

# SYBA: Bayesian estimation of synthetic accessibility of organic compounds

**Milan Voršilák**

Ustav molekularni genetiky Akademie Ved Ceske Republiky <https://orcid.org/0000-0002-8923-1627>

**Michal Kolář**

Ustav molekularni genetiky Akademie Ved Ceske Republiky <https://orcid.org/0000-0002-4593-1525>

**Ivan Čmelo**

University of Chemistry and Technology Prague <https://orcid.org/0000-0001-7787-8653>

**Daniel Svozil** (✉ [daniel.svozil@gmail.com](mailto:daniel.svozil@gmail.com))

University of Chemistry and Technology <https://orcid.org/0000-0003-2577-5163>

---

## Research article

**Keywords:** synthetic accessibility, Bayesian analysis, Bernoulli naïve Bayes

**Posted Date:** May 20th, 2020

**DOI:** <https://doi.org/10.21203/rs.2.22597/v3>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Journal of Cheminformatics on May 20th, 2020. See the published version at <https://doi.org/10.1186/s13321-020-00439-2>.

# SYBA: Bayesian estimation of synthetic accessibility of organic compounds

Milan Voršilák<sup>1,2</sup>, Michal Kolář<sup>3,4</sup>, Ivan Čmelo<sup>1</sup>, Daniel Svozil<sup>1,2,\*</sup>

<sup>1</sup> CZ-OPENSREEN: National Infrastructure for Chemical Biology, Department of Informatics and Chemistry, Faculty of Chemical Technology, University of Chemistry and Technology Prague, Technická 5, 166 28 Prague 6, Czech Republic.

<sup>2</sup> CZ-OPENSREEN: National Infrastructure for Chemical Biology, Institute of Molecular Genetics of the Czech Academy of Sciences, Vídeňská 1083, 142 20 Prague 4, Czech Republic.

<sup>3</sup> Laboratory of Genomics and Bioinformatics, Institute of Molecular Genetics of the Czech Academy of Sciences, Vídeňská 1083, 142 20 Prague 4, Czech Republic.

<sup>4</sup> Department of Informatics and Chemistry, Faculty of Chemical Technology, University of Chemistry and Technology Prague, Technická 5, 166 28 Prague 6, Czech Republic.

\* corresponding author, e-mail: Daniel.Svozil@vscht.cz

Email addresses:

MV: milan.vorsilak@img.cas.cz, ORCID: 0000-0002-8923-1627

MK: kolarmi@img.cas.cz, ORCID: 0000-0002-4593-1525

IČ: ivan.cmelo@vscht.cz, ORCID: 0000-0001-7787-8653

DS: daniel.svozil@vscht.cz, ORCID: 0000-0003-2577-5163

## Abstract

SYBA (SYnthetic Bayesian Accessibility) is a fragment-based method for the rapid classification of organic compounds as easy- (ES) or hard-to-synthesize (HS). It is based on a Bernoulli naïve Bayes classifier that is used to assign SYBA score contributions to individual fragments based on their frequencies in the database of ES and HS molecules. SYBA was trained on ES molecules available in the ZINC15 database and on HS molecules generated by the Nonpher methodology. SYBA was

compared with a random forest, that was utilized as a baseline method, as well as with other two methods for synthetic accessibility assessment: SAScore and SCScore. When used with their suggested thresholds, SYBA improves over random forest classification, albeit marginally, and outperforms SAScore and SCScore. However, upon the optimization of SAScore threshold (that changes from 6.0 to ~4.5), SAScore yields similar results as SYBA. Because SYBA is based merely on fragment contributions, it can be used for the analysis of the contribution of individual molecular parts to compound synthetic accessibility. SYBA is publicly available at <https://github.com/lich-uct/syba> under the GNU General Public License.

## Keywords

synthetic accessibility – Bayesian analysis – Bernoulli naïve Bayes

## Background

Chemical space available for the generation of new molecules is huge [1-4], making the synthesis and testing of all possible compounds impractical. Therefore chemists, both experimental and computational, developed tools and approaches for the exploration of chemical space with the aim to identify new compounds with desirable physico-chemical, biological and pharmacological properties [5-12]. A major *in silico* method for chemical space exploration is *de novo* molecular design in which new virtual molecules are assembled from scratch [13-18]. An essential requirement for *de novo* designed compounds is their synthetic accessibility. Synthetic accessibility is commonly incorporated into *de novo* design programs by employing chemical strategies that guide an assembly process. For example, the connections between certain atom types can be disallowed [19], established chemical reactions can be used to connect individual molecular building blocks [20, 21] or the retrosynthetic rules can be directly incorporated into the assembly process [22, 23].

The latest development in *de novo* molecular design are molecular generators based on deep learning approaches [24-26]. These typically construct new molecules not by assembling the building blocks, but by producing chemically feasible SMILES strings [27-32]. The generators are able to produce millions of virtual compounds, synthetic accessibility of which has to be quickly and efficiently assessed. Quick synthetic accessibility assessment can be based [33] on molecule's complexity that is

typically calculated [34-37] from the number of atoms, bonds, rings, and/or hard-to-synthesize motifs, such as chiral centers or uncommon ring fusions. However, the definition of molecular complexity is ambiguous and context dependent [38, 39]. The structural complexity is not equivalent to the synthetic one as complexity-based metrics do not incorporate any information about starting materials and tend to remove molecules that can be synthesized from already existing complex precursors [40, 41]. A better way of synthetic accessibility assessment is to use the complexity of the synthetic route [42]. Based on this principle, SCScore, a data-driven metric designed to describe real syntheses, was developed recently [43]. SCScore is based on the idea that reaction products are synthetically more complex than reactants. To quantify this, a deep feed-forward neural network, that assigns a synthetic complexity score between 1 and 5, was trained on 22 million reactant-product pairs from the Reaxys database [44]. Using the hinge loss objective function, that supports the separation between scores in each reactant-product pair, the model learns synthetic complexity score that correlates with the number of reaction steps, but does not rely on the availability of reaction database or organic chemist ranking.

SAScore [45], another popular and rapid method for synthetic accessibility assessment, is based on the analysis of ECFP4 [46] fragments obtained from one million compounds randomly selected from the PubChem database [47]. The main idea of SAScore is that when a molecular fragment occurs often in the PubChem database, it contributes to the synthetic accessibility of a molecule more than a less frequently occurring fragment. Each fragment is assigned a numerical score, frequent fragments have positive scores and less frequent fragments have negative scores. In addition to the fragment score, SAScore consists of a complexity penalty and symmetry bonus. These terms penalize nonstandard structural motives such as macrocycles, stereo centers, spiro and bridge atom, but reward the symmetry of a structure. SAScore acquires values between 1 (easy to make) and 10 (very difficult to make), where 6.0 is suggested by the authors [45] as a threshold to distinguish between easy- and hard-to-synthesize compounds. SAScore is a popular high-throughput measure and proved to be a very useful tool in many cheminformatics applications [27, 48-50].

In the present work, we further expand on main concepts of SAScore construction. We developed SYBA (SYnthetic Bayesian Accessibility), a rapid fragment-based score derived using Bayesian probabilistic modeling. Fragment contributions to SYBA are calculated not only from fragments present in synthetically accessible molecules, but also from fragments appearing in hard-to-synthesize molecules.

## Methods

### SYBA score derivation

SYBA is a Bernoulli naïve Bayes classifier based on the frequency of molecular fragments that are present in the database of easy-to-synthesize (ES) and hard-to-synthesize (HS) molecules and on the assumption of the independence of molecular fragments. Though such assumption is bold, it was shown to provide surprisingly good results in many cheminformatics studies [51-55].

Each compound is represented by a binary fingerprint  $\mathbf{F} = [f_1, f_2, \dots, f_M]$  of length  $M$  where  $f_i$  indicates the presence ( $f_i = 1$ ) or absence ( $f_i = 0$ ) of the specific fragment  $i$  in the compound. SYBA uses this fingerprint to assign the molecule to a class  $C \in \langle \text{ES}, \text{HS} \rangle$ . The calculation is based on the Bayes theorem

$$\text{Equation 1} \quad p(C|\mathbf{F}) = \frac{p(\mathbf{F}|C) p(C)}{p(\mathbf{F})},$$

where  $p(C|\mathbf{F})$  is the posterior probability that a compound with a certain set of molecular fragments  $\mathbf{F}$  belongs to the class  $C$ . The likelihood  $p(\mathbf{F}|C)$  is the conditional probability that a compound from the class  $C$  contains a set of molecular fragments  $\mathbf{F}$ . The marginal probabilities  $p(\mathbf{F})$  and  $p(C)$  express our belief to observe a set of molecular fragments  $\mathbf{F}$  and the molecule that belongs to the class  $C$ .

The SYBA score is defined as the logarithm of the ratio of the posterior probabilities that the molecule belongs to the ES and HS classes,

$$\text{Equation 2} \quad \text{SYBA}(\mathbf{F}) = \ln \left( \frac{p(\text{ES}|\mathbf{F})}{p(\text{HS}|\mathbf{F})} \right).$$

Using Equation 1, the SYBA score can be expressed as

Equation 3 
$$\text{SYBA}(\mathbf{F}) = \ln\left(\frac{p(\text{ES})}{p(\text{HS})}\right) + \ln\left(\frac{p(\mathbf{F}|\text{ES})}{p(\mathbf{F}|\text{HS})}\right).$$

In the data set SYBA was derived from (further referred to as the training data set S), ES and HS compounds are represented evenly, the priors  $p(\text{ES})$  and  $p(\text{HS})$  are thus equal and the term  $\ln\left(\frac{p(\text{ES})}{p(\text{HS})}\right)$  becomes zero:

Equation 4 
$$\text{SYBA}(\mathbf{F}) = \ln\left(\frac{p(\mathbf{F}|\text{ES})}{p(\mathbf{F}|\text{HS})}\right).$$

Assuming the independence of molecular fragments, the conditional probability  $p(\mathbf{F}|C)$  factorizes to  $p(\mathbf{F}|C) = \prod_{i=1}^M p(f_i|C)$  and the SYBA score simplifies to

Equation 5 
$$\text{SYBA}(\mathbf{F}) = \sum_{i=1}^M s_i(f_i)$$

where  $s_i(f_i)$  is the score contribution from the fragment  $i$  (SYBA fragment score) given as

Equation 6 
$$s_i(f_i) = \ln\left(\frac{p(f_i|\text{ES})}{p(f_i|\text{HS})}\right).$$

Considering that  $p(f_i|\text{ES}) = 1 - p(f_i|\text{HS})$ , the fragment scores  $s_i(f_i)$  in Equation 6 represent logits and can be expressed using the fragment frequencies in the training data set S as

Equation 7 
$$s_i(f_i) = \ln\frac{N_{\text{HS}}+2}{N_{\text{ES}}+2} + f_i \ln\frac{(n_{\text{ES},i}+1)}{(n_{\text{HS},i}+1)} + (1 - f_i) \ln\frac{(N_{\text{ES}}-n_{\text{ES},i}+1)}{(N_{\text{HS}}-n_{\text{HS},i}+1)},$$

where  $N_{\text{HS}}$  is the number of HS and  $N_{\text{ES}}$  the number of ES molecules in the training data set S,  $n_{\text{HS},i}$  is the number of HS molecules in the training data set S that contain the fragment  $i$ , and  $n_{\text{ES},i}$  is the number of ES molecules in the training data set S that contain the fragment  $i$ . See Additional file 2 for a detailed derivation. Positive  $s_i(f_i)$  means that the presence/absence of the fragment  $i$  is more probable in ES than in HS class and vice versa. Positive SYBA means that the compound belongs more likely to the ES class, while negative SYBA means that the compound belongs more likely to the HS class. The higher the absolute value of SYBA, the more evidence for the class membership is present in the molecule.

## Training set construction

The training data set  $S$  consists of two subsets:  $S_+$  contains ES structures and  $S_-$  contains HS structures (Figure 1, Additional file 1). While ES molecules can be readily obtained, for example, from the ZINC database of purchasable compounds [56, 57], no equivalent database of HS molecules exists. However, HS molecules can be designed by Nonpher [58], a method based on a molecular morphing approach [59]. In Nonpher, a starting molecule is gradually transformed into a more complex compound using small structural perturbations, such as the addition or removal of an atom or a bond. To prevent the creation of overly complex structures, four complexity indices (Bertz [34], Whitlock [35], BC [36] and SMCM [37]) are monitored and once their respective thresholds (Additional file 2 – Table S1) are exceeded, Nonpher is stopped.

Using Nonpher, 693 353 HS molecules were generated and they form the  $S_-$  data set. The  $S_+$  data set, containing ES compounds, is formed by the same number of molecules randomly chosen (excluding natural products) from the ZINC15 database [57] so that their distribution of the number of heavy atoms is the same as in the  $S_-$  data sets. Every  $S_+$  and  $S_-$  molecule was fragmented using the Morgan fingerprint function in the RDKit toolkit [60]. Fragments with the radius of 4 and smaller, corresponding to radial ECFP8 [46] fragments, were used. This type of fragments consists of a central atom and atoms distant from the central atom up to four bonds. Besides ECFP8 fragments, the number of stereocenters was also included into SYBA as the molecules with more stereocenters are typically more difficult to synthesize. The stereo score is based, similarly to the fragment score, on the analysis of the number of stereocenters in the training set  $S$ . To obtain the stereo score, molecules were divided into 6 bins differing by the number of stereocenters (0, 1, 2, 3, 4 and 5+) and individual score contributions were calculated from Equation 8.

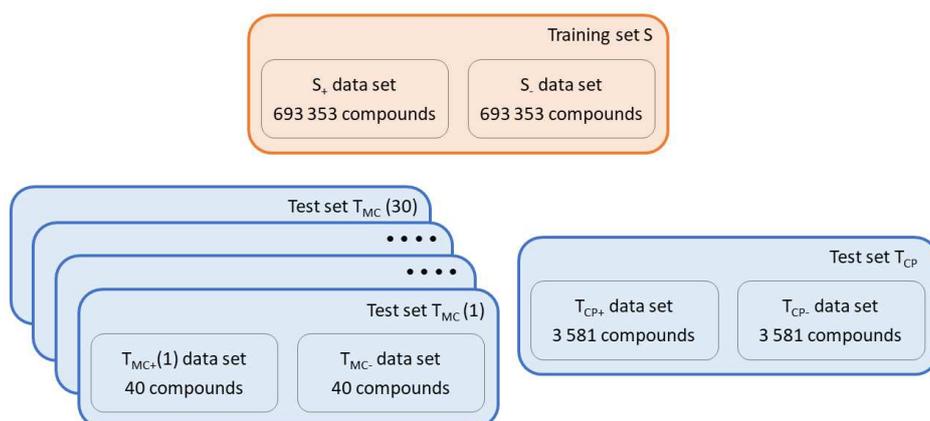
## Test set construction

SYBA performance could have been assessed using a test set created in a similar way as the training set  $S$ , i.e. using HS compounds generated by Nonpher. However, such test set would be clearly biased towards chemical space covered by Nonpher. Therefore, two test sets were constructed in a conceptually different manner. First test set, further denoted as  $T_{MC}$ , was manually curated from the

literature, second test set, referred to as  $T_{CP}$ , was computationally picked from the ZINC15 [57] and GDB17 databases [61].

HS compounds in  $T_{MC}$  (denoted as  $T_{MC-}$ ) were obtained by the analysis [58] of 296 published compounds assessed by experienced medicinal chemists [41, 45, 62, 63]. Based on original chemists' scores, the final  $T_{MC-}$  data set of 40 HS compounds was assembled. A complementary  $T_{MC+}$  data set consists of 40 ES compounds selected from the ZINC15 database [57] in such a way that the distribution of the number of their heavy atoms is the same as in the  $T_{MC-}$  data set. Because small  $T_{MC}$  size may bias the results, 30 different  $T_{MC}$  data set instances were generated using the same 40  $T_{MC-}$  compounds, but different 40  $T_{MC+}$  compounds (Additional file 3).

HS compounds in the  $T_{CP}$  test set (Additional file 4), denoted as  $T_{CP-}$ , were obtained by the analysis of the publicly available subset of 50M molecules from the GDB-17 database [61]. Only molecules exceeding thresholds (Additional file 2 – Table S1) of all monitored complexity indices (Bertz [34], Whitlock [35], BC [36] and SMCM [37]) were considered to be HS. In total, 3 581 molecules form the  $T_{CP-}$  data set. A complementary  $T_{CP+}$  data set consists of the same number of compounds randomly selected from the ZINC15 database [57] that follow the same size distribution as HS compounds and that, in addition, do not exceed any of the aforescribed complexity indices. Data sets used in the present work are summarized in Figure 1.



**Figure 1** Data set summary. Training set was used to derive SYBA scores, as well as to train a random forest classifier. Training set consists of 693 353 molecules randomly selected from the ZINC15 database [57] that are considered to be ES ( $S_+$  data set) and of the same number of HS molecules generated by Nonpher [58] ( $S_-$  data set). Two test sets were used to compare the performance of SYBA, a random forest, SAScore [45] and SCScore [43]. Manually curated test set ( $T_{MC}$ ) contains 40 compounds ( $T_{MC-}$  data set) considered to be HS by experienced medicinal chemists [58] supplemented by 40 ES

compounds randomly selected from the ZINC15 database ( $T_{MC+}$  data set). 30  $T_{MC}$  data set instances differing in  $T_{MC+}$  compounds were constructed. Computationally picked test set ( $T_{CP}$ ) consists of 3 581 HS compounds that were obtained from the GDB-17 database [61] ( $T_{CP-}$  data set) complemented by the same number of compounds randomly selected from the ZINC15 database ( $T_{CP+}$  data set).

## Performance evaluation

The performance of classification models studied in the present work was assessed by four different metrics: the classification accuracy ( $Acc$ ), sensitivity ( $SN$ ), specificity ( $SP$ ) and area under the ROC curve ( $AUC$ ).  $Acc$  gives the percentage of correctly classified samples regardless of their class.

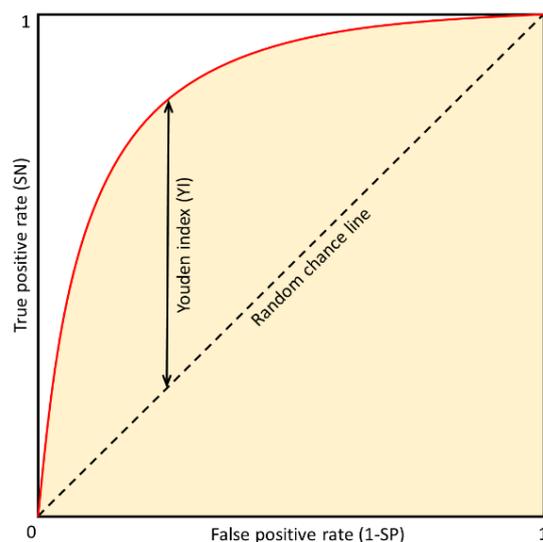
$$\text{Equation 9} \quad \text{Accuracy } (Acc) = \frac{TP+TN}{TP+TN+FN+FP}$$

where true positives ( $TP$ ) are ES compounds predicted by a model to be ES, true negatives ( $TN$ ) are HS compounds predicted to be HS, false positives ( $FP$ ) are HS compounds predicted to be ES and false negatives ( $FN$ ) are ES compounds predicted to be HS. The accuracy can also be evaluated for positive and negative classes independently leading to  $SN$  and  $SP$ .  $SN$  is the percentage of correctly predicted positive class compounds, while the percentage of correctly predicted negative class compounds is known as  $SP$ .

$$\text{Equation 10} \quad \text{Sensitivity } (SN) = \frac{TP}{TP+FN}$$

$$\text{Equation 11} \quad \text{Specificity } (SP) = \frac{TN}{TN+FP}$$

$SN$  and  $SP$  can be combined in the receiver operating characteristic (ROC) curve that is the graphical representation of the trade-off between true positive rate (given as  $SN$ ) and false positive rate (given as  $1 - SP$ ) over all possible thresholds (Figure 2). The area under the ROC curve ( $AUC$ ) is the quantitative measure of the performance of a classifier and is equal to the probability that a classifier will rank a randomly chosen positive example higher than a randomly chosen negative example. A random classifier has  $AUC$  of 0.5, while  $AUC$  for a perfect classifier is equal to 1.



**Figure 2** ROC curve and Youden index. The ROC curve (red line) is the dependency of true positive rate (it equals to  $SN$ ) on false positive rate (it equals to  $1 - SP$ ) at various thresholds. The random chance line represents a classifier that assigns examples into individual classes randomly. Orange shaded area represents the area under the ROC curve ( $AUC$ ). The larger the  $AUC$ , the better is the overall performance of the classifier. Youden index ( $YI$ ) is the point on the ROC curve that is farthest from the random chance line along the  $SN$  axis.

## Random forest classification, SAScore and SCScore

Because of its wide adoption in various cheminformatics applications [58, 64-66], the random forest (RF) classifier with compounds encoded by 1024-bits long Morgan fingerprint with radius 2 was used as a baseline method with which SYBA, SAScore [45] and SCScore [43] were compared. The RF classifier was implemented in Scikit-learn [67]. Two RF hyperparameters were optimized: the number of trees (50, 100, 300 and 500) and the maximum number of features considered when looking for the best split (10% out of 1024 = 102, 25% = 256, 50% = 512, 75% = 768, 100% = 1024,  $\sqrt{1024} = 32$  and  $\log_2(1024) = 10$ ). For each pair of hyperparameters, RF model was trained using the training set S and the prediction accuracy was evaluated on the test set  $T_{CP}$  (Additional file 2 – Table S2, Figure S2-S8). The setting used in this work (100 trees and 32 features) represents the best trade-off between computational efficiency and prediction accuracy [64]. RF was trained using the training set S (Figure 1). SAScore was calculated by the RDKit toolkit [60]. SCScore code was downloaded from the public GitHub repository [68].

## Classification thresholds

In SYBA, more positive value means a higher probability that the compound is ES and more negative value indicates a higher probability that the compound is HS (Equation 4). The threshold value of zero is used to distinguish between ES and HS compounds. For SAScore, the recommended value of 6.0 [45] was used as a threshold. In RF, the final prediction is based on a number of decision trees that predict either of classes. Here, 0.5 is used as a threshold, i.e., if more decision trees predict ES than HS class, a compound is classified as ES and vice versa. For SCScore [43], no threshold was suggested by the authors. In such case, the threshold can be identified by the analysis of the ROC curve. A frequently used measure that enables the selection of an optimal threshold is the Youden index (*YI*) [69, 70]. *YI* is defined as

$$\text{Equation 12} \quad YI = \max (SN + SP - 1)$$

and ranges between 0 and 1 (Figure 2). The optimal threshold value is selected by maximizing *YI*, i.e., by maximizing the sum of *SN* and *SP*.

## Statistical comparison of model performance

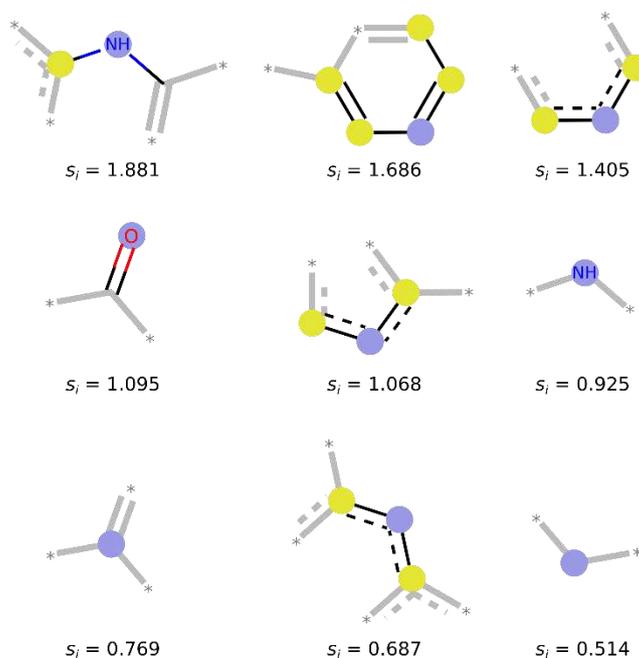
The performance of studied classification models was compared using non-parametric Cochran's Q test [71], an omnibus test for testing for differences between three or more machine learning models. In the case of the statistically significant result of Cochran's Q test, differing pairs of classification models were identified by McNemar's post-hoc paired test [72] with Benjamini-Hochberg false discovery rate adjustment [73]. McNemar's test checks if the distribution of disagreements between two methods is imbalanced. The statistical significance for all tests in the present work was assessed at the significance level  $\alpha = 0.05$ .

## Results and discussion

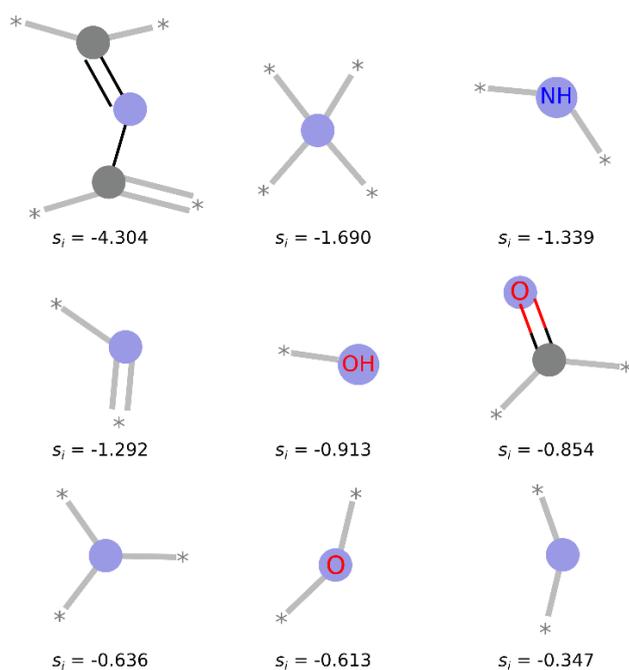
### Chemical space covered by SYBA data sets

The examples of training set compounds are given in Additional file 2 – Figure S9-S12. In total, 3 439 074 ECFP8 fragments were obtained for ES compounds and 23 447 524 fragments for HS

compounds. 458 040 fragments are common for both S+ and S- subsets. 55.0 % of S+ fragments and 91.7 % of S- fragments are present only once in the whole data set S (singletons). Typical ES and HS fragments are shown in Figure 3 and Figure 4, fragments with very low SYBA in Additional file 2 – Figure S13 and compounds containing these fragments in Additional file 2 – Figure S14.

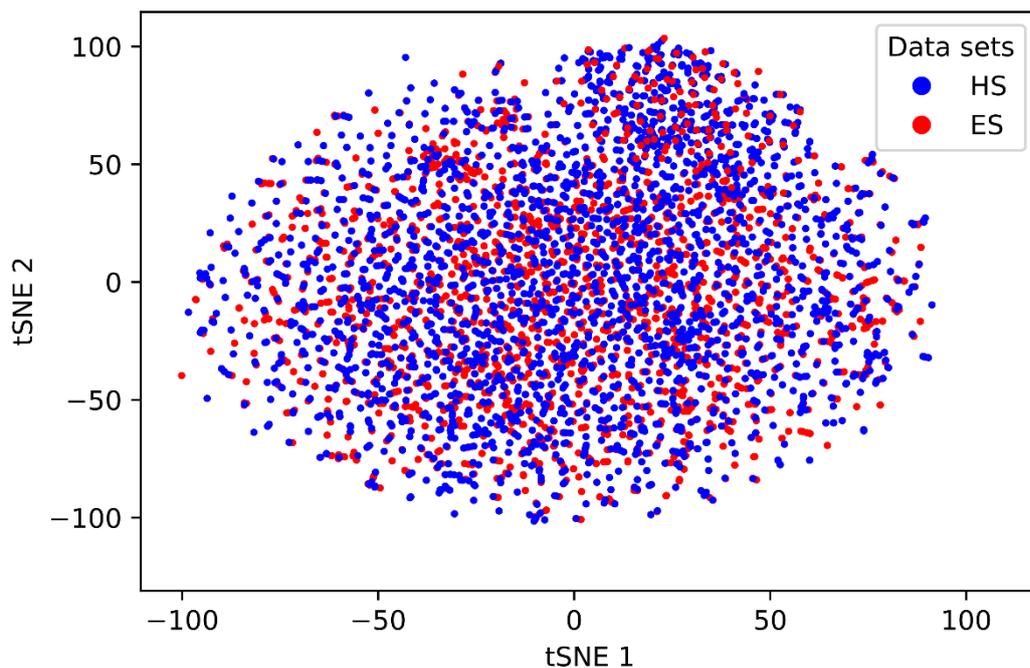


**Figure 3** ES fragments enriched in the S<sub>+</sub> data set. Nine fragments that are most frequent in the S<sub>+</sub> data set and, at the same time, least frequent in the S<sub>-</sub> data set.  $s_i$  is SYBA fragment score. Blue circles represent each fragment central atom, yellow circles represent aromatic atoms. Fragment images were generated by the RDKit function DrawMorganEnvs().



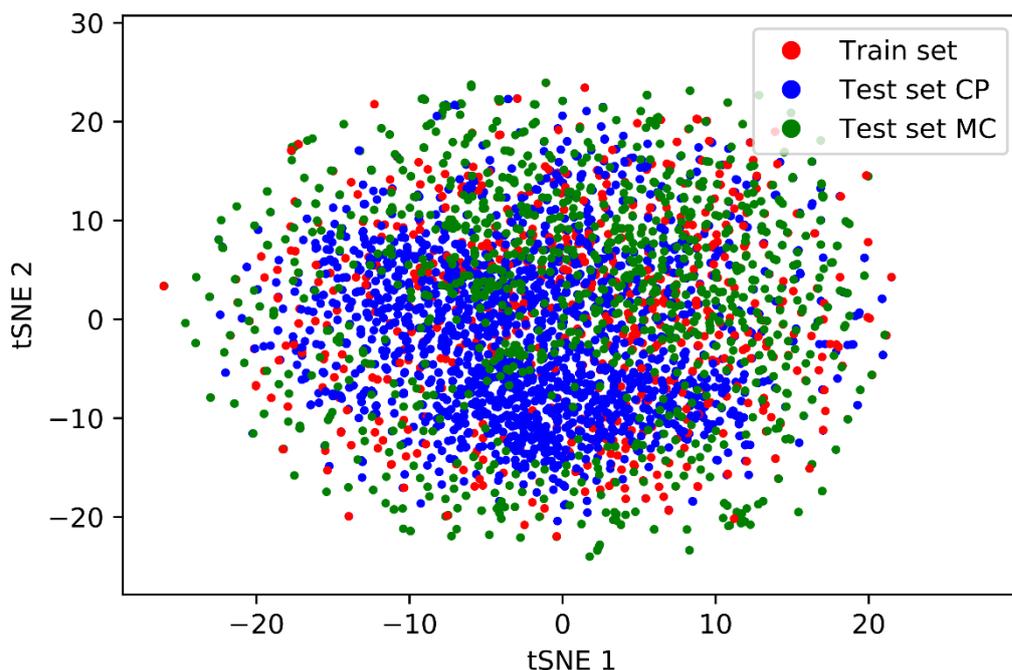
**Figure 4** HS fragments enriched in the S- data set. Nine fragments that are most frequent in the S- data set and, at the same time, least frequent in the S<sub>+</sub> data set.  $s_f$  is SYBA fragment score. Blue circles represent fragment central atom, gray circles represent aliphatic ring atoms. Fragment images were generated by the RDKit function DrawMorganEnvs().

Though the number of fragments in HS compounds is much larger than in ES compounds, chemical space is equally covered by both HS and ES molecules and there is no bias towards HS compounds (Figure 5).



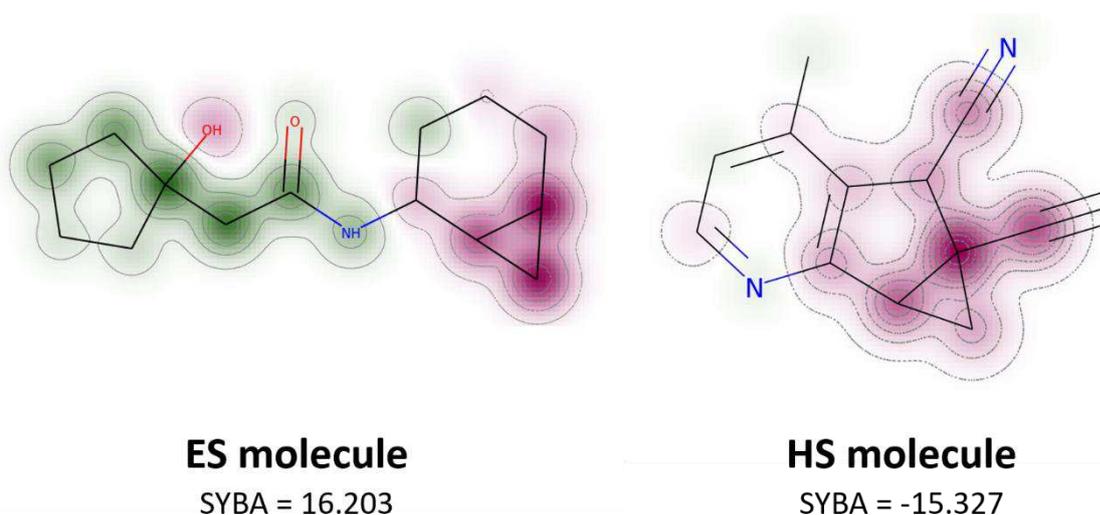
**Figure 5** Chemical space coverage by ES and HS training set compounds. 3000 ES and 3000 HS compounds were randomly selected from the training set and each compound was encoded by 1024 bits long ECFP4 fingerprint. The dimensionality of the input space was reduced by SVD to 500 components that explain 85% of the variance in the data.

The visualization of chemical space (Figure 6) covered by S, T<sub>CP</sub> and T<sub>MC</sub> compounds shows that test set compounds lie within chemical space of training compounds. The examples of T<sub>CP</sub> and T<sub>MC</sub> compounds are given in Additional file 2 – Figure S15-S21.



**Figure 6** Chemical space coverage by training set  $S$  and test sets  $T_{CP}$  and  $T_{MC}$ .  $T_{MC}$  data set consists of 40 HS compounds and 1200 ES compounds, from  $S$  and  $T_{CP}$  data sets random samples of 1240 compounds were generated. Each compound was encoded by 1024 bits long ECFP4 fingerprint. The dimensionality of the input space was reduced by SVD to 500 components that explain 88% of the variance in the data.

SYBA also enables the visualization and interpretation of fragment score contributions. Each SYBA fragment score can be projected to the corresponding fragment root atom and this projection can be used to analyze which fragments contribute unfavorably to molecule synthetic accessibility (Figure 7).



**Figure 7** SYBA fragment score visualization. Fragment score is projected on the fragment root atom and the whole molecule is visualized as a similarity map [74]. The more frequent the fragment is in the  $S_+$  data set compared to the  $S_-$  data set, the greener is its central atom. Similarly, the more frequent the fragment is in the  $S_-$  data set compared to the  $S_+$  data set, the

redder is its central atom. This visualization enables to analyze the contributions of the individual parts of the molecule to its synthetic accessibility. In the HS molecule, the quaternary carbon is most problematic. Another substructure decreasing compound synthetic accessibility is a fused cyclopropane ring as can be observed both in ES and HS compounds.

## Classifier performance on manually curated test set

The differences in classification of  $T_{MC}$  (Figure 1) compounds using SYBA, SAScore and RF with default thresholds are statistically significant. The Cochran's Q test p-value is  $2 \times 10^{-5}$  for the  $T_{MC}$  test set with the smallest SYBA AUC. The results of the corresponding McNemar's paired tests [72] are summarized in Table 1. While RF and SYBA do not differ significantly, both RF and SYBA yield significantly better results than SAScore.

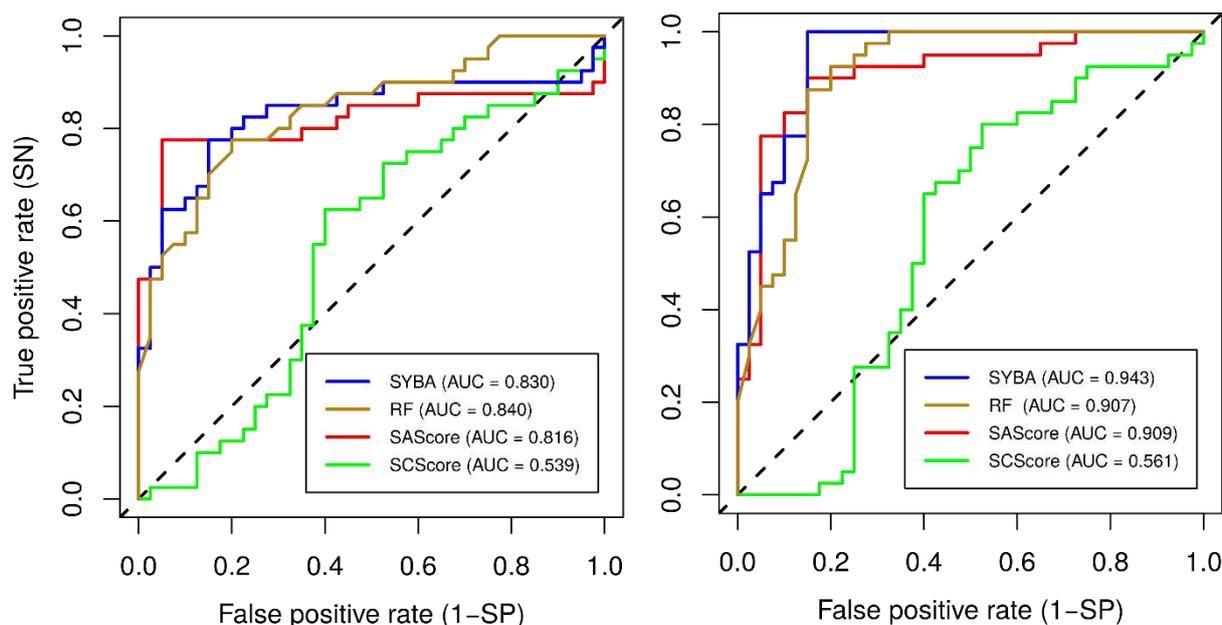
|                  | Adjusted p-value |
|------------------|------------------|
| RF vs. SAScore   | 0.002            |
| RF vs. SYBA      | 1.000            |
| SAScore vs. SYBA | 0.001            |

**Table 1** The results of McNemar's two-sided paired tests for the  $T_{MC}$  test set with the smallest SYBA AUC (AUC = 0.830). The default threshold values were used (0.0 for SYBA, 0.5 for RF, and 6.0 for SAScore). The p-values were adjusted using Benjamini-Hochberg method.

The quality measures of the classification of the compounds in the manually curated  $T_{MC}$  test set, averaged over 30  $T_{MC}$  instances, are summarized in Table 2. The corresponding confusion matrices are reported in Additional file 2 – Panel S1 and Panel S2 and ROC curves are shown in Figure 8.

| Model               | AUC          | Acc          | SN           | SP           | Threshold |
|---------------------|--------------|--------------|--------------|--------------|-----------|
| Default threshold   |              |              |              |              |           |
| SYBA                | <b>0.903</b> | <b>0.844</b> | 0.913        | <b>0.775</b> | 0.0       |
| SAScore             | 0.865        | 0.617        | <b>0.934</b> | 0.300        | 6.0       |
| RF                  | 0.875        | 0.819        | 0.863        | 0.775        | 0.5       |
| Optimized threshold |              |              |              |              |           |
| SYBA                | <b>0.903</b> | <b>0.871</b> | <b>0.902</b> | 0.840        | 19.1      |
| SAScore             | 0.865        | 0.859        | 0.799        | <b>0.919</b> | 3.9       |
| SCScore             | 0.528        | 0.601        | 0.707        | 0.496        | 3.7       |
| RF                  | 0.875        | 0.842        | 0.855        | 0.828        | 0.6       |

**Table 2** The performance of classification models for the manually curated  $T_{MC}$  test set. Quality measures AUC, Acc, SN and SP, as well as thresholds, are reported as their average values over 30  $T_{MC}$  instances.



**Figure 8** The ROC curves of classification models for the manually curated  $T_{MC}$  test set. Out of 30 possible  $T_{MC}$  instances, ROC curves of the  $T_{MC}$  test set with the smallest (left) and largest (right) SYBA AUC are shown.

In terms of *Acc*, the best performing model is SYBA followed by RF and SAScore. While SYBA and RF sensitivity and specificity are well balanced, SAScore shows high sensitivity ( $SN = 0.934$ , i.e., on average 93.4 % of *ES* compounds are predicted as *ES*), while its specificity (i.e., the ability to correctly classify *HS* compounds) is rather low ( $SP = 0.300$ ). The observed high sensitivity of SAScore is not surprising as only 0.2 % of ZINC structures have SAScore greater than 6.0 and out of these, only lower units were selected into the  $T_{MC+}$  set.

For optimized thresholds, the differences between SYBA, SAScore and RF are again statistically significant (Cochran's Q test p-value is  $2 \times 10^{-14}$  for the  $T_{MC}$  test set with the smallest SYBA AUC). However, McNemar's paired test (Table 3) identifies significant differences only between SCScore and other methods meaning that SAScore results improve significantly upon threshold optimization.

|                     | Adjusted p-value   |
|---------------------|--------------------|
| RF vs. SAScore      | 0.164              |
| RF vs. SCScore      | $2 \times 10^{-6}$ |
| RF vs. SYBA         | 0.687              |
| SAScore vs. SCScore | $1 \times 10^{-8}$ |
| SAScore vs. SYBA    | 0.466              |
| SCScore vs. SYBA    | $2 \times 10^{-6}$ |

**Table 3** The results of McNemar's two-sided paired tests for the  $T_{MC}$  test set with the smallest SYBA AUC (AUC = 0.830). The optimized threshold values were used (27.7 for SYBA, 0.5 for RF, 3.7 for SAScore, and 4.0 for SCScore). The p-values were adjusted using Benjamini-Hochberg method.

In terms of performance measures (Table 2), the most accurate classifiers are SYBA, RF and SAScore followed afar by SCScore. Notable is the improvement of SAScore *SP* by 0.619 compared to the default threshold. The increase in SAScore *SP* comes, however, at the cost of *SN* that decreases by 0.135. The worst performing model, SCScore, is only slightly better (*AUC* = 0.528) than a random model. However, because  $T_{MC+}$  and  $T_{MC-}$  data sets consist of only 40 compounds each, the results must be interpreted with caution as small changes in confusion matrices lead to relatively large changes in reported metrics.

### Classifier performance on computationally picked test set

Even stronger evidence of the differences between the models is provided by the classification of compounds in the large computationally picked  $T_{CP}$  test set (Figure 1). Using the default thresholds, the differences between the classifiers are statistically significant (Cochran's Q test p-value <  $10^{-16}$ ) and all classifiers differ significantly (Table 4).

|                  | Adjusted p-value |
|------------------|------------------|
| RF vs. SAScore   | < $10^{-16}$     |
| RF vs. SYBA      | < $10^{-16}$     |
| SAScore vs. SYBA | < $10^{-16}$     |

**Table 4** The results of McNemar's two-sided paired tests for the  $T_{CP}$  test set. The default threshold values were used (0.0 for SYBA, 0.5 for RF, and 6.0 for SAScore). The p-values were adjusted using Benjamini-Hochberg method.

When used with their default thresholds, both SYBA and RF are more accurate than SAScore (Table 5, Additional file 2 – Panel S3). Low observed SAScore accuracy (*Acc* = 0.665) is caused by its low specificity when almost 70 % of HS compounds are predicted to be ES (Table 5).

| Model               | <i>AUC</i>   | <i>Acc</i>   | <i>SN</i>    | <i>SP</i>    | Threshold |
|---------------------|--------------|--------------|--------------|--------------|-----------|
| Default threshold   |              |              |              |              |           |
| SYBA                | 0.903        | <b>0.962</b> | 0.925        | <b>1.000</b> | 0.0       |
| SAScore             | <b>0.999</b> | 0.665        | <b>0.999</b> | 0.331        | 6.0       |
| RF                  | 0.995        | 0.892        | 0.784        | 0.999        | 0.5       |
| Optimized threshold |              |              |              |              |           |
| SYBA                | 0.998        | 0.988        | 0.985        | 0.991        | -18.6     |
| SAScore             | <b>0.999</b> | <b>0.990</b> | <b>0.986</b> | <b>0.993</b> | 4.5       |
| SCScore             | 0.641        | 0.612        | 0.499        | 0.725        | 3.1       |
| RF                  | 0.995        | 0.973        | 0.960        | 0.986        | 0.2       |

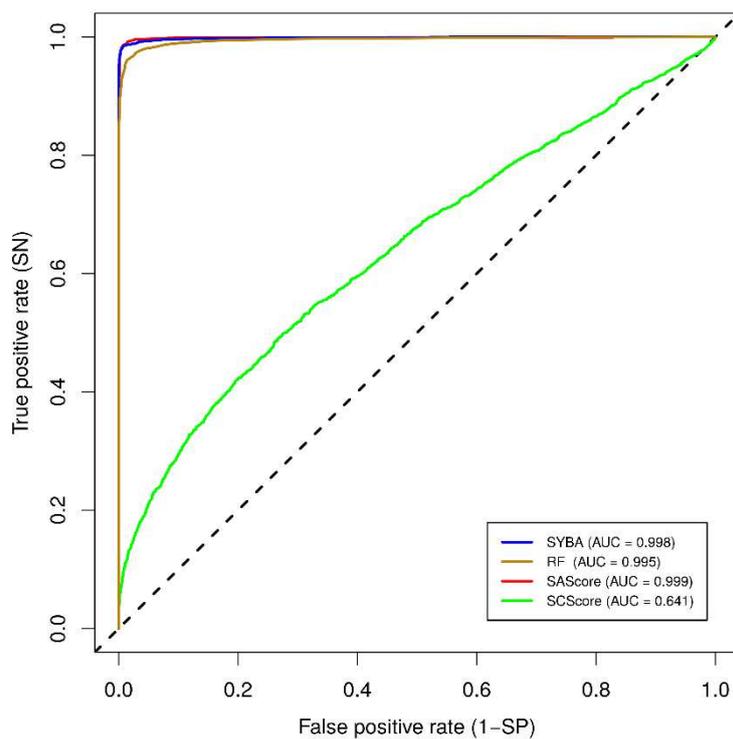
**Table 5** The performance of classification models for the computationally picked  $T_{CP}$  test set.

Cochran’s Q test identifies significant differences (p-value  $<10^{-16}$ ) also if classifiers are used with the optimized thresholds. While statistically significant differences were detected within RF/SAScore and RF/SYBA pairs (Table 6), the effect size is still rather small (Table 5). However, the difference between SCScore and all other methods is statistically significant and large (Table 5, Table 6). No statistically significant difference was observed between SYBA and SAScore meaning that these two methods yield comparable results.

|                     | Adjusted p-value    |
|---------------------|---------------------|
| RF vs. SAScore      | $4 \times 10^{-15}$ |
| RF vs. SCScore      | $<10^{-16}$         |
| RF vs. SYBA         | $2 \times 10^{-14}$ |
| SAScore vs. SCScore | $<10^{-16}$         |
| SAScore vs. SYBA    | 0.474               |
| SCScore vs. SYBA    | $<10^{-16}$         |

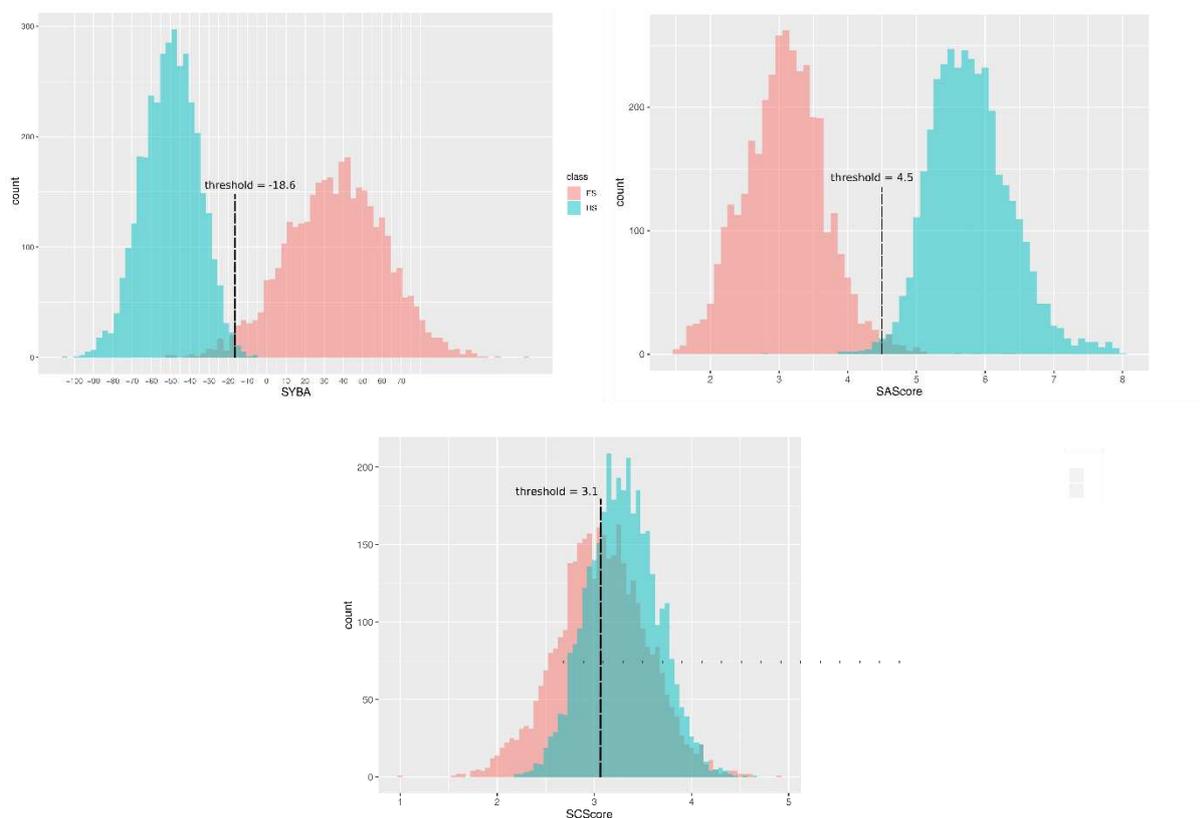
**Table 6** The results of McNemar’s two-sided paired tests for the  $T_{CP}$  test set. The optimized values of threshold (-18.6 for SYBA, 0.2 for RF, 4.5 for SAScore, 3.1 for SCScore) were used. p-values were adjusted using Benjamini-Hochberg method.

Compared to its default threshold of 6.0, SAScore specificity increases by 0.662 (Table 5) when the threshold is shifted to the optimal value of 4.5 (Figure 10). At this threshold, SAScore is on par with SYBA and RF methods (Table 5). However, SYBA retains its high performance over much broader range of threshold values than SAScore (Additional file 2 – Figure S1).



**Figure 9** ROC curves of classification models for the TCP test set.

High performance of SYBA, RF and SAScore is also evident from their *AUC* that is close to one (Figure 9, Table 5). On the other hand, SCScore fails to distinguish between ES and HS compounds as can be deduced from its ROC curve (Figure 9). In its optimal threshold of 3.1, SCScore predicts a majority of TCP compounds as HS (Figure 10, Additional file 2 – Panel S3).

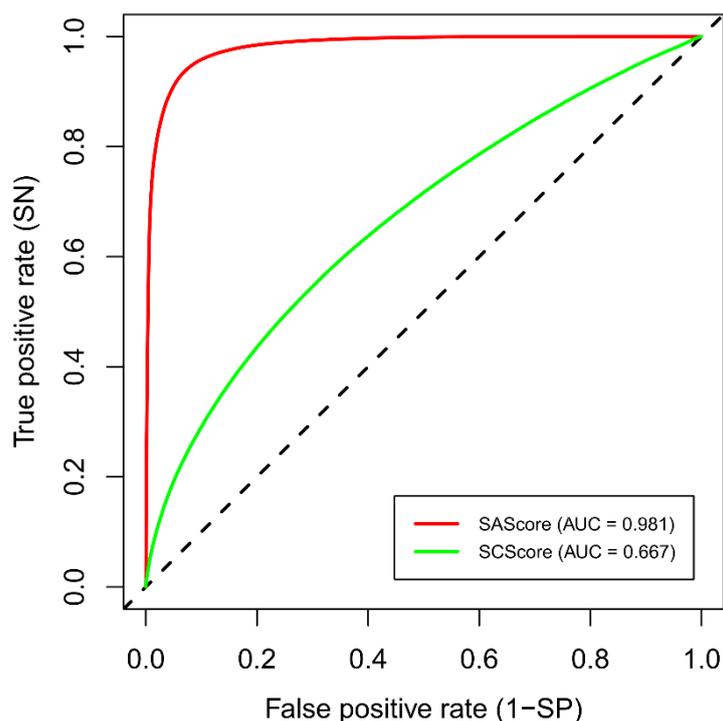


**Figure 10** SYBA, SAScore and SCScore histograms of ES and HS compounds in the computationally picked  $T_{CP}$  test set. The positions of optimal thresholds are shown. SAScore recommended threshold of 6.0 leads to a large number of FP (Additional file 2 – Panel S3). If the threshold is moved to its optimal value of 4.5, SAScore specificity increases from 0.317 to 0.994, i.e. by 0.677.

In addition to the  $T_{MC}$  and  $T_{CP}$  test sets, the performance of SAScore and SCScore was also assessed using the training set S, as this data set was not used for their parametrization. Classification results are shown in Table 7 and Figure 11, confusion matrices are available in Additional file 2 – Panel S4.

| Model               | <i>AUC</i>   | <i>Acc</i>   | <i>SN</i>    | <i>SP</i>    | Threshold |
|---------------------|--------------|--------------|--------------|--------------|-----------|
| Default threshold   |              |              |              |              |           |
| SAScore             | 0.981        | 0.767        | 0.998        | 0.536        | 6.0       |
| Optimized threshold |              |              |              |              |           |
| SAScore             | <b>0.981</b> | <b>0.933</b> | <b>0.935</b> | <b>0.932</b> | 4.4       |
| SCScore             | 0.667        | 0.623        | 0.564        | 0.682        | 3.7       |

**Table 7** The performance of classification models for the training set S.



**Figure 11** ROC curves of classification models for the training set S.

In agreement with previous experiments on the  $T_{MC}$  and  $T_{CP}$  test sets, SAScore is able to distinguish between ES and HS compounds more accurately than SCScore. However, to achieve the best performance, SAScore classification threshold must be shifted from its default value of 6.0 to the optimal value of 4.4. At this threshold, SAScore is both highly sensitive and specific, while SCScore is, using its optimal threshold of 3.7, sensitive and specific only moderately.

The observed poor performance of SCScore in all data sets may follow from the fact that SCScore differs conceptually from other methods tested in the present work. In SCScore, the problem of predicting synthetic complexity is reformulated as the analysis of reactions consisting of reactant-product pairs and SCScore correlates with a number of reaction steps. It means that, contrary to SYBA and SAScore, SCScore captures synthetic feasibility of compounds, not structural complexity. In SCScore derivation, each molecule is analyzed as a whole in the context of all molecules and reactions as they appear in the Reaxys database [44]. Thus, SCScore is biased [43] by the types of reactants and products in the Reaxys database. Therefore, we hypothesize that the unsatisfactory results of SCScore are caused by the fact that compounds in our data sets come from chemical subspace insufficiently covered by Reaxys compounds.

## Conclusions

In the present work, SYBA method for the classification of organic compounds as easy- and hard-to-synthesize is described. SYBA is an additive fragment-based approach meaning that the compound is decomposed into individual substructure fragments, each fragment is assigned its respective SYBA fragment score and these are summed to obtain the final SYBA score. The fragment scores were derived by the Bayesian analysis of the frequency of ECFP8 fragments occurring in the database of ES compounds, that were randomly chosen from the ZINC15 database [57], and HS compounds, that were generated using the Nonpher approach [58]. Because SYBA was derived from ECFP8 fragments that utilize only molecular connectivity and no 3D information, the influence of stereochemistry on synthetic accessibility is not accounted for. However, apart from ECFP8 fragments, the number of stereocenters is included in SYBA calculation and the compounds with many stereocenters are penalized. If the SYBA score is positive, the compound is considered to be ES and vice versa. While SYBA score can theoretically assume values between plus and minus infinity, a majority of compounds will have SYBA score between -100 and +100 in real applications. It must be stressed here that the absolute value of the SYBA score is the measure of the confidence of the prediction and not of the degree of the synthetic accessibility.

SYBA was compared with other two recent classification methods, SAScore [45] and SCScore [43]. As a baseline for the comparison, RF/ECFP4 classifier was used due to its wide adoption in many cheminformatics applications. All methods were assessed using accuracy, sensitivity, specificity and area under the ROC curve. While SYBA and RF provide similar performance, we recommend to use SYBA due to its smaller complexity, lower computational demands and more straightforward analysis of the individual fragment contributions. SYBA outperforms SAScore when this is used with the threshold of 6.0 proposed by the authors [45]. However, if the SAScore threshold is changed to the value of ~4.5, the accuracy of SYBA and SAScore becomes comparable. Therefore, to reduce the number of false positives in workflows that utilize SAScore, we recommend to decrease the SAScore threshold to ~4.5. SYBA, RF and SAScore substantially outperform SCScore. The observed weak performance of SCScore can be, in our opinion, attributed to the fact that our test set compounds

come from a part of chemical space that is insufficiently covered by Reaxys compounds used to derive SCScore. The SYBA fragment scores can be mapped [74] onto a molecule and used for the analysis of the contribution of its individual substructures to the overall synthetic accessibility.

SYBA can be used to quickly rank large molecular data sets that originate, for example, from *de novo* molecular design. However, SYBA is conceptually based on the notion that a compound can be categorized as easy- and hard-to-synthesize. As the synthetic accessibility is a vaguely defined term, SYBA's simplifying approach, though accurate enough, cannot compete with more sophisticated synthetic path-reconstruction methods that enable the incorporation of other factors such as the availability of starting materials, reaction yields or a price aspect. In the end, the definitive assessment of synthetic accessibility is in the hands of experienced organic chemists.

## Authors' contributions

MV, MK and DS conceptualized the problem. MV was responsible for method development, implementation and validation. MV also maintains the SYBA GitHub repository. MK derived SYBA and IČ took part in method testing. DS supervised the study, performed statistical analyses and prepared the manuscript with the active participation of MV, MK and IČ. All authors read and approved the final manuscript.

## Acknowledgements

Computational resources were supplied by the project "e-Infrastruktura CZ" (e-INFRA LM2018140) provided within the program Projects of Large Research, Development and Innovations Infrastructures.

## Competing interests

The authors declare that they have no competing interests.

## Funding

This work was supported from the Ministry of Education of the Czech Republic (RVO 68378050-KAV-NPUI and LM2018130). IČ research was further supported by a specific university research (MSMT No. 20/2015). MK was supported by the project "Center for Tumor Ecology - Research of the Cancer Microenvironment Supporting Cancer Growth and Spread" (reg. No.

CZ.02.1.01/0.0/0.0/16\_019/0000785) within the Operational Programme Research, Development and Education.

## Availability of data and materials

The code used to train and benchmark SYBA model is available from <https://github.com/lich-uct/syba> repository. Nonpher is available from <https://github.com/lich-uct/nonpher>.

## Abbreviations

Acc – accuracy

AUC – area under the ROC curve

ES – easy-to-synthesize

HS – hard-to-synthesize

MW – molecular weight

RF – Random forest

ROC – receiver operating characteristic

S – training set

SN – sensitivity

SP – specificity

SVD – singular value decomposition

SYBA – SYnthetic Bayesian Accessibility

T<sub>CP</sub> – computationally picked test set

T<sub>MC</sub> – manually curated test set

YI – Youden index

## Additional files

Additional file 1 – Training set S. It consists of 693 353 ES compounds selected from the ZINC15 database and of 693 353 ES compounds generated by Nonpher.

Additional file 2 – This supporting document contains the detailed description of the SYBA score derivation, threshold values of complexity indices, RF hyperparameter optimization results, the dependence of the accuracy of the classification of  $T_{CP}$  compounds on SYBA and SAScore thresholds, examples of correctly predicted and mispredicted S,  $T_{MC}$  and  $T_{CP}$  compounds, fragments with very low SYBA contributions and HS compounds containing these fragments, and confusion matrices of the classification of  $T_{MC}$ ,  $T_{CP}$  and S data sets.

Additional file 3 – Manually curated test set ( $T_{MC}$ ). It consists of 40 HS compounds manually selected from scientific papers and of 30 ES sets, each of them contains 40 compounds selected from the ZINC15 database.

Additional file 4 – Computationally picked test set ( $T_{CP}$ ). It consists of 3 581 HS compounds that were obtained from the GDB-17 database complemented by the same number of compounds randomly selected from the ZINC15 database.

## References

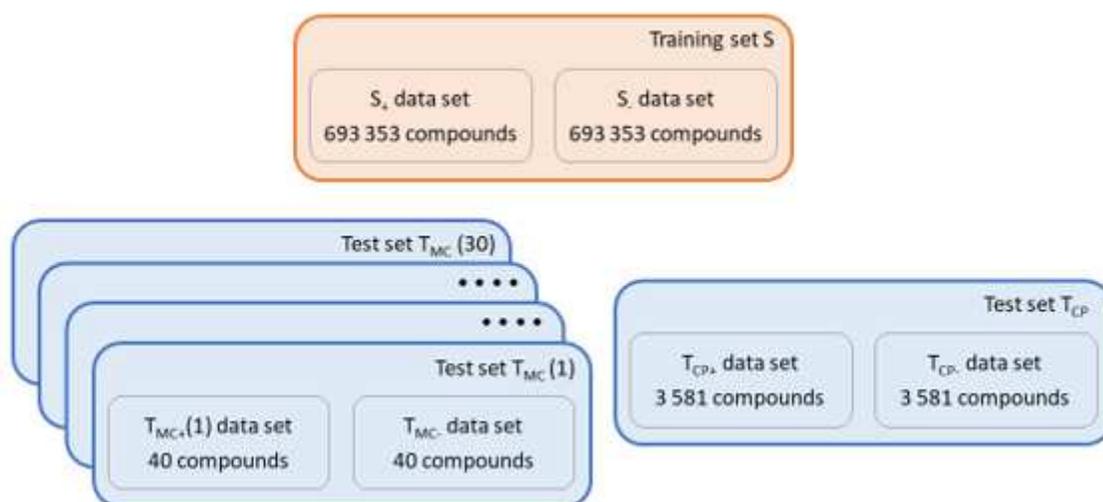
1. Bohacek RS, McMartin C, Guida WC (1996) The art and practice of structure-based drug design: a molecular modeling perspective. *Med Res Rev* 16(1):3-50.
2. Polishchuk PG, Madzhidov TI, Varnek A (2013) Estimation of the size of drug-like chemical space based on GDB-17 data. *J Comput Aided Mol Des* 27(8):675-679.
3. Ertl P (2003) Cheminformatics analysis of organic substituents: identification of the most common substituents, calculation of substituent properties, and automatic identification of drug-like bioisosteric groups. *J Chem Inf Comput Sci* 43(2):374-380.
4. Reymond JL, van Deursen R, Blum LC, Ruddigkeit L (2010) Chemical space as a source for new drugs. *Medchemcomm* 1(1):30-38.
5. Llanos EJ, Leal W, Luu DH, Jost J, Stadler PF, Restrepo G (2019) Exploration of the chemical space and its three historical regimes. *Proc Natl Acad Sci U S A* 116(26):12660-12665.
6. Karlov DS, Sosnin S, Tetko IV, Fedorov MV (2019) Chemical space exploration guided by deep neural networks. *Rsc Advances* 9(9):5151-5157.
7. Gromski PS, Henson AB, Granda JM, Cronin L (2019) How to explore chemical space using algorithms and automation. *Nature Reviews Chemistry* 3(2):119-128.
8. Walters WP (2019) Virtual Chemical Libraries. *J Med Chem* 62(3):1116-1124.
9. Franzini RM, Neri D, Scheuermann J (2014) DNA-encoded chemical libraries: advancing beyond conventional small-molecule libraries. *Acc Chem Res* 47(4):1247-1255.
10. Lopez-Vallejo F, Caulfield T, Martinez-Mayorga K, Giulianotti MA, Nefzi A, Houghten RA, Medina-Franco JL (2011) Integrating virtual screening and combinatorial chemistry for accelerated drug discovery. *Comb Chem High Throughput Screen* 14(6):475-487.
11. Hoffmann T, Gastreich M (2019) The next level in chemical space navigation: going far beyond enumerable compound libraries. *Drug Discov Today* 24(5):1148-1156.
12. van Hilten N, Chevillard F, Kolb P (2019) Virtual Compound Libraries in Computer-Assisted Drug Discovery. *J Chem Inf Model* 59(2):644-651.
13. Schneider G, Fechner U (2005) Computer-based de novo design of drug-like molecules. *Nat Rev Drug Discov* 4(8):649-663.
14. Loving K, Alberts I, Sherman W (2010) Computational approaches for fragment-based and de novo design. *Curr Top Med Chem* 10(1):14-32.
15. Medina-Franco JL, Martinez-Mayorga K, Meurice N (2014) Balancing novelty with confined chemical space in modern drug discovery. *Expert Opin Drug Discov* 9(2):151-165.

16. Schneider P, Schneider G (2016) De Novo Design at the Edge of Chaos. *J Med Chem* 59(9):4077-4086.
17. Kutchukian PS, Shakhnovich EI (2010) De novo design: balancing novelty and confined chemical space. *Expert Opinion on Drug Discovery* 5(8):789-812.
18. Hartenfeller M, Schneider G (2011) De novo drug design. *Methods Mol Biol* 672:299-323.
19. Hartenfeller M, Proschak E, Schuller A, Schneider G (2008) Concept of combinatorial de novo design of drug-like molecules by particle swarm optimization. *Chem Biol Drug Des* 72(1):16-26.
20. Vinkers HM, de Jonge MR, Daeyaert FF, Heeres J, Koymans LM, van Lenthe JH, Lewi PJ, Timmerman H, Van Aken K, Janssen PA (2003) SYNOPSIS: SYNthesize and OPTimize System in Silico. *J Med Chem* 46(13):2765-2773.
21. Hartenfeller M, Zettl H, Walter M, Rupp M, Reisen F, Proschak E, Weggen S, Stark H, Schneider G (2012) DOGS: reaction-driven de novo design of bioactive compounds. *PLoS Comput Biol* 8(2):e1002380.
22. Schneider G, Lee ML, Stahl M, Schneider P (2000) De novo design of molecular architectures by evolutionary assembly of drug-derived building blocks. *J Comput Aided Mol Des* 14(5):487-494.
23. Fechner U, Schneider G (2006) Flux (1): a virtual synthesis scheme for fragment-based de novo design. *J Chem Inf Model* 46(2):699-707.
24. Chen H, Engkvist O, Wang Y, Olivecrona M, Blaschke T (2018) The rise of deep learning in drug discovery. *Drug Discov Today* 23(6):1241-1250.
25. Hessler G, Baringhaus KH (2018) Artificial Intelligence in Drug Design. *Molecules* 23(10).
26. Xu Y, Lin K, Wang S, Wang L, Cai C, Song C, Lai L, Pei J (2019) Deep learning for molecular generation. *Future Med Chem* 11(6):567-597.
27. Popova M, Isayev O, Tropsha A (2018) Deep reinforcement learning for de novo drug design. *Sci Adv* 4(7):eaap7885.
28. Olivecrona M, Blaschke T, Engkvist O, Chen H (2017) Molecular de-novo design through deep reinforcement learning. *J Cheminform* 9(1):48.
29. Blaschke T, Olivecrona M, Engkvist O, Bajorath J, Chen H (2018) Application of Generative Autoencoder in De Novo Molecular Design. *Mol Inform* 37(1-2).
30. Segler MHS, Kogej T, Tyrchan C, Waller MP (2018) Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Cent Sci* 4(1):120-131.
31. Gupta A, Muller AT, Huisman BJH, Fuchs JA, Schneider P, Schneider G (2018) Generative Recurrent Networks for De Novo Drug Design. *Mol Inform* 37(1-2).
32. Merk D, Friedrich L, Grisoni F, Schneider G (2018) De Novo Design of Bioactive Small Molecules by Artificial Intelligence. *Mol Inform* 37(1-2).
33. Mendez-Lucio O, Medina-Franco JL (2017) The many roles of molecular complexity in drug discovery. *Drug Discov Today* 22(1):120-126.
34. Bertz SH (1981) The first general index of molecular complexity. *Journal of the American Chemical Society* 103(12):3599-3601.
35. Whitlock HW (1998) On the Structure of Total Synthesis of Complex Natural Products. *The Journal of Organic Chemistry* 63(22):7982-7989.
36. Barone R, Chanon M (2001) A New and Simple Approach to Chemical Complexity. Application to the Synthesis of Natural Products. *Journal of Chemical Information and Computer Sciences* 41(2):269-272.
37. Allu TK, Oprea TI (2005) Rapid Evaluation of Synthetic and Molecular Complexity for in Silico Chemistry. *Journal of Chemical Information and Modeling* 45(5):1237-1243.
38. Selzer P, Roth HJ, Ertl P, Schuffenhauer A (2005) Complex molecules: do they add value? *Curr Opin Chem Biol* 9(3):310-316.
39. Sheridan RP, Zorn N, Sherer EC, Campeau LC, Chang CZ, Cumming J, Maddess ML, Nantermet PG, Sinz CJ, O'Shea PD (2014) Modeling a crowdsourced definition of molecular complexity. *J Chem Inf Model* 54(6):1604-1616.
40. Gillet VJ, Myatt G, Zsoldos Z, Johnson AP (1995) SPROUT, HIPPO and CAESA: Tools for de novo structure generation and estimation of synthetic accessibility. *Perspectives in Drug Discovery and Design* 3:34-50.

41. Huang Q, Li L-L, Yang S-Y (2011) RASA: A Rapid Retrosynthesis-Based Scoring Method for the Assessment of Synthetic Accessibility of Drug-like Molecules. *J Chem Inf Model* 51(10):2768-2777.
42. Li J, Eastgate MD (2015) Current complexity: a tool for assessing the complexity of organic molecules. *Org Biomol Chem* 13(26):7164-7176.
43. Coley CW, Rogers L, Green WH, Jensen KF (2018) SCScore: Synthetic Complexity Learned from a Reaction Corpus. *Journal of Chemical Information and Modeling* 58(2):252-261.
44. Reaxys. <https://www.reaxys.com>. Accessed 24th January 2020.
45. Ertl P, Schuffenhauer A (2009) Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *J Cheminf* 1:1-11.
46. Rogers D, Hahn M (2010) Extended-Connectivity Fingerprints. *J Chem Inf Model* 50(5):742-754.
47. Kim S, Chen J, Cheng T, Gindulyte A, He J, He S, Li Q, Shoemaker BA, Thiessen PA, Yu B *et al* (2019) PubChem 2019 update: improved access to chemical data. *Nucleic Acids Research* 47(D1):D1102-D1109.
48. Besnard J, Ruda GF, Setola V, Abecassis K, Rodriguiz RM, Huang XP, Norval S, Sassano MF, Shin AI, Webster LA *et al* (2012) Automated design of ligands to polypharmacological profiles. *Nature* 492(7428):215-220.
49. Yang X, Zhang J, Yoshizoe K, Terayama K, Tsuda K (2017) ChemTS: an efficient python library for de novo molecular generation. *Sci Technol Adv Mater* 18(1):972-976.
50. Chevillard F, Kolb P (2015) SCUBIDOO: A Large yet Screenable and Easily Searchable Database of Computationally Created Chemical Compounds Optimized toward High Likelihood of Synthetic Tractability. *J Chem Inf Model* 55(9):1824-1835.
51. Clark AM, Dole K, Coulon-Spektor A, McNutt A, Grass G, Freundlich JS, Reynolds RC, Ekins S (2015) Open Source Bayesian Models. 1. Application to ADME/Tox and Drug Discovery Datasets. *J Chem Inf Model* 55(6):1231-1245.
52. Xia X, Maliski EG, Gallant P, Rogers D (2004) Classification of kinase inhibitors using a Bayesian model. *J Med Chem* 47(18):4463-4470.
53. Bender A (2011) Bayesian methods in virtual screening and chemical biology. *Methods Mol Biol* 672:175-196.
54. Vogt M, Bajorath J (2007) Introduction of an information-theoretic method to predict recovery rates of active compounds for Bayesian in silico screening: theory and screening trials. *J Chem Inf Model* 47(2):337-341.
55. Koutsoukas A, Lowe R, Kalantarmotamedi Y, Mussa HY, Klaffke W, Mitchell JB, Glen RC, Bender A (2013) In silico target predictions: defining a benchmarking data set and comparison of performance of the multiclass Naive Bayes and Parzen-Rosenblatt window. *J Chem Inf Model* 53(8):1957-1966.
56. Irwin JJ, Sterling T, Mysinger MM, Bolstad ES, Coleman RG (2012) ZINC: A Free Tool to Discover Chemistry for Biology. *J Chem Inf Model* 52(7):1757-1768.
57. Sterling T, Irwin JJ (2015) ZINC 15 – Ligand Discovery for Everyone. *J Chem Inf Model* 55(11):2324-2337.
58. Voršilák M, Svozil D (2017) Nonpher: computational method for design of hard-to-synthesize structures. *J Cheminf* 9(1):1-20.
59. Hoksza D, Skoda P, Vorsilak M, Svozil D (2014) Molpher: a software framework for systematic chemical space exploration. *J Cheminf* 6:1-13.
60. RDKit: Open-source cheminformatics. <http://www.rdkit.org>. Accessed 24th January 2020.
61. Ruddigkeit L, van Deursen R, Blum LC, Reymond J-L (2012) Enumeration of 166 Billion Organic Small Molecules in the Chemical Universe Database GDB-17. *J Chem Inf Model* 52(11):2864-2875.
62. Boda K, Seidel T, Gasteiger J (2007) Structure and reaction based evaluation of synthetic accessibility. *J Comput-Aided Mol Des* 21(6):311-325.
63. Fukunishi Y, Kurosawa T, Mikami Y, Nakamura H (2014) Prediction of Synthetic Accessibility Based on Commercially Available Compound Databases. *J Chem Inf Model* 54(12):3259-3267.

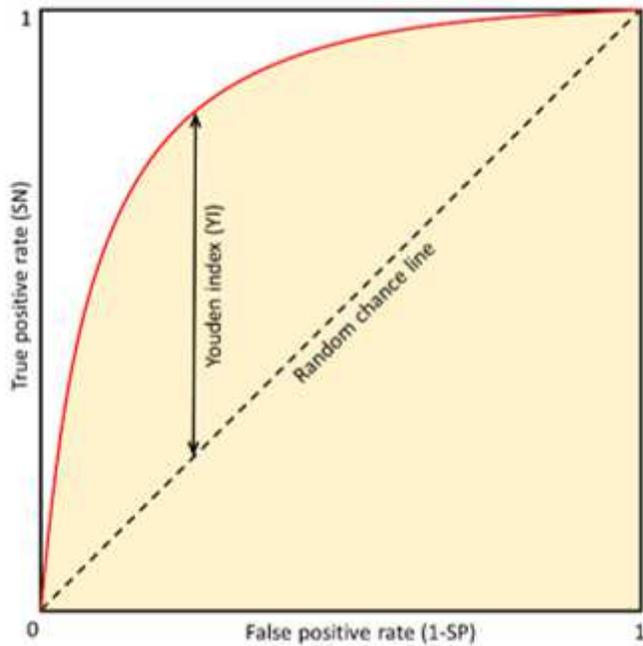
64. Sheridan RP (2013) Using random forest to model the domain applicability of another random forest model. *J Chem Inf Model* 53(11):2837-2850.
65. Kensert A, Alvarsson J, Norinder U, Spjuth O (2018) Evaluating parameters for ligand-based modeling with random forest on sparse data sets. *J Cheminform* 10(1):49.
66. Svetnik V, Liaw A, Tong C, Culberson JC, Sheridan RP, Feuston BP (2003) Random forest: a classification and regression tool for compound classification and QSAR modeling. *J Chem Inf Comput Sci* 43(6):1947-1958.
67. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V *et al* (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research (J Mach Learn Res)* 12:2825-2830.
68. SCScore GitHub. <https://github.com/connorcoley/scscore>. Accessed 24th January 2020.
69. Youden WJ (1950) Index for rating diagnostic tests. *Cancer* 3(1):32-35.
70. Fluss R, Faraggi D, Reiser B (2005) Estimation of the Youden Index and its associated cutoff point. *Biom J* 47(4):458-472.
71. Looney SW (1988) A Statistical Technique for Comparing the Accuracies of Several Classifiers. *Pattern Recognition Letters* 8(1):5-9.
72. Westfall PH, Troendle JF, Pennello G (2010) Multiple McNemar tests. *Biometrics* 66(4):1185-1191.
73. Benjamini Y, Hochberg Y (1995) Controlling the False Discovery Rate - a Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-Statistical Methodology* 57(1):289-300.
74. Riniker S, Landrum GA (2013) Similarity maps - a visualization strategy for molecular fingerprints and machine-learning methods. *J Cheminform* 5(1):43.

# Figures



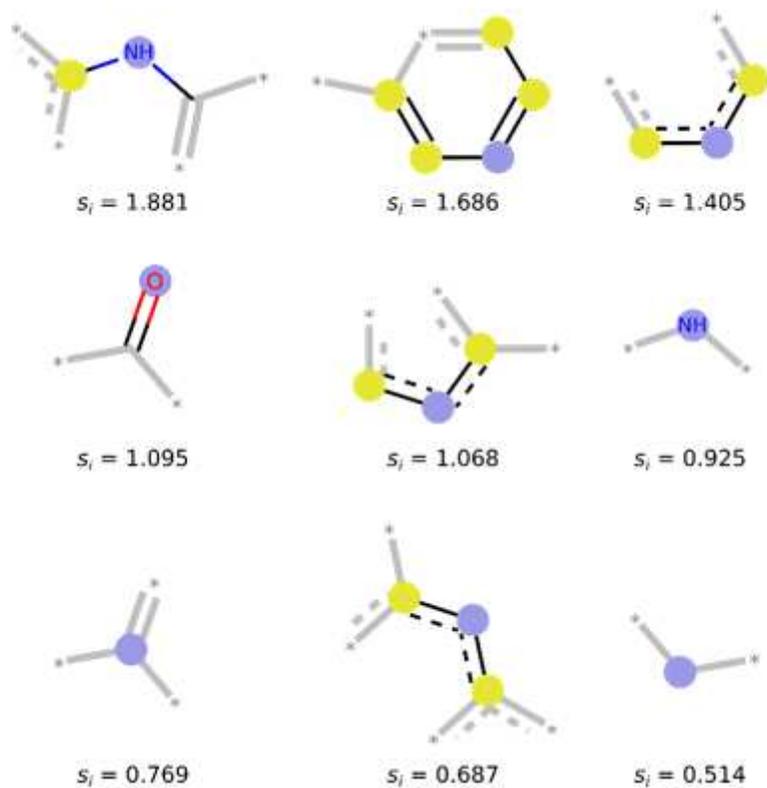
**Figure 1**

Data set summary. Training set was used to derive SYBA scores, as well as to train a random forest classifier. Training set consists of 693 353 molecules randomly selected from the ZINC15 database [57] that are considered to be ES (S+ data set) and of the same number of HS molecules generated by Nonpher [58] (S- data set). Two test sets were used to compare the performance of SYBA, a random forest, SAScore [45] and SCScore [43]. Manually curated test set (TMC) contains 40 compounds (TMC- data set) considered to be HS by experienced medicinal chemists [58] supplemented by 40 ES compounds randomly selected from the ZINC15 database (TMC+ data set). 30 TMC data set instances differing in TMC+ compounds were constructed. Computationally picked test set (TCP) consists of 3 581 HS compounds that were obtained from the GDB-17 database [61] (TCP- data set) complemented by the same number of compounds randomly selected from the ZINC15 database (TCP+ data set).



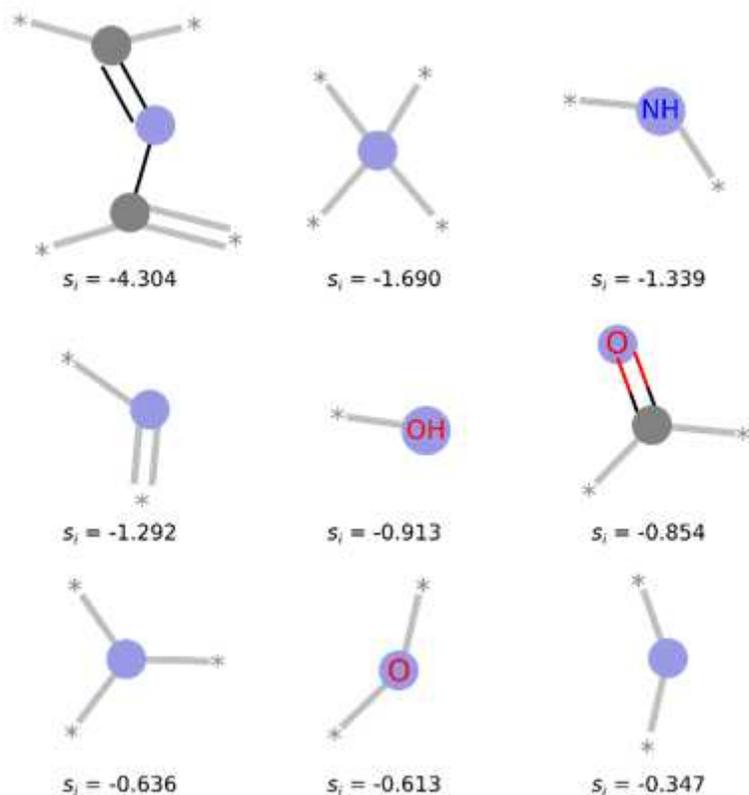
**Figure 2**

ROC curve and Youden index. The ROC curve (red line) is the dependency of true positive rate (it equals to SN) on false positive rate (it equals to  $1 - SP$ ) at various thresholds. The random chance line represents a classifier that assigns examples into individual classes randomly. Orange shaded area represents the area under the ROC curve (AUC). The larger the AUC, the better is the overall performance of the classifier. Youden index (YI) is the point on the ROC curve that is farthest from the random chance line along the SN axis.



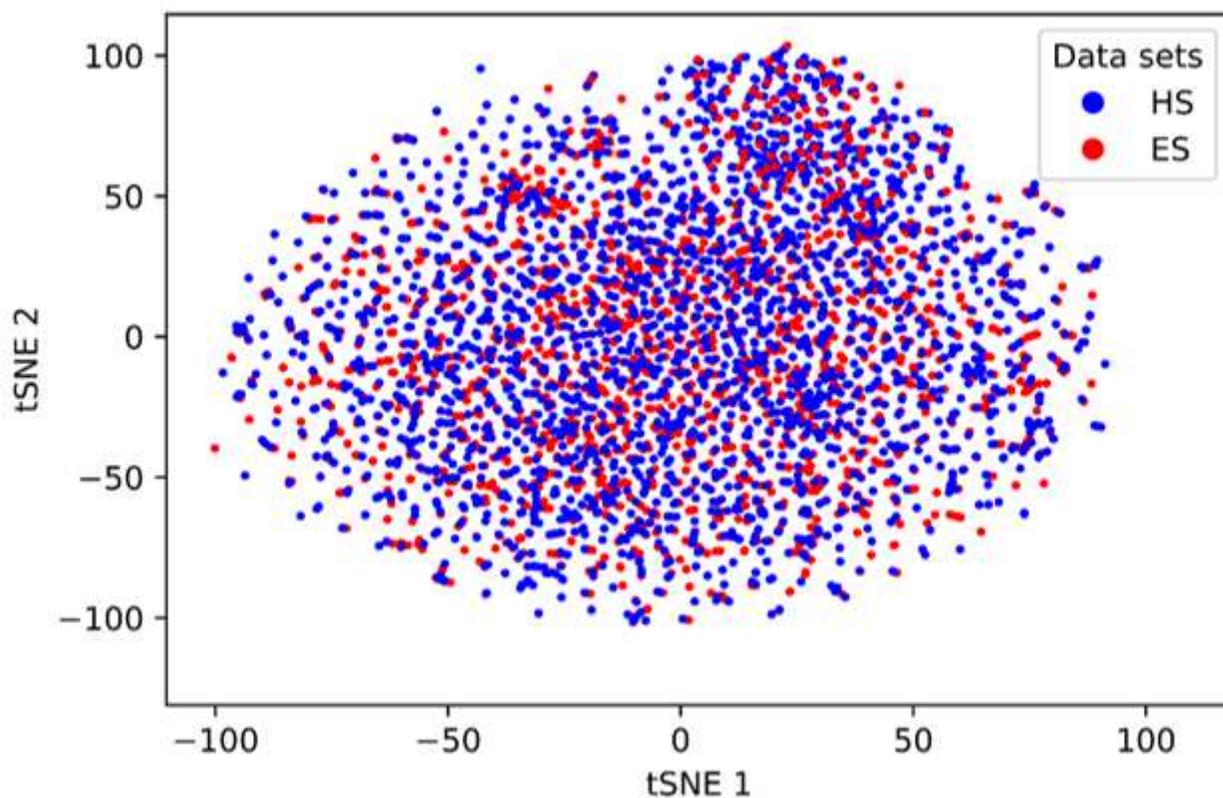
**Figure 3**

ES fragments enriched in the S+ data set. Nine fragments that are most frequent in the S+ data set and, at the same time, least frequent in the S- data set.  $s_i$  is SYBA fragment score. Blue circles represent each fragment central atom, yellow circles represent aromatic atoms. Fragment images were generated by the RDKit function `DrawMorganEnvs()`.



**Figure 4**

HS fragments enriched in the S- data set. Nine fragments that are most frequent in the S- data set and, at the same time, least frequent in the S+ data set.  $s_i$  is SYBA fragment score. Blue circles represent fragment central atom, gray circles represent aliphatic ring atoms. Fragment images were generated by the RDKit function `DrawMorganEnvs()`.



**Figure 5**

Chemical space coverage by ES and HS training set compounds. 3000 ES and 3000 HS compounds were randomly selected from the training set and each compound was encoded by 1024 bits long ECFP4 fingerprint. The dimensionality of the input space was reduced by SVD to 500 components that explain 85% of the variance in the data.

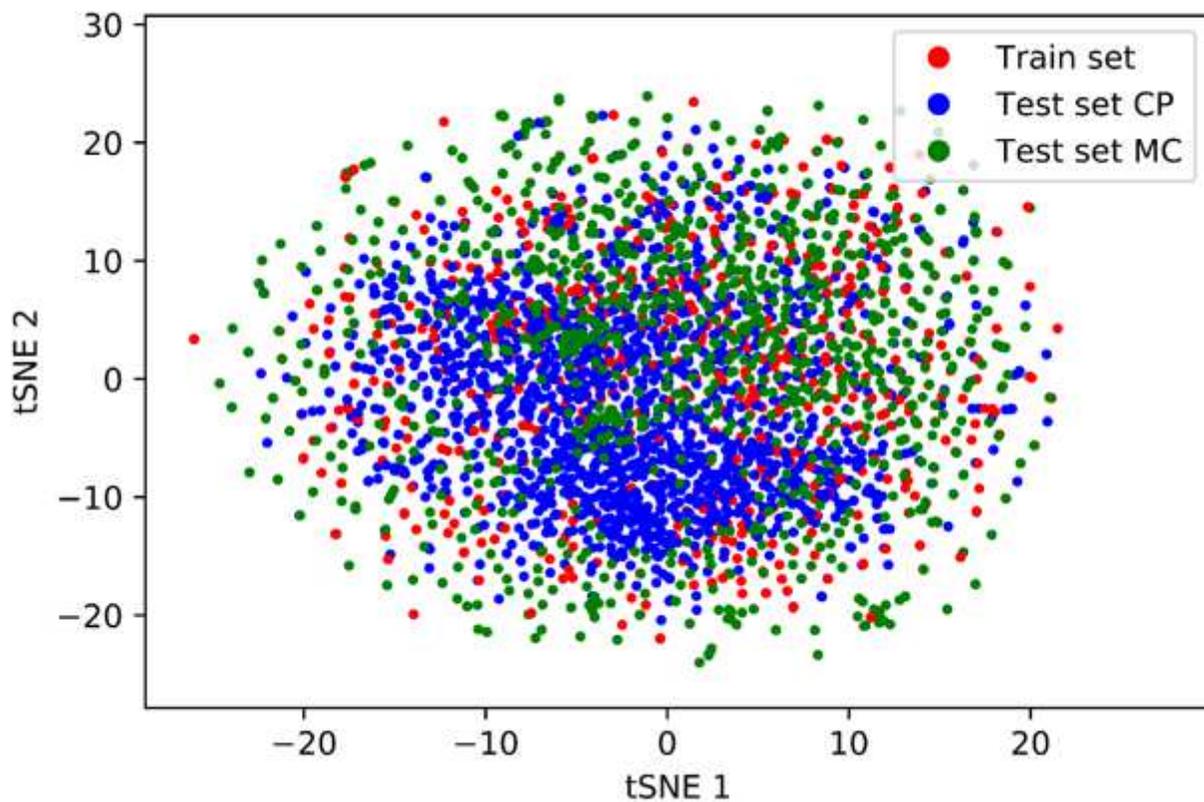


Figure 6

Chemical space coverage by training set S and test sets TCP and TMC. TMC data set consists of 40 HS compounds and 1200 ES compounds, from S and TCP data sets random samples of 1240 compounds were generated. Each compound was encoded by 1024 bits long ECFP4 fingerprint. The dimensionality of the input space was reduced by SVD to 500 components that explain 88% of the variance in the data.

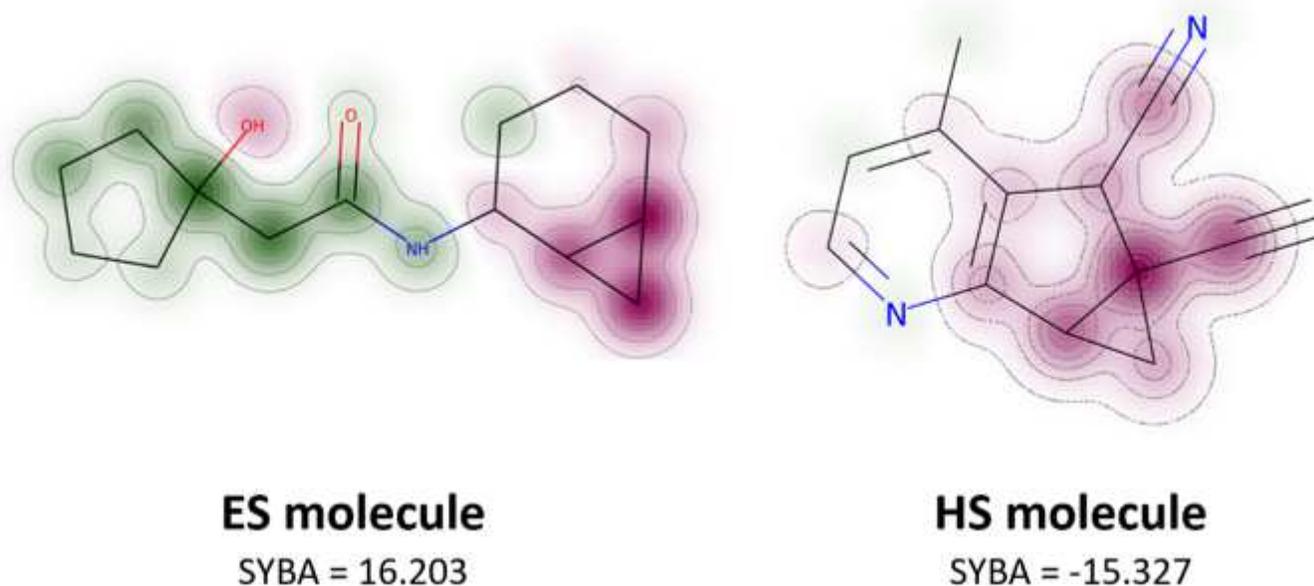
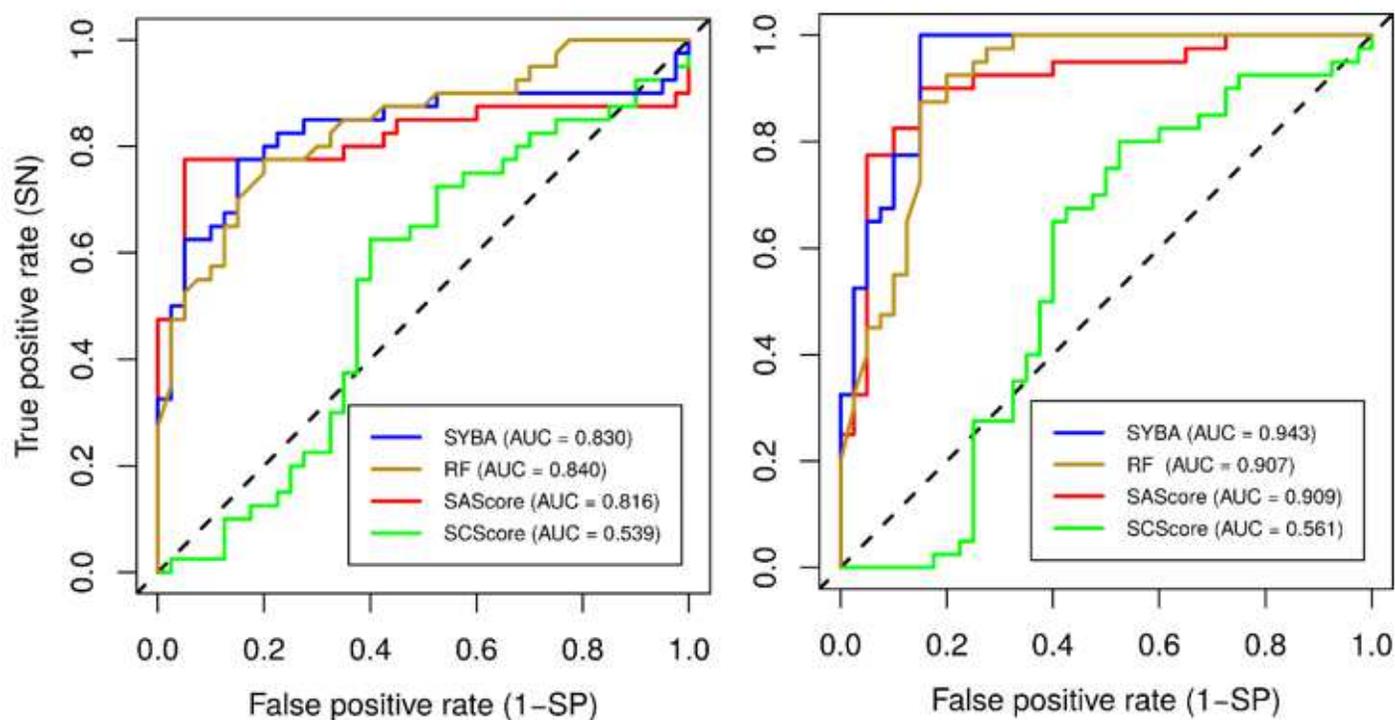


Figure 7

SYBA fragment score visualization. Fragment score is projected on the fragment root atom and the whole molecule is visualized as a similarity map [74]. The more frequent the fragment is in the S+ data set compared to the S- data set, the greener is its central atom. Similarly, the more frequent the fragment is in the S- data set compared to the S+ data set, the redder is its central atom. This visualization enables to analyze the contributions of the individual parts of the molecule to its synthetic accessibility. In the HS molecule, the quaternary carbon is most problematic. Another substructure decreasing compound synthetic accessibility is a fused cyclopropane ring as can be observed both in ES and HS compounds.



**Figure 8**

The ROC curves of classification models for the manually curated TMC test set. Out of 30 possible TMC instances, ROC curves of the TMC test set with the smallest (left) and largest (right) SYBA AUC are shown.

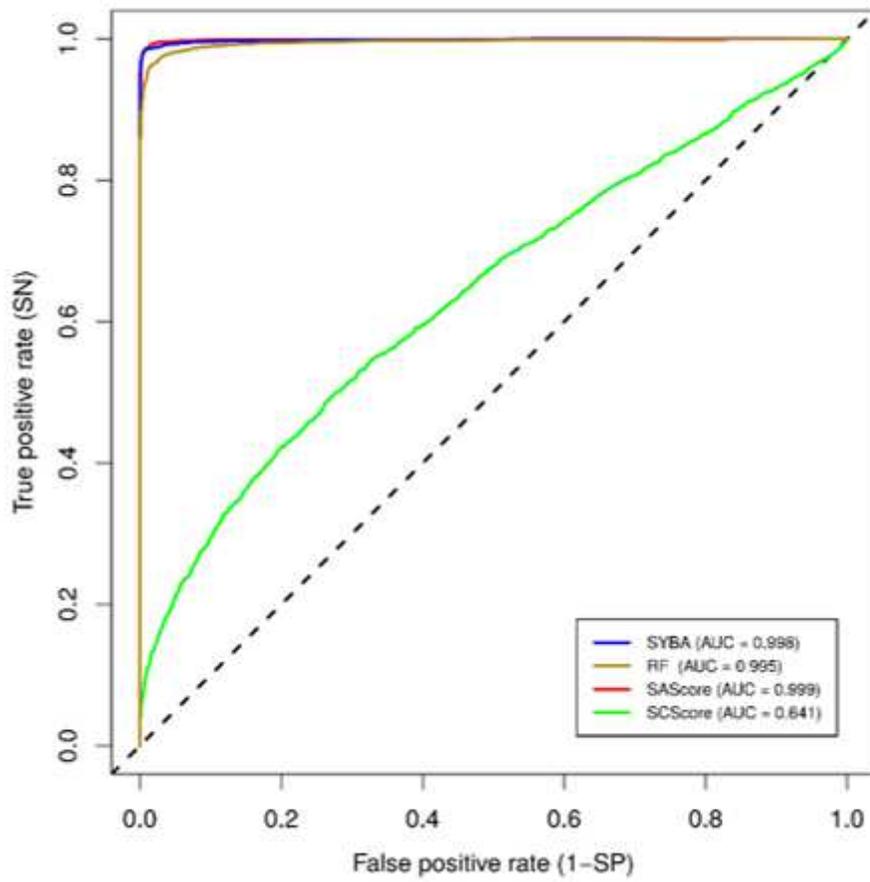
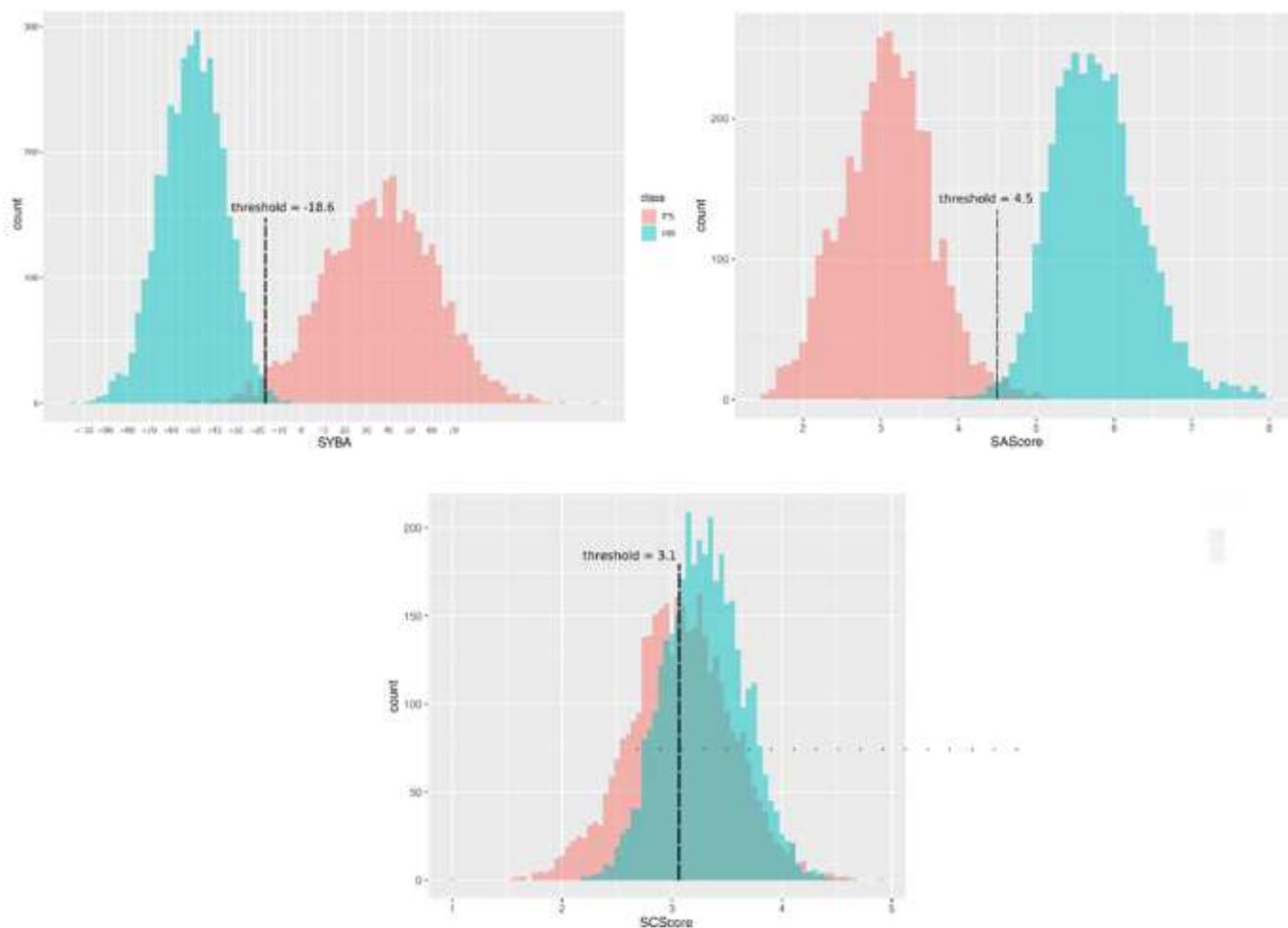


Figure 9

ROC curves of classification models for the TCP test set.



**Figure 10**

SYBA, SAScore and SCScore histograms of ES and HS compounds in the computationally picked TCP test set. The positions of optimal thresholds are shown. SAScore recommended threshold of 6.0 leads to a large number of FP (Additional file 2 – Panel S3). If the threshold is moved to its optimal value of 4.5, SAScore specificity increases from 0.317 to 0.994, i.e. by 0.677.

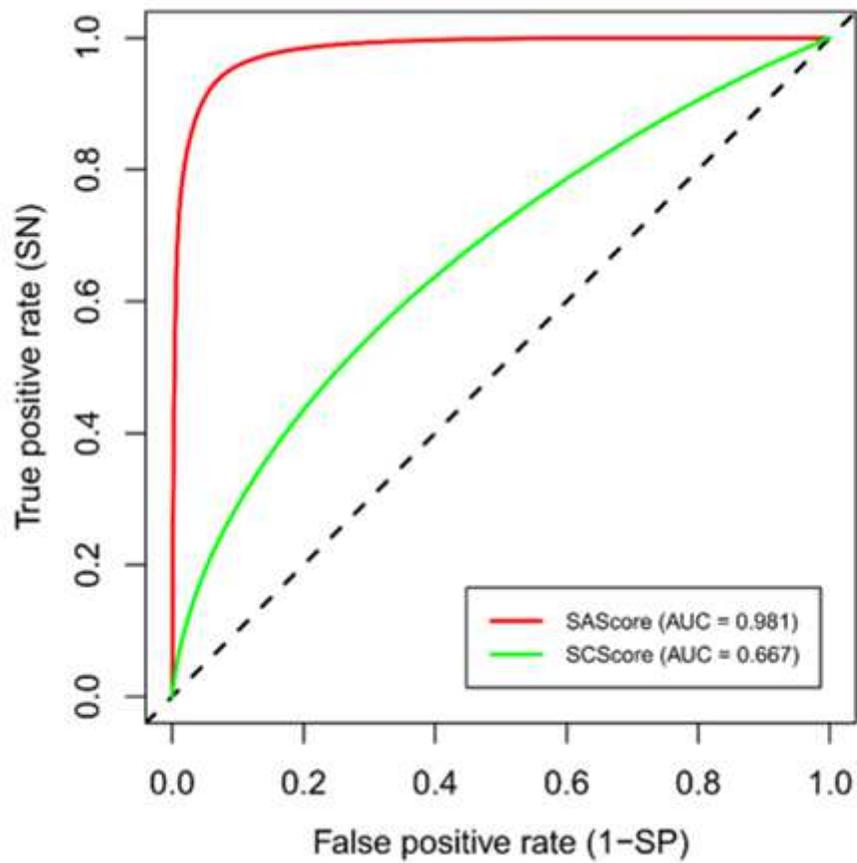


Figure 11

ROC curves of classification models for the training set S.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile2.docx](#)
- [Additionalfile1.xlsx](#)
- [Additionalfile4.xlsx](#)
- [Additionalfile3.xlsx](#)