

Systems Biology

APPAGATO: an APproximate PArallel and stochastic GrAph querying TOol for biological networks

Vincenzo Bonnici¹, Federico Busato¹, Giovanni Micale², Nicola Bombieri¹,
Alfredo Pulvirenti³, Rosalba Giugno^{1,3,*}

¹Department of Computer Science, University of Verona, Strada le Grazie 15 - 37134 Verona and

²Department of Math and Computer Science, University of Catania, Viale A. Doria 6 - 95125 Catania and

³Department of Clinical and Experimental Medicine, University of Catania, via Palermo, 636 - 95122 Catania.

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Biological network querying is a problem requiring a considerable computational effort to be solved. Given a target and a query network, it aims to find occurrences of the query in the target by considering topological and node similarities (i.e. mismatches between nodes, edges, or node labels). Querying tools that deal with similarities are crucial in biological network analysis since they provide meaningful results also in case of noisy data. In addition, since the size of available networks increases steadily, existing algorithms and tools are becoming unsuitable. This is rising new challenges for the design of more efficient and accurate solutions.

Results: This paper presents *APPAGATO*, a stochastic and parallel algorithm to find approximate occurrences of a query network in biological networks. *APPAGATO* handles node, edge, and node label mismatches. Thanks to its randomic and parallel nature, it applies to large networks and, compared to existing tools, it provides higher performance as well as statistically significant more accurate results. Tests have been performed on protein-protein interaction networks annotated with synthetic and real gene ontology terms. Case studies have been done by querying protein complexes among different species and tissues.

Availability and implementation: *APPAGATO* has been developed on top of CUDA-C++ Toolkit 7.0 framework. The software is available at <http://profs.sci.univr.it/~bombieri/APPAGATO>.

Contact: rosalba.giugno@univr.it

Supplementary information: Supplementary data is available at *Bioinformatics* online.

1 Introduction

Technological advances have led to the inference and the validation of structured interaction networks involving genes, proteins, drugs, phenotype, and diseases (Kelley *et al.*, 2003; Panni and Rombo, 2015; Barabasi and Oltvai, 2004). According to the data type, such networks are referred to as: (i) protein-protein interaction (PPI) networks representing either physical or functional interactions among proteins; (ii) gene regulatory networks that express how the activity of genes is regulated; (iii)

metabolic networks describing biochemical reactions between chemical compound of cells; and (iv) signalling networks representing inner/outer cell communications.

A typical example that highlights the advantages and possibilities of analysing interaction relationships is protein function prediction. Although sequence homology is commonly used to functionally annotate proteins, a great amount of them remained uncharacterised (Yu *et al.*, 2013). In this context, different algorithms and tools that compare biological networks have been applied to predict novel protein functions (Jiang *et al.*, 2011; Wang *et al.*, 2013; Malod-Dognin and Pržulj, 2015).

In disease studies, genes showing similar phenotypes tend to be neighbours in protein interaction networks and their aggregation in connected sub-networks is effective to detect biomarkers (Creixell *et al.*, 2015; Fortney *et al.*, 2010; Ideker *et al.*, 2002). Also, finding similar functional and topological sub-networks helps analyzing the conservation among species (Lim *et al.*, 2006). In all these applications, graphs serve as the underlying structures for representing biological networks¹ and graph algorithms solve problems such as network alignment, network querying, motif extractions and network perturbation (Panni and Rombo, 2015; Ma and Gao, 2012; Ciriello *et al.*, 2012; Malod-Dognin and Pržulj, 2015).

In this paper we address the problem of *approximate network querying*, which finds, in a *target* network, *similar* occurrences of a so-called *query* network. The notion of similarity takes into account both the similarities between target nodes and query nodes, and a cost measuring the differences of nodes and their connections. An approximate network querying algorithm has to find the query occurrences, among all possible, with the maximum combined similarity.

Querying tools that deal with similarities are effective in biological network analysis since they provide results also in case of noisy data. They are also suitable in the case of partial knowledge of users when formulating queries. Furthermore, they can be used to compare data from different species where some fundamental and functional structures are partially preserved.

Solving approximate network querying implies applying instances of subgraph isomorphism, which is a NP-complete problem (Dost *et al.*, 2008). In literature, several heuristics have been proposed to solve such a problem in reasonable running time. Examples include, restricting the topology of queries to paths or trees (Dost *et al.*, 2008; Kelley *et al.*, 2004; Shlomi *et al.*, 2006; Pinter *et al.*, 2005), applying network alignment strategies (Gulsoy and Kahveci, 2011; Yuanyuan and Patel, 2008; Tian *et al.*, 2007, 2008), dealing with node similarities and ignoring the query topology (Bruckner *et al.*, 2010; Blin *et al.*, 2010), fixing the topology and computing differences of node labels (Liang *et al.*, 2015). Other methods consist of building indexes to reduce the query time (Khan *et al.*, 2013; Zhang *et al.*, 2009); filtering the set of possible similar target data (Sahraeian and Yoon, 2012; Hong *et al.*, 2015; Pienta *et al.*, 2014); to find only exact occurrences of the query in the network (Bonnici *et al.*, 2013; Cordella *et al.*, 2004; Sun *et al.*, 2012; Bonnici and Giugno, 2016); finding the largest part of the query exactly contained in the target graph and replace the query edges not present in the target with paths (Pienta *et al.*, 2014).

We have created *APPAGATO*, a tool that relies on an iterative sampling method (Lawrence *et al.*, 1993; Micale *et al.*, 2014) to compute functional and topological similarities between a query and a target network. Through a matching probability matrix and a weighted sampling procedure, it selects a seed from which the query-target matching starts. Then, by associating a cost to each approximation, it iteratively extends the match by selecting the approximations with the lowest possible cost. The algorithm runs K times and returns a set of K approximate matches. *APPAGATO* performs approximate network querying by considering the topology of query, taking into account node and edge deletions together with differences on node labels.

To speed-up the querying process in large biological networks, *APPAGATO* has been implemented to run on graphics processing units (GPUs). Due to their low cost, high-performance, and easy integration to any personal computer, GPUs have been increasingly applied to accelerate bioinformatics problems (Dematté and Prandi, 2010; Zhao and Chu, 2014; Vouzis and Sahinidis, 2011). Our aim is to handle large biological networks

in a reasonable time yielding accurate results. We compare *APPAGATO* with *RESQUE* (Sahraeian and Yoon, 2012) and *NeMa* (Khan *et al.*, 2013) since, to the best of our knowledge, they are the most efficient and stable tools in literature very close to *APPAGATO* on both the problem they address and on the approximation concept they assume. We run the tools with different PPI networks as input and compared nodes by using similarities of protein sequences and functional gene ontology annotations. We extensively compare the tools in terms of running time, costs of returned matches, and accuracy in finding protein complexes among different species. The results show that *APPAGATO* outperforms the other two tools yielding more accurate results on large PPI networks.

2 Materials and methods

2.1 Definitions and notations

A graph G is a pair (V, E) , where V is the set of nodes and $E \subseteq (V \times V)$ is the set of edges. If $(u, v) \in E$, we say that v is a neighbour of u . G is *undirected* iff $\forall (u, v) \in E$, then $(v, u) \in E$, i.e. u is a neighbour of v and vice-versa. The degree of a node u , $Deg(u)$, is the number of its neighbours. Given a set of labels A , the function $Lab : V \rightarrow A$ assigns a label to each node of G . We assume that graphs are undirected and labelled only on nodes.

2.1.1 Exact Subgraph Isomorphism

Let $Q = (V, E)$ and $T = (V', E')$ two graphs, named *query* and *target*, respectively. The *exact SubGraph Isomorphism* problem (SubGI) aims to find an injective function, $M : V \rightarrow V'$, which maps each node in Q to a unique node in T , such that $\forall (u, v) \in E$: (i) $(M(u), M(v)) \in E'$; (ii) $Lab(u) = Lab(M(u))$; (iii) $Lab(v) = Lab(M(v))$. A solution of the SubGI problem can be represented as the set $m = \{(v_1, M(v_1)), (v_2, M(v_2)), \dots, (v_{|V|}, M(v_{|V|}))\}$, called a *match* of Q in T . Q may have different maps m_i in T .

2.1.2 Inexact Subgraph Isomorphism and matching costs

In this paper, we deal with the *Inexact SubGraph Isomorphism* problem (ISubGI)², which is a variant of the SubGI problem, and in which we admit node and edge mismatches. A mismatch occurs when (i) two nodes with different labels are mapped through a similarity function, or (ii) a query edge or (iii) a query node is missing in the target graph. The absence of a node implies mismatches for all its edges. A cost c is associated to each mismatch. For the sake of simplicity, the same cost $c = 1$ is associated to each of the three types of mismatch.

We denote with $C = \sum c$ the total cost of mismatches between Q and T . The goal of the ISubGI problem is to find an injective function $M : V \rightarrow V'$, such that C is minimized. In this case, a solution for the ISubGI, $m = \{(v_1, M(v_1)), (v_2, M(v_2)), \dots, (v_k, M(v_k))\}$ with $k \leq |V|$, is called an *approximate match* with a cost $C \geq 0$.

Let $Q_m = (V_m, E_m)$ be the subgraph of query Q that has been mapped in the match m , that is, $V_m = \{v \in V : (v, M(v)) \in m\}$ and $E_m = \{(u, v) \in E : (u, M(u)) \in m, (v, M(v)) \in m, (M(u), M(v)) \in E'\}$. We define $V_{\bar{m}} = V \setminus V_m$ and $E_{\bar{m}} = E \setminus E_m$, the nodes and the edges in Q , respectively, that have not been matched in m . Let $S_{|V| \times |V'|}$ be the label similarities between each node $q \in Q$ and $t \in T$. The *label similarity* values belong to the interval $[0, 1]$. The computation of S is application dependent. In the case of PPI networks, the similarity can be based on sequences, functional, or structural protein similarity.

¹ For the sake of clarity, in this article, we use the terms *graph* and *network* indistinctly.

² Here called also approximate subgraph querying

For example, establishing the conservation of a protein-complex CO of the species A within the species B , consists of searching the subgraph Q_{CO} , extracted from the PPI of A (named G_A), into the PPI of B (named G_B). The two PPIs may have different proteins (i.e., nodes with different names), but with similar function, detectable by looking at sequence similarities. An ISubGI algorithm must search for occurrences of Q_{CO} in G_B that minimize sequences and topology differences. We conclude that CO is conserved in B if we find highly similar occurrences.

The total matching cost C is obtained by summing all node and edge costs and by normalizing them over the number of query elements, as follows:

$$C = \frac{\sum_{q \in V_m} (1 - S(q, M(q)) + |V_m| + |E_m|)}{|V| + |E|} \quad (1)$$

2.2 The APPAGATO algorithm

The method consists of the following three main phases.

2.2.1 Phase 1: Computation of matching probability matrix

Before starting the search, APPAGATO computes a matrix P of matching probabilities between all possible node pairs $\langle q, t \rangle$ ($q \in Q$ and $t \in T$), by combining (i) the label similarity $S(q, t)$, (ii) the degree similarity $D(q, t)$, and (iii) the breadth-first similarity $BFS_{Sim}(q, t)$. The label similarity has been defined in Section 2.1.2. In APPAGATO the label similarity matrix, S , may be provided as input by the user. Alternatively, APPAGATO computes a boolean similarity function to compare node labels. It assigns 1 if labels are identical, 0 otherwise. The degree similarity is a binary function $D(q, t) = 1$ if $Deg(q) \leq Deg(t)$, otherwise it is 0. $BFS_{Sim}(q, t)$ is computed by performing breadth-first visits (BFSs) of the query and target graphs by starting from q and t and evaluating label and degree similarities of the visited nodes, level by level. The maximum depth of the BFS visits is a user-defined parameter l_{max} , with $l_{max} \geq 1$. Given a node x , and a level $l \leq l_{max}$ we denote with $BFS_l(x)$ the set of nodes at level l in the BFS tree rooted at x . An edge $e = (u, v)$ in the BFS tree of q is defined *matchable* iff there exists an edge $e' = (u', v')$ in the BFS tree of t such that $S(u, u')$ and $S(v, v')$ are not 0 and $D(u, u') = D(v, v') = 1$. We denote with $MaxMatch(BFS_l(q), BFS_l(t))$ a maximal set of matchable edges in the BFS tree of q at level l , with respect to the BFS tree of level l rooted in t . The BFS similarity between q and t assumes values in $[0, 1]$ and is defined as follows:

$$BFS_{Sim}(q, t) = \frac{\sum_{l=1}^{l_{max}} l \times |MaxMatch(BFS_l(q), BFS_l(t))|}{\sum_{l=1}^{l_{max}} l \times |BFS_l(q)|} \quad (2)$$

Matching probability matrix. The three similarity values are linearly combined in $MScore(q, t) = S(q, t) + D(q, t) + BFS_{Sim}(q, t)$ and normalized to get the matching probability:

$$P(q, t) = \frac{MScore(q, t)}{\sum_{z \in T} MScore(q, z)} \quad (3)$$

Equation 3 ensures that $\sum_{t \in T} P(q, t) = 1$. In phase 2, the probability matrix is used as a transition matrix within an iterative sampling to extract the best possible matches. The upper side of Figure 1 shows an example of such a matrix computation.

2.2.2 Phase 2: Seed selection

APPAGATO searches the first pair of nodes to be matched by randomly selecting q and t according to the probabilities defined in Equation 3 (see the example of Figure 1).

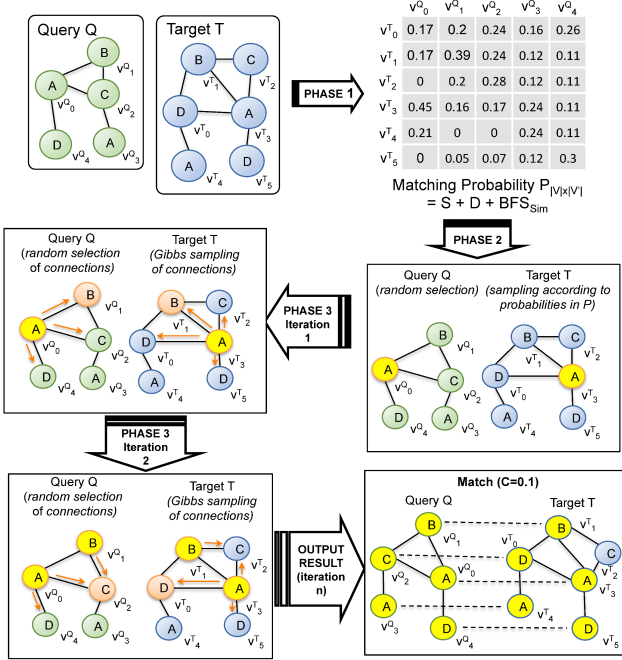


Fig. 1. The APPAGATO approximate matching algorithm.

2.2.3 Phase 3: Extension

Gibbs sampling is used to navigate within a Markov chain, where each state represents a possible query-target node match. The initial state corresponds to the seed selected in phase 2. The sampling method iteratively performs a transition from a state to another, by replacing the query-target nodes pair with a new one, according to a properly defined transition probability. As an example, Figure 1 shows the first two iterations of the extension phase. Transition probabilities are defined by starting from similarity scores, and by taking into account the connections of candidate nodes with already matched nodes. Let Q_m and T_m be the set of query-target matched nodes at a certain step of the extension process. We denote with $Q_m[i]$ ($T_m[i]$) the i -th query (target) node added to the partial match. Let q be a query node neighbour to at least one node in Q_m and t be a target node neighbour to at least one node in T_m . We represent the set of connections between q and the nodes in Q_m through a bit vector $CP(q)$ of $|Q_m|$ elements, called *connection profile* of q , where the i -th element is defined as follows:

$$CP(q)[i] = \begin{cases} 1 & \text{if } (q, Q_m[i]) \in E \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

We define $CP(t)$ in the same way. The *connection profile similarity* between q and t is the corresponding number of equal bits in the connection profiles of q and t :

$$CP_{Sim}(q, t) = \frac{|\{1 \leq i \leq |CP(q)| : CP(q)[i] = CP(t)[i]\}|}{|CP(q)|} \quad (5)$$

The overall similarity scores is $MScoreExt(q, t) = S(q, t) \times CP_{Sim}(q, t)$. The result value is normalized to obtain the final transition probability³:

$$P_T(q, t) = \frac{MScoreExt(q, t)}{\sum_{z \in T} MScoreExt(q, z)} \quad (6)$$

After a number of iterations, n , which is a user-defined parameter, the algorithm returns the reached match between the query and the target

³ Notice that $MScore$ is not used in the extension phase. $MScoreExt$ strongly influences the convergence of the approach (Lawrence *et al.*, 1993; Micale *et al.*, 2014).

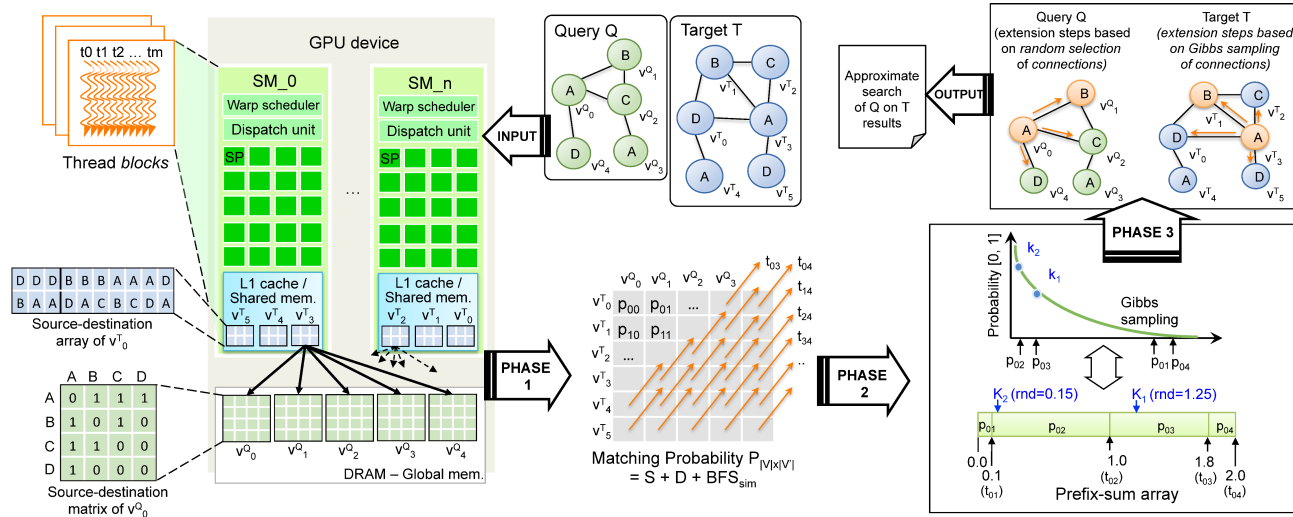


Fig. 2. The parallel search of APPAGATO on the GPU device.

node. The quality of such a match is evaluated by summing the costs of node and edge mismatches between Q and T . APPAGATO does not require any user-defined threshold for the maximum allowed cost of a match. In Figure 1, the approximate match has only a label mismatch, v_2^Q whose label C is mapped with v_0^T having label D, and the cost of the match is $C = 0.1$, computed by applying equation 1. APPAGATO iterates K times phases 2 and 3 and, in each iteration, it starts the sampling procedure from a different seed. Each run of APPAGATO always returns K solutions (approximate matches), each one with the corresponding cost.

2.3 The APPAGATO parallel implementation for GPUs

APPAGATO has been implemented to take advantage of massively parallel GPU architectures. All the processing phases presented in Section 2.2 have been implemented through different CUDA kernels⁴, which are invoked by the host CPU. This allows performing the most compute-intensive tasks of the search algorithm on the GPU device. As for the parallel implementation paradigm for GPUs, each kernel is executed in parallel by several blocks of threads. Thread blocks spread and run concurrently and independently over streaming multiprocessors (SMs). Threads of the same block efficiently cooperate through fast shared memory and by synchronizing their execution through extremely fast (i.e., HW implemented) barriers. Groups of 32 threads of the same block are called warps. Each warp executes one kernel instruction at time in parallel on different data (i.e., single instruction multiple data-SIMD architecture) over the many stream processors (cores) of the GPU device. A warp scheduler efficiently switches between warps with the aim of hiding the latency of thread accesses to the memory.

Given the query and the target graphs, Q and T , the three phases have been implemented as follows (see Figure 2).

2.3.1 Phase 1: Parallel computation of matching probability matrix.

Computing the matching probability matrix is one of the most computation-intensive part of the whole algorithm. It requires $|V| \times |V'|$ computations of Equation 3 and, in particular, $O(|V| + |V'|)$ BFSs over Q and T and the corresponding comparisons between the visited edges (Equation (2)).

APPAGATO implements such a phase through a customized version of *BFS-4K* (Busato and Bombieri, 2015), a parallel implementation of BFS for GPU architectures. *BFS-4K* relies on the concept of *frontier* (Cormen

et al., 2009) (i.e., a FIFO queue that contains the nodes to be visited at each BFS iteration) to implement the graph visit. Through the frontier-based visiting, *BFS-4K* allows equation (3) to be performed over two levels of parallelism: Each parallel warp of a block is mapped to each node of the frontier, and, each parallel thread of a warp is mapped to each outgoing edge from a frontier node.

APPAGATO extends the BFS visit over a third level of parallelism, by running a total number of $|V| + |V'|$ independent BFSs in parallel, one for each node of Q and T . This is done by allocating one block of threads per BFS. The block allocation is automatically done at runtime. A total number of $|V|$ thread blocks perform, in parallel, $|V|$ BFSs (of depth l_{max}) for the query graph. The result consists of *source-destination matrices*, one per node, which are stored in the global memory (the left-most side of Figure 2 shows an example, assuming $l_{max} = 2$). Each matrix contains information on the labels of such edges visited during the BFS from the node along l_{max} levels. In the example of Figure 2, the V_0^Q matrix contains information on the edges of the first level BFS ($A - B$, $A - C$, $A - D$) as well as the edges of the second level BFS ($B - A$, $B - C$, $C - A$, $C - B$, $D - A$).

Similarly, and concurrently, a total number of $|V'|$ thread blocks perform the BFSs for the target graph. The result consists of a set of *source-destination arrays*, one per node, which are stored in the device *shared memory*. This allows an extremely fast memory access for the following comparisons between the generated node structures. The array data structure has been chosen as it allows to represent in a more compact way the source-destination information of T in the limited shared memory. In contrast, the matrix data structure has been chosen as it guarantees a faster access to the source-destination information of Q , to be stored in the larger global memory.

Finally, $|V'|$ thread blocks compare, in parallel, their own source-destination array stored in the local shared memory with all the source-destination matrices in global memory. Such a data structure organization over the GPU memory hierarchy allows the complexity of equation (3) to be reduced from $O(|V| \times |V'|)$ as for the sequential algorithm, to a parallel complexity of $O(1)$. The result of Phase 1 is the matrix $P_{|V| \times |V'|}$, which is stored in the device global memory (see center part of Figure 2).

2.3.2 Phase 2: Parallel seed selection.

APPAGATO emulates the Gibbs sampling to select the K seeds for the successive extension phase. The emulation relies on two parallel primitives, *prefix-sum* (Billeter et al., 2009; Mark Harris, 2008) and

⁴ <http://www.nvidia.co.uk/object/cuda-parallel-computing-uk.html>

weighed random number generation⁵, which are efficiently implemented in the literature for GPUs. Given the similarity value of each query-target node pair p_{xy} of $P_{|V| \times |V'|}$, *APPAGATO* performs the parallel prefix-sum of such values through $|V| \times |V'|$ threads (i.e., one thread per similarity value). The result is a prefix-sum array, in which each element is associated to a thread and the corresponding similarity value. As an example, Figure 2 shows the prefix-sum array of four threads, $t_{01}, t_{02}, t_{03}, t_{04}$ having similarity value 0.1, 0.9, 0.8, and 0.2, respectively. The array elements have been depicted through different sizes to better represent the corresponding similarity values. Then, all the threads generate a random sequence of K values in the interval $[0, \sum p_{xy}]$ (i.e., $[0, 2]$ in the example). The parallel primitive for the random number generation allows the threads to share the generation seed and, as a consequence, to generate the same sequence of random values. This allows the threads to concurrently recognize whether the own boundaries in the prefix-sum array include any randomly generated value. In the example, the sequence of random values $K_1 = 1.25$ and $K_2 = 0.15$ leads to the pair of nodes (v_0^Q, v_3^T) and (v_0^Q, v_2^T) associated to threads t_{03} and t_{02} , respectively, to be selected for the extension phase.

2.3.3 Phase 3: Parallel extension.

The extension phase has been implemented through primitives of BFS, prefix-sum, weighed random number generation over different levels of parallelism. As a first level, the K query-target nodes selected in phase 2 are mapped to thread blocks (i.e., one pair of query-target nodes per block). They are concurrently processed as follows. Given a node pair (e.g., (v_0^Q, v_3^T) in Figure 2) the two nodes are processed in parallel by two thread warps (second level of parallelism). The two warps perform a one-step parallel BFS (third level of parallelism) on Q and T , respectively, to visit the neighbour nodes (i.e., candidate connections) of v_0^Q and v_3^T . The result is two frontiers of neighbours $(\{v_1^Q, v_2^Q, v_4^Q\})$ and $(\{v_0^T, v_1^T, v_2^T, v_5^T\})$ in the example). One step of extension over Q performs through a random selection of a node (connection) from the first frontier (v_1^Q in the example). For such a node, *APPAGATO* generates the connection profile through a one-step parallel BFS. Such a connection profile strongly affects the extension over T , which is performed as follows. Starting from all the nodes of the second frontier, *APPAGATO* (i) runs one step of parallel BFS (one per node), (ii) generates the connection profiles of the visited nodes, and (iii) generates the connection profile similarity of each of such nodes with the connection of Q . Through an emulation of the Gibbs sampling similar to that implemented in phase 2, *APPAGATO* selects the new connection for T . The algorithm iterates over the new pair of nodes (i.e., connection of Q and connection T) for a total number $n = |V|$ iterations.

2.4 Datasets

Physical Interaction Networks We used the PPI networks taken from the STRING v10.0 databases (Szklarczyk *et al.*, 2011) of three species: *Mus musculus*, *Homo sapiens*, and *Danio rerio*. These networks differ significantly in size (number of nodes and edges) and density (i.e., the average number of neighbours per node). For each network, we used up to 250 synthetic labels and gene ontologies annotation downloaded from *BioDbNet*⁶. This yielded 12 different PPIs (i.e., 3 species, each one labelled in 4 different ways). We constructed the queries by randomly extracting sets of 100 connected subgraphs, from each network, by varying the size of the queries up to 128 nodes. In this dataset, the similarities matrix $S_{|V| \times |V'|}(q, t) = 1$ if $Lab(q) = Lab(t)$ otherwise is set to 0.

Functional Interaction Networks The STRING database reports, among two proteins and beside the direct physical interactions used above, indirect functional relations such as structural similarity, similarity between the transcript sequences encoding them, and functional correlations. It gives a score, ranging from 0 (namely no relation is known) to 999, which combines physical and functional (i.e., co-expression data analysis) interactions. We constructed a second dataset by taking into account such a combined score. We extracted 4 PPI networks related to the species *Mus musculus*, *Homo sapiens*, *Danio rerio* and *Saccharomyces cerevisiae*. We fixed the interaction score threshold at 998 to get few but highly functional related interactions within each network. As queries, we used 10 human protein complexes taken from the CORUM database (Ruepp *et al.*, 2010). Since CORUM only reports the set of proteins belonging to a given complex, and not their interactions, we reconstructed the topology of the complex by taking into account the interactions reported in the full STRING database with respect to the *Homo sapiens* species. Finally, we labeled target and query nodes with the protein sequences. We computed the query-target node similarities matrix $S_{|V| \times |V'|}$, by making use of *CUDASW*⁷, which implements a parallel version for GPUs of the Smith-Waterman algorithm for local alignment of sequences. We normalized the matrix by row in order to set to 1 the maximum similarity of the target and query node. We used this dataset to investigate the biological significance of the results. The approximate subgraph matching algorithms were capable to identify functional conservation of protein complexes among different species. We refer the reader to Section 1 and Tables S1-S2 of the Supplementary for more details.

3 Results and discussion

We compared *APPAGATO* with *NeMA* (Khan *et al.*, 2013) and *RESQUE* (Sahraeian and Yoon, 2012) on both the physical and functional datasets described in Section 2.4. All the tools solve ISubGI by taking into account the query topology. Unless differently specified, with the term *APPAGATO* we refer to its implementation on top of CUDA. In the Supplementary, Section 2, we report details on the *APPAGATO* implementation and tuning of parameters (Fig. S1-S3), we assess the robustness of *APPAGATO* over query construction (Fig. S4-S5) and the efficiency of both sequential and parallel versions of *APPAGATO* (Fig. S6-S7).

3.1 Performance

For the physical interaction networks, we report the comparison results only between *APPAGATO* and *NeMA*, since *RESQUE* does not support such a large dataset. Fig. 3 shows the average running times of the two tools on the *Danio rerio* network. In the total running time of *NeMA*, we distinguish the target preprocessing and the querying time. Note that *APPAGATO* does not perform any preprocessing step. The results show that *APPAGATO* is at least three times faster than *NeMA* in case of very small queries (i.e., 4, 8, 16 nodes). The performance difference sensibly increases with larger queries. The plots clearly show that the *APPAGATO* running time is almost constant when increasing the query size and the number of labels. We do not report the comparison results on *Mus musculus* and *Homo sapiens* since, in those networks, the running time difference is even more evident (i.e., *NeMA* requires more than 10,000 seconds for the preprocessing phase and more than 6,000 seconds for the execution phase, while *APPAGATO* always requires around two seconds). Fig. S8 in Section 3 of Supplementary reports the details on the *APPAGATO* running time in all the physical interaction networks, by showing its efficiency varying the number of labels, query size, and network size. Fig. 4 reports the

⁵ <https://developer.nvidia.com/curand>

⁶ <http://biodbnet.abcc.ncifcrf.gov>

⁷ <http://cudasw.sourceforge.net>

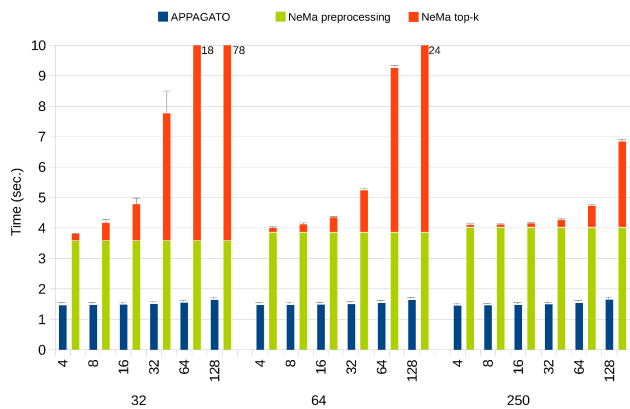


Fig. 3. The running time comparison between APPAGATO and NeMa on the *Danio rerio* PPI network, randomly labelled with 32, 64 and 250 labels. Chart values report the average time on 100 queries. Queries are grouped with respect to the number of nodes, namely 4, 8, 16, 32, 64, 128. For each query, the tools have been run to find 10, 50 and 100 matches.

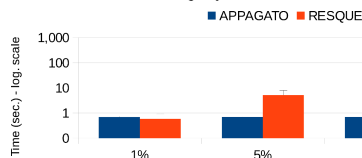


Fig. 4. Running times of APPAGATO and RESQUE on the functional interaction networks. Results are grouped by the similarity thresholds. The running time of RESQUE highly depends on the number of target nodes that can be matched with a query node (i.e., on the similarity threshold t).

comparison of APPAGATO with RESQUE on the functional interaction networks. For the sake of clarity, we do not include the NeMa results in the comparison since in this kind of networks, RESQUE outperforms NeMa. The performance of RESQUE mainly depends on the size of query and target and on the number of possible candidates for each query node. RESQUE requires, as an input, a similarity matrix between query and target nodes. Such a matrix can be partially defined and this affects the quality of the results. If the similarity matrix is fully defined, then the algorithm execution becomes infeasible (i.e., RESQUE takes hours for a single query run). Therefore, we run several tests by changing the percentage of target nodes that can match to a specific query node. Given a threshold t , we set all entries in the similarity matrix with values less than t to 0 (i.e., making them not possible candidates). We then normalized each row by the row maximum value. We chose the percentages 10%, 5% and 1% to obtain reasonable RESQUE running times (i.e., 14, 5, 1 seconds, respectively). APPAGATO always requires around 0.69 seconds. The RESQUE running time rapidly rises as the t threshold increases. In contrast, the APPAGATO running times are always below 1 second.

3.2 Quality measurements of matches

Fig. 5 shows a comparison of the average response costs of APPAGATO and NeMa on the *Danio rerio* physical PPI network. We removed the duplicated matches from the results of APPAGATO to avoid the bias coming from low cost matches. Both algorithms are executed to return the best 10, 50, 100 matches. As expected, both algorithms are highly dependent on the query size. However there is a clear difference in their output quality. The cost of NeMa results are often close to 1, which means they involve a high number of mismatches. In contrast, the averages of the APPAGATO costs range from 0.1 to 0.55. Fig. S9-S10 in Section 3 of Supplementary confirm the accuracy of APPAGATO also on *Homo sapiens* and *Mus musculus*. We measured the statistical significance of the differences between the APPAGATO and NeMa performance. We computed the p-values with a Wilcoxon rank-sum test together with a FDR-correction (false discovery rate) for multiple testing. Fig. S11 in Section

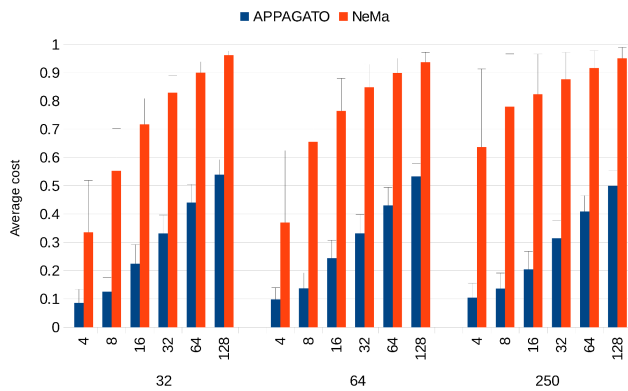


Fig. 5. Average costs (and their standard deviations) by taking into account the set of distinct output matches. Analysis have been performed on the physical interaction PPI of *Danio rerio*. Results are grouped with respect to the number of target labels and query size.

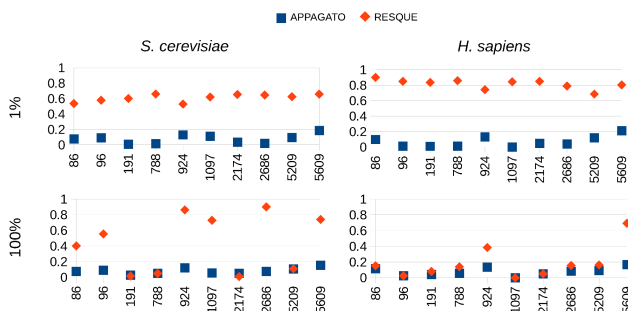


Fig. 6. A chart showing the costs of the 10 protein complexes over the *S. cerevisiae* and *H. sapiens* networks. The CORUM ID of the protein complexes is reported on the x-axis. In the top charts, the similarity threshold is equal to 1%. For those reported in the bottom side the similarity matrix has not been filtered.

3 of Supplementary shows that APPAGATO significantly outperforms NeMa. The number of tested queries having lower p-values increases as the output size becomes larger, particularly when the number of required output matches increases.

3.3 Querying protein complexes among different species.

We compared APPAGATO and RESQUE using 10 human protein complexes taken from CORUM and queried on the functional interaction dataset composed by *Mus musculus*, *Homo sapiens*, *Drosophila melanogaster* and *Saccharomyces cerevisiae* networks (see Fig. 6 and Fig. S12 Supplementary). We test RESQUE using two similarity threshold values, 1% and 100%. RESQUE shows the main performance limitation with a similarity threshold equal to 1% on every target network, while it provides better performance by increasing the cut-off. In all cases, APPAGATO outperforms RESQUE even on the quality of the results. To confirm this, we run the Wilcoxon rank-sum tests (see Fig. S13 in Supplementary). For low similarity thresholds (from 1% to 10%), APPAGATO provides p-values close to 1×10^{-12} . Better p-values (between 1×10^{-5} and 1×10^{-6}) are shown when we defined the whole similarity matrix. Nevertheless, this turned out to be unfeasible from the running time point of view. Fig. S14 in Section 4 of Supplementary shows the functional coherence of results with respect to gene ontology. We computed the average p-value for both algorithms obtained by querying the 10 protein complexes for each of the four species. APPAGATO outperforms RESQUE on every type of target networks and similarity threshold. We refer the reader to Sections 4-5 (Fig. S15-S16-S17) of the Supplementary for details and further application of APPAGATO to compare disease modules over tissue specific protein interaction networks.

4 Conclusions

We have developed *APPAGATO*, a stochastic and parallel algorithm to find approximate occurrences of a query in biological networks. *APPAGATO* deals with node, edge, and node label mismatches. It is implemented for GPUs. The choice of such devices is motivated by their accessible costs, high-performance, and widespread availability on any personal computer. All above features allow *APPAGATO* to compute efficiently functional and topological node similarity together with fast searching of a large number of query matching within the target graph. The results show that *APPAGATO* outperforms the existing tools in terms of running time and result accuracy and, unlike competitors, it scales also on very large PPI networks.

Acknowledgement

We thank S Mohammad E Sahraeian and Byung-Jun Yoon for all their help to use and test their software RESQUE. We thank the authors of NeMA, Arijit Khan, Yinghui Wu, Charu C. Aggarwal and Xifeng Yan, for distributing their software and their prompt support to evaluate it. We thank Dr Anna Privitera for her helpful discussion on *APPAGATO* application.

References

- Barabasi, A.-L. and Oltvai, Z. N. (2004). Network biology: understanding the cell’s functional organization. *Nature Reviews Genetics*, **5**(2), 101–113.
- Billeter, M., Olsson, O., and Assarsson, U. (2009). Efficient stream compaction on wide SIMD many-core architectures. In *Proceedings of the Conference on High Performance Graphics 2009*, pages 159–166.
- Blin, G., Sikora, F., and Viallette, S. (2010). Querying graphs in protein-protein interactions networks using feedback vertex set. *Computational Biology and Bioinformatics, IEEE/ACM Transactions on*, **7**(4), 628–635.
- Bonnici, V. and Giugno, R. (2016). On the variable ordering in subgraph isomorphism algorithms. *Computational Biology and Bioinformatics, IEEE/ACM Transactions on*, (99).
- Bonnici, V., Giugno, R., Pulvirenti, A., Shasha, D., and Ferro, A. (2013). A subgraph isomorphism algorithm and its application to biochemical data. *BMC bioinformatics*, **14**(Suppl 7), S13.
- Bruckner, S., Hübner, F., Karp, R., Shamir, R., and Sharan, R. (2010). Topology-free querying of protein interaction networks. *J Comput Biol*, **17**(3), 237–52.
- Busato, F. and Bombieri, N. (2015). BFS-4K: an efficient implementation of BFS for kepler GPU architectures. *IEEE Transactions on Parallel Distributed Systems*, **26**(7), 1826–1838.
- Ciriello, G., Cerami, E., Sander, C., and Schultz, N. (2012). Mutual exclusivity analysis identifies oncogenic network modules. *Genome Res.*, **22**(2), 398–406.
- Cordella, L. P., Foggia, P., Sansone, C., and Vento, M. (2004). A (sub) graph isomorphism algorithm for matching large graphs. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, **26**(10), 1367–1372.
- Cormen, T., Leiserson, C., Rivest, R., and Stein, C. (2009). *Introduction to Algorithms*. MIT press.
- Creixell, P., Reimand, J., Haider, S., Wu, G., Shibata, T., Vazquez, M., Mustonen, V., Gonzalez-Perez, A., Pearson, J., Sander, C., Raphael, B., Marks, D., Ouellette, B., Valencia, A., Bader, G., Boutros, P., Stuart, J., Linding, R., Lopez-Bigas, N., and Stein, L. (2015). Pathway and network analysis of cancer genomes. *Nat Methods*, **12**(7), 615–621.
- Dematté, L. and Prandi, D. (2010). Gpu computing for systems biology. *Briefings in Bioinformatics*, **11**(3), 323–333. cited By 56.
- Dost, B., Shlomi, T., Gupta, N., Rupp, E., Bafna, V., and Sharan, R. (2008). Qnet: a tool for querying protein interaction networks. *J Comput Biol.*, **15**(7), 913–25.
- Fortney, K., Kotlyar, M., and Jurisica, I. (2010). Method inferring the functions of longevity genes with modular subnetwork biomarkers of caenorhabditis elegans aging.
- Gulsoy, G. and Kahveci, T. (2011). RINQ: Reference-based indexing for network queries. *Bioinformatics*, **27**(13), i149–i158.
- Hong, L., Zou, L., Lian, X., and Yu, P. (2015). Subgraph matching with set similarity in a large graph database. *Knowledge and Data Engineering, IEEE Transactions on*, **27**(9), 2507–2521.
- Ideker, T., Ozier, O., Schwikowski, B., and Siegel, A. F. (2002). Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics*, **18**(suppl 1), S233–S240.
- Jiang, X., Gold, D., and Kolaczyk, E. D. (2011). Network-based Auto-probit Modeling for Protein Function Prediction. *Biometrics*, **67**(3), 958–966.
- Kelley, B., Sharan, R., Karp, R., Sittler, T., Root, D., Stockwell, B., and Ideker, T. (2003). Conserved pathways within bacteria and yeast as revealed by global protein network alignment. *PNAS*, **100**(20), 11394–11399.
- Kelley, B., Yuan, B., Lewitter, F., Sharan, R., Stockwell, B., and Ideker, T. (2004). PathBLAST: a tool for alignment of protein interaction networks. *Nucleic Acids Res.*, **1**(32), W83–8.
- Khan, A., Wu, Y., Aggarwal, C. C., and Yan, X. (2013). NeMa: fast graph search with label similarity. In *Proceedings of the 39th international conference on Very Large Data Bases*, PVLDB’13, pages 181–192. VLDB Endowment.
- Lawrence, C., Altschul, S., Boguski, M., Liu, J., Neuwald, A., and Wootton, J. (1993). Detecting subtle sequence signals: a gibbs sampling strategy for multiple alignment. *Science*, **262**(5131), 208–214.
- Liang, H., Lei, Z., Xiang, L., and Philip S., Y. (2015). Subgraph matching with set similarity in a large graph database. *IEEE Transactions on Knowledge and Data Engineering*, **27**(9), 2507–2521.
- Lim, J. et al. (2006). A protein–protein interaction network for human inherited ataxias and disorders of purkinje cell degeneration. *Cell*, **125**(4), 801–814.
- Ma, X. and Gao, L. (2012). Biological network analysis: insights into structure and functions. *Briefings in Functional Genomics*, **11**(6), 434–442.
- Malod-Dognin, N. and Pržulj, N. (2015). L-GRAAL: Lagrangian graphlet-based network aligner. *Bioinformatics*, **31**.
- Mark Harris, Shubhabrata Sengupta, J. D. O. (2008). *GPU Gems 3: Parallel Prefix Sum (Scan) with CUDA*, chapter 3. Addison Wesley Professional.
- Micale, G., Pulvirenti, A., Giugno, R., and Ferro, A. (2014). GASOLINE: a greedy and stochastic algorithm for optimal local multiple alignment of interaction networks. *PLoS ONE*, **9**(6), e98750.
- Panni, S. and Rombo, S. E. (2015). Searching for repetitions in biological networks: methods, resources and tools. *Briefings in Bioinformatics*, **16**(1), 118–136.
- Pienta, R., Tamersoy, A., Tong, H., and Chau, D. H. (2014). MAGE: matching approximate patterns in richly-attributed graphs. In *2014 IEEE International Conference on Big Data, Big Data 2014, Washington, DC, USA, October 27-30, 2014*, pages 585–590.
- Pinter, R. Y., Rokhlenko, O., Yeger-Lotem, E., and Ziv-Ukelson, M. (2005). Alignment of metabolic pathways. *Bioinformatics*, **21**(16), 3401–3408.
- Ruepp, A., Waegle, B., Lechner, M., Brauner, B., Dunger-Kaltenbach, I., Fobo, G., Frishman, G., Montrone, C., and Mewes, H.-W. (2010). CORUM: the comprehensive resource of mammalian protein complexes. *Nucleic Acids Research*, **38**(suppl 1), D497–D501.
- Sahraeian, S. M. E. and Yoon, B.-J. (2012). RESQUE: Network reduction using semi-markov random walk scores for efficient querying of biological networks. *Bioinformatics*, **28**(16), 2129–2136.
- Shlomi, T., Segal, D., Rupp, E., and Sharan, R. (2006). QPath: a method for querying pathways in a protein-protein interaction network. *BMC Bioinformatics*, **10**(7), 199.
- Sun, Z., Wang, H., Wang, H., Shao, B., and Li, J. (2012). Efficient subgraph matching on billion node graphs. *Proc. VLDB Endow.*, **5**(9), 788–799.
- Szklarczyk, D., Franceschini, A., Kuhn, M., Simonovic, M., Roth, A., Minguéz, P., Doerks, T., Stark, M., Müller, J., Bork, P., Jensen, L., and von Mering, C. (2011). The STRING database in 2011: functional interaction networks of proteins, globally integrated and scored. *Nucleic Acids Res.*, **39**(Database issue), D561–8.
- Tian, Y., McEachin, R., Santos, C., States, D., and JM, P. (2007). SAGA: a subgraph matching tool for biological graphs. *Bioinformatics*, **15**(23), 232–9.
- Tian, Y., Patel, J. M., Nair, V., Martini, S., and Kretzler, M. (2008). Periscope/gq: A graph querying toolkit. *Proc. VLDB Endow.*, **1**(2), 1404–1407.
- Vouzis, P. D. and Sahinidis, N. V. (2011). GPU-BLAST: using graphics processors to accelerate protein sequence alignment. *Bioinformatics*, **27**(2), 182–188.
- Wang, H., Huang, H., and Ding, C. (2013). Function–function correlated multi-label protein function prediction over interaction networks. *Journal of Computational Biology*, **20**(4), 322–343.
- Yu, D., Kim, M., Xiao, G., and Hwang, T. H. (2013). Review of Biological Network Data and Its Applications. *Genomics & informatics*, **11**(4), 200–210.
- Yuan, T. and Patel, J. (2008). Tale: A tool for approximate large graph matching. In *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, pages 963–972.
- Zhang, S., Li, S., and Yang, J. (2009). Gaddi: Distance index based subgraph matching in biological networks. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology, EDBT ’09*, pages 192–203, New York, NY, USA. ACM.
- Zhao, K. and Chu, X. (2014). G-BLASTN: accelerating nucleotide alignment by graphics processors. *Bioinformatics*, **30**(10), 1384–1391.