

1 A comparison of nine machine learning models based on an
2 expanded mutagenicity dataset and their application for
3 predicting pyrrolizidine alkaloid mutagenicity

4 Christoph Helma^{*1}, Verena Schöning⁴, Philipp Boss³, and Jürgen Drewe²

5 ¹in silico toxicology gmbh, Rastatterstrasse 41, 4057 Basel, Switzerland

6 ²Zeller AG, Seeblickstrasse 4, 8590 Romanshorn, Switzerland

7 ³Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular
8 Medicine in the Helmholtz Association, Robert-Rössle-Strasse 10, Berlin, 13125, Germany

9 ⁴Clinical Pharmacology and Toxicology, Department of General Internal Medicine,
10 Bern University Hospital, University of Bern, Inselspital, 3010 Bern, Switzerland

11 ^{*} Correspondence: Christoph Helma <helma@in-silico.ch>

12 Random forest, support vector machine, logistic regression, neural
13 networks and k-nearest neighbor (**lazar**) algorithms, were applied to new
14 *Salmonella* mutagenicity dataset with 8309 unique chemical structures. The
15 best prediction accuracies in 10-fold-crossvalidation were obtained with
16 **lazar** models and MolPrint2D descriptors, that gave accuracies (%) similar
17 to the interlaboratory variability of the Ames test.

18 **TODO:** PA results

19 **Introduction**

20 **TODO:** rationale for investigation

21 The main objectives of this study were

- 22 • to generate a new mutagenicity training dataset, by combining the most compre-
23 hensive public datasets
- 24 • to compare the performance of MolPrint2D (*MP2D*) fingerprints with Chemistry
25 Development Kit (*CDK*) descriptors
- 26 • to compare the performance of global QSAR models (random forests (*RF*), support
27 vector machines (*SVM*), logistic regression (*LR*), neural nets (*NN*)) with local
28 models (*lazar*)
- 29 • to apply these models for the prediction of pyrrolizidine alkaloid mutagenicity

30 **Materials and Methods**

31 **Data**

32 **Mutagenicity training data**

33 An identical training dataset was used for all models. The training dataset was compiled
34 from the following sources:

- 35 • Kazius/Bursi Dataset (4337 compounds, Kazius, McGuire, and Bursi (2005)):
36 http://cheminformatics.org/datasets/bursi/cas_4337.zip
- 37 • Hansen Dataset (6513 compounds, Hansen et al. (2009)): [http://doc.ml.tu-berlin.](http://doc.ml.tu-berlin.de/toxbenchmark/Mutagenicity_N6512.csv)
38 [de/toxbenchmark/Mutagenicity_N6512.csv](http://doc.ml.tu-berlin.de/toxbenchmark/Mutagenicity_N6512.csv)
- 39 • EFSA Dataset (695 compounds EFSA (2016)): [https://data.europa.eu/euodp/](https://data.europa.eu/euodp/data/storage/f/2017-0719T142131/GENOTOX%20data%20and%20dictionary.xls)
40 [data/storage/f/2017-0719T142131/GENOTOX%20data%20and%20dictionary.xls](https://data.europa.eu/euodp/data/storage/f/2017-0719T142131/GENOTOX%20data%20and%20dictionary.xls)

41 Mutagenicity classifications from Kazius and Hansen datasets were used without further
42 processing. To achieve consistency with these datasets, EFSA compounds were classified
43 as mutagenic, if at least one positive result was found for TA98 or T100 Salmonella

44 strains.

45 Dataset merges were based on unique SMILES (*Simplified Molecular Input Line Entry*
46 *Specification*) strings of the compound structures. Duplicated experimental data with
47 the same outcome was merged into a single value, because it is likely that it originated
48 from the same experiment. Contradictory results were kept as multiple measurements
49 in the database. The combined training dataset contains 8309 unique structures.

50 Source code for all data download, extraction and merge operations is pub-
51 licly available from the git repository <https://git.in-silico.ch/mutagenicity-paper>
52 under a GPL3 License. The new combined dataset can be found at <https://git.in-silico.ch/mutagenicity-paper/tree/data/mutagenicity.csv>.
53

54 **Pyrrolizidine alkaloid (PA) dataset**

55 The testing dataset consisted of 602 different PAs.

56 The PA dataset was created from five independent, necine base substructure searches in
57 PubChem (<https://pubchem.ncbi.nlm.nih.gov/>) and compared to the PAs listed in the
58 EFSA publication EFSA (2011) and the book by Mattocks Mattocks (1986), to ensure,
59 that all major PAs were included. PAs mentioned in these publications which were
60 not found in the downloaded substances were searched individually in PubChem and,
61 if available, downloaded separately. Non-PA substances, duplicates, and isomers were
62 removed from the files, but artificial PAs, even if unlikely to occur in nature, were kept.
63 The resulting PA dataset comprised a total of 602 different PAs.

64 The PAs in the dataset were classified according to structural features. A total of 9
65 different structural features were assigned to the necine base, modifications of the necine
66 base and to the necic acid:

67 For the necine base, the following structural features were chosen:

- 68 • Retronecine-type (1,2-unsaturated necine base)
- 69 • Otonecine-type (1,2-unsaturated necine base)
- 70 • Platynecine-type (1,2-saturated necine base)

71 For the modifications of the necine base, the following structural features were chosen:

- 72 • N-oxide-type
- 73 • Tertiary-type (PAs which were neither from the N-oxide- nor DHP-type)
- 74 • DHP-type (pyrrolic ester)

75 For the necic acid, the following structural features were chosen:

- 76 • Monoester-type
- 77 • Open-ring diester-type
- 78 • Macrocyclic diester-type

79 The compilation of the PA dataset is described in detail in Schöning et al. (2017).

80 **Descriptors**

81 **MolPrint2D (*MP2D*) fingerprints**

82 MolPrint2D fingerprints (O’Boyle et al. (2011)) use atom environments as molecular
83 representation. They determine for each atom in a molecule, the atom types of its
84 connected atoms to represent their chemical environment. This resembles basically the
85 chemical concept of functional groups.

86 In contrast to predefined lists of fragments (e.g. FP3, FP4 or MACCs fingerprints) or
87 descriptors (e.g CDK) they are generated dynamically from chemical structures. This
88 has the advantage that they can capture substructures of toxicological relevance that
89 are not included in other descriptors.

90 Chemical similarities (e.g. Tanimoto indices) can be calculated very efficiently with Mol-

Print2D fingerprints. Using them as descriptors for global models leads however to huge, sparsely populated matrices that cannot be handled with traditional machine learning algorithms. In our experiments none of the R and Tensorflow algorithms was capable to use them as descriptors.

MolPrint2D fingerprints were calculated with the OpenBabel cheminformatics library (O’Boyle et al. (2011)).

Chemistry Development Kit (*CDK*) descriptors

Molecular 1D and 2D descriptors were calculated with the PaDEL-Descriptors program (<http://www.yapcwssoft.com> version 2.21, Yap (2011)). PaDEL uses the Chemistry Development Kit (*CDK*, <https://cdk.github.io/index.html>) library for descriptor calculations.

As the training dataset contained over 8309 instances, it was decided to delete instances with missing values during data pre-processing. Furthermore, substances with equivocal outcome were removed. The final training dataset contained 8080 instances with known mutagenic potential.

During feature selection, descriptors with near zero variance were removed using ‘*NearZeroVar*’-function (package ‘*caret*’). If the percentage of the most common value was more than 90% or when the frequency ratio of the most common value to the second most common value was greater than 95:5 (e.g. 95 instances of the most common value and only 5 or less instances of the second most common value), a descriptor was classified as having a near zero variance. After that, highly correlated descriptors were removed using the ‘*findCorrelation*’-function (package ‘*caret*’) with a cut-off of 0.9. This resulted in a training dataset with 516 descriptors. These descriptors were scaled to be in the range between 0 and 1 using the ‘*preProcess*’-function (package ‘*caret*’). The scaling routine was saved in order to apply the same scaling on the testing dataset. As these

three steps did not consider the dependent variable (experimental mutagenicity), it was decided that they do not need to be included in the cross-validation of the model. To further reduce the number of features, a LASSO (*least absolute shrinkage and selection operator*) regression was performed using the ‘*glmnet*’-function (package ‘*glmnet*’). The reduced dataset was used for the generation of the pre-trained models.

CDK descriptors were used in global (RF, SVM, LR, NN) and local (**lazar**) models.

Algorithms

lazar

lazar (*lazy structure activity relationships*) is a modular framework for read-across model development and validation. It follows the following basic workflow: For a given chemical structure **lazar**:

- searches in a database for similar structures (neighbours) with experimental data,
- builds a local QSAR model with these neighbours and
- uses this model to predict the unknown activity of the query compound.

This procedure resembles an automated version of read across predictions in toxicology, in machine learning terms it would be classified as a k-nearest-neighbour algorithm.

Apart from this basic workflow, **lazar** is completely modular and allows the researcher to use arbitrary algorithms for similarity searches and local QSAR (*Quantitative structure–activity relationship*) modelling. Algorithms used within this study are described in the following sections.

Neighbour identification

Utilizing this modularity, similarity calculations were based both on MolPrint2D finger-

138 prints and on CDK descriptors.

139 For MolPrint2D fingerprints chemical similarity between two compounds a and b is
140 expressed as the proportion between atom environments common in both structures
141 $A \cap B$ and the total number of atom environments $A \cup B$ (Jaccard/Tanimoto index).

$$sim = \frac{|A \cap B|}{|A \cup B|}$$

142 For CDK descriptors chemical similarity between two compounds a and b is expressed
143 as the cosine similarity between the descriptor vectors A for a and B for b .

$$sim = \frac{A \cdot B}{|A||B|}$$

144 Threshold selection is a trade-off between prediction accuracy (high threshold) and the
145 number of predictable compounds (low threshold). As it is in many practical cases
146 desirable to make predictions even in the absence of closely related neighbours, we follow
147 a tiered approach:

- 148 • First a similarity threshold of 0.5 is used to collect neighbours, to create a local
149 QSAR model and to make a prediction for the query compound. This are predic-
150 tions with *high confidence*.
- 151 • If any of these steps fails, the procedure is repeated with a similarity threshold
152 of 0.2 and the prediction is flagged with a warning that it might be out of the
153 applicability domain of the training data (*low confidence*).
- 154 • Similarity thresholds of 0.5 and 0.2 are the default values chosen by the software
155 developers and remained unchanged during the course of these experiments.

156 Compounds with the same structure as the query structure are automatically eliminated

157 from neighbours to obtain unbiased predictions in the presence of duplicates.

158 **Local QSAR models and predictions**

159 Only similar compounds (neighbours) above the threshold are used for local QSAR
160 models. In this investigation, we are using a weighted majority vote from the neigh-
161 bour’s experimental data for mutagenicity classifications. Probabilities for both classes
162 (mutagenic/non-mutagenic) are calculated according to the following formula and the
163 class with the higher probability is used as prediction outcome.

$$p_c = \frac{\sum \text{sim}_{n,c}}{\sum \text{sim}_n}$$

164 p_c Probability of class c (e.g. mutagenic or non-mutagenic)

165 $\sum \text{sim}_{n,c}$ Sum of similarities of neighbours with class c

166 $\sum \text{sim}_n$ Sum of all neighbours

167 **Applicability domain**

168 The applicability domain (AD) of **lazar** models is determined by the structural diver-
169 sity of the training data. If no similar compounds are found in the training data no
170 predictions will be generated. Warnings are issued if the similarity threshold had to be
171 lowered from 0.5 to 0.2 in order to enable predictions. Predictions without warnings
172 can be considered as close to the applicability domain (*high confidence*) and predictions
173 with warnings as more distant from the applicability domain (*low confidence*). Quantita-
174 tive applicability domain information can be obtained from the similarities of individual
175 neighbours.

176 **Availability**

- 177 • **lazar** experiments for this manuscript: <https://git.in-silico.ch/mutagenicity-paper>
178 (source code, GPL3)
- 179 • **lazar** framework: <https://git.in-silico.ch/lazar> (source code, GPL3)
- 180 • **lazar** GUI: <https://git.in-silico.ch/lazar-gui> (source code, GPL3)
- 181 • Public web interface: <https://lazar.in-silico.ch>

182 **R Random Forest, Support Vector Machines, and Deep Learning**

183 The RF, SVM, and DL models were generated using the R software (R-project for
184 Statistical Computing, <https://www.r-project.org/>; version 3.3.1), specific R packages
185 used are identified for each step in the description below.

186 **Random Forest (*RF*)**

187 For the RF model, the ‘*randomForest*’-function (package ‘*randomForest*’) was used. A
188 forest with 1000 trees with maximal terminal nodes of 200 was grown for the prediction.

189 **Support Vector Machines (*SVM*)**

190 The ‘*svm*’-function (package ‘*e1071*’) with a *radial basis function kernel* was used for the
191 SVM model.

192 **TODO: Verena, Phillip** Sollen wir die DL Modelle ebenso wie die Tensorflow als
193 Neural Nets (NN) bezeichnen?

194 **Deep Learning**

195 The DL model was generated using the ‘*h2o.deeplearning*’-function (package ‘*h2o*’). The
196 DL contained four hidden layer with 70, 50, 50, and 10 neurons, respectively. Other
197 hyperparameter were set as follows: $l1=1.0E-7$, $l2=1.0E-11$, $\epsilon = 1.0E-10$, $\rho =$

198 0.8, and quantile_alpha = 0.5. For all other hyperparameter, the default values were
199 used. Weights and biases were in a first step determined with an unsupervised DL model.
200 These values were then used for the actual, supervised DL model.

201 To validate these models, an internal cross-validation approach was chosen. The training
202 dataset was randomly split in training data, which contained 95% of the data, and
203 validation data, which contain 5% of the data. A feature selection with LASSO on the
204 training data was performed, reducing the number of descriptors to approximately 100.
205 This step was repeated five times. Based on each of the five different training data,
206 the predictive models were trained and the performance tested with the validation data.
207 This step was repeated 10 times.

208 **Applicability domain**

209 **TODO: Verena:** Mit welchen Deskriptoren hast Du den Jaccard index berechnet?
210 Fuer den Jaccard index braucht man binaere Deskriptoren (zB MP2D), mit PaDEL
211 Deskriptoren koennte man zB eine euklidische oder cosinus Distanz berechnen.

212 The AD of the training dataset and the PA dataset was evaluated using the Jaccard
213 distance. A Jaccard distance of '0' indicates that the substances are similar, whereas a
214 value of '1' shows that the substances are different. The Jaccard distance was below 0.2
215 for all PAs relative to the training dataset. Therefore, PA dataset is within the AD of
216 the training dataset and the models can be used to predict the genotoxic potential of
217 the PA dataset.

218 **Availability**

219 R scripts for these experiments can be found in [https://git.in-silico.ch/mutagenicity-](https://git.in-silico.ch/mutagenicity-paper/tree/scripts/R)
220 [paper/tree/scripts/R](https://git.in-silico.ch/mutagenicity-paper/tree/scripts/R).

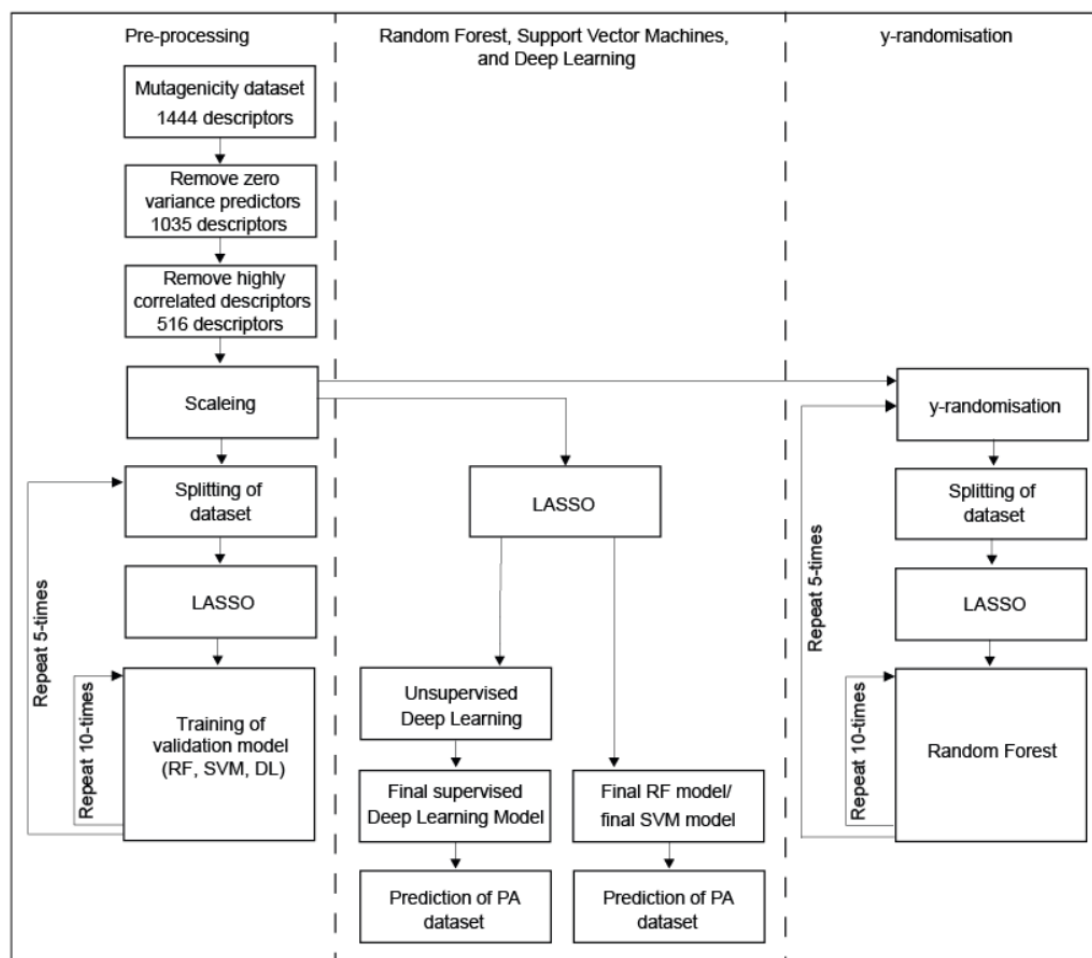


Figure 1: Flowchart of the generation and validation of the models generated in R-project

221 **Tensorflow models**

222 Data pre-processing was done by rank transformation using the ‘*QuantileTransformer*’
223 procedure. A sequential model has been used. Four layers have been used: input layer,
224 two hidden layers (with 12, 8 and 8 nodes, respectively) and one output layer. For the
225 output layer, a sigmoidal activation function and for all other layers the ReLU (‘*Rectified*
226 *Linear Unit*’) activation function was used. Additionally, a L^2 -penalty of 0.001 was used
227 for the input layer. For training of the model, the ADAM algorithm was used to minimise
228 the cross-entropy loss using the default parameters of Keras. Training was performed
229 for 100 epochs with a batch size of 64. The model was implemented with Python 3.6
230 and Keras.

231 **TODO: Philipp** Ich hab die alten Ergebnisse mit feature selection weggelassen, ist das
232 ok? Dann muesste auch dieser Absatz gestrichen werden, oder?

233 **TODO: Philipp** Kannst Du bitte die folgenden Absaetze ergaenzen

234 **Random forests (*RF*)**

235 **Logistic regression (SGD) (*LR-sgd*)**

236 **Logistic regression (scikit) (*LR-scikit*)**

237 **TODO: Philipp, Verena** DL oder NN?

238 **Neural Nets (*NN*)**

239 Alternatively, a DL model was established with Python-based Tensorflow program ([https:](https://www.tensorflow.org/)
240 [//www.tensorflow.org/](https://www.tensorflow.org/)) using the high-level API Keras ([https://www.tensorflow.org/](https://www.tensorflow.org/guide/keras)
241 [guide/keras](https://www.tensorflow.org/guide/keras)) to build the models.

Tensorflow models used the same CDK descriptors as the R models.

Validation

10-fold cross-validation was used for all Tensorflow models.

Availability

Jupyter notebooks for these experiments can be found in <https://git.in-silico.ch/mutagenicity-paper/tree/scripts/tensorflow>.

Results

10-fold crossvalidations

Crossvalidation results are summarized in the following tables: Table ?? shows **lazar** results with MolPrint2D and CDK descriptors, Table ?? R results and Table ?? Tensorflow results.

Table 1: Summary of crossvalidation results with MolPrint2D descriptors

	lazar-HC	lazar-all	RF	LR-sgi	LR-scikit	NN	SVM
Accuracy	84	82	80	84	84	84	84
True positive rate	89	85	81	84	84	85	85
True negative rate	79	79	79	84	84	83	84
Positive predictive value	83	80	78	83	83	83	83
Negative predictive value	86	84	82	84	85	85	86
Nr. predictions	5808	7790	8290	8290	8290	8290	8290

Table 2: Summary of crossvalidation results with CDK descriptors

	lazar-HC	lazar-all	RF	LR-sgi	LR-scikit	NN	SVM
Accuracy	38	38	84	80	80	85	82
True positive rate	17	17	85	79	80	85	82
True negative rate	60	60	82	80	80	85	82
Positive predictive value	31	31	81	81	80	85	82
Negative predictive value	42	42	86	78	80	85	82
Nr. predictions	8083	8083	8065	8065	8065	8065	8065

Figure 2 depicts the position of all crossvalidation results in receiver operating characteristic (ROC) space.

Confusion matrices for all models are available from the git repository <https://git.in-silico.ch/mutagenicity-paper/tree/10-fold-crossvalidations/confusion-matrices/>, individual predictions can be found in <https://git.in-silico.ch/mutagenicity-paper/tree/10-fold-crossvalidations/predictions/>.

The most accurate crossvalidation predictions have been obtained with standard **lazar** models using MolPrint2D descriptors (for predictions with high confidence, for all predictions). Models utilizing CDK descriptors have generally lower accuracies ranging from (R deep learning) to (R/Tensorflow random forests). Sensitivity and specificity is generally well balanced with the exception of **lazar**-CDK (low sensitivity) and R deep learning (low specificity) models.

Pyrrolizidine alkaloid mutagenicity predictions

Mutagenicity predictions from all investigated models for 602 pyrrolizidine alkaloids (PAs) are shown in Table 4. A CSV table with all predictions can be downloaded from

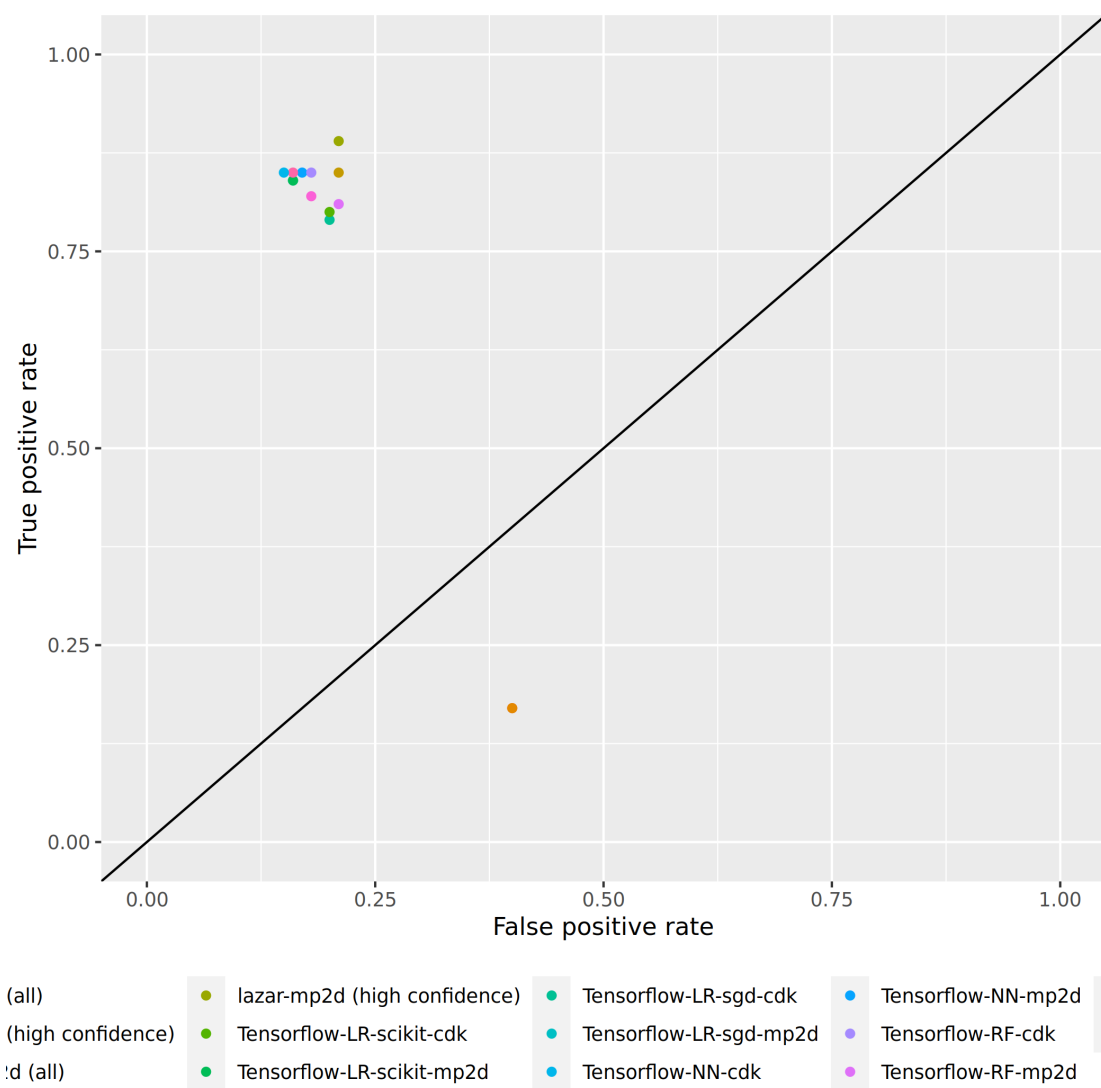
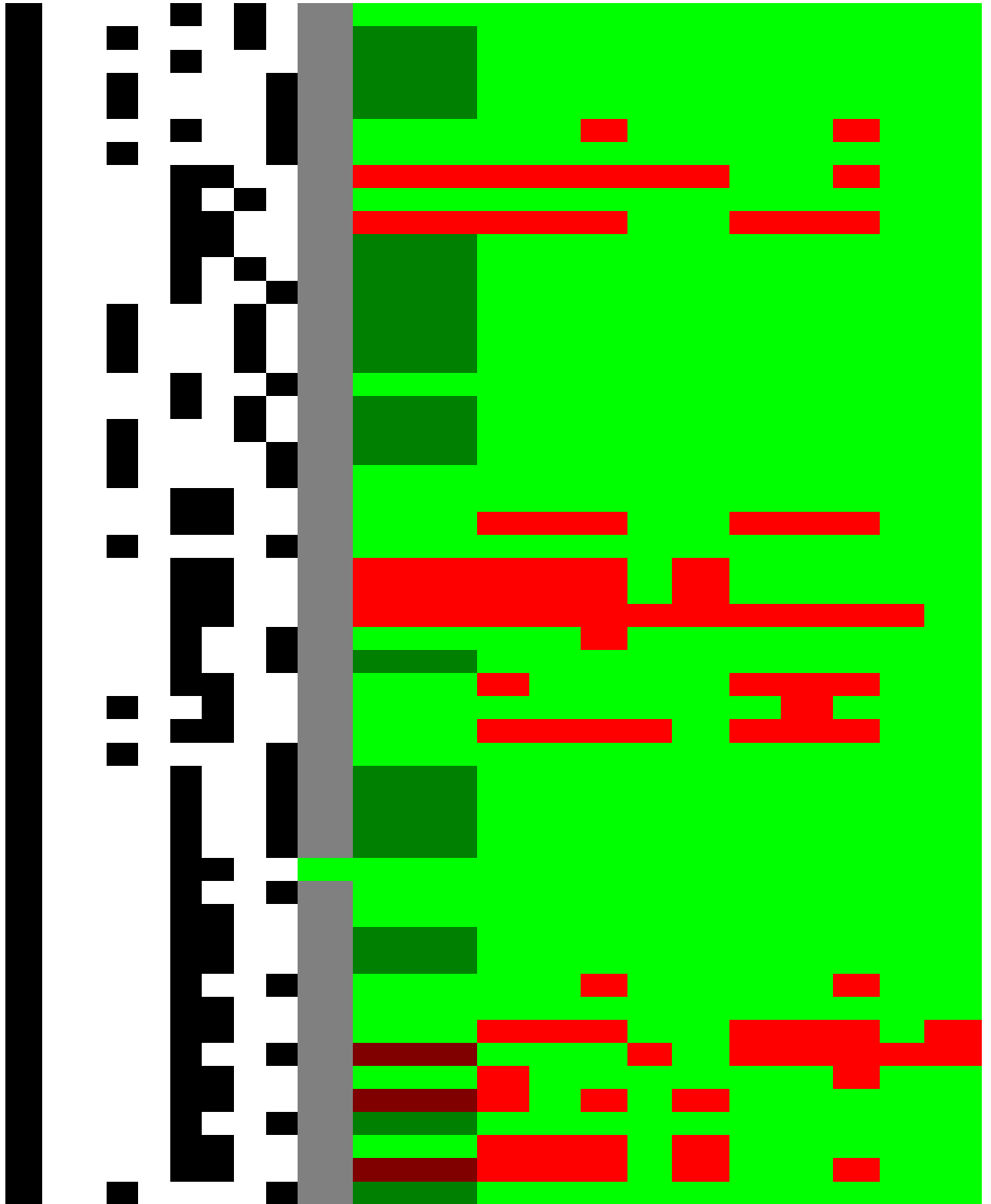
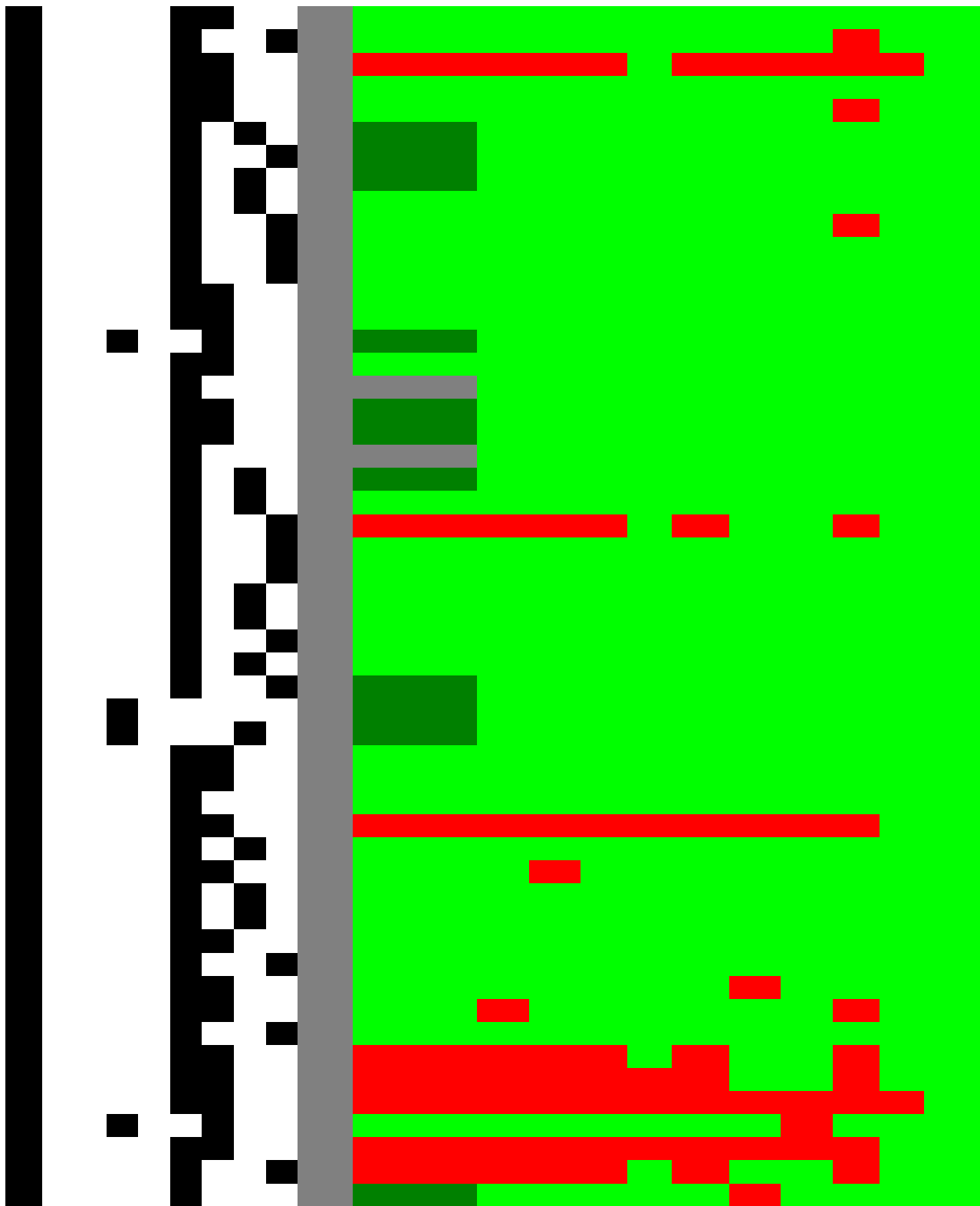
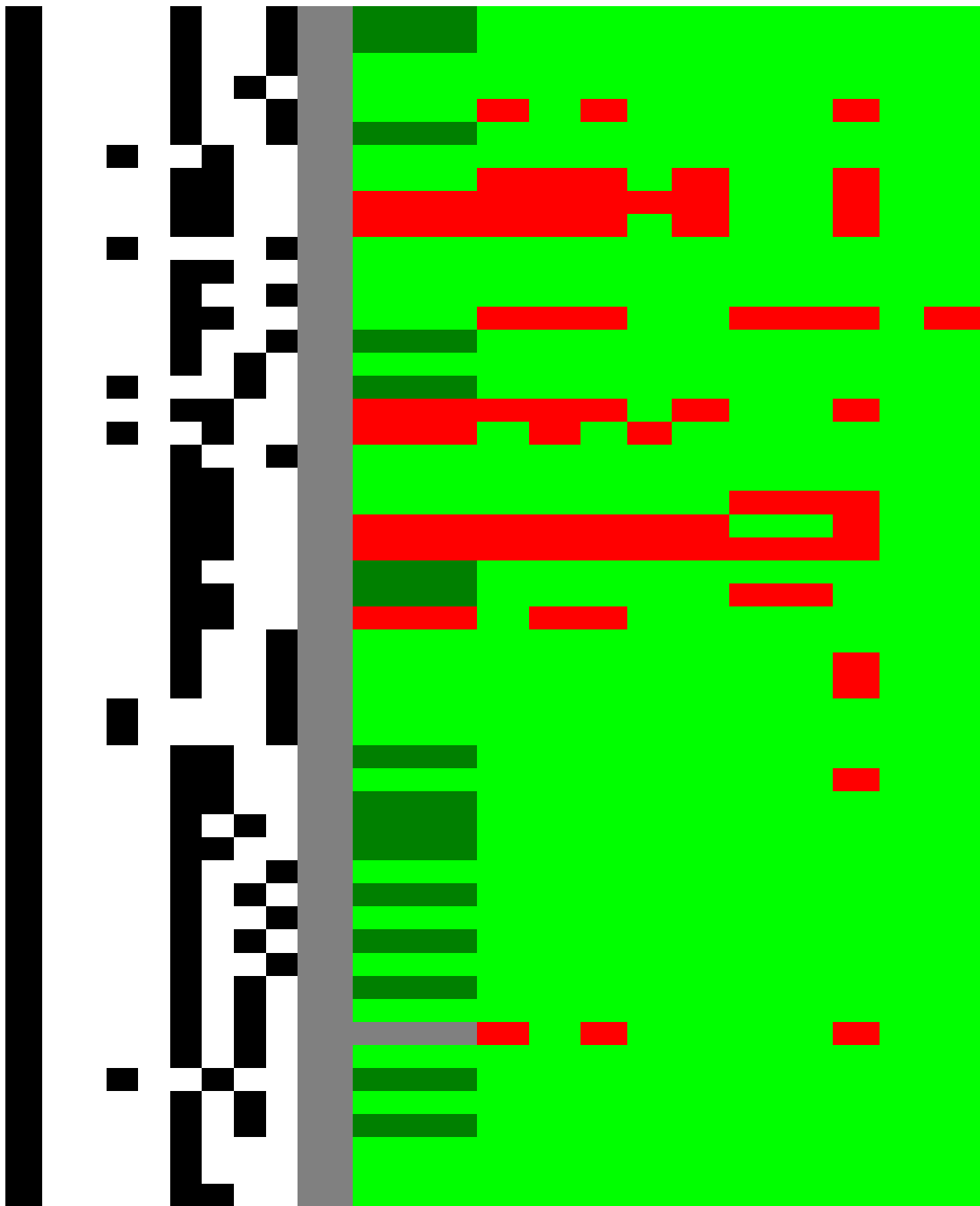
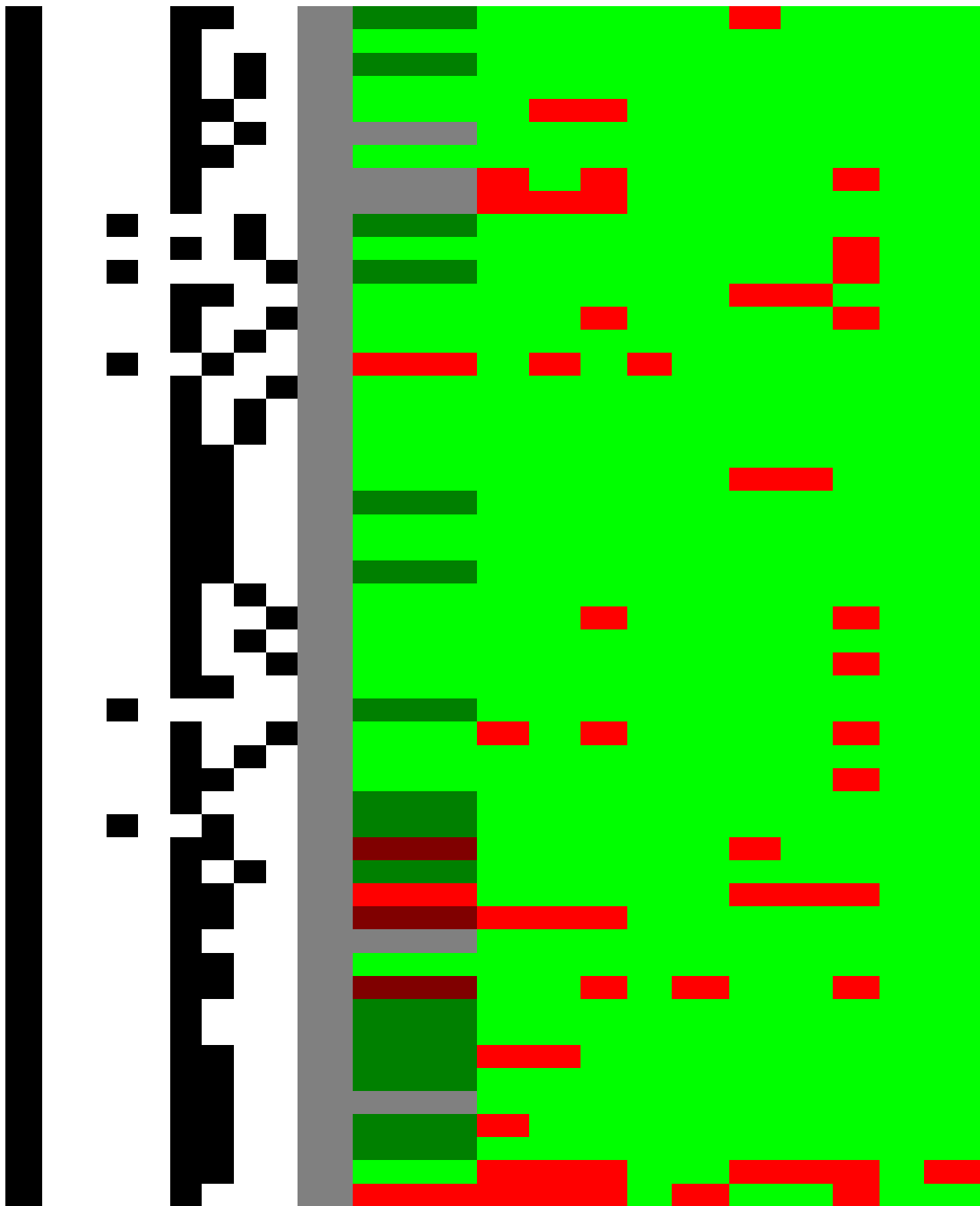


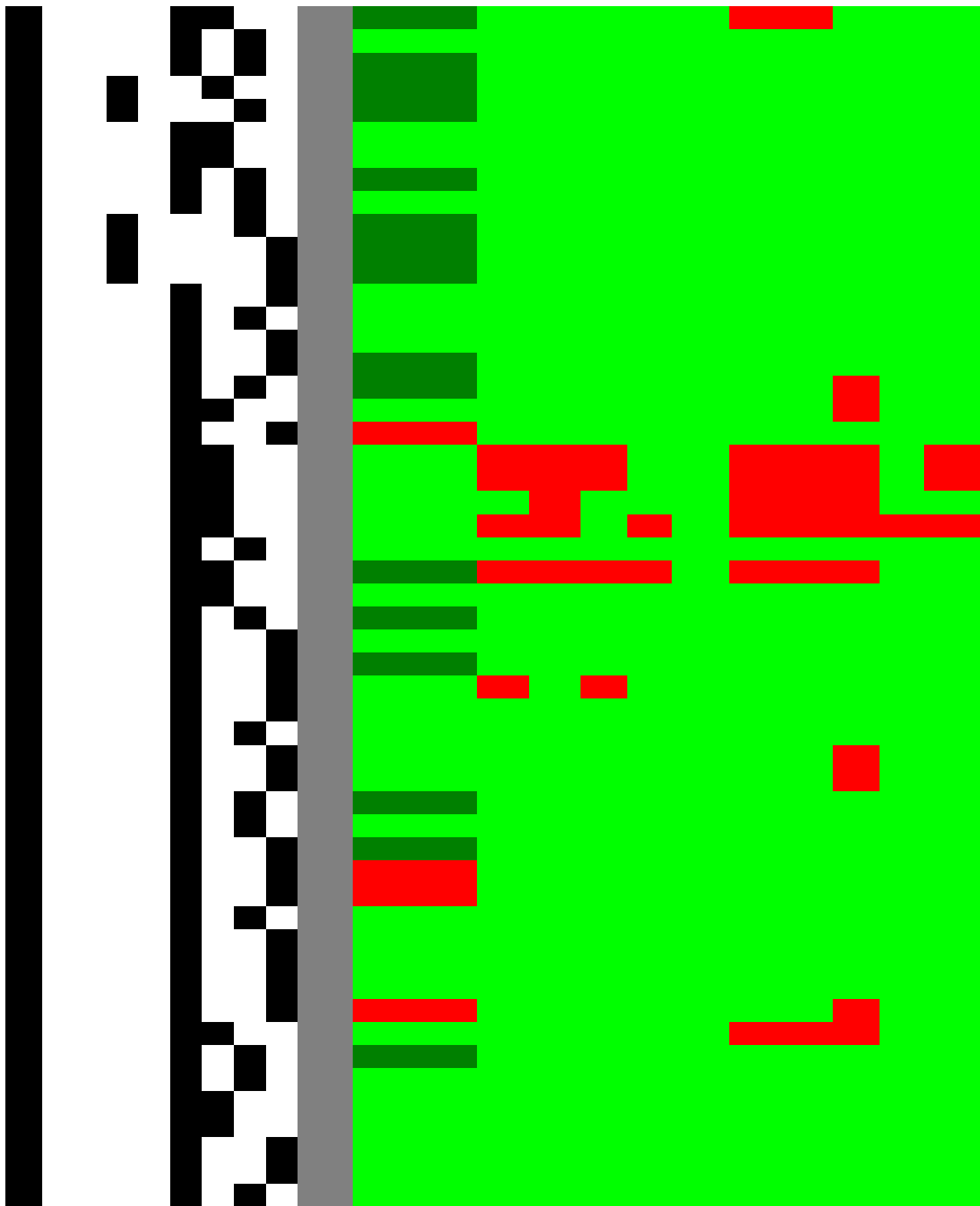
Figure 2: ROC plot of crossvalidation results.

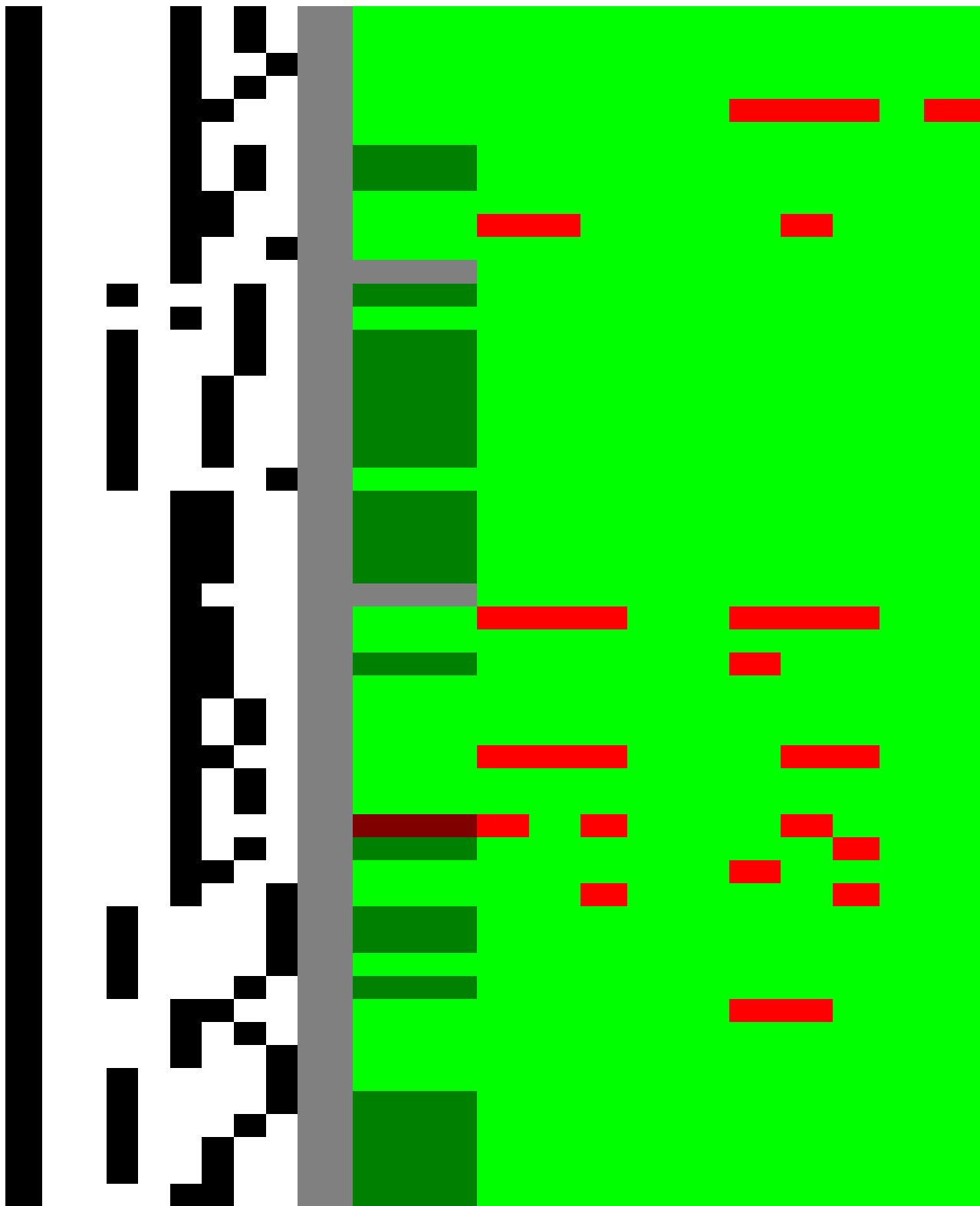


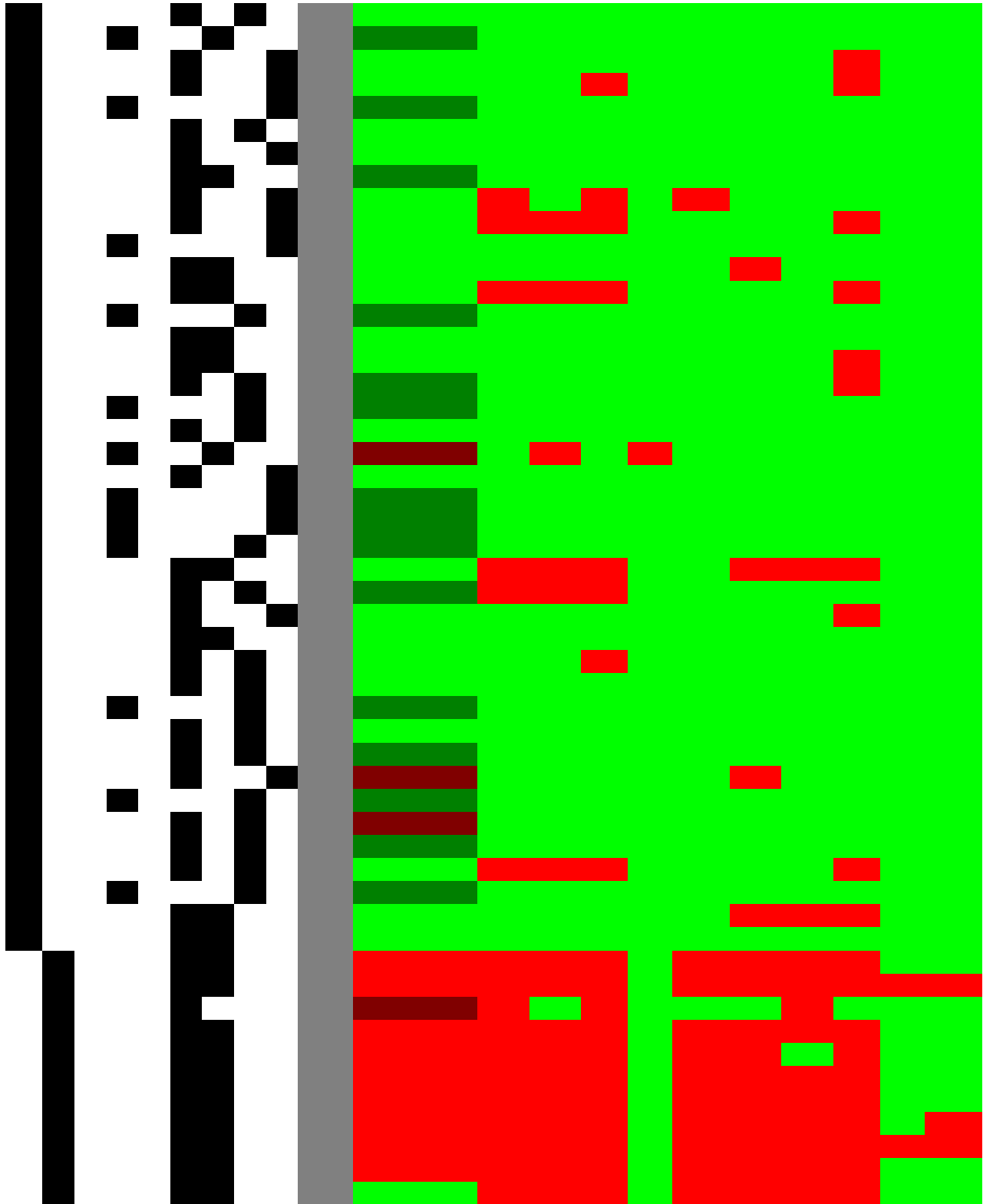


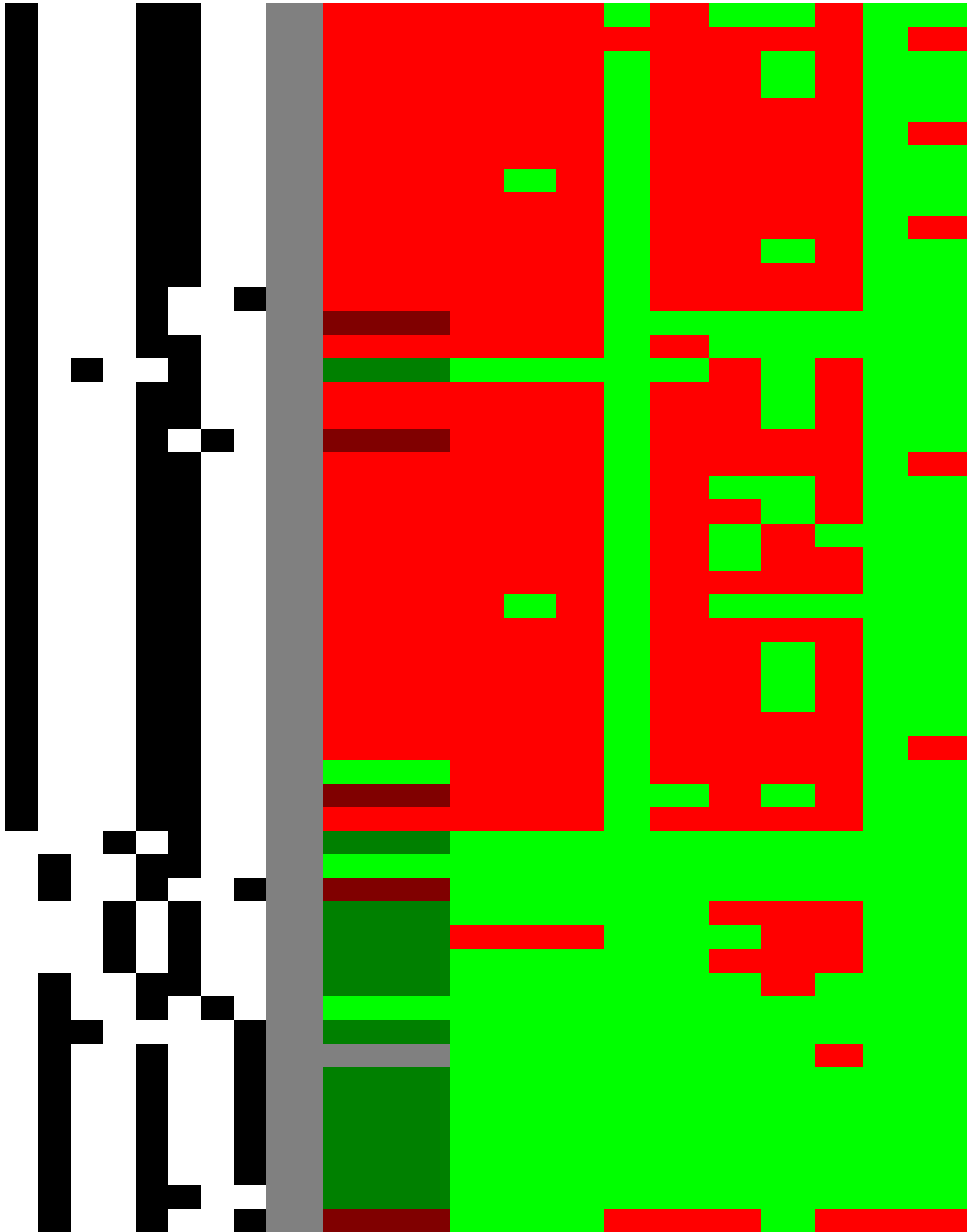


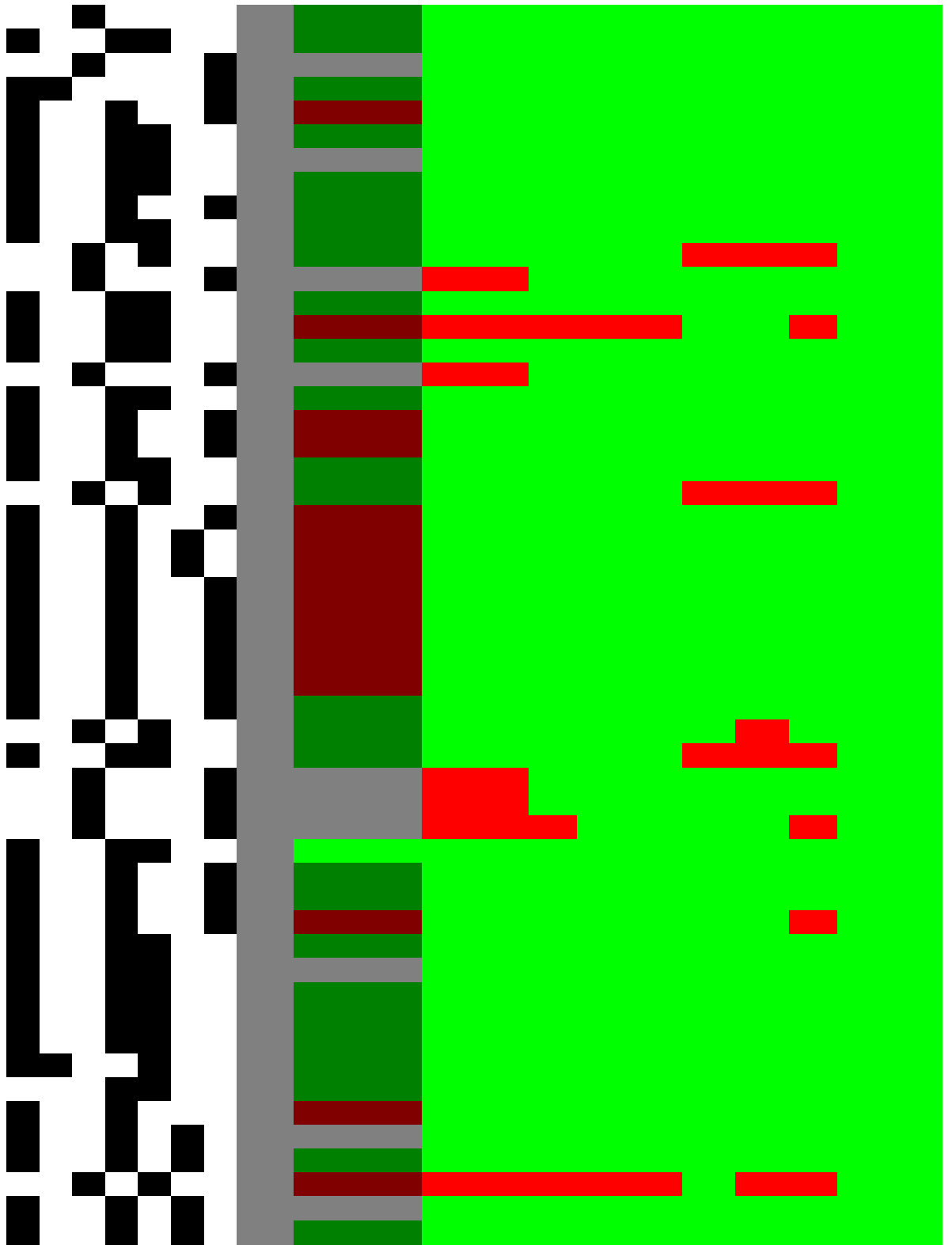


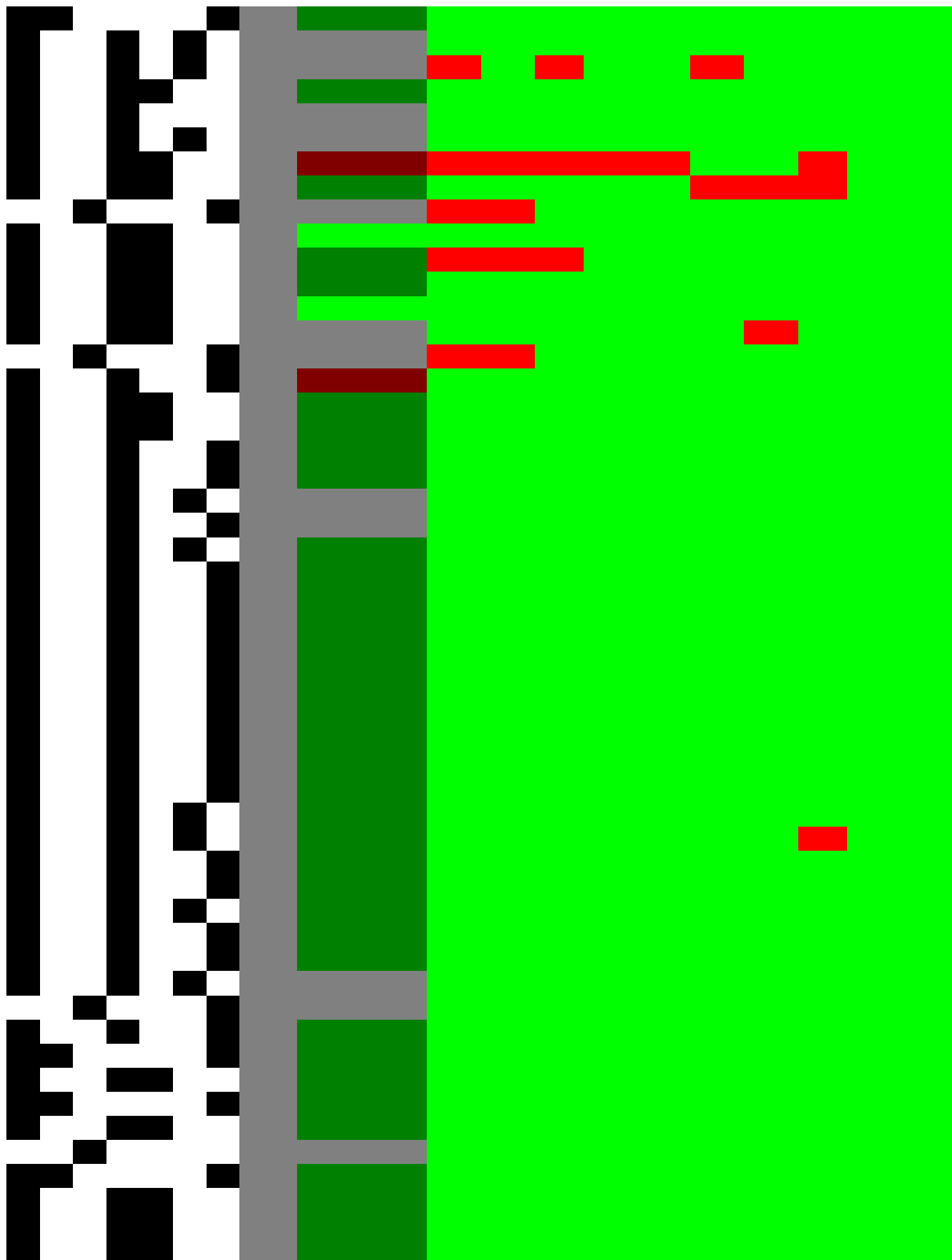


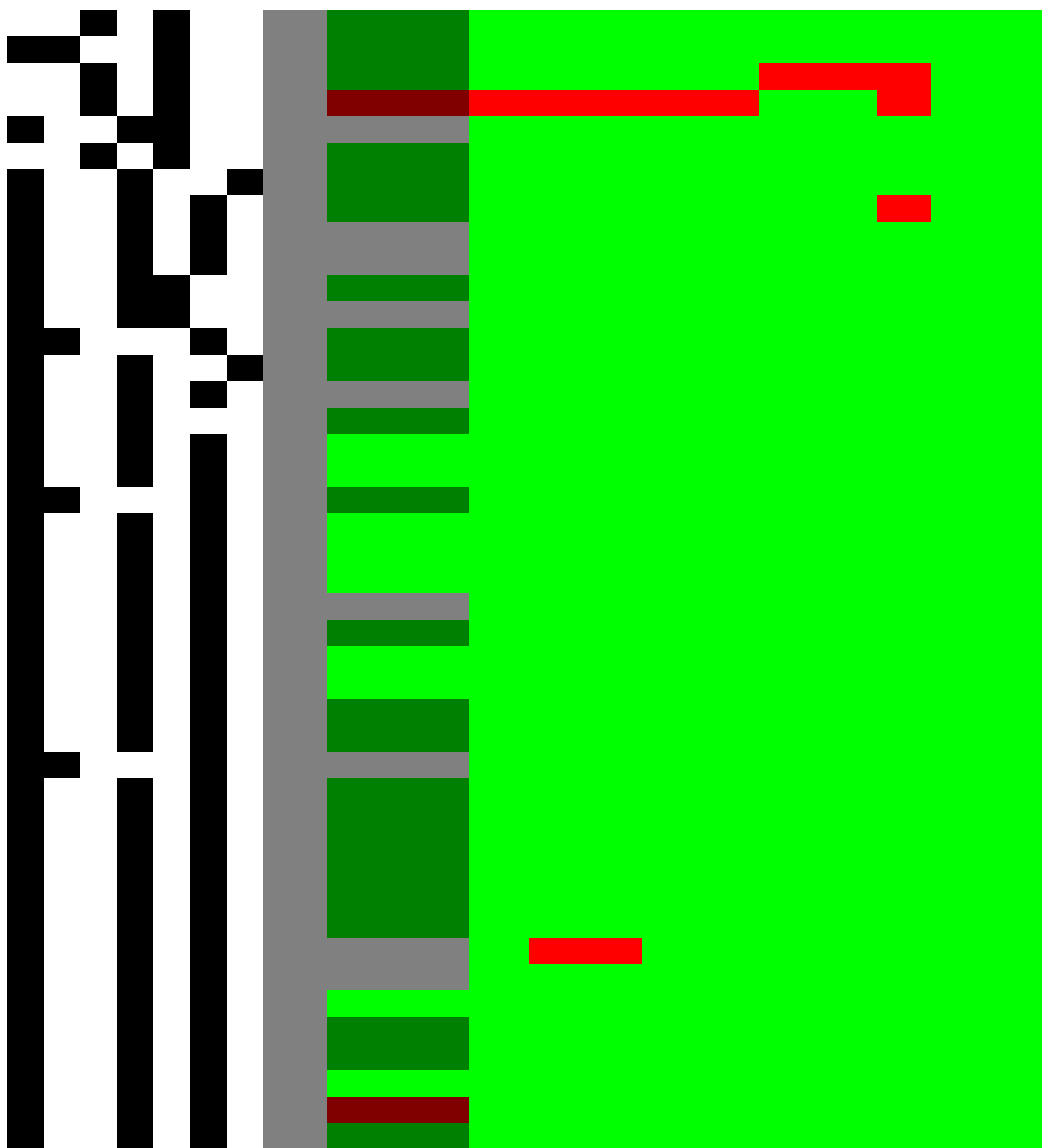












271 Table ?? summarises the number of positive and negative mutagenicity predictions for
 272 all investigated models.

273 For the visualisation of the position of pyrrolizidine alkaloids in respect to the train-

ing data set we have applied t-distributed stochastic neighbor embedding (t-SNE, Maaten and Hinton (2008)) for MolPrint2D and CDK descriptors. t-SNE maps each high-dimensional object (chemical) to a two-dimensional point, maintaining the high-dimensional distances of the objects. Similar objects are represented by nearby points and dissimilar objects are represented by distant points.

Figure 12 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training data in MP2D space (Tanimoto/Jaccard similarity).

Figure 13 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training data in CDK space (Euclidean similarity).

Discussion

Data

A new training dataset for *Salmonella* mutagenicity was created from three different sources (Kazius, McGuire, and Bursi (2005), Hansen et al. (2009), EFSA (2016)). It contains 8309 unique chemical structures, which is according to our knowledge the largest public mutagenicity dataset presently available. The new training data can be downloaded from <https://git.in-silico.ch/mutagenicity-paper/tree/data/mutagenicity.csv>.

Model performance

Table ??, Table ??, Table ?? and Figure 2 show that the standard **lazar** algorithm (with MP2D fingerprints) give the most accurate crossvalidation results. R Random Forests, Support Vector Machines and Tensorflow models have similar accuracies with balanced sensitivity (true position rate) and specificity (true negative rate). **lazar** models with CDK descriptors have low sensitivity and R Deep Learning models have low specificity. The accuracy of **lazar** *in-silico* predictions are comparable to the interlaboratory vari-

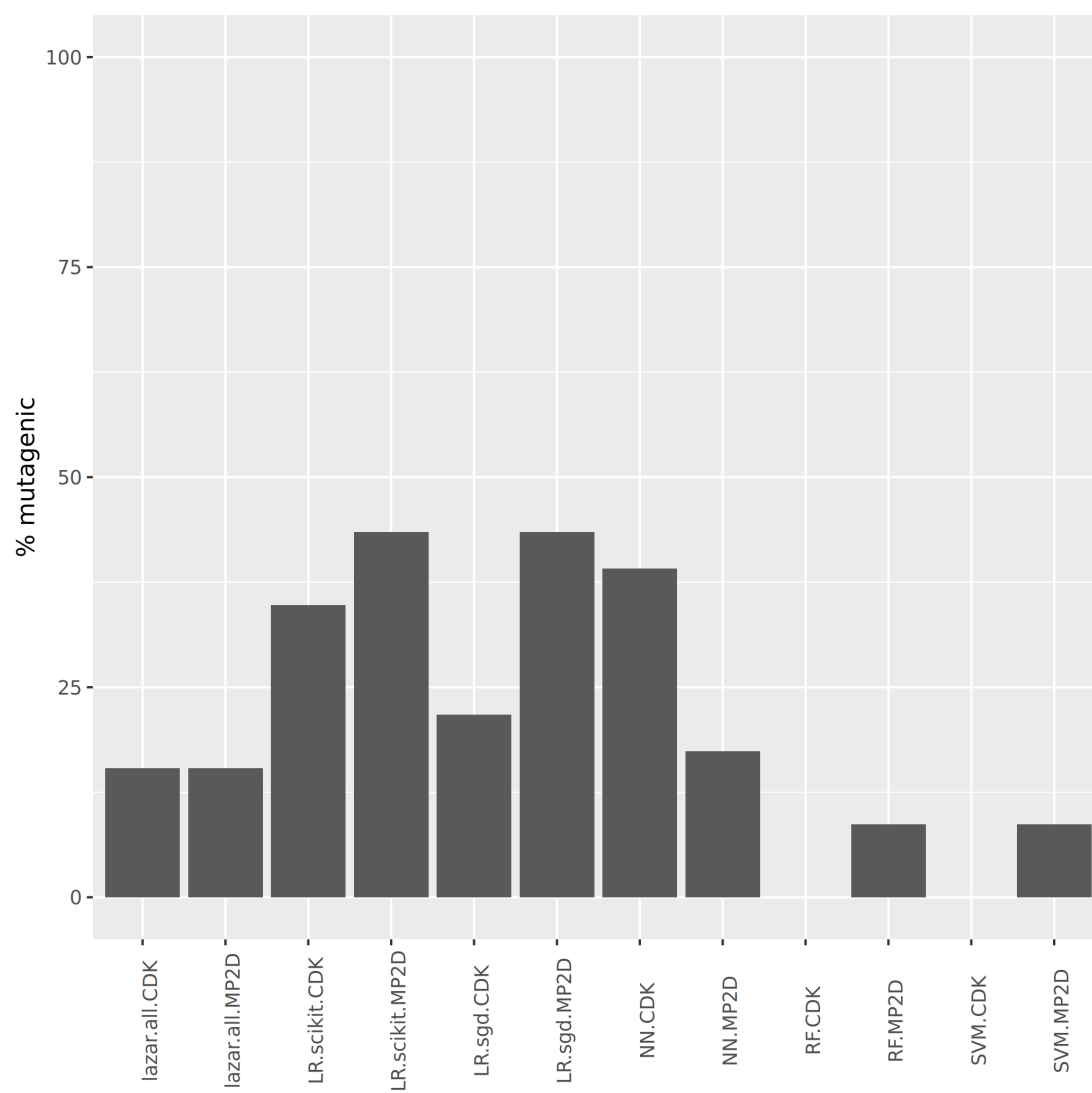


Figure 3: Summary of Dehydropyrrolizidine predictions

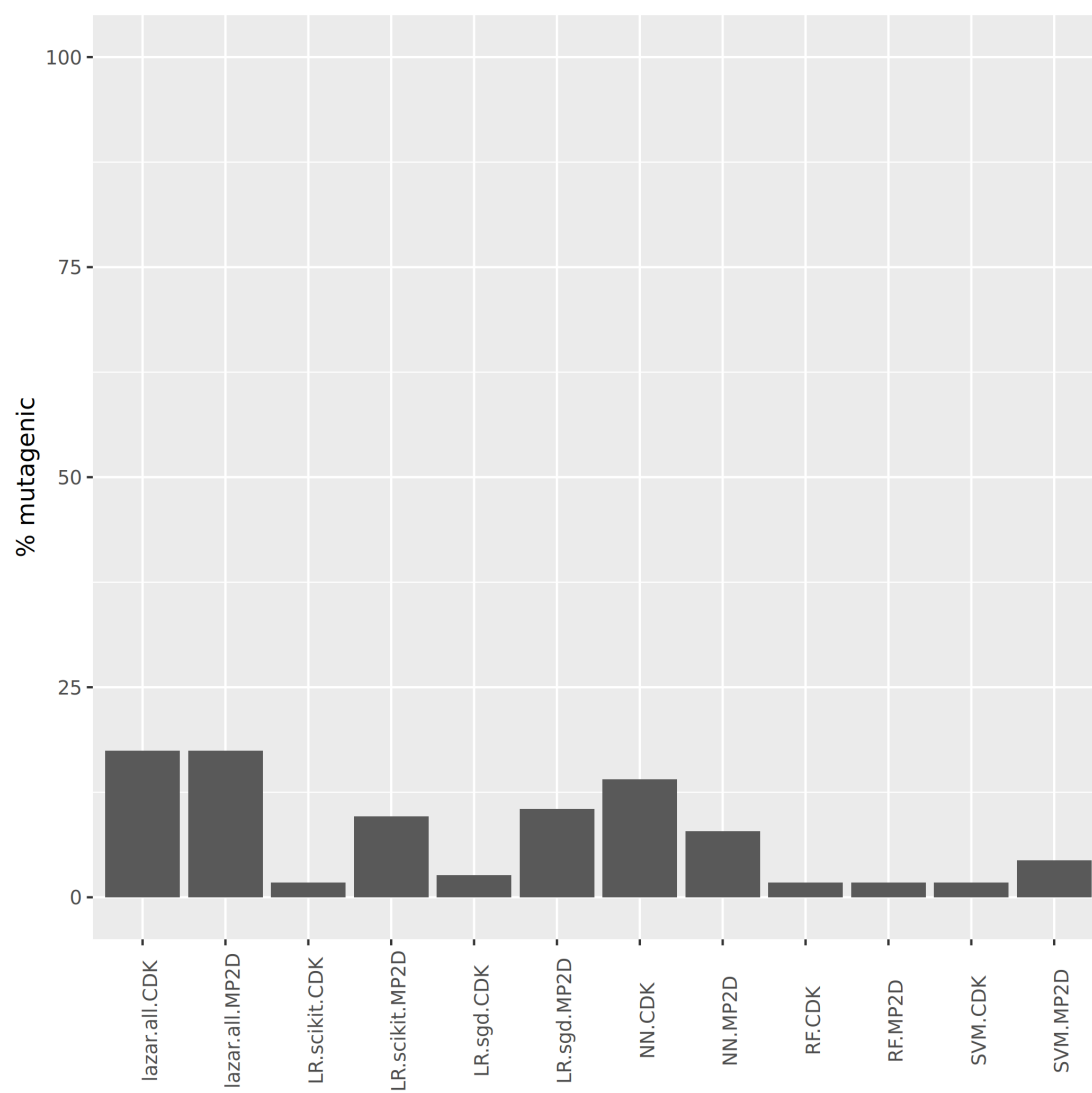


Figure 4: Summary of Diester predictions

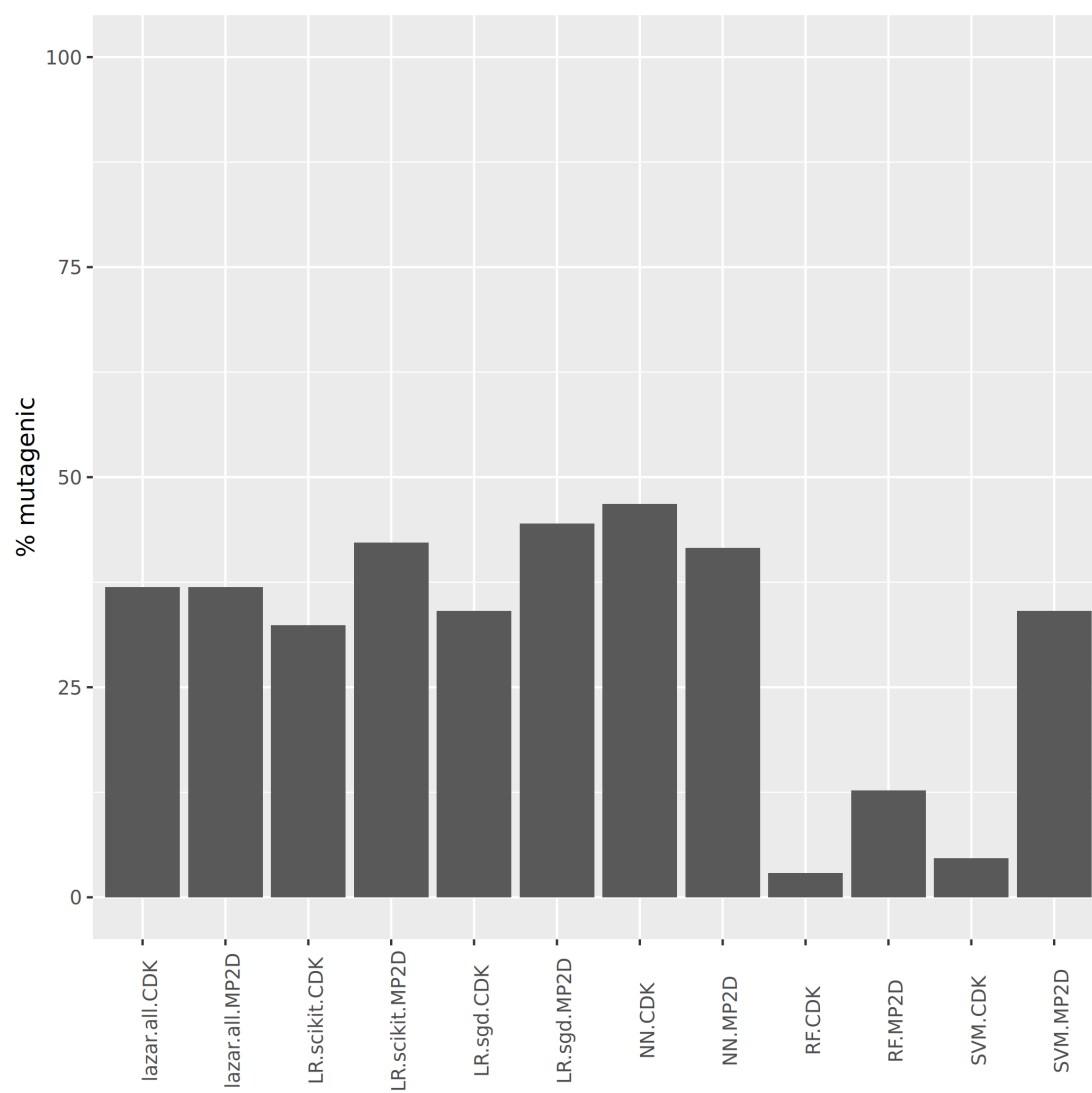


Figure 5: Summary of Macrocytic-diester predictions

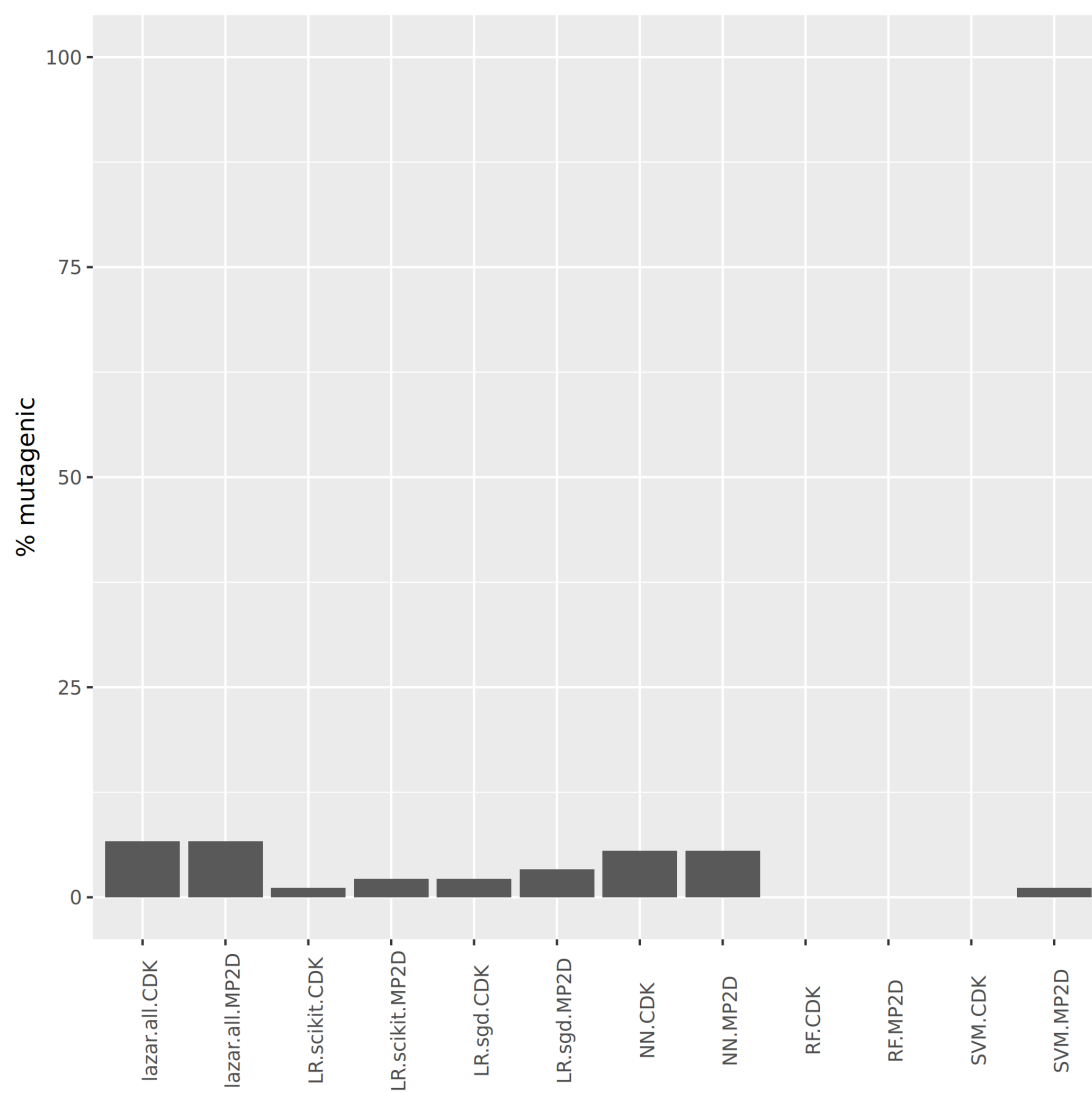


Figure 6: Summary of Monoester predictions

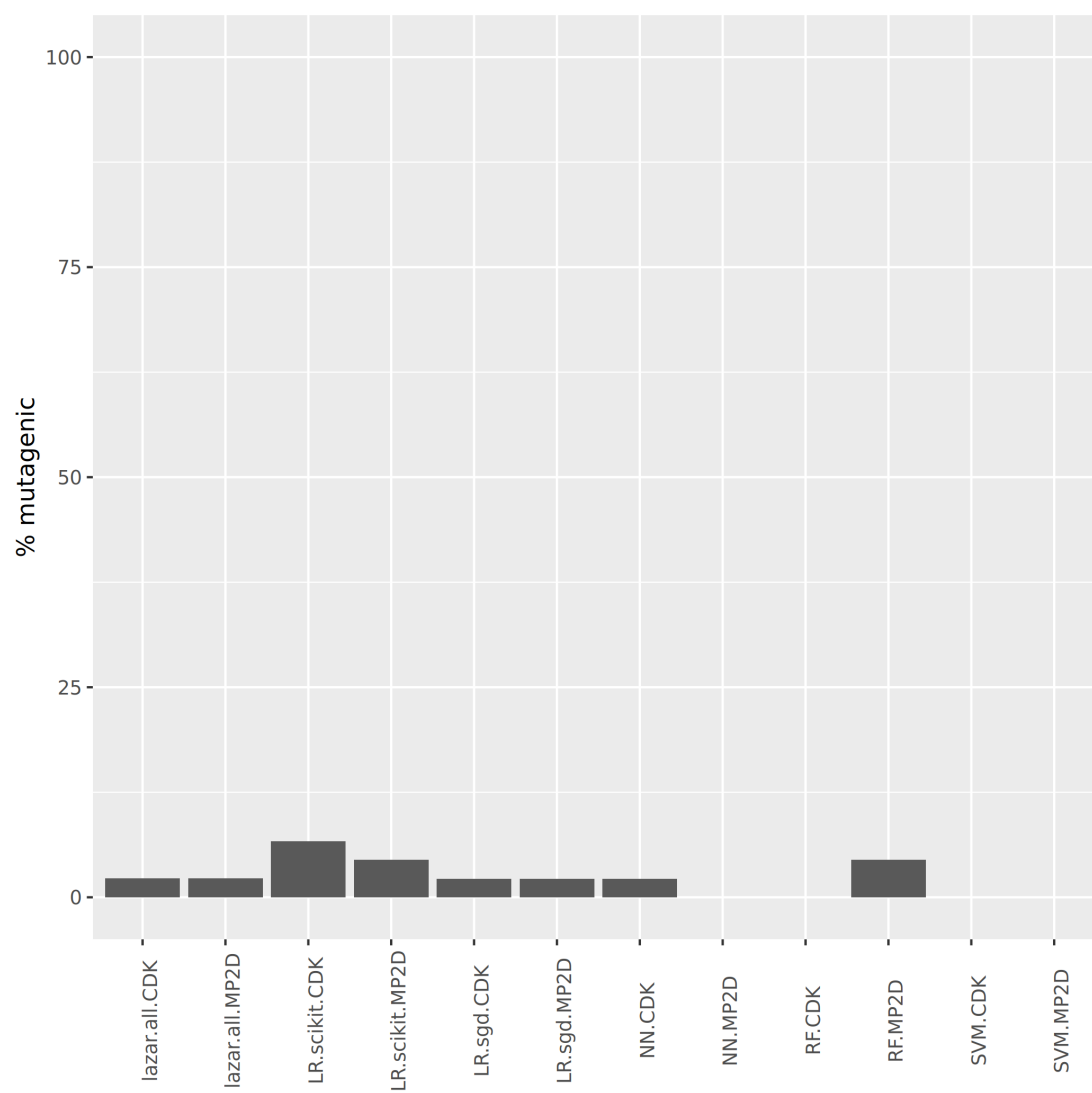


Figure 7: Summary of N-oxide predictions

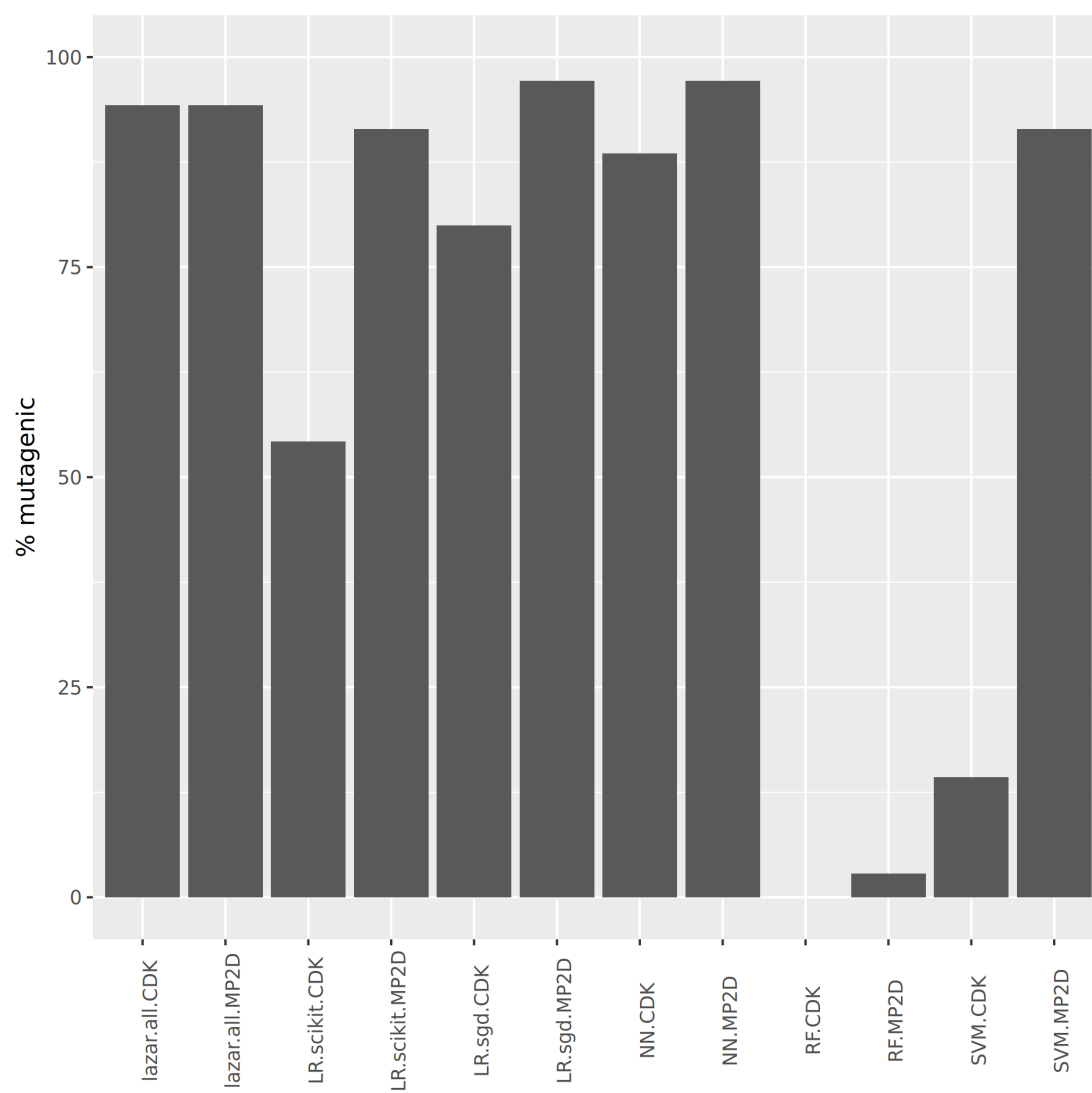


Figure 8: Summary of Otonecine predictions

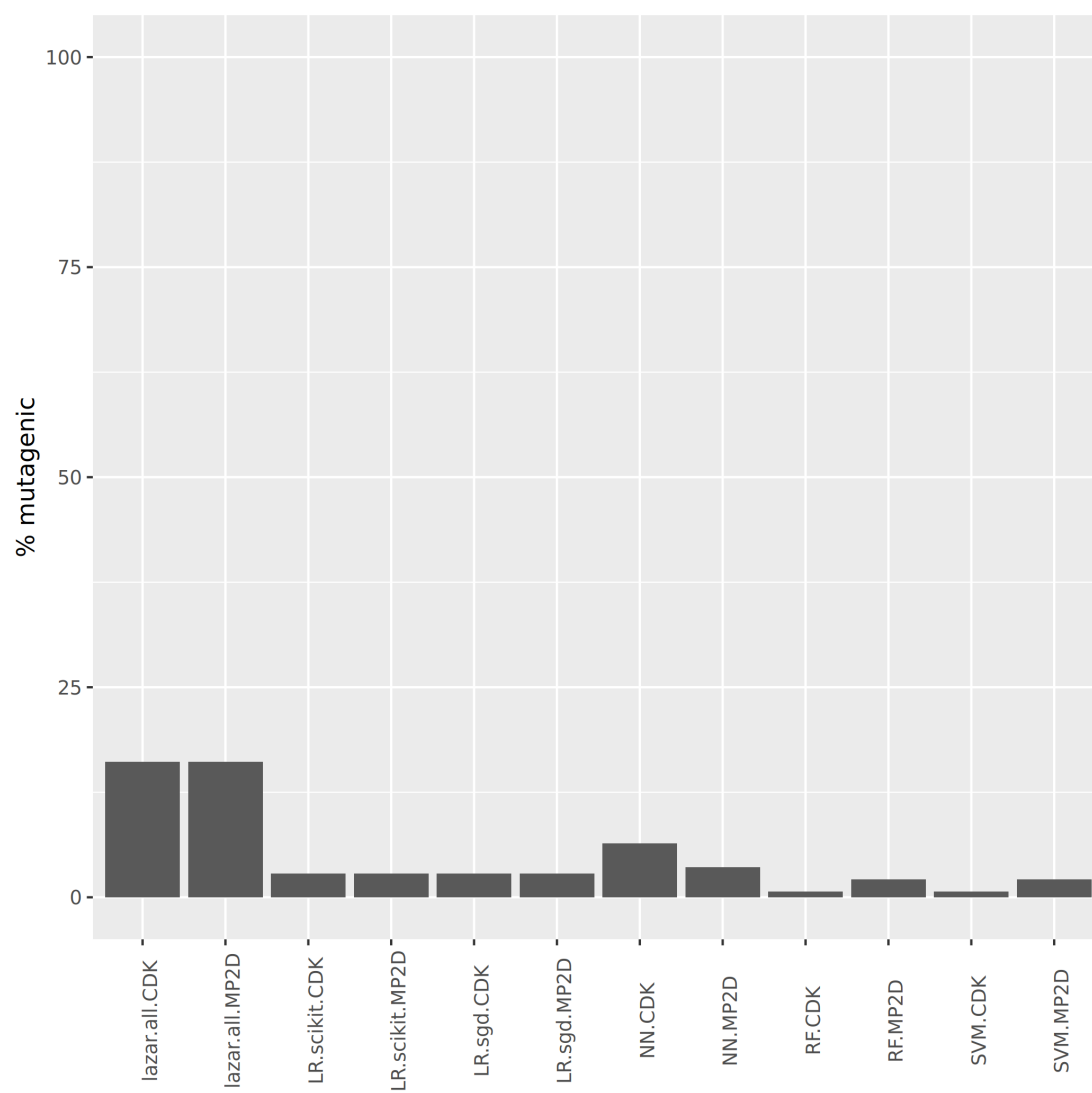


Figure 9: Summary of Platynecine predictions

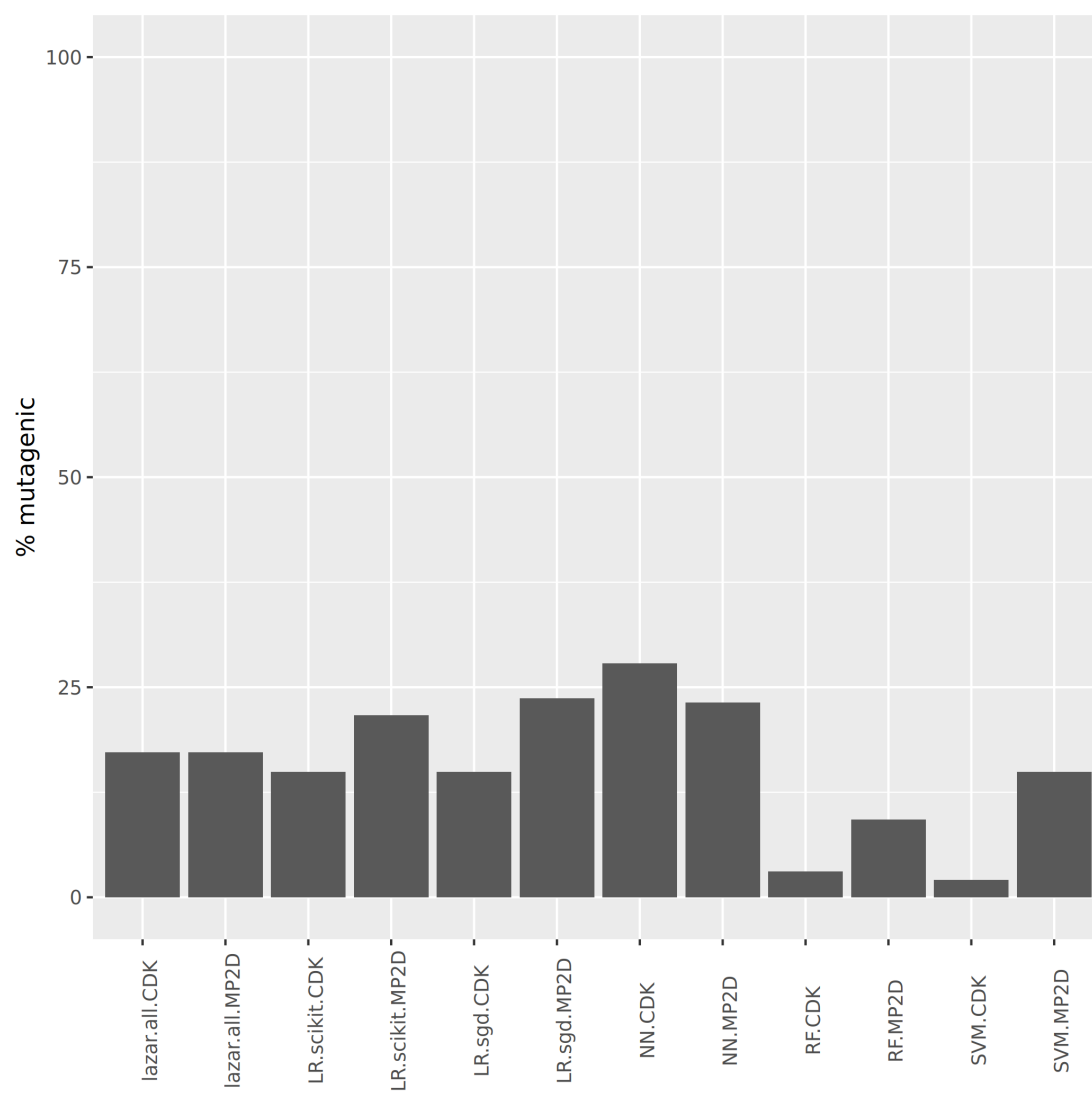


Figure 10: Summary of Retronecine predictions

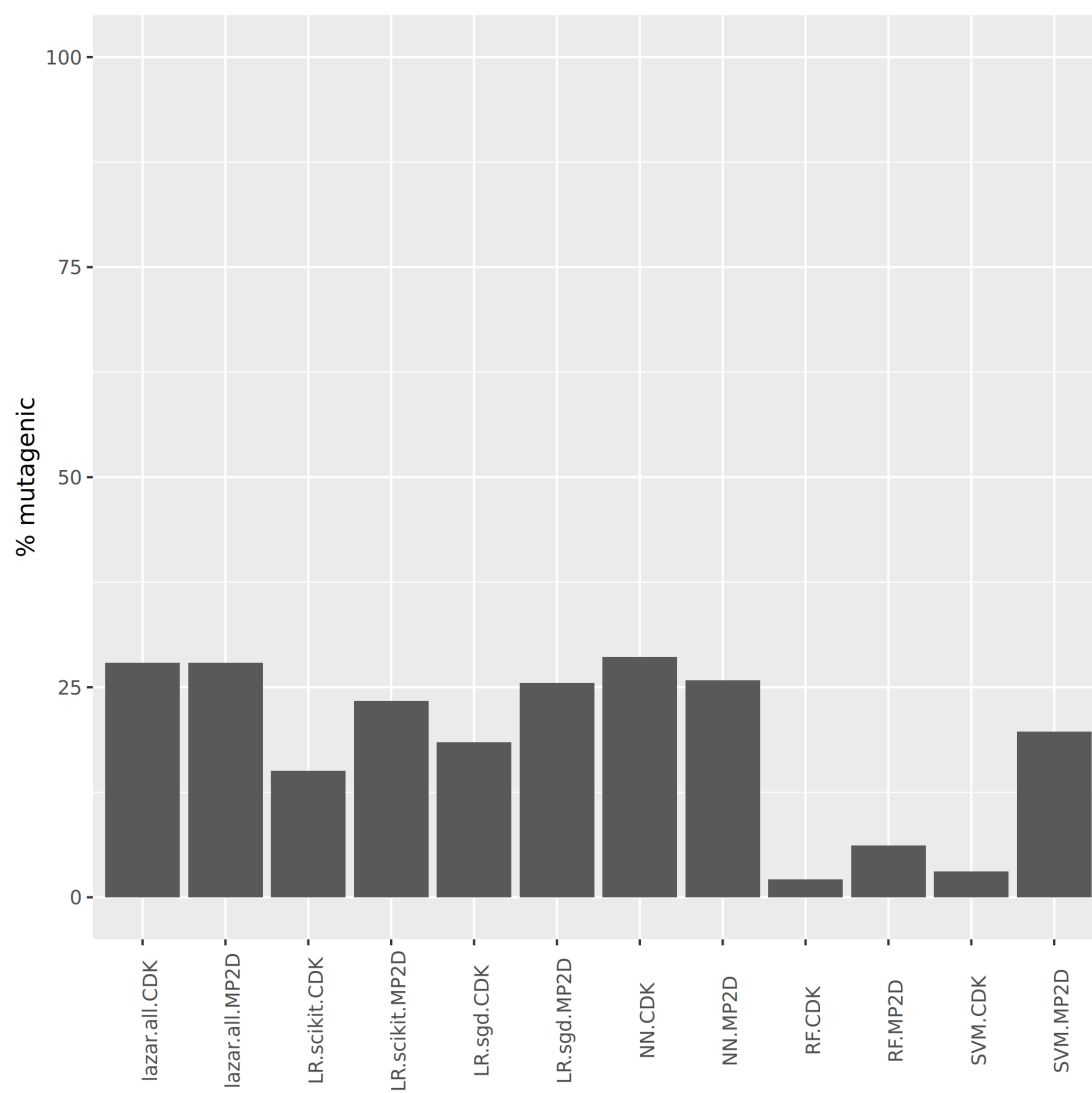


Figure 11: Summary of Tertiary PA predictions



Figure 12: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA)



Figure 13: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA)

ability of the Ames test (80-85% according to Benigni and Giuliani (1988)), especially for predictions with high confidence (%). This is a clear indication that *in-silico* predictions can be as reliable as the bioassays, if the compounds are close to the applicability domain. This conclusion is also supported by our analysis of **lazar** lowest observed effect level predictions, which are also similar to the experimental variability (Helma et al. (2018)).

The lowest number of predictions () has been obtained from **lazar**-CDK high confidence predictions, the largest number of predictions comes from Tensorflow models (). Standard **lazar** give a slightly lower number of predictions () than R and Tensorflow models. This is not necessarily a disadvantage, because **lazar** abstains from predictions, if the query compound is very dissimilar from the compounds in the training set and thus avoids to make predictions for compounds out of the applicability domain.

Descriptors

This study uses two types of descriptors for the characterisation of chemical structures: *MolPrint2D* fingerprints (MP2D, Bender et al. (2004)) use atom environments (i.e. connected atom types for all atoms in a molecule) as molecular representation, which resembles basically the chemical concept of functional groups. MP2D descriptors are used to determine chemical similarities in the default **lazar** settings, and previous experiments have shown, that they give more accurate results than predefined fragments (e.g. MACCS, FP2-4).

In order to investigate, if MP2D fingerprints are also suitable for global models we have tried to build R and Tensorflow models, both with and without unsupervised feature selection. Unfortunately none of the algorithms was capable to deal with the large and sparsely populated descriptor matrix. Based on this result we can conclude, that MolPrint2D descriptors are at the moment unsuitable for standard global machine learning

322 algorithms.

323 **lazar** does not suffer from the size and sparseness problem, because (a) it utilizes inter-
324 nally a much more efficient occurrence based representation and (b) it uses fingerprints
325 only for similarity calculations and not as model parameters.

326 CDK calculates topological and physical-chemical descriptors.

327 **TODO: Verena** kannst Du bitte die Deskriptoren nochmals kurz beschreiben

328 *CDK* descriptors were used for **lazar**, R and Tensorflow models. All models based
329 on CDK descriptors had similar crossvalidation accuracies that were significantly lower
330 than **lazar** MolPrint2D results. Direct comparisons are available only for the **lazar**
331 algorithm, and also in this case CDK accuracies were lower than MolPrint2D accuracies.

332 Based on **lazar** results we can conclude, that CDK descriptors are less suited for chemi-
333 cal similarity calculations than MP2D descriptors. It is also likely that CDK descriptors
334 lead to less accurate predictions for global models, but we cannot draw any definitive
335 conclusion in the absence of MP2D models.

336 Algorithms

337 **lazar** is formally a *k-nearest-neighbor* algorithm that searches for similar structures
338 for a given compound and calculates the prediction based on the experimental data
339 for these structures. The QSAR literature calls such models frequently *local models*,
340 because models are generated specifically for each query compound. R and Tensorflow
341 models are in contrast *global models*, i.e. a single model is used to make predictions
342 for all compounds. It has been postulated in the past, that local models are more
343 accurate, because they can account better for mechanisms, that affect only a subset of
344 the training data. Our results seem to support this assumption, because standard **lazar**
345 models with MolPrint2D descriptors perform better than global models. The accuracy

of **lazar** models with CDK descriptors is however substantially lower and comparable to global models with the same descriptors.

This observation may lead to the conclusion that the choice of suitable descriptors is more important for predictive accuracy than the modelling algorithm, but we were unable to obtain global MP2D models for direct comparisons. The selection of an appropriate modelling algorithm is still crucial, because it needs the capability to handle the descriptor space. Neighbour (and thus similarity) based algorithms like **lazar** have a clear advantage in this respect over global machine learning algorithms (e.g. RF, SVM, LR, NN), because Tanimoto/Jaccard similarities can be calculated efficiently with simple set operations.

Pyrrolizidine alkaloid mutagenicity predictions

lazar models with MolPrint2D descriptors predicted % of the pyrrolizidine alkaloids (PAs) (% with high confidence), the remaining compounds are not within its applicability domain. All other models predicted 100% of the 602 compounds, indicating that all compounds are within their applicability domain.

Mutagenicity predictions from different models show little agreement in general (table 4). 42 from 602 PAs have non-conflicting predictions (all of them non-mutagenic). Most models predict predominantly a non-mutagenic outcome for PAs, with exception of the R deep learning (DL) and the Tensorflow Scikit logistic regression models (and % positive predictions).

R RF and SVM models favor very strongly non-mutagenic predictions (only and % mutagenic PAs), while Tensorflow models classify approximately half of the PAs as mutagenic (RF %, LR-sgd %, LR-scikit:, LR-NN:%). **lazar** models predict predominately non-mutagenicity, but to a lesser extend than R models (MP2D:, CDK:).

It is interesting to note, that different implementations of the same algorithm show little

371 accordance in their prediction (see e.g R-RF vs. Tensorflow-RF and LR-sgd vs. LR-scikit
372 in Table 4 and Table ??).

373 **TODO Verena, Philipp** habt ihr eine Erklaerung dafuer?

374 Figure 12 and Figure ?? show the t-SNE of training data and pyrrolizidine alkaloids. In
375 Figure 12 the PAs are located closely together at the outer border of the training set.
376 In Figure ?? they are less clearly separated and spread over the space occupied by the
377 training examples.

378 This is probably the reason why CDK models predicted all instances and the MP2D
379 model only PAs. Predicting a large number of instances is however not the ultimate
380 goal, we need accurate predictions and an unambiguous estimation of the applicabil-
381 ity domain. With CDK descriptors *all* PAs are within the applicability domain of the
382 training data, which is unlikely despite the size of the training set. MolPrint2D descrip-
383 tors provide a clearer separation, which is also reflected in a better separation between
384 high and low confidence predictions in **lazar** MP2D predictions as compared to **lazar**
385 CDK predictions. Crossvalidation results with substantially higher accuracies for MP2D
386 models than for CDK models also support this argument.

387 Differences between MP2D and CDK descriptors can be explained by their specific prop-
388 erties: CDK calculates a fixed set of descriptors for all structures, while MolPrint2D
389 descriptors resemble substructures that are present in a compound. For this reason
390 there is no fixed number of MP2D descriptors, the descriptor space are all unique sub-
391 structures of the training set. If a query compound contains new substructures, this is
392 immediately reflected in a lower similarity to training compounds, which makes appli-
393 cability domain estimations very straightforward. With CDK (or any other predefined
394 descriptors), the same set of descriptors is calculated for every compound, even if a
395 compound comes from an completely new chemical class.

396 From a practical point we still have to face the question, how to choose model predictions,

if no experimental data is available (we found two PAs in the training data, but this number is too low, to draw any general conclusions). Based on crossvalidation results and the arguments in favor of MolPrint2D descriptors we would put the highest trust in **lazar** MolPrint2D predictions, especially in high-confidence predictions. **lazar** predictions have a accuracy comparable to experimental variability (Helma et al. (2018)) for compounds within the applicability domain. But they should not be trusted blindly. For practical purposes it is important to study the rationales (i.e. neighbors and their experimental activities) for each prediction of relevance. A freely accessible GUI for this purpose has been implemented at <https://lazar.in-silico.ch>.

TODO: Verena Wenn Du **lazar** Ergebnisse konkret diskutieren willst, kann ich Dir ausfuehrliche Vorhersagen (mit aehnlichen Verbindungen und deren Aktivitaet) fuer einzelne Beispiele zusammenstellen

Conclusions

A new public *Salmonella* mutagenicity training dataset with 8309 compounds was created and used it to train **lazar**, R and Tensorflow models with MolPrint2D and CDK descriptors. The best performance was obtained with **lazar** models using MolPrint2D descriptors, with prediction accuracies (%) comparable to the interlaboratory variability of the Ames test (80-85%). Models based on CDK descriptors had lower accuracies than MolPrint2D models, but only the **lazar** algorithm could use MolPrint2D descriptors.

TODO: PA Vorhersagen

References

Bender, Andreas, Hamse Y. Mussa, Robert C. Glen, and Stephan Reiling. 2004. "Molecular Similarity Searching Using Atom Environments, Information-Based Feature Selec-

tion, and a Naïve Bayesian Classifier.” *Journal of Chemical Information and Computer Sciences* 44 (1): 170–78. <https://doi.org/10.1021/ci034207y>.

Benigni, R., and A. Giuliani. 1988. “Computer-assisted Analysis of Interlaboratory Ames Test Variability.” *Journal of Toxicology and Environmental Health* 25 (1): 135–48. <https://doi.org/10.1080/15287398809531194>.

EFSA. 2011. “Scientific Opinion on Pyrrolizidine Alkaloids in Food and Feed.” *EFSA Journal*, no. 9: 1–134.

———. 2016. “Guidance on the Establishment of the Residue Definition for Dietary Assessment: EFSA Panel on Plant Protect Products and Their Residues (PPR).” *EFSA Journal*, no. 14: 1–12.

Hansen, Katja, Sebastian Mika, Timon Schroeter, Andreas Sutter, Antonius ter Laak, Thomas Steger-Hartmann, Nikolaus Heinrich, and Klaus-Robert Müller. 2009. “Benchmark Data Set for in Silico Prediction of Ames Mutagenicity.” *Journal of Chemical Information and Modeling* 49 (9): 2077–81. <https://doi.org/10.1021/ci900161g>.

Helma, Christoph, David Vorgrimmler, Denis Gebele, Martin Gütlein, Barbara Engeli, Jürg Zarn, Benoit Schilter, and Elena Lo Piparo. 2018. “Modeling Chronic Toxicity: A Comparison of Experimental Variability with (Q)SAR/Read-Across Predictions.” *Frontiers in Pharmacology*, no. 9: 413.

Kazius, J., R. McGuire, and R. Bursi. 2005. “Derivation and Validation of Toxicophores for Mutagenicity Prediction.” *J Med Chem*, no. 48: 312–20.

Maaten, L. J. P. van der, and G. E. Hinton. 2008. “Visualizing Data Using T-Sne.” *Journal of Machine Learning Research*, no. 9: 2579–2605.

Mattocks, AR. 1986. *Chemistry and Toxicology of Pyrrolizidine Alkaloids*. Academic Press.

444 O’Boyle, Noel, Michael Banck, Craig James, Chris Morley, Tim Vandermeersch, and
445 Geoffrey Hutchison. 2011. “Open Babel: An open chemical toolbox.” *J. Cheminf.* 3 (1):
446 33. <https://doi.org/doi:10.1186/1758-2946-3-33>.

447 Schöning, Verena, Felix Hammann, Mark Peinl, and Jürgen Drewe. 2017. “Editor’s
448 Highlight: Identification of Any Structure-Specific Hepatotoxic Potential of Different
449 Pyrrolizidine Alkaloids Using Random Forests and Artificial Neural Networks.” *Toxicol.*
450 *Sci.*, no. 160: 361–70.

451 Yap, CW. 2011. “PaDEL-Descriptor: An Open Source Software to Calculate Molecular
452 Descriptors and Fingerprints.” *Journal of Computational Chemistry*, no. 32: 1466–74.