

Machine-learning models able to predict phage-bacteria interactions

Diogo Leite¹, Grégory Resch², Yok-Ai Que³, Xavier Brochet¹, Carlos Peña¹

¹School of Business and Engineering Vaud (HEIG-VD), University of Applied Sciences Western Switzerland (HES-SO), Swiss Institute of Bioinformatics (SIB), Switzerland

²Department of Fundamental Microbiology, University of Lausanne, Lausanne, Switzerland

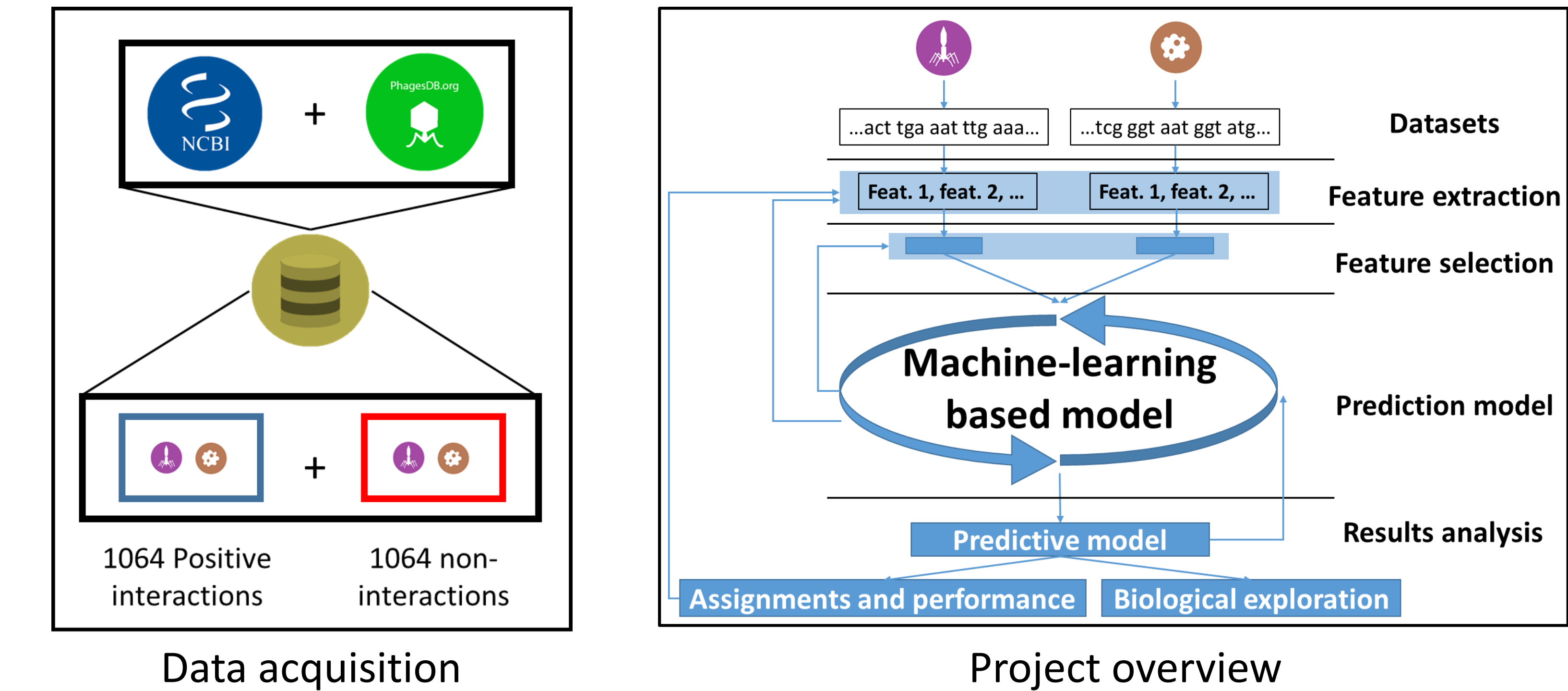
³Department of Intensive Care Medicine, Bern University Hospital (Inselspital), Bern, Switzerland



Abstract

Phage-therapy, a promising alternative to antibiotic-resistance, uses phages to infect and kill pathogenic bacteria. It requires finding perfectly matching phage-bacterium pairs, a time and money-consuming task, currently achieved empirically in laboratory. Our project aims at improving this task by predicting, *in-silico*, if a given phage-bacterium pair would interact. Predictions are performed on the base of public genomic data combined with machine-learning algorithms. With such an approach we have obtained around 90% of predictive power. In order to improve these results, we will extend our methodology and we will validate it with newly-generated clinically-relevant data.

Overview: prediction of phage-bacteria interactions



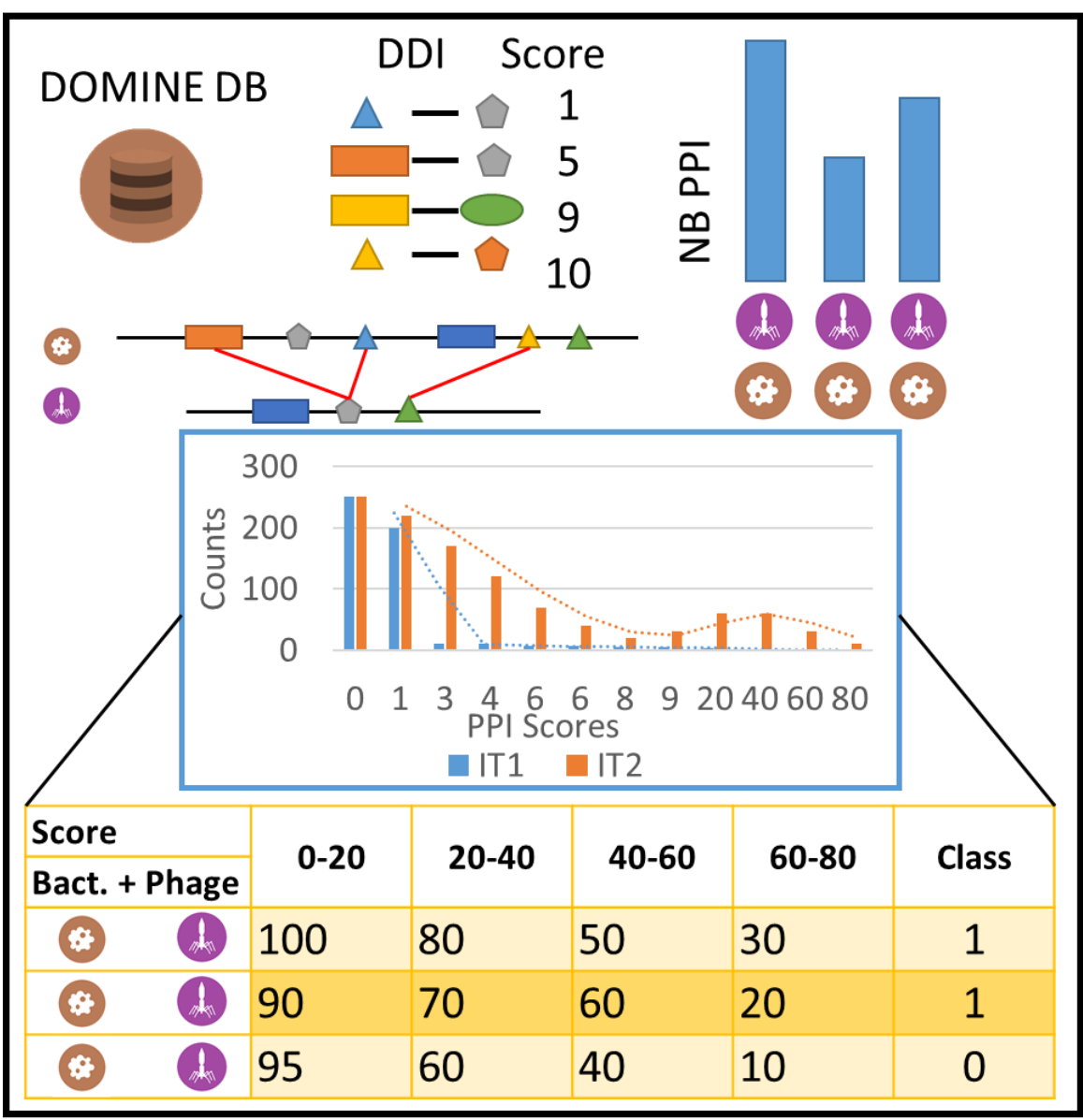
For our project we:

- Acquired data from public databases as NCBI [1] and phagesdb.org [2]. From GeneMarkS [3] we predict genes in bacterial and phages genomes.
- Constituted a positive dataset with 1064 phage-bacteria interactions pairs.
- Developed a model able to predict phage-bacteria interactions

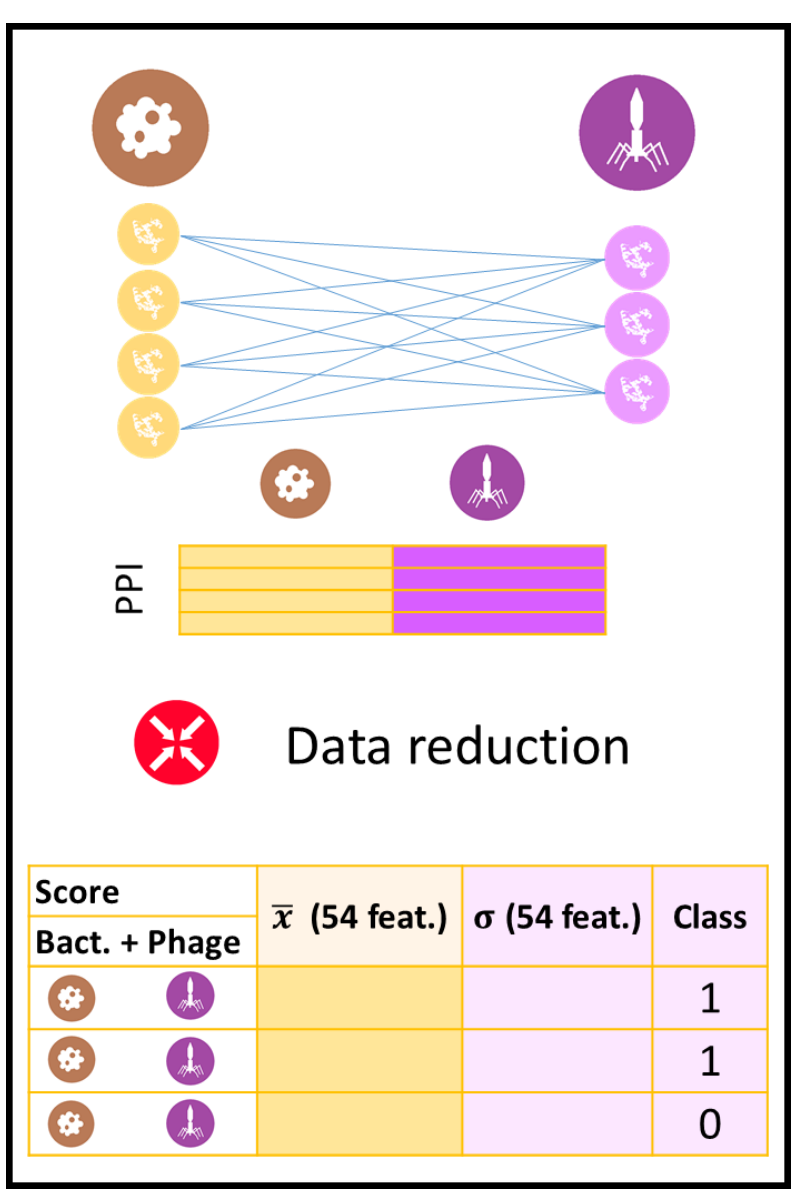
Feature engineering: obtention of informative features

On our data we:

- Extracted two kinds of features based on:
 - Protein interactions (PFAM Domain-domain interactions) - 18 datasets
 - Genomic sequences (% of amino acids, chemical components and weight) - 1 dataset
- Corrected the over-representation of the Mycobacterium smegmatis mc2 155 bacterium, reducing it from 86% to 14%

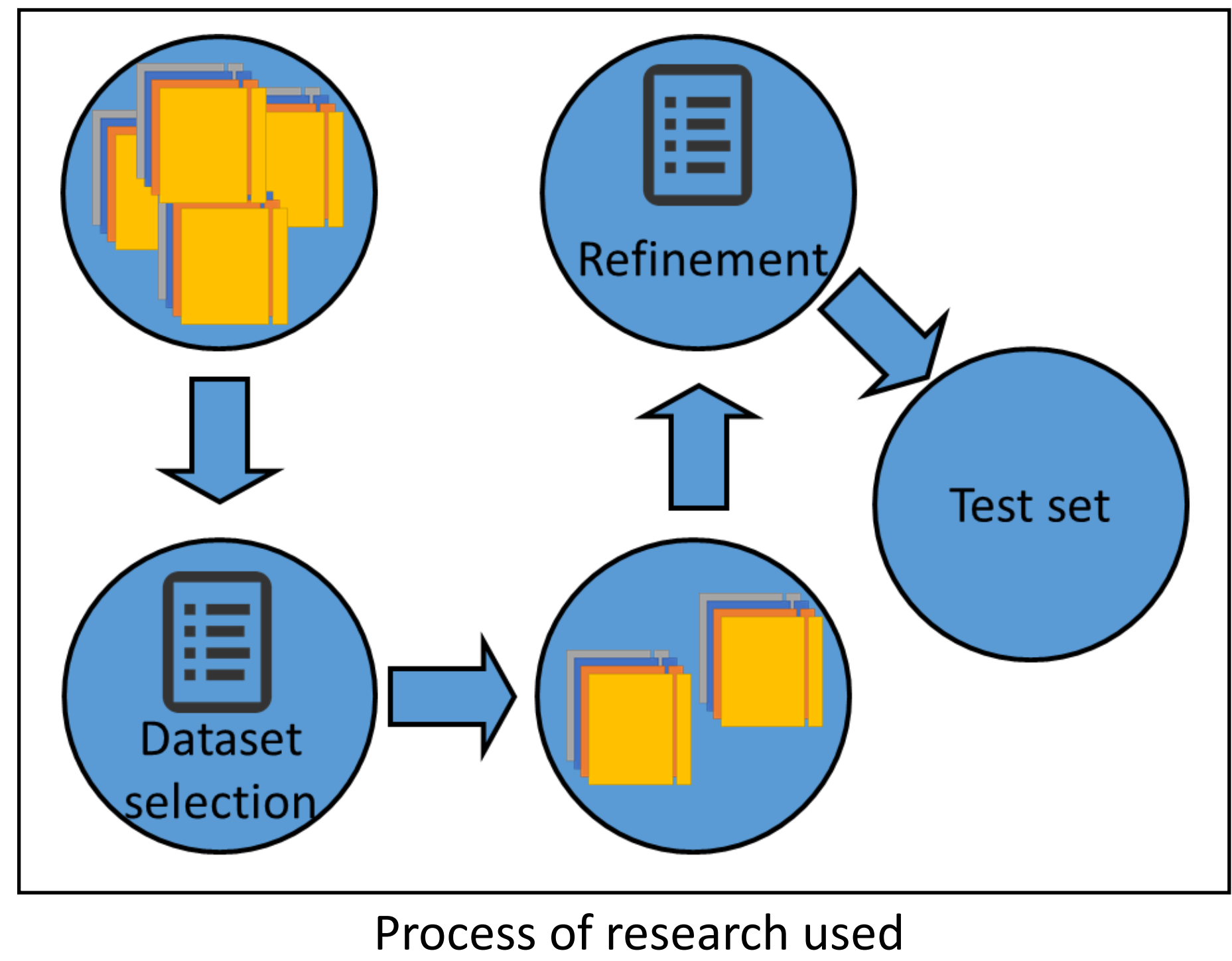


Domain-domain scoring



Sequence quantification

Machine-Learning based modeling



In our machine-learning search we used:

- 19 datasets (12'918 samples)
- Predictive model building with 10-fold cross-validation
- Four approaches: K-Nearest Neighbors (kNN), Random Forests (RF), Support Vector Machines (SVM), and Artificial Neural Networks (ANN)

In order to identify the best datasets and configurations of algorithms, we performed the following:

1. Selected the datasets exhibiting the highest scores across the four approaches
2. Further modelling to refine the algorithms' hyperparameters

Main results and conclusions

The best results on the test, obtained by ANN with 9 neurons and 50 epochs are presented below:

Datasets	Accuracy	F-Score	Sensitivity	Specificity
NB50	89.78%	90.13%	89.56%	90.12%
NBN50	89.79%	90.13%	89.56%	90.12%
S1e ⁻⁶	85.79%	86.24%	85.43%	86.35%

Future work:

- Explore other types of features to further increase model performance.
- Improve the predictivity of the dataset by pruning redundant/correlated variables that may perturb modeling
- Search for new relevant interactions allowing to predict interactions for different strains of a given bacterial species
- Improve domain interactions scores