

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



**Barcelona  
Supercomputing  
Center**

Centro Nacional de Supercomputación

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

Introduction

Related work and context

Analysis, design and implementation

Experimental results

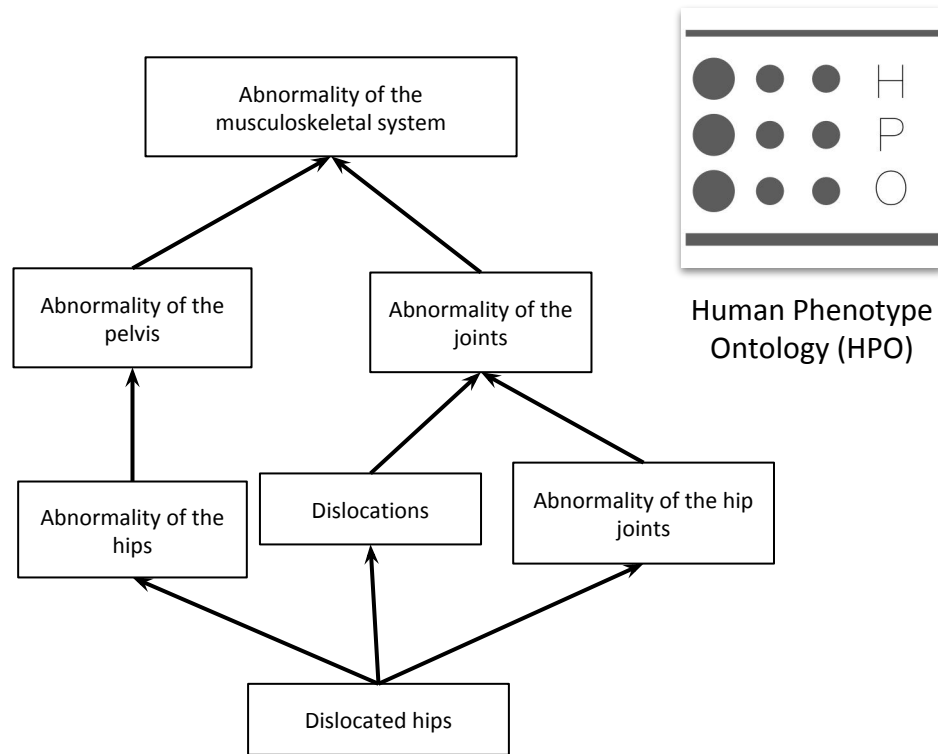
Conclusions

Future work

## 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**.



## Open Biomedical and Biological Ontologies (OBO) library

- best practices
- curated corpora of ontologies

Phenotypic and genotypic relationships studies → Identify molecular drivers underlying human diseases

## Genome-wide association studies (GWAS)

- biological complexity
- lack of consensus
- susceptibility

**Limiting factors**

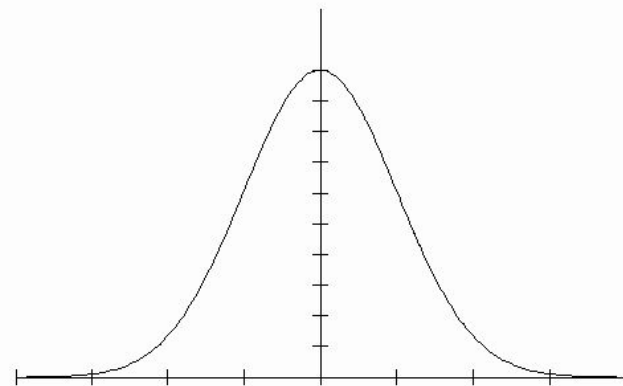
**HPC:** Message Passing Interface (*MPI*), Open Multi-Processing (*OpenMP*), and OpenMP SuperScalar (*OmpSs*)

**Probability**, assuming a gaussian distribution,...

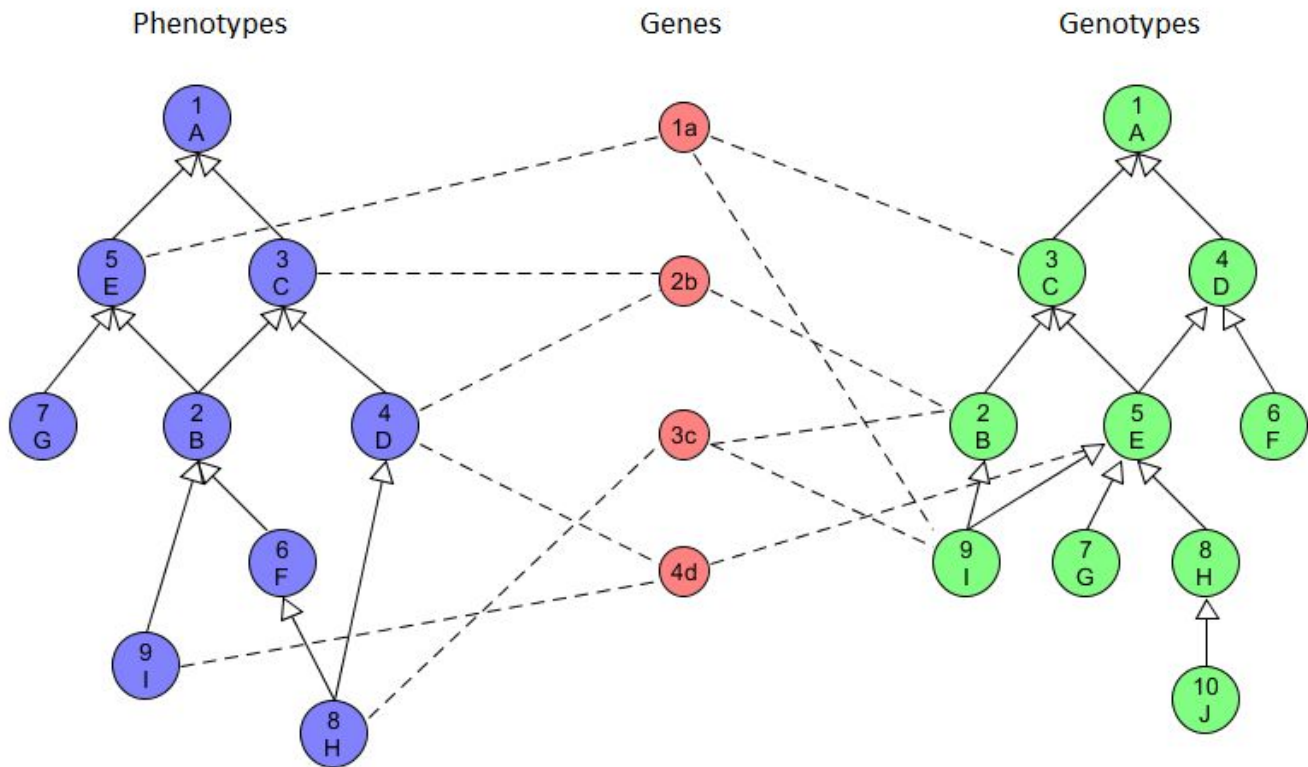
$$P(c|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

... and computing the **odds**

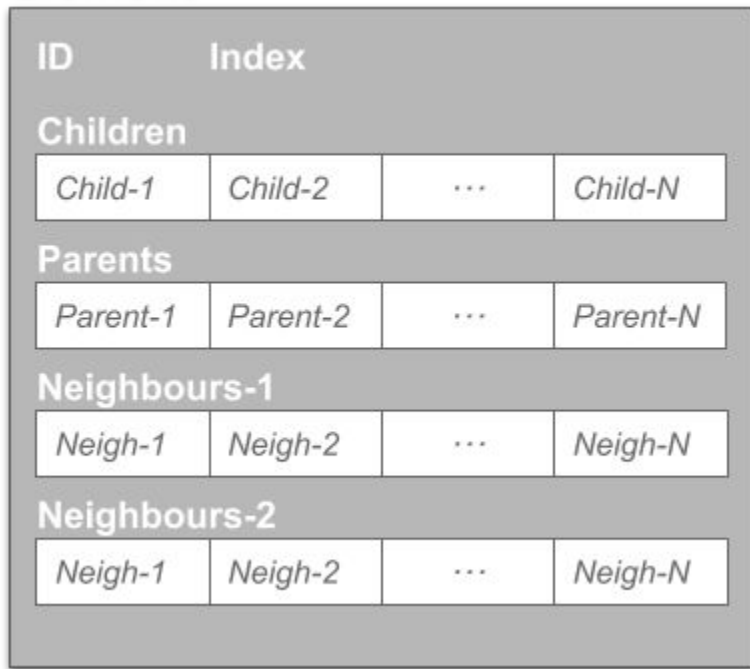
$$\text{Odds}(\text{connected}|\text{Counts}) = \prod_{pt \in \text{path types}} \frac{P(\text{Count}(pt)|\mu_{\text{con}}(pt), \sigma_{\text{con}}(pt))}{P(\text{Count}(pt)|\mu_{\text{discon}}(pt), \sigma_{\text{discon}}(pt))}$$



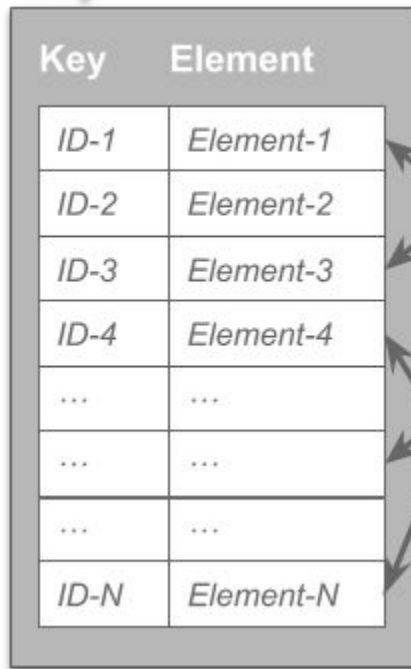
# Implementation: the graph



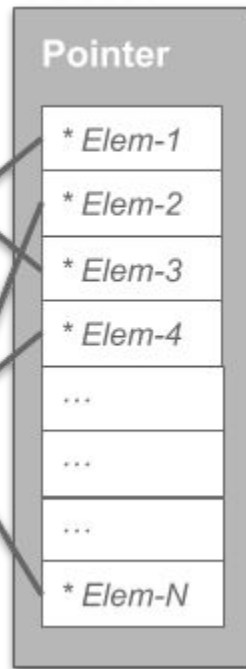
## Element



## Map



## Vector



## Map

Key	Element
ID-1	Element-1
ID-2	Element-2
ID-3	Element-3
ID-4	Element-4
...	...
...	...
...	...
ID-N	Element-N

## Vector

Pointer
* Elem-1
* Elem-2
* Elem-3
* Elem-4
...
...
...
* Elem-N

## Algorithm 2 MPI parallelization with communication

```
1:  $chSize \leftarrow \text{Ontology 1 size} / \text{num Ranks}$ 
2:  $(start, end) \leftarrow (my Rank * chSize, my Rank * chSize + chSize)$ 
3:  $(nPaths, nPairs) \leftarrow (0, 0)$ 
4: for  $i \leftarrow start, end$  do
5:   for  $j \leftarrow 0, \text{Ontology 2 size}$  do
6:      $nPaths \leftarrow nPaths + \text{Search Paths}(\text{Ontology 1}(i), \text{Ontology 2}(j), \text{path type})$ 
7:      $nPairs \leftarrow nPairs + 1$ 
8: MPI_AllReduce ( $nPaths, nPairs$ )
9:  $(average, st dev) \leftarrow (nPaths/nPairs, \text{Calculate Local St Dev})$ 
10: MPI_Reduce ( $st dev$ )
11: if  $my Rank = 0$  then
12:    $st dev \leftarrow \text{Calculate Overall St Dev}$ 
```

Even number of elements distributed among MPI ranks



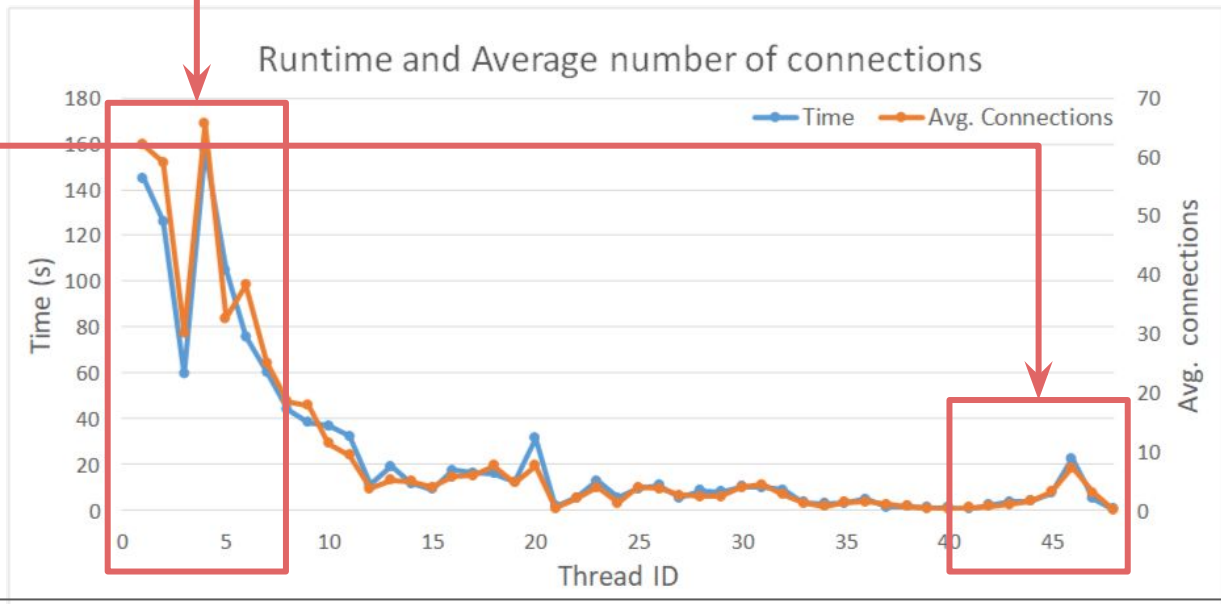
# Implementation: the load imbalance problem

Map

Key	Element
ID-1	Element-1
ID-2	Element-2
ID-3	Element-3
ID-4	Element-4
...	...
...	...
...	...
ID-N	Element-N

Vector

Pointer
* Elem-1
* Elem-2
* Elem-3
* Elem-4
...
...
...
* Elem-N



## Map

Key	Element
ID-1	Element-1
ID-2	Element-2
ID-3	Element-3
ID-4	Element-4
...	...
...	...
...	...
ID-N	Element-N

## Vector

Pointer
* Elem-1
* Elem-2
* Elem-3
* Elem-4
...
...
...
* Elem-N

## Algorithm 3 MPI parallelization with new work distribution

```
1:  $total\ work \leftarrow 0$ 
2: for  $i \leftarrow 0, Ontology\ 1\ size$  do
3:    $total\ work \leftarrow total\ work + connections(i)$ 
4:  $chSize \leftarrow total\ work / num\ Ranks$ 
5:  $(start, end) \leftarrow (begin(my\ Rank, chSize), end(my\ Rank, chSize))$ 
6:  $(nPaths, nPairs) \leftarrow (0, 0)$ 
7: for  $i \leftarrow start, end$  do
8:   for  $j \leftarrow 0, Ontology\ 2\ size$  do
9:      $nPaths \leftarrow nPaths + Search\ Paths(Ontology\ 1(i), Ontology\ 2(j), path\ type)$ 
10:     $nPairs \leftarrow nPairs + 1$ 
11:  $MPI\_AllReduce(nPaths, nPairs)$ 
12:  $(average, st\ dev) \leftarrow (nPaths/nPairs, Calculate\ Local\ St\ Dev)$ 
13:  $MPI\_Reduce(st\ dev)$ 
14: if  $my\ Rank = 0$  then
15:    $st\ dev \leftarrow Calculate\ Overall\ St\ Dev$ 
```

Different number of elements, according with the number of connections, distributed among MPI ranks

## Map

Key	Element
ID-1	Element-1
ID-2	Element-2
ID-3	Element-3
ID-4	Element-4
...	...
...	...
...	...
ID-N	Element-N

## Vector Ontology 1

Pointer
*Elem-1
*Elem-2
*Elem-3
*Elem-4
...
...
...
*Elem-N

## Map

Key	Element
ID-1	Element-1
ID-2	Element-2
ID-3	Element-3
ID-4	Element-4
...	...
...	...
...	...
ID-N	Element-N

## Vector Ontology 2

Pointer
*Elem-1
*Elem-2
*Elem-3
*Elem-4
...
...
...
*Elem-N

## Algorithm 4 OmpSs parallelization with reduction

```
1:  $(nPaths, nPairs) \leftarrow (0, 0)$ 
2: for  $i \leftarrow start, end$  do
3:   for  $j \leftarrow 0, \text{Ontology 2 size}$  do
4:     # pragma omp task reduction(+, nPaths)
5:      $nPaths \leftarrow nPaths + \text{Search Paths}(\text{Ontology 1}(i), \text{Ontology 2}(j), \text{path type})$ 
6:      $nPairs \leftarrow nPairs + 1$ 
7:   # pragma omp taskwait / barrier
8:   MPI All Reduce (num paths, num pairs)
9:    $(average, st\ dev) \leftarrow (nPaths/nPairs, 0)$ 
10:  # pragma omp parallel for reduction(+, st dev)
11:  for  $i \leftarrow 0, start - end$  do
12:     $st\ dev \leftarrow st\ dev + (nPaths(i) - average)^2$ 
13:  MPI Reduce (st dev)
```

## Use case

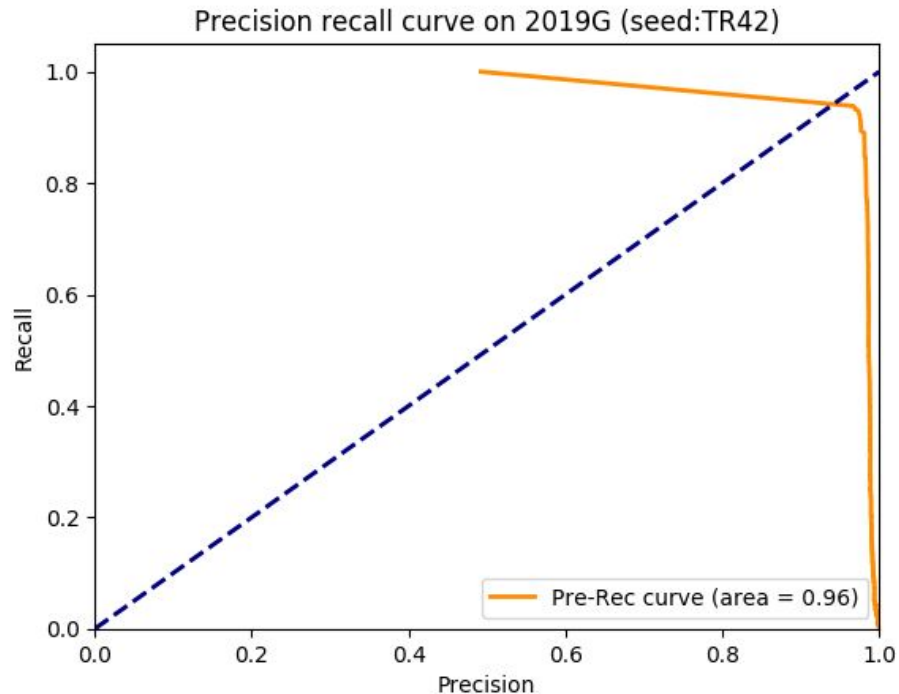
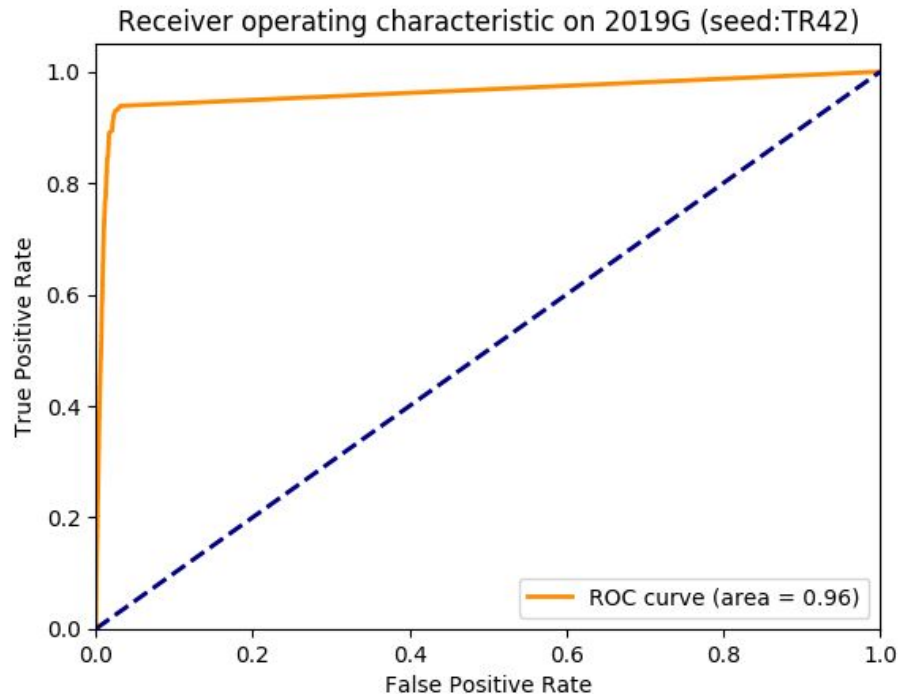
- **Ontologies:** GO and HPO, only human, version: January 2019.
- **Path size:** all types of paths with size of 4 or 5 elements.
- **Pairs nature:** from phenotypes to genotypes.
- **Samples:** 85,750 randomly sampled pairs of both types.
- **Direct edge removal:** yes

## Methods:

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

**Algorithm:** *Each path type possible of up to length 5*

# Experimental results: model validation



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

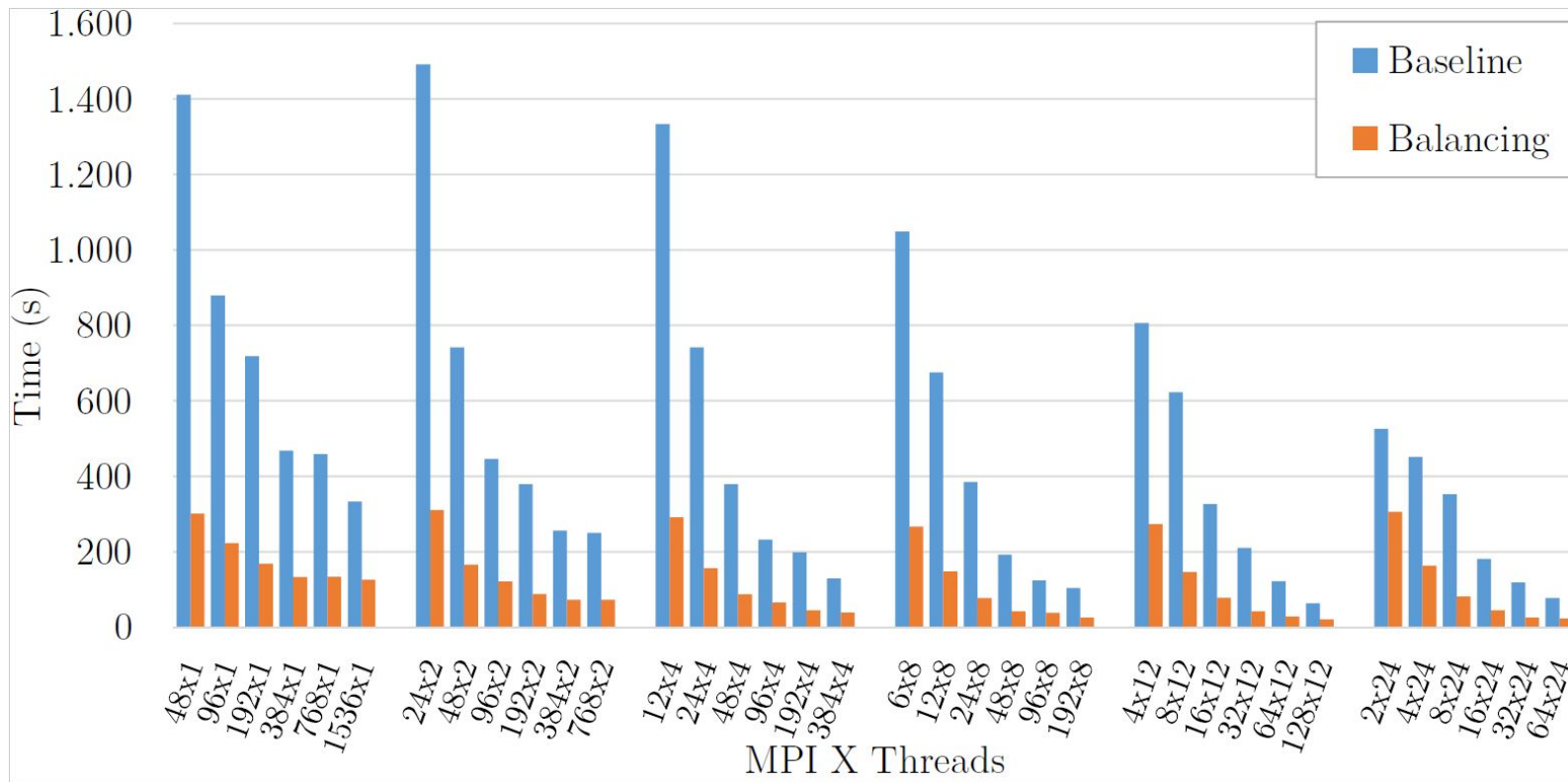
- 2 Intel Xeon Platinum 8160, running at 2.1 GHz
- 48 cores (i.e., 24 per processor) and 33 MB L3 Cache
- 2 NUMA sockets (i.e., 1 socket per processor), 192GB per socket

## 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



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 → balancing schedule → scalability plots



## 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:

<https://www.linkedin.com/in/xteruel>



**Barcelona  
Supercomputing  
Center**

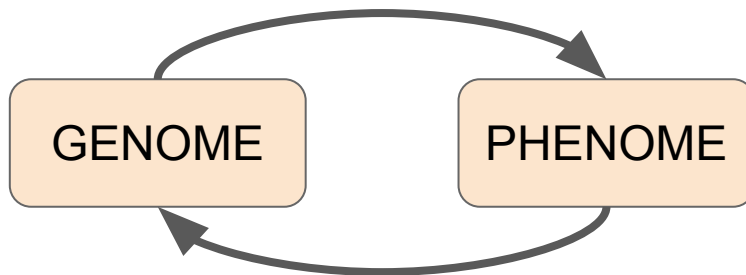
Centro Nacional de Supercomputación

## Open Biomedical and Biological Ontologies (OBO) Foundry

- best practices
- curated corpora of ontologies



[www.obofoundry.org](http://www.obofoundry.org)



Phenotypic and genotypic relationships studies → Identify molecular drivers underlying human diseases

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