ROSSINI, A. J., SAWITZKI, G., SMITH, C., SMYTH, G., TIERNEY, L., YANG, J. Y. & ZHANG, J. (2004) Bioconductor: open software development for computational biology and bioinformatics. *Genome Biol*, 5, R80.
- GODFREY, D. I. & KRONENBERG, M. (2004) Going both ways: immune regulation via CD1d-dependent NKT cells. *J Clin Invest*, 114, 1379-88.
- HENRICK, K. & THORNTON, J. M. (1998) PQS: a protein quaternary structure resource. *Trends Biochem Sci*, 23, 358-361.
- JIANG, H. & CHESS, L. (2004) An integrated model of immunoregulation mediated by regulatory T cell subsets. *Adv Immunol*, 83, 253-88.
- JONES, S. & THORNTON, J. M. (1996) Principles of protein-protein interactions. *Proc Natl Acad Sci U S A*, 93, 13-20.
- KAMPIS, G. (1991) *Self-modifying systems in Biology and Cognitive Science*, Oxford, New York, Pergamon Press.
- KIM, C. H., ROTT, L., KUNKEL, E. J., GENOVESE, M. C., ANDREW, D. P., WU, L. & BUTCHER, E. C. (2001) Rules of chemokine receptor association with T cell polarization in vivo. *J Clin Invest*, 108, 1331-9.
- KLEYWEGT, G. J. & JONES, D. T. (1998) Databases in protein crystallography. *Acta Cryst. D*, 54, 119-1131.
- KNUTH, D. E. (1998) *The Art of Computer Programming*, Addison-Wesley.
- KONOPKA, A. K. (1994) Fundamentals of Biomolecular Cryptology. IN SMITH, D. W. (Ed.) *Biocomputing: Informatics and Genome Projects*. San Diego, New York, Boston, Academic Press.

- LUSCOMBE, N. M., AUSTIN, S. E., BERMAN, H. M. & THORNTON, J. M. (2000) An overview of the structures of protein-DNA complexes. *Genome Biol*, 1.
- MAYNARD SMITH, J. & SZATHMÁRY, E. (1995) *The Major Transitions in Evolution*, Oxford, New York, Heidelberg, W.H. Freeman.
- MURZIN, A. G., BRENNER, S. E., HUBBARD, T. & CHOTHIA, C. (1995) SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J Mol Biol*, 247, 536-40.
- ORENGO, C. A., MICHIE, A. D., JONES, S., JONES, D. T., SWINDELLS, M. B. & THORNTON, J. M. (1997) CATH--a hierarchic classification of protein domain structures. *Structure*, 5, 1093-108.
- PINKER, S. (2001) *The Language Instinct: How the Mind Creates Language*, Perennial.
- POLJAK, R. J. (1994) An idiotope--anti-idiotope complex and the structural basis of molecular mimicking. *Proc Natl Acad Sci U S A.*, 91, 1599-1600.
- PONGOR, S. (1988) Novel databases for molecular biology. *Nature*, 332, 24.
- PONGOR, S., SKERL, V., CSERZŐ, M., HÁTSÁGI, Z., SIMON, G. & BEVILACQUA, V. (1993) The SBASE domain library: a collection of annotated protein segments. *Protein Eng*, 6, 391-5.
- RICE, P., LONGDEN, I. & BLEASBY, A. (2000) EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet*, 16, 276-7.
- RIPLEY, B. D. (1999) *Pattern Recognition and Neural Networks*, Cambridge, Cambridge University Press.
- SADRA, A., CINEK, T. & IMBODEN, J. B. (2004) Translocation of CD28 to lipid rafts and costimulation of IL-2. *Proc Natl Acad Sci U S A*, 101, 11422-7.
- SANDER, C. & SCHNEIDER, R. (1991) Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins*, 9, 56-68.
- SCHULZE-KREMER, S. (1997) Adding semantics to genome databases: towards an ontology for molecular biology. *Proc Int Conf Intell Syst Mol Biol*, 5, 272-5.
- SHANNON, C. E. (1948a) Communication theory of secrecy systems. *Bell Syst. Tech. J.*, 29.
- SHANNON, C. E. (1948b) A mathematical theory of communication II. *Bell Syst. Tech. J.*, 27.
- SHANNON, C. E. (1948c) Prediction and entropy of printed English. *Bell Syst. Tech. J.*, 30.
- SIDDIQUI, A. S., DENGLER, U. & BARTON, G. J. (2001) 3Dee: A database of protein structural domains. *Bioinformatics*, 17, 200-201.
- SIM, K. L., UCHIDA, T. & MIYANO, S. (2001) ProDDO: a database of disordered proteins from the Protein Data Bank (PDB). *Bioinformatics*, 17, 379-380.
- STAJICH, J. E., BLOCK, D., BOULEZ, K., BRENNER, S. E., CHERVITZ, S. A., DAGDIGIAN, C., FUELLEN, G., GILBERT, J. G., KORF, I., LAPP, H., LEHVASLAIHO, H., MATSALLA, C., MUNGALL, C. J., OSBORNE, B. I., POCOCK, M. R., SCHATTNER, P., SENGER, M., STEIN, L. D., STUPKA, E., WILKINSON, M. D. & BIRNEY, E. (2002) The Bioperl toolkit: Perl modules for the life sciences. *Genome Res*, 12, 1611-8.
- SZATHMÁRY, E. & SMITH, J. M. (1995) The major evolutionary transitions. *Nature*, 374, 227-32.
- WALSH, P. T., TAYLOR, D. K. & TURKA, L. A. (2004) Tregs and transplantation tolerance. *J Clin Invest*, 114, 1398-1403.
- WANG, E. & FALUS, A. (2004) Changing paradigm through a genome-based approach to clinical and basic immunology. *J Transl Med*, 2, 2.
- WESTBROOK, J. D. & BOURNE, P. E. (2000) STAR/mmCIF: an ontology for macromolecular structure. *Bioinformatics*, 16, 159-68.
- WILSON, M. J., TORKAR, M., HAUDE, A., MILNE, S., JONES, T., SHEER, D., BECK, S. & TROWSDALE, J. (2000) Plasticity in the organization and sequences of human KIR/ILT gene families. *Proc Natl Acad Sci U S A*, 97, 4778-83.
- WITTGENSTEIN, L. (1922) *Tractatus Logico Philosophicus*, London, Routledge and Kegan Paul.
- ZHANG, X., SCHWARTZ, J. C., ALMO, S. C. & NATHENSON, S. G. (2003) Crystal structure of the receptor-binding domain of human B7-2: insights into organization and signaling. *Proc Natl Acad Sci U S A*, 100, 2586-91.