



AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Flexible case-based retrieval for comparative genomics

This is the author's manuscript					
Original Citation:					
Availability:					
This version is available http://hdl.handle.net/2318/135681 since 2016-08-05T11:14:17Z					
Published version:					
DOI:10.1007/s10489-012-0399-z					
Terms of use:					
Open Access					
Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.					

(Article begins on next page)



UNIVERSITÀ DEGLI STUDI DI TORINO

This is an author version of the contribution published on: Questa è la versione dell'autore dell'opera: [APPLIED INTELLIGENCE, Vol. 39 (1), 2013, DOI 10.1007/s10489-012-0399-z] ovvero Stefania Montani, Giorgio Leonardi, Stefano Ghignone, Luisa Lanfranco. Vol. 39 (1), 2013, pag. 144–152

> *The definitive version is available at:* La versione definitiva è disponibile alla URL:

[http://link.springer.com/article/10.1007/s10489-012-0399-z]

Flexible case-based retrieval for comparative genomics

Stefania Montani¹, Giorgio Leonardi¹, Stefano Ghignone² and Luisa Lanfranco^{2,3} ¹DISIT - Computer Science Institute, University of Piemonte Orientale, Alessandria, Italy ²IPP – CNR, UAS Torino, Turin, Italy

³Department of Life Science and Systems Biology, University of Turin, Turin, Italy

<u>stefania.montani@unipmn.it</u> (corresponding author); <u>giorgio.leonardi@unipmn.it</u>, <u>stefano.ghignone@unito.it</u>, luisa.lanfranco@unito.it

ABSTRACT

Background: Comparative genomics represents a key instrument to discover or validate phylogenetic relationships, to give insights on genome evolution, and to infer metabolic functions of a given organism. A tool for properly supporting comparative genomics is of paramount importance in several application domains.

Results: We have developed a tool for supporting customized comparative genomics searches. The tool is based on the retrieval step of the Case-Based Reasoning methodology. It takes advantage of an abstraction technique similar to Temporal Abstractions, thus allowing to neglect un-relevant details.

Conclusions: By means of our tool, retrieval is made flexible by the use of abstractions, and efficient by the use of proper taxonomical index structures. Moreover, end-users are allowed to progressively relax or refine their queries, in an interactive way. The tool functionalities are exemplified referring to the study of endobacteria living in arbuscular mycorrhizal fungi.

Keywords: Comparative Genomics; Case-Based Reasoning; Flexible Retrieval; Index-based Retrieval; Interactivity

1 BACKGROUND

Comparative genomics represents a key instrument to answer fundamental questions concerning the biology, ecology and evolutionary history of an organism, or of a biological system and of its composing elements. In particular, the comparative genomics approach is extremely useful when biochemical and physiological data are not available and/or hard to obtain.

The introduction of massive parallel technologies used in next-generation sequencing (NGS) are revolutionizing genomic research [1]. The manipulation and interpretation of the millions of nucleotide sequences generated by NGS present significant computational challenges from the reads assembly to comparisons of many entire genomes. Therefore there is urgent need of platforms able to handle such large-scale genomic data.

In this frame, we are developing a modular architecture for comparative genomics. We have chosen to exploit as much as possible the available tools built around GMOD, the Generic Model Organism Database project [2]. Indeed, GMOD brought to the development of a whole, and still under expansion, collection of open source software tools for creating and managing genome-scale biological databases.

In particular, we are extending the functionalities offered by GMOD, in order to properly meet the needs of a specific comparative genomics study we are interested in, within the BIOBITS project (see section 1.1).

In this paper, we will focus on the characteristics of one single extension we are providing, namely a flexible genome search and retrieval tool, based on the Case-Based Reasoning (CBR) methodology [3]. Such a tool allows to search in the genomes database by expressing queries at different levels of abstraction detail, resorting to a technique similar to Temporal Abstractions (TA) [4; 5]. Moreover, end-users are allowed to progressively relax or refine their queries, in an interactive way. Finally, retrieval is made efficient by the use of multidimensional orthogonal index structures, which allow for early pruning and focusing.

We will present the tool referring to a case study, in which we compare genomes of bacterial species belonging to the Burkholderiaceae family. This family includes a number of well described free-living species and the still enigmatic obligate endosymbionts such those living in arbuscular mycorrhizal fungi (AMF).

The paper is organized as follows. Section 1.1 quickly describes the application we will refer to in order to illustrate our approach. Section 2 first illustrates the general architecture we are implementing for genome analysis in the project, and then focuses on our retrieval tool. Section 3 presents a case study. Section 4 discusses related works and section 5 summarizes our conclusions.

1.1 The application domain

Bacterial endosymbionts in the animal kingdom have been and are excellent models for investigating important biological events, such as organelle evolution, genome reduction, and transfer of genetic information among host lineages [6]. By contrast, the knowledge on endobacteria living in fungi are limited [7; 8]. Arbuscular mycorrhizal fungi (AMFs) are a crucial component of soil microbial communities and exert positive impacts on plants health and productivity in natural and agricultural systems. They establish a symbiotic association with the root of land plants, providing a better mineral nutrition and an increased tolerance to stress conditions. AMFs are thus a significant resource for sustainable agriculture. AMFs are often in further symbiosis with uncultivable bacteria which are hosted in hyphae of AMFs, leading to a complex tripartite system (i.e. endobacterium-AMF-plant roots) [9; 10].

In our current project BIOBITS (see acknoledgments), we focus on the *Candidatus* Glomeribacter gigasporarum [11], endobacterium of the AMF *Gigaspora margarita* (isolate BEG34), currently used as a model system to investigate endobacteria-AMFs interactions.

2 METHODS

2.1 System Architecture

The system architecture we are developing within the BIOBITS project has been engineered exploiting the standard modules and interfaces offered by the GMOD project [2], and completed with custom modules to provide new functionalities (see Figure 1).

The main module of the system contains the database, which provides all the data needed to perform the in-silico activities. We adopted the GMOD Chado database schema [12], to take advantage of its completeness and of its support for controlled vocabularies and ontologies. Furthermore, Chado is the standard database for most of the GMOD modules. Therefore, we can reuse these modules to support the main activities of the project, and extend the system incrementally as the researchers' needs evolve. The database in this module stores and provides all the information about the organisms to be studied (mainly bacteria), their genomes, their known annotations, their proteins and metabolic pathways. It also includes the newly discovered annotations, which can be stored and managed locally until they are confirmed and published.

As explained, our database contains information to be used and stored locally, but we have added the possibility to populate and update it with information retrieved from the biological databases accessible through the Internet. This feature is provided by the set of modules in the *Import modules* section (see Figure 1). The main module (*RRE - Queries*), which is built on the basis of a previously published tool [13], performs queries to different biological databases through the Internet (e.g. the GenBank [14]) and converts the results into a standard format. Afterwards, the *Import module* inserts or updates the retrieved information into the Chado database. This process can be started ondemand, or performed automatically on a regular basis, in order to maintain the local database up-to-date.

Chado also acts as the data interface for the software layers implementing the functionalities and tools used by the researchers. From the architectural point of view, we offer two types of services: the services implemented through existing modules of GMOD (*GMOD modules* section in Figure 1), and new services implemented through new modules, developed ad-hoc (*New applications* section).

The existing GMOD modules we are exploiting are the following:

- *CMap*, which allows users to view comparisons of genetic and physical maps. The package also includes tools for maintaining map data;
- *GBrowse*, which is a genome viewer, and also permits the manipulation and the display of annotations on genomes;
- *GBrowse_syn*, which is a GBrowse-based synteny browser designed to display multiple genomes, with a central reference species compared to two or more additional species;
- SyBil, which is a system for comparative genomics visualizations.

All the GMOD tools exploit a web-based interface to be more user-friendly and easy to use. Moreover, they can be reused as they are, but they can also be customized to meet the researchers' recommendations before being integrated in our software architecture.

As for the new modules, they consist of:

- a set of clustering and other data mining modules (whose description is outside the scope of this paper). Such tools rely on *BioMart* [15], a GMOD facility able to reorganize the information stored in the Chado database into a data warehouse;
- the *flexible retrieval module*, which is the main topic of this paper, and which will be described in section 2.2.



Fig. 1. System architecture.

Every new module added in the *New applications* section of our architecture, or every customized module in the *GMOD modules* section, connects to the other modules of our architecture using GMOD standard interfaces. Therefore, each of them can be published to the GMOD community, in order to extend and enrich the functionalities of this platform.

2.2 A flexible retrieval tool

The flexible retrieval module we have designed in the project architecture implements the retrieval step of the Case Based Reasoning (CBR) [3] cycle.

CBR is a reasoning paradigm that exploits the knowledge collected on previously experienced situations, known as *cases*. The CBR cycle operates by:

- 1. retrieving past cases that are similar to the current one and by
- 2. reusing past successful solutions after, if necessary, properly
- 3. *revising* them; the current case can then be
- 4. *retained* and put into the system knowledge base, called the case base.

Purely retrieval systems, leaving to the user the completion of the reasoning cycle (steps 2 to 4), are however very valuable decision support tools [16], especially when automated adaptation strategies can hardly be identified, as in biology and in medicine [17]. In BIOBITS we are following exactly this research line.

In the flexible retrieval module, the information stored in cases is related to genomes expressed as sequences of nucleotides, each one taken from a different organism, and properly aligned with the genome of a reference organism.

Completing the alignment task is therefore a prerequisite for representing cases in our library.

Details of our alignment strategy, of case representation and of case retrieval are discussed in the next subsections.

2.2.1 Case representation

As previously observed, our first step towards case representation requires alignment. To this end, we rely on BLAST [18], a state-of-the-art local alignment algorithm, specifically designed for bioinformatics applications.

From a typical BLAST output (Figure 2) one can extract basic information (e-value, score and identity percentage) that can be easily plotted as represented in Figure 3, providing a piecewise constant function, which graphically represents the alignment itself. Specifically, Figure 3 refers to e-value.

Query	657724	CGGTGCCGCCGTTCGTGCCCGAAACCGATGAGCAGAAGGTGCTCGCGAAATATGCGGCCG	657783
Sbjct	2790220	CGGTGCCGCCGTTCGTGCCCGTCAACGACGAGCAGAAGCGGCTCGCGCAGTACGCGATGG	2790161
Query	657784	ACGTGCGGGCCGCGCTCGACAAGATCGTCGAGATGAAGCCGGAAGAGGTogcgaagggcg	657843
Sbjct	2790160	ACGTGCGTGCGGCGCTCGACAAGATCGTCGAGCTCAAGCCGGAGGAAGTCGCGAAAGGGG	2790101
Query	657844	cggtctgagcg 657854 Score	
Sbjct	2790100	AAGTGTGAGCG 2790090 E-value	
Score	e = 604 b ities = 3	115 (22), Expect = (1e-169) 81/408 (33), Gaps = 07405 (03) Micromototat that and and (in pucketider)	e
Ouerv	2351954	Alighment start and end (in nucleotides)	2352013
Shict	635050		634991
Oueru	2352014	GATGACTACATCCACTACCACCACCACCACCACCACCACCACCA	2352073
Sbjct	634990	GATCACCATGCTCGACTGCACGCATGCGCCGGGCGGGCGG	634931
Query	2352074	GAACACCTTGTTGAACGACGCGAAATCGCGCGCGTCATCGAGCCACACGCCGCAGCGCAC	2352133
Sbjct	634930	GAAGACCTTGTTGAACGACGCGGAAATCGCGCGCGCGCGC	634871
Query	2352134	GACGIGCICGAGICCGIAGCCGGCITCCTICAGGAICGCGAICACGITCICGAICGICIG	2352193
Sbjct	634870	GACGTGCTCGAGGCCGTAGCCCGCTTCCTTCAGGATCGCGATCACGTTCTCGATGGCCTG	634811
Query	2352194	CTICGACTGCGTGACGATCCCGCCTTCGACGACCTCGCCGTTTACCATCGGCGTCTGGCC	2352253
Sbjct	634810	CTTCGACTGCGTGACGATCCCGCCCTCGACGACCTCGCCGTTCACCATCGGCGTCTGGCC	634751
Query	2352254	CGACACGTACAGCCAGCCGTCGGCCTCGACCGCCCGTGCAAACGGCATCACCTGGCCGCC	2352313
Sbjct	634750	GEACACGTACAGCCACCCGTCGGCCTCGACCGCGCGCGCGAACGGCATCACCTGGCCGCC	
Query	2352314	CGTGCCCTTCGCTTCGCCTACGCCATATCGCTTCATCGTTCACTCCT 2352361	
Sbjct	634690	CGTGCCCTTCGCTTCGCCCACGCCATATCGCTTCATCGTTTCACTCCT 634643	

Fig. 2. BLAST sequence alignment.

Quantitative e-values provided by BLAST can be converted into a set of qualitative levels (e.g. low, high similarity), thus providing a "higher level" view of the information, able to abstract from unnecessary details. To perform such a conversion, we propose a semantic-based abstraction process, similar to the Temporal Abstractions (TA) techniques [4; 5].

TA is an Artificial Intelligence methodology which allows to move from a point-based to an interval-based representation of a time series, where time series points are converted into intervals (episodes), aggregating adjacent points sharing a common behavior, persistent over time. In particular, state abstractions [5] allow to extract episodes associated with qualitative levels of the variable represented by the time series. In state abstractions, the mapping between qualitative abstractions and quantitative values has to be parameterized on the basis of domain knowledge.



Fig. 3. A graphical visualization of sequence alignment (x-axis: nucleotide position of the alignment with respect to the reference string; y-axis: e-values).

In our framework, we adopt this methodology with a main difference: the independent variable is the symbol position in the aligned strings, instead of time. The values to be converted into qualitative levels are then the e-values calculated between the two strings themselves. An example is provided in Figure 3, where e.g. a $10^{-32} < e$ -value $< 10^{-5}$ is considered to be moderately high (Hm).

Indeed, on the basis of domain knowledge, we define a whole taxonomy of qualitative levels, where high values can be further specialized (see Figure 4). Abstractions of e-values at different levels of detail can support different requests of the researchers.

We also allow to abstract the information about the "localization" of the aligned substrings along the nucleotide sequences at different granularities (see Figure 5). This allows to visualize genome information at the level of groups of nucleotides, genes, regions, chromosomes or even complete genomes¹.

This granularity change makes sense from a biological point of view: consider e.g. that a region may be conserved among relative organisms, while a specific gene within the region may not. Thus, a high similarity at the region level might be difficultly identified at the level of single genes (see the example in section 3).

In particular, as shown in Figure 5, we abstract the row data at a minimum granularity of 500-nucleotides-long sequences. This choice is motivated by the fact that the average gene length in our application domain is 1000 nucleotides – but our framework is parametric with respect to this choice.

Indeed, the whole taxonomy of granularities is strongly influenced by domain semantics. Domain knowledge also strongly influences the conversion of a string of symbols from a given granularity to a different one, as required for flexible retrieval (see section 2.2.2).

To summarize, in our framework case representation is obtained as follows. First, a pair of nucleotide strings, optimally aligned as calculated by BLAST, is taken. In particular, for each subsequence of nucleotides, the e-value with respect to the aligned nucleotides in the paired string is provided. Abstractions on such quantitative levels are then calculated, and allow to convert these values into qualitative ones. Abstractions are calculated at the ground level in the symbol taxonomy (and operate also at the ground level in the granularity taxonomy, since they work on nucleotides, see Figure 5). The resulting string of symbols is finally stored in the case library as a case.

Despite the fact that cases are stored as abstractions at the ground level, they could be easily converted at coarser levels in both dimensions (i.e. in the dimension of the taxonomy of symbols, and in the one of granularities). Such a conversion is actually the means by which we support flexible case retrieval, and will be described in the next section.



Fig. 4. Taxonomy of state abstraction symbols. The qualitative symbol low (L) corresponds to e-value> 10^{-5} , while high (H) corresponds to e-value $\leq 10^{-32}$; Further, very high (Hv) corresponds to e-value $< 10^{-32}$; moderately high (Hm) corresponds to $10^{-32} \leq e$ -value $\leq 10^{-5}$.

¹ Since, in our application, the full genome of an organism is typically subdivided into one or more chromosome/plasmids, similarity between two genomes has to be calculated by applying an aggregation operation to the similarities at the single chromosome/plasmid level. We are currently using the arithmetic mean as an aggregation operator.



Fig. 5. Taxonomy of sequence granularities. Observe that, given the lowest level granularity, corresponding to 500-nucleotides-length sequences, it is possible to move up both to 1000-nucleotides-length sequences (and additional levels could be added), and/or, resorting to annotations, to the gene level.

2.2.2 Case retrieval

As described in section 2.2.1, our retrieval framework allows for abstractions according to two dimensions, namely a taxonomy of state abstraction symbols, and a variety of granularities.

Taking advantage from this data representation, we support flexible retrieval.

In particular, we allow users to express their queries at any level of detail, both in the dimension of the taxonomy of symbols and in the dimension of granularity. Obviously, since cases are stored at the ground level in both dimensions, in order to identify the cases that match a specific query, we have to provide a function for scaling up (*up* henceforth) two or more symbols expressed at a specific granularity level to a single symbol expressed at a coarser one. Moreover, we need to define a proper distance function for retrieval.

The data structures described in section 2.2.1, as well as the *up* and the distance functions, have to be detailed on the basis of the semantics of the specific application domain.

In particular, scaling up between the lower levels in the taxonomy of granularities requires the adoption of specific domain-dependent choices. For instance, in our domain it makes sense that an H (high) and an L (low) symbols at the 500-nucleotides-length level are converted into an H symbol at the 1000-nucleotides-length one. Scaling up e.g from the nucleotides level to the genes one is even more application-dependent, and requires to parse and exploit annotations. Annotations in fact allow to know where the gene is located on the nucleotide sequence. In this case the *up* function also has to specify how to scale up the two or more symbols in the 500-nucleotides-length sequences intersecting the gene to the single symbol labeling the gene itself. Similar issues have to be dealt with when scaling up to regions. However, the framework is parametric with respect to these choices, and can be quickly adapted to different domains.

Moreover, we have identified a set of general (domain independent) "consistency" constraints, that any meaningful choice must satisfy, in order to avoid ambiguous or meaningless situations. For instance, distance "preserves" ordering also in case *isa* relationships between symbols are involved (e.g. the distance between L (low) and Hm (moderately high) is smaller than the distance between L (low) and Hv (very high)). The exhaustive presentation of such constraints can be found in [19].

In order to increase efficiency, our framework also takes advantage of multi-dimensional orthogonal index structures, which allow for early pruning and focusing in query answering. Indexes are built on the basis of the data structures described in section 2.2.1. The root node of each index is a string of symbols, defined at the highest level in the symbol taxonomy (i.e. the children of "Any", see Figure 4) and in the granularity taxonomy. A (possibly incomplete) index stems from each root, describing refinements along the granularity and/or the symbol dimension. An example multi-dimensional index, rooted in the H symbol, is represented in Figure 6. Note that, in the figure, granularity has been chosen as the leading dimension, i.e. the root symbol is first specialized in the granularity dimension. From each node of the resulting index, the sequence of symbols of the node itself is then orthogonally specialized in the secondary (i.e. the symbol) dimension, while keeping granularity fixed. However, the opposite choice for instantiating the leading and the secondary dimensions would also be possible.

Technically speaking, to answer a query, in order to enter the more proper index structure, we first progressively generalize the query along the secondary dimension (i.e. the symbol taxonomy), while keeping the leading dimension (i.e. granularity) fixed. Then, we generalize the query in the other dimension as well. Following the generalization steps backwards, we can enter the index from its root, and descend along it, until we reach the node which fits the original query leading dimension level. If an orthogonal index stems from this node, we can descend along it, always following the query generalization steps backwards. We stop when we reach the same detail level in the secondary dimension as in the original query. If the query detail level is not represented in the index, because the index is not complete, we stop at the most detailed possible level. We then return all the cases indexed by the selected node. It is worth noting that indexes may be incomplete with respect to the taxonomies. Index refinement can be automatically triggered by the memorization of new cases in the case base, and by the types of queries which have been issued so far. In particular, if queries have often involved e.g. a symbol taxonomy level which is not yet represented in the index(es), the corresponding level can be created. A proper frequency threshold for counting the queries has to be set to this end. This policy allows to augment the indexes discriminating power only when it is needed, while keeping the memory occupancy of the index structures as limited as possible.

Flexibility and interactivity are also supported by a user-friendly graphical interface, which has been designed by following software engineering principles, in order to enhance usability and user friendliness in the interaction with the system. Through the interface, we provide a graphical representation of the indexes (conceptually depicted as in Figure 6), whose nodes can be exploded or iconified, facilitating index navigation (see Figure 7). Moreover, the graphical interface can support the user in selecting the proper navigation direction, providing him/her with quantitative and qualitative information about the cases indexed by sons and siblings of the currently visited node. For instance, we provide information about the number of indexed cases, the sequence of abstractions representing the cases, and the distance from the sequence of abstractions representing the node currently visited by the user.

Very encouraging experimental results have already been obtained by resorting to the same framework, in the field of haemodialysis [19].



Fig. 6. An example multi-dimensional orthogonal index

🔝 NAVIGAZIONE STRUTTURA DI INDXIZZAZIONE					
We Alter subspace is probability We Alter subspace is probability * Ch * * *	CI SONO I CASI NEL INCOO III SONI ATI TIROVATI Cano a Costor: guarante: 8: Costor: guarante: 9: Costor: 9				

Fig. 7. A snapshot of index navigation as rendered by the system graphical interface.

3 RESULTS

In this section we illustrate the retrieval mechanisms of our approach by means of a specific case study. We take as a reference organism a bacterium belonging to the Burkholderia genus. All bacteria belonging to this family share a region, called DCW cluster, which is involved in the synthesis of peptidoglycan precursors and cell division. The DCW cluster is composed by 14 genes: FtsA, FtsI, FtsL, FtsQ, FtsW, FtsZ, mraW, mraZ, murC, murD, murE, murF, murG. The prominent feature of this cluster is that it is conserved with a high similarity in many bacterial genomes over a broad taxonomic range. Notwithstanding some bacteria belonging to the studied family simply miss one of the 14 genes (specifically the third), while all of the others maintain a high similarity at the DCW region level with their rela-

tives. Therefore, it makes sense to define the *up* function as follows: up(HHLHHHHHHHHHHHH)=H (where the absence of a gene is identified by a low similarity value in the gene position).

Suppose that, more precisely, the user expresses the query HvHvLvHvHvHvHvHvHvHvHvHvHvHvHv, aiming at retrieving the specific bacteria missing the third gene, but very similar to the reference one as regards the other genes.

Our system will first generalize the query in the symbol taxonomy dimension, providing the string HHLHHHHHHHHHHHHH (see Figure 6), and then in the granularity dimension, providing the query H at the region level. This allows to enter the index in Figure 6 from its root. Then, following the generalization step backwards, a node identical to the query can be found, and the ground cases indexed by it can be retrieved. The index search steps are highlighted in the figure.

4 DISCUSSION

In 1993, Aaronson [21] suggested that analogical reasoning (which includes CBR) is particularly applicable to the biological domain, because biological systems are often homologous, and because biologists often design and perform experiments based on the similarity between features of the new system to be investigated, and already known ones. As a matter of fact, since then, some CBR applications in biology and bioinformatics have been published. The paper in [22] is an interesting survey on the topic. The surveyed papers are mostly related to experimental design in protein crystallization and protein structure prediction. However, one contribution [21] also makes the hypothesis of using a CBR approach for predicting unknown regulatory regions. More recently, a hybrid method (resorting to Bayesian techniques and CBR) for feature selection in microarray data analysis has been presented [23]. Except for the work in [21], however, we are not aware of CBR works in genomic comparison. Moreover, the work in [21] does not support any flexible and interactive case retrieval, as we are able to do by means of the abstraction mechanism.

As stated in section 2.2.1, our abstraction mechanisms resembles the one of TA [4; 5]. In fact, the present work has been developed starting from our previous experience on TA-based time series retrieval [19; 20]. With respect to those works, here we have properly adapted the characteristics of the existing framework to the biological domain. Actually, such a framework was designed in a modular and domain independent way. We have realized the adaptation to the biological domain by acquiring the specific domain knowledge, which is the basis for a proper definition of the taxonomies and of the distance and up functions.

As regards TA, they have been extensively resorted to in the literature, especially in the medical field, from diabetes mellitus [24; 25], to artificial ventilation of intensive care units patients [26] (see also the survey in [27]). However, typically they were exploited with the aim to solve a data interpretation task [4], and not to support flexible retrieval.

The goal of our proposal is to try to fill this gap, by exploiting an abstraction mechanism for supporting data interpretation, as well as case exploration and retrieval; this idea thus appears to be significantly innovative in the recent literature panorama.

As previously observed, one of our main goals is also interactivity. It is worth noting that, in classical CBR systems, interactivity is typically not supported: the user is asked to input the entire, precise problem description as a query for case retrieval. This means s/he must know the relevance of every case feature - which is not always straightforward in practice. A research direction meant to overcome this limitation indeed exists in the CBR literature, and is known as Conversational CBR (CCBR) (see e.g. [28; 29; 30]). In CCBR, the user is allowed to input just a brief free text description of the case, to start. The system then supports a progressive query refinement through a conversation, in which best matching cases are listed, and further questions meant to reduce and specialize the retrieval set are asked by the system. Our framework thus loosely resembles CCBR. However, CCBR is characterized by some strong challenges, mainly related to case authoring, and dialog inference. Both aspects are non-trivial, and should be solved by specific modules (based e.g. on machine learning [30] or model-based reasoning [31] techniques). Such an additional effort is not required in our framework.

Interestingly, a number of (non CBR) tools to support bioinformatics applications are available in the literature (see e.g. [32; 33; 34; 35; 36; 37; 38; 39; 40]). Most of them are only loosely related to our work. On the other hand, the approaches in [32; 39; 40] deal with comparative genomics. In particular, the works in [32; 40] afford the problem of multiple sequence alignment, also discussed in [41]. The work in [32] is a methodological contribution that introduces a genetic algorithm to explore the search space for the multiple sequence alignment task. The approach also refines the search through local search optimization. The work in [40] describes VISTA, a tool which allows the visualization of pre-computed pairwise and multiple alignments of whole genome assemblies. With respect to these approaches, our tool provides additional interesting features. Namely, it allows to mine genomes at multiple levels: customized searches can be performed, to retrieve genomes and/or genomic segments matching specific features as described by the query at the desired granularity. Furthermore, queries can be performed efficiently, and potentially on very large databases.

As a last remark, the work in [33] introduces a technique for pattern matching with regular expressions, which has been tested on biological applications. Indeed, as a future work, we would like to extend our querying capabilities, including the definition of a more powerful query language. Such a language could involve the use of regular expressions (see also [19]). Indeed, queries with regular expressions could capture in an abstract and concise way a set of specific "ground" queries. Therefore, works like the one in [33] will be the object of our study in the next future.

5 CONCLUSIONS

In this paper, we have described a modular architecture for supporting in-silico comparative genomics analysis, being developed within the BIOBITS project. In particular, we have focused on the main features of a genome search tool, which implements the retrieval step of the CBR cycle. Such a tool provides researchers with flexible retrieval capabilities, in an interactive fashion. Flexibility and interactivity are also supported by a user-friendly graphical interface. Moreover, retrieval performances are optimized by resorting to multi-dimensional orthogonal index structures, allowing for a quick query answering.

Our future work will be devoted to prove this statement, by means of an extensive experimental work. Future enhancements, such as the definition of a richer query language, are also foreseen.

ACKNOWLEDGEMENTS

Funding: The work is supported by the BIOBITS project, a grant of Regione Piemonte, under the Converging Technologies Call, which involves the University of Turin, the University of Piemonte Orientale, the IPP-CNR and the companies ISAGRO Ricerca s.r.l., GEOL Sas, Etica s.r.l.

REFERENCES

[1] J. Zhang, R. Chiodini, A. Badr, G. Zhang. The impact of next-generation sequencing on genomics, Journal of Genetics and Genomics 38: 95-109, 2011.

- [2] B. Osborne and GMOD Community. GMOD, 2000.
- [3] A. Aamodt and E. Plaza. Case-based reasoning: foundational issues, methodological variations and systems approaches. AI Communications, 7:3959, 1994.
- [4] Y. Shahar. A framework for knowledge-based temporal abstractions. Artificial Intelligence, 90:79133, 1997.
- [5] R. Bellazzi, C. Larizza, and A. Riva. Temporal abstractions for interpreting diabetic patients monitoring data. Intelligent Data Analysis, 2:97122, 1998.
- [6] N.A. Moran, A.J. McCutcheon, and P. Nakabachi. Genomics and evolution of heritable bacterial symbionts. Annu. Rev. Genet., 42:165190, 2008.
- [7] E. Lumini, S. Ghignone, V. Bianciotto, and P. Bonfante. Endobacteria or bacterial endosymbionts? to be or not to be. New Phytol, 170:205208, 2006.
- [8] G. Lackner, N. Moebius, L. Partida-Martinez, C. Hertweck. Complete genome sequence of Burkholderia rhizoxinica, an endosymbiont of Rhizopus microsporus. J Bacteriol 193: 783–784, 2011.
- [9] P. Bonfante P, I. Anca. Plants, mycorrhizal fungi, and bacteria: a network of interactions. Annu Rev Microbiol 63: 363-38, 2009.
- [10] M. Naumann, A. Schusler, P. Bonfante. The obligate endobacteria of arbuscular mycorrhizal fungi are ancient heritable components related to the Mollicutes. ISME J 4: 862– 871, 2010.
- [11] S. Ghignone, A. Salvioli, I. Anca, E. Lumini, G. Ortu, L. Petiti, S. Cruveiller, V. Bianciotto, P. Piffanelli, L. Lanfranco, P. Bonfante. The genome of the obligate endobacterium of an AM fungus reveals an interphylum network of nutritional interactions. The ISME J 6:136-145, 2012.
- [12] C.J. Mungall, D.B. Emmert. The FlyBase Consortium. A chado case study: an ontology-based modular schema for representing genome-associated biological information. Bioinformatics, 23(13):i33746, 2007.
- [13] F. Lazzarato, G. Franceschinis, M. Botta, F. Cordero, and R. Calogero. Rre: a tool for the extraction of non-coding regions surrounding annotated genes from genomic datasets. Bioinformatics, 1;20(16): 2848-50,2004.
- [14] http://www.ncbi.nlm.nih.gov/Genbank/
- [15] http://www.biomart.org/
- [16] I. Watson. Applying Case-Based Reasoning: techniques for enterprise systems. Morgan-Kaufmann, 1997.
- [17] S. Montani. Exploring new roles for case-based reasoning in heterogeneous AI systems for medical decision support. Applied Intelligence, 28:275285, 2008.
- [18] http://blast.ncbi.nlm.nih.gov/Blast.cgi
- [19] S. Montani, G. Leonardi, A. Bottrighi, L. Portinale, P. Terenziani. Supporting flexible, efficient and user-interpretable retrieval of similar time series, IEEE Transactions on Knowledge and Data Engineering 2012 (DOI: <u>http://doi.ieeecomputersociety.org/10.1109/TKDE.2011.264</u>).
- [20] A. Bottrighi, G. Leonardi, S. Montani, L. Portinale, P. Terenziani. Intelligent data interpretation and case base exploration through Temporal Abstractions, Proc. International Conference on Case-Based Reasoning (ICCBR), LNCS 6176, I. Bichindaritz, S. Montani eds., Springer, Berlin, pp. 36-50, 2010.
- [21] J.S. Aaronson, G. Juergen, G.C. Overton. Knowledge discovery in GENBANK, Proc. International Conference on Intelligent Systems for molecular biology, L. Hunter, D. Searls, U. Shavlik eds., AAAI Press, pp. 3-11, 1993.
- [22] I. Jurisica, J. Glasgow. Applications of Case-Based Reasoning in molecular biology, AI Magazine, vol. 25(1), pp. 85-95, 2004.
- [23] I. Bichindaritz, A. Annest. Case based reasoning with Bayesian model averaging: an improved method for survival analysis on microarray data, Proc. International Conference on Case-Based Reasoning (ICCBR), LNCS 6176, I. Bichindaritz, S. Montani eds., Springer, Berlin, pp. 346-359, 2010.

- [24] Y. Shahar, M. Musen. Knowledge-based temporal abstraction in clinical domains, Artificial Intelligence in Medicine 8:267-298, 1996.
- [25] R. Bellazzi, C. Larizza, P. Magni, S. Montani, M. Stefanelli. Intelligent analysis of clinical time series: an application in the diabetes mellitus domain, Artificial Intelligence in Medicine 20: 37–57, 2000.
- [26] S. Miksch, W. Horn, C. Popow, F. Paky. Utilizing temporal data abstractions for data validation and therapy planning for artificially ventilated newborn infants. Artificial Intelligence in Medicine 8: 543–576, 1996.
- [27] P. Terenziani, E. German, and Y. Shahar. The temporal aspects of clinical guidelines, in Computer-based Medical Guidelines and Protocols: A Primer and Current Trends, A. T. Teije, S. Miksch, and P. Lucas, Eds., 2008.
- [28] D. Aha and H. Munoz-Avila. Introduction: interactive case-based reasoning. Applied Intelligence, 14:78, 2001.
- [29] M. Manago, K.-D. Althoff, E. Auriol, R. Traphoner, S. Wess, N. Conruyt, F. Maurer. Induction and reasoning from cases, in Proc. European Workshop on CBR. Springer, pp. 313–318, 1993.
- [30] D. Aha and L. Breslow, Refining conversational libraries, in Proc International Conference on Case Based Reasoning. Springer, pp. 267–278, 1997.
- [31] D. Aha, T. Maney, L. Breslow. Supporting dialogue inferencing in conversational case-based reasoning, in Proc European Workshop on Case Based Reasoning. Springer, pp.262–273, 1998.
- [32] F.J. Mateus da Silva, J. M. Sánchez Pérez, J. A. Gómez Pulido, M. A. Vega Rodríguez, AlineaGA—a genetic algorithm with local search optimization for multiple sequence alignment, Applied Intelligence, 32: 164-172, 2010
- [33] B. J. Ross, Probabilistic Pattern Matching and the Evolution of Stochastic Regular Expressions, Applied Intelligence, 13: 285-300, 2000
- [34] W. Zhang, J. Liu, Y. Niu, Quantitative prediction of MHC-II peptide binding affinity using relevance vector machine, Applied Intelligence, 31: 180-187, 2009
- [35] X. Jin, A. Xu, R. Bie, Cancer classification from serial analysis of gene expression with event models, Applied Intelligence, 29: 35-46, 2008
- [36] H. González-Vélez, M. Mier, M. Julià-Sapé, T. N. Arvanitis, J.M. García-Gómez, HealthAgents: distributed multi-agent brain tumor diagnosis and prognosis, Applied Intelli-
- gence, 30: 191-202, 2009
- [37] K. Webb, T. Whyte, Cell modeling with reusable agent-based formalisms, Applied Intelligence, 24: 169-181, 2006
- [38] S. B. Cho, H. Won, Cancer classification using ensemble of neural networks with multiple significant gene subsets, Applied Intelligence, 26: 243-250, 2007
- [39] C. Wu, H. Chen, S. Chen, Counter-Propagation Neural Networks for Molecular Sequence Classification: Supervised LVQ and Dynamic Node Allocation, Applied Intelli-
- gence, 7: 27-38, 1997
- [40] http://genome.lbl.gov/vista/index.shtml
- [41] T. Lassmann, E.L.L. Sonnhammer, Quality assessment of multiple alignment programs, FEBS Lett 529:126-130, 2002