



UNIVERSITY OF LEEDS

This is a repository copy of *Comparison of the Predictive Performance and Interpretability of Random Forest and Linear Models on Benchmark Datasets*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/119210/>

Version: Accepted Version

Article:

Marchese Robinson, RL, Palczewska, A, Palczewski, JA orcid.org/0000-0003-0235-8746 et al. (1 more author) (2017) Comparison of the Predictive Performance and Interpretability of Random Forest and Linear Models on Benchmark Datasets. *Journal of Chemical Information and Modeling*, 57 (8). pp. 1773-1792. ISSN 1549-9596

<https://doi.org/10.1021/acs.jcim.6b00753>

© 2017 American Chemical Society. This document is the Accepted Manuscript version of a Published Work that appeared in final form in *Journal of Chemical Information and Modeling*, copyright © American Chemical Society after peer review and technical editing by the publisher. To access the final edited and published work see <https://doi.org/10.1021/acs.jcim.6b00753>. Uploaded in accordance with the publisher's self-archiving policy.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Comparison of the Predictive Performance and Interpretability of Random Forest and Linear Models on Benchmark Datasets

Richard Liam Marchese Robinson, Anna Palczewska, Jan Palczewski, and Nathan Kidley

J. Chem. Inf. Model., **Just Accepted Manuscript** • DOI: 10.1021/acs.jcim.6b00753 • Publication Date (Web): 17 Jul 2017

Downloaded from <http://pubs.acs.org> on July 18, 2017

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

Comparison of the Predictive Performance and Interpretability of Random Forest and Linear Models on Benchmark Datasets

Richard L. Marchese Robinson,^{a,b,c,} Anna Palczewska,^{d,e} Jan Palczewski,^f Nathan Kidley^a*

- a. Syngenta Ltd, Jealott's Hill International Research Centre, Bracknell, Berkshire, RG42 6EY, United Kingdom
- b. School of Pharmacy and Biomolecular Sciences, Liverpool John Moores University, James Parsons Building, Byrom Street, Liverpool, L3 3AF, United Kingdom
- c. Present Address: School of Chemical and Process Engineering, University of Leeds, Leeds LS2 9JT, United Kingdom
- d. Department of Computing, University of Bradford, BD7 1DP Bradford, United Kingdom
- e. Present Address: School of Geography, Leeds Institute for Data Analytics, University of Leeds, Leeds LS2 9JT, United Kingdom
- f. School of Mathematics, University of Leeds, Leeds LS2 9JT, United Kingdom

*Corresponding author. E-mail: R.L.MarcheseRobinson@leeds.ac.uk

KEYWORDS

quantitative structure-activity relationships; model interpretation; Machine Learning; Heat Map; Random Forest; Partial Least Squares; Support Vector Machines; Support Vector Regression

ABSTRACT

The ability to interpret the predictions made by quantitative structure activity relationships (QSARs) offers a number of advantages. Whilst QSARs built using non-linear modelling approaches, such as the popular Random Forest algorithm, might sometimes be more predictive than those built using linear modelling approaches, their predictions have been perceived as difficult to interpret. However, a growing number of approaches have been proposed for interpreting non-linear QSAR models in general and Random Forest in particular. In the current work, we compare the performance of Random Forest to two widely used linear modelling approaches: linear Support Vector Machines (SVM), or Support Vector Regression (SVR), and Partial Least Squares (PLS). We compare their performance in terms of their predictivity as well as the chemical interpretability of the predictions, using novel scoring schemes for assessing Heat Map images of substructural contributions. We critically assess different approaches to interpreting Random Forest models as well as for obtaining predictions from the forest. We assess the models on a large number of widely employed, public domain benchmark datasets corresponding to regression and binary classification problems of relevance to hit identification and toxicology. We conclude that Random Forest typically yields comparable or possibly better predictive performance than the linear modelling approaches and that its predictions may also be interpreted in a chemically and biologically meaningful way. In contrast to earlier work looking at interpreting non-linear QSAR models, we directly compare two methodologically distinct approaches for interpreting Random Forest models. The approaches for interpreting Random

1
2
3 Forest assessed in our article were implemented using Open Source programs, which we have
4 made available to the community. These programs are the *rfFC* package [[https://r-forge.r-](https://r-forge.r-project.org/R/?group_id=1725)
5 project.org/R/?group_id=1725] for the R Statistical Programming Language, along with a
6 Python program *HeatMapWrapper* [<https://doi.org/10.5281/zenodo.495163>] for Heat Map
7 generation.
8
9
10
11
12
13
14
15
16

17 INTRODUCTION

18
19 The ability to interpret the predictions made by quantitative structure-activity relationships
20 (QSARs), in terms of the size and sign of the influence of the underlying descriptor values, is
21 valuable for a number of reasons. Mechanistic interpretation is viewed as particularly important
22 in a regulatory context¹ and could yield novel mechanistic insight or consistency with established
23 understanding could complement statistical validation, by conferring greater confidence in the
24 predictions.² Furthermore, interpretable predictions could guide structural modifications required
25 for desired changes in molecular properties.² Analysis of descriptor contributions may also help
26 to rationalize outliers.³
27
28
29
30
31
32
33
34
35
36
37
38

39 Traditionally, QSARs were linear models.^{2,4} For linear models, both the relative size and sign of
40 the influence of an individual descriptor on the predictions is straightforward to obtain in terms
41 of the corresponding coefficient.^{4,5} The product of the coefficient and descriptor value for a
42 given chemical will yield the component of an individual prediction arising from a particular
43 descriptor.⁶ In principle, although complications may arise in practice,^{6,7} this makes the
44 predictions obtained from linear QSARs relatively straightforward to interpret.
45
46
47
48
49
50
51
52
53

54 A number of algorithms for building linear QSARs, either for prediction of continuous
55 (regression models) or categorical (classification models) measures of biological activities, are
56
57
58
59
60

1
2
3 used by modern QSAR practitioners. These include partial least squares (PLS),⁸ for regression
4 or classification⁹ as well as, when used with a linear kernel function, Support Vector Machines
5 (SVMs)^{10,11} and Support Vector Regression (SVR)¹² for classification¹³ and regression¹⁴
6 respectively.
7
8
9

10
11
12
13 However, in recent decades, QSAR models have been increasingly developed using non-linear
14 modelling approaches.^{2,15} Interpreting the predictions made by these models is typically harder
15 than for linear models. Indeed, it has been argued that there is an inherent trade-off between
16 predictive performance and interpretability, with non-linear models suggested to generally be
17 more predictive but less interpretable.¹⁶ Nonetheless, in addition to approaches for estimating the
18 overall “importance” of various descriptors in a QSAR model,¹⁶ a growing number of studies has
19 advocated approaches which evaluate how specific molecular characteristics contribute towards
20 individual predictions made by a non-linear QSAR model.^{2,17-25}
21
22
23
24
25
26
27
28
29
30
31
32

33
34 One non-linear modelling approach which is popular within the QSAR community^{15,16} is
35 Random Forest.^{9,26} Approaches for estimating descriptor “importance” for Random Forest
36 models are well established.^{9,26-28} However, more recently, proposals for estimating the relative
37 magnitude and sign of the influence of a given descriptor on a given prediction, as is the focus of
38 our article, have been presented,^{2,18,19,23,29-31} including those which are specifically designed for
39 Random Forest models.^{19,23,29-31}
40
41
42
43
44
45
46
47

48
49 In our current article, we consider QSAR models developed using Random Forest, PLS and
50 SVM (or Support Vector Regression), using a linear kernel, for binary classification and
51 regression modelling tasks which are relevant for hit identification or predictive toxicology. We
52 compare the predictive performance of these algorithms, as well as the extent to which their
53
54
55
56
57
58
59
60

1
2
3 predictions for individual molecules can be interpreted in a chemically and biologically
4 meaningful manner according to various established approaches for model interpretation, using
5 well-known, publicly available benchmark datasets. As well as evaluating different algorithms
6 for interpreting individual predictions, we further investigate the degree to which obtaining a
7 chemically and biologically meaningful interpretation depends upon the quality of the
8 corresponding prediction for a given molecule.
9

10
11 Specifically, we compare two different approaches for interpreting the same Random Forest
12 prediction. The first approach^{19,23,29–32} was originally proposed for Random Forest regression
13 models¹⁹ and was then extended to binary classification models.^{29,30} The second approach is
14 based on estimating local gradients (i.e. vectors of partial derivatives) of the prediction with
15 respect to the descriptors and is analogous to algorithms recently used to interpret various kinds
16 of non-linear QSAR models.^{2,18,20,33} We evaluate these approaches via means of novel scoring
17 schemes for assessing Heat Map images representing the influence of molecular substructures on
18 model predictions. We extend the Heat Map approach introduced by Rosenbaum et al.¹³ to work
19 with any linear or non-linear classification or regression model, built using the same descriptors
20 as per their earlier work, and present an Open Source implementation of our extension. Our
21 extension of their approach includes novel “symmetrized” normalization schemes for translating
22 descriptor contributions associated with substructural features into Heat Map molecular images.
23 We present the first evaluation of our approach for interpreting binary classification Random
24 Forest models^{29,30} in cheminformatics, using suitable benchmark classification datasets. To the
25 best of our knowledge, we present the first comparative assessment of methodologically distinct
26 approaches to interpreting predictions obtained from a non-linear QSAR model.
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

METHODS AND DATA

Modelling approaches

More expansive descriptions of these modelling algorithms and, for those algorithms yielding linear models, the equations generating the predictions, are provided in the Supporting Information.

Partial least squares (PLS) regression^{8,34} develops a linear regression model based upon so-called “latent variables”, which correspond to a selected set of orthogonal, linear combinations of the original descriptors. A simple modification can transform PLS regression into a binary classifier: the two class labels are mapped onto two real numbers (1 and 0), PLS regression is used to build a model from the training set and the numeric predictions (or score) made by the regression model are converted into predicted class labels via comparison to some suitable threshold.⁹ This so-called “probit” treatment yields *Probit-PLS*.³⁵ Note that this model differs from a classical probit regression, for which a random outcome is observed with a probability described by the model.³⁶

Support Vector Machines (SVMs)^{10,11} build a binary classification model with a functional form which is contingent upon the selected “kernel function”. *Support Vector Regression (SVR)*^{12,14} is an analogous approach to SVM. Using a linear “kernel function”^{11,13} yields linear SVM and SVR models.

Random Forest classification^{9,26} builds classification models, meaning binary classification models in our current article, via generating a forest of decision trees using the training data. Each tree is grown using independent random samples of the training data. Whilst “bootstrap” samples are commonly employed,^{9,15,26} we employed sampling without replacement. This was

1
2
3 required in order to be able to calculate descriptor contributions via the “Kuz'min/Palczewska”
4
5 approach (see below) using the *rfFC* software package.³⁷
6
7

8
9 By default, binary classification predictions are made for a new molecule by passing it down
10
11 each of the decision trees in the forest, each tree calculating a score corresponding to the fraction
12
13 of molecules belonging to “class 1” in the terminal node to which the new molecule is assigned.
14
15 If that fraction is greater than 0.5, i.e. if a majority of molecules belong to “class 1”, the tree
16
17 predicts the new molecule to belong to “class 1”. The forest then returns an overall score
18
19 corresponding to the proportion of trees predicting the molecule to belong to “class 1”. If this
20
21 overall score is greater (not greater) than the threshold of 0.5, i.e. if a majority of trees (do not)
22
23 vote for “class 1”, the predicted class is “class 1” (“class 0”). We also considered a non-default
24
25 “averaged predictions” approach, which is analogous to Random Forest regression, where the
26
27 overall score is calculated as the arithmetic mean of the scores generated by each tree and
28
29 classification is based on the same threshold of 0.5.
30
31
32
33
34
35

36 *Random Forest regression*^{9,26} model development is virtually identical to Random Forest
37
38 classification. The bioactivity predicted for a molecule is the arithmetic mean of the bioactivity
39
40 of all molecules assigned to the relevant terminal node, for a given tree, averaged across all trees
41
42 in the forest.
43
44
45

46 **Hyperparameters selection**

47
48 A variety of “hyperparameters” can affect the performance of the modelling algorithms used for
49
50 our current article. In addition, the degree of random “downsampling”,³⁸ to address class
51
52 imbalance, may be considered another hyperparameter for the binary classification datasets.
53
54 Hence, a variety of hyperparameter combinations were evaluated, which are summarized in
55
56
57
58
59
60

1
2
3 Supporting Information Table S1. The best hyperparameter options were selected via iterative
4 application of stratified, 2-fold cross-validation on the training set.³⁹ Different combinations of
5 hyperparameters were evaluated in terms of the cross-validated Matthews Correlation
6 Coefficient (MCC)^{40,41} or the Coefficient of Determination (R^2)^{42,43} for a binary classification or
7 regression modelling approach respectively. Further details are provided in the Supporting
8 Information.

18 **Descriptors employed**

20 The descriptors employed in our current work correspond to a set of “features” representing
21 molecular substructures. The corresponding descriptor takes the value 1, or 0, if a feature
22 encoded substructure is present, or not present, in a particular molecule. Specifically, the widely
23 used extended connectivity fingerprint (ECFP)⁴⁴ features were employed. Here, we calculated
24 these features as per Rosenbaum et al.¹³ In brief, for each molecule, features are grown to
25 represent circular substructures centered on each of the heavy atoms in the molecule. For each
26 central atom, an initial identifier was assigned based on the Daylight invariant information
27 (including atomic number and number of connections)⁴⁵ and ring type information. The
28 assignment of this initial identifier is considered iteration number 0. In subsequent iterations, this
29 initial identifier is supplemented with the identifiers of attached heavy atoms. Features are
30 generated at each iteration by hashing the combinations of identifiers. The unique features,
31 corresponding to molecular substructures of increasing size, obtained at each iteration are saved.
32 This procedure stops when the maximum number of iterations is reached, which is controlled by
33 the depth of the fingerprint: the longest path between atoms in the fingerprint. As per Rosenbaum
34 et al., we used a depth of 4.¹³

Feature selection

Prior to selecting the most appropriate hyperparameters and building a final model on a given training set, the training set was used to select the most biologically relevant 200 features.

Further details are provided in the Supporting Information.

Benchmark datasets

The datasets used in the current work were the protein binding regression datasets of Sutherland et al.,⁴⁶ along with several binary classification datasets containing molecules assessed with respect to various endpoints: protein binding (MUV846, MUV548, MUV713 datasets),⁴⁷ Ames test (Kazius dataset)⁴⁸ and Chromosome aberration test (CA dataset)¹⁷ toxicity tests. (A summary of these datasets is provided in **Table 1**.) Whilst possible limitations of the Kazius⁴⁸ dataset have been noted and expanded Ames datasets described in the literature,^{22,49} the Kazius,⁴⁸ Sutherland⁴⁶ and MUV⁴⁷ datasets are accepted QSAR benchmarks^{13,50–53} and the Kazius, MUV and CA datasets have previously been used to illustrate approaches for interpreting QSAR predictions.^{13,17}

Table 1. Benchmark datasets

Type	Dataset Name	Endpoint	Size ^a	Unique ECFP features before feature selection	Reference
Binary classification	MUV846	Coagulation factor Xia (protease) inhibition	15030 (30 actives, 15000 decoys)	175146	Rohrer and Baumann ⁴⁷
	MUV548	Protein kinase A	15030 (30 actives,	154311	

		inhibition	15000 decoys)		
	MUV713	estrogen receptor alpha coactivator binding inhibition	15030 (30 actives, 15000 decoys)	150857	
	Kazius	Mutagenicity (Ames test)	4337 (2401 mutagens, 1936 non-mutagens)	46282	Kazius et al. ⁴⁸
	CA	Chromosome aberration test	939 (351 positive, 588 negative)	15846	Mohr et al. ¹⁷
Regression	DHFR	dihydrofolate reductase inhibition (pIC ₅₀ values)	397	4792	Sutherland et al. ⁴⁶
	COX2	cyclooxygenase-2 inhibition (pIC ₅₀ values)	322	3434	
	BZR	benzodiazepine receptor activity (pIC ₅₀ values)	163	2755	
	AchE	acetylcholinesterase inhibition (pIC ₅₀ values)	111	1587	
	ACE	angiotensin converting enzyme inhibition	114	1637	Sutherland et al. & Priest et

		(pIC ₅₀ values)			al. ^{46,54}
	THR	thrombin inhibition (pK _i values)	88	1084	Sutherland et al. & Böhm et al. ^{46,55}
	THERM	thermolysin inhibition (pK _i values)	76	1017	Sutherland et al. & Gohlke
	GPB	glycogen phosphorylase b inhibition (pK _i values)	66	817	et al. ^{46,56}

- a. One CA dataset¹⁷ entry (“CA847”) was automatically removed during the preparation of datasets for calculations.

For all MUV datasets, the actives were treated as “class 1” and the decoys were “class 0”. For the Kazius dataset, mutagens were “class 1” and non-mutagens were “class 0”. For the CA dataset, positives, i.e. compounds determined to cause chromosome aberrations based on the chromosome aberration test,¹⁷ were “class 1” and negatives were “class 0”. Not all the MUV datasets⁴⁷ were modelled, due to the considerable computational overhead.

Preparation of datasets for calculations

All datasets were initially retrieved as SDF files and the structures processed using Pipeline Pilot⁵⁷ prior to calculating ECFP features. This was designed to generate neutral, single molecule representations, in standard tautomeric forms, with consistently represented functional groups

1
2
3 including explicit hydrogens, as well as remove any obviously erroneous structures. Further
4
5 details are provided in the Supporting Information.
6
7

8 9 **Statistical evaluation of predictive performance**

10
11 The performance of all modelling approaches was assessed via R repetitions of K -fold “external
12 cross-validation”.⁵⁸ Feature selection was carried out, independently, using the entirety of each
13 cross-validation training set and hyperparameters selection was carried out using “internal”⁵⁸
14 cross-validation on those training sets. This avoided optimistically biasing the results obtained on
15 the “external cross-validation” test sets due to their having been used for model optimization.
16
17
18
19
20
21
22
23

24 For all datasets except the CA dataset, for which the single set of 10 folds defined by Mohr et al.
25 was used,¹⁷ the partitioning into folds was carried out randomly and independently for each
26 repetition, with stratified sampling employed for the classification datasets. Two repetitions of
27 five-fold cross-validation were used for the Kazius⁴⁸ and MUV⁴⁷ datasets, as per Rosenbaum et
28 al.,¹³ whilst the performance of the regression approaches on the Sutherland datasets⁴⁶ was
29 evaluated using 20 repetitions of 10-fold cross-validation as per Hinselmann et al.⁵⁰
30
31
32
33
34
35
36
37
38

39 Finally, it should be recalled that all model development protocols were dependent upon random
40 selections. In the case of Random Forest models, this is intrinsic to the algorithm. In the case of
41 binary classification models, this is true due to the random selection of subsets of the majority
42 class when downsampling. In the case of all models, except for Random Forest regression
43 models, this is further due to the random selection of internal cross-validation folds used to guide
44 hyperparameter selection (see “Hyperparameters selection”). Hence, as well as averaging test set
45 results across multiple repetitions of K -fold “external cross-validation”, performance statistics
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 were averaged across five repetitions of the entire modelling and evaluation procedure: each
4
5 repetition employed a different random number generator (RNG) seed for all random operations.
6
7

8
9 For binary classification datasets, the overall predictive performance of the models was
10
11 evaluated in terms of Matthew's correlation coefficient (*MCC*),^{40,41} Cohen's Kappa (*Kappa*),^{35,59}
12
13 the area under the ROC curve (*AUC*),^{60,61} and balanced accuracy (*BA*).⁶² For regression datasets,
14
15 the overall predictive performance of the models was evaluated in terms of the coefficient of
16
17 determination (R^2), Pearson's correlation coefficient (r) and root mean squared error
18
19 (*RMSE*).^{8,40,42,43} These were calculated, and may be interpreted, as explained in the Supporting
20
21 Information.
22
23
24

25 26 27 **Descriptor contributions**

28
29 By "descriptor contributions", we mean a vector of estimates of signed influences of each
30
31 descriptor value upon the prediction made by the model for a given molecule. Here, we first
32
33 provide an overview of the different kinds of descriptor contributions which were calculated, and
34
35 their interpretation, followed by a detailed explanation of how these different kinds of descriptor
36
37 contributions were calculated.
38
39

40
41 For the linear modelling approaches (PLS regression, Probit-PLS, Support Vector Machines and
42
43 Support Vector Regression with linear kernel functions), the coefficients were treated as
44
45 descriptor contributions. For both Random Forest regression and Random Forest classification,
46
47 two different kinds of methodology were employed to estimate descriptor contributions. The first
48
49 approach was developed by Kuz'min et al.¹⁹ for Random Forest regression and extended by
50
51 Palczewska et al. to binary classification.^{29,30} The second approach was to estimate the local
52
53 gradient, i.e. a vector of partial derivatives of the predicted value with respect to each of the
54
55
56
57
58
59
60

1
2
3 descriptors. The local gradients approach is analogous to the approaches presented in a number
4
5 of recent publications for interpreting non-linear QSARs.^{2,18,20,33} N.B. For brevity, we refer to the
6
7 first and second approach as the “Kuz’min/Palczewska” approach and the “local gradients”
8
9 approach respectively.
10

11
12
13 Whilst this would not generally be the case, all analyzed descriptor contributions presented in
14
15 our current work are estimates of the extent to which the corresponding molecular substructures
16
17 promoted biological activity. This is because we employed binary descriptors for which the
18
19 descriptors took values of 1 or 0, denoting the presence or absence of a corresponding
20
21 substructure, and we only analyzed the descriptor contributions for molecules where the
22
23 corresponding substructures were present. (See sections “Descriptors employed” and “Heat Map
24
25 interpretation of individual predictions” for details.) However, for the linear models, the
26
27 significance of given substructures in specific, local contexts (e.g. specific chemical classes), as
28
29 opposed to their global significance, is not estimated.
30
31
32
33
34
35

36 For Random Forest classification, only the descriptor contributions calculated based on the
37
38 “local gradients approach” were dependent upon the whether the default (majority voting) or
39
40 “averaged predictions” approach was employed. However, as explained under
41
42 “Kuz’min/Palczewska approach” below, only the predictions made using the “averaged
43
44 predictions” approach can be guaranteed to be consistent with the descriptor contributions
45
46 calculated according to the “Kuz’min/Palczewska” approach. Hence, we only consider the
47
48 “averaged predictions” when evaluating the interpretability of Random Forest classification
49
50 predictions based on the “Kuz’min/Palczewska” approach.
51
52
53
54

55
56 *Linear coefficients approach*
57
58
59
60

1
2
3 The output generated for all linear modelling approaches can be expressed using a common
4 functional form. This allows coefficients to be obtained for the descriptors. The output generated
5 by each model, for any given molecule, is either the predicted bioactivity, for regression models,
6 or the score used, in combination with a threshold, to make predictions for classification models.
7
8 The common functional form is presented in equation (1): o_i is the model output for the i th
9 molecule, d_{ji} is the value of the j th descriptor (out of J descriptors in total) for the i th molecule,
10 c_j is the coefficient for the j th descriptor and b is an offset constant. N.B. Consideration of
11 Supporting Information equations (i)-(vi) demonstrates that the values of c_j in equation (1) are
12 calculated from different model-specific parameters.
13
14
15
16
17
18
19
20
21
22
23
24
25

$$o_i = \sum_j^J c_j d_{ji} + b \quad (1)$$

26
27
28
29
30 In the current work, these coefficients (c_j) were treated as “descriptor contributions” for all linear
31 models. In general, treating the coefficients as “descriptor contributions” does not quantify the
32 total component of a given prediction arising from a given descriptor, which is given by the
33 product of the descriptor and coefficient values (i.e. $c_j d_{ji}$).⁶ Nonetheless, in the current work, the
34 “descriptor contributions” analyzed for linear models were directly equivalent to the components
35 of a given prediction arising from individual descriptors. (This is because, as described under
36 “Descriptors employed”, we employed binary descriptors, for which $d_{ji} = 1$ or 0 , and we only
37 analyzed the descriptor contributions where $d_{ji} = 1$ – corresponding to the presence of molecular
38 substructures, as explained under “Heat Map interpretation of individual predictions”.)
39
40 Specifically, they corresponded to an estimate of the extent to which the corresponding
41 molecular substructures promoted biological activity. Since the coefficients calculated for linear
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 models are constant for all molecules, the significance of given substructures in specific, local
4 contexts (e.g. specific chemical classes), as opposed to their global significance, is not estimated.
5
6
7

8 9 *Kuz'min/Palczewska approach*

10
11
12 This approach starts with a trained Random Forest model for which the training set molecules,
13 and their measured bioactivities, assigned to each node are known. Subsequent to treating the
14 class labels as numeric bioactivity values (i.e. “class 1” is treated as 1 and “class 0” is treated as
15 0), the approaches considered in our current article are identical for both binary classification and
16 regression i.e. binary classification “descriptor contributions” are calculated towards “class 1”,
17 such that negative values indicate contributions towards “class 0”. The steps required to obtain
18 the descriptor contribution for a single descriptor and for a single molecule (X) are as follows.
19
20
21
22
23
24
25
26
27
28
29

- 30 1. For a single tree in the forest, start at the “root node” corresponding to the complete
31 “local training set”: this “root node” is the first “parent node”.
32
33
34
35
- 36 2. Obtain the arithmetic mean bioactivity for all training set molecules in the current “parent
37 node”. For binary classification, this is equivalent to calculating the fraction of “local
38 training set” molecules belonging to “class 1” (see the preceding discussion).
39
40
41
42
43
- 44 3. Go to the “child node” into which the molecule X is directed based on its value for the
45 descriptor used to partition the “parent node” compared to the “split value” (see
46 “Modelling approaches”).
47
48
49
50
51
- 52 4. *If* the descriptor used to split the “parent node” corresponds to the descriptor of interest,
53 obtain the corresponding “local increment”, as per equation (2). In equation (2), $LI_j^{C,T}$
54 denotes the “local increment” associated with the j th descriptor calculated for the *relevant*
55
56
57
58
59
60

“child node” C in tree T ; $Y_{mean}^{C,T}$ ($Y_{mean}^{P,T}$) corresponds to the arithmetic mean bioactivity of the “local training set” molecules assigned to the “child node” C (“parent node” P) in tree T .

$$LI_j^{C,T} = Y_{mean}^{C,T} - Y_{mean}^{P,T} \quad (2)$$

5. Repeat steps [2]-[4], with each “child node” from the previous iteration becoming the “parent node” in the next iteration, until a “terminal node” is reached.
6. Repeat steps [2]-[5] for all trees in the forest.
7. Calculate the “descriptor contribution” from the “local increments” as per equation (3):
 dc_j^X denotes the descriptor contribution obtained for the j th descriptor for molecule X ; T denotes a particular tree in the forest; n_{tree} denotes the total number of trees in the forest; C denotes a particular “child node” in a given tree. N.B. The set of “child nodes” for which “local increments” ($LI_j^{C,T}$) are summed in equation (3) will be, as explained for step [4], contingent upon the specific molecule and descriptor for which a “descriptor contribution” is calculated.

$$dc_j^X = \frac{1}{n_{tree}} \sum_T^{n_{tree}} \sum_C LI_j^{C,T} \quad (3)$$

It should be noted that these “descriptor contributions” are equivalent to the descriptor specific components of the prediction generated by a Random Forest regression model. This is also true for the score for “class 1” generated by a binary classification model under either of the following scenarios. The first scenario is that the “averaged predictions” predictions approach (see “Modelling approaches”) of averaging the proportion of terminal node “local training set

1
2
3 molecules” in “class 1”, across all trees in the forest, is employed. The second scenario is that the
4
5 default approach of each tree predicting the terminal node majority class, with the model
6
7 returning the proportion of trees voting for “class 1”, is employed and, for each tree, all the
8
9 terminal node molecules belong to the same class. (This default approach is equivalent to
10
11 predicting the class with the majority of votes once this proportion is compared to the threshold
12
13 of 0.50.) Indeed, under the second scenario, the “averaged predictions” and default predictions
14
15 approaches are identical. Under these circumstances, the relationship between the Random
16
17 Forest model output and the descriptor contributions is given in equation (4).^{19,29,30} In equation
18
19 (4), o_X is the model output for molecule X (either the predicted bioactivity for a regression model
20
21 or the score for “class 1” for a binary classification model); dc_j^X denotes the descriptor
22
23 contribution obtained for the j th descriptor for molecule X ; A_0 is a constant, for a given model.
24
25
26
27
28
29
30

$$o_X = A_0 + \sum_j^J dc_j^X \quad (4)$$

31 32 33 34 *Local gradients approach*

35
36
37 For this approach, an estimate of the value, for a given molecule X , of the partial derivative of
38
39 the model output with respect to a given descriptor was calculated as per equation (5). In
40
41 equation (5), dc_j^X denotes the descriptor contribution obtained for the j th descriptor for molecule
42
43 X ; $o_X(d_{jX}, \{d_{j^*X}\}_{j^* \neq j})$ denotes the model output calculated for molecule X based on its values
44
45 for the j th descriptor (d_{jX}) and all other descriptors; $o_X(d_{jX} + \Delta, \{d_{j^*X}\}_{j^* \neq j})$ denotes the model
46
47 output when the value of the j th descriptor is adjusted by Δ and all other descriptor values are
48
49 kept the same.
50
51
52
53
54
55
56
57
58
59
60

$$dc_j^X = \frac{o_X(d_{jX+\Delta, \{d_{j^*X}\}_{j^* \neq j}}) - o_X(d_{jX, \{d_{j^*X}\}_{j^* \neq j}})}{\Delta} \quad (5)$$

N.B. (1) The model output is the predicted bioactivity for a regression model, or score for “class 1” for a binary classification model. (2) Since binary descriptors were used for the work presented in the current article (see “Descriptors employed”), i.e. the only possible descriptor values were 1 or 0, Δ would be +1 if $d_{jX} = 0$ and -1 if $d_{jX} = 1$. (3) In practice, since only descriptor contributions for descriptors where $d_{jX} = 1$ were used as the basis for analysis in our current article (see “Heat Map interpretation of individual predictions”), $\Delta = -1$ for all calculations performed.

Heat Map interpretation of individual predictions

The high dimensionality (see Table 1) of ECFP fingerprints and the fact that their corresponding substructures overlap^{13,44} means it makes sense^{13,24} to aggregate the descriptor contributions obtained for all substructures containing a given bond (or atom) and color that bond (or atom) according to the aggregate descriptor contributions. This allows for the contributions of different regions of the molecule to the output generated by the model to be visualized in a manner that is transparent to chemists.

Here, we employed variations of the Heat Map coloring scheme of Rosenbaum et al.,¹³ where bonds (or atoms) are colored via calculating raw scores for those bonds (or atoms), based on the descriptor contributions for the corresponding ECFP feature encoded substructures present in the molecule, and then normalizing the raw bond (or atom) scores between 0 and 1. The original, raw score for a bond (or atom) is computed by summing the descriptor contribution values for all substructures containing the bond (or atom). This means that bonds (or atoms) associated with

1
2
3 features making a positive (negative) contribution to the model prediction, according to the
4
5 descriptor contributions, are assigned positive (negative) raw scores. The normalized bond (or
6
7 atom) score is transformed to a color, which runs from strong red (normalized score of 0), to
8
9 orange (normalized score above 0.25), through yellow (normalized score close to 0.5), to green
10
11 (normalized score at least 0.54), and ultimately to strong green (normalized score of 1). The
12
13 more positive the raw bond (or atom) score, the closer to 1 the normalized score will lie and,
14
15 hence, the more green the assigned color will be.¹³
16
17
18
19

20
21 In this work, various approaches to normalizing the raw bond (or atom) scores were considered:
22
23 (a) Single Molecule Normalization, (b) Full Dataset normalization, (c) Symmetrized Single
24
25 Molecule Normalization and (d) Constant Symmetrized Normalization. Approaches (a – d)
26
27 correspond to equations (6 – 9) respectively. Approaches (a) and (b) were proposed by
28
29 Rosenbaum et al.,¹³ whilst the symmetrized schemes (c) and (d) were introduced in the current
30
31 work. In contrast to the previously proposed normalization schemes, Symmetrized Single
32
33 Molecule Normalization and Constant Symmetrized Normalization ensure that negative and
34
35 positive raw scores are matched to normalized scores less than 0.5 and greater than 0.5
36
37 respectively. This ensures that regions of the molecule contributing negatively (positively) to the
38
39 prediction are colored from yellow, to orange, to red (from yellow to green). In principle,
40
41 Symmetrized Single Molecule Normalization, as opposed to Constant Symmetrized
42
43 Normalization, may result in weaker green coloring being assigned to equally positive regions of
44
45 the molecule upon moving to a molecule with large magnitude negative raw scores. However,
46
47 both Full Dataset Normalization and Constant Symmetrized Normalization may result in small
48
49 magnitude positive or negative raw bond (or atom) scores being assigned normalized scores
50
51
52
53
54
55
56
57
58
59
60

close to 0.5, such that they appear yellow i.e. their sign is indeterminate from the Heat Map image.

$$s_n = \frac{(s_r - s_r^{min,mol})}{(s_r^{max,mol} - s_r^{min,mol})} \quad (6)$$

$$s_n = \frac{(s_r - s_r^{min,ds})}{(s_r^{max,ds} - s_r^{min,ds})} \quad (7)$$

$$s_n = \frac{(s_r + |s|_r^{max,mol})}{(2 \times |s|_r^{max,mol})} \quad (8)$$

$$s_n = \frac{(s_r + |s|_r^{max,set})}{(2 \times |s|_r^{max,set})} \quad (9)$$

In equations (6 – 9), s_r denotes a raw bond (or atom) score and s_n denotes a normalized score. $s_r^{min,mol}$ denotes the minimum (i.e. most negative if there are negative raw scores) raw score for the molecule, whilst $s_r^{max,mol}$ denotes the maximum (i.e. most positive if there are positive raw scores) raw score for the molecule. $s_r^{min,ds}$ and $s_r^{max,ds}$ denote the minimum and maximum raw bond (or atom) scores respectively across a specified dataset of molecules; in keeping with Rosenbaum et al.,¹³ this dataset was the training set plus the molecule of interest. $|s|_r^{max,mol}$ denotes the largest raw bond (or atom) score *magnitude* for the molecule of interest. Finally, $|s|_r^{max,set}$ denotes the largest raw bond (or atom) score *magnitude* corresponding to a set of model predictions and associated descriptor contributions. In the current work, this set was the complete set of descriptor contributions calculated for all molecules of interest (see “Molecules considered for detailed analysis”) for a given dataset with all combinations of modelling,

1
2
3 descriptor contribution and, where relevant to the descriptor contributions, predictions
4
5 approaches.
6
7

8
9 Due to software implementation challenges, it was not possible to generate Heat Map images
10
11 corresponding to Full Dataset Normalization, i.e. equation (7), for Random Forest models. (See
12
13 “Computational details” in the Supporting Information.)
14
15

16
17 The interpretability of predictions was assessed via evaluating the biological plausibility of the
18
19 Heat Map colors. Further details are provided under “Quality assessment of Heat Map images”.
20
21

22 23 **Molecules considered for detailed analysis**

24
25 Heat Map images and corresponding “external, leave-one-out” predictions were generated and
26
27 evaluated for specific molecules of interest. The models used to generate these Heat Maps and
28
29 predictions were built upon the entirety of the corresponding dataset, save for the molecule of
30
31 interest where applicable, using the same procedure employed to build models evaluated via
32
33 “external cross-validation”. For the Kazius and CA datasets, the same molecules subjected to
34
35 detailed analysis by Rosenbaum et al.¹³ and Mohr et al.¹⁷ respectively were considered and
36
37 removed, in turn, from the datasets. For all protein interaction datasets, corresponding Protein
38
39 Data Bank (PDB)^{63,64} protein-ligand crystal structures, with ligands not found in the original
40
41 datasets, were identified. Further details are provided in the Supporting Information.
42
43
44
45
46

47
48 In addition to these molecules, for which manual evaluation of the Heat Map images was
49
50 performed, “external, leave-one-out” predictions and corresponding Heat Map images were
51
52 generated for an additional set of molecules from the Kazius dataset matching a set of
53
54 toxicophore encoding SMARTS patterns. These SMARTS patterns are provided in the
55
56
57
58
59
60

1
2
3 Supporting Information file “Kazius_2005_SI_SA_SMARTS_subset_v2.xls” and were derived
4
5 from Kazius et al.,⁴⁸ as explained therein.
6
7

8 9 **Quality assessment of Heat Map images**

10
11 The Heat Map images generated for the molecules of interest were assessed in terms of (a)
12
13 available mechanistic knowledge concerning (potential) toxicophores and detoxifying groups or,
14
15 for the models trained on the protein interaction datasets, (b) the corresponding PoseView⁶⁵
16
17 images of protein-ligand interactions available from the PDB.
18
19

20
21
22 Primarily, these images were assessed by considering how the molecular substructures
23
24 highlighted as most positive (i.e. green), suggesting a significant positive contribution towards a
25
26 prediction of membership of the active (toxic) class (“class 1”) or the predicted continuous
27
28 biological activity, corresponded to those substructures expected to be directly responsible for
29
30 promoting biological activity. In the case of molecules from the Kazius⁴⁸ and CA datasets,¹⁷ this
31
32 entailed comparing the most positive regions of the molecule to the known (or hypothesized)
33
34 toxicophores for Ames mutagenicity or chromosome aberration respectively.^{13,17,48} In the case of
35
36 those images which corresponded to PDB ligands, the most positive regions of the molecule
37
38 were compared to the interactions indicated by the PoseView images⁶⁵ provided by the PDB.⁶³
39
40
41
42
43

44
45 An important point to note is that, by assessing the correspondence between the Heat Map and
46
47 the known influence of the biologically relevant substructures, this assessment is not necessarily
48
49 equivalent to an assessment of how well the interpretation approach provides the basis for the
50
51 model’s prediction. Rather, this assessment evaluates the extent to which the Heat Map images
52
53 can provide insights into the biological significance of different molecular substructures. For a
54
55 prediction made for an untested chemical, this is useful information to obtain from a QSAR
56
57
58
59
60

1
2
3 model in its own right, hence an assessment of how well the Heat Map images can provide this
4 information on chemicals for which the most biologically relevant substructures are known is
5 valuable. (This assumes the ability of the Heat Map images to provide this information on those
6 chemicals is representative of how well they would provide this information on chemicals for
7 which the most biologically significant substructures are unknown.) Furthermore, if it can be
8 assumed that the model makes correct predictions for the right reason, this assessment indicates
9 how well a Heat Map image explains the basis for the model's prediction.
10
11
12
13
14
15
16
17
18
19

20
21 Assessment of the correspondence between green coloring, suggesting a significant positive
22 contribution to predicted bioactivity, on the Heat Map and the substructures expected to be
23 directly responsible for promoting biological activity took account of two important
24 considerations. The first consideration was the extent to which regions of the molecule known to
25 directly promote biological activity were highlighted as green. The second consideration was the
26 extent to which biologically irrelevant and, for the minority of molecules for which this was
27 applicable, known activity reducing regions of the molecule were highlighted as green. Both
28 failing to (fully) identify those regions of the molecule directly responsible for promoting
29 biological activity, as well as wrongly identifying biologically irrelevant or activity reducing
30 regions of the molecule as promoting biological activity with respect to the endpoint of interest,
31 would mislead scientists seeking to interpret the images, e.g. medicinal chemists seeking ideas
32 for lead optimization. Hence, the following ranking scheme was initially used to assess the
33 images (**Table 2**).
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51
52 For some molecules for which Heat Map images were analyzed, detoxifying substructures were
53 identified which were known to reduce (or even abolish) the biological activity expected to be
54 conferred by the corresponding toxicophore.^{13,48} For Heat Map images corresponding to
55
56
57
58
59
60

1
2
3 predictions for PDB ligands, hydrogen bond donors/acceptors (e.g. tertiary amines) in regions
4 suggested to be involved in hydrophobic interactions by the corresponding PoseView⁶⁵ images
5 were penalized when evaluating positive (green) coloring according to the ranking scheme
6
7
8
9
10
11 (Table 2). A failure to satisfy ligand hydrogen bond acceptors and, possibly even more so,
12 donors located within a hydrophobic pocket is enthalpically unfavorable.⁶⁶
13
14
15

16 **Table 2.** Summary of the ranking scheme initially used to denote the quality of the Heat Map
17 images based on considering the biological relevance of the most positively highlighted (i.e.
18 green) regions of the molecule
19
20
21
22
23

Biological relevance (positive coloring): Rank ^a	Biological relevance (positive coloring): Label	Criteria used for assignment
1	Strongest	All regions of the molecule contributing directly towards biological activity were highlighted as most positive in their entirety, without any biologically irrelevant substructures, or substructures expected to reduce biological activity, highlighted as most positive.
2	Strong	The regions of the molecule contributing directly towards biological activity were partially highlighted as most positive, without any biologically irrelevant substructures, or substructures expected to reduce biological activity,

		highlighted as most positive.
3	Weak	All regions of the molecule contributing directly towards biological activity were highlighted as most positive in their entirety, but some biologically irrelevant substructures, or substructures expected to reduce biological activity, were also highlighted as most positive.
4	Weaker	The regions of the molecule contributing directly towards biological activity were partially highlighted as most positive, but some biologically irrelevant substructures, or substructures expected to reduce biological activity, were also highlighted as most positive.
5	None	No match between the molecular regions highlighted as most positive and those known to directly promote biological activity was observed.

a. A lower rank denotes greater biological relevance i.e. higher Heat Map image quality.

For the minority of molecules with known activity reducing substructures, an analogous analysis was performed based upon assessing the correspondence between those substructures highlighted as most negative (i.e. red) and the known activity reducing substructures.

Ultimately, these Heat Map images will be useful for chemists inspecting them to gain insights into the basis for a given model prediction and/or the indicated biological significance of molecular substructures, according to a given interpretation approach and QSAR model.

1
2
3 However, manual assessment of Heat Map images is not only time consuming but is also prone
4
5 to inconsistent categorization, i.e. the scheme outlined in **Table 2** does not fully remove all
6
7 potential subjectivity or inconsistency associated with manual Heat Map evaluation.
8
9

10 Nonetheless, assessment according to a systematized scheme is necessary for effective
11
12 evaluation of the degree to which different interpretation approaches yield chemically and
13
14 biologically meaningful interpretations of QSAR predictions. To overcome the limitations of
15
16 manual assessment, an automated, numeric “quality score” scheme was implemented. This
17
18 automated scheme, as illustrated by **Figure 1**, was only implemented in the current work for
19
20 Heat Map images based on atom coloring and either Symmetrized Single Molecule
21
22 Normalization, or Constant Symmetrized Normalization. Furthermore, it was only applied to a
23
24 subset of Kazius dataset molecules for which toxicity promoting substructures and, where
25
26 applicable, corresponding detoxifying substructures were encoded using SMARTS patterns
27
28 derived from the work of Kazius et al.⁴⁸
29
30
31
32
33

34
35 Each of the relevant Heat Map images is assigned a single, numeric “quality score”, as
36
37 illustrated by **Figure 1**, designed to estimate the degree to which the image provides a
38
39 chemically and biologically meaningful interpretation of the corresponding model prediction.
40
41 For a molecule with one or more SMARTS pattern identified toxicophores, the quality score
42
43 measures the degree to which the green coloring of atoms fully and precisely identifies the atoms
44
45 lying inside the biologically active substructure(s). (Only atoms estimated to make a relatively
46
47 significant positive contribution to the predicted bioactivity will be colored green if symmetrized
48
49 normalization and atom coloring is used.) This is quantified in terms of the F-measure, which
50
51 takes into account the proportion of true positives to false “negatives” and false positives.⁶⁷ Here,
52
53 the true positives are all atoms with a normalized atom score ≥ 0.54 , which was visually
54
55
56
57
58
59
60

1
2
3 estimated to be the point at which an atom starts to appear green, and which were matched by
4 toxicophore SMARTS. False positives have normalized atom scores ≥ 0.54 , yet are not
5
6
7
8 matched by a toxicophore SMARTS. False “negative” atoms are those matched by the
9
10
11 toxicophore SMARTS, but which are not selected based on this threshold. This calculation is
12
13 illustrated, for an example molecule, on the LHS of **Figure 1**.

14
15
16 For a molecule which also contains substructures that are known to be deactivating, an analogous
17
18 calculation is performed with regards to those atoms highlighted as sufficiently red. (This
19
20 calculation is illustrated, for an example molecule, on the RHS of **Figure 1**.) For Heat Map
21
22 images generated via symmetrized normalization and atom coloring, this means those atoms are
23
24 known to make a significant negative contribution to the predicted bioactivity. The atoms
25
26 indicated to significantly reduce toxicity, based on the Heat Map, were those with a normalized
27
28 score ≤ 0.25 , which was visually estimated to be the point at which an atom becomes
29
30
31 discernably red. In this latter case, the overall quality score is calculated as the arithmetic mean
32
33
34
35 of the two F-measures.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

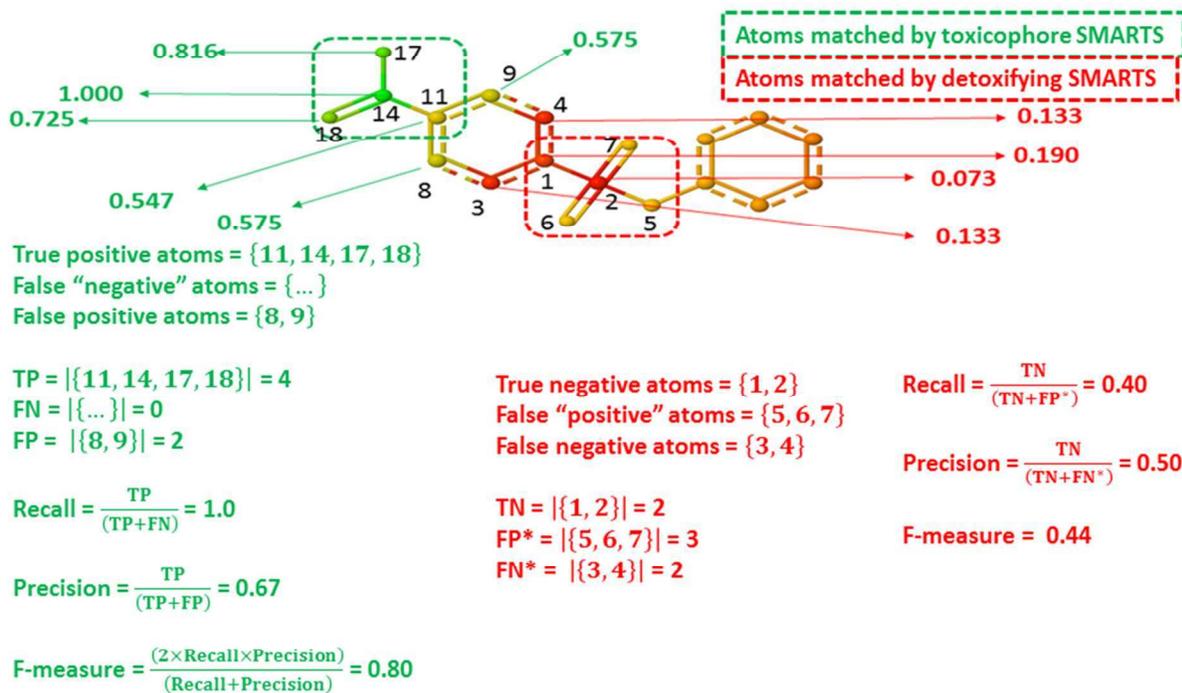


Figure 1 The procedure for calculating numeric Heat Map quality scores for an image obtained via symmetrized normalization (Symmetrized Single Molecule Normalization in this example) and atom coloring. The relevant normalized atom score estimates are shown matched to the atom IDs. In this case, the overall quality score would be the arithmetic mean of the two F-measures i.e. 0.62. Further details are provided in the associated text. SMARTS matching of toxicophores and detoxifying substructures was implemented using OpenBabel (version 2.3.2) & Pybel,^{68,69} based on SMARTS adapted from the work of Kazius et al.⁴⁸ and reported in the Supporting Information. The atom IDs, used to identify the corresponding atom score estimates for this image, were identified using OpenBabel⁶⁸, via this command: "obabel [SDF containing molecule of interest] -O [name of image].svg -xi -xu".

Normalization of binary classification scores for molecules subjected to detailed analysis

The quality of the Heat Map images was compared to the quality of the corresponding predictions (see below). The raw prediction scores of the binary classification models were normalized to facilitate comparison across methods. The raw prediction scores for “class 1”, the “toxic” class, are of different scales. Normalization entailed transforming these values, where required, to lie between 0 and 1. Subsequent to normalization, a score greater than 0.50 resulted in predicted assignment to “class 1” for all methods. Details of this normalization are provided in the Supporting Information.

Correlation between prediction quality and Heat Map image quality

Heat Map images have the potential to inform chemists in the analysis and design of molecules. The quality of those images, in terms of correctly highlighting active regions of the molecule, is therefore of paramount importance. It is useful to understand the extent to which the quality of a Heat Map image depends upon the quality/strength of the corresponding prediction. If a strong correlation between Heat Map quality and prediction quality is observed, this can help chemists judge whether a Heat Map image offers meaningful interpretation in the absence of knowing the true determinants of biological activity. Firstly, in the case that predictions are made for a biologically untested molecule, statistical estimates of the prediction quality (e.g. as obtained from cross-validation on the training set) may suggest whether the Heat Map image can be trusted. Secondly, in the case that the prediction is experimentally tested yet the mechanism of action remains unknown, the known quality of the prediction could be used to judge whether the Heat Map image offers useful insights.

1
2
3 In order to determine whether the quality of these Heat Map images, calculated robustly as per
4
5 **Figure 1**, corresponded to the quality of the corresponding predictions, it was necessary to rank
6
7 the quality of these predictions. Since the scoring scheme described in **Figure 1** could only be
8
9 applied to external leave-one-out binary classification predictions for the Kazius dataset, high
10
11 quality predictions were defined as correct predictions, with all incorrect predictions being
12
13 considered of low quality.
14
15

16 17 18 **Computational details** 19

20
21 Scripts and software configuration files, used for evaluation of the modelling approaches
22
23 considered in our article, have been made freely available.⁷⁰ These files are described, along with
24
25 information about software and hardware dependencies used to generate our results, in the
26
27 Supporting Information. Two key dependencies, which we have made freely available to the
28
29 cheminformatics community, are the *rfFC* package³⁷ for the R Statistical Programming
30
31 language⁷¹ and the *HeatMapWrapper* Python program.⁷² The *rfFC* package implements
32
33 algorithms for interpreting predictions made using the *randomForest* package. The
34
35 *HeatMapWrapper* program allows for any set of descriptor contributions, associated with binary
36
37 descriptors corresponding to fingerprints calculated as per our work, to be used as the basis for
38
39 Heat Map images. Hence, the former program supports interpretation of Random Forest
40
41 predictions, whilst the latter program supports the interpretation of any QSAR prediction built
42
43 using the same fingerprints as our models, as long as suitable descriptor contributions can be
44
45 calculated. These dependencies were used to generate the Heat Map images analyzed in our
46
47 current article, as well as provide estimates of the underlying normalized and corresponding raw
48
49 atom scores. Additional details are presented in the Supporting Information.
50
51
52
53
54
55
56
57
58
59
60

RESULTS & DISCUSSION

Predictive performance

The overall performance estimates for the investigated modelling approaches, from “external cross-validation”, are presented in **Table 3** for key binary classification datasets. The performance of all modelling approaches on all remaining binary classification datasets was worse, in terms of all figures of merit, than for the Kazius dataset and typically worse than for the CA dataset. Subsequent investigations of prediction interpretations, based on additional leave-one-out predictions, are discussed in detail for the Kazius dataset, with some illustrative examples also taken from the CA dataset. Since modelling results on the other datasets are not considered in greater detail in our article, further detailed performance estimates are only presented in the Supporting Information (Tables S5 and S6 under “Predictive performance: additional result summaries”).

Nonetheless, in summary, Random Forest classification, with “averaged predictions”, is the best performing approach, in terms of the arithmetic mean AUC, for all but one (MUV713) of the five classification datasets. (For MUV713, Random Forest classification using the default majority vote predictions is marginally better, with identical mean AUC at 3dp. However, the average performance of all models on the MUV713 dataset is close to that expected from random guessing.)⁶¹ Random Forest classification, using the “averaged predictions” approach, also typically yields the highest predictive performance in terms of the arithmetic mean values for all other performance metrics. However, it should be noted that the differences in performance between methods are sometimes marginal and may not be robust.

Likewise, Random Forest regression is the best performing approach according to the arithmetic mean R^2 for four (ACE, DHFR, GPB, THR) out of the eight regression datasets. Again, it should be noted that the differences in performance are sometimes marginal and may not be robust.

Regarding the robustness of the results, an approximate assessment of statistical significance was carried out for the pairwise differences in mean performance metrics. For both the classification and regression datasets, over half of differences were statistically significant according to a standard paired, two-tail t-test, even upon applying multiple testing corrections.⁷³ However, this test suffers from elevated Type 1 error in our evaluation framework of averaged K-fold testing,^{74,75} as explained in the discussion of “Statistical Significance” in the Supporting Information. Therefore, it is likely that the results are not trustworthy and, hence, no detailed analysis is presented.

Table 3. Overall performance estimates of binary classification modelling approaches obtained from "external cross-validation" for key datasets: arithmetic mean performance statistics (3dp)

Dataset	Modelling Approach	AUC	MCC	Kappa	BA
CA	Random Forest classification (AP) ^a	0.762	0.369	0.352	0.683
	Random Forest classification	0.748	0.358	0.333	0.676
	Support Vector Machine	0.747	0.346	0.327	0.670
	Probit-PLS	0.746	0.339	0.326	0.664
Kazius	Random Forest classification (AP) ^a	0.876	0.600	0.598	0.802

	Random Forest classification	0.862	0.595	0.592	0.799
	Probit-PLS	0.855	0.570	0.566	0.787
	Support Vector Machine	0.854	0.568	0.564	0.785

- a. “AP” denotes Random Forest binary classification predictions made using the, non-default, “averaged predictions” approach described under the Random Forest classification subsection of “Modelling approaches”.

Summary of Heat Map interpretations of individual predictions

A variety of Heat Map images were generated, corresponding to (a) leave-one-out binary classification predictions for the Kazius and CA datasets, for 39 and 15 molecules respectively, or (b) regression predictions made for eight PDB ligands corresponding to one of the Sutherland datasets. A summary of these images and the assessment made of the correspondence between the Heat Map and the known biologically relevant substructures is provided in the Supporting Information (see “FinalImagesAnalysis_withImages_resub.xlsx”).

As well as different approaches to generating descriptor contributions, different Heat Map images were also obtained for the same molecule via translating these descriptor contributions into Heat Map images according to (a) different normalization schemes and (b) atom or bond coloring based on the descriptor contributions. These different schemes are explained in “Heat Map interpretation of individual predictions” under “Methods & Data” above. Due to software implementation limitations explained under “Computational details” in the Supporting Information, it was only possible to generate “Full Dataset Normalization” Heat Map images for

1
2
3 the linear modelling methods: Probit-PLS, PLS regression and, with a linear kernel, SVM and
4
5 Support Vector Regression.
6
7

8
9 Systematic manual assessments, according to **Table 2**, were made for all Heat Map images
10
11 corresponding to the examples from the Kazius and CA datasets (two and 15 molecules
12
13 respectively) which were previously used for assessment of interpretation approaches,^{13,17} along
14
15 with all PDB ligands. In addition, it was possible to calculate automated numeric Heat Map
16
17 quality scores, as per **Figure 1**, for all 39 molecules in the Kazius dataset, identified via
18
19 toxicophore SMARTS pattern matches, for which Heat Map images were generated. Whilst 2026
20
21 Kazius dataset molecules were identified via toxicophore SMARTS pattern matches and, in
22
23 principle, automated numeric Heat Map quality scores could have been calculated for all of
24
25 these, the currently available software for generating Heat Map images⁷² requires laborious
26
27 manual intervention via a GUI to obtain each image. This limited the number of Heat Map
28
29 images it was practical to generate for the current paper.
30
31
32
33
34
35

36 **Effect of different coloring schemes upon Heat Map interpretations**

37
38 The choice of atom or bond coloring, as well as the choice of normalization scheme, can yield
39
40 discernably different Heat Map images. This is illustrated by **Figure 2**.
41
42

43
44 Only the symmetrized normalization schemes introduced in this paper (Symmetrized Single
45
46 Molecule Normalization and Constant Symmetrized Normalization) are guaranteed to ensure that
47
48 positive raw bond (atom) scores, calculated from the underlying descriptor contributions, are
49
50 translated into yellow to green coloring, becoming increasingly green with increasing magnitude,
51
52 and negative raw scores are translated into yellow to orange to red coloring. With the previously
53
54 proposed normalization schemes,¹³ Full Dataset Normalization and especially Single Molecule
55
56
57
58
59
60

1
2
3 Normalization, green or red coloring might be assigned which would give a misleading
4
5 impression as to the sign of the raw bond (or atom) scores indicating whether a particular
6
7 molecular substructure was responsible for a negative or positive effect on the model's
8
9 prediction. This can be seen by comparing the images in **Figure 2** with the estimated
10
11 corresponding raw atom scores for this molecule in **Figure 3**. N.B. All normalized and raw atom
12
13 scores we present are estimates, as discussed under "Computational Details" in the Supporting
14
15 Information.
16
17
18
19

20
21 As illustrated by **Figure 2**, both Full Dataset Normalization and Constant Symmetrized
22
23 Normalization can make the magnitude of all normalized bond (or atom) scores close to 0.5. This
24
25 either reduces the colors to yellow, which conveys no information about the sign of the raw bond
26
27 (or atom) scores and allows for no visual interpretation, or makes them only weakly green or
28
29 orange, which makes differences between different parts of the molecule harder to discern.
30
31
32

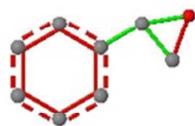
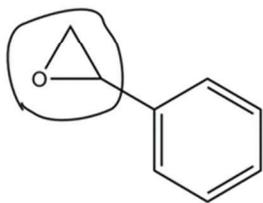
33
34 In principle, Constant Symmetrized Normalization ensures that only substructures contributing
35
36 negatively to the prediction and substructures contributing positively would appear red and green
37
38 respectively, in keeping with Symmetrized Single Molecule Normalization, as well as offering
39
40 certain theoretical advantages. Firstly, it allows for the variation in the strength of negative and
41
42 positive coloring between molecules to be evaluated, due to the use of consistent normalization
43
44 across all Heat Map images. However, as long as the regions are sufficiently positive or negative
45
46 to appear green and red respectively under Symmetrized Single Molecule Normalization or
47
48 Constant Symmetrized Normalization, this would have no bearing on the outcome of the
49
50 evaluation schemes applied in our work. Secondly, the use of Constant Symmetrized
51
52 Normalization avoids a hypothetical problem with Symmetrized Single Molecule Normalization:
53
54 the possibility that a positive region of the molecule may switch from green to yellow due to the
55
56
57
58
59
60

1
2
3 inclusion of a large magnitude negative raw atom score associated with a red substructure.

4
5 However, this does not appear to have happened in practice, as consideration of both **Figure 4**
6
7 and **Figure 5** indicates. Rather, as illustrated by **Figure 5**, the effect of Constant Symmetrized
8
9 Normalization appears to typically be to reduce the normalized atom scores close to 0.5, making
10
11 the corresponding atoms (or bonds) appear yellow and, hence, impossible to obtain any
12
13 interpretation from. For this reason, the use of Symmetrized Single Molecule Normalization is
14
15 recommended and Heat Map images generated according to this scheme were used to draw
16
17 conclusions from our work.
18
19
20
21

22
23 Regarding the choice of atom or bond coloring, atom coloring is arguably more appropriate for
24
25 discerning whether a hydrogen bond donor or acceptor atom⁶⁶ is making a significant
26
27 contribution towards the prediction. Different bond coloring either side of this atom could lead to
28
29 an ambiguous interpretation. However, it *may* be possible for some pi-bonds to act as weak
30
31 hydrogen bond acceptors,⁷⁶ in which case bond coloring would be more appropriate, as was
32
33 previously suggested in the literature on different grounds.¹³ However, the Pybel software⁶⁹ used
34
35 for SMARTS matching, hence automated numeric Heat Map quality scoring, for the current
36
37 paper identified the parts of the molecule matched by the SMARTS patterns in terms of atom
38
39 indices. Hence, this automated scoring was only applicable to Heat Map images based on atom
40
41 coloring. For this reason, only Heat Map images corresponding to atom coloring are considered
42
43 for the rest of our article.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Highlighted: epoxide moiety

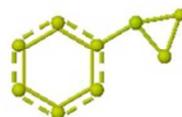


(A) Expert reasoning assigned toxicophore

(B) Bond color, Single Molecule Normalization

(C) Bond color, Full Dataset Normalization

(D) Bond color, Symmetrized Single Molecule Normalization



(E) Atom color, Symmetrized Single Molecule Normalization

(F) Atom color, Single Molecule Normalization

(G) Atom color, Full Dataset Normalization

(H) Atom color, Constant Symmetrized Normalization

Figure 2 Heat Map analysis of a binary classification leave-one-out SVM prediction made for molecule “CA23” in the CA dataset (c.f. Figure 4(b) in Mohr et al.), determined to be toxic (positive result) in the chromosome aberration test.¹⁷ (A) Expert reasoning assigned toxicophore (epoxide functionality), based on the rationale provided in Mohr et al.:¹⁷ this substructure was flagged by an Ames test alerts knowledgebase, with 80% of Ames positives being chromosome aberration positives, and is known to exhibit electrophilic reactivity with DNA. (B – H): Heat Map images based on descriptor contributions associated with the same SVM prediction, generated according to either bond or atom coloring obtained via different normalization schemes.

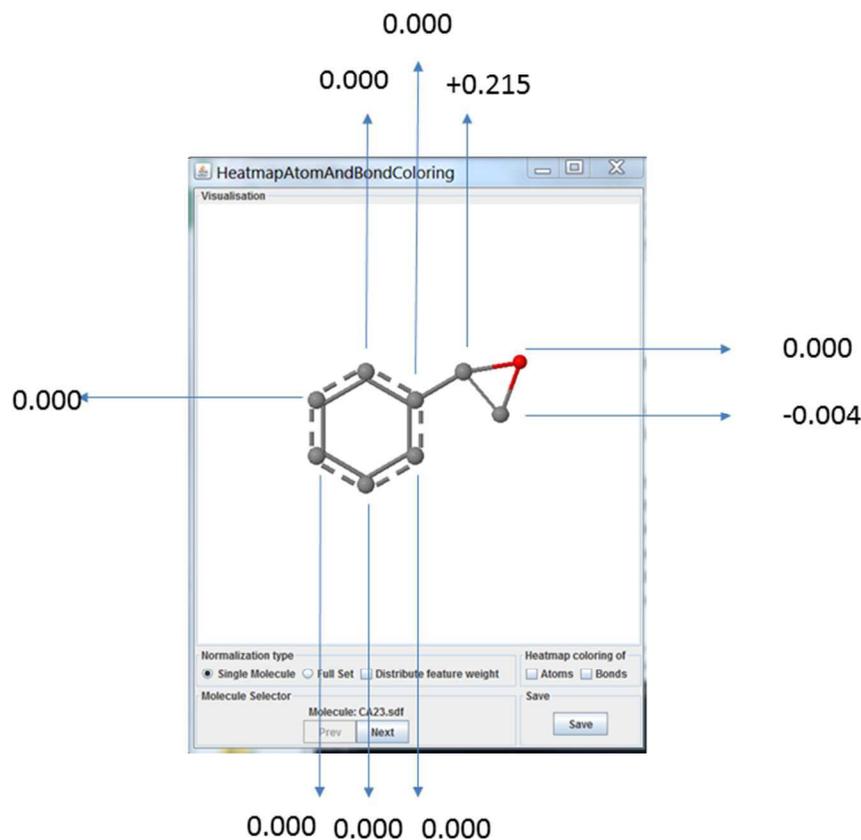


Figure 3 Image of the molecule “CA23” in the CA dataset¹⁷ displayed in HeatmapViewer.jar GUI,¹³ as generated by running the HeatMapWrapper tool,⁷² with both atom and bond coloring switched off. The intermediate output generated via running the HeatMapWrapper tool in atom color, Symmetrized Single Molecule Normalization mode, based on leave-one-out SVM descriptor contributions, was processed to reveal estimated raw atom scores ranging from weakly negative in one case (-0.004), through zero in most cases, to positive (+0.215) in one case. The estimated raw atom scores are rounded to 3dp. This image corresponds to all Heat Map images shown in **Figure 2**, with different normalization schemes generating the different normalized scores used for coloring. Via switching off atom coloring, based on normalized scores, the epoxide oxygen (in red) and carbon atoms (in black) can readily be discerned. N.B. The specific

1
2
3 atom IDs, used to identify the corresponding atom score estimates, were identified as per **Figure**
4
5
6 **1**.

9 **Strengths and weaknesses of different Heat Map quality assessment schemes**

10
11 Ultimately, the Heat Map images are designed to provide chemically and biologically
12
13 meaningful interpretations of model predictions to chemists, hence require manual consideration.
14
15 However, a systematized evaluation scheme is required to allow for a reasonable assessment of
16
17 the degree to which different modelling and prediction interpretation approaches, underpinning
18
19 these Heat Map images, tend to yield chemically and biologically meaningful interpretations.
20
21

22
23
24 Whilst a manual, qualitative evaluation scheme (**Table 2**) was initially devised, this has a
25
26 number of weaknesses. These prompted the development of an automated, numeric scoring
27
28 scheme (**Figure 1**). The following figures (**Figure 4**, **Figure 5**, **Figure 6**) compare and contrast
29
30 the application of both schemes to two Kazius dataset molecules, with Heat Map images
31
32 generated according to various scenarios, for which external leave-one-out predictions were
33
34 obtained.
35
36

37
38
39 Arguably the greatest weakness of the manual ranking scheme (**Table 2**), is the potential for
40
41 subjectivity or inconsistency in assignments to be made. Inconsistency is likely to be greater if a
42
43 larger number of images are considered, with the laborious requirement to visually consider each
44
45 image restricting the number of images which can be considered and, hence, limiting the
46
47 robustness of any observed trends. The application of the automated, numeric scoring scheme
48
49 (**Figure 1**) allows for greater objectivity and for a large number of Heat Map images to be
50
51 assessed with ease, allowing more robust conclusions to be drawn. Of course, this assumes a
52
53 large number of Heat Map images have been generated and the molecular features directly
54
55
56
57
58
59
60

1
2
3 responsible for promoting (or reducing) activity have been suitably encoded as SMARTS
4
5 patterns.
6
7

8
9 Nonetheless, it could be argued that the imposition of hard boundaries between those atoms
10 involved in biologically relevant and irrelevant molecular substructures, via SMARTS matching,
11 whilst removing subjectivity, might sometimes be problematic. Firstly, it should be
12 acknowledged that this sharp distinction between biologically relevant and irrelevant molecular
13 substructures is a simplification e.g. whole molecule physicochemical properties, such as logP,
14 may affect biological activity and arise from molecular features not directly involved in
15 interactions at the site of biological action. Secondly, the SMARTS patterns may fail to account
16 for all biologically relevant knowledge. For example, consider the molecule shown in **Figure 7**.
17 Whilst the putative toxicophore¹⁷ encompasses the entire arylhydrazide substructure, it *may* be
18 the case that biological activity arises from metabolism of this substructure to yield the aromatic
19 amine substructure.^{17,48} Hence, it can reasonably be argued that application of the qualitative,
20 manual quality assessment scheme should only assign credit for green coloring of that
21 substructure. A different example, which further illustrates the potential problem associated with
22 applying hard boundaries between biologically relevant and irrelevant sets of atoms as supposed
23 by the automated scoring scheme, is **Figure 4**. According to the SMARTS patterns used for
24 application of the automated scoring scheme, **Figure 4** (B) and (E) should be penalized for
25 highlighting the ortho positions, to the sulfonamide attached to the aromatic ring, as red.
26 However, one might reasonably argue that, at least in this context, this still serves to indicate that
27 the specific effect of having a sulfonamide attachment at this end of the ring is detoxifying. In
28 principle, one could adjust the SMARTS patterns to accommodate these scenarios when
29 applying the numeric scoring scheme e.g. by expanding the SMARTS patterns used for
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 automated assessment of the molecule in **Figure 4** to accommodate more of the aromatic ring.
4
5 However, the most appropriate modifications would still need to be decided with care. For
6
7 example, consider simply extending these SMARTS patterns to encompass a larger portion of
8
9 the aromatic ring. This could risk giving credit to Heat Map images which highlighted a
10
11 significant portion of the aromatic ring, or other attachment points, without highlighting the
12
13 attachments which were actually responsible for promoting (or attenuating) biological activity.
14
15 Likewise, with reference to this example, the exact point at which to stop the SMARTS pattern
16
17 indicating that a sulfonamide attachment to the ring was detoxifying and a para-nitro attachment
18
19 was a toxicophore is also debatable.
20
21
22
23
24

25
26 In spite of this, it should still be concluded that the numeric, automated scoring scheme is a more
27
28 appropriate means for systematic assessment of Heat Map images than the qualitative, manual
29
30 scheme. As well as the potential for inconsistency and difficulty of applying to large numbers of
31
32 images already discussed, consider scenario in which there are two toxicophores that can be
33
34 encoded as SMARTS patterns (**Figure 6**). The qualitative scheme would fail to assign more
35
36 credit to a Heat Map image which fully highlighted one toxicophore and partially highlighted
37
38 another than to a Heat Map image which only partially highlighted one of them. The numeric,
39
40 automated scheme overcomes this limitation.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

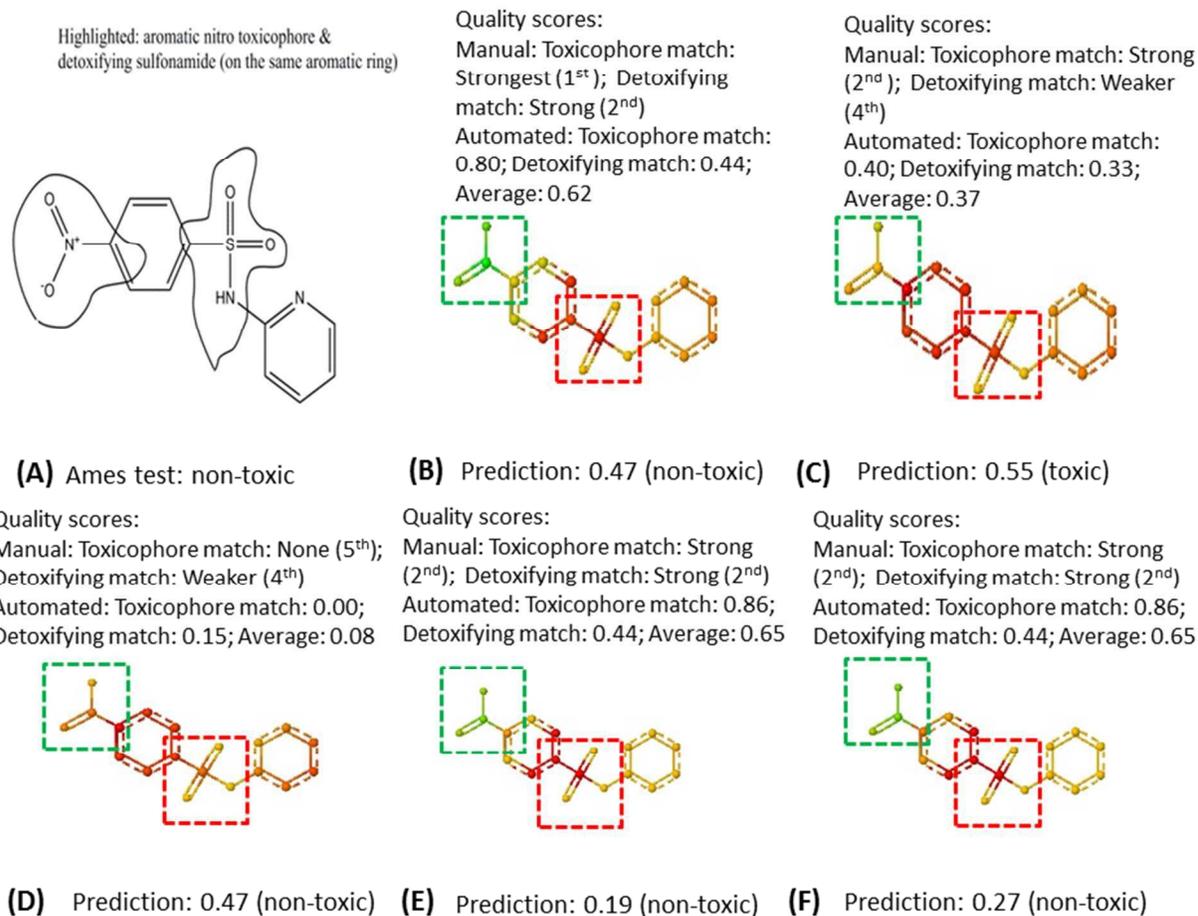


Figure 4 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule “1028-11-1” in the Kazius dataset (c.f. Figure 3(A) in Rosenbaum et al.).¹³ (A) Molecule with biologically significant substructures shown, based on a combination of mechanistic reasoning, expert knowledge and data analysis reported previously in the literature,^{13,48} and experimental Ames test assignment reported.¹³ (B – F) Heat Map images: (B) Random Forest classification (Kuz’min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. The predictions shown are the normalized scores (transformed into values between 0 and 1) for the toxic class, with predicted assignment to the toxic class if this score was greater than

0.50. The green dotted lines show the atoms matched to toxicophore SMARTS and the red dotted lines the atoms matched to corresponding detoxifying SMARTS. SMARTS matching and identification of matched atoms were carried out as per **Figure 1**. The manual and automated assessments of Heat Map quality are based on the schemes explained under “Quality assessment of Heat Map images”.

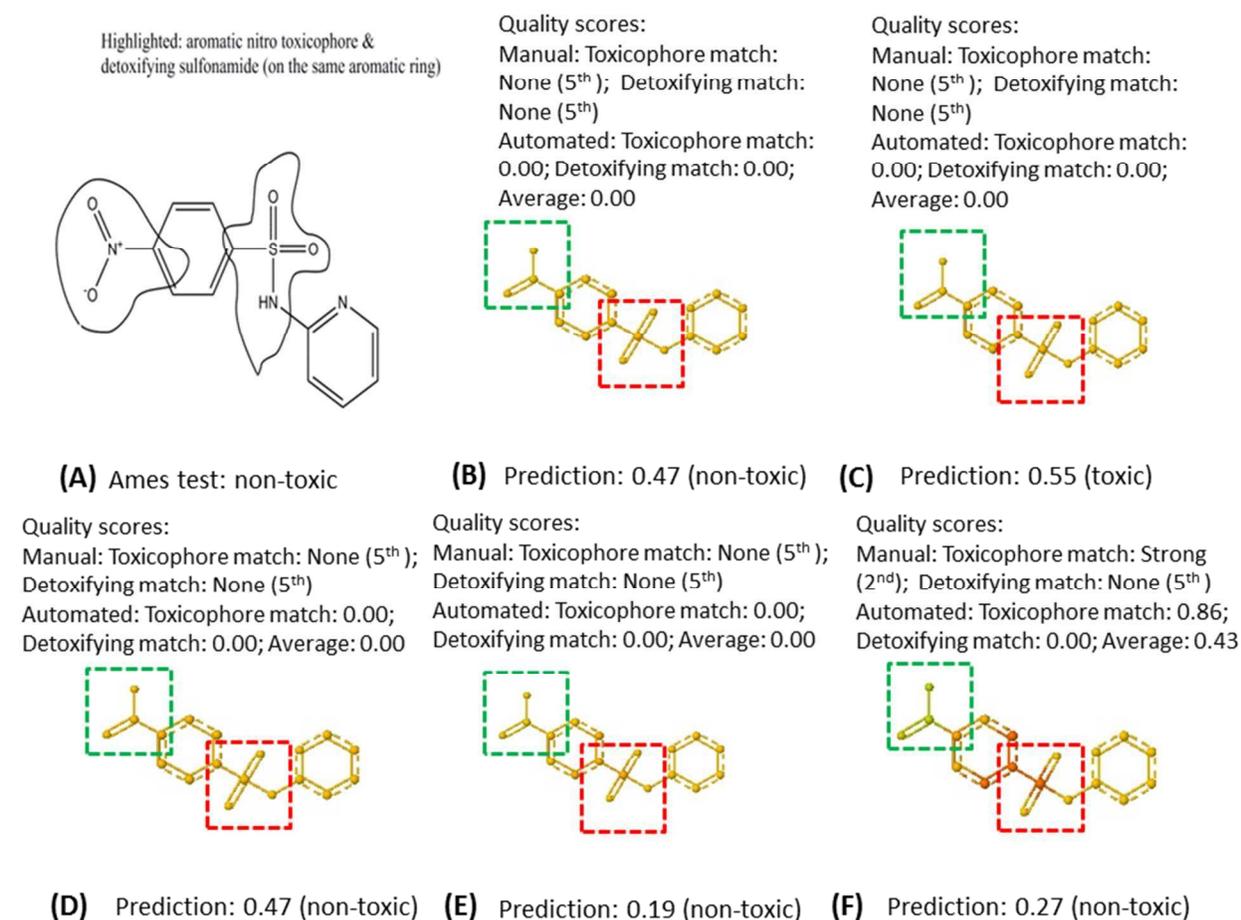


Figure 5 Heat Map (atom coloring & Constant Symmetrized Normalization) analysis of binary classification leave-one-out predictions made for molecule “1028-11-1” in the Kazius dataset (c.f. Figure 3(A) in Rosenbaum et al.).¹³ (A) Molecule with biologically significant substructures shown, based on a combination of mechanistic reasoning, expert knowledge and data analysis

1
2
3 reported previously in the literature,^{13,48} and experimental Ames test assignment reported.¹³ (B –
4
5 F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged
6
7 predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D)
8
9 Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are
10
11 reported as per **Figure 4**. The green dotted lines show the atoms matched to toxicophore
12
13 SMARTS and the red dotted lines the atoms matched to corresponding detoxifying SMARTS.
14
15 SMARTS matching and identification of matched atoms were carried out as per **Figure 1**. The
16
17 manual and automated assessments of Heat Map quality are based on the schemes explained
18
19 under “Quality assessment of Heat Map images”.
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

28166-06-5: SMARTS matches



Quality scores:
 Manual: Toxicophore match:
 Strong (2nd)
 Automated: Toxicophore match:
 0.73

Quality scores:
 Manual: Toxicophore match:
 None (5th)
 Automated: Toxicophore match:
 0.00



(A) Ames test: toxic

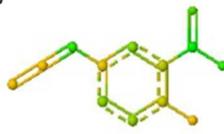
(B) Prediction: 0.48 (non-toxic)

(C) Prediction: 0.54 (toxic)

Quality scores:
 Manual: Toxicophore match:
 None (5th)
 Automated: Toxicophore match:
 0.00

Quality scores:
 Manual: Toxicophore match:
 Weaker (4th)
 Automated: Toxicophore match:
 0.63

Quality scores:
 Manual: Toxicophore match:
 Weaker (4th)
 Automated: Toxicophore match:
 0.59



(D) Prediction: 0.48 (non-toxic)

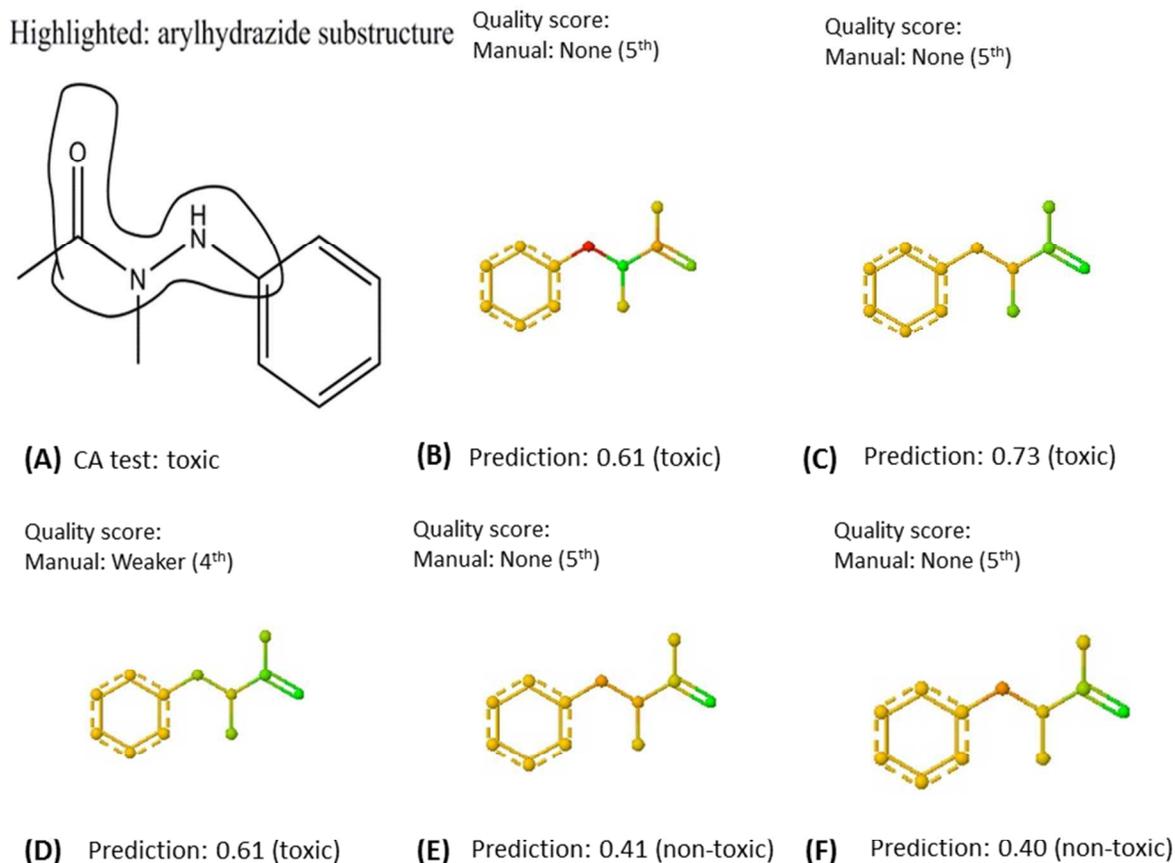
(E) Prediction: 0.82 (toxic)

(F) Prediction: 0.87 (toxic)

Figure 6 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule “28166-06-5” in the Kazius dataset. (A) Molecule with SMARTS matches shown in green for specific aromatic nitro (RHS) and azide (LHS) toxicophores, derived from the work of Kazius et al.⁴⁸ based on a combination of mechanistic reasoning, expert knowledge and data analysis reported, and experimental Ames test assignment reported in the Kazius dataset SDF file.¹³ The SMARTS matches (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per

Figure 4. The SMARTS matches show the atoms matched to SMARTS patterns. SMARTS

1
2
3 matching and identification of matched atoms were carried out as per **Figure 1**. The manual and
4
5 automated assessments of Heat Map quality are based on the schemes explained under “Quality
6
7 assessment of Heat Map images”. Quality scores, not to be confused with prediction scores,
8
9 based on Heat Map correspondence to detoxifying groups were not assigned as no detoxifying
10
11 groups were identified via SMARTS matching.
12
13
14
15
16



48
49 **Figure 7** Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of
50
51 binary classification leave-one-out predictions made for molecule CA31 in the CA dataset (c.f.
52
53 Figure 6(A) in Mohr et al.).¹⁷ (A) Molecule with biologically significant substructures shown,
54
55 based on a combination of mechanistic reasoning and expert knowledge reported previously in
56
57
58
59
60

1
2
3 the literature,¹⁷ and experimental chromosome aberration test assignment reported.¹⁷ (B – F)
4
5 Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged
6
7 predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D)
8
9 Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are
10
11 reported as per **Figure 4**. The manual assessments of Heat Map quality are based on the scheme
12
13 explained under “Quality assessment of Heat Map images”. Here, the assignments were not
14
15 simply based on matching the entire putative toxicophore,¹⁷ but were informed by consideration
16
17 of the likely toxic substructure revealed via metabolism, which is expected to be the
18
19 arylhydrazine substructure and may even be the aromatic amine substructure arising from further
20
21 metabolism.^{17,48}
22
23
24
25
26
27

28 **Distributions of numeric Heat Map quality scores for different interpretation approaches**

29
30 The distributions of overall quality scores, calculated as the average of scores based on Heat Map
31
32 correspondence to SMARTS pattern assigned toxicophores and detoxifying substructures as per
33
34 **Figure 1**, for all Kazius dataset Heat Map images generated according to atom coloring and
35
36 either Symmetrized Single Molecule Normalization or Constant Symmetrized Normalization are
37
38 shown in **Figure 8** (A) and (B) respectively. As indicated by **Figure 5**, Constant Symmetrized
39
40 Normalization tended to produce yellow coloring, yielding no information, which results in low
41
42 quality scores. Hence, as previously noted, Symmetrized Single Molecule Normalization is
43
44 recommended and no further consideration is given to the results obtained from Constant
45
46 Symmetrized Normalization.
47
48
49
50

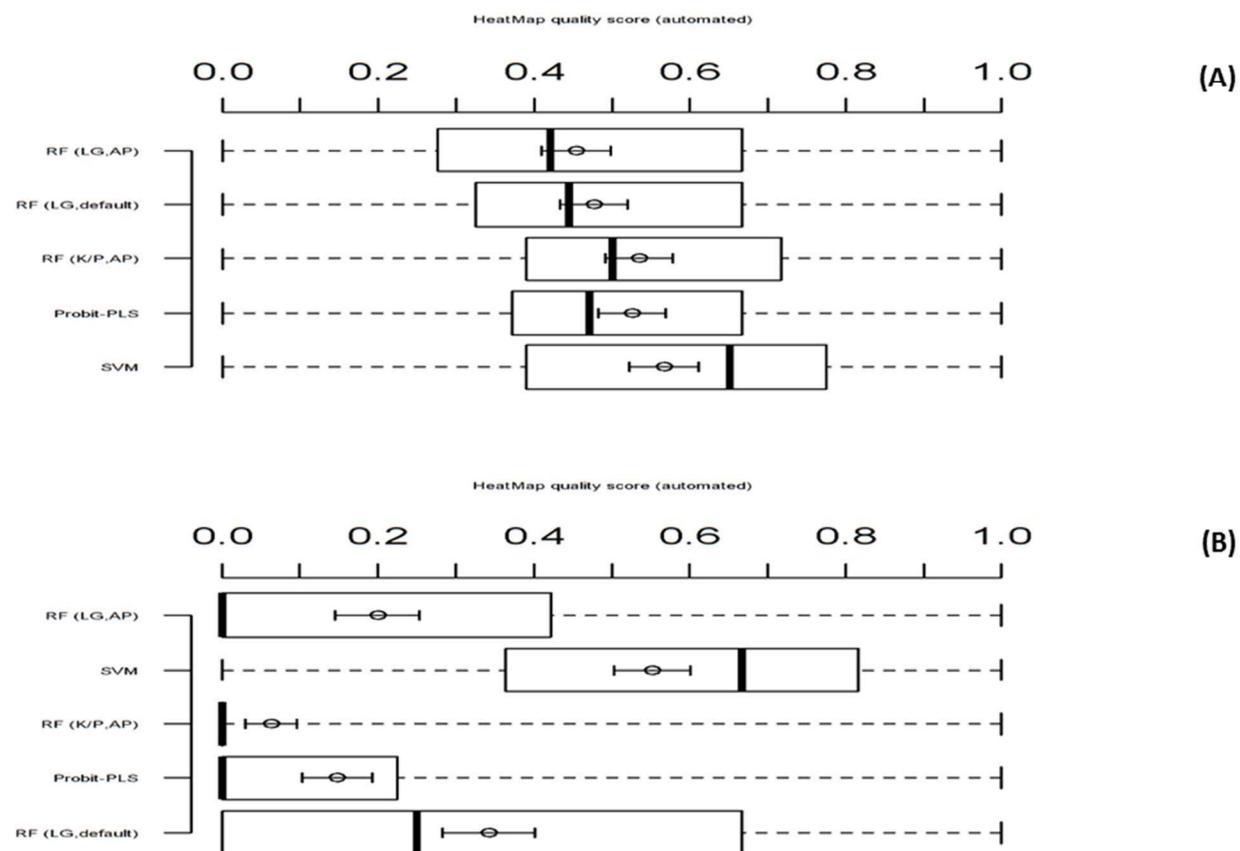
51
52 Furthermore, since only one of the molecules for which numeric Heat Map quality scores were
53
54 assigned was associated with a detoxifying SMARTS pattern match, the use of the average of the
55
56 scores calculated as per **Figure 1** is not representative of the distribution. Hence, this single
57
58
59
60

1
2
3 molecule was excluded from further analysis and all subsequent distributions reflect the quality
4 scores calculated based on Heat Map correspondence to the identified toxicophores for the
5 remaining 38 molecules, or relevant subsets thereof. As expected, removing this single molecule
6 has little impact upon the overall distribution, as can be seen by comparing **Figure 8 (A)** and
7
8
9
10
11
12
13 **Figure 9 (A)**.

14
15
16 Consider the median quality scores presented in **Figure 9 (A)**. These suggest that the degree to
17 which Heat Map images, based on atom coloring and Symmetrized Single Molecule
18 Normalization, can be meaningfully interpreted was typically somewhat higher for those
19 corresponding to linear SVM predictions compared to all other interpretation approaches. They
20 also suggest the Kuz'min/Palczewska approach to calculating descriptor contributions for
21 Random Forest classification predictions yielded Heat Map images with slightly higher median
22 quality scores than the local gradients approach. The trend in arithmetic mean quality scores
23 followed the same pattern.
24
25
26
27
28
29
30
31
32
33
34
35

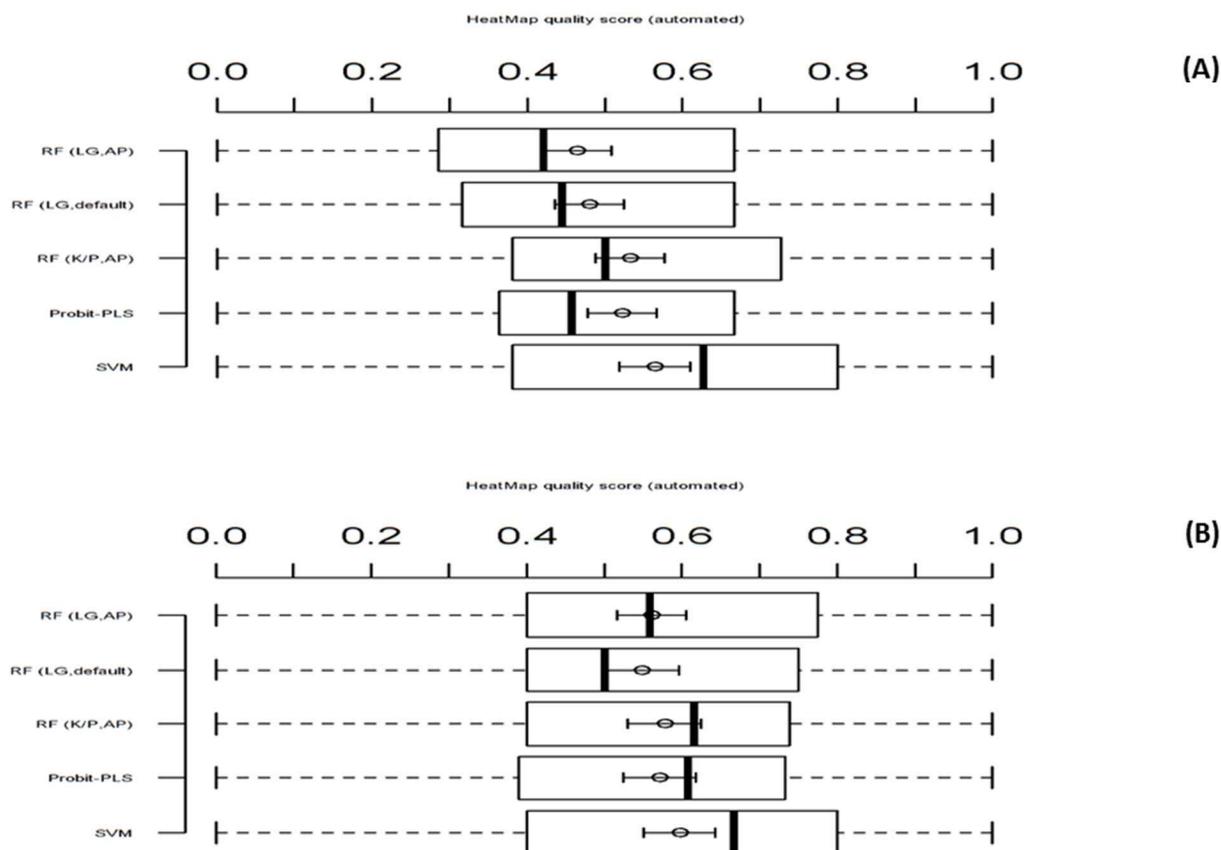
36 However, it might be argued that the quality scores of Heat Map images associated with toxic
37 predictions are of greater relevance, as these indicate the extent to which toxicophores are
38 correctly identified when a toxic prediction is made. A Heat Map image which fully and
39 precisely highlighted the toxicophores in combination with a prediction of toxicity, as per **Figure**
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
11 (F), would be valuable to medicinal chemists seeking to design out a predicted toxic liability
during lead optimization. (Of course, this is also reliant on the prediction of toxicity being
reliable as per **Figure 11 (F)**.) Consider the distributions of quality scores, reflecting the degree
to which Heat Map images associated with toxic predictions correctly identified the
toxicophore(s), in **Figure 10 (A)**. These indicate roughly the same rank ordering as per **Figure 9**
(A), except that the ordering of quality scores associated with Heat Map images corresponding to

1
2
3 different kinds of Random Forest classification local gradients descriptor contributions is
4
5
6 reversed.
7
8
9



41 **Figure 8** Automated Heat Map numeric quality scores, calculated as the average of F-measures
42 for their correspondence to toxicophore and detoxifying substructure SMARTS matches as
43 explained in **Figure 1**, for Heat Map images generated using atom coloring and (A)
44 Symmetrized Single Molecule Normalization or (B) Constant Symmetrized Normalization.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 underlying descriptor contributions: RF (LG, AP) = Random Forest classification (local
4 gradients, averaged predictions); RF (LG, default) = Random Forest classification (local
5 gradients, averaged predictions); RF (LG, default) = Random Forest classification (local
6 gradients, majority votes predictions); RF (K/P, AP) = Random Forest classification
7 (Kuz'min/Palczewska, averaged predictions); Probit-PLS (linear coefficients); SVM (linear
8 coefficients). The distributions are represented as boxplots, with the bold lines showing the
9 medians and the superimposed circle the arithmetic means, with error bars denoting the standard
10 errors in the means.
11
12
13
14
15
16
17
18
19
20
21



53 **Figure 9** Automated Heat Map numeric quality scores, based on their correspondence to
54 toxicophore substructure SMARTS matches as explained in **Figure 1**, for Heat Map images
55
56
57
58
59
60

1
2
3 generated using atom coloring and Symmetrized Single Molecule Normalization. In sub-figure
4
5 (A), quality score distributions correspond to 38 leave-one-out predictions made for a subset of
6
7 Kazius dataset molecules identified via toxicophore SMARTS matches, excluding the single
8
9 molecule (“1028-11-1”) for which a detoxifying substructure was also identified via SMARTS
10
11 pattern matches. In sub-figure (B), the distributions correspond to the subsets of those 38
12
13 molecules which were correctly predicted. The number of molecules correctly predicted to be
14
15 toxic or non-toxic varied across methods: Random Forest classification (averaged predictions) =
16
17 28; Random Forest classification (majority vote predictions) = 29; Probit-PLS = 28; SVM = 31.
18
19 The different quality scores distributions are annotated according to the different approaches
20
21 used to obtain the underlying descriptor contributions: RF (LG, AP) = Random Forest
22
23 classification (local gradients, averaged predictions); RF (LG, default) = Random Forest
24
25 classification (local gradients, majority votes predictions); RF (K/P, AP) = Random Forest
26
27 classification (Kuz’min/Palczewska, averaged predictions); Probit-PLS (linear coefficients);
28
29 SVM (linear coefficients). The distributions are represented as boxplots, with the bold lines
30
31 showing the medians and the superimposed circle the arithmetic means, with error bars denoting
32
33 the standard errors in the means.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

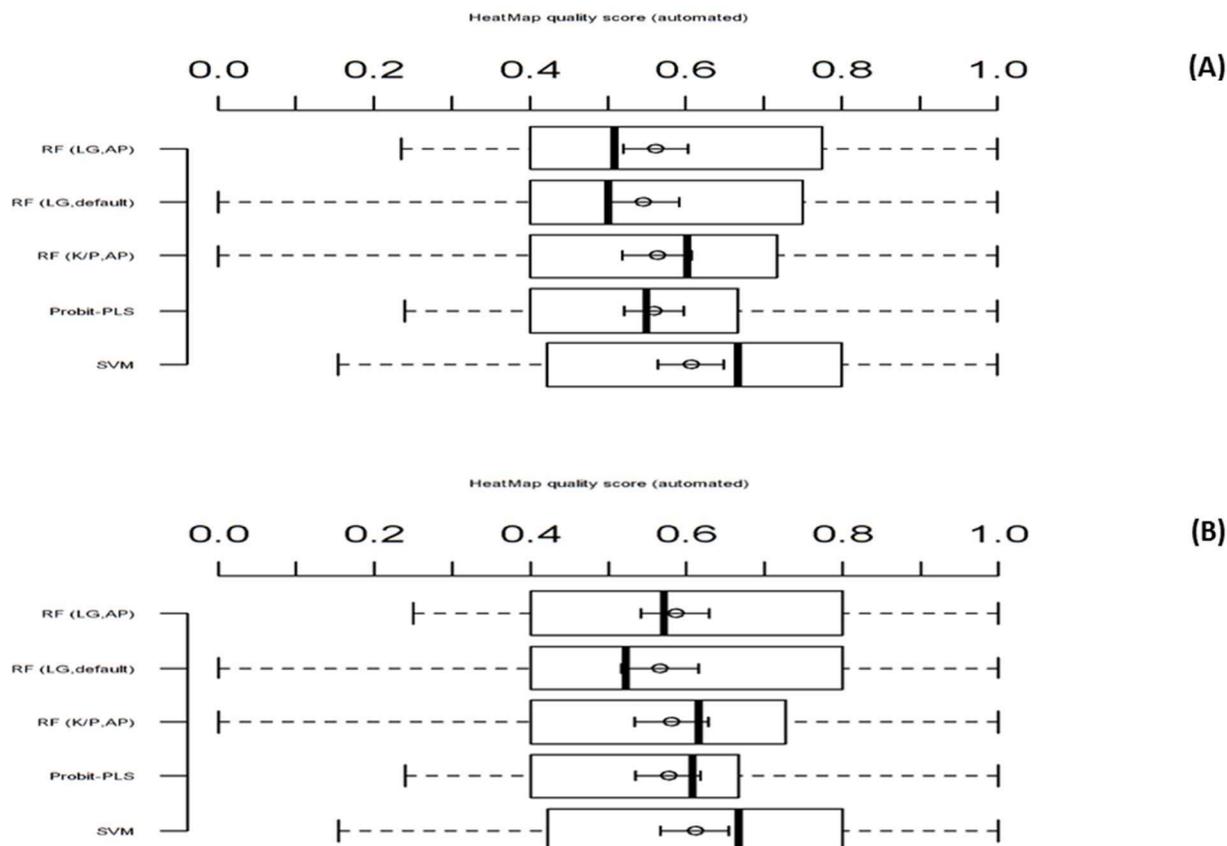


Figure 10 Automated Heat Map numeric quality scores, calculated based on their correspondence to toxicophore SMARTS matches as explained in **Figure 1**, for Heat Map images generated using atom coloring and Symmetrized Single Molecule Normalization. In sub-figure (A), the distributions of quality scores correspond to all toxic leave-one-out predictions made for a subset of Kazius dataset molecules identified via toxicophore SMARTS matches, excluding the single molecule (“1028-11-1”) for which a detoxifying substructure was also identified via SMARTS pattern matches. In sub-figure (B), the distributions of quality scores correspond to the subsets of those toxic predictions which were correct. The number of molecules (correctly) predicted to be toxic, hence the number of Heat Map images for which these distributions were generated, varied across methods: Random Forest classification

1
2
3 (averaged predictions) = 28 (25 correct); Random Forest classification (majority vote
4
5 predictions) = 29 (26 correct); Probit-PLS = 30 (26 correct); SVM = 31 (28 correct). The
6
7 different quality scores distributions are annotated according to the different approaches used to
8
9 obtain the underlying descriptor contributions: RF (LG, AP) = Random Forest classification
10
11 (local gradients, averaged predictions); RF (LG, default) = Random Forest classification (local
12
13 gradients, majority votes predictions); RF (K/P, AP) = Random Forest classification
14
15 (Kuz'min/Palczewska, averaged predictions); Probit-PLS (linear coefficients); SVM (linear
16
17 coefficients). The distributions are represented as boxplots, with the bold lines showing the
18
19 medians and the superimposed circle the arithmetic means, with error bars denoting the standard
20
21 errors in the means.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

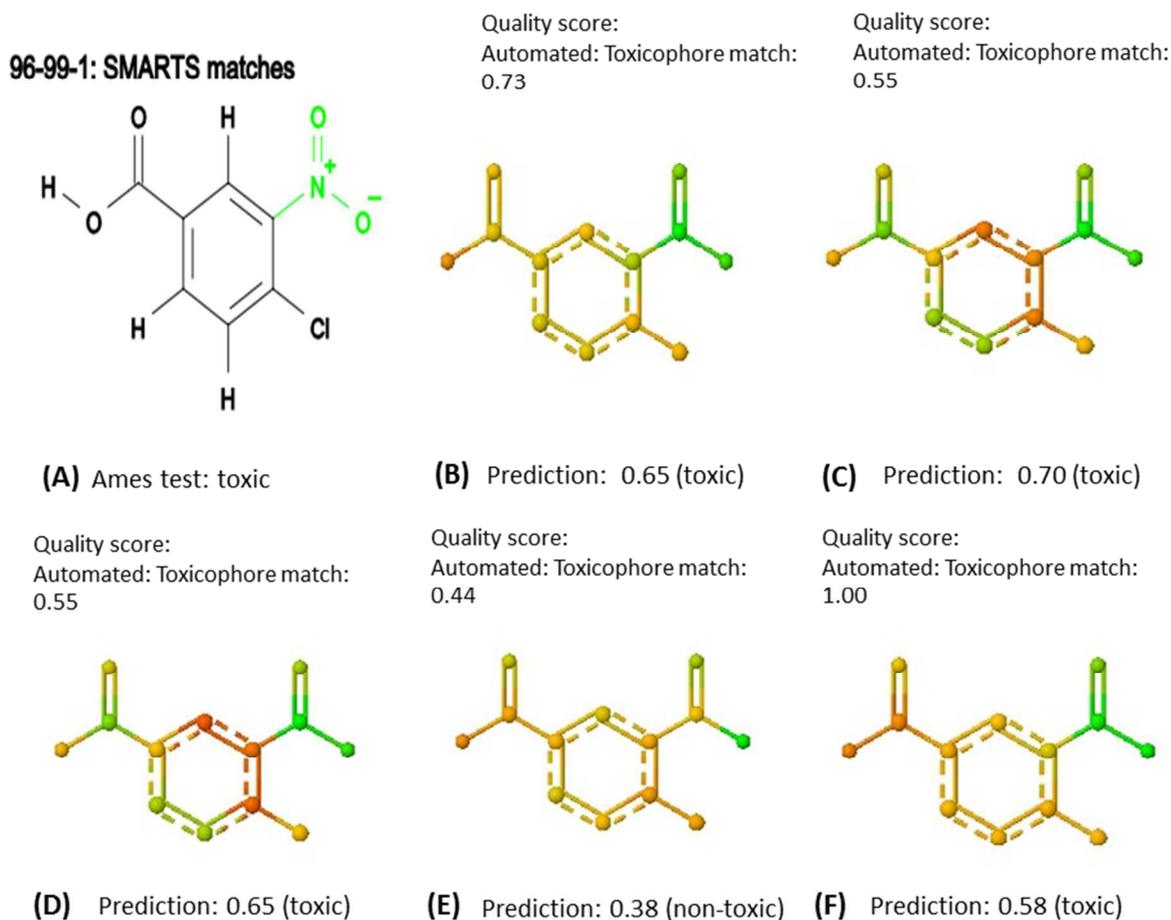


Figure 11 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule “96-99-1” in the Kazius dataset. (A) Molecule with SMARTS matches shown in green for specific aromatic nitro toxicophore, derived from the work of Kazius et al.⁴⁸ based on a combination of mechanistic reasoning, expert knowledge and data analysis, and experimental Ames test assignment reported in the Kazius dataset SDF file.¹³ The SMARTS matches (B – F) Heat Map images: (B) Random Forest classification (Kuz’min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per **Figure 4**. The SMARTS matches show the atoms matched to SMARTS patterns. SMARTS matching and

1
2
3 identification of matched atoms were carried out as per **Figure 1**. The automated assessments of
4 Heat Map quality are based on the scheme explained in **Figure 1**. Quality scores, not to be
5 confused with prediction scores, based on Heat Map correspondence to detoxifying groups were
6 not assigned as no detoxifying groups were identified via SMARTS matching.
7
8
9
10
11
12

13 14 **Effect of prediction quality on Heat Map quality**

15
16 **Figure 9 (B)** and **Figure 10 (B)** correspond to **Figure 9 (A)** and **Figure 10 (A)** respectively,
17 after removing incorrect predictions. The corresponding average quality scores are summarized
18 in Supporting Information tables S7 (corresponding to **Figure 9 (A)** and **Figure 9 (B)**) and S8
19 (corresponding to **Figure 10 (A)** and **Figure 10 (B)**). With the exception of SVM toxic
20 predictions, there is an increase in the average numeric quality score, based on Heat Map
21 correspondence to toxicophores, upon moving from all to correct and all toxic to correct toxic
22 predictions. However, especially in the case of toxic predictions, this change is often quite small
23 and the numbers of incorrect predictions are rather small, which may mean this finding is not
24 robust. Nonetheless, these findings do suggest that in the case of an experimentally confirmed
25 prediction, or in light of strong statistical grounds for believing a specific prediction is correct,
26 the substructural cause of toxicity suggested by the (atom coloring, Symmetrized Single
27 Molecule Normalization) Heat Map interpretation of the predictions may be considered more
28 trustworthy. Hence, under these circumstances, the Heat Map image may be considered more
29 useful for guiding lead optimization or suggesting the mechanism of toxic action.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Differences between Random Forest interpretation approaches

Our work revealed that different approaches for interpreting Random Forest predictions may sometimes yield different interpretations. This can be observed by contrasting the Heat Map images (atom coloring, Symmetrized Single Molecule Normalization) corresponding to the same Random Forest classification (averaged predictions) prediction using either the Kuz'min/Palczewska (sub-figure B) or local gradients approach (sub-figure D) in the following figures: **Figure 4**, **Figure 6**, **Figure 7**, **Figure 11**. Consideration of **Figure 12**, which presents raw atom score estimates corresponding to **Figure 11**, emphasizes that this is not merely an artefact of the Symmetrized Single Molecule Normalization scheme: the raw atom scores, derived from the different kinds of descriptor contributions, are clearly not identical.

Consideration of **Figure 12** shows that the estimated influence of specific atoms on the prediction may, indeed, change sign. This could simply reflect the fact that the local gradients approach is not as good as the Kuz'min/Palczewska approach for estimating descriptor influences on the prediction, as may be indicated by consideration of the somewhat lower median Heat Map quality scores shown in **Figure 8 (A)**, **Figure 9**, **Figure 10**. As to why this should be the case, we speculate that the manner in which local gradients are estimated in our work (see equation (5)) fails to take account of the possibility that switching the value of one binary descriptor might affect the manner in which the molecule of interest is passed down the trees of the forest. This could affect how other descriptors end up contributing to the prediction. This is a potential weakness compared to the Kuz'min/Palczewska approach, which merely takes account of the influence of the specific descriptor of interest on the prediction.

Nonetheless, in spite of these *potential* limitations of the local gradients approach compared to the Kuz'min/Palczewska approach, it should be acknowledged that the correlation between

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

estimated raw atom scores was typically higher than observed for the example presented in (Figure 12). This can be seen by comparing the Pearson's correlation coefficient (0.74) between these two sets of estimated raw atom scores and the distribution of correlation coefficients summarized in Figure 13 (arithmetic mean = 0.79, median = 0.93).

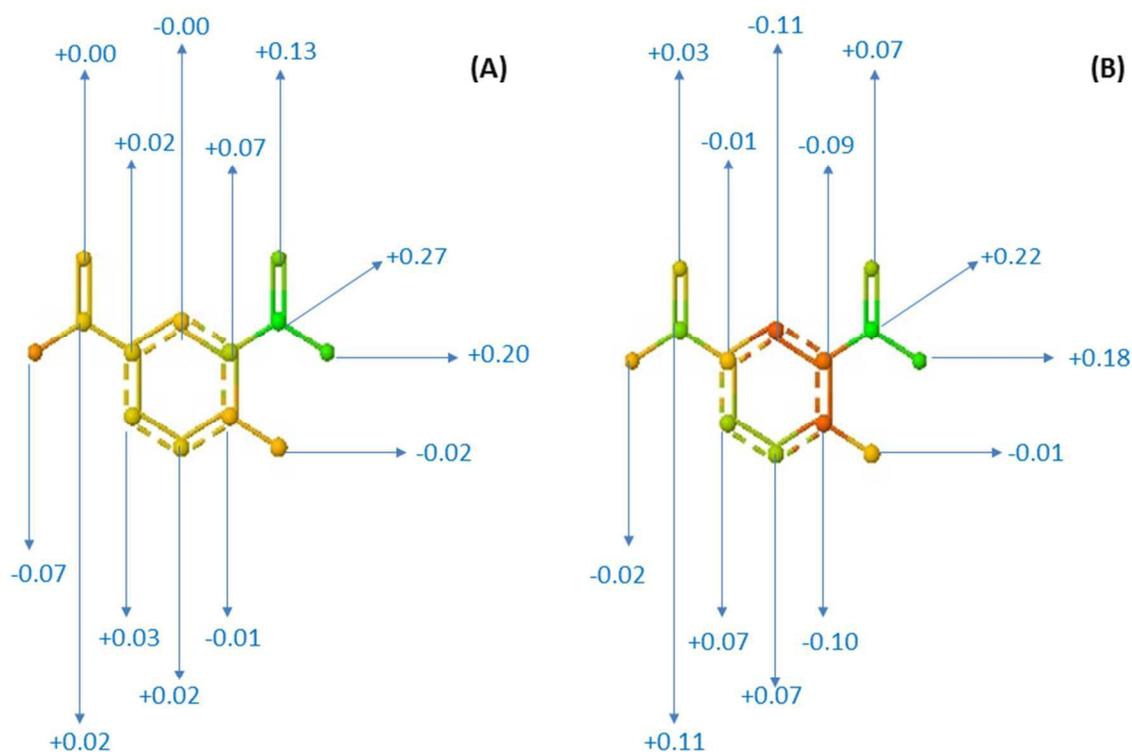
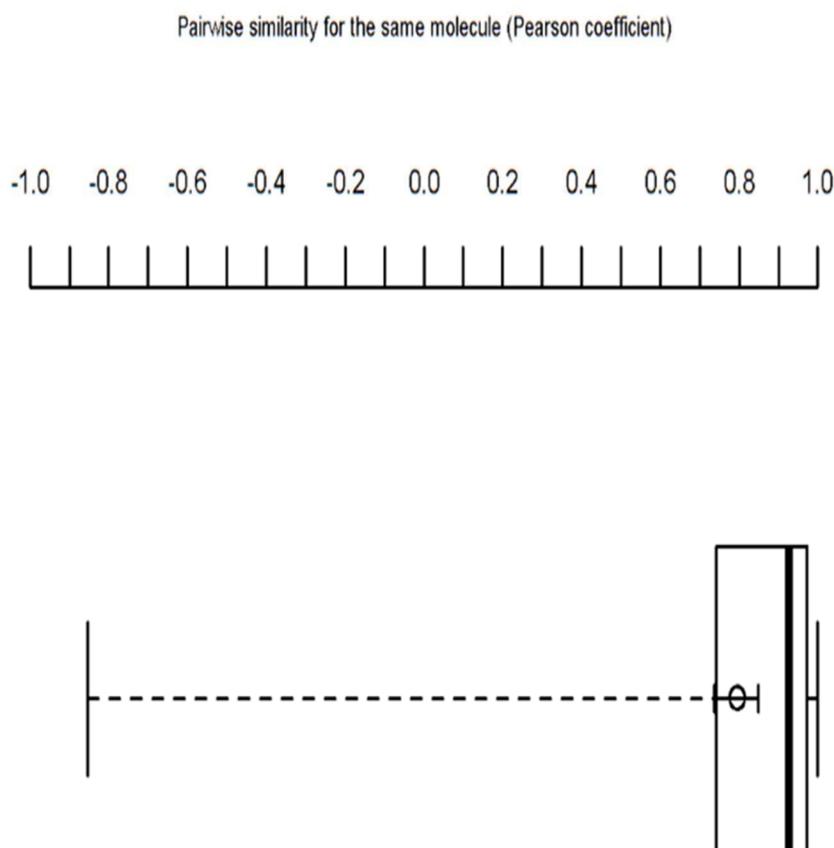


Figure 12 Comparison between the estimated raw atom scores (2dp), for the molecule "96-99-1" in the Kazius dataset, derived from the descriptor contributions corresponding to a single leave-one-out prediction made using Random Forest classification (averaged predictions) according to the Kuz'min/Palczewska (A) and local gradients (B) approaches. The raw scores were estimated

1
2
3 via processing output from the HeatMapWrapper tool⁷² in atom coloring and Symmetrized
4
5 Single Molecule Normalization mode. The correlation between the corresponding raw atom
6
7 score estimates, in terms of the Pearson's correlation coefficient, is 0.74. The raw atom score
8
9 estimates are used to annotate the corresponding atoms in the Heat Map (atom coloring &
10
11 Symmetrized Single Molecule Normalization) image. N.B. The specific atom IDs, used to
12
13 identify the corresponding atom score estimates, were identified as per **Figure 1**.
14
15
16
17
18
19



51 **Figure 13** The distribution of pairwise similarities, calculated via the Pearson's correlation
52
53 coefficient using estimates of the raw atom scores derived from descriptor contributions
54
55 generated via the Kuz'min/Palczewska and local gradients approaches, for all leave-one-out
56
57
58
59
60

1
2
3 Random Forest classification predictions (averaged predictions approach) made for a subset of
4 molecules in the Kazius dataset identified via toxicophore SMARTS pattern matches. (The raw
5 atom score estimates were obtained via processing the output of the HeatMapWrapper
6 program⁷² in atom coloring and Symmetrized Single Molecule Normalization mode.) The
7 distribution is represented as a boxplot, with the bold line showing the median and the
8 superimposed circle the arithmetic mean, with error bars denoting the standard error in the mean.
9
10
11
12
13
14
15
16
17
18
19

20 CONCLUSIONS

21
22 In this paper, a comparison was presented between the widely used non-linear Random Forest
23 algorithm and well-known linear modelling approaches: Probit-PLS, PLS-regression, Support
24 Vector Machines and Support Vector Regression, with the latter two approaches yielding linear
25 models when employed with a linear kernel function. Specifically, these algorithms were
26 compared in terms of their ability to build predictive models, based on commonly employed
27 extended connectivity fingerprints, for numerical and categorical measures of biological activity
28 on a variety of established benchmark datasets. As well as comparing the predictive performance
29 of the models, using external cross-validation, the interpretability of individual predictions was
30 also assessed.
31
32
33
34
35
36
37
38
39
40
41
42
43

44 For Random Forest, different approaches for generating binary classification predictions and for
45 interpreting individual predictions were investigated. To the best of our knowledge, our work
46 represents the first time that methodologically distinct approaches for interpreting predictions
47 obtained from a non-linear QSAR model have been compared. Predictions were interpreted via
48 converting the corresponding substructural descriptor contributions, obtained via different
49 algorithms, into colored Heat Map molecular images, as previously proposed. We demonstrated
50
51
52
53
54
55
56
57
58
59
60

1
2
3 that different ways of translating descriptor contributions into Heat Map images can yield
4
5 different interpretations and we recommend the use of a novel approach to Heat Map generation:
6
7 atom coloring combined with Symmetrized Single Molecule Normalization. Another novelty of
8
9 our work is the introduction of systematic assessment schemes for assessing the quality of Heat
10
11 Map molecular images as a means of assessing the extent to which the corresponding predictions
12
13 can be interpreted in a chemically and biologically meaningful fashion. These images are
14
15 ultimately intended to provide information to chemists, hence require manual inspection.
16
17
18 However, we advocate the use of our automated, numeric scoring scheme for systematic
19
20 evaluation of trends in the degree to which Heat Maps generated according to different protocols
21
22 can be meaningfully interpreted. A manual ranking scheme, applied to an initial set of Heat Map
23
24 prediction interpretations corresponding to models based upon a variety of datasets, was
25
26 determined to be too problematic for systematic evaluation of these images. However, due to
27
28 implementation challenges and the requirement for molecular substructures responsible for
29
30 (attenuating) biological activity to be encoded as SMARTS patterns, it was only practical to
31
32 apply the automated scoring scheme to a subset of leave-one-out binary classification predictions
33
34 obtained for an Ames mutagenicity dataset.
35
36
37
38
39

40
41
42 Random Forest classification was observed to produce better predictive performance, on
43
44 average, for most of the analyzed datasets. The non-default “averaged predictions” approach to
45
46 generating classification predictions was typically observed to produce better predictive
47
48 performance, on average, compared to the default “majority votes” approach. Random Forest
49
50 regression was observed to less consistently outperform the other regression approaches.
51
52
53 However, differences in performance were sometimes marginal and may not be robust.
54
55
56
57
58
59
60

1
2
3 This work emphasizes that predictions obtained from both the non-linear (Random Forest) and
4 linear modelling algorithms considered here *can* be meaningfully interpreted via identifying
5 regions of the molecule expected to be responsible for biological activity. A detailed evaluation
6 was performed of Heat Map images obtained via a variety of interpretation approaches, each
7 interpretation approach corresponding to a different combination of modelling algorithm and, as
8 applicable, method for generating predictions and corresponding descriptor contributions from
9 the model. Our evaluation was based on the distribution of numeric Heat Map quality scores
10 obtained for leave-one-out predictions of Ames mutagenicity, denoting the extent to which those
11 images correspond to the known toxicophores. Considering those distributions, the usefulness of
12 Heat Map images corresponding to linear SVM predictions appeared somewhat higher on
13 average than all other interpretation approaches and the Kuz'min/Palczewska approach to
14 calculating descriptor contributions for Random Forest classification predictions (averaged
15 predictions) yielded Heat Map images with slightly higher average quality scores than the local
16 gradients approach. At least in terms of the median quality scores, a similar pattern was observed
17 upon consideration of toxic predictions, for which Heat Map ability to identify toxicophores
18 would be useful for medicinal chemists seeking to design out a toxic liability during lead
19 optimization. These trends were also observed when only correct, or correct toxic, predictions
20 were considered.

21
22 Improved prediction quality seemed to typically correspond to improved interpretability, in terms
23 of the correspondence between the Heat Map images and known toxicophores, based upon
24 considering the same set of leave-one-out predictions of Ames mutagenicity. However, increases
25 in average Heat Map quality scores were typically small, upon moving from all (toxic) to correct
26 (toxic) predictions and only a small number of incorrect (toxic) predictions were considered.
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Nonetheless, these findings suggest that, if the corresponding predictions were experimentally verified, there were even stronger grounds for believing the Heat Map images could be trusted to offer at least some insight into an unknown mechanism of action or opportunities for lead optimization. This also suggests that, where there are statistical grounds for high confidence in the prediction of an untested molecule, there are stronger grounds for drawing conclusions from the corresponding Heat Map image.

Overall, we can advocate Random Forest as a suitable option for QSAR modelers interested in both accurate and interpretable predictions. In keeping with the literature, our results suggest Random Forest may yield better predictive performance than linear methods. They also indicate that our Open Source *rfFC* and *HeatMapWrapper* software tools may allow for chemically and biologically meaningful interpretations of Random Forest predictions. For the purpose of generating Heat Map molecular interpretations of those predictions, we advocate the use of atom coloring, in combination with our Symmetrized Single Molecule Normalization scheme, and tentatively suggest the Kuz'min/Palczewska approach may be preferable to the local gradients approach considered in our work.

Supporting Information

Supporting Information Available: (1) "SupportingInformation_doc_JCIM_rev.2.ret.pdf" contains three sections: (A) describes the methods used and computational analysis in greater detail; (B) explains how we generated our results; (C) summarizes additional results, including raw results provided in "additional.results.resub.fixed.zip". (2) "FinalImagesAnalysis_withImages_resub.xlsx" presents the full analysis of all Heat Map images and corresponding predictions. (3) "repeating_ligand_bioactivity_searches.xlsx" documents how

1
2
3 experimental pK_i and pIC_{50} values were obtained to derive the values reported in
4
5 “FinalImagesAnalysis_withImages_resub.xlsx”. (4) “ligands_released.zip” presents copies of the
6
7 ligand files used to generate Heat Map images. (5) “additional.results.resub.fixed.zip” provides
8
9 additional raw results files. (6) “Kazius_2005_SI_SA_SMARTS_subset_v2.xls” provides the
10
11 SMARTS patterns for toxicophores and detoxifying moieties used for automated, numeric
12
13 scoring of Kazius LOO Heat Map images. This material is available free of charge via the
14
15 Internet at <http://pubs.acs.org>.
16
17
18
19
20
21

22 Author Contributions

23
24 RLMR conceived of the project, designed the computational experiments, wrote and executed
25
26 the necessary scripts, analyzed the results and wrote the first draft of the manuscript. NK helped
27
28 design the Pipeline Pilot workflow for standardizing molecular structures and prepared all ligand
29
30 SDF files from PDB files for analysis using Schrödinger’s Protein Preparation Wizard. AP and
31
32 JP provided feedback which improved the analysis carried out by RLMR and assisted with
33
34 writing some scripts. AP developed the *rfFC* package, with some conceptual input from RLMR
35
36 and JP. RLMR wrote the *HeatMapWrapper* Python tool used for generating Heat Map images
37
38 analyzed in the current publication. JP assisted with the design and implementation of the
39
40 Symmetrized Single Molecule Normalization scheme. The final manuscript was written through
41
42 contributions of all authors. All authors have given approval to the final version of the
43
44 manuscript.
45
46
47
48
49
50
51

52 Funding Sources

53
54 N/A
55
56
57
58
59
60

Notes

N/A

ACKNOWLEDGEMENT

RLMR thanks Kim Travis (Syngenta Ltd.), Mark Earll (Syngenta Ltd.) and Mark Forster (Syngenta Ltd.) for useful discussions. RLMR thanks Lars Rosenbaum (University of Tübingen) for providing guidance on how to generate the Heat Map images reported in this publication. RLMR also thanks Andreas Zell (University of Tübingen) for further helpful correspondence regarding the Heat Map software. RLMR thanks Johannes Mohr (Technische Universität Berlin) for valuable correspondence regarding the CA dataset. RLMR thanks Michael Mackay (Liverpool John Moores University) for providing access to a virtual machine used to run cross-validation calculations. RLMR thanks Rachel Kramer Green (RCSB PDB Leadership Team), Heather Carlson (University of Michigan) and Renxiao Wang (Shanghai Institute of Organic Chemistry) for helpful correspondence regarding the PDB, BindingMOAD and PDBbind databases respectively. RLMR thanks Knut Baumann (Technische Universität Braunschweig) for providing access to the MUV datasets. The authors thank Lucy Entwistle (Syngenta Ltd.) for writing the Pipeline Pilot workflow used to process all molecules prior to fingerprint calculations.

ABBREVIATIONS

N/A

REFERENCES

- (1) *Guidance Document on the Validation of (Quantitative)Structure-Activity Relationships [(Q)SAR] Models (ENV/JM/MONO(2007)2)*; OECD Series on Testing and Assessment; 69; Organisation for Economic Co-operation and Development, 2007.
- (2) Cherkasov, A.; Muratov, E. N.; Fourches, D.; Varnek, A.; Baskin, I. I.; Cronin, M.; Dearden, J.; Gramatica, P.; Martin, Y. C.; Todeschini, R.; et al. QSAR Modeling: Where Have You Been? Where Are You Going To? *J. Med. Chem.* **2014**, *57*, 4977–5010.
- (3) Gavaghan, C. L.; Arnby, C. H.; Blomberg, N.; Strandlund, G.; Boyer, S. Development, Interpretation and Temporal Evaluation of a Global QSAR of hERG Electrophysiology Screening Data. *J. Comput.-Aided Mol. Des.* **2007**, *21* (4), 189–206.
- (4) Hansch, C.; Fujita, T. P- σ - π Analysis. A Method for the Correlation of Biological Activity and Chemical Structure. *J. Am. Chem. Soc.* **1964**, *86* (8), 1616–1626.
- (5) Baskin, I. I.; Ait, A. O.; Halberstam, N. M.; Palyulin, V. A.; Zefirov, N. S. An Approach to the Interpretation of Backpropagation Neural Network Models in QSAR Studies. *SAR QSAR Environ. Res.* **2002**, *13* (1), 35–41.
- (6) Johnson, S. R.; Chen, Q. X.; Murphy, D.; Gudmunsson, O. A Computational Model for the Prediction of Aqueous Solubility That Includes Crystal Packing, Intrinsic Solubility, and Ionization Effects. *Mol. Pharmaceutics* **2007**, *4* (4), 513–523.
- (7) Dearden, J. C.; Cronin, M. T. D.; Kaiser, K. L. E. How Not to Develop a Quantitative Structure–activity or Structure–property Relationship (QSAR/QSPR). *SAR QSAR Environ. Res.* **2009**, *20* (3–4), 241–266.
- (8) Jensen, F. Chapter 17: Statistics and QSAR. In *Introduction to Computational Chemistry*; John Wiley & Sons Ltd, 2007; pp 547–561.
- (9) Svetnik, V.; Liaw, A.; Tong, C.; Culberson, J. C.; Sheridan, R. P.; Feuston, B. P. Random Forest: A Classification and Regression Tool for Compound Classification and QSAR Modeling. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (6), 1947–1958.
- (10) Muller, K.; Mika, S.; Ratsch, G.; Tsuda, K.; Scholkopf, B. An Introduction to Kernel-Based Learning Algorithms. *IEEE Transactions on Neural Networks* **2001**, *12* (2), 181–201.
- (11) Hsu, C.-W.; Chang, C.-C.; Lin, C.-J. *A Practical Guide to Support Vector Classification*; Department of Computer Science, National Taiwan University: Taiwan, 2010.
- (12) Smola, A. J.; Schölkopf, B. A Tutorial on Support Vector Regression. *Statistics and Computing* **2004**, *14* (3), 199–222.
- (13) Rosenbaum, L.; Hinselmann, G.; Jahn, A.; Zell, A. Interpreting Linear Support Vector Machine Models with Heat Map Molecule Coloring. *J. Cheminf.* **2011**, *3* (1), 11.
- (14) Song, M.; Clark, M. Development and Evaluation of an in Silico Model for hERG Binding. *J. Chem. Inf. Model* **2006**, *46* (1), 392–400.
- (15) Mitchell, J. B. O. Machine Learning Methods in Chemoinformatics. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2014**, *4*, 468–481.
- (16) Guha, R. On the Interpretation and Interpretability of Quantitative Structure–activity Relationship Models. *J. Comput.-Aided Mol. Des.* **2008**, *22* (12), 857–871.
- (17) Mohr, J.; Jain, B.; Sutter, A.; Laak, A. T.; Steger-Hartmann, T.; Heinrich, N.; Obermayer, K. A Maximum Common Subgraph Kernel Method for Predicting the Chromosome Aberration Test. *J. Chem. Inf. Model.* **2010**, *50* (10), 1821–1838.
- (18) Carlsson, L.; Helgee, E. A.; Boyer, S. Interpretation of Nonlinear QSAR Models Applied to Ames Mutagenicity Data. *J. Chem. Inf. Model.* **2009**, *49* (11), 2551–2558.

- 1
2
3 (19) Kuz'min, V. E.; Polishchuk, P. G.; Artemenko, A. G.; Andronati, S. A. Interpretation of
4 QSAR Models Based on Random Forest Methods. *Mol. Inf.* **2011**, *30* (6-7), 593–603.
- 5 (20) An, Y.; Sherman, W.; Dixon, S. L. Kernel-Based Partial Least Squares: Application to
6 Fingerprint-Based QSAR with Model Visualization. *J. Chem. Inf. Model.* **2013**, *53*, 2312–
7 2321.
- 8 (21) Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. Universal Approach
9 for Structural Interpretation of QSAR/QSPR Models. *Mol. Inf.* **2013**, *32* (9–10), 843–853.
- 10 (22) Webb, S. J.; Hanser, T.; Howlin, B.; Krause, P.; Vessey, J. D. Feature Combination
11 Networks for the Interpretation of Statistical Machine Learning Models: Application to
12 Ames Mutagenicity. *J. Cheminf.* **2014**, *6* (1), 8.
- 13 (23) Welling, S. H.; Clemmensen, L. K. H.; Buckley, S. T.; Hovgaard, L.; Brockhoff, P. B.;
14 Refsgaard, H. H. F. In Silico Modelling of Permeation Enhancement Potency in Caco-2
15 Monolayers Based on Molecular Descriptors and Random Forest. *Eur. J. Pharm.*
16 *Biopharm.* **2015**, *94*, 152–159.
- 17 (24) Balfer, J.; Bajorath, J. Visualization and Interpretation of Support Vector Machine
18 Activity Predictions. *J. Chem. Inf. Model.* **2015**, *55* (6), 1136–1147.
- 19 (25) Polishchuk, P.; Tinkov, O.; Khristova, T.; Ognichenko, L.; Kosinskaya, A.; Varnek, A.;
20 Kuz'min, V. Structural and Physico-Chemical Interpretation (SPCI) of QSAR Models and
21 Its Comparison with Matched Molecular Pair Analysis. *J. Chem. Inf. Model.* **2016**, *56* (8),
22 1455–1469.
- 23 (26) Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45* (1), 5–32.
- 24 (27) Marchese Robinson, R. L.; Glen, R. C.; Mitchell, J. B. O. Development and Comparison
25 of hERG Blocker Classifiers: Assessment on Different Datasets Yields Markedly
26 Different Results. *Mol. Inf.* **2011**, *30*, 443–458.
- 27 (28) Strobl, C.; Boulesteix, A.-L.; Zeileis, A.; Hothorn, T. Bias in Random Forest Variable
28 Importance Measures: Illustrations, Sources and a Solution. *BMC Bioinf.* **2007**, *8* (1), 25.
- 29 (29) Palczewska, A.; Palczewski, J.; Robinson, R. M.; Neagu, D. Interpreting Random Forest
30 Models Using a Feature Contribution Method. In *2013 IEEE 14th International*
31 *Conference on Information Reuse and Integration (IRI)*; 2013; pp 112–119.
- 32 (30) Palczewska, A.; Palczewski, J.; Robinson, R. M.; Neagu, D. Interpreting Random Forest
33 Classification Models Using a Feature Contribution Method. In *Integration of Reusable*
34 *Systems*; Bouabana-Tebibel, T., Rubin, S. H., Eds.; Advances in Intelligent Systems and
35 Computing; Springer International Publishing, 2014; pp 193–218.
- 36 (31) Welling, S. H.; Refsgaard, H. H. F.; Brockhoff, P. B.; Clemmensen, L. H. Forest Floor
37 Visualizations of Random Forests. *arXiv.org, e-Print Arch., Statistics* **2016**.
- 38 (32) Cassano, A.; Marchese Robinson, R. L.; Palczewska, A.; Puzyn, T.; Gajewicz, A.; Tran,
39 L.; Manganelli, S.; Cronin, M. T. D. Comparing the CORAL and Random Forest
40 Approaches for Modelling the in Vitro Cytotoxicity of Silica Nanomaterials. *ATLA,*
41 *Altern. Lab. Anim.* **2016**, *44* (6), 533–556.
- 42 (33) Stålring, J.; Almeida, P. R.; Carlsson, L.; Helgee Ahlberg, E.; Hasselgren, C.; Boyer, S.
43 Localized Heuristic Inverse Quantitative Structure Activity Relationship with Bulk
44 Descriptors Using Numerical Gradients. *J. Chem. Inf. Model.* **2013**, *53* (8), 2001–2017.
- 45 (34) Mevik, B.-H.; Wehrens, R. The Pls Package: Principal Component and Partial Least
46 Squares Regression in R. *J Stat Softw* **2007**, *18* (2).
- 47 (35) Obrezanova, O.; Segall, M. D. Gaussian Processes for Classification: QSAR Modeling of
48 ADMET and Target Activity. *J. Chem. Inf. Model.* **2010**, *50* (6), 1053–1061.
- 49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- (36) Agresti, A. Generalized Linear Models for Binary Data. In *An Introduction to Categorical Data Analysis*; Wiley Series in Probability and Statistics; WILEY-INTERSCIENCE, 2007; pp 68–72.
- (37) R-Forge: randomForest Feature Contribution: R Development Page (rffC package repository) https://r-forge.r-project.org/R/?group_id=1725 (accessed Apr 30, 2016).
- (38) Liu, Y.; Chawla, N. V.; Harper, M. P.; Shriberg, E.; Stolcke, A. A Study in Machine Learning from Imbalanced Data for Sentence Boundary Detection in Speech. *Comput. Speech Lang.* **2006**, *20* (4), 468–494.
- (39) Braga-Neto, U. M.; Dougherty, E. R. Is Cross-Validation Valid for Small-Sample Microarray Classification? *Bioinformatics* **2004**, *20* (3), 374–380.
- (40) Baldi, P.; Brunak, S.; Chauvin, Y.; Andersen, C. A. F.; Nielsen, H. Assessing the Accuracy of Prediction Algorithms for Classification: An Overview. *Bioinformatics*. **2000**, *16* (5), 412–424.
- (41) Gorodkin, J. Comparing Two K-Category Assignments by a K-Category Correlation Coefficient. *Comput. Biol. Chem.* **2004**, *28* (5–6), 367–374.
- (42) Konovalov, D. A.; Llewellyn, L. E.; Vander Heyden, Y.; Coomans, D. Robust Cross-Validation of Linear Regression QSAR Models. *J. Chem. Inf. Model.* **2008**, *48* (10), 2081–2094.
- (43) Alexander, D. L. J.; Tropsha, A.; Winkler, D. A. Beware of R²: Simple, Unambiguous Assessment of the Prediction Accuracy of QSAR and QSPR Models. *J. Chem. Inf. Model.* **2015**, *55* (7), 1316–1322.
- (44) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model* **2010**, *50* (5), 742–754.
- (45) Weininger, D.; Weininger, A.; Weininger, J. L. SMILES. 2. Algorithm for Generation of Unique SMILES Notation. *J. Chem. Inf. Comput. Sci.* **1989**, *29*, 97–101.
- (46) Sutherland, J. J.; O'Brien, L. A.; Weaver, D. F. A Comparison of Methods for Modeling Quantitative Structure–Activity Relationships. *J. Med. Chem.* **2004**, *47* (22), 5541–5554.
- (47) Rohrer, S. G.; Baumann, K. Maximum Unbiased Validation (MUV) Data Sets for Virtual Screening Based on PubChem Bioactivity Data. *J. Chem. Inf. Model.* **2009**, *49* (2), 169–184.
- (48) Kazius, J.; McGuire, R.; Bursi, R. Derivation and Validation of Toxicophores for Mutagenicity Prediction. *J. Med. Chem.* **2005**, *48* (1), 312–320.
- (49) Przybylak, K. R.; Madden, J. C.; Cronin, M. T. D.; Hewitt, M. Assessing Toxicological Data Quality: Basic Principles, Existing Schemes and Current Limitations. *SAR QSAR Environ. Res.* **2012**, *23* (5–6), 435–459.
- (50) Hinselmann, G.; Rosenbaum, L.; Jahn, A.; Fechner, N.; Zell, A. jCompoundMapper: An Open Source Java Library and Command-Line Tool for Chemical Fingerprints. *J. Cheminf* **2011**, *3* (1), 3.
- (51) Hinselmann, G.; Rosenbaum, L.; Jahn, A.; Fechner, N.; Ostermann, C.; Zell, A. Large-Scale Learning of Structure–Activity Relationships Using a Linear Support Vector Machine and Problem-Specific Metrics. *J. Chem. Inf. Model* **2011**, *51* (2), 203–213.
- (52) Maunz, A.; Gütlein, M.; Rautenberg, M.; Vorgrimmler, D.; Gebele, D.; Helma, C. Lazar: A Modular Predictive Toxicology Framework. *Front. Pharmacol* **2013**, *4*, 38.
- (53) Manchester, J.; Czermiński, R. CAUTION: Popular “Benchmark” Data Sets Do Not Distinguish the Merits of 3D QSAR Methods. *J. Chem. Inf. Model.* **2009**, *49* (6), 1449–1454.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- (54) DePriest, S. A.; Mayer, D.; Naylor, C. B.; Marshall, G. R. 3D-QSAR of Angiotensin-Converting Enzyme and Thermolysin Inhibitors: A Comparison of CoMFA Models Based on Deduced and Experimentally Determined Active Site Geometries. *J. Am. Chem. Soc.* **1993**, *115* (13), 5372–5384.
- (55) Böhm, M.; Stürzebecher, J.; Klebe, G. Three-Dimensional Quantitative Structure–Activity Relationship Analyses Using Comparative Molecular Field Analysis and Comparative Molecular Similarity Indices Analysis To Elucidate Selectivity Differences of Inhibitors Binding to Trypsin, Thrombin, and Factor Xa. *J. Med. Chem.* **1999**, *42* (3), 458–477.
- (56) Gohlke, H.; Klebe, G. DrugScore Meets CoMFA: Adaptation of Fields for Molecular Comparison (AFMoC) or How to Tailor Knowledge-Based Pair-Potentials to a Particular Protein. *J. Med. Chem.* **2002**, *45* (19), 4153–4170.
- (57) BIOVIA Pipeline Pilot | Scientific Workflow Authoring Application for Data Analysis <http://accelrys.com/products/collaborative-science/biovia-pipeline-pilot/> (accessed May 2, 2016).
- (58) Low, Y.; Uehara, T.; Minowa, Y.; Yamada, H.; Ohno, Y.; Urushidani, T.; Sedykh, A.; Muratov, E.; Kuz'min, V.; Fourches, D.; et al. Predicting Drug-Induced Hepatotoxicity Using QSAR and Toxicogenomics Approaches. *Chem. Res. Toxicol.* **2011**, *24*, 1251–1262.
- (59) Cohen, J. Weighted Kappa: Nominal Scale Agreement Provision for Scaled Disagreement or Partial Credit. *Psychol. Bull.* **1968**, *70* (4), 213–220.
- (60) Sing, T.; Sander, O.; Beerenwinkel, N.; Lengauer, T. ROCR: Visualizing Classifier Performance in R. *Bioinformatics* **2005**, *21* (20), 3940–3941.
- (61) Berrar, D.; Flach, P. Caveats and Pitfalls of ROC Analysis in Clinical Microarray Research (and How to Avoid Them). *Briefings Bioinf.* **2012**, *13* (1), 83–97.
- (62) Anger, L. T.; Wolf, A.; Schleifer, K.-J.; Schrenk, D.; Rohrer, S. G. Generalized Workflow for Generating Highly Predictive in Silico Off-Target Activity Models. *J. Chem. Inf. Model.* **2014**, *54* (9), 2411–2422.
- (63) RCSB Protein Data Bank - RCSB PDB <http://www.rcsb.org/pdb/home/home.do> (accessed Feb 4, 2016).
- (64) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucl. Acids Res.* **2000**, *28* (1), 235–242.
- (65) Stierand, K.; Rarey, M. Drawing the PDB: Protein–Ligand Complexes in Two Dimensions. *ACS Med. Chem. Lett.* **2010**, *1* (9), 540–545.
- (66) Bissantz, C.; Kuhn, B.; Stahl, M. A Medicinal Chemist's Guide to Molecular Interactions. *J. Med. Chem.* **2010**, *53* (14), 5061–5084.
- (67) Cannon, E. O.; Nigsch, F.; Mitchell, J. B. A Novel Hybrid Ultrafast Shape Descriptor Method for Use in Virtual Screening. *Chem. Cent. J.* **2008**, *2*, 3–3.
- (68) O'Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An Open Chemical Toolbox. *J. Cheminf.* **2011**, *3* (1), 33.
- (69) O'Boyle, N.; Morley, C.; Hutchison, G. Pybel: A Python Wrapper for the OpenBabel Cheminformatics Toolkit. *Chem. Cent. J.* **2008**, *2* (1), 5.
- (70) Marchese Robinson, Richard. Scripts for “Comparison of the predictive performance and interpretability of Random Forest and linear models on benchmark datasets” [revision 2.0] <https://doi.org/10.5281/zenodo.824245> (accessed Jul 7, 2017).
- (71) The R Project for Statistical Computing <http://www.r-project.org/> (accessed Oct 18, 2016).

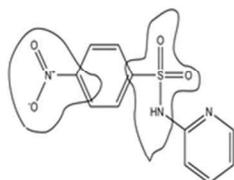
- 1
2
3 (72) Marchese Robinson, R. L. HeatMapWrapper (version 0.16)
4 <https://doi.org/10.5281/zenodo.495163> (accessed Apr 5, 2017).
5
6 (73) Benjamini, Y.; Yekutieli, D. The Control of the False Discovery Rate in Multiple Testing
7 under Dependency. *The Annals of Statistics* **2001**, *29* (4), 1165–1188.
8 (74) Bengio, Y.; Grandvalet, Y. No Unbiased Estimator of the Variance of K-Fold Cross-
9 Validation. *J. Mach. Learn. Res.* **2004**, *5*, 1089–1105.
10 (75) Nadeau, C.; Bengio, Y. Inference for the Generalization Error. *Mach. Learn.* **2003**, *52* (3),
11 239–281.
12 (76) Steiner, T. The Hydrogen Bond in the Solid State. *Angew. Chem., Int. Ed.* **2002**, *41* (1),
13 48–76.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

for Table of Contents use only

Comparison of the predictive performance and interpretability of Random Forest and linear models on benchmark datasets

Richard L. Marchese Robinson, Anna Palczewska, Jan Palczewski, Nathan Kidley

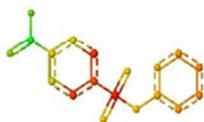
Non-toxic molecule with toxicophore (LHS) and detoxifying substructure (RHS) highlighted



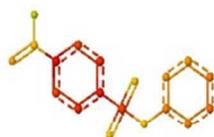
**How good are these QSAR predictions?
How meaningful are their interpretations?**

Random Forest

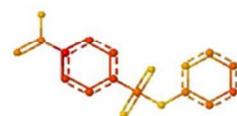
Interpretation approach 1



Interpretation approach 2(a)

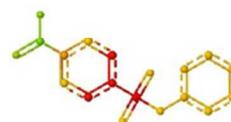


Interpretation approach 2(b)

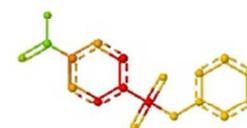


Linear models

Probit-PLS



SVM (linear kernel)



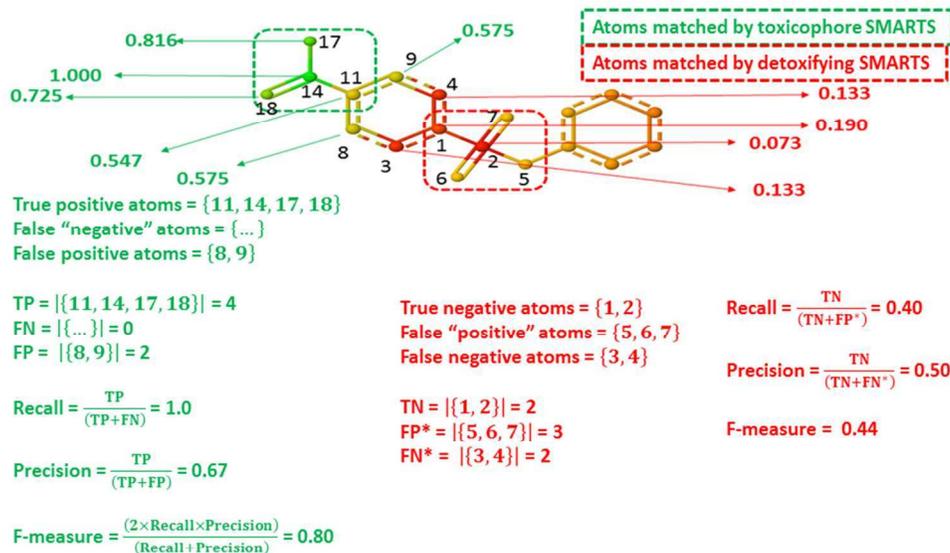
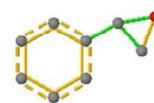
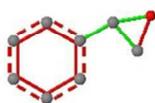
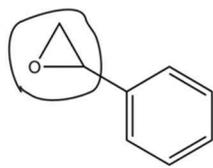


Figure 1 The procedure for calculating numeric Heat Map quality scores for an image obtained via symmetrized normalization (Symmetrized Single Molecule Normalization in this example) and atom coloring. The relevant normalized atom score estimates are shown matched to the atom IDs. In this case, the overall quality score would be the arithmetic mean of the two F-measures i.e. 0.62. Further details are provided in the associated text. SMARTS matching of toxicophores and detoxifying substructures was implemented using OpenBabel (version 2.3.2) & Pybel,^{68,69} based on SMARTS adapted from the work of Kazius et al.⁴⁸ and reported in the Supporting Information. The atom IDs, used to identify the corresponding atom score estimates for this image, were identified using OpenBabel⁶⁸, via this command: "obabel [SDF containing molecule of interest] -O [name of image].svg -xi -xu".

254x190mm (96 x 96 DPI)

Highlighted: epoxide moiety

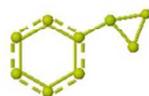
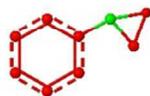


(A) Expert reasoning assigned toxicophore

(B) Bond color, Single Molecule Normalization

(C) Bond color, Full Dataset Normalization

(D) Bond color, Symmetrized Single Molecule Normalization



(E) Atom color, Symmetrized Single Molecule Normalization

(F) Atom color, Single Molecule Normalization

(G) Atom color, Full Dataset Normalization

(H) Atom color, Constant Symmetrized Normalization

Figure 2 Heat Map analysis of a binary classification leave-one-out SVM prediction made for molecule "CA23" in the CA dataset (c.f. Figure 4(b) in Mohr et al.), determined to be toxic (positive result) in the chromosome aberration test.¹⁷ (A) Expert reasoning assigned toxicophore (epoxide functionality), based on the rationale provided in Mohr et al.:¹⁷ this substructure was flagged by an Ames test alerts knowledgebase, with 80% of Ames positives being chromosome aberration positives, and is known to exhibit electrophilic reactivity with DNA. (B – H): Heat Map images based on descriptor contributions associated with the same SVM prediction, generated according to either bond or atom coloring obtained via different normalization schemes.

254x190mm (96 x 96 DPI)

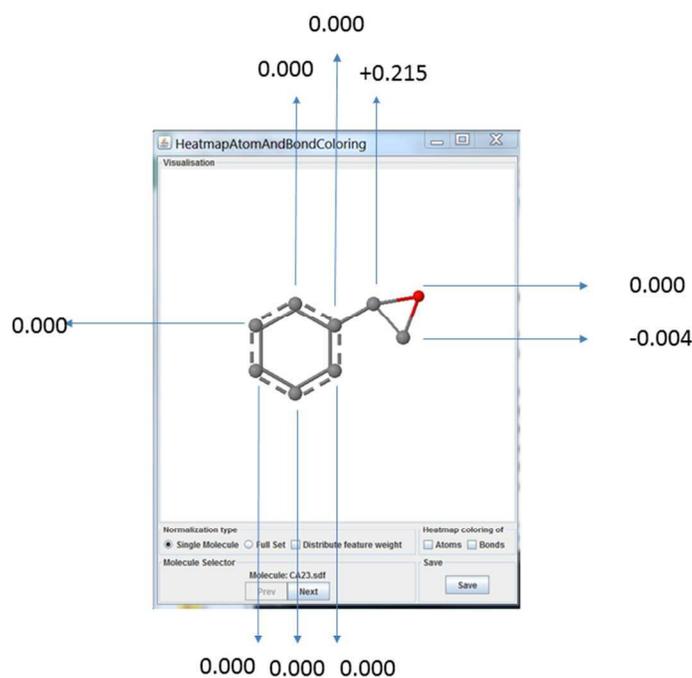


Figure 3 Image of the molecule "CA23" in the CA dataset¹⁷ displayed in HeatmapViewer.jar GUI,¹³ as generated by running the HeatMapWrapper tool,⁷² with both atom and bond coloring switched off. The intermediate output generated via running the HeatMapWrapper tool in atom color, Symmetrized Single Molecule Normalization mode, based on leave-one-out SVM descriptor contributions, was processed to reveal estimated raw atom scores ranging from weakly negative in one case (-0.004), through zero in most cases, to positive (+0.215) in one case. The estimated raw atom scores are rounded to 3dp. This image corresponds to all Heat Map images shown in Figure 2, with different normalization schemes generating the different normalized scores used for coloring. Via switching off atom coloring, based on normalized scores, the epoxide oxygen (in red) and carbon atoms (in black) can readily be discerned. N.B. The specific atom IDs, used to identify the corresponding atom score estimates, were identified as per Figure 1.

254x190mm (96 x 96 DPI)

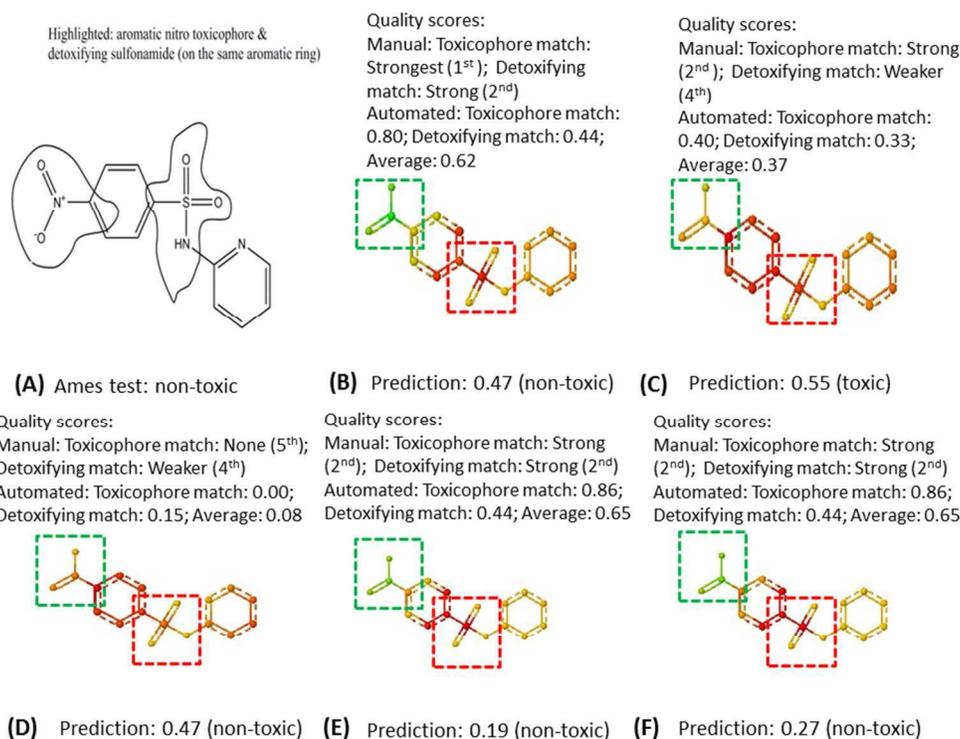


Figure 4 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule "1028-11-1" in the Kazius dataset (c.f. Figure 3(A) in Rosenbaum et al.).¹³ (A) Molecule with biologically significant substructures shown, based on a combination of mechanistic reasoning, expert knowledge and data analysis reported previously in the literature,^{13,48} and experimental Ames test assignment reported.¹³ (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. The predictions shown are the normalized scores (transformed into values between 0 and 1) for the toxic class, with predicted assignment to the toxic class if this score was greater than 0.50. The green dotted lines show the atoms matched to toxicophore SMARTS and the red dotted lines the atoms matched to corresponding detoxifying SMARTS. SMARTS matching and identification of matched atoms were carried out as per Figure 1. The manual and automated assessments of Heat Map quality are based on the schemes explained under "Quality assessment of Heat Map images".

254x190mm (96 x 96 DPI)

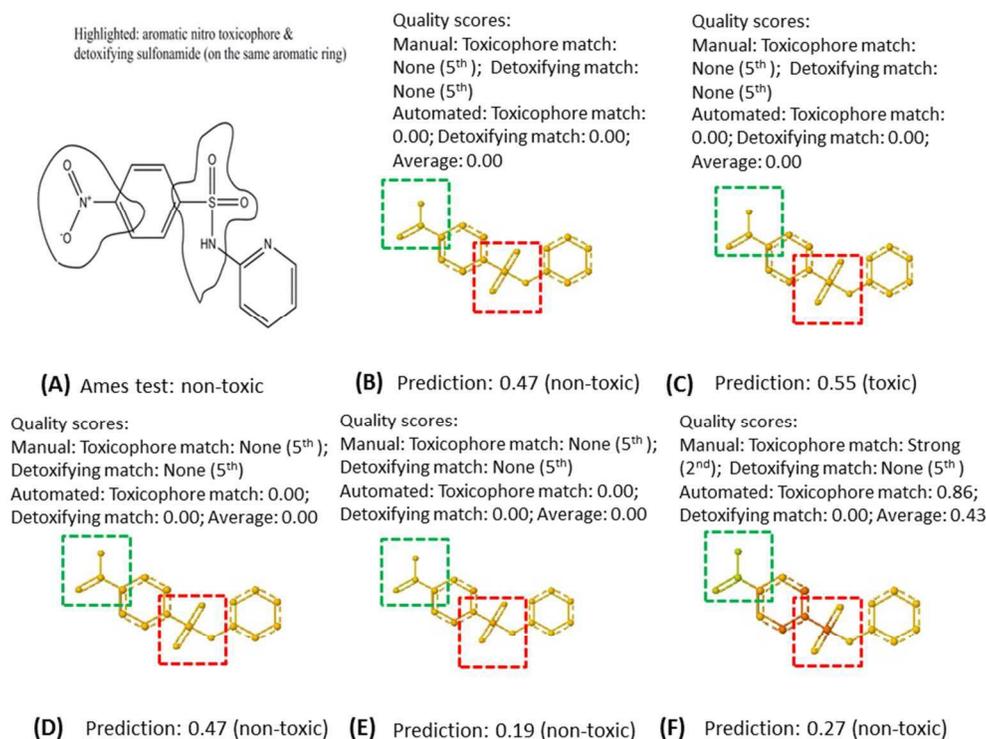
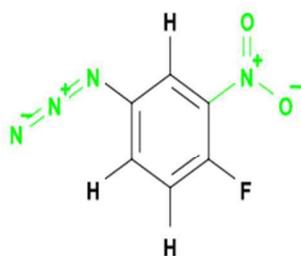


Figure 5 Heat Map (atom coloring & Constant Symmetrized Normalization) analysis of binary classification leave-one-out predictions made for molecule "1028-11-1" in the Kazius dataset (c.f. Figure 3(A) in Rosenbaum et al.).¹³ (A) Molecule with biologically significant substructures shown, based on a combination of mechanistic reasoning, expert knowledge and data analysis reported previously in the literature,^{13,48} and experimental Ames test assignment reported.¹³ (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per Figure 4. The green dotted lines show the atoms matched to toxicophore SMARTS and the red dotted lines the atoms matched to corresponding detoxifying SMARTS. SMARTS matching and identification of matched atoms were carried out as per Figure 1. The manual and automated assessments of Heat Map quality are based on the schemes explained under "Quality assessment of Heat Map images".

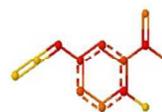
254x190mm (96 x 96 DPI)

28166-06-5: SMARTS matches



Quality scores:
Manual: Toxicophore match:
Strong (2nd)
Automated: Toxicophore match:
0.73

Quality scores:
Manual: Toxicophore match:
None (5th)
Automated: Toxicophore match:
0.00



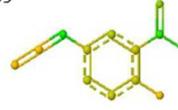
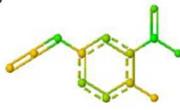
(A) Ames test: toxic

(B) Prediction: 0.48 (non-toxic) (C) Prediction: 0.54 (toxic)

Quality scores:
Manual: Toxicophore match:
None (5th)
Automated: Toxicophore match:
0.00

Quality scores:
Manual: Toxicophore match:
Weaker (4th)
Automated: Toxicophore match:
0.63

Quality scores:
Manual: Toxicophore match:
Weaker (4th)
Automated: Toxicophore match:
0.59



(D) Prediction: 0.48 (non-toxic) (E) Prediction: 0.82 (toxic) (F) Prediction: 0.87 (toxic)

Figure 6 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule "28166-06-5" in the Kazius dataset. (A) Molecule with SMARTS matches shown in green for specific aromatic nitro (RHS) and azide (LHS) toxicophores, derived from the work of Kazius et al.⁴⁸ based on a combination of mechanistic reasoning, expert knowledge and data analysis reported, and experimental Ames test assignment reported in the Kazius dataset SDF file.¹³ The SMARTS matches (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per Figure 4. The SMARTS matches show the atoms matched to SMARTS patterns. SMARTS matching and identification of matched atoms were carried out as per Figure 1. The manual and automated assessments of Heat Map quality are based on the schemes explained under "Quality assessment of Heat Map images". Quality scores, not to be confused with prediction scores, based on Heat Map correspondence to detoxifying groups were not assigned as no detoxifying groups were identified via SMARTS matching.

254x190mm (96 x 96 DPI)

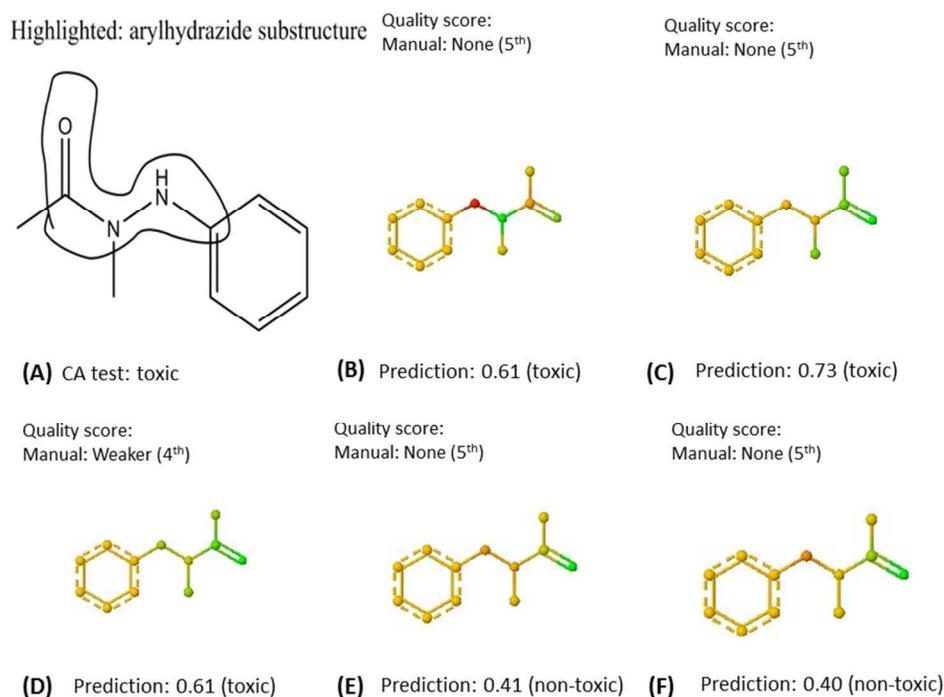


Figure 7 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule CA31 in the CA dataset (c.f. Figure 6(A) in Mohr et al.).¹⁷ (A) Molecule with biologically significant substructures shown, based on a combination of mechanistic reasoning and expert knowledge reported previously in the literature,¹⁷ and experimental chromosome aberration test assignment reported.¹⁷ (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per Figure 4. The manual assessments of Heat Map quality are based on the scheme explained under "Quality assessment of Heat Map images". Here, the assignments were not simply based on matching the entire putative toxicophore,¹⁷ but were informed by consideration of the likely toxic substructure revealed via metabolism, which is expected to be the arylhydrazine substructure and may even be the aromatic amine substructure arising from further metabolism.^{17,48}

254x190mm (96 x 96 DPI)

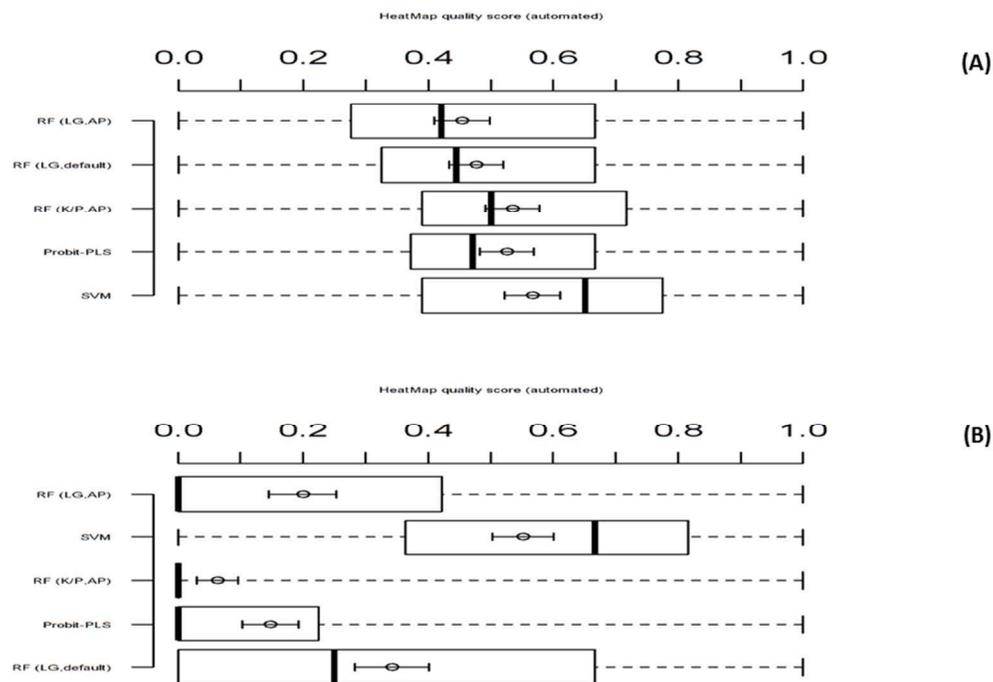


Figure 8 Automated Heat Map numeric quality scores, calculated as the average of F-measures for their correspondence to toxicophore and detoxifying substructure SMARTS matches as explained in Figure 1, for Heat Map images generated using atom coloring and (A) Symmetrized Single Molecule Normalization or (B) Constant Symmetrized Normalization. Quality score distributions correspond to all 39 leave-one-out predictions made for a subset of Kazius dataset molecules identified via toxicophore SMARTS matches. The different quality score distributions are annotated according to the different approaches used to obtain the underlying descriptor contributions: RF (LG, AP) = Random Forest classification (local gradients, averaged predictions); RF (LG, default) = Random Forest classification (local gradients, majority votes predictions); RF (K/P, AP) = Random Forest classification (Kuz'min/Palczewska, averaged predictions); Probit-PLS (linear coefficients); SVM (linear coefficients). The distributions are represented as boxplots, with the bold lines showing the medians and the superimposed circle the arithmetic means, with error bars denoting the standard errors in the means.

254x190mm (96 x 96 DPI)

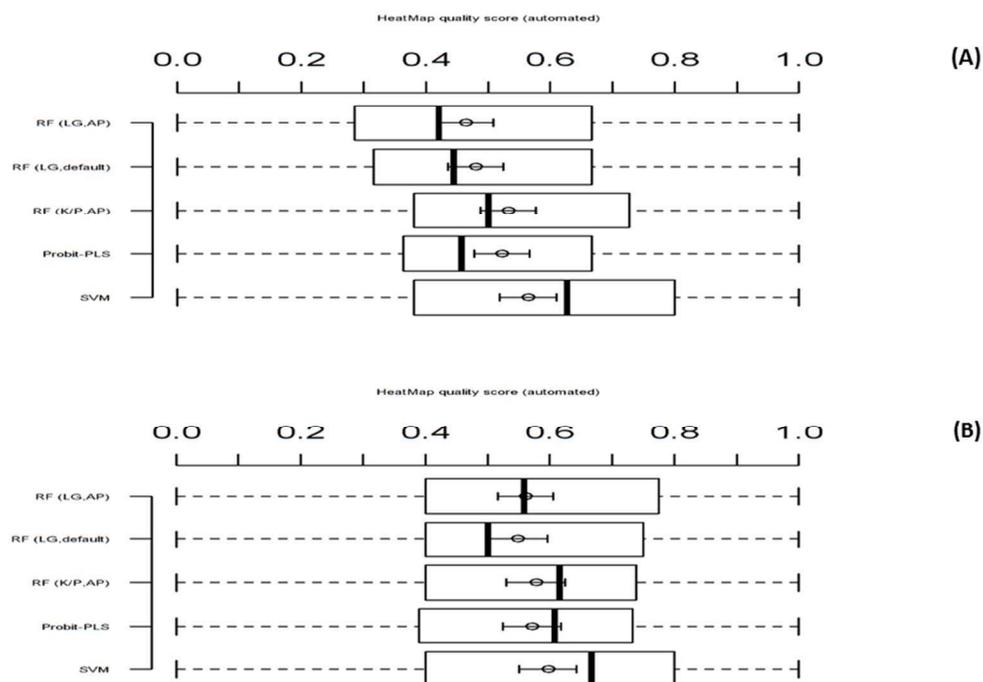


Figure 9 Automated Heat Map numeric quality scores, based on their correspondence to toxicophore substructure SMARTS matches as explained in Figure 1, for Heat Map images generated using atom coloring and Symmetrized Single Molecule Normalization. In sub-figure (A), quality score distributions correspond to 38 leave-one-out predictions made for a subset of Kazius dataset molecules identified via toxicophore SMARTS matches, excluding the single molecule ("1028-11-1") for which a detoxifying substructure was also identified via SMARTS pattern matches. In sub-figure (B), the distributions correspond to the subsets of those 38 molecules which were correctly predicted. The number of molecules correctly predicted to be toxic or non-toxic varied across methods: Random Forest classification (averaged predictions) = 28; Random Forest classification (majority vote predictions) = 29; Probit-PLS = 28; SVM = 31. The different quality scores distributions are annotated according to the different approaches used to obtain the underlying descriptor contributions: RF (LG, AP) = Random Forest classification (local gradients, averaged predictions); RF (LG, default) = Random Forest classification (local gradients, majority votes predictions); RF (K/P, AP) = Random Forest classification (Kuz'min/Palczewska, averaged predictions); Probit-PLS (linear coefficients); SVM (linear coefficients). The distributions are represented as boxplots, with the bold lines showing the medians and the superimposed circle the arithmetic means, with error bars denoting the standard errors in the means.

254x190mm (96 x 96 DPI)

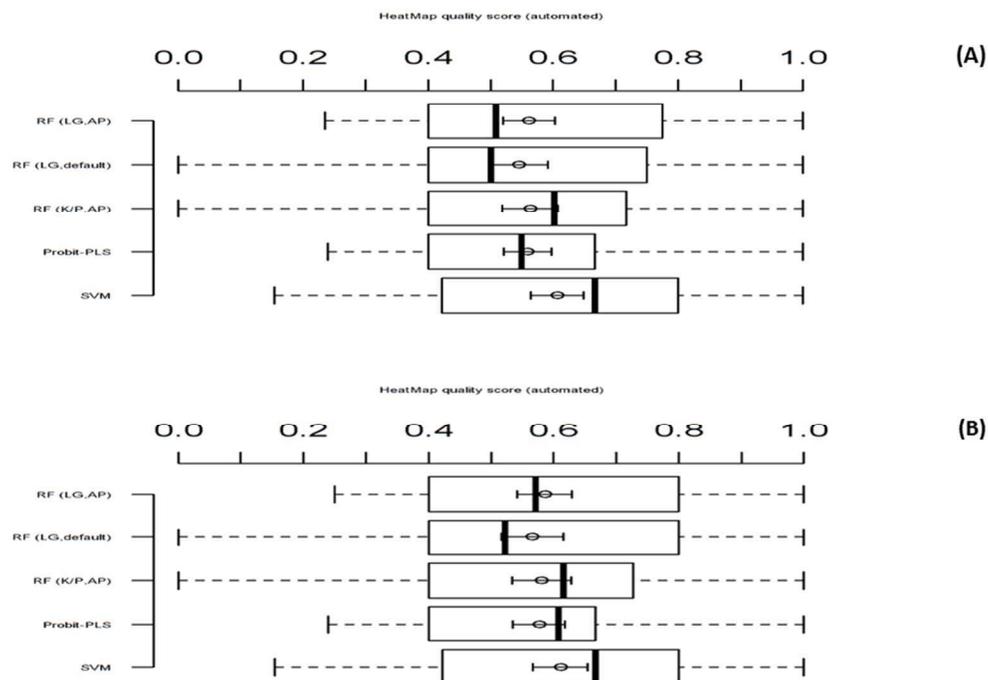


Figure 10 Automated Heat Map numeric quality scores, calculated based on their correspondence to toxicophore SMARTS matches as explained in Figure 1, for Heat Map images generated using atom coloring and Symmetrized Single Molecule Normalization. In sub-figure (A), the distributions of quality scores correspond to all toxic leave-one-out predictions made for a subset of Kazius dataset molecules identified via toxicophore SMARTS matches, excluding the single molecule ("1028-11-1") for which a detoxifying substructure was also identified via SMARTS pattern matches. In sub-figure (B), the distributions of quality scores correspond to the subsets of those toxic predictions which were correct. The number of molecules (correctly) predicted to be toxic, hence the number of Heat Map images for which these distributions were generated, varied across methods: Random Forest classification (averaged predictions) = 28 (25 correct); Random Forest classification (majority vote predictions) = 29 (26 correct); Probit-PLS = 30 (26 correct); SVM = 31 (28 correct). The different quality scores distributions are annotated according to the different approaches used to obtain the underlying descriptor contributions: RF (LG, AP) = Random Forest classification (local gradients, averaged predictions); RF (LG, default) = Random Forest classification (local gradients, majority votes predictions); RF (K/P, AP) = Random Forest classification (Kuz'min/Palczewska, averaged predictions); Probit-PLS (linear coefficients); SVM (linear coefficients). The distributions are represented as boxplots, with the bold lines showing the medians and the superimposed circle the arithmetic means, with error bars denoting the standard errors in the means.

254x190mm (96 x 96 DPI)

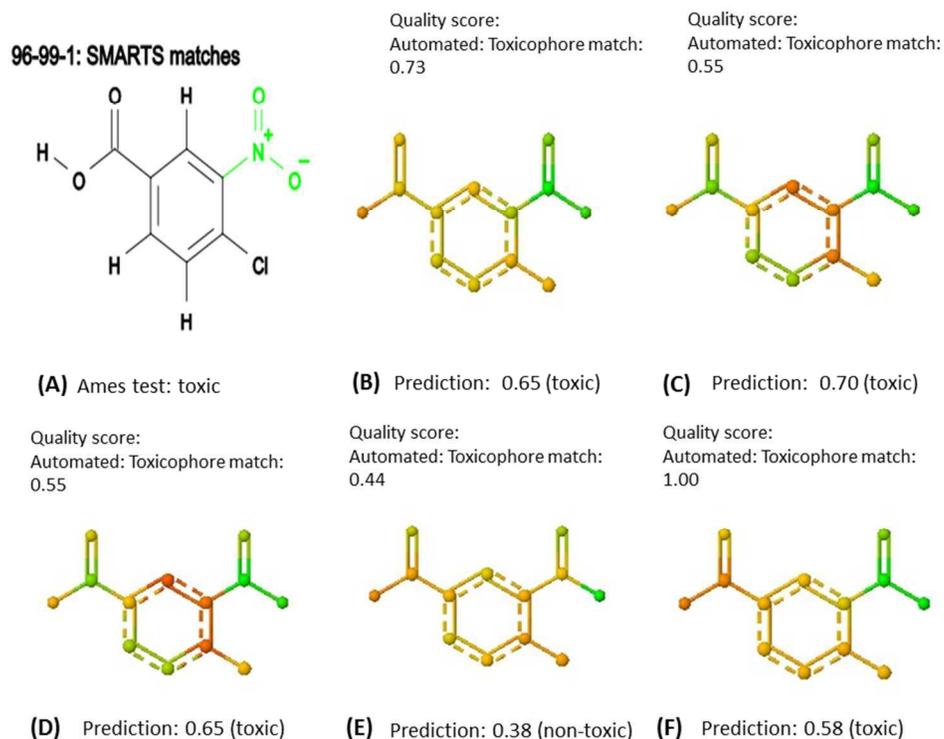


Figure 11 Heat Map (atom coloring & Symmetrized Single Molecule Normalization) analysis of binary classification leave-one-out predictions made for molecule "96-99-1" in the Kazius dataset. (A) Molecule with SMARTS matches shown in green for specific aromatic nitro toxicophore, derived from the work of Kazius et al.⁴⁸ based on a combination of mechanistic reasoning, expert knowledge and data analysis, and experimental Ames test assignment reported in the Kazius dataset SDF file.¹³ The SMARTS matches (B – F) Heat Map images: (B) Random Forest classification (Kuz'min/Palczewska, averaged predictions); (C) Random Forest classification (local gradients, majority vote predictions); (D) Random Forest (local gradients, averaged predictions); (E) Probit-PLS; (F) SVM. Predictions are reported as per Figure 4. The SMARTS matches show the atoms matched to SMARTS patterns. SMARTS matching and identification of matched atoms were carried out as per Figure 1. The automated assessments of Heat Map quality are based on the scheme explained in Figure 1. Quality scores, not to be confused with prediction scores, based on Heat Map correspondence to detoxifying groups were not assigned as no detoxifying groups were identified via SMARTS matching.

254x190mm (96 x 96 DPI)

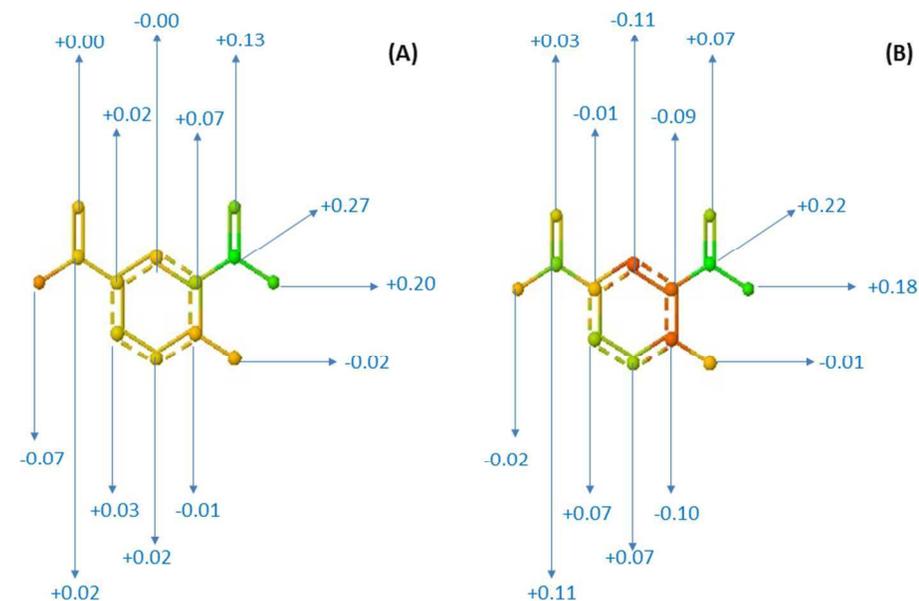


Figure 12 Comparison between the estimated raw atom scores (2dp), for the molecule "96-99-1" in the Kazius dataset, derived from the descriptor contributions corresponding to a single leave-one-out prediction made using Random Forest classification (averaged predictions) according to the Kuz'min/Palczewska (A) and local gradients (B) approaches. The raw scores were estimated via processing output from the HeatMapWrapper tool⁷² in atom coloring and Symmetrized Single Molecule Normalization mode. The correlation between the corresponding raw atom score estimates, in terms of the Pearson's correlation coefficient, is 0.74. The raw atom score estimates are used to annotate the corresponding atoms in the Heat Map (atom coloring & Symmetrized Single Molecule Normalization) image. N.B. The specific atom IDs, used to identify the corresponding atom score estimates, were identified as per Figure 1.

254x190mm (96 x 96 DPI)

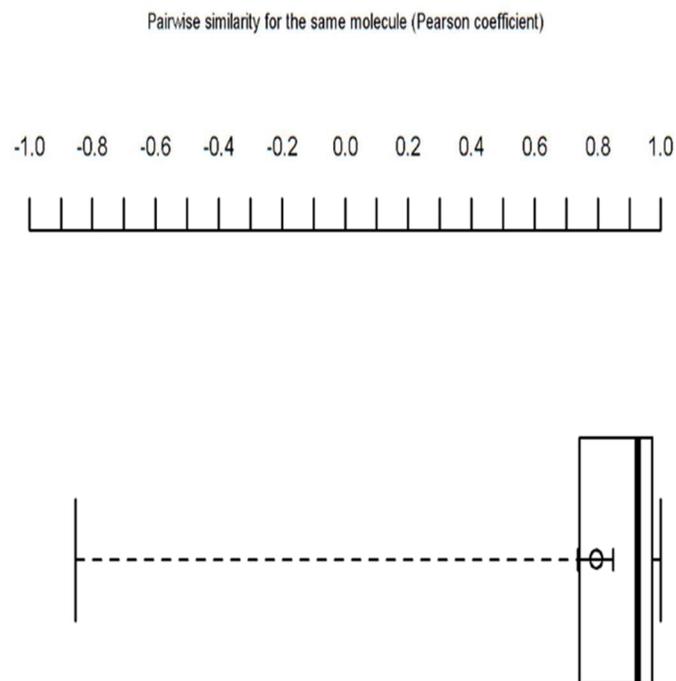


Figure 13 The distribution of pairwise similarities, calculated via the Pearson's correlation coefficient using estimates of the raw atom scores derived from descriptor contributions generated via the Kuz'min/Palczewska and local gradients approaches, for all leave-one-out Random Forest classification predictions (averaged predictions approach) made for a subset of molecules in the Kazius dataset identified via toxicophore SMARTS pattern matches. (The raw atom score estimates were obtained via processing the output of the HeatMapWrapper program⁷² in atom coloring and Symmetrized Single Molecule Normalization mode.) The distribution is represented as a boxplot, with the bold line showing the median and the superimposed circle the arithmetic mean, with error bars denoting the standard error in the mean.

254x190mm (96 x 96 DPI)