

1 Many thanks to all reviewers for their useful comments, which will improve the final version of the paper. We are very  
 2 happy that the novelty, simplicity, effectiveness, significance, and generality of our method has been appreciated. We  
 3 will address all the comments about the related work, the tables, and the details of the datasets in the final paper.

4 **@R1: Motivation of loss function?** **@R2: Slightly different toy problems?** **@R3: Using max has already been**  
 5 **proposed by other methods (e.g., DeepGO).** **@R3: Why MCLoss is better than BCELoss ( $\mathcal{L}$ )?** While the max  
 6 function has already been used, nobody so far has shown how to deploy it *effectively*, which is what the com-  
 7 bination of MCM and MCLoss does. Experimentally, the ablation studies in Table 4 show that MCLoss allows  
 8 for *consistent* higher performance than  $\mathcal{L}$ . In general, consider a class  $A$  with ancestors  $A_1 \dots A_n$  in the hierar-  
 9 chy. The higher  $n$  the more likely it is that a neural network (NN) with MCM trained with  $\mathcal{L}$  will remain stuck  
 10 in bad local optima. Indeed, consider a datapoint such that  $h_A > h_{A_1} \dots h_{A_n}$ ,  $y_A = 0$ ,  $y_{A_1} \dots y_{A_n} = 1$ : then  
 11  $\mathcal{L} = -\ln(1 - h_A) - \sum_{i=1}^n \ln(h_{A_i}) = -\ln(1 - h_A) - n \ln(h_{A_i})$ , while  $\text{MCLoss} = -\ln(1 - h_A) + \sum_{i=1}^n \text{MCLoss}_{A_i}$ .  
 12 Notice that, in the GO datasets, it is common to have  $n > 10$ , given that the hierarchies

13 have 13 levels, and each class can have more than one par-  
 14 ent. To visualize the negative impact of using  $\mathcal{L}$  instead  
 15 of MCLoss, consider the leftmost figure with 9 rectan-  
 16 gles named 1 . . . 9. Further, assume (i) we have classes  
 17  $A_1 \dots A_9$ , (ii) that a datapoint belongs to  $A_i$  if it belongs  
 18 to the  $i$ th rectangle, and (iii) that  $A_5$  (resp.,  $A_3$ ) is an an-  
 19 cestor (resp., descendant) of every class. Thus, all points  
 20 in rectangle 3 belong to all classes, and if a datapoint  
 21 belongs to a rectangle, then it also belongs to class  $A_5$ .

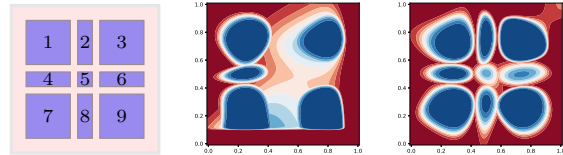


Figure 1: From left to right: (i) rectangles disposition, (ii)  
 decision boundaries for  $A_5$  of  $h + \text{MCM}$  trained with  $\mathcal{L}$ ,  
 and (iii) decision boundaries for  $A_5$  of  $\text{C-HMCNN}(h)$ .

22 Let  $h$  be a NN with a single hidden layer with 7 neurons. Then, the average  $AU(\overline{PRC})$  (and std) over 10 runs for  
 23  $h + \text{MCM}$  trained with  $\mathcal{L}$  is 0.938 (0.038), while  $h + \text{MCM}$  trained with MCLoss ( $\text{C-HMCNN}(h)$ ) is 0.974 (0.007).  
 24 Notice that not only  $h + \text{MCM}$  performs worse, but also, due to the convergence to bad local optima, the std obtained  
 25 with  $h + \text{MCM}$  is 5 times higher than the one of  $\text{C-HMCNN}(h)$ : the (min, median, max)  $AU(\overline{PRC})$  for  $h + \text{MCM}$  are  
 26 (0.871, 0.945, 0.990) while for  $\text{C-HMCNN}(h)$  are (0.964, 0.975, 0.990). The figure shows the decision boundaries of  
 27 the 6th best performing networks for class  $A_5$  (given the even number of runs, we could not take the exact median). We  
 28 will include this example and a general discussion in the paper. Thanks also to **R2** for the suggestion on how to create  
 29 alternative synthetic examples: we will explore such direction in the future.

30 **@R2: The need for this method is not fully demonstrated.** **@R2 It is hard to believe that NNs can outperform**  
 31 **more classical methods on such datasets.** As pointed out by **R1**, the paper deals with a problem with a lot of  
 32 significance and relevance to real-world applications, while we achieved SOTA results on most benchmark datasets used  
 33 in the literature of HMC. As shown by [2], who established the current SOTA on these datasets, NNs can beat more  
 34 classical methods. In the paper, we compared ourselves with the current SOTA models (one is Clus-Ens: an ensemble  
 35 of predictive clustering trees), which have already proved to beat famous tree-based and kernel-based methods.

36 **@R2: What is the added value of using a NN here?** NNs have shown the ability of performing very well in different  
 37 scenarios. We wanted to create a model that was broadly applicable and easy to use, and hence NNs were our first choice.  
 38 Notice though that  $h$  can be any model that outputs a probability for each class and can be trained with backpropagation.

39 **@R2: Is there a danger to get stuck in bad local optima?** Our results show our model is very stable. As stated at  
 40 line 205 in the paper, the std over 10 runs is very small (in the range  $[0.5 \times 10^{-3}, 2.6 \times 10^{-3}]$ ). This is not surprising:  
 41 while, in theory, gradient descent can be performed only on differentiable functions, in practice, it performs well also on  
 42 non-differentiable functions (e.g., on  $\text{ReLU}()$  and  $\text{max}()$ ) leading to very good performances [1]. Considering the  $\text{max}()$   
 43 function, it is differentiable almost everywhere, and software implementations of NN training (e.g., Adam in Pytorch)  
 44 usually return one of the one-sided derivatives. This may be heuristically justified by observing that gradient-based  
 45 optimization on a digital computer is subject to numerical error anyway [1]. Hence, in practice, the non-differentiability  
 46 of a small set of points does not affect the learning algorithm.

47 **@R3: Comparative analysis of computation time?** In addition to the inference time per batch in Table 4 (appendix),  
 48 we will include a table with the training times of each model. Due to lack of space, here, we report the ones for  
 49  $\text{CELLCYCLE FUN}$  (in minutes, averaged over 10 runs): (i)  $\text{C-HMCNN}(h)$ :  $\sim 6\text{m}$ , (ii) Clus-Ens:  $\sim 20\text{m}$ , and (iii)  
 50  $\text{HMC-LMLP}$ :  $\sim 51\text{m}$ . Since we do not have the code, we could not measure the times for  $\text{HMCN-R}$  and  $\text{HMCN-F}$ .

51 **@R4, R2 Plug in our method on top of [2] and/or CNNs?** We still have to try this; we will do so in future work.

52 **@R4: Did you optimize them over the validation sets?** Yes, we optimized the hyperparameters on the validation  
 53 sets. We used the test sets only to report the results. We will make it clearer in the appendix.

54 **@R5: The results are different between tables, I would expect them to be the same.** We conducted the ablation  
 55 studies on the validation set (see caption of Table 3), while we report the results (after re-training on training+validation  
 56 set) on the test set. We will specify this in the caption of Table 2.

57 [1] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*, pages 188–189. MIT Press, 2016. <http://www.deeplearningbook.org>.  
 58 [2] J. Wehrmann, R. Cerri, and R. C. Barros. Hierarchical multi-label classification networks. In *Proc. of ICML*, 2018.