

# Protein function prediction for poorly annotated species

# **Introduction**

### **Computational methods for PFP and network-based methods**

Protein function prediction (PFP) is one of the most important aims of modern biology. In crop and livestock species, PFP is conventionally based on annotation transfer from the few well-studied species, such as *Arabidopsis* and humans. While successful, these methods rely on the assumption that homologous proteins share function, which has been proved wrong in many cases [1]. It is thus desirable to complement the orthology-based methods with other approaches.

Network-based methods, for instance, infer the function of proteins exploiting the principle of guilt-by-association. Based on this, proteins that interact are likely to have similar function [2]. These methods have a lot of potential because they can utilize the information generated by high-throughput biological experiments, such as co-expression [3] or protein-protein-interactions [4] to construct networks from which to infer function.

### **Networks in poorly annotated species**

Because a wide range of data can be combined in these networks, the network approaches seem particularly relevant for poorly annotated species, such as agricultural species, where the validated data of a particular kind (i.e. coexpression, protein-protein-interactions...) is scarce [12]. The identification of putative transcription factors or hub genes via networks, for instance, may be more relevant for poorly annotated species in which knowledge about key regulatory elements is more limited. Network approaches, however, are also more challenging for these species because the data may be insufficient to carry a network analysis. A previous study, has shown that it is possible to develop network-based methods that can utilize the limited network resources of some crop species like rice, poplar, soybean and tomato, and achieve accurate PFP [5] by combining different data sources. In their approach, they combined data co-expression and protein-protein-interaction data, as well as information from other well annotated species, such as *Arabidopsis Thaliana*.

In livestock species, there is increasing interest for functional annotation. Efforts such as The Functional Annotation of Animal Genomes consortium (FAANG) [6], are currently generating functional annotations for some relevant species such as pig and chicken. Although network data for livestock species may be even more limited than for the aforementioned crop species, there are some studies that have used network data from livestock species. Stanley et al., [12], for instance, used a co-expression networks in chickens to infer function via defining GO-enrichment-modules. [In their approach, they stressed the importance of network methods to identify gene modules, as well as regulatory genes that are relevant for a set of functions or functional cascades.](#) This approach based on GO-enrichment, however, is indirect in that it first identify identify functional modules in the network. Sharan et. al. [16] judged the direct methods (directly predict the function of a protein) as slightly superior to the indirect ones. Furthermore, the method does not enjoy the advantage of the direct statistical-learning-based methods. Methods based on statistical learning have more potential than other methods because they can identify combinations of features that correlate with certain functions [20]. An interesting question therefore is whether it is possible to use some of the direct statistical-learning-based methods that are efficient for PFP, in livestock species. [Furthermore, it is interesting to investigate whether the network defined for PFP can be utilized to gain some understanding on the co-expression cascades of genes that redeem functions at different levels of specificity.](#)

### **BMRF**

In order to develop a statistical-learning-based network method that is efficient for livestock species, a logical approach is to utilize one of the state-of-art methods used for crop species. Bayesian Markov Random Field (BMRF) [1] is a prediction method that was developed with the purpose of achieving accurate predictions when the data was far from complete.

Add here a paragraph about methods predecessors of BMRF?

??

??

The prediction ability of BMRF was compared to other methods' in the Critical Assessment of protein Function Annotation (CAFA)[X] experiment and its prediction performance was high for some species (Illustration 1). BMRF is particularly efficient for poorly annotated species for two reasons mainly: First, it can synthesizes heterogeneous data into one network. For instance, it can integrate co-expression from the same species, as well as from related species; and second, though Gibbs sampling, BMRF can take into account unlabeled proteins to estimate the parameters of the model.

From a biological perspective, we would expect that the genes that are involved in certain specific functions are more prone to be involved in co-expression cascades and therefore have a larger number of neighbors in the network with whom they share a functional role. BMRF is a network method that directly exploits the information from the neighbors and therefore we would expect that predictions with BMRF will be more accurate for those genes that are co-expressed with genes that share a functional role. Consequently, it is interesting to investigate whether BMRF could be useful for the task of identifying relevant genes in the co-expression cascades. Questions to be addressed are whether the more specific GO terms have a higher degree of connection in the network and whether this is translated into a better accuracy of PFP with BMRF. Also, whether genes that are associated with a large number of GO terms of different levels of specificity are better predicted with BMRF.

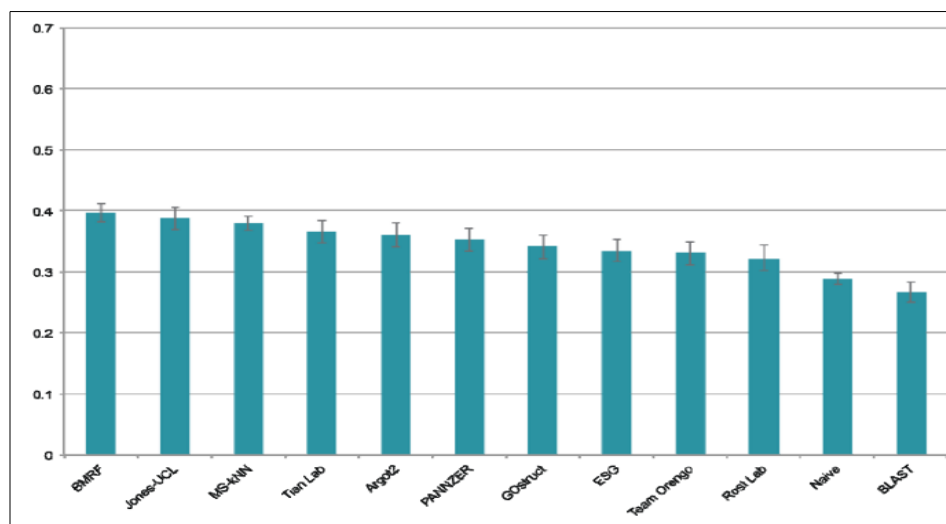


Illustration 1: Comparison of BMRF with other PFP methods. Evaluation for the Biological process category in *H. Sapiens*. Source: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3584181/f>

### **Positive-unlabeled learning**

We discovered an important potential problem with BMRF. The learning process may be biased because BMRF attempts to solve a one class classification problem where the annotated data consists solely on positive cases. This is because it is very difficult to be certain that a protein does not perform a function. In fact, this problem applies to many situations in bioinformatics because from a biological perspective, the lack of evidence for a connection does not imply that such a connection does not exist. As an example, for rice 415 proteins have experimental evidence for a biological process, but not a single protein has a validated proof of no-connection with a function. Since only positive associations are reported, the negative set is composed of all unlabeled data. This leads to some bias in the prediction because the unlabeled data may contain some positive cases. This problem increases with increasing numbers of unannotated proteins, such as in the case of network data from livestock species. To overcome this problem, a new type of machine learning has emerged called Positive Unlabeled learning (PU). With PU it is possible to identify the proteins that are more unlikely to have a given function. Hence, the number of unlabeled cases can be minimized by extracting some “reliable negatives” proteins from the set of unlabeled. It has been proved theoretically that, by identifying sets of reliable negatives, PU improves the performance of machine learning algorithms in situations where only positive labels are known [7]. PU has been successfully applied to a variety of problems related to PFP [7-10,13]. In [7], for instance, the authors extracted a set of reliable-negatives proteins from the unlabeled dataset by defining a threshold of similarity based on the euclidean distance between a set of positive proteins and the set of unlabeled. Yang et al. [8] developed a multi-label version of PU learning to identify genes associated with diseases; [9] developed two novel approaches to identify reliable negatives that can be applied in different algorithms. Jiang et al. [10] applied PU on a support vector machine and outperformed all pre-existing methods for pupylation sites prediction, and then Nan et al. [13] improved their method by adding as a first step to the algorithm in [10], the method described in [7]. Lastly, in another recent study, Nusrath et al. [14] used a Self organizing map to extract reliable negatives from unlabeled data-set of drug-drug interactions. None of these studies, however, have applied PU learning on a BMRF. BMRF

is one of the few methods that can handle networks with a large portion of unknown cases, and PU can extract reliable information from these unknowns. Our hypothesis therefore is that a PU implementation of BMRF will be particularly effective for PFP in species with limited network resources, such as pig and chicken. Furthermore, since with PU it is possible to create an alternative database with negative examples, the task of identification of regulatory elements is, to some extent, simplified.

The aim of this study is to develop a PU implementation of an existing Bayesian Markov Random Field algorithm that can efficiently assign proteins to GO term categories using network data from chickens. The methodology used could potentially be helpful in identifying genes that play relevant roles in the co-expression cascades.

### **Thesis setup**

Illustration 1, shows a diagram of the thesis setup.

<b>Part 1</b>	<div>a) Choice of network data for chickens</div> <div>b) Tune of the model parameters</div>	
<b>Part 2</b>	<div>Differences in prediction performance based on data available (using BMRF)</div>	<div>Section 2A: differences between species</div> <div>Section 2B: Impact of the quality of the data</div> <div>Section 2C: Impact of the characteristics of the co-expression analysis</div> <div>Choice of the improvement strategy</div>
<b>Part 3</b>	<div>Differences in prediction performance between GO-terms (using BMRF)</div>	
<b>Part 4</b>	<div>a) PU-BMRF development</div> <div>b) PU-BMRF performance evaluation</div>	
<b>Part 5</b>	<div>Biological support of the approach</div>	<div>Is there a biological explanation for the conclusions in parts 2 and 3?</div> <div>The results in part 4?</div>

Illustration 2: Thesis setup

## **Material and methods**

## **Data preparation**

Analyses were carried with co-expression data from three species: yeast, humans and chickens. In the case of yeast, the analyses were also performed with ppi data as well as co-expression data. Data sources in the different cases are given in Table X in Appendix II-Data overview.

Most of the analyses were carried in the three species and yeast ppi. However, some analysis were carried on yeast data because it was easily accessible, whereas some other analysis were carried on human as we thought that chicken data would not become available and it resembles more the situation in chickens.

In the GO-terms file, the associations were coded as "valid" if, for at least one of the association available in data, they correspond to Experimental evidence scores (precisely: 'EXP', 'IDA', 'IEP', 'IMP', 'IPI', 'IGI' ) and to the category of Biological process (BP), and as "non-valid" otherwise. Note that "valid" and "NONvalid" are just a level in a binary class that is defined for each associations and should not be misled with the label of a gene with regard to a GO term (0,1), as explained in the section "Validation in BMRF".

In order to maintain this distinction throughout the analysis between "valid" and "non-valid" associations these two groups of associations were up-propagated independently and then both sets of associations were combined into a single file, adding a third variable depending on whether the association corresponds to "valid or NONvalid" (table XX in Appendix II-Overview\_data ).

Domain and GO-terms files were pruned to exclude genes that are not available in the network file, as a requirement for the BMRF code. Then, the GO-size filter was applied to exclude the GO-terms that were too general or whose number of known associated genes was excessively low for the BMRF computations (Appendix I). The GO-size filter is based solely on the "valid" associations for two reasons: (1) the non-valid associations are not used in the validation and (2) the GO-size filter allows to make sure that there are enough number of genes in the validation. Analogous to the GO-size filter, BMRF uses another filter to exclude from the analysis the domains whose number of genes is below a certain threshold. This filter is named as 'DF-size filter' and is defined in table XX in Appendix I-Concepts.

Table XX in Appendix II-Data-overview shows one example of the first rows of each file type after data preparation.

## **Markov random fields**

Markov random fields are random field that satisfies the markov properties. In a markov random field the most likely discrete class of an element can be predicted by the joint probability distribution of its neighbors. A typical example is image restoration, where the value of a pixel (color) can be predicted based on the value of the neighbors pixels. The strength of markov random fields to infer the class of an element is that they allow for simultaneous predictions of many elements, and therefore they may achieve accurate predictions even when non of the neighbors are unknown. MRF are successful, in general, for prediction problems in which the principle of guilt-by-association holds, like in image restoration or PFP.

In a MRF the probability of a certain assignment of discrete states  $x=x_1,...,x_N$  is (as explained in [21]):

$$P(x) = \frac{1}{Z} \exp(-H(x)) = \frac{1}{Z} \exp\left(-\sum_{c \in C} H_c(x_c)\right)$$

Where N is the total number of variables, Z is the normalizing constant, C is the set of all the cliques in the network,  $H_c$  is a potential function associated with clique c and  $x_c$  is the assignment of states to the members of c.

Computing this is hard and often the equation is kept only until the second order. Also, equation XX is often homogenize by defining the same potential function for all cliques of the same size. Thus, we have:

$$H(x) = \sum_{v \in V} H_1(x_v) + \sum_{(u,v) \in E} H_2(x_{u,v})$$



In the PFP context, the proteins are represent as nodes in a network and the edges represent a similarity relationship between the two nodes they are connecting (for instances, either protein-protein-interactions, domain information pr co-expression between genes). The annotation is 1/0, either the protein has the function or not. Deng et al., 2003 adapted MRF to the PFP problem and stated that the probability over the enter network is proportional to  $\exp(\alpha N_{01} + \beta N_{11} + N_{00})$ , where  $N_{01}$ ,  $N_{11}$ ,  $N_{00}$ , correspond to the number of pairs of proteins that, while interacting, only one has the function, both of them have the function, and non of them has the function, respectively. And  $\alpha$  and  $\beta$  are weighting the contribution of each of these classes of pars of proteins. Then, combining the a priori probability of an assignment with  $N_1$  '1's, which depends on the frequency of the function and is proportional to  $(f/(1-f))^{N_1}$ , they obtain a homogeneous second order MRF and the following equation to estimate the probability that protein  $v$  is assigned with the function given the annotations of its neighbors  $N(v)$

$$P(X_{(v)}=1|X_{N(v)}=1)=\text{logit}\left(\log\frac{f}{1-f}+\beta N(v,1)+\alpha(N(v,1)-N(v,0))-N(v,0)\right)$$

where  $N(v,i)$  is the number of neighbors of  $v$  that are assigned with  $i \in \{0,1\}$  and  $\text{logit}$  is the logistic function  $\text{logit}(x)=1/(1+e^{-x})$ . [Deng et al \(2003\)](#) estimate the two parameters of the model using a quasi-likelihood method and apply Gibbs sampling to infer the unknown functional annotations. The approach has two steps: first, they estimate the parameters and, second, they infer the label using gibb-smapling. The parameters are estimated by maximizing the PLF with logistic regression. In this, each protein is a statistical unit, the predictors are  $N(v,1)$  and  $N(v,0)$ , and the response is the assignment. However, because some proteins are not annotated, the response will be missing for these and there will be uncertainty within the predictors. Deng eta la, overcame this by simply ignoring the unannotated proteins in the parameter estimation step. This can be problematic when the number of unknown proteins is large because by ignoring the unknowns, the neighbors are pruned and they may no longer express the complexity of the network.

### **Bayesian markov random fields**

Janis et al., 2010 developed a Bayesian Markov Random field (BMRF) to overcame the aforementioned problem. In BMRF, a MCMC algorithm is used to sample from the joint posterior density of  $\alpha$  and  $\beta$ . Thus, the label of the proteins is iteratively updated conditionally on the parameters  $\alpha$  and  $\beta$  though Gibb sampling. Thus, a candidate point  $\theta'=(\alpha', \beta')$  is obtained using the equation:

$$\theta'=\theta+\gamma(Z_{R1}-Z_{R2})+\epsilon$$

where  $\theta$  denotes the current state of the parameter vector, is the scaling parameter  $\gamma \sim U(\gamma/2, \gamma')$  is the scaling parameter and  $\gamma'=2.38/(\sqrt{2d})$  is the optimal step size[41], where  $d$  is the parameter dimension ( $d=2$ , in this case).  $Z_{R1}$  and  $Z_{R2}$  are randomly selected form past samples of the Markov Chains stored in the matrix  $Z$  and  $\epsilon \sim MVN(0, 10^{-4})$ .  $\Theta^*$  is accepted using a Metropolis step, with probability:

$$r=\min\left(1, \frac{PLF(x^{(t)}|\theta')}{PLF(x^{(t)}|\theta)}\right)$$

The labeling vector  $x$  is initialized using the output of the MRF-Deng, as explained in [plos1]

### **Validation in BMRF**

BMRF [1] was used to do PFP and learn about the impact of different method and network parameters on the prediction performance. In this framework, predictions are made individually for each GO-term that passes the GO-size filter (see Appendix I - Concepts). The predictions, however are not completely independent for each GO-term in the dataset because the genes that are not associated with any of the GO-terms in the dataset are treated differently. These genes are coded as "unknowns", and the number of unknown genes depends on the number of GO-terms in the database (which can be regulated through the GO-size filter). Each gene in the network file will enter the BMRF code with one of three possible labels: '1', '0' or '-1', where '1' stands for positive,

'0' stands for unlabeled (it is not known whether the gene has the function) and '-1' stands for unknown (it is not known whether the gene has the function and its label in the training set will be '0' or '1' based on Gibbs-sampling taking into account the label of its neighbors). Thus, the training set consists on those genes that enter the BMRF prediction function with a '1', or a '0', and the test set consists on those genes which are labeled as a '-1'. Note each gene have one label for each GO term (1 or 0) depending on whether the gene is associated or not with the GO term. However, that the unknown are labeled as '-1', for all the GO terms.

Then, the function provides a vector of posteriors corresponding to the probability that each gene has the function, and these posteriors are transformed into a '1' if the is above a certain threshold or '0' otherwise. The threshold is chosen in such a way that AUC for the test set is maximized. Note that both, the train set and the test set will be predicted but the accuracy of prediction was solely based on the test set.

The advantage of treating the unknown genes (genes known to be associated with zero GO-terms) differently is that genes that have never been predicted as positives are less likely to be non-associated for a given function than those genes that have been found as positives for some function but not for the function of interest. This has to do with the fact that some function are more difficult to predict than others. Thus, if based on data, a gene has never been identified as positive for any function, it is more fair to assume that the function is particularly difficult to predict using experimental approaches, that assuming that a very low number of genes have the function. In other words, a large portion of unlabeled genes ("0") implies that the function is rare, (almost never observed), whereas a large portion of -1s implies that the function is difficult to predict. This distinction between unlabeled and unknowns genes, additionally, allows to do predictions in a fairer fashion. By treating the gene of interest as unknown (-1) instead of as unlabeled (0), the prediction of the gene-GO association we are interested in, will be more free from the prediction of the GO-term. If the portion of "-1" however becomes very large with respect to the portion of "0", however, Gibbs-sampling will fail in the relabeling, because it will expect that the portion of genes that have the function is very large.

The labeling at a given fold in the k-fold CV is as follows:

Gene classes in BMRF	gene class as defined in Appendix I-concepts	Label in BMRF input	Expected label in BMRF output	fold
Associated with the function but the association is non-validated	Positive "NONvalid"	1*		
Not known to be associated with any function	unkown	-1		
Associated with the function. Assoc is validated but it is masked in this fold by labeling as '-1'	Positive-train	-1	1	a
Not associated with the function. This information is masked in this fold by labeling as '-1'	Unlabeled-train	-1	0	b
Associated with the function. Assoc is validated and it is not masked in this fold	Positive-test	1		a
Not associated with the function. This information is not masked in this fold	Unlabeled-test	0		b

Table 1: Labeling of genes at a given fold in the k-fold CV

\* Only take label '1' if the parameter *Only\_EES* is set to 'False', otherwise these genes are excluded from the analysis.

*Fold a: k-fold CV for the positives. The positive genes (this includes all genes that are associated with the function and whose association is validated) are divided into k sets until the end of the analysis. Genes in one of these k sets will be in Positive-test, the rest will be in Positive-train. Within the same analysis, this is repeated k times. Each time corresponds to a fold with a new labeling configuration such as represented in Table 1. Note that only the label of those gene classes that have 'a' or 'b' in the fold column change with each fold. For each fold, a new k set takes place as Positive-test and the rest are Positive-train, until each of the the k sets that were defined at the start of the analysis have been assigned once to the Positive-test.*

*Fold b: k-fold CV for the unlabeled. Same as Fold abut for unlabeled genes instead of positives.*

### **Part 1- Choosing the data and tuning the model parameters for prediction performance using BMRF.**

A first step involved the choice of a co-expression threshold for chickens given a conditionally independent co-expression network. Conventionally, a pearson correlation of 0.7 is used as a threshold for co-expression. However, in the case of chickens, the number of co-expression experiments is limited and we a pearson correlation of 0.7 would lead to a scarcity of network data in BMRF. Furthermore, the assoictaions file would need to be heavily pruned because BMRF does not allow for any gene in the association file sthat is not in the network. Thus, we computed AUC for several pearson correlations thersholds. Note that for each person correlation, the GO file needs to be uppropagated. Details about which pearson correlation cutoffs were considered are given in section 1A in Appendix III-Additional results.

Since the assignments of genes to folds in the k-fold CV is a random process, we investigated how many replicartes were required to achieve reproducible results. Then, we tunned the follwoig model paramateres by changing values in each of them at a time: GO-size filter, Number of folds in k-validation, Number of iteration in Gibb-sampling. Lastly we computed AUC with and without domain information and non-validated associations. Details about which values were considered for each model parameter are given in section 1AB in Appendix III-Additional results.

### **Part 2- Differences in prediction performance based on data available (using BMRF)**

Brief explanation. Part 2 and part3 aim the selection of an strategy of improvement.. Part 2 consits of three sections:

Section 2A): Investigate the differences in prediction performance between species

Section 2B): Investigate the impact of the quality of the data on the prediction performance

Section 2C): Investigate the Impact of the nature of the network on the prediction performance

#### **Section 2A)**

Investigate the differences in prediction performance between species

We analyzed the differences in data between the species considered and we investigated how this difference is translated in a different prediction performance. Differences in network data that we considered were:  $\#te$ ,  $\#epp$ ,  $\#epn$ ,  $\#enn$ ,  $\#epp \cdot 100 / \#te$ ,  $\#epp / tpepp$ ; and differences in annotation data considered were:  $\#genes$  per GO-term,  $\#GO$ -terms per gene,  $\#edges$  per gene, and  $\#epp$  per GO term. These differences can be found in Appendix II-data overview, as well as in section 3a in Appendix III-Additional results. Appendix II provides a more graphical representation of these differences whereas Appendix III focuses on how these differences may be affecting the prediction performance. Formal definitions of these parameters can be found in Appendix I-concepts.

#### **Section 2B)**

Investigate the impact of the quality of the data on the prediction performance

In order to investigate the impact that the quality of the data has on the predictions, we computed AUC after random removal of associations and edges from the data. We distinguished between 2 types of associations: association of the GO-term of interest and associations of other GO-terms; and four types of edges: edges positive-positive,  $epn$ ,  $enn$ ,  $te$ . Portions subtracted were 0, 10, 30, 50, 90 and 95% when yeast data was used; and 0, 0.05, 0.1, 0.2, 0.4, 0.6, 0.8, 0.99 when human data was used. Then the correlation between AUC and the percentage of removal was computed in order to assess the impact of the amount of data considered and the prediction performance.

#### **Section 2C)**

Investigate the impact of the nature of the network on the prediction performance

We investigated how the conditions of the co-expression analyses influence the prediction performance. For this we took different subsets of co-expression data and assessed the prediction performance in each case. The source of expression data for yeast is organized based on experiments and a brief description of these is provided. Thus, we searched for key words, like "oxidation" or "stress" and created subsets of networks with data from those experiments that have those words in the description.

In the case of humans, data is organized by tissues. Thus, each tissue was a subset. We homogenized the subsets based on  $epp / tpepp$  instead of  $\#te$ , because in part 1 we learnt that the former parameter has barely any impact on the prediction performance. In order to homogenize, we removed randomly edges from the network. Then, we computed AUC with the expression data from the different tissue experiments.

We investigated the global effect that the 'nature' of the network has on the prediction performance, and we searched for evidence of biological support. For instance, from a biological perspective, we would expect that a co-expression analysis carried for one specific tissue will allow for better predictions in those GO terms whose function is more relevant for that particular tissue.

Based on the results from these 4 types of analyses, we chose an approach to improve the BMRF PFP method. All analyses were carried out with the network data and the values for model parameters chosen in part 1. In order to estimate the reproducibility of the approach, the standard

deviation across folds within the same replicate, as well as across replicates, were computed.

### **Part 3- Differences in prediction performance between GO-terms (using BMRF)**

We Investigated the impact of different GO-term-properties on the prediction performance. A total of 9 GO-term-properties were defined: epp/tpEpp, eppV/tpEppV, #genesV, sepc, teV/tpV, depth, AUC and sdAUC. We computed the correlation between these in order to investigate how they may be affecting the prediction performance. Definitions of these GO-term-properties are also given in Appendix I-concepts.

Given the definition of depth (appendix I-Concepts), the most general GO terms have a lower depth than the most specific GO terms. However, it should be considered that depth is not a good estimate of the specificity of the GO terms, because in the GO hierarchy, different branches have different levels. Thus GO terms with similar specificity can have different depth, depending on the branch. Moreover, some GO terms appear in more than one branch. For this reason another parameter was defined called “spec”, that corresponds to the inverse of the sum of all the validated genes from the 4 species considered.

### **Part 4a- PU-BMRF development**

Similarly to Part 1, in Part 2, the folds for the set of positives were created solely based on the validated associations, and the same holds for the unlabeled genes. We also used 10fold -CV. However, due to time constraints we performed the analysis only with 4 replicates (instead of 20 in the conventional approach of BMRF). Also, for simplicity, in part 2 the non-validated positive associations were always excluded from the analysis. This can be done by simply setting the “only\_EES” parameter to “True”.

Steps 1 to 6 aim the computation 86 features including 70 non-GO-specific features and 16 GO-specific features. These features were computed for each of the 1,656,000 total gene-associations (138 GO-terms x 12,000 genes). Note that, for chickens, only 138 GO-terms passed the GO-size filter.

Based on these features, we computed the euclidean distance between each of the genes in the set of known genes that were not masked (positive-train) and the set of unknown genes that were not masked (unlabelled train), and we estimated the average euclidean distance between the known unknown genes and the known positive genes. Then, we classified the genes in the test set as positive or negative based on whether their distance to the positive genes was larger or smaller than the distance from the known unlabeled genes to the positive set, as explained in the algorithm introduced by [1] (Illustration XX):

1.  $RN = \emptyset$ ;
2. Represent each gene  $g_i$  in  $P$  and  $U$  as a vector  $V_{g_i}$ ;
3.  $pr = \sum_{i=1}^{|P|} V_{g_i} / |P|$ ;
4.  $Ave\_dist = 0$ ;
5. **For each**  $g_i \in U$  **do**
6.      $Ave\_dist += dist(pr, V_{g_i}) / |U|$ ;
7. **For each**  $g_i \in U$  **do**
8.     **If**  $(dist(pr, V_{g_i}) > Ave\_dist)$
9.          $RN = RN \cup \{g_i\}$

*Illustration 3: Algorithm for extraction of RN.*

Source: [1]

Note that step 6 in the algorithm allows to introduce a constant to specify how strict we want to be. A value of 1 means that genes in the test set that are closer to the positive set than the average unlabeled genes will be classified as positive, but a constant of 1.5 would mean that we are more strict and only genes that are 50% closer to the positive set than the unlabeled genes will be classified as positive.

➤ **Step 1 – Similarity Matrix:**

A similarity matrix between the GO-terms is computed using the R package “GOSim”. The use of computing this matrix is two-folded: First, it allows to extract a set of unrelated GO-terms, in case we cannot carry the analysis for all the GO-terms that pass the filter (for instance, due to time constraints); and second, this matrix will be used in the computation of the features (step 5).

➤ **Step 2 – Creating the folds:**

The training and test-sets are created by randomly sampling genes among the positive associations for each GO-term. Then, for each fold, one GO\_file is created, in which the associations in the test-set have been excluded. Also, the set of genes that, after “hiding” the test set, are associated with the GO-terms are stored, and also their neighbors. Thus, for each GO-term, two objects are stored: the set of positive and the neighbors of the genes that are in the positives. Finally, another object is extracted with the neighbors of each gene. These objects are different for each fold and will be used in steps 4 and 5.

➤ **Step 3 – Network features:**

Transitivity, closeness and betweenness are computed for each gene-GO combination in three different networks:

- A network for all edges that have at least one node in the set of Positives. This is expected to be useful because two types of nodes are included (positive and unlabeled), and genes that are positive (either discovered or to be discovered) are expected to be more interconnected in this network.
- A network for all edges that have at least one node in the set of Positives, and all the edges that have both their nodes in the set of neighbors of positives (step 2). In this network it is also expected that the positive (either discovered or to be discovered) are more interconnected than the genes that are not associated with the GO-term.
- A network of all edges except those that link to genes in the positive set. We expect that in this network, the genes that are positive are less interconnected.

➤ **Step 4 – non-GO-specific features:**

The following features are computed for each of the 12,000 genes:

- features f1-f4 refer to the number of GO-terms of the gene. It is expected that the genes that are associated with a large number of GO-terms are more likely to be associated with a novel GO-term. We expect this probability to be higher for genes that are associated with a large number of specific GO-terms (before up-propagating) because they may be involved in a large number of functions. Genes that are associated with a large number of GO-terms only after up-propagating, however, may be associated with these simply because they are associated with more general GO-terms.

f1) The number of GO-terms the gene is associated with.

f2) The number of GO-terms the gene is validated-associated with

f3) The number of GO-terms the gene is associated with, in a GO file before up-propagating

f4) The number of GO-terms the gene is associated with, in a GO file before up-propagating

- features f5 and f6 refer to the number of GO-terms of the neighbors of the gene.

f5) The sum of the GO-terms that are associated with the genes that are co-expressed with the target gene.

f6) The sum of the number of GO-terms that are associated with the genes that are co-expressed with the target gene in a database where only validated GO terms are considered.

- features f7-f9 refer to the number of neighbor.

f7) The number of genes that are co-expressed with the gene of interest. BMRF accounts for this information, but only when the network data was extracted with a Pearson Correlation threshold of 0.35. In this step, we can provide information from other networks with other co-expression thresholds as well.

f8) The number of genes that are co-expressed with the gene of interest and are associated with at least 2 GO-terms.

f9) The number of genes that are co-expressed with the gene of interest and are associated with at least 5 GO-terms.

We expect that genes that are co-expressed with genes that have multiple functions are more likely to have multiple functions and therefore are more likely to be associated with a novel GO-terms. The thresholds of 2 GO-terms and 5 GO-terms in f8 and f9 were chosen based on the variability of the feature. We are interested in features that are highly variable across the genes.

- features f10 to f70 are same as f1-f7 but for different features Pearson correlation thresholds. Mainly: 0.1, 0.2, 0.35, 0.5, 0.6, 0.7 and 0.8. Note that as the network database changes, so does the GO-file because BMRF does not allow for any gene that is not in the network. Due to the constraints in the data that come after the different correlation thresholds, it is expected that if a gene is associated with a very large number of GO-terms when the Pearson correlation was high (i.e. 0.6), it must be a gene involved in many different functions, whereas it may be that other genes are associated with more GO-terms when the Pearson correlation is lower, and this should also be considered.

The possibilities of restricting the data-set based on the Pearson correlation cutoff greatly increases the information available. This is because based on different correlation cutoffs we may be able to observe different patterns in data. Note that this is particularly useful since we are using a condition-independent network where the data of different experiments is combined.

### ➤ Step 5 – GO-specific features:

For each gene, we computed up to 16 GO-specific features. First, we defined four intersection for each gene and most of the features defined in step 5 will be computed for each of these intersection (for each gene).

Intersections of genes:

(1) Whether the gene of interest is in the set of positives

(2) Genes that are found in the neighbors of the gene of interest and in the set of positives

(3) Genes that are found in the gene of interest and the neighbors of the genes in the set of positives

(4) Genes that are found in the neighbors of the gene of interest and the neighbors of the genes in the set of positives.

F1-f4) Following from step 3, for each gene, we compute the sum of the betweenness, transitivity and closeness of the GO-terms that are in the interactions 1-4.

f5-f9) The number of genes in intersections 1-4 and the portions of neighbors of the genes of interest that are in the intersection s1-4.

F10-f13) The number of domains of the genes in intersections 1-4 the number of genes that share domains in intersections 1-4, and the number or unique domains that are sheared between genes in intersections 1-4.

f14-f16) The number of genes in interaction 1-2 (neighbors of the genes in the positive set are not considered here) weighted by the degree of similarity between the GO-terms they have in common and the GO-term we are interested in.

Note that these features vary depending on the train and test set, and therefore they need to be computed k-times per replicate.

#### ➤ **Step 6 – extraction of RN:**

The databases with features information obtained in steps 4 and 5 were combined and the values of the features were scaled by dividing each value by the square root of the squared sum of all the values. Thus, for each GO-term we have a database with 86 features per gene. Checking the label of the genes that are in the training set, we apply the algorithm in illustration XX:

We check whether any of the RN is in the set of non-validated positive cases and we removed from the analysis those that did so. We tried different thresholds in step 6 of the algorithm in Illustration 1 (default value is 1) and through an iterative process we adjusted the threshold to the highest value that allows to extract a maximum number of RN. We gradually increased the threshold by 0.05 if the criteria was not satisfied (thus, if the number of RN was excessively high for our purpose). We proceeded the analysis for different values of “maximum number of RN”, mainly, for 1000, 2000... and 8000.

In step 6, we also extracted the same amount of RN but through random extraction and stored them separately. Two approaches were also used to extract RN. Instead of defining a fixed number of RN, we allowed the threshold to change according to a desired value of AUC in the process of extraction of RN. For instance, we can specify that we want to extract as many RN as possible as long as the AUC is equal to 1 in the process of extraction. The genes that are used to evaluate the performance of extraction are those in the test-set. These genes were excluded from the analysis in step 1 of the PU-BMRF and therefore were not considered to define the threshold in step 6 of the algorithm in illustration 1.

Finally, in order to be more reliable about the RN not being positives, we excluded from the set of RN those genes that while being in the set of RN they were and also in the set of positive-NONvalid.

Note that a different set of RN will be extracted per fold, per replicate and per GP term.

After the extraction of RN we train the BMRF classifier (RN vs positives), instead of (unlabeled vs positives) in the conventional BMRF approach.

#### **Part 4b- PU-BMRF performance evaluation**



- Evaluation of the process of extraction of RN

We computed the accuracy of extraction of the RN. Note that, in principle, only unlabelled genes should be classified as RN. For AUC in BMRF, the expected label for the positives-train is 1, and the unlabeled-test as 0. And the predicted label was '0' if the gene was classified as RN, and '1', otherwise. This means that a gene whose expected label is '1' and its predicted label is a '0' is a false positive, because it is a gene that is positive and therefore it should not have been classified as RN. Similarly, a gene whose expected label is 0 and whose predicted label is 1, may have been properly classified (it is and unknown gene that was not extracted as RN); a gene whose expected label is '1' and whose predicted label is '1', has been surely properly classified because it was a positive gene and it has not been classified as RN; and lastly, a gene whose expected label is '0' and prediction label is '0', has been properly classified as well, because it is an unlabelled gene that was classified as RN.

Since the focus was on the accuracy of extraction of RN we extracted from the computation of AUC those genes for which we expect a '0'; and we observed a '1'. So the unlabelled genes that were not classified as RN. Note that if these were included in the computation of AUC, AUC could be inflated based on this very large subset of genes when in fact we cannot tell if they were properly classified. This is an important aspect to consider, especially since we want to compare the accuracy of prediction for different number of RN and this class of genes may become very large when the number of RN is very low (for instance setting a maximum of 1000).

- Reproducibility of the process of extraction of RN

In order to assess the reproducibility of the extraction of RN across replicates, we investigated which portion of the RN extracted were common in all the replicates of the analysis. Further, we assessed the reproducibility across folds within the same replicate. For this, we calculated how many different RN were extracted per replicate (combining the extraction of each fold).

- Prediction performance using the RN

The positive genes and the genes in the set of RN were used to train the BMRF classifier. It is expected that the classification will be more accurate now that the set of unlabelled genes have been replaced by a smaller set of RN. Apart from that, the RN were treated in a similar way than the unlabelled genes in the conventional approach of BMRF.

## **Part5-Biological support of the approach**

In principle, PU-learning should contribute in improving the epp/tpepp by removing edges from the data. From a biological perspective we would expect that the genes that are involved in more specific functions have a higher epp/tpepp. The aim of this part is to investigate whether we observe that principle in the data, whether the ratio improves more in genes involved in specific function or in general functions, and finally how this related to the increase in accuracy with PU-BMRF with respect to BMRF in both types of genes.

To investigate whether the principle of more specific-better connects, we followed two approaches. First we computed the correlation between specificity and epp/tpepp, and second, we compared epp/tpepp in two groups of genes depending on whether these genes are associated with specific or general GO terms.

To compare epp/tpepp in these two groups, we looked into genes that are exclusive from general or from specific GO terms. For this, we subsetted GO terms with a specificity larger than 5 and less than 3000 associated genes. The first requirement is to make sure that the GO terms considered have at least a minimum of edges. Otherwise we would be taking as specific GO terms GO terms that are simply associated with almost none of the genes. The

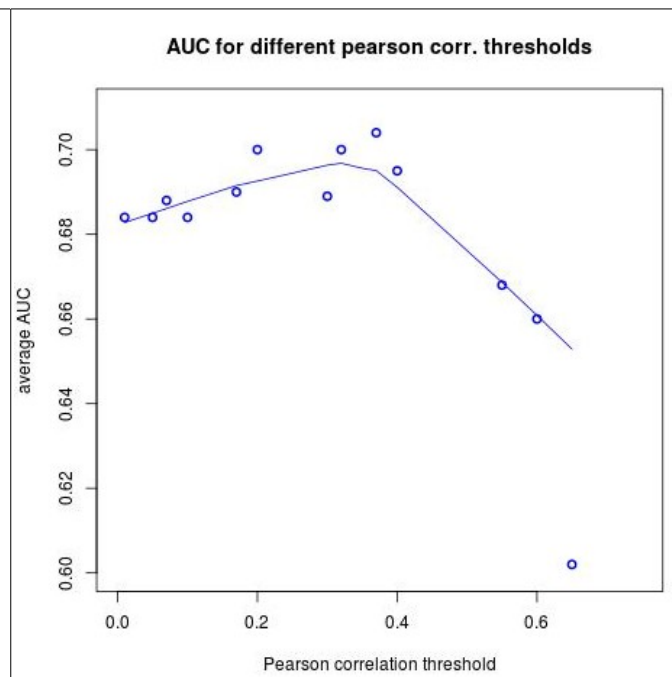
second requirement is guarantee that there are some genes in the specific GO term that are not in the specific GO terms. For instance, if we choose the GO-term XXX-Biological process, which is associated with XXX genes, it would be very difficult to find genes that are associated with a more specific GO term but not with XXX.

Then we distinguish between two groups of GO terms, one specific and one of general GO terms. We decided that the specific GO terms should have a maximum of 7 genes whereas the general should have 2700. The choice of these values is trivial because we want a large difference between the two groups, at the same time that enough number of representatives in both groups. Furthermore, if the groups are too large, then it will be very difficult to find genes that are exclusive only of one of the groups. Then we sampled 200 genes that are exclusive from each group and we compared the number of edges that connect the genes of each of these sets.

## Results

### **Part 1- Choice of the network data for chickens and tune of the model parameters.**

The prediction performance was higher for chickens when a Pearson correlation of 0.35 was used as a co-expression threshold, as it can be observed in Illustration 4. We therefore carried the next analysis on chickens with this co-expression cutoff.



*Illustration 4: AUC for different Pearson correlation thresholds for the chickens' co-expression network.*

*The Y-axis corresponds to the average AUC for all GO-terms considered in the analysis. The line corresponds to the lowess line. Values corresponding to the plot are in table 1.*

After tuning the model parameters we come with the following conclusions about how to carry the next analyses:

- 10-fold CV will be used, because the accuracy of prediction is significantly larger than for 5-fold CV.
- Each analysis will be carried with 20 replicates, because the sd across replicates is low ( $<0.02$ ) on human data.
- The number of iterations for the Gibbs-sampling in BMRF will be set to 30, since this number is enough even when the number of unlabeled genes is large
- A GO-size filter of (20 , 0.1), for minGOsize and maxGOsize, respectively. This filter specifies that no predictions will be made for the most general GO terms (since they are less relevant), and it guarantees that there will be enough number of positive cases in each fold when BMRF is used.
- We will use domain information, as well as non-validated associations because they seem to improve the performance.

*Table 2: Conclusions from tuning the model parameters.*

## **Part 2- Differences in prediction performance based on data available (using BMRF)**

The conclusions from this part are:

### **Section 2.a) Impact data differences between the species on the prediction performance**

Table 5 shows the prediction performance for the species considered using BMRF.

	yeast	yeast_ppi	humans	Chicken
# GO terms	1,102	1,019	1,982	138
mean AUC (sd AUC)	0.764 (0.081)	0.714(0.09)	0.7 (0.077)	0.726 (0.08)
median AUC	0.762	0.711	0.701	0.721
mean sd across replicates	0.016	0.02	0.017	0.03

Table 3: Overall prediction performance for the different species using BMRF

Predictions were more accurate for yeast, then for chickens, yeast\_ppi and finally, for humans.

The fact that we achieve better prediction performance for chickens than for yeast\_ppi, is most probably linked to the fact that a much lower number of GO terms are predicted when the chicken data was used, than in the other cases. It may be the case that the GO terms that were predicted with chickens data are more easy to predict than the overall set of GO terms that are predicted in yeast, yeast\_ppi and humans.

Therefore, we aimed to compare the prediction performance of the GO terms that were predicted in the 4 cases. We observed that of the 138 GO terms predicted with chicken data, 20, were also predicted in the other 3 cases. Illustration 7 shows, a comparison of the prediction performance using data from the different species:

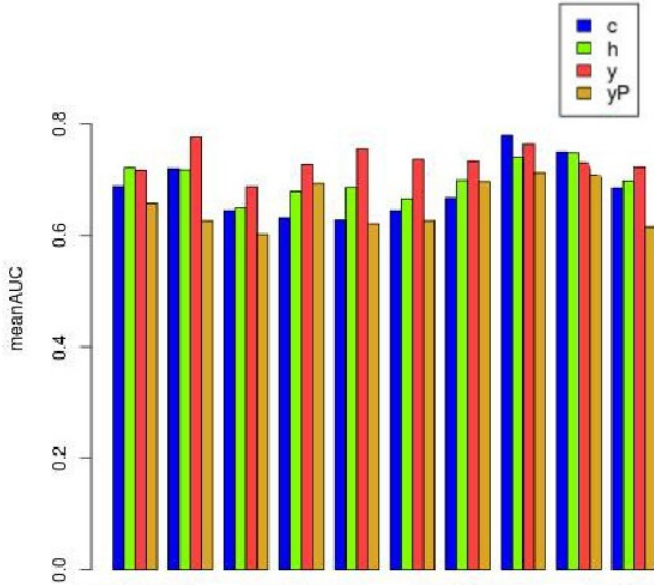


Illustration 5:  
GO\_0009605,GO\_0009653,GO\_0009719,GO\_0009892,GO\_0009966,GO\_0010033,GO\_0010605,GO\_0010628,GO\_0010629,GO\_0030154

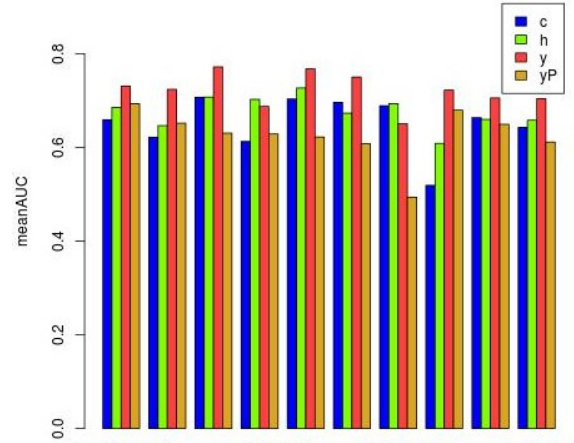


Illustration 6:  
GO\_0031324,GO\_0032879,GO\_0048468,GO\_0048584,GO\_0048646,GO\_0048869,GO\_0050793,GO\_0051128,GO\_0070887,GO\_1901700

Illustration 7 revealed that the prediction performance is generally larger for yeast, the n for humans, the chickens and finally yeast\_ppi. Then mean AUC (and sd) were, 0.729(0.03), 0.688(0.03), 0.668(0.05) and 0.642(0.05), respectively for the 4 species.

The fact that predictions are still high for chickens even though a Pearson correlation of 0.35 was used instead of 0.7 in the other 3 cases, suggests that BMRF has better prediction performance when the network has more edges even at the expenses of a lower reliability of the edges. Furthermore we can conclude that BMRF had a high prediction performance on chickens, which, to our knowledge is the most poorly annotated species ever considered for PFP via networks.

We then aimed to investigate whether the fact that we achieve better predictions in some species and in others is related to the fact that the state of annotation is larger for these.

### Part3) Differences in prediction performance between GO-terms (using BMRF)

With the purpose of identifying which GO-term-properties have a more direct effect on the predictions performance, we computed the correlation between AUC and the different GO-term properties, The results are shown in Table 4.

	yeast	yeast_PPI	humans	chickens
<b>epp_V/tpEppV</b>	0.62	0.373	0.462	.
<b>epp/tpEpp</b>	0.582	0.34	0.378	.
<b>sdAUC</b>	-0.408	.	-0.337	.
<b>teV/tpV</b>	0.394	0.241	-0.13	-0.263
<b>depth</b>	0.265	0.104	.	0.478
<b>spec</b>	0.095	.	.	0.518
<b>#genesV</b>	.	.	.	-0.518

Table 4: Correlations with AUC\_comparison species

Only significant correlations ( $pvalue < 0.05$ ) were added to the table.

From table 4 we learn that epp/tpEpp and sd are the parameters that affect more the prediction performance. Epp/tpEpp shows a favorable correlation with AUC, meaning that for GO terms whose associated genes are interconnected in the network (coexpressed), the method has more chances to distinguish genes that are associated with the GO term from genes that are not.

In order to achieve high PFP accuracy, epp/tpEpp should be as large as possible. Epp depends on data available and cannot be increased with methods, however tpEpp could be reduced by PU-learning. Note that, from theory, PU-learning improves two ratios: (1) it reduces the portion of epn within enn, and (2) it reduces the epn, as some unknown genes are dropped from the analysis. By improving the second ratio, PU-learning would be increasing the epp/tpEpp ratio.

Sd shows a negative correlation with AUC, meaning that for those GO terms whose AUC fluctuates more from replicate to replicate are overall worse predicted. A possible explanation for this is that the sd is high when epp/tpEpp is low. Thus, indirectly, high sd means low overall AUC. This is because if only a few of the associated genes are interconnected with each other, the results will depend on whether these interconnected labeled genes enter the training or the test set in the cross-validation.

In Material and methods it was explained that depth is not a reliable estimate of the specificity of the GO terms. Thus, it is not very surprising that we observed positive correlations between AUC and depth.

Illustration X shows a graphical representation of how these correlations change across the species considered.

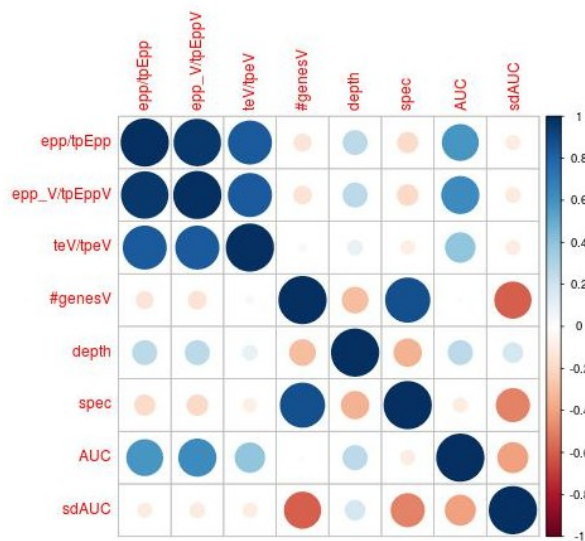


Illustration 8: Correlations yeast

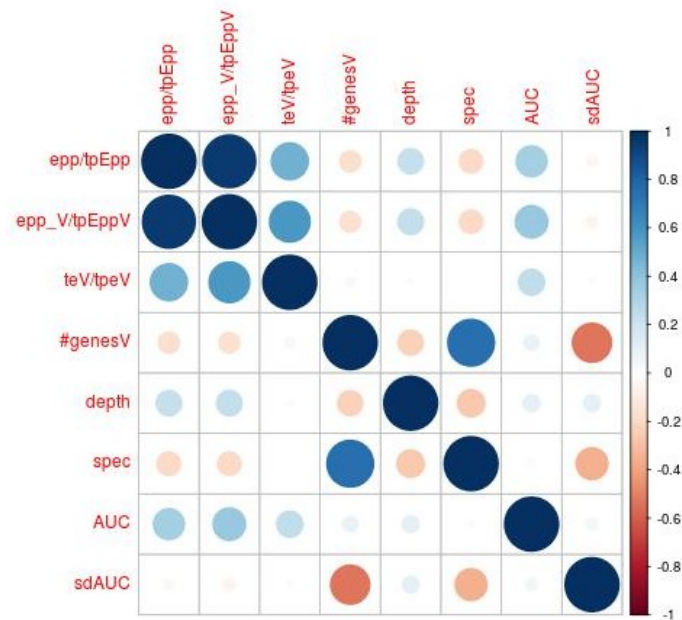


Illustration 9: Correlations chickens

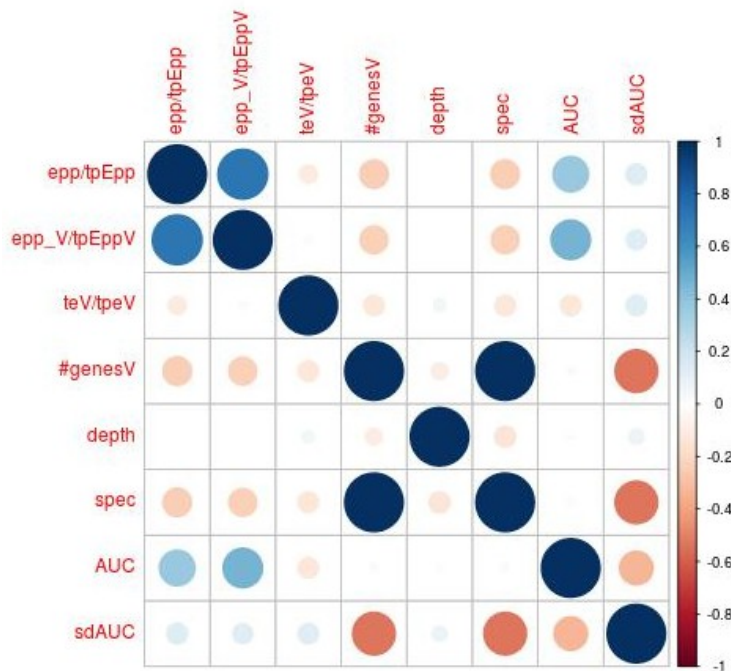
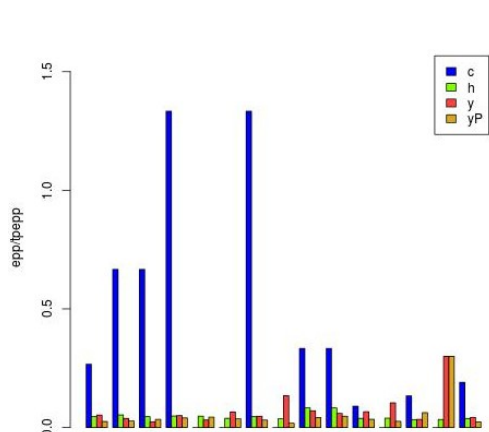


Illustration 10: Correlations humans

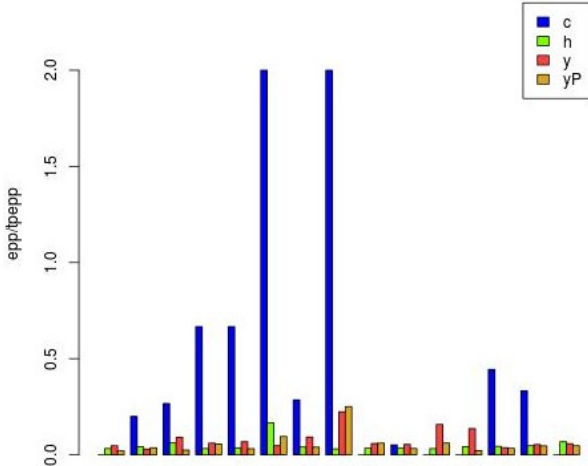


The analysis with chickens data showed a different pattern of correlations. The strong correlation between specificity and AUC is surprising, as well as the correlation between genesV and AUC (-0.52).

Compare the ratio  $epp/tpepp$  at the level of individual GO terms. (Illustration X) for 30 randomly chosen GO terms that are predicted in the four cases

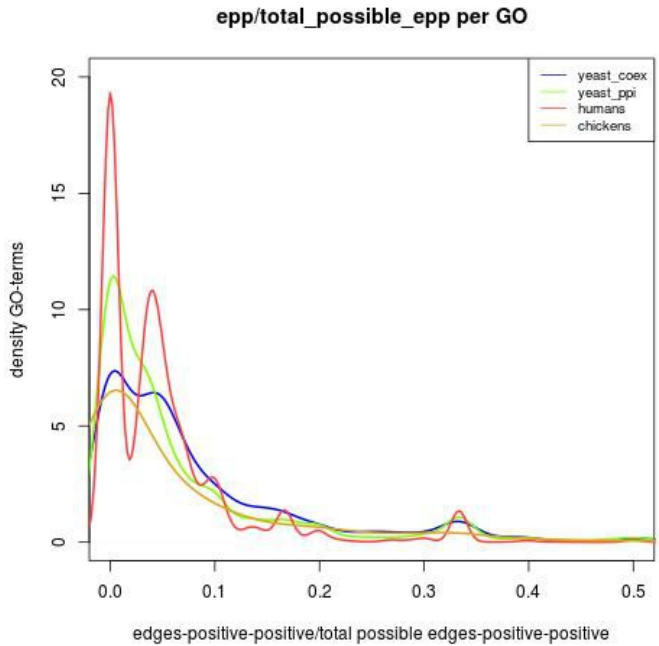


*Illustration 12:*  
GO\_0000003,GO\_0000075,GO\_0000077,GO\_0000086,GO\_0000122,GO\_0000165,GO\_0000278,GO\_0000302,GO\_0000724,GO\_0000725,GO\_0000902,GO\_0001101,GO\_0001558,GO\_0001676,GO\_0001932



*Illustration 11:*  
GO\_0001933,GO\_0001934,GO\_0003006,GO\_0005975,GO\_0006066,GO\_0006071,GO\_0006082,GO\_0006090,GO\_0006109,GO\_0006139,GO\_0006163,GO\_0006164,GO\_0006259,GO\_0006260,GO\_0006261

The ratio is considerably larger for chickens, although this could be an artifact of the scarce data that is available, as explained in Table XX of Appendix II-Data\_overview. After that, the GO terms seem to have a higher ratio for yeast than for the other species, which is i



	Average $epp/tpepp$ (sd)
yeast	0.122 (0.204)
yeast_PPI	0.106 (0.204)
humans	0.066 (0.138)
chickens	0.255 (0.55)

Table 5: average  $epp/tpepp$  for the different species

*Illustration 13:  $epp/tpepp$ . Value sin table x appendix I*

epp/tpepp seem to be the network parameter with a larger impact in the prediction performance.

In part 2 we have seen that prediction sin indiviudal GO terms were higher when human data was used than when chiocken data was used, because human data is

we conlsude that a method the predcition performace with BMRF could be iomproved with PU-learning because it will allow to improv ethe ratioepp/tpepp. More

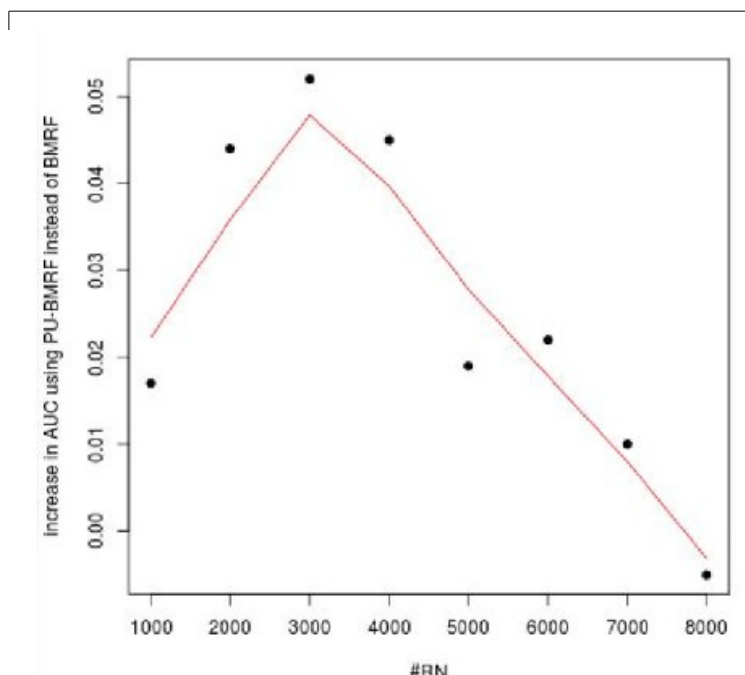
#### **Part 4- PU-BMRF performance evaluation.**

Table 5 illustrates the accuracy of prediction sin the two steps of PU-BMRF, first the RN are extracted from the set of unlabeled genes, and second BMRF is trained on the set of positives and the set of RN.

	Max # of RN extracted							
	1000	2000	3000	4000	5000	6000	7000	8000
<b>Accuracy of extraction of RN</b>	0.981 (0.01)	0.973 (0.015)	0.966 (0.019)	0.965 (0.02)	0.966 (0.021)	0.961 (0.023)	0.959 (0.023)	0.958 (0.024)
<b>PFP AUC using PU-BMRF (sd)</b>	0.723 (0.08)	0.75 (0.072)	0.758 (0.084)	0.751 (0.086)	0.725 (0.094)	0.728 (0.089)	0.716 (0.092)	0.701 (0.095)

*Table 6: Accuracy in the two steps of PU-BMRF, for different choices of max #RN extracted. Results from 30 GO-terms*

The prediction accuracy for these GO terms using BMRF was 0.706 (0.026). Thus, except when a maximum of 8000 RN were extracted, PU-BMRF outperformed BMRF. The maximum improvement (+0.052 AUC) was when a maximum of RN was set to 3000. This is also illustrated in Illustration 7.



*Illustration 14: Increase in AUC PU-BMRF vs BMRF, for different choices of max #RN extracted*

A disappointing results was the increase in standard deviation across replicates (0.03 with BMRF vs 0.076 with PU-BMRF). This increase is surprising considering that, as shown below, the reproducibility in the process of extraction of RN was high. The most likely explanation for this is that the network used was smaller in the case of PU-BMRF. In fact, it was shown that the sd was even higher when a set of 3000 RN was extracted randomly (sd was 0.1).

The average AUC in the process of extraction of 3000 RN for the 30 GO-terms was 0.966 (0.019), and the average sd across folds was (0.0434). However, in order to test the reproducibility in the process of extraction of RN, we calculated which portion of the RN were extracted in the 4 replicates. We carried the analysis in the situation where a maximum of 3000 were chosen. We observed than one average 96.7% (0.0159) of the RN were extracted in the 4 replicates. Results for 30 GO terms are given in table XX in Appendix III-Additional\_results and in illustration XX, part 3. Furthermore, we observed that On average, 3033.583 (104.67) different RNN were extracted in the 1- folds of each replicate, which is very close to the 3000 minimum RN per fold, indicating that the large majority of the RN were common in the 10 folds.

Another important result is that, as expected, the accuracy of prediction was lower when the RN were extracted randomly than when BMRF, or PU-BMRF was used, as it is shown illustration 9

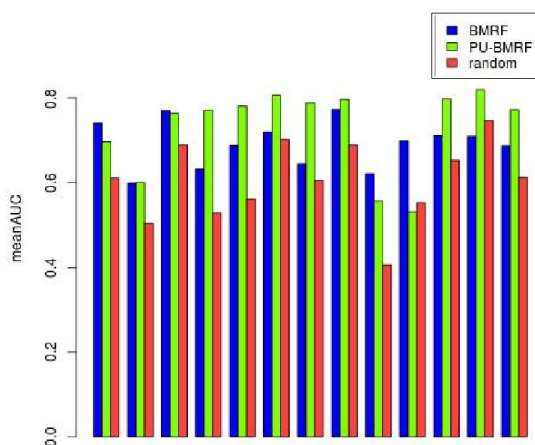


Illustration 15:

GO\_0006928,GO\_0006950,GO\_0007423,  
GO\_0008283,GO\_0009605,GO\_0009653,  
GO\_0009719,GO\_0009887,GO\_0032879,  
GO\_0033554,GO\_0040011,GO\_0044699,  
GO\_0044700

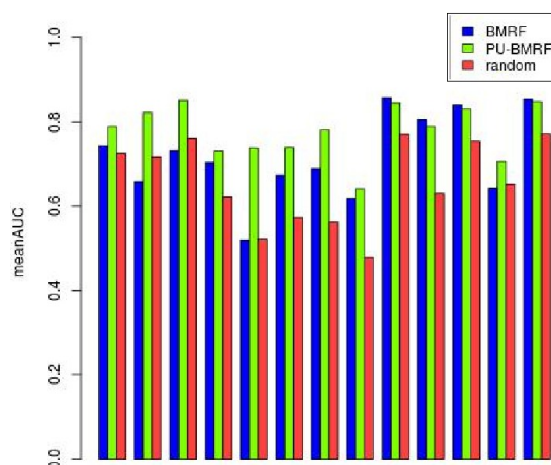
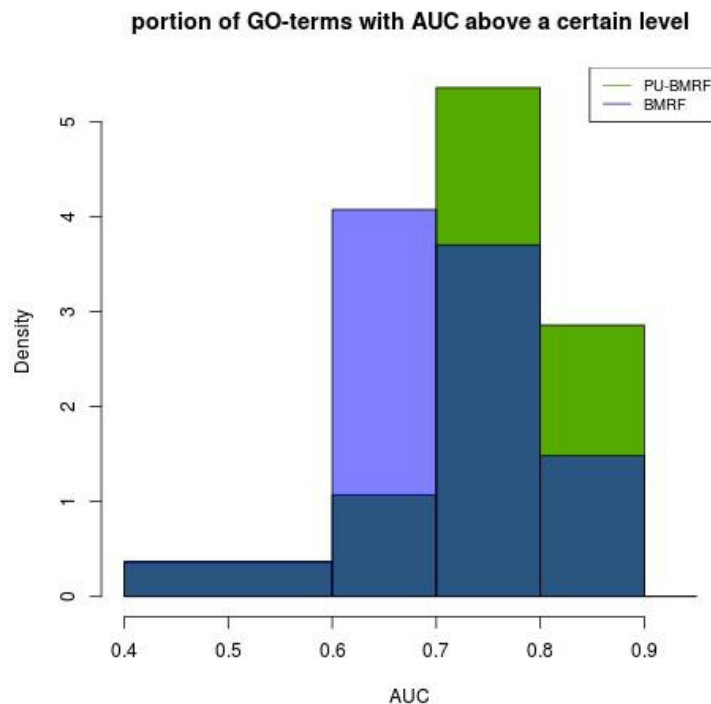


Illustration 16:

GO\_0044710,GO\_0044763,GO\_0044767,  
GO\_0048646,GO\_0051128,GO\_0051240,  
GO\_0060485,GO\_0065008,GO\_0097659,  
GO\_1901360,GO\_1901362,GO\_1901700,  
GO\_1903506

Illustration XX shows to which extend PU-BMRF changed the distribution of AUC across the 30 GO terms used. It is shown that the portion of GO terms for which AUC higher than 70 or 80% increases, whereas there are less GO terms with an AUC of 60-70% when PU-BMRF was used. Further, even when PU-BMRF was used, the same portion of GO terms with AUC<60% remained.



We then investigated whether the increase in AUC with PU-BMRRF vs BMRF was correlated with any of the GO-term properties seen in part 2. We observed that non of the GO-term properties was significantly correlated with the increase in AUC. However, the strongest correlation was with epp/tpepp (0.25), with a p-value of 0.21. We would expect that this p-value will be lower if more than 30 GO terms were used to compute the correlation.

## **Part 5- Biological support of the approach**

Questions to be addressed are:

1. whether the more specific GO terms have a higher degree of connection in the network.
2. whether this is translated into a better accuracy of PFP with BMRF.
3. whether genes that are associated with a large number of GO terms of different levels of specificity are better predicted with BMRF.

We are basically interested in the relationship between 4 variables (it would be good to provide the partial corrs as well as the corrs:

AUC

increase in AUC (PU-BMRF vs BMRF)

specificity

degree of connection

However, there are a few ways to compute degree of connection. And AUC should concern individual genes rather than GO-terms

1. The answer is “yes, to some extent”. It would be good to do other experiment and test for significance. As it can be seen below I tried to link to PU.
2. 3. I will try to do before Saturday. It requires changing the code from “predictions accuracy per GO-term” to “prediction accuracy per gene”

In part 3 we observe that PU contributes more in GO terms whose associated genes show a high degree of connection (high epp/tpepp). Moreover, in part 2 we observed that epp/tpepp shows a negative corr with 'generality' (-0.23 in humans). Thus, PU may contribute more in specific GO terms. Furthermore, since PU contributes more when the ratio epp/tpepp is high, we would expect that PU contributes more in species for which the overall ratio epp/tpepp is high. From table XX in Appendix II, we learnt that the overall ratio epp/tpepp was higher for chickens than for the other more annotated species. Thus, it would be interesting to check whether the ratio epp/tpepp is generally larger for poorly annotated species and if PU contributes more to increase the prediction performance in poorly annotated species than in other more annotated species. Table XX shows the correlation between 'generality' and epp/tpepp.

Genes that are exclusive	corr(spec-epp/tpepp) (p_value)
yeast	-0.192 (0)
yeast_PPI	-0.199 (0)
humans	-0.235(0)
chickens	-0.108 (0.21)

In table XX we observe that the correlation between 'generality' and epp/tpepp ranges between -0.1 and -0.24 depending on the species, and is lower for poorly annotated species. In addition, further analysis revealed that this correlation changes depending on whether we include in the computation of the correlation GO-terms with a minimum number of associated genes (Table XX). This can be due to the fact that epp/tpepp tends to be larger when more genes are considered. The correlation

reaches its highest magnitude when only Go terms with a minim of 2000 associated genes. In other words, we observed that within a group of very general GO terms, the principle that “genes of specific GO terms are highly connected” holds to a larger extend. This makes, sense since when the number of positive genes is very large it is more difficult that they all are co-expressed at the same time. Whereas, when the number of genes is low, the genes within the group have higher chances to be co-expressed with other genes that are not in the group. More so I f we take into account that some of the genes that while being outside the group of positives, are co-expressed with positive genes, are actually positive genes not yest discovered. And we would expect that the chances of making false positives decreases as we look into more general GO terms. Table XX shows the effect of minGOsize on the correlation between generality and epp/tpepp.

<b>Genes that are exclusive</b>	<b>cor(corr(spec-epp/tpepp) – minGOsize)</b>
yeast	-0.1503 (0.006)
yeast_PPI	-0.87 (0)
humans	-0.901 (0)
chickens	-0.402 (0.06)

In order to correct for the size of the group off genes, we carried another experiment. We compared the degree of connection between a group of gens that are only involved in general GO term, and a group of genes that are only involved in general GO terms. We homogenized the size of the two groups and we observed that in out of 10000 replicates of genes sampling for the two groups, in 8394 replicates the genes in the group of specific GO terms appeared to be more interconnected than the genes in the set of general GO terms. On average, these genes appeared to be 6.36% more highly connected than the genes exclusive of the =general GO terms.

## **Discussion**

Development of computational methods for PFP based on network data is a challenging problem in poorly annotated species. Here, we expand upon an existing BMRF and develop a PU implementation (PU-BMRF) that is more accurate than predecessor for PFP in poorly annotated species such as chicken. The efficiency of BMRF to infer function of proteins in poorly annotated had been previously noted [5,6]. Nevertheless, the efficiency of BMRF is hampered because the algorithm attempts to solve a two class classification problem when in fact the annotated data is from one single class (positive class). PU-BMRF tackles this problem by adding a previous step to the BMRF algorithm, in which a set of reliable negative are extracted. Subsequently, the BMRF classifier can be trained with a representative set of genes of each class (negative and positives) and prediction become more accurate.

The basis of PU is in extracting the genes within the set of unlabeled, that show strong differences with respect to the set of positives. It is a must that the features that are used to investigate these differences are as unrelated as possible to the features that the classifier, BMRF in this case, uses afterward. BMRF uses neighborhood information but neglects other important features like, for instance, whether the gene is associated with a related GO term, or the degree of connection between the neighbors of the gene and the genes that are associated with the GO-term of interest. We developed a PU implementation in which 64 features were taken into account to identify a set of reliable negatives (RN) for each GO term.

### **Lack of reproducibility**

Overall, 76% of the GO-terms were more accurately predicted when PU-BMRF was used with respect to when only BMRF was used. PU-BMRF, however, have some limitations. We observed that the extraction of RN has low reproducibility. To our knowledge this lack of consistency in the extraction of RN can be explained by the low number of folds that were used in the cross validation (only two folds). In fact, probably the main drawback of PU-BMRF is that the computational time is around 200 time larger than for BMRF and this may force the user to choose a low value of k-folds in the validation. The increase in computational time was, to some extent, expected given that PU-BMRF requires the computation for 64 features for each gene within each GO term; and some of which are complex like, for instance, the sum of the closeness of a set of gene in the network. As a consequence, one may have to choose a low value of k, at the expenses of a lack of reproducibility in the extraction of RN.

The lack of reproducibility, had been previously noted in PU methods. [13], for instance extracted 20,099 and 4066 RN with two different approaches, respectively and only 589 of the RN were common between the two approaches. [7] explained that in some cases the information from GO alone is not enough to predict a good set of negative examples, because some proteins defy the conventional annotation patterns.

In line with this, it should be considered that the main goal of PU applied to PFP is to eventually predict the function of genes that do not fall in the set of positives or reliable negatives. These predictions, however, cannot



be evaluated with the same method that was used to do the predictions because in most PU approaches, the genes that do not fall in the set of positives or reliable negatives are extracted from the analysis. It could be, thus, argued that the lack of reproducibility is advantageous in that the predictions can be evaluated for a larger number of genes. Since in PU-BMRF, the set of RN that are extracted differs slightly from run to run, we would expect that given a sufficient number of runs, most of the genes will be included in the set of RN at least once. Thus, due to the low reproducibility, we may be able to evaluate the predictions for nearly every gene in the database.

Both, the reproducibility and the number of genes for which the predictions are made, can to some extent, be regulated in PU-BMRF by specifying the number of RN that we want to extract in each run. The user can choose to extract a very large number of RN, although it would come at the cost of a lower increase in accuracy. In [1], for instance, they set a value of 1 as cutoff in equation XX that regulates how many RN will be extracted, whereas in [2] they chose 1.05. In our case, we changed this value according to the GO term as we aimed to extract a fixed number of RN. We decided that this was a reasonable option since we observed that the accuracy did not increase when the number of RN was very large or very small. Another common threshold to separate RN from the unlabeled data is that specificity is equal to one. Thus, the cutoff in equation XX or in any other equation that aims to separate RN from the unlabeled data, could be defined in such a way that none of the known positives would fall in the set of RN.

### **Factors that determine the prediction performance**

Overall, since PU approaches consist on adding previous step to the classification algorithms, one may say that several PU approaches can be considered (either combining or applying one after another), in order to extract a more reliable set of RN. [2], for instance, extracted a set of likely negatives (LN) based on the approach introduced by [1] and then they extracted a set of RN using the approach developed in [3]. [13], for instance extracted RN based in two sets of features, one that consider each features individually and one that considered only the mean of the features of each group of features, and then extracted those RN that were extracted with both approaches. An important aspect to consider is that the quality of the method will not depend so much on the number of features that are defined or on how many PU approaches are integrated, but rather on two main factors: (1) How different the Positives are from the negatives for the features considered (data properties and quality of the features) and (2) Which portion of the unlabeled genes will be extracted as RN. Regarding to the first factor, [7] referred to the so-called “moonlighting” proteins to those proteins that having a unique combination of features cannot be predicted from the conventional annotations.

In addition to this, there is another important factor that determined not only the accuracy with which the RN are extracted, but also the accuracy with which the final predictions are made via the final classifier. This factor is the extent to which the principle of guilt by association holds. It was previously shown that this “guilt by association” heuristic is universal and preserved beyond organism boundaries [12, 13]. However, the same greatly depends on the quality of the data. In particular, in the case of the co-expression networks, it should be considered that even if the quality of the data is good, the phenotypic variation is controlled at many levels, some of which are independent of transcript abundance.

### **Avenues of improvement**

The current method can integrate protein-protein-interaction (ppi) data with the co-expression data, as well as data from some well-annotated related species, like, *Mus musculus* or *Homo sapiens* in the case of chicken, and this could lead to an improvement in accuracy of prediction as explained in [5,6]. The main aspect to be improved regarding PFP in poorly annotated species like chicken, however, is not the accuracy of prediction but the number of GO-terms for which the predictions can be made. This number is very low even when the GO-size filters were at its minimum (minGOsize:9, 321 GO-terms). This problem, in-fact, may be even more severe in PU-BMRF than in BMRF, because since the computational time is much larger for PU-BMRF, the value of k may be lower, and subsequently a larger number of positives cases is required.

The number of GO terms for which make predictions can be made, could be improved, for instance, by expanding upon the existing R function “glmnet” that BMRF uses to train the classifier, in such a way that the matrix are less sparse or that sparser matrix can be solved. In fact, probably the best avenue of improvement for PU-BMRF is to extend the “glmnet” function to more than two classes. This would allow to do predictions within the set

of unlabeled genes, taking simultaneously into account the information from the set of positives and the set of negatives. More importantly, the matrix would become less sparse and predictions could be made for a lower number of positives.

In order to reduce the computational time of PU-BMRF, a first step would be to identify and discard the features that are weakly contributing to extraction of RN. Also, to discard those features whose values change depending on which genes are in the training set. Were these features discarded, the computation of features would need to be computed only once. This could lead to a significant decrease of the running time given that the computation of features accounts for nearly 82% of the running time of PU-BMRF.

Further aspects to be improved are, for instance, a more accurate extraction of RN, for instance, using Self organizing maps as explained in [13], or the Rochio technique [1]; accounting for the magnitude of coexpression; or including an option to make predictions for individual genes rather than for each GO-term.

## **Conclusions**

### PU-BMRF prediction performance

1. PU-BMRF outperformed BMRF using chicken data: 0.758 vs 0.706 AUC. However, the standard deviation across replicates was more than double (0.078 vs 0.026)
  - The best performance of PU-BMRF was when a maximum of 3000 RN were extracted
  - The reproducibility in the process of extraction of RN was high, with an average of 96.7% of the RN extracted being extracted in the four replicates considered.
  - When the RN were extracted randomly, PU-BMRF achieved a low prediction performance (0.652 using chicken data), lower than when the conventional BMRF was used (0.706).
  - The increase in accuracy when PU-BMRF was used instead of BMRF was higher for the most specific GO terms (0.126, pvalue=0.53\*)

### BMRF and PU-BMRF for poorly annotated species

- Both, the conventional BMRF and PU-BMRF achieved a good performance on chickens. However, predictions can only be made by a maximum of 364 GO terms due to the GO-size-filters that BMRF uses.
- the most important limitation of BMRF and PU-BMRF for poorly annotated species is that a minimum of genes annotated with the GO terms is required (minimum 8) genes.
- The increase in accuracy when PU-BMRF was used instead of BMRF was higher for those GO terms with a high epp/tpepp ratio (correlation of 0.22, pvalue 0.21\*). We observed that the ratio epp/tpepp for a given GO term was considerably higher in the chicken data, therefore we would expect that PU-learning is particularly effective in poorly annotated species where the proportion of unlabeled genes is large.

### Data properties that determine the prediction performance

- In the trade-off between the extension of the network and the quality of it, BMRF seemed to perform better when the co-expression threshold was set to 0.35.
- The GO term property most strongly correlated with the prediction performance was epp/tpepp
- Decreasing the number of edges in the network seemed to improve the prediction performance considerably
- The size of the network did not seem to have a strong impact on the accuracy of prediction

### Biology behind co-expression networks

- The characteristics of the co-expression analysis (tissue, experimental conditions...) had a very low impact on the prediction performance.
- The genes from the most specific GO terms seemed to be more interconnected, with a correlation between epp/tpepp of 0.235 in humans. Furthermore these correlations increased as we stick to the most general GO terms.
- Genes that are exclusive from specific GO terms seem to have a higher degree of connection than genes that are exclusive from general GO terms

### Others

- An in-depth overview of the differences in data between the species considered, as well as the current state of annotation is provided in Appendix II-Data-overview.

\* This correlation was estimated with only 30 GO terms. This could explain the high p\_value.



## **References**

- 1.KOURMPETIS, Y. A. I., VAN DIJK, A. D. J., BINK, M., VAN HAM, R. & TER BRAAK, C. J. F. 2010. Bayesian Markov Random Field Analysis for Protein Function Prediction Based on Network Data. *Plos One*, 5.
- 2.SCHWIKOWSKI, B., UETZ, P. AND FIELDS, S.. 2000. A network of protein-protein interactions in yeast. *Nature Biotechnology* 18:1257 – 1261.
- 3.MINGHUA DENG, KUI ZHANG, SHIPRA MEHTA, TING CHEN & SUN., F. 2004. Prediction of Protein Function Using Protein-Protein Interaction Data. *Journal of Computational Biology*, 10, 6.
- 4.YAN, K. K., WANG, D. F., ROZOWSKY, J., ZHENG, H., CHENG, C. & GERSTEIN, M. 2014. OrthoClust: an orthology-based network framework for clustering data across multiple species. *Genome Biology*, 15.
- 5.JOACHIM W. BARGSTENA, EDOUARD I. SEVERINGB, J.-P. N., GABINO F. SANCHEZ-PEREZA & DIJK, A. D. J. V. 2014. Biological process annotation of proteins across the plant kingdom. *Current Plant Biology*, 1.
- 6.Consortium TF, Andersson L, Archibald AL, et al. Coordinated international action to accelerate genome-to-phenome with FAANG , the Functional Annotation of Animal Genomes project. *Genome Biol.* 2015:4-9. doi:10.1186/s13059-015-0622-4.
- 7.BHARDWAJ, N., GERSTEIN, M. & LU, H. 2010. Genome-wide sequence-based prediction of peripheral proteins using a novel semi-supervised learning technique. *Bmc Bioinformatics*, 11.
- 8.YANG, S. K. 2012. Positive-unlabeled learning for disease gene identification. *Bioinformatics*, 28, 2640-2647.
- 9.YOUNGS, N., PENFOLD-BROWN, D., BONNEAU, R. & SHASHA, D. 2014. Negative Example Selection for Protein Function Prediction: The NoGO Database. *Plos Computational Biology*, 10.10.YANG, P., LI, X. L., MEI, J. P., KWOH, C. K. & 10.JIANG, M. & CAO, J. Z. 2016. Positive-Unlabeled Learning for Pupylation Sites Prediction. *Biomed Research International*.
- 11.XIAOLI, LI & BING, LIU. 2003. *IJCAI'03 Proceedings of the 18th international joint conference on Artificial intelligence*. 587-592.
- [12]<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2164-14-13>
- [13]<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1546-7>
16. Sharan R, Ulitsky I, Shamir R (2007) Network-based prediction of protein function. *Molecular Systems Biology* 3: 1–13.R. SharanI. UlitskyR. Shamir2007Network-based prediction of protein function.*Molecular Systems Biology*3113
19. <https://www.frontiersin.org/articles/10.3389/fpls.2016.00444/full>
20. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3584181/f>
21. <http://msb.embopress.org/content/3/1/88.long>