

CCS/ITALY

ITALIAN REGIONAL  
CONFERENCE ON  
COMPLEX SYSTEMS



# Classification of Genome Wide Association data by Belief Propagation Neural network

*Daniele Dall'Olio<sup>1</sup>, Nico Curti<sup>1,2</sup>, Armando Bazzani<sup>1,2</sup>,  
Daniel Remondini<sup>1,2</sup>, Gastone Castellani<sup>1,2</sup>*

<sup>1</sup>Department of Physics and Astronomy, University of Bologna

<sup>2</sup>INFN Bologna

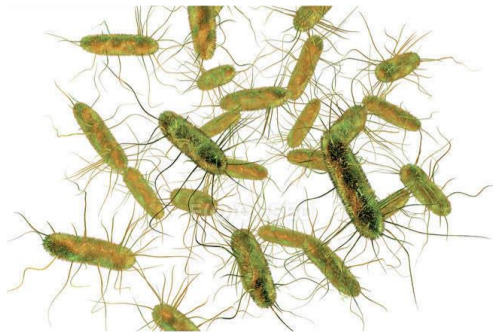
**Presented by:**  
Daniele Dall'Olio

Trento, Italy  
July 1-3, 2019

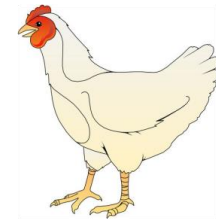
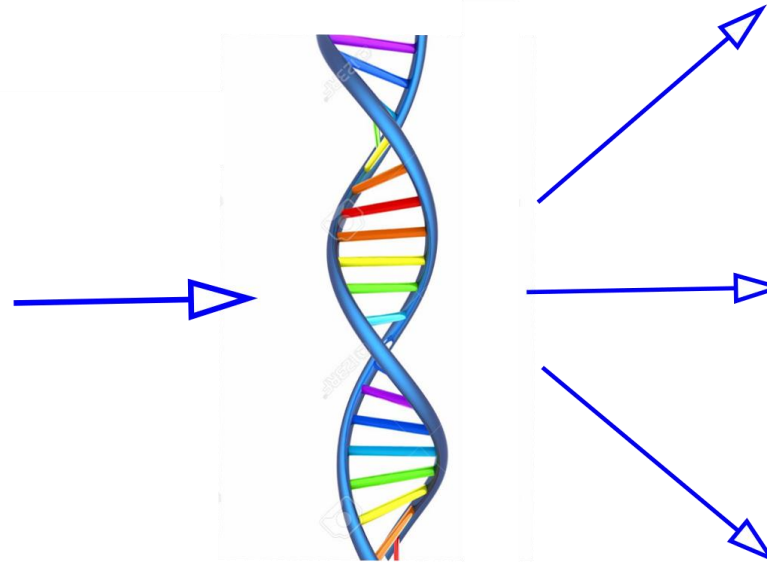
**Objective:** *Identification, containment and mitigation of emerging infectious diseases and foodborne outbreaks.*

Source Attribution

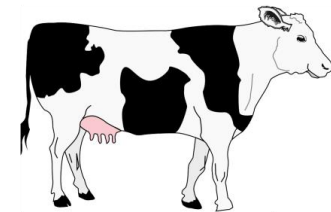
Bacteria



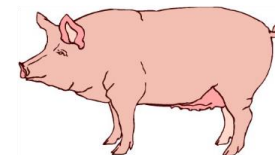
Genome Wide Association



?



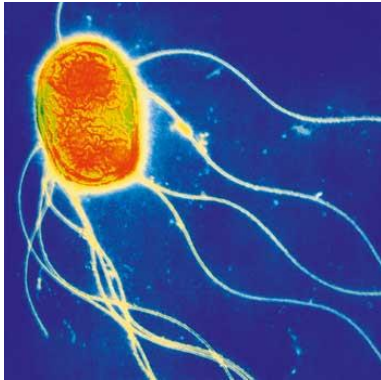
?



?

## SNPs

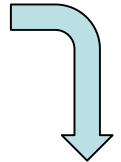
### *Salmonella Enterica*



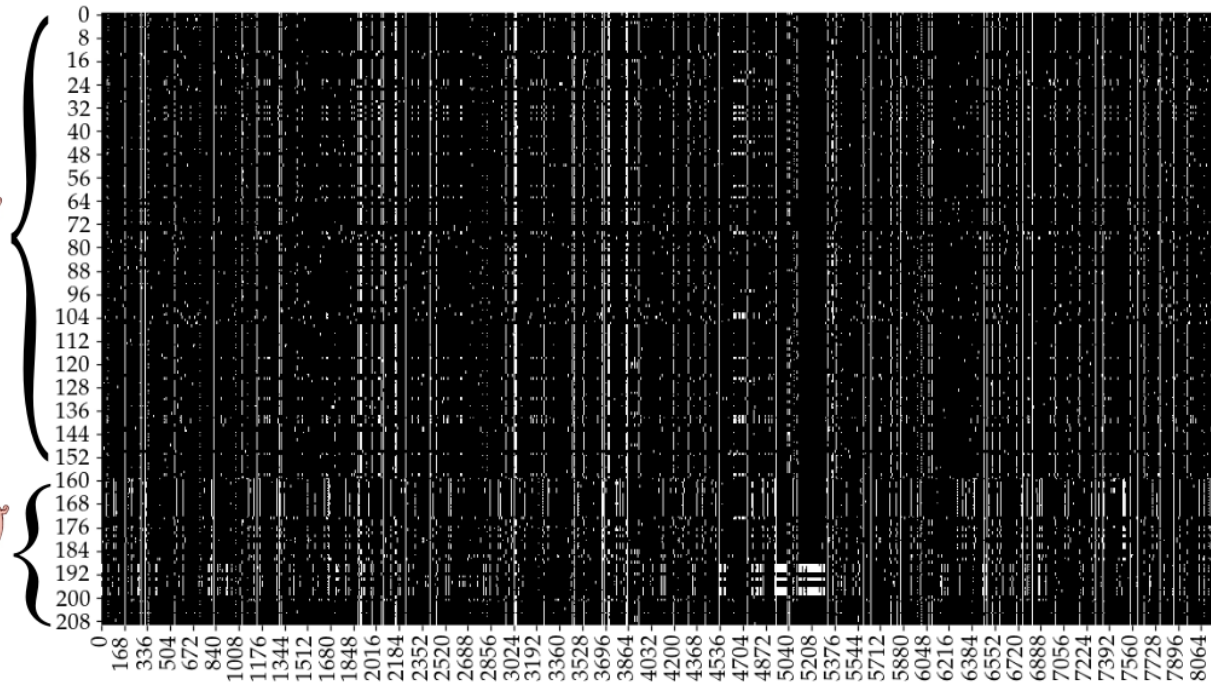
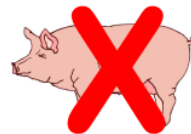
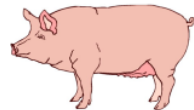
Reference : CCGTTAGAGTTACAATTCGA

Sample : CGGTTAGAGTAACTATTCCA

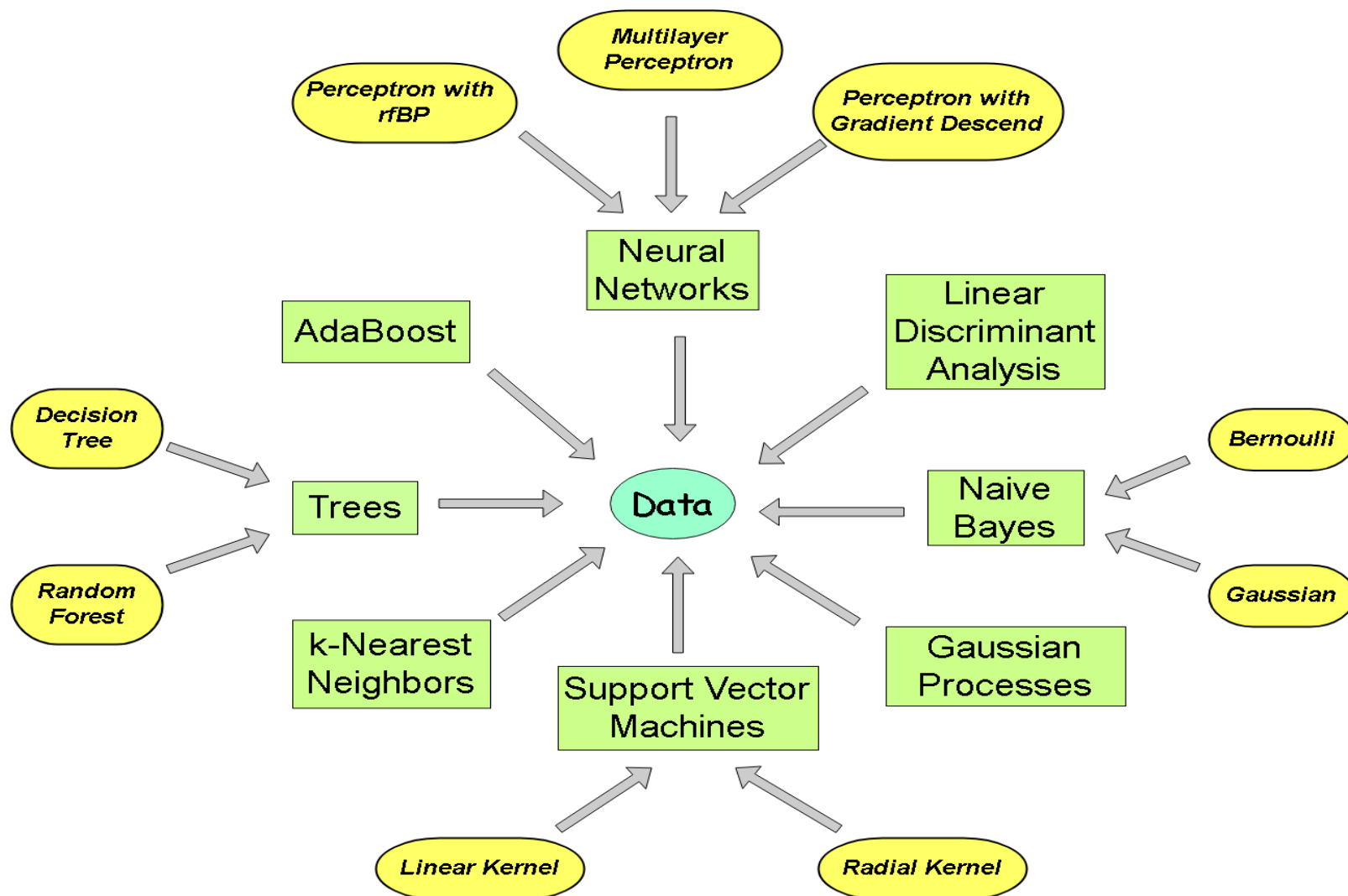
Binary SNPs : 01000000001001000010



Samples	210
Pigs	159
no-Pigs	51
Filtered Bases	8189



Filtered Bases



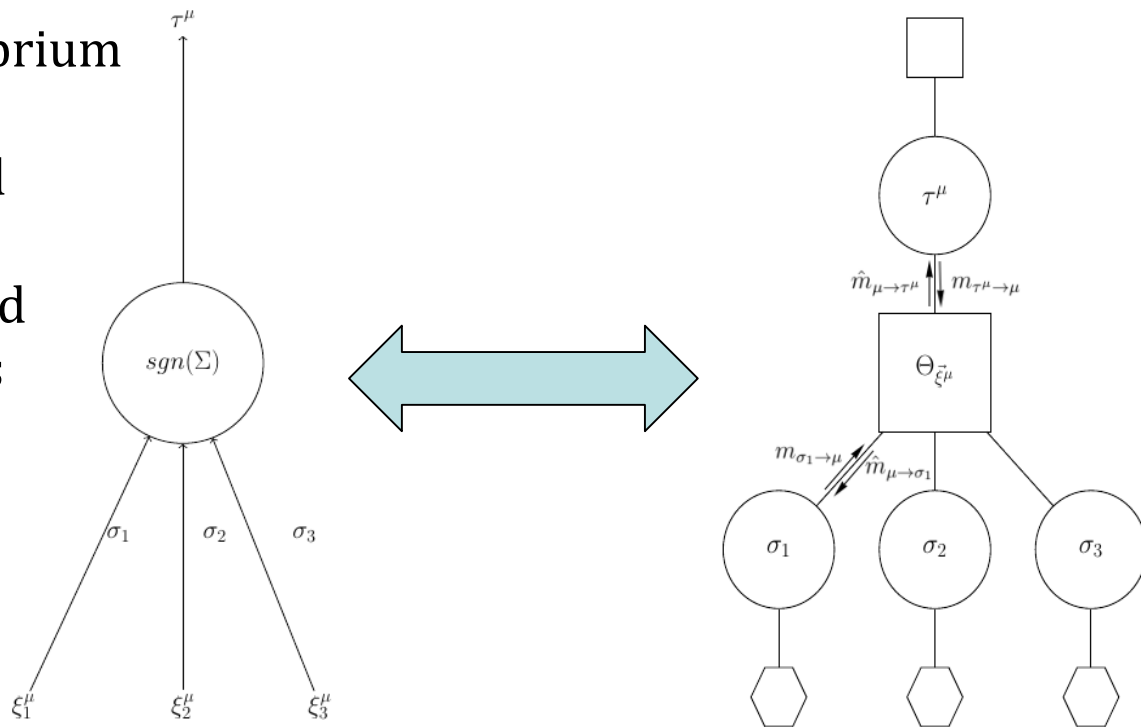
# replicated focusing Belief Propagation<sup>1</sup>

- Derived from a out-of-equilibrium distribution
- Entropy-maximization based learning rule
- Parameterized and reinforced Belief Propagation equations

## New implementation:

- C++ library
- Python wrapper
- Optimized for parallel computing
- Integrated with *scorer* library

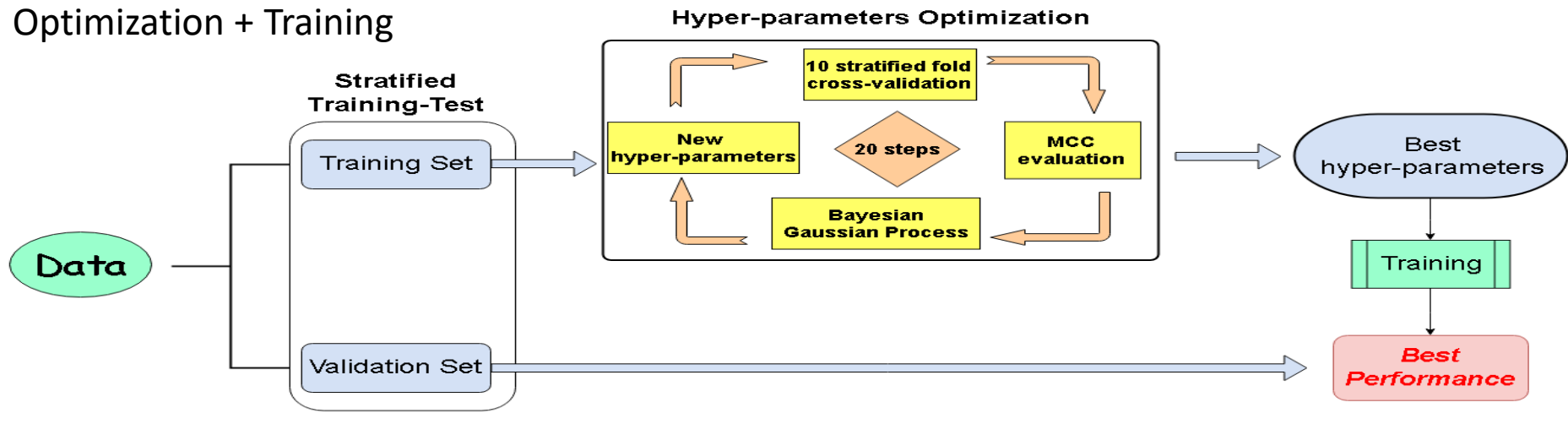
- Link: <https://github.com/Nico-Curti/rFBP>



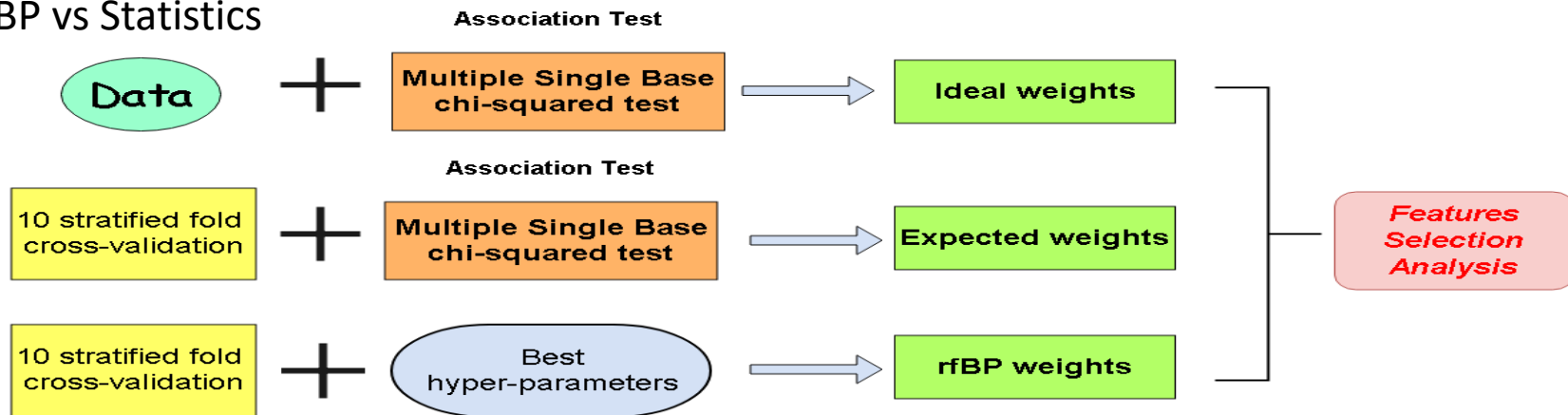
Binary Perceptron  $\longleftrightarrow$  Perceptron graph

<sup>1</sup>C. Baldassi et al. *Unreasonable Effectiveness of Learning Neural Networks: From Accessible States and Robust Ensembles to Basic Algorithmic Schemes*, 2016.

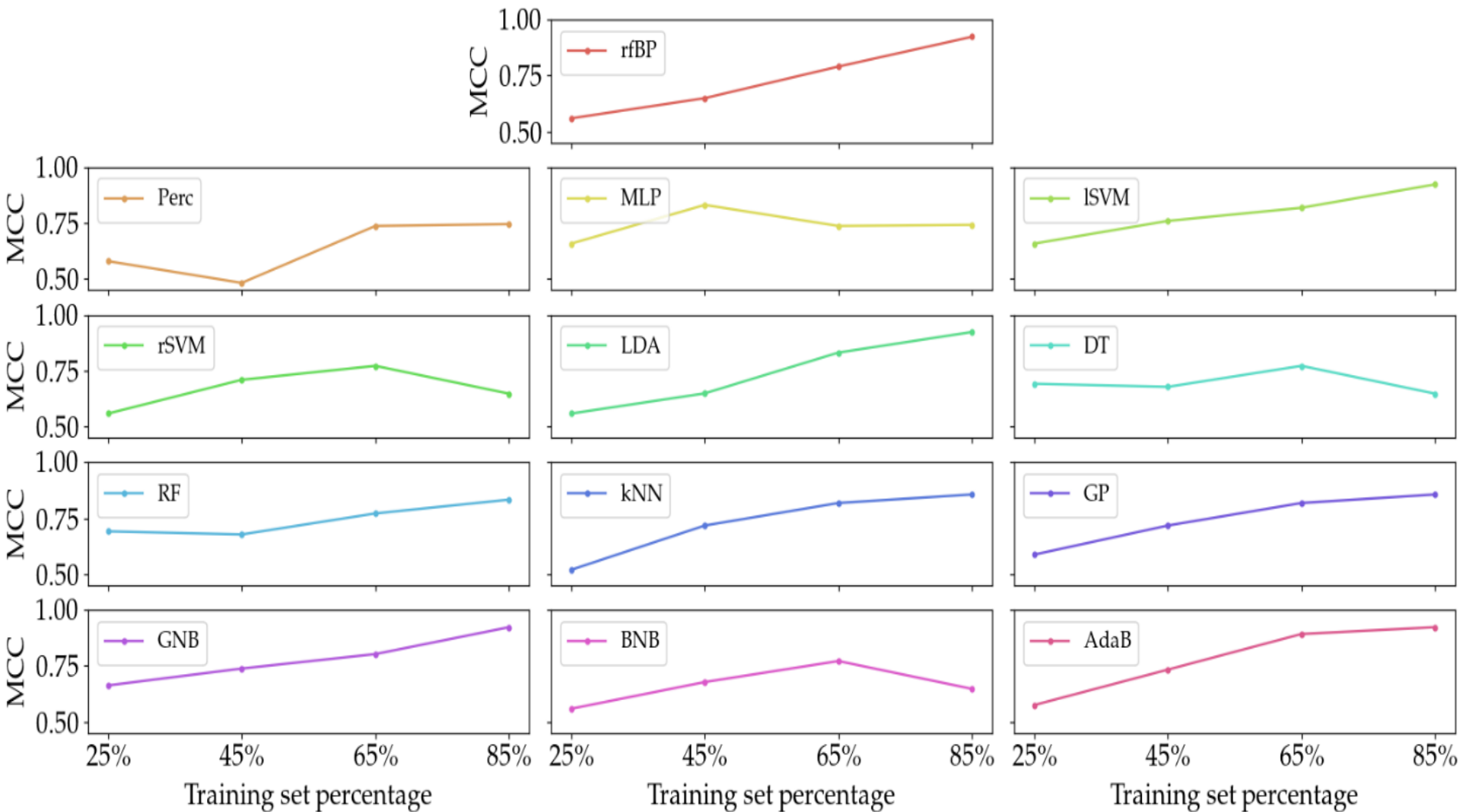
## Optimization + Training

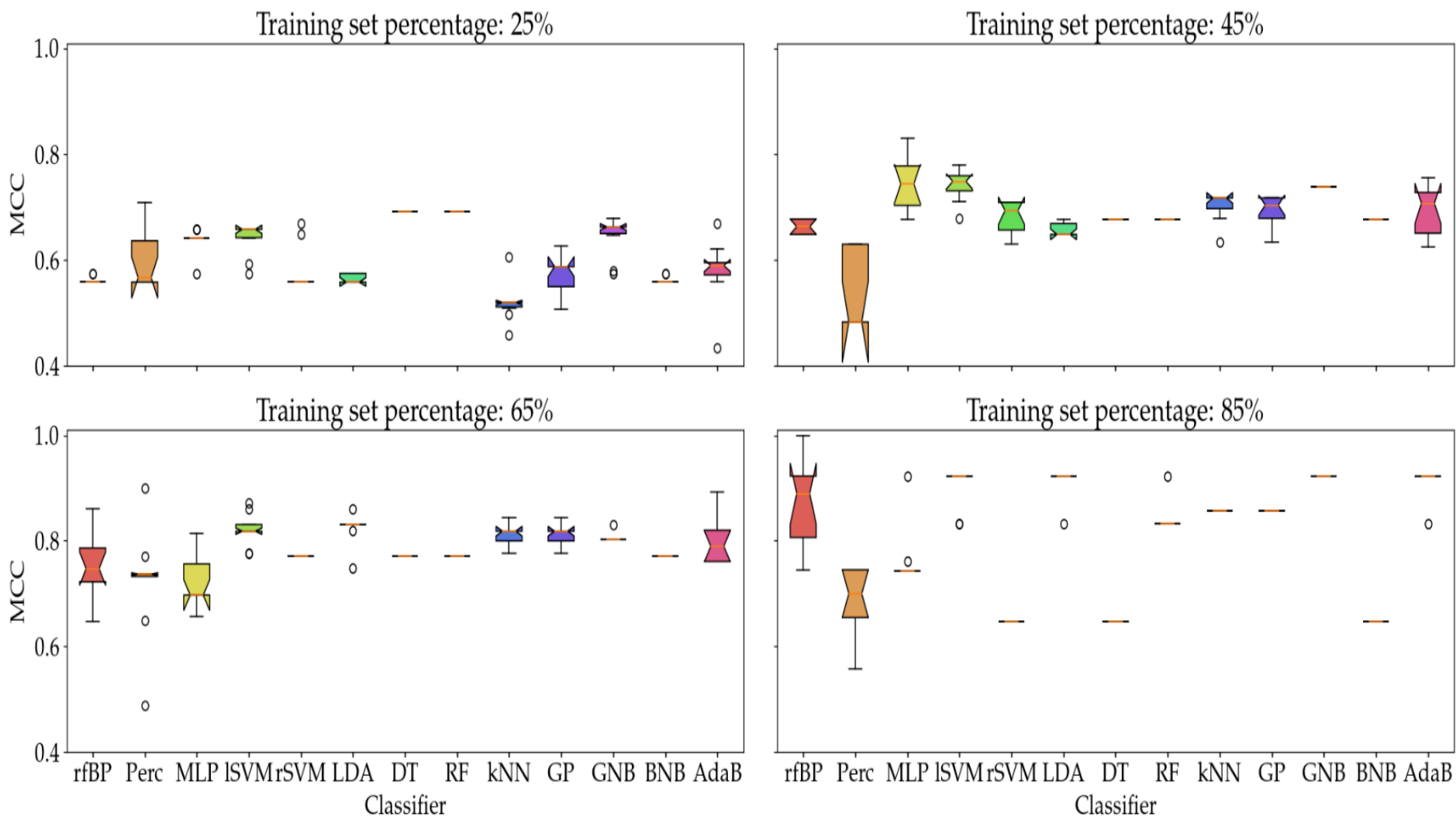


## rfBP vs Statistics



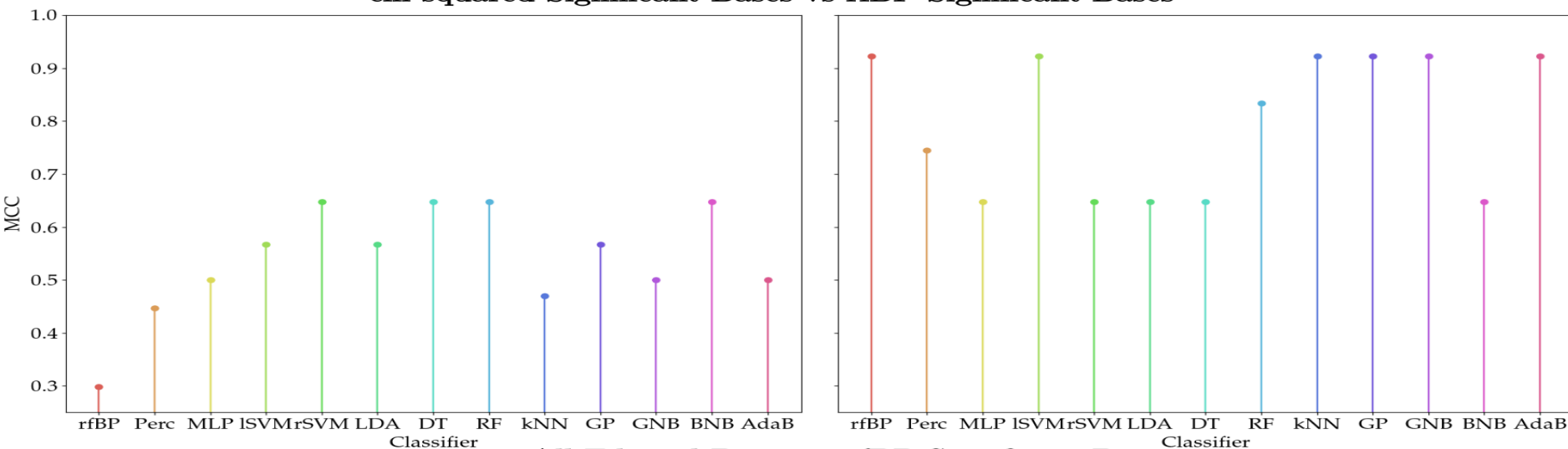




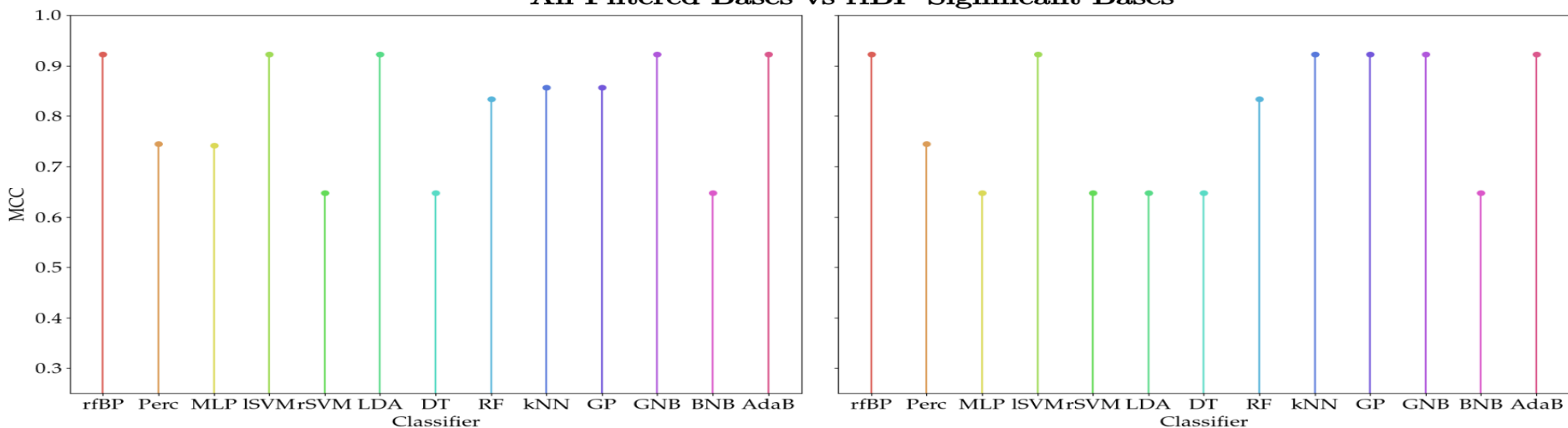




**chi-squared Significant Bases vs rfBP Significant Bases**

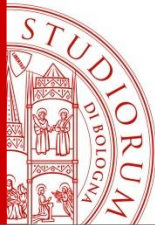


**All Filtered Bases vs rfBP Significant Bases**



- Classification of binary SNPs for Source Attribution is optimally performed by Entropy-maximization based algorithm called rfBP
- SNPs binary nature seems to favour rfBP binary properties
- rfBP significant bases are much better to classify than test association significant bases → can be actively used for Features Selection
- New C++ and Python (scikit-learn format) implementations are user-friendly and can be efficiently run on real data

Link: <https://github.com/Nico-Curti/rFBP>



# Acknowledgement



Horizon 2020  
European Union funding  
for Research & Innovation

