

# Immersive Graph-based Visualization and Exploration of Biological Data Relationships

N. Férey, P.-E. Gros, J. Hérisson, R. Gherbi

Bioinformatics team, Human-Computer Communication Department  
LIMSI-CNRS, University Paris-Sud, BP 133 F-91403 ORSAY CEDEX (France)

phone: +33 (0)1 69 85 81 64 fax: +33 (0)1 69 85 80 88

genoteam@limsi.fr

## Abstract

*Genomic information shows some characteristics that make them very difficult to interpret and to exploit. Such data constitute an important factual resource (GenBank, SwissProt, GeneOntology, or Decryphon...), are heterogeneous, huge in quantity, and are geographically distributed. They are also recorded in structured or semi-structured formats within public or private databanks. Nevertheless, genome knowledge could not be limited to DNA or protein annotated sequences. Indeed, there is a significant quantity of information relating to these genes, recorded in an unstructured format within millions publications (PubMed, Medline). This paper presents Genome3DExplorer, a new modeling and software solution to explore textual and factual genomic data based on adapted federator description language. Moreover, Genome3DExplorer offers us a user-friendly visualization of this information within an immersive environment. The visualization is based on a well-adapted graphical paradigm that automatically helps to build a graph-based representation. This solution allows biologist to explore suitably huge sets of genomic data, but it could be applied to other application fields. This kind of graphical exploration has the advantage to highlight some global topological characteristics, which are uneasily visible using traditional exploration tools. Finally, some results produced by Genome3DExplorer software on various sets of biological data, are presented.*

## Introduction

As in *ADN-Viewer* [1] and *SequenceWord* [2], our objective is to elaborate new solution in order to explore virtually various kinds of genomic data. These data come from the many databases, such as *GenBank* [12], *SwissProt* [11] or *Decryphon* [10]. Our approach is mainly based on the definition of a genomic data federator language, answering the requirements and specificities of genomic databases. Then we explain the representation methods to view these data within an immersive framework. Finally, we present some results produced by our *Genome3DExplorer* software on various sets of biological data.

In order to visualize efficiently biological data, we need to define a common data description language that must accommodate and represent knowledge resulting from structured but heterogeneous databanks. We describe in this section how we used the specific characteristics of the genomic data to find an adapted description format for this kind of data.

### 1. Biological databank specificities

Although the genomic databases are very heterogeneous (format or quality), they involve some specific characteristics. Indeed, they are often focused on biological object of interest (protein, gene...), described by an attribute set. Moreover, these objects are often compared one to another by a measurement (sequence alignment score, functional similarity...). For instance, *GenBank* contains annotated DNA

sequences, and provides *BLAST* tools in order to compare these sequences, as *SwissProt*, which deals with annotated protein sequences.

## 1.1 Definition of a genomic data representation language

In the most commonly cases, different kinds of biological objects (protein sequences, DNA sequences, biological terms...) are often connected by binary relationships. For example, using text corpora, biologists could extract co-occurrence relationship between two biological terms, as in *BioBiblioMetrics* [3], or more specific semantic relationships, coming from text information extraction processing [4] (like protein inhibition...). In databank case, biologist can extract alignment measurements between two DNA or protein sequences. These binary relationships can be valuated, by alignments or co-occurrence measurements (numerical values), or by semantic relations extracted from texts (symbolic values). Biological objects can also be valuated, by their biological properties. The values can be of string type (label, sequence...), numerical one (like co-occurrence score, alignment measurements), or symbolic one (type of interaction, like positive retroaction...). Taking account into these characteristics, we define a XML-based data representation language, based on the concept of multi-valuated objects and relationships, which is particularly adapted to describe biological data. In this uncompleted example, the studied biological objects are yeast genes. Two values characterize these genes. The name (value0 of object tag) and the number of co-regulator factors (value1 of object tag). Two values characterize the relationships between two genes, the kind of relationship (value1 of relation tag), and a numerical value (value0 of relation tag).

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<data>
<object id="1" value0="YCR107W" value1="3"/>
<object id="2" value0="YCR106W" value1="3"/>
<object id="3" value0="YCR105W" value1="3"/>
<object id="4" value0="YCR104W" value1="1"/>
<object id="5" value0="YCR102W-A" value1="2"/>
<object id="6" value0="YCR102C" value1="2"/>
<object id="7" value0="YCR101C" value1="1"/>
...
<relation id1="1" id2="6" value0="6,76" value1="corlink"/>
<relation id1="7" id2="5" value0="6,76" value1="corlink"/>
<relation id1="4" id2="6" value0="6,77" value1="corlink"/>
<relation id1="6" id2="1" value0="6,81" value1="corlink"/>
...
<relation id1="1" id2="2" value0="1" value1="physlink"/>
<relation id1="2" id2="3" value0="1" value1="physlink"/>
<relation id1="3" id2="4" value0="1" value1="physlink"/>
</data>
```

## 2. Representation modalities in an immersive framework

The characteristics of the data representation language allow us to describe biological data, but it remains to define an immersive visualization paradigm. This paradigm must be adapted at the same time to this language and to the user's needs. We present in this section how we map the data representation language defined in the second section, to a visual representation.

### 2.1 Graph visualization

The selected federator language describes a list of valued objects with their binary valued relationships. We choice to consider this data as a graph, where biological objects are nodes and relationships between

us are edges. Visualizing and exploring data with 3D graph has no reality references (on the contrary to metaphoric representation) and is independent from data.

## 2.2 Visualization description format

Visualizing data by a 3D graph results from the following motivation: invent a system that ensures independence between data and their representation. This motivation transformed into requirements because of the format heterogeneity. Indeed, we did not want to import this heterogeneity in the visualization system. However it remains to choice how graphically represents each value (both object and relationship values) within a 3D graph. In order to map data to their 3D graph, we defined and XML-based data visualization language. To do that, we started to build a short inventory of the graphic 3D object characteristics for 3D graph nodes and edges: Taking into account the description data format, these graphic characteristics can be classified in the three following groups: symbolic, numeric, or both. So within 3D graph visualization, numerical object values may be represented by numerical graphic properties, like node position (x, y, z), node size, node color (r, g, b components) or node transparency (alpha component). Numerical relationship values may be visualized by edge length, edge weight, edge color, or edge transparency. Each symbolic values may be represented by a predefined shape (cube, sphere for nodes, and cylinder, line for edges) or predefined color (red, pink, blue), according to the kind of value (object or relationship value). Finally, string values may be visualized by a 3D text label within the 3D graph representation.

## 2.3 Node placement problem with 3D weight graph visualization

However, we are faced of the following problem: mapping numerical value (correlation) to a distance (edge length) in the 3D space between graph nodes has often no graphic solution, due mainly to 3D Euclidian space constraints.

We use an approach, proposed by Eades [6] simulating two kinds of force between each node. To place two nodes that are in relations respecting distance constraints, he proposed to apply them an attraction force, in order to minimize the global energy  $E$  of this spring system,

$$E = \sum_{\substack{i < j \\ \text{Node}}} k \left( |p_i - p_j| - l_{ij} \right)^2$$

where  $p_i$  is the position of a node  $i$ ,  $l_{ij}$  the optimal distance between the node  $i$  and the node  $j$ , and  $k$  a constant factor.

Moreover, a repulsion force is applied on two close nodes, which are not connected. After several iterations ( $nbiter$ ), this dynamic property allows system to converge into a satisfactory solution where all the distances are as closed more possible than desired edges length. The main disadvantage of above approach is its complexity about:

$$\theta \left( |Node|^2 + |Edge| \right) \times nbiter$$

Indeed, each node reacts to the presence of all its connected neighbors by an attraction force, and moves according to the presence of all the other nodes per repulsion. The complexity strongly decreases by applying a visibility threshold on the not connected nodes. The nodes too much far according to this fixed threshold do not repulse. Moreover, we can remark that, the superposition of node problem on 2D space is not a problem in 3D space, because user cans easy turn around his data. We have just to

approximate distance between connected nodes, but not repulsion between not-connected nodes. We can so decrease the complexity:

$$\theta(|Edge|) \times nbiter$$

On very huge graph, it however was necessary to use segmentation process in order to make some partition of the graph, using MCL clustering [8], Karipys [7] or Fiduccia algorithm [9].

### 3. Results

#### 3.1 Factual data: Yeast gene block duplications

This system was used firstly in order to gene duplications in the yeast chromosomes. In this experiment, each object is one of the yeast chromosome arms. For each object, the values are chromosome name, chromosome size, chromosome side (right or left arm). In each binary relationship between chromosomes, the values are the number of same gene shared by two chromosomes arms.

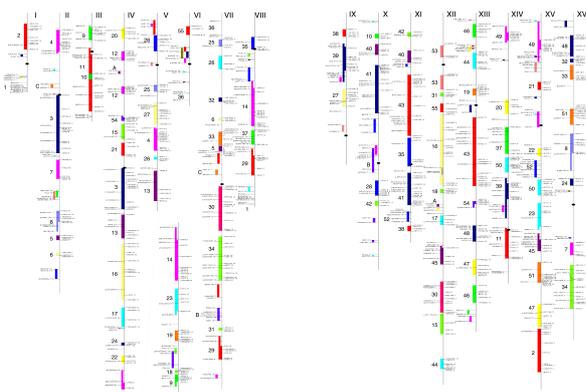


Figure 1: Traditional 2D visualization (16 Yeast chromosomes)

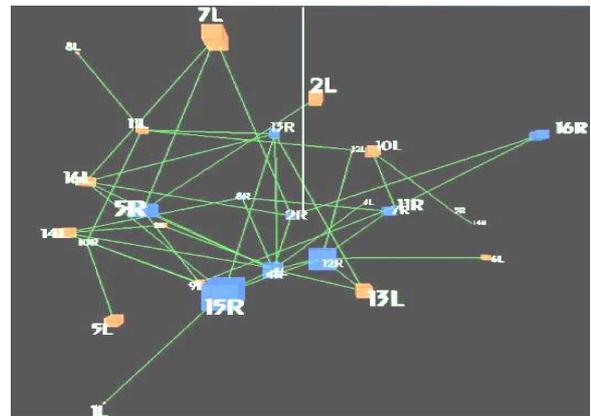


Figure 2: Immersive synthetic visualization (32 Yeast chromosome arms)

In the synthetic representation showed in Figure 2, we can directly see that several chromosome arms, like 4R, are in the centre of the global gene duplications, whereas other chromosome arms take placed in periphery. This representation helps biologist to launch a work on correlation between chromosome placement in the cells and the gene duplication between Yeast chromosomes during evolution.

#### 3.2 Decryphon: A huge protein-to-protein sequences alignment dataset

The Decryphon database [10] contains the results of an exhaustive comparison of all known proteins from living organisms (animals, plants and humans), including the coding sequences from 76 completely sequenced genomes. There are currently two ways for Decryphon analyzing: biologists can use the Decryphon browser to query which proteins are homologous to a targeted protein. However, Decryphon browser does not allow request on a set of proteins. In order to request Decryphon on a set of protein, users must follow the second way: download raw data.



In order to position gene profile nodes, the edge length between is inversely proportional to their correlation score.

Biologists validate this exploration method partially, because many known genes in the same sub network are the same known pathway, for example SSA2 and SSA3 in Figure 9.

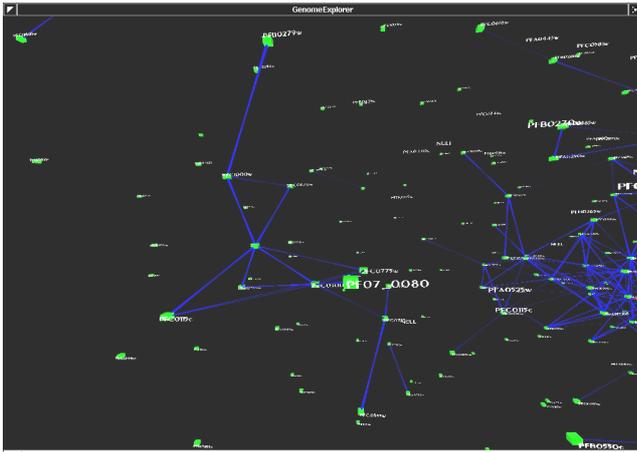


Figure 6 : Plasmodium falciparum correlation expression profile network (1)

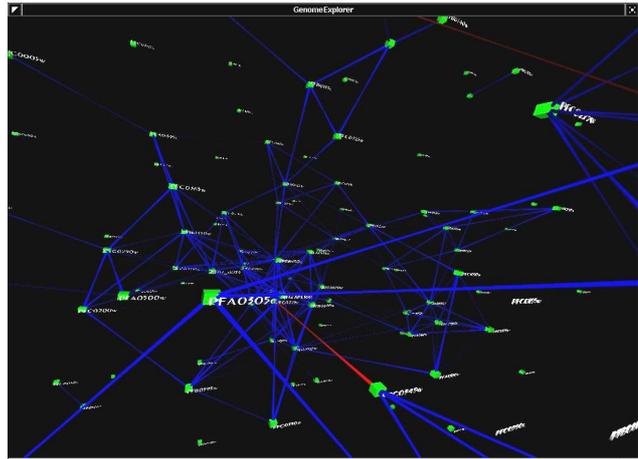


Figure 7 : Plasmodium falciparum correlation expression profile network (2)

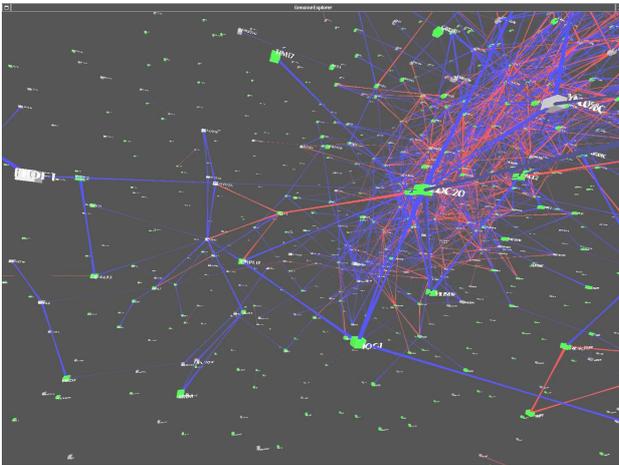


Figure 8: Yeast correlation expression profile network during elutriation phase (1)

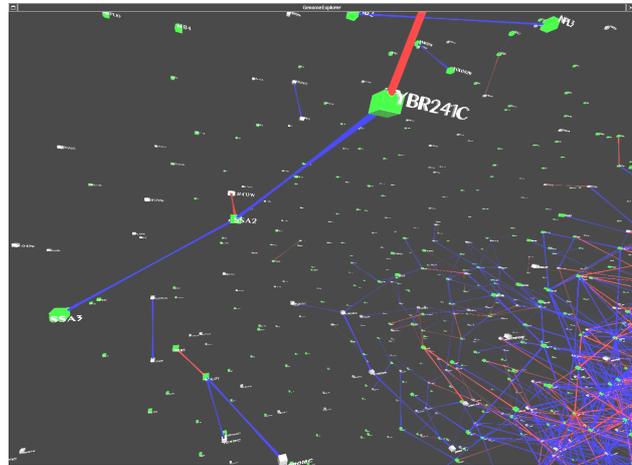


Figure 9: Yeast correlation expression profile network during elutriation phase (2)

#### 4. Conclusion and future work

In this paper, our objective was to elaborate new solutions in order to explore both biological genomic data. This approach is mainly based on the definition of a genomic data representation language, answering the requirements and specificities of biological databases. The representation methods to view these data within an immersive environment, like in *SequenceWorld* [2], were presented and approved successively on various sets of biological data. We try to represent these biological data by 3D graph, using force directed placement algorithm, like *BioBiblioMetrics* [3]. Compared to this work, the data description language offers biologist to precise the semantic of the edge in the representation, and the semantic of all the others graphic characteristics. Moreover, we can represent any biological objects (like protein sequence, biological terms, chromosome arm...), on the contrary to the *Sequence World*, which only deals with genetic sequences. The immersive aspect gives the possibility of exploring huge data in a synthetic way and so constitute the strong points of our system, because it offers a global point of view of the data subjacent structure. These characteristics are particularly interesting when biologists wish to explore a mass of data without precisely knowing what they seek. For example, the partial analysis of *Decryphon* data [10] shows directly several clusters within the representation. This study was concretized by a software development, named *Genome3DExplorer*, which was used to generate the results presented in this paper.

## Reference

1. Hérisson, J. Gros, P.-E. Férey, N. Magneau, N. and Gherbi, R. DNA in Virtuo: Visualization and Exploration of 3D Genomic Structures. 3rd ACM International Conference on Virtual Reality, Computer Graphics, Visualization and Interaction (2004).
2. Rojdestvenski, I. Pettersson, F. Modjeska, D. Sequence World: A Genetics Database in Virtual Reality. Proceedings of the International Conference on Information Visualization (2000).
3. Stapley, B.J. & Benoit, G. Biobibliometrics: Information Retrieval and Visualization from Co-occurrences of Genes Names in Medline Abstracts. International Pacific Symposium on Biocomputing, 5 (2000) pp. 526-537.
4. Pustejovsky, J. Castano, J. & Zhang J. Robust Relational Parsing over Biomedical Literatures: Extracting Inhibit Relations. Proceedings of Pacific Symposium on Biocomputing, (2002).
5. Bozdech, Z. Llinas, M. Pulliam, B.L. Wong, E.D. Zhu, J. DeRisi, J. The Transcriptome of the Intraerythrocytic Developmental Cycle of *Plasmodium Falciparum*. PLoS Biol .1(1) (2003).
6. Eades, P. A Heuristic for Graph Drawing. *Congressus Nutnerantiunt*. 42 (1984) pp. 149–160.
7. Karypis, G. and Kumar, V. Multilevel Algorithms for Multi-Constraint Graph Partitioning. In Proceedings of the IEEE/ACM Conference (SC98), (1998) pp 28.
8. Enright A.J., Van Dongen S., Ouzounis C.A. An Efficient Algorithm for Large-Scale Detection of Protein Families. *Nucleic Acids Research* 30(7) (2002) pp. 1575-1584.
9. Fiduccia, C.M. and Mattheyses, R.M. A Linear-Time Heuristic for Improving Network Partitions. In Proceedings of the 19th ACM/IEEE Design Automation Conference (DAC'82), (1982) pp. 175–181.
10. Decryphon Project (AFM, IBM, GENOMINING) <http://www.infobiogen.fr/services/decryphon/index.html>
11. SwissProt : <http://us.expasy.org/sprot/>
12. GenBank : <http://www.psc.edu/general/software/packages/genbank/genbank.html>
13. Bozdech, Z. Llinas, M. Pulliam, B.L. Wong, E.D. Zhu, J. DeRisi, J. The Transcriptome of the Intraerythrocytic Developmental Cycle of *Plasmodium Falciparum*. PLoS Biol. 1(1) (2003).