# GOPHER

an HPC framework for large scale graph exploration and inference (MLHPCS)



Barcelona M2 Supercomputing Center Center Centro Nacional de Supercomputación

Marc Josep, Xavier Teruel, Victor Giménez, et al.

# Outline



#### Introduction

- Related work and context
- Analysis, design and implementation
- **Experimental results**
- Conclusions
- Future work

# Introduction

# **Ontologies** are widely used in Biomedicine

- Data integration & interoperability
- Interpretation of high-throughput experiments and clinical information
- E.g., HPO describes phenome abnormalities (symptoms)

**HPC and AI solutions** are needed to traverse and model the interconnectivity of **multiple ontologies**.



### Related work and context



Open Biomedical and Biological Ontologies (OBO) library

- best practices
- curated corpora of ontologies

Phenotypic and genotypic relationships studies  $\rightarrow$  Identify molecular drivers underlying human diseases

Genome-wide association studies (GWAS)

- biological complexity
- lack of consensus
- susceptibility

Limiting factors

**HPC:** Message Passing Interface (*MPI*), Open Multi-Processing (*OpenMP*), and OpenMP SuperScalar (*OmpSs*)

# Analysis and design (methods)



#### Probability, assuming a gaussian distribution,...

$$P(c|\mu,\sigma)=rac{1}{\sigma\sqrt{2\pi}}e^{-(x-\mu)^2/2\sigma^2}$$

... and computing the odds



$$Odds(connected|Counts) = \prod_{pt \in pathtypes} \frac{P(Count(pt)|\mu_{con}(pt), \sigma_{con}(pt))}{P(Count(pt)|\mu_{discon}(pt), \sigma_{discon}(pt))}$$

# Implementation: the graph





#### Implementation: data structures



D	Index	
Children		
Child-1	Child-2	 Child-N
Parents		
Parent-1	Parent-2	 Parent-N
Neighbo	urs-1	
Neigh-1	Neigh-2	 Neigh-N
Neighbo	urs-2	
Neigh-1	Neigh-2	 Neigh-N



# Implementation: MPI Implementation (baseline)





### Implementation: the load imbalance problem

![](_page_8_Figure_1.jpeg)

**GOPHER**, an HPC framework for large scale graph exploration and inference (MLHPCS'20)

June, 25th 2020

9

BSC

# Implementation: MPI Implementation (balancing)

![](_page_9_Picture_1.jpeg)

![](_page_9_Figure_2.jpeg)

connections, distributed among MPI ranks

# Implementation: OpenMP/OmpSs implementation

![](_page_10_Picture_1.jpeg)

![](_page_10_Figure_2.jpeg)

#### Experimental results: model use case

![](_page_11_Picture_1.jpeg)

#### Use case

- **Ontologies:** GO and HPO, only human, version: January 2019.
- Path size: all types of paths with size of 4 or 5 elements.
- Pairs nature: from phenotypes to genotypes.
- Samples: 85,750 randomly sampled pairs of both types.
- Direct edge removal: yes

#### Methods:

- Receiver Operating Characteristic (ROC)
- Precision-Recall (PR) curves

Algorithm: Each path type possible of up to length 5

#### Experimental results: model validation

![](_page_12_Picture_1.jpeg)

![](_page_12_Figure_2.jpeg)

## Experimental results: performance use case

![](_page_13_Picture_1.jpeg)

#### Environment - MareNostrum IV cluster, located at BSC, each node:

- 2 Intel Xeon Platinum 8160, running at 2.1 GHz
- 48 cores (i.e., 24 per processor) and 33 MB L3 Cache
- 2 NUMA sockets (i.e., 1 socket per processor), 192GB per socket

#### Use case

- **Ontologies:** GO and HPO, only human, version: January 2019.
- Path size: all types of paths up to a size of 5 elements.
- **Pairs nature:** from phenotypes to genotypes.
- **Samples:**100,000 randomly sampled pairs (constant seed).
- Direct edge removal: yes

Algorithm: Number of paths for each path type

#### Experimental results: scaling factors

![](_page_14_Figure_1.jpeg)

![](_page_14_Figure_2.jpeg)

### Conclusions

![](_page_15_Picture_1.jpeg)

#### Introduce the GOPHER framework for large graph exploration and inference

*"estimate the likelihood that two ontology terms are associated when missing a direct connection through a co-annotated gene"* 

#### An interdisciplinary work:

- A Biological topic;
- A Machine Learning approach;
- By means of High Performance Computing technology

#### **Preliminary results**

- Model analysis: obtaining an AUC score of 0.96 over 1.
- Performance: load imbalance problem  $\rightarrow$  balancing schedule  $\rightarrow$  scalability plots

### Future Work

![](_page_16_Picture_1.jpeg)

#### The HPC approach

- Study GOPHER behaviour in other architectures
- Further performance analysis (explore other metrics, trace analysis)
- Optimisation opportunities
  - Improve intra-node balance: use Dynamic Load Balancing library
  - Improve inter-node balance: use OmpSs@Cluster
- The ML approach
  - Further study (and validation) of the proposed thesis
- The Biological approach
  - Other actionable use cases: anticancer treatment recommendations
  - Other biological ontologies: mouse and fruit fly

# Thanks!

# Further information at: <a href="https://www.linkedin.com/in/xteruel">https://www.linkedin.com/in/xteruel</a>

![](_page_17_Picture_2.jpeg)

Barcelona Supercomputing Center Centro Nacional de Supercomputación

## Related work and context

Open Biomedical and Biological Ontologies (OBO) Foundry

- best practices
- curated corpora of ontologies

![](_page_18_Picture_4.jpeg)

![](_page_18_Picture_5.jpeg)

www.obofoundry.org

Phenotypic and genotypic relationships studies  $\rightarrow$  Identify molecular drivers underlying human diseases

**HPC:** Message Passing Interface (*MPI*), Open Multi-Processing (*OpenMP*), and OpenMP SuperScalar (*OmpSs*)

![](_page_18_Picture_11.jpeg)