



Contents lists available at ScienceDirect

# Biochemical and Biophysical Research Communications

journal homepage: [www.elsevier.com/locate/ybbrc](http://www.elsevier.com/locate/ybbrc)



## Review

# Integration, visualization and analysis of human interactome



Chiara Pastrello<sup>a,1</sup>, Elisa Pasini<sup>b,1</sup>, Max Kotlyar<sup>a</sup>, David Otasek<sup>a</sup>, Serene Wong<sup>a,c</sup>, Waheed Sangrar<sup>a</sup>, Sara Rahmati<sup>a,d</sup>, Igor Jurisica<sup>a,d,e,\*</sup>

<sup>a</sup> Princess Margaret Cancer Centre, University Health Network and TECHNA Institute for the Advancement of Technology for Health, TMDT, Room 11–314, 101 College Street, Toronto, ON M5G 1L7, Canada

<sup>b</sup> CRO Aviano National Cancer Institute, Cancer Bio-Immunotherapy Unit, Via Franco Gallini 2, 33081 Aviano, PN, Italy

<sup>c</sup> York University, Department of Computer Science and Engineering, Toronto, ON, Canada

<sup>d</sup> University of Toronto, Department of Medical Biophysics, Toronto, ON, Canada

<sup>e</sup> University of Toronto, Department of Computer Science, Toronto, ON, Canada

## ARTICLE INFO

### Article history:

Available online 1 February 2014

### Keywords:

Omics  
Interactome  
Network analysis  
Gastric cancer  
Visual data mining

## ABSTRACT

Data integration and visualization are crucial to obtain meaningful hypotheses from the diversity of 'omics' fields and the large volume of heterogeneous and distributed data sets. In this review we focus on network analysis as a key technique to integrate, visualize and extrapolate relevant information from diverse data. We first describe challenges in integrating different types of data and then focus on systematically exploring network properties to gain insight into network function. We also describe the relationship between network structures and function of elements that form it. Next, we highlight the role of the interactome in connecting data derived from different experiments, and we stress the importance of network analysis to recognize interaction context-specific features. Finally, we present an example integration to demonstrate the value of the network approach in cancer research, and highlight the importance of dynamic data in the specific context of signaling pathways.

© 2014 Elsevier Inc. All rights reserved.

## Contents

1. Omics .....	758
2. Integration .....	758
3. Accurate representation of omics data .....	758
3.1. Data collection and storage .....	758
3.2. Data exchange .....	759
3.3. Example difficulties in data exchange .....	759
4. Network visualization .....	759
5. Example integration .....	760
5.1. New era in pharmacology .....	760
5.2. Gastric cancer .....	760
5.3. Illustrative network for gastric cancer .....	761
5.4. Investigating the relationship between tetracycline and gastric cancer through network analysis .....	761
6. Network structure–function relationship .....	762
6.1. Network properties .....	762
6.1.1. Global network properties .....	763
6.1.2. Local network properties .....	763
6.2. Computational challenges .....	764
7. Importance of the interactome in network analysis .....	764
7.1. Expansion of the known human interactome .....	764

\* Corresponding author at: Princess Margaret Cancer Centre, University Health Network, TMDT, Room 11–314, 101 College Street, Toronto, ON M5G 1L7, Canada.

E-mail address: [juris@ai.utoronto.ca](mailto:juris@ai.utoronto.ca) (I. Jurisica).

<sup>1</sup> These authors contributed equally to this work.

7.2.	Interactome expansion .....	764
7.3.	Benefits of interactome expansion .....	764
7.4.	Interactome biases .....	764
7.5.	Interactome accuracy .....	765
7.6.	Comparisons of human interactomes across time .....	765
8.	Power of network analysis to identify context-specific information .....	767
8.1.	Context-specificity of interactomes .....	767
8.2.	A context application: modeling of signaling pathways .....	768
9.	Discussion .....	769
	Acknowledgments .....	769
	References .....	769

## 1. Omics

The suffix “omics” is appended to words describing a field of study and usually involves large scale, comprehensive and systematic techniques. The first use of omics in this manner was genomics [1], and the first international project to collect the most complete data set, a building block for genomics was the Human Genome Project, launched in 1990 and completed in 2003 (<http://www.genome.gov/10001772>). The list of completed genomes includes 170 Eukaryota organisms, almost 3500 Virus and more than 2500 Bacteria (as listed in <http://www.ebi.ac.uk/genomes/>; last accessed 13th Dec. 2013). Comparably, there are 277 completed proteomes for Eukaryota organisms, more than 1600 Bacteria and more than 1100 Virus (as listed in <http://www.uniprot.org/taxonomy/complete-proteomes>; last accessed 13th Dec. 2013). Moreover, the structure of more than 94,000 proteins across different organisms has been described so far (as listed in <http://www.rcsb.org/pdb/home/home.do>; last accessed 13th Dec. 2013; the organism with the highest number of structures being *Homo Sapiens*).

Increasing data collections enable more comprehensive analyses, but coping with this data deluge is not trivial. For example, one mass spectrometry experiment can result in thousands to hundreds of thousands of spectra for one sample [2]. Likewise, next-generation sequencing produces millions of reads per sample [3]. The amount of data being generated calls for uniformity, standardization and optimized workflows [4]. Even if these are very basic concepts, they are not as widespread as one would expect. For example, there are 338 protein–protein interaction (PPI) databases, 243 metabolic pathway databases and 202 signaling pathway databases (as listed in <http://www.pathguide.org/>; last accessed 13th Dec. 2013) of which only some are in a format that supports data interchange across databases. Navigating through these vast resources can be challenging but integrating such data is both beneficial and increasingly necessary.

## 2. Integration

The emergence of high-throughput (HT) assays shifted research from hypothesis-driven exploration to data-driven hypothesis generation. However, generating substantially more data, HT methods in turn led to shifting from predominantly using statistical tools to depending on computational biology approaches, especially data mining and machine learning algorithms, to aid data analysis and interpretation [5,6]. As the number of omics disciplines grows, and with them the amount of data, the combination of an increasing number of different perspectives can give the scientist a more complete (and more realistic) view of the system they are studying. Now, the challenge is data integration, and in turn integrative data analysis [7]. For example, integrating gene expression with copy number variation data, mutation status, methylation profile and microRNA targeting can highlight the key players in a specific disease. Complex, multifactorial diseases can only be fully

investigated with this type of approach. Relationships between these data and entities can effectively be represented as graphs. Thus, network visualization and analysis is becoming one of the key tools for integrative analysis.

## 3. Accurate representation of omics data

Data integration requires immense attention to information representation, annotation and support for accurate data exchange. Integrative computational biology supports modeling biological processes using data integrated across many omics fields. To address this, numerous data architectures have been established to effectively and efficiently collect, store, annotate and exchange data. These architectures vary in scope, intent, and standards they use. They are continuously being updated to represent the most current knowledge, and as such will contain inconsistencies and incompleteness. Often, researchers rely on one or many such architectures to integrate pre-existing research with their own, or to share their own results. Understanding the nature of the architectures available as well as being able to accurately specify which have been used is critical to reducing ambiguities in this process, and improves the quality and utility of published results.

### 3.1. Data collection and storage

Information in omics data changes frequently, arriving in the form of peer-reviewed studies. These can be small-scale, hypothesis-based studies with a small number of results, or wider scale HT studies with thousands of results. Collecting and storing these results has necessitated the use of many diverse omics databases. These databases vary greatly in scope; for our purposes, they address entities, relationships between them, and annotations. Entity databases cover proteins, genes, small molecules, or other biologically relevant objects, and include for example UniProt [8], GeneCards [9], RefSeq [10], and DrugBank [11,12]. Relationship databases describe how these entities relate to each other. These include, but are not limited to, protein or gene interactions, drug targeting, and biological pathways. Some examples include IntAct [13], I2D [14–16], BioGRID [17], and KEGG [18]. IMEx consortium [19] and PSICQUIC registry [20] are notable collections of interaction databases that share representation and curation workflow standards [19]. Annotation databases attempt to create indices of terms and definitions. These terms are intended to unambiguously describe entities or relationships between them, often structured as an ontology, which further describes how terms relate to each other in a standard way. Some examples include Gene Ontology [21] and the numerous controlled vocabularies included in the PSI-MI standard [22]. Some databases may rely on each other's standards: for example, UniProt and other sequence databases may contain references to relevant Gene Ontology terms in their records. Resources such as GeneCards [23] integrate heterogeneous

Download English Version:

<https://daneshyari.com/en/article/10756906>

Download Persian Version:

<https://daneshyari.com/article/10756906>

[Daneshyari.com](https://daneshyari.com)