

GLOLLOC: MIXTURE OF GLOBAL AND LOCAL EXPERTS FOR MOLECULAR ACTIVITY PREDICTION

Hélène A. Gaspar

BenevolentAI
4-8 Maple St, Bloomsbury
London W1T 5HD
helena.gaspar@benevolent.ai

Matthew P. Seddon

BenevolentAI
4-8 Maple St, Bloomsbury
London W1T 5HD

ABSTRACT

Quantitative structure-activity relationships (QSAR) models have been used for decades to predict the activity of small molecules, using encodings of the molecular structure, for which simple 2D descriptors of the molecular graph are still most commonly used. One of the recurrent problems of QSAR is that relationships observed for a specific scaffold (pruned molecular skeleton) are often not transferable to another; this is often addressed by building several local models from subsets of the chemical space. Similarly, single task models sometimes outperform large multi-task models in predicting the activity of small molecules against specific proteins. In this paper, we introduce Glolloc, a global-local MoE-QSAR architecture, based on a Mixture of Experts (MoE) framework. Glolloc combines predictions from global and local experts, provides a built-in model introspection tool, can enhance model performance, and removes the need to maintain several local models.

1 INTRODUCTION

Mixtures of experts (MoE) are a subset of ensemble learning algorithms, consisting in individual learners or experts, and a gate that assigns weights to different learners. They were relatively popular in the 90s (Jordan & Jacobs, 1993; Jacobs et al., 1991), and have recently regained interest in the deep learning community, notably for language and vision (Ruiz et al., 2021; Shazeer et al., 2017). MoE can be seen as a type of meta-learning approach as it "learns to learn" from an ensemble of experts, which can be explicitly pre-defined e.g. by using clustering or subsets of input features (Guttag et al., 2000; Tang et al., 2002).

In quantitative structure-activity relationships (QSAR) models, where the aim is to predict activity from molecular structures (e.g. activity against a protein of interest), local models trained on subsets of molecules often provide better performances than global models trained on the whole dataset (Yuan et al., 2006). In drug discovery projects, a local model is often constructed for each molecular series of interest (set of molecules with a similar skeleton or "scaffold"). At a given stage, several chemical series might be investigated, and several models are separately optimised and maintained. However, these local models discard potentially valuable information, and do not always perform better than global models. Furthermore, some molecules share characteristics from different series and would benefit from combining several local model predictions. Methods that are able to combine advantages from both local and global models should therefore be of particular interest in the field.

We thought a MoE architecture could be a way to combine predictions from global and local experts, each local expert being specialised in a chemical series of interest. Local experts for QSAR can also be defined in other ways. When the goal is to predict the activity of a small molecule against several proteins at the same time, one local expert can be assigned to each protein task. Multi-protein models can perform better or worse than single task learning, depending on the protein (Gaspar et al., 2021) - subdividing the problem and training one expert per protein or protein family, in addition to a global expert trained on everything, could help to better adapt predictions to specific proteins.

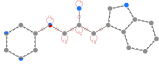
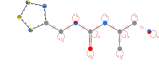
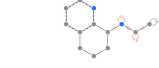
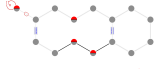
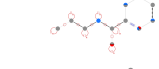
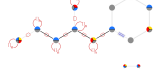
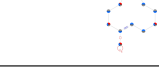
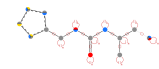
The MoE literature for QSAR is very scarce - MoEs were recently investigated (Dörgő et al., 2020) as a QSAR ensemble method across different feature spaces. In methods introduced in the following sections, the problem is naively partitioned by input instances rather than input features, although combining both approaches would be an interesting follow-up. Partitioning the input space based on chemistry expert knowledge, and then using a MoE architecture to choose between local and global experts has not yet been attempted. The main aim of the following experiments is to explore different ways to partition the problem in a chemically relevant way, and come up with a model able to focus on local aspects whilst not discarding valuable information.

The main contribution of this paper is a global-local ("GloLoc") MoE-QSAR architecture for single or multi-task learning that: (a) combines advantages from both local and global models, which are routinely used separately in QSAR pipelines, (b) is structured using expert knowledge (e.g. medicinal chemist partitioning the space by chemical series), (c) removes the need to maintain many local models, (d) determines which experts are more appropriate for a given molecule, (e) provides a built-in model introspection tool via weights given to experts, and (f) can improve performance over a single global model or an ensemble of random experts.

2 METHODS

2.1 DATA

Table 1: ChEMBL datasets and SMARTS pattern defining chemical compounds in validation and test set. The images were generated using SMARTSplus (Ehrt et al., 2020)

abbreviation	protein name(s)	train:valid:test	valid/test chemical series
<i>single task modelling</i>			
akt1	RAC-alpha serine/threonine-protein kinase	642:35:35	
cp3a4	Cytochrome P450 3A4	1316:35:35	
cxcr4	C-X-C chemokine receptor type 4	145:20:20	
gcr	Glucocorticoid receptor	1118:70:70	
kif11	Kinesin-like protein KIF11	531:21:20	
protease	HIV-1 protease	1753:281:281	
revtrans	HIV-1 reverse transcriptase	1266:104:104	
<i>multi-task modelling</i>			
cyp	CP3A4, CP2D6, CP1A2, CP2C9, CP2CJ	2915:568:35	

In the single task experiments, the 7 datasets (Table 1) correspond to 7 diverse proteins from the diverse subset of the DUD-E data collection (Mysinger et al., 2012). Note that whilst the diverse subset from DUD-E was used to choose the proteins, the ligand data was retrieved from a publicly available benchmark dataset (Lenselink et al., 2017b) curated by Lenselink et al. (Lenselink et al., 2017a), which compiles activities of small molecules extracted from the ChEMBL database (Gaulton et al., 2016). For our purposes, instances are small molecules, and the endpoint is pChEMBL activity (-Log(molar IC50, XC50, EC50, AC50, Ki, Kd or Potency)). For each of the 7 datasets, a tree-based "autosmarts" algorithm (Algorithm 1) was used on molecules with pChEMBL value > 5.5 to find

a top chemical series of interest - its corresponding SMARTS pattern visualisation is reported in Table 1.

Compounds matching that pattern were split into train, valid and test using time splits. This train set was augmented with data not matching the pattern, resulting in a train set that contains both the local data and other diverse molecules. The time splits were defined using publication year, to mimic what happens in a drug discovery project, where the aim is to do better on more recent data. It should be noted that publication year is not an ideal proxy and does not reflect the actual experiment date.

The cp3a4 dataset was supplemented by other cytochrome data for the multi-task challenge (CP2D6, CP1A2, CP2C9, and CP2CJ) to construct the cyp dataset - the test dataset being the same as for cp3a4.

ECFP4_2048 descriptors were used to featurise molecules: Morgan fingerprints with a radius of 2 and size 2048.

2.2 PROPOSED SINGLE TASK AND MULTI-TASK NETWORKS

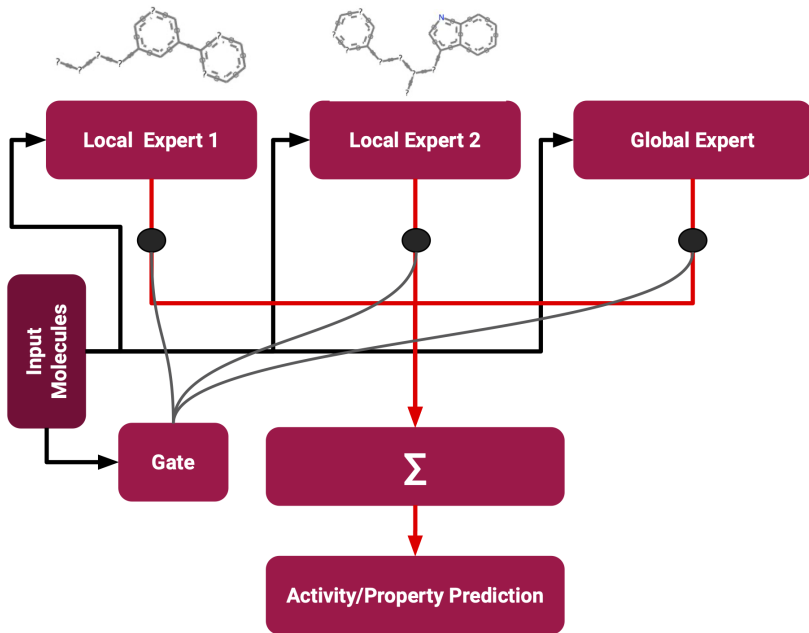


Figure 1: Mixture of Experts (MoE) for molecular activity modelling: the gate assigns weights to predictions from a global expert and individual experts specialised in specific chemical series

Two architectures were investigated: a single-task network with local experts specialising in different regions of the chemical space (Figure 1), and a multi-task network to predict the activity of multiple proteins, where local experts specialise in predicting specific proteins (Figure 2). In each case, in addition to the local networks (chemical series or protein experts), a global expert is trained on the whole data set, and is expected to perform better for molecules outside of the applicability domain of the local models.

The gate network, trained at the same time as the experts, assigns weights to each of them. The final predicted activity for a molecule n and task t is the sum of outputs O from a set of experts E weighted by gate weights G (Equation 1).

$$\hat{Y}_{nt} = \sum_{i=1}^{|E|} G_{nti} O_{nti} \quad (1)$$

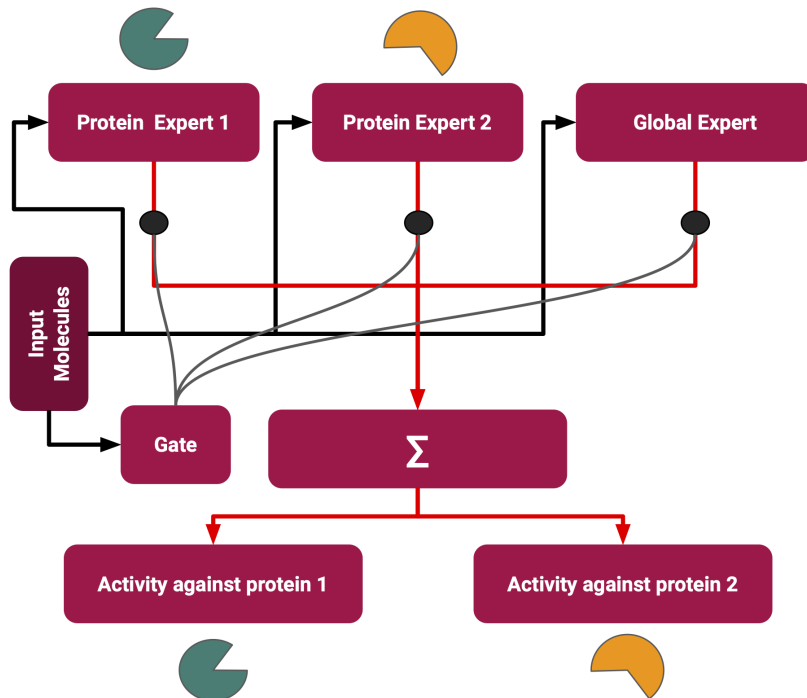


Figure 2: Mixture of Experts (MoE) for multi-task ligand-protein modelling: the gate assigns weights to predictions from a global expert and individual protein-specific experts

Experts and gate are trained at the same time, although each expert is kept independent by using expert-specific optimisers and losses (cf. Appendix). The gate is also assigned its own optimiser, but is made "aware" of expert predictions in the gate-specific loss (Equation 2), which is used to update the gate parameters.

$$loss_{gate} = \frac{1}{NT} \sum_n \sum_t (Y_{nt} - \hat{Y}_{nt})^2 \quad (2)$$

2.3 AUTOMATED METHODS TO DEFINE CHEMICAL SERIES

In a drug discovery project, chemical series are typically identified by medicinal chemists. Chemoinformaticians encode them as SMARTS - a type of regex language, which allows to query molecules containing specific substructures.

However, for benchmarks on public data sets, an automated way to partition the chemical space is needed. Note that in a real-world scenario, project chemists would help in the definition of the series to add their domain knowledge to the model structure. Glolloc implements 3 methods to generate chemical series experts from a chemical dataset: cluster-based, tree-based, and neighbourhood-based.

The cluster-based method is self-explanatory: a k-means clusterer using scikit-learn's implementation (Pedregosa et al., 2011) is trained on the training set, and cluster labels are used to construct cluster-specific experts.

The tree-based and neighbourhood-based methods both rely on finding fuzzy maximum substructures (MCS) in groups of molecules. They are described in further detail in the Appendix. We used the tree-based method in this paper to define our chemical series, as we obtained similar results with the neighbourhood-based method. The advantage of MCS-based methods is the automatically generated SMARTS patterns that can be visualised and interpreted easily - this type of approach is therefore coined "autosmarts" in this paper.

Table 2: (a) Mean squared error (MSE) and (b) Spearman’s rank correlation on single test set (using best model determined from validation set) for a single global expert, and different global + local expert setups: global expert with 6 random experts (global+6xrandom), 6 clusters (global+6clusters) or an automatically determined number of chemical series (global+autosmarts)

(a) MSE	global	global+6xrandom	global+6xclusters	global+autosmarts
akt1	1.048	1.036	0.985	1.109
cp3a4	0.142	0.107	0.117	0.118
cxcr4	0.148	0.145	0.153	0.148
gcr	1.442	1.472	1.348	1.264
kif11	0.533	0.603	0.534	0.530
protease	1.025	1.002	0.980	0.894
revtrans	0.833	0.854	0.788	0.885
(b) Spearman’s ρ	global	global+6xrandom	global+6xclusters	global+autosmarts
akt1	0.634	0.624	0.688	0.603
cp3a4	0.624	0.657	0.692	0.630
cxcr4	0.208	0.150	0.338	0.190
gcr	0.320	-0.147	0.349	0.361
kif11	0.703	0.615	0.650	0.664
protease	0.440	0.443	0.457	0.499
revtrans	0.281	0.253	0.353	0.213

3 RESULTS

3.1 SINGLE TASK: PARTITION PROBLEM BY CHEMICAL SERIES

This section compares Random Forests and Glocloc architectures using different partitions of the chemical space to define experts (Figure 1).

Two control experiments were designed for the Glocloc framework: a model using a single global expert, and a model using one global expert and 6 random experts. The random experts divide the chemical space into 6 random sections. The control experiments are run against two structured Glocloc models: one with the global expert and 6 cluster experts, the other with the global expert and autosmarts (Algorithm 1) experts. The more structured experts (based on clustering or autosmarts algorithms) yield lower mean squared errors (MSE) for 5 out of 7 datasets, and higher Spearman rank correlation in 6 out of 7 (Table 2).

When comparing the best glocloc expert combination (whether random or structured) against a single expert or Random Forests (Figure 3), Random Forests wins only for the gcr dataset and has the worst performance in all other cases.

3.2 MULTI-TASK: PARTITION PROBLEM BY PROTEIN

Two control experiments were run on the cyp dataset in a multi-task setting (Figure 2), with the goal of predicting CP3A4 activity in the test set. The first control consists in a single task model with only one global expert, the second in a multi-task model with also only one global expert. These controls are compared to a multi-task model with both a global expert and 5 protein experts, one for each of the proteins in the dataset (CP2D6, CP1A2, CP2C9, CP2CJ, CP3A4). Results are reported in Figure 4.

Predictions from the single task model are generally higher than from the multi-task model, and lower-activity molecules are not well predicted. The multi-task baseline pushes those predictions down but the range of predicted values stays narrow; adding protein experts improves the MSE in this particular instance.

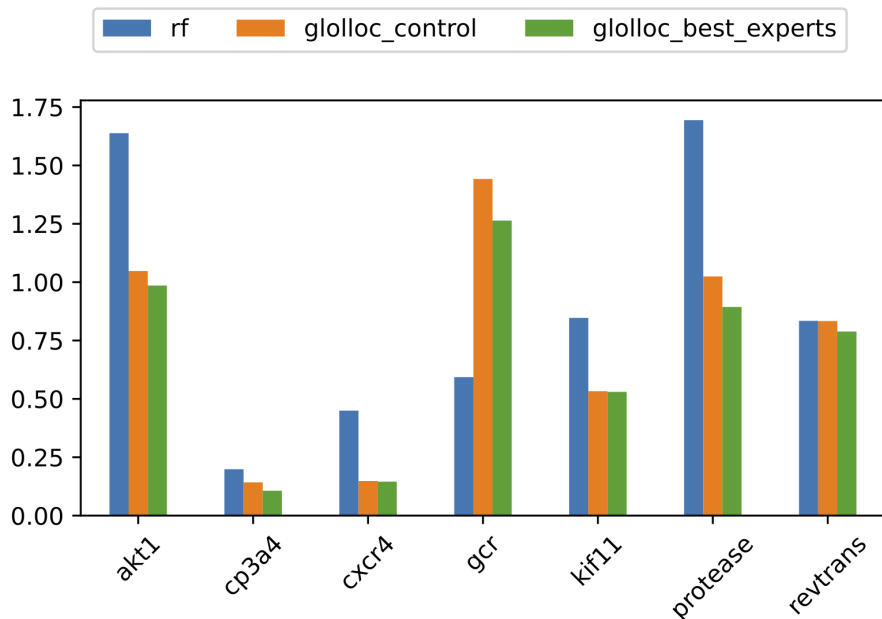


Figure 3: Test set mean squared error (MSE) for Random Forests (rf), a single global expert (gloalloc_control) and best combination of experts (gloalloc_best_experts).

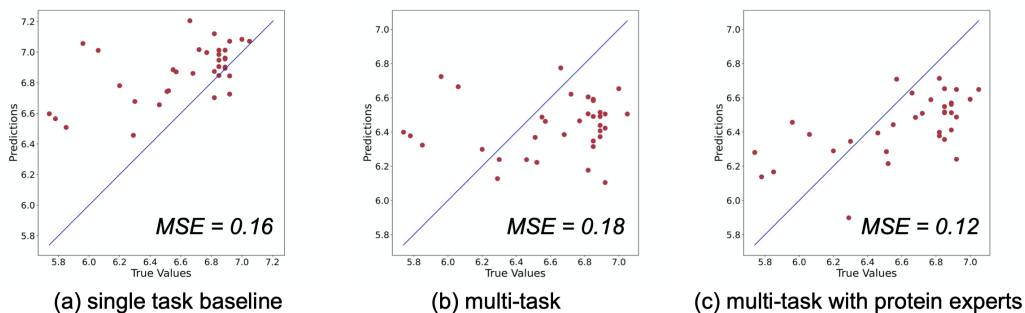


Figure 4: Test set mean squared error (MSE) on the cyp test set for (a) a single task global expert, (b) multi-task global expert and (c) multi-task with both global expert and 5 protein experts - one for each of the proteins in the cyp training set.

3.3 MODEL INTROSPECTION

An important advantage of partitioning a QSAR problem is downstream model introspection. In Figure 5, the weights assigned to different local experts are shown for two molecules in the akt1 validation set. If the model was trained correctly, most of the gate weights should go to experts specialised in the relevant chemical series and the global expert. Molecules belonging to chemical series which are not covered by local experts are expected to be mostly assigned to the global expert.

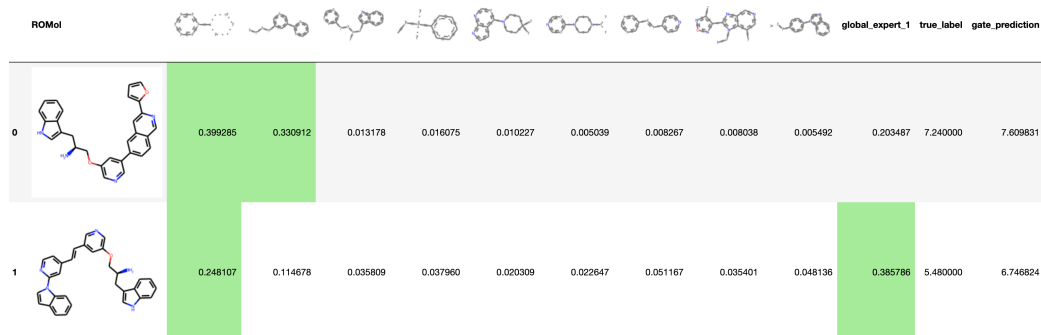


Figure 5: Gate weights visualisation for two molecules in the akt1 validation set: columns with SMARTS patterns symbolise local experts specialised in chemical series. For each molecule, the two experts with the largest weight are highlighted.

The same type of visualization can be produced for the multi-task Glolloc architecture (Figure 2): in Figure 6, weights assigned to different protein experts are shown for two molecules in the cyp validation set. The network predicts the activity against 5 proteins (CP3A4, CP2D6, CP1A2, CP2C9, CP2CJ), but only weights for the CP3A4 task are shown in the figure. For that task, most predictions are dominated by the CP3A4 expert and the global expert.



Figure 6: Gate weights visualisation for two molecules in the multi-task cytochrome dataset, querying predictions for CP3A4 activity. CP2D6, CP1A2, CP2C9, CP2CJ, and CP3A4 are protein-specific experts. For each molecule, the two experts with the highest gate weight are highlighted.

4 CONCLUSION AND NEXT STEPS

Glolloc adds local structure to QSAR models, which can be used for downstream model introspection by visualising how local or global structures contribute to the final prediction.

In the preliminary investigations shown in this paper, models adding more structure, based on chemical series, clusters, or single proteins, can often perform better than controls (RF/single experts/random experts). However, these models need more careful optimisation, and more datasets and splits should be investigated in follow-up experiments to build more confidence in performance improvement.

There is still progress to be made on the side of expert training, as experts are defined using data masks and users have to monitor a series of losses. Arguably, the discrete, masked multi-expert setting could be replaced by a more continuous approach. Alternatively, an attention mechanism could be used instead of experts, to give more attention to predefined dataset structures.

Different 2D and 3D depictions of molecules could also be used to define additional experts, as well as potential groupings of related tasks, such as protein families in the multi-task setting. Glolloc-like architectures could also be used to model different endpoint types or experimental settings.

REFERENCES

- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001. doi: 10.1023/a:1010933404324. URL <https://doi.org/10.1023/a:1010933404324>.
- Gyula Dörgő, Omar Péter Hamadi, Tamás Varga, and János Abonyi. Mixtures of QSAR models: Learning application domains of pka predictors. *Journal of Chemometrics*, 34(4), April 2020. doi: 10.1002/cem.3223. URL <https://doi.org/10.1002/cem.3223>.
- Christiane Ehrt, Bennet Krause, Robert Schmidt, Emanuel S. R. Ehmki, and Matthias Rarey. SMARTS.plus – a toolbox for chemical pattern design. *Molecular Informatics*, 39(12):2000216, October 2020. doi: 10.1002/minf.202000216. URL <https://doi.org/10.1002/minf.202000216>.
- Hélène A. Gaspar, Mohamed Ahmed, Thomas Edlich, Benedek Fabian, Zsolt Varszegi, Marwin Segler, Joshua Meyers, and Marco Fiscato. Proteochemometric models using multiple sequence alignments and a subword segmented masked language model. *ChemRxiv*, May 2021. doi: 10.26434/chemrxiv.14604720.v1. URL <https://doi.org/10.26434/chemrxiv.14604720.v1>.
- Anna Gaulton, Anne Hersey, Michał Nowotka, A. Patrícia Bento, Jon Chambers, David Mendez, Prudence Mutowo, Francis Atkinson, Louisa J. Bellis, Elena Cibrián-Uhalte, Mark Davies, Nathan Dedman, Anneli Karlsson, María Paula Magariños, John P. Overington, George Papadatos, Ines Smit, and Andrew R. Leach. The ChEMBL database in 2017. *Nucleic Acids Research*, 45(D1):D945–D954, November 2016. doi: 10.1093/nar/gkw1074. URL <https://doi.org/10.1093/nar/gkw1074>.
- Srinivas Gutta, Jeffrey R. J. Huang, P. Jonathon Phillips, and Harry Wechsler. Mixture of experts for classification of gender, ethnic origin, and pose of human faces. *IEEE Transactions on Neural Networks*, 11(4):948–960, July 2000. doi: 10.1109/72.857774. URL <https://doi.org/10.1109/72.857774>.
- Kazunari Hattori, Hiroaki Wakabayashi, and Kenta Tamaki. Predicting key example compounds in competitors' patent applications using structural information alone. *Journal of Chemical Information and Modeling*, 48(1):135–142, January 2008. doi: 10.1021/ci7002686. URL <https://doi.org/10.1021/ci7002686>.
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3:79–87, 1991.
- Michael I. Jordan and Robert A. Jacobs. Hierarchical mixtures of experts and the EM algorithm. In *Proceedings of 1993 International Conference on Neural Networks (IJCNN-93-Nagoya, Japan)*.

- IEEE, 1993. doi: 10.1109/ijcnn.1993.716791. URL <https://doi.org/10.1109/ijcnn.1993.716791>.
- Franziska Kruger, Nikolas Fechner, and Nikolaus Stiefl. Automated identification of chemical series: Classifying like a medicinal chemist. *Journal of Chemical Information and Modeling*, 60(6):2888–2902, May 2020. doi: 10.1021/acs.jcim.0c00204. URL <https://doi.org/10.1021/acs.jcim.0c00204>.
- Eelke B. Lenselink, Niels ten Dijke, Brandon Bongers, George Papadatos, Herman W. T. van Vlijmen, Wojtek Kowalczyk, Adriaan P. IJzerman, and Gerard J. P. van Westen. Beyond the hype: deep neural networks outperform established methods using a ChEMBL bioactivity benchmark set. *Journal of Cheminformatics*, 9(1), August 2017a. doi: 10.1186/s13321-017-0232-0. URL <https://doi.org/10.1186/s13321-017-0232-0>.
- Eelke B. Lenselink, Niels ten Dijke, Brandon Bongers, George Papadatos, Herman W. T. van Vlijmen, Wojtek Kowalczyk, Adriaan P. IJzerman, and Gerard J. P. Van westen. Beyond the hype: Deep neural networks outperform established methods using a chembl bioactivity benchmark set [version 1]. 4TU.ResearchData <https://doi.org/10.4121/uuid:547e8014-d662-4852-9840-c1ef065d03ef>, Jul 2017b. URL https://data.4tu.nl/articles/dataset/Beyond_the_Hype_Deep_Neural_Networks_Outperform_Established_Methods_Using_A_ChEMBL_Bioactivity_Benchmark_Set_version_1_/12694478/1.
- Michael M. Mysinger, Michael Carchia, John. J. Irwin, and Brian K. Shoichet. Directory of useful decoys, enhanced (DUD-e): Better ligands and decoys for better benchmarking. *Journal of Medicinal Chemistry*, 55(14):6582–6594, July 2012. doi: 10.1021/jm300687e. URL <https://doi.org/10.1021/jm300687e>.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12: 2825–2830, 2011.
- RDKit, online. RDKit: Open-source cheminformatics. <http://www.rdkit.org>, 2021. [v. 2021.3.3].
- Carlos Riquelme Ruiz, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=FrIDgjDOHlu>.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer, 2017. URL <https://arxiv.org/abs/1701.06538>.
- Bin Tang, Malcolm I. Heywood, and Michael Shepherd. Input partitioning to mixture of experts. In *Proceedings of the 2002 International Joint Conference on Neural Networks. IJCNN'02 (Cat. No.02CH37290)*. IEEE, 2002. doi: 10.1109/ijcnn.2002.1005474. URL <https://doi.org/10.1109/ijcnn.2002.1005474>.
- Hua Yuan, Yongyan Wang, and Yiyu Cheng. Local and global quantitative structure-activity relationship modeling and prediction for the baseline toxicity. *Journal of Chemical Information and Modeling*, 47(1):159–169, December 2006. doi: 10.1021/ci600299j. URL <https://doi.org/10.1021/ci600299j>.

5 APPENDIX

5.1 DESCRIPTORS

We used the rdkit (RDKit, online) implementation of Morgan fingerprints.

5.2 HYPERPARAMETER SEARCH

For all experiments, parameters were optimised on the valid set without cross-validation. For Random Forests, combinations of hyperparameters were systematically screened from maximum depth = {5, 25, 50, 75, 100} and number of estimators = {100, 200}. We used the scikit-learn (Pedregosa et al., 2011) implementation of Random Forests (Breiman, 2001) regression. For Glolloc networks, hidden layer sizes were the only variable hyperparameter (two hidden layers with [64, 64] or [1024, 1024] units). Only fully connected layers were used, with fixed parameters for dropout (0.25 for experts, 0.025 for the gate), early stopping patience (15), batch size (32), learning rate (0.0005), and temperature (1) for the gate softmax. For the datasets in this paper, which are relatively small, the networks are easily trainable on a single CPU within a few minutes.

5.3 GLOLLOC EXPERTS TRAINING

Expert losses are defined by the following equation:

$$loss_{expert_i} = \frac{1}{\sum_n^N \sum_t^T \mathbf{1}_{\{n,t\} \in expert_i}} \sum_n^N \sum_t^T \mathbf{1}_{\{n,t\} \in expert_i} (Y_{nt} - O_{nti})^2 \quad (3)$$

Where $expert_i$ is a set of pairs of molecules and tasks $\{n, t\}$ relevant to expert i , and the indicator function $\mathbf{1}_{\{n,t\} \in expert_i}$ returns 1 if $\{n, t\}$ is relevant to the i th expert, 0 if not.

In practice, experts are defined using masks, so that the losses for irrelevant data points are not taken into account during the training process: e.g. if an expert focuses on pyrazole-containing molecular compounds, all losses for molecules not containing a pyrazole substructure will be zeroed out. Experts are independent, each having their own optimiser and loss. This independence makes training the network challenging, as validation sets need to be carefully defined and loss curves monitored for each expert - local experts are in effect trained on subsets of the training data, and they are prone to overfitting. This is only partially addressed by adding learning rate decay for each optimiser.

5.3.1 THE GATE NETWORK

The gate is a fully connected network with a softmax function, outputting a gate weight matrix of dimension [batch size, number of tasks, number of experts], the softmax function being applied on the last dimension, so that weights across experts sum to 1 for each task.

5.4 AUTOMATED METHODS TO DEFINE CHEMICAL SERIES

Glolloc can use automated methods to define chemical series ("autosmarts"): tree-based, or neighbourhood based. They both rely on finding a set of maximum common substructures in the dataset, using hierarchical clustering or neighbourhood membership, respectively. We used the RDKit (RDKit, online) implementation of fuzzy MCS - with a minimum of 20 molecules per MCS.

The tree-based method (Algorithm 1) is similar to the one described recently by Kruger et al. (Kruger et al., 2020): building fuzzy maximum common substructures (MCS) from nodes in a tree obtained using an agglomerative clustering algorithm, going from the root to the leaves: when a series (MCS) is found in a parent node, children are discarded. The difference between the Glolloc implementation and the one by Kruger et al. is the criterion to validate the MCS: they use a probability that a random molecule matches the scaffold, whereas Glolloc sets a minimum number of atoms for the scaffold, to guarantee a large enough scaffold (default value: 11 atoms). Strict MCS ring matching rules and loose rules for atom and bond matching were arbitrarily used - the impact of different MCS algorithm settings is not explored in this paper.

Algorithm 1 Tree-based autosmarts algorithm: find fuzzy maximum common substructures (MCS) in different regions of the chemical space using hierarchical clustering

```

min_atoms ← user parameter                                ▷ minimum number of atoms in MCS
min_node_size ← user parameter                          ▷ minimum number of molecules to construct an MCS from
T ← dendrogram(D)                                     ▷ construct dendrogram from the molecule data matrix D
N ← node_dict(T)                                     ▷ dictionary mapping tree nodes to molecules
N' ← sort(N)                                         ▷ sort nodes from root to leaves
MCS ← ∅                                                ▷ the MCS set stores the maximum common substructures
for node, molecules_in_node in N' do
  if count(molecules_in_node) ≥ min_node_size then
    mcs ← fuzzy_maximum_substructure(molecules_in_node)
    ▷ if a large enough MCS is found in parent node
    if mcs exists and mcs.n_atoms ≥ min_atoms then
      MCS ← MCS ∪ {mcs}
      set children of node in N' to ∅                ▷ discard children nodes
    end if
  end if
end for

```

Other methods of detecting chemical series in a data set can be used beyond tree-based, e.g. varying the number of molecules that can match an MCS in a data set, fragment-based methods, scaffold tree methods, or nearest neighbours algorithms. The neighbourhood-based autosmarts method for series identification also implemented in Glolloc (Algorithm 2) is based on the algorithm described by Hattori et al. (Hattori et al., 2008) and adapted for MCS series definition. Hattori et al.’s algorithm is based on iteratively finding representative molecules with the largest neighbourhood and removing them and their neighbours from the pool of molecules. This method is dependent on a similarity threshold parameter to define the neighbours. For Glolloc, neighbourhoods are also iteratively defined, but at each iteration a fuzzy MCS from the top neighbourhood is found and all molecules matching that MCS are removed from the pool used at the next iteration. The neighbourhood-based approach is not used in this paper for benchmarks as it seemed to provide similar results to the tree-based method when varying the similarity threshold for our data sets.

Algorithm 2 Neighbourhood autosmarts algorithm: find fuzzy maximum common substructures (MCS) in different regions of the chemical space using nearest neighbours

```

min_similarity ← user parameter                        ▷ similarity cutoff
min_molecules ← user parameter                      ▷ minimum number of molecules per MCS
bag_of_molecules ← set of all molecules indices
S ← similarity_matrix(D)                          ▷ similarity matrix from molecule data matrix D
MCS ← ∅                                             ▷ the MCS set stores the maximum common substructures
while bag_of_molecule is not empty do
  ▷ find neighbourhood of molecule with max number of neighbours
  top_neighbourhood ← f(S, min_similarity)
  if count(top_neighbourhood) ≥ min_molecules then
    mcs ← fuzzy_maximum_substructure(top_neighbourhood)
    if mcs exists then                                ▷ if MCS is found in neighbourhood
      MCS ← MCS ∪ {mcs}
      mcs_matches ← matches(bag_of_molecule, mcs)
      ▷ discard molecules matching MCS for next iteration
      bag_of_molecule ← bag_of_molecule \ mcs_matches
      S[mcs_matches] ← 0.                             ▷ set similarities of discarded molecules to 0
    end if
  end if
end while

```
