

# Big Data for Disease Prevention and Precision Medicine

Marco Moscatelli<sup>1</sup>, Matteo Gnocchi<sup>1</sup>, Andrea Manconi<sup>1</sup>, and Luciano Milanese<sup>1</sup>

<sup>1</sup> Institute for Biomedical Technologies – National Research Council (CNR-ITB)

Corresponding Author:

Marco Moscatelli<sup>1</sup>

Via Fratelli Cervi, 93 - 20090 Segrate (MI) -Italy

Email address: marco.moscatelli@itb.cnr.i

## 1 Motivation

2 Nowadays, advances in technology has arisen in a huge amount of data in both biomedical  
3 research and healthcare systems. This growing amount of data gives rise to the need for new  
4 research methods and analysis techniques. Analysis of these data offers new opportunities to  
5 define novel diagnostic processes. Therefore, a greater integration between healthcare and  
6 biomedical data is essential to devise novel predictive models in the field of biomedical  
7 diagnosis. In this context, the digitalization of clinical exams and medical records is becoming  
8 essential to collect heterogeneous information. Analysis of these data by means of big data  
9 technologies will allow a more in depth understanding of the mechanisms leading to diseases,  
10 and contextually it will facilitate the development of novel diagnostics and personalized  
11 therapeutics. The recent application of big data technologies in the medical fields will offer new  
12 opportunities to integrate enormous amount of medical and clinical information from population  
13 studies. Therefore, it is essential to devise new strategies aimed at storing and accessing the  
14 data in a standardized way. Moreover, it is important to provide suitable methods to manage  
15 these heterogeneous data.

16

## 17 Methods

18 In this work, we present a new information technology infrastructure devised to efficiently  
19 manage huge amounts of heterogeneous data for disease prevention and precision medicine. A  
20 test set based on data produced by a clinical and diagnostic laboratory has been built to set up  
21 the infrastructure.

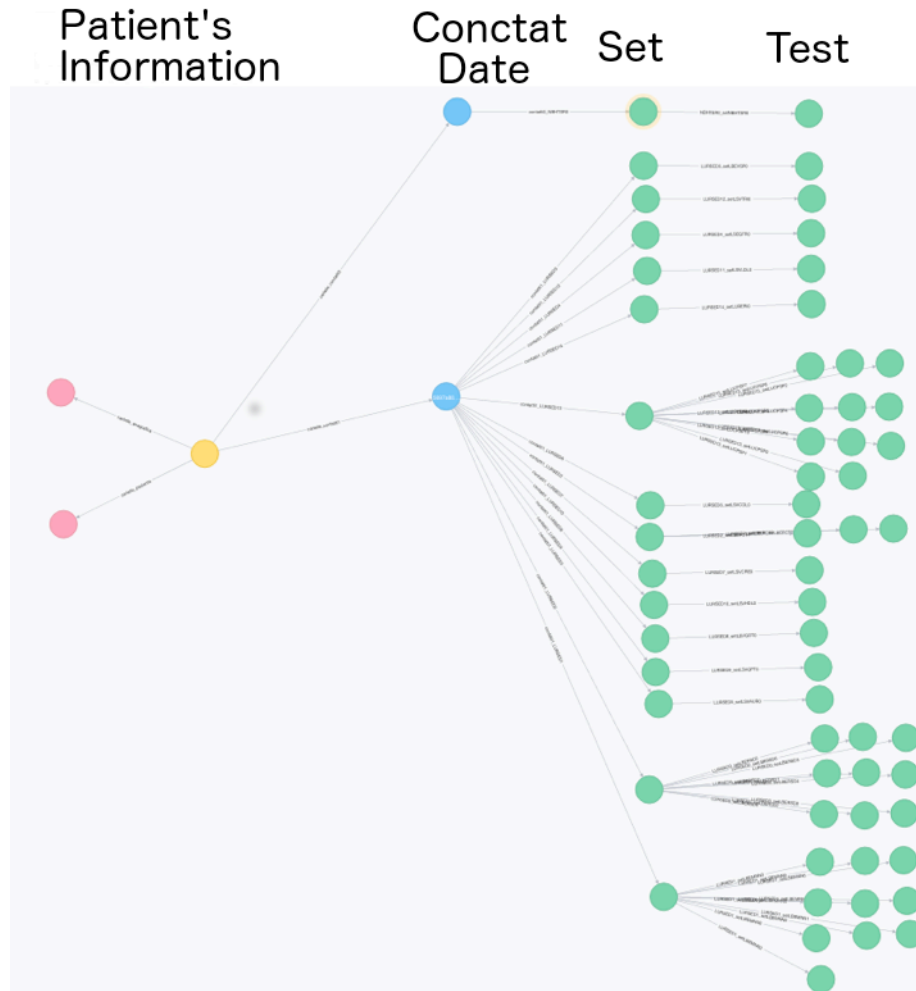
22 When working with clinical data is essential to ensure the confidentiality of sensitive patient  
23 data. Therefore, the set up phase has been carried out using "anonymous data". To this end,  
24 specific techniques have been adopted with the aim to ensure a high level of privacy in the  
25 correlation of the medical records with important secondary information (e.g., date of birth, place  
26 of residence).

27 It should be noted that the rigidity of relational databases does not lend to the nature of these  
28 data. In our opinion, better results can be obtained using non-relational (NoSQL) databases.  
29 Starting from these considerations, the infrastructure has been developed on a NoSQL  
30 database with the aim to combine scalability and flexibility performances. In particular,  
31 MongoDB [1] has been used as it fits better to manage different types of data on large scale. In  
32 doing so, the infrastructure is able to provide an optimized management of huge amounts of  
33 heterogeneous data, while ensuring high speed of analysis.

34

## 35 Results

36 The presented infrastructure exploits big data technologies in order to overcome the limitations  
37 of relational databases when working with large and heterogeneous data. The infrastructure  
38 implements a set of interface procedures aimed at preparing the metadata for importing data in  
39 a NOSQL DB. Moreover, data can also be represented as a graph using Neo4j [2]; The Neo4J  
40 DB allows you to emphasize and enhance the connections between the data and facilitate the  
41 retrieve and navigation of data (Fig 1).



42

43

44 Experimental tests on huge amount of data show that our infrastructure exhibits performances  
 45 in terms of speed and scalability unachievable with relational databases. These performances  
 46 are mainly related to ability of the infrastructure to index any type of field as well as to customize  
 47 the queries. In particular, the high flexibility to customize the queries increases the search  
 48 performance and specificity of the results.

49 As for future work, we planned to implement new functions and operators to perform specialized  
 50 statistics analysis on big data.

51

## 52 References

53 [1] <http://www.mongodb.org>

54 [2] <http://neo4j.com/>

55

## 56 Acknowledgement

57 This work has been supported by the “Fondazione Bracco”, the Italian Ministry of Education and  
58 Research Flagship (PB05) “InterOmics” and the European “MIMOMICS” (305280) projects.