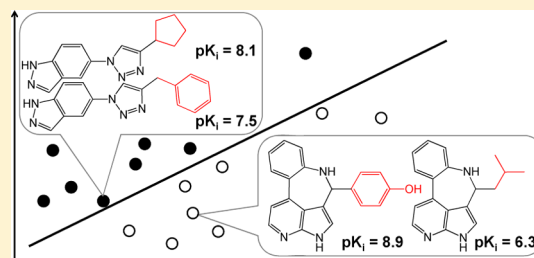


# Prediction of Activity Cliffs Using Support Vector Machines

Kathrin Heikamp,<sup>†,§</sup> Xiaoying Hu,<sup>†,‡,§</sup> Aixia Yan,<sup>‡</sup> and Jürgen Bajorath<sup>\*,†</sup><sup>†</sup>Department of Life Science Informatics, B-IT, LIMES Program Unit Chemical Biology and Medicinal Chemistry, Rheinische Friedrich-Wilhelms-Universität, Dahlmannstr. 2, D-53113 Bonn, Germany<sup>‡</sup>State Key Laboratory of Chemical Resource Engineering, Department of Pharmaceutical Engineering, P.O. Box 53, Beijing University of Chemical Technology, 15 BeiSanHuan East Road, Beijing 100029, People's Republic of China

**ABSTRACT:** Activity cliffs are formed by pairs of structurally similar compounds that act against the same target but display a significant difference in potency. Such activity cliffs are the most prominent features of activity landscapes of compound data sets and a primary focal point of structure–activity relationship (SAR) analysis. The search for activity cliffs in various compound sets has been the topic of a number of previous investigations. So far, activity cliff analysis has concentrated on data mining for activity cliffs and on their graphical representation and has thus been descriptive in nature. By contrast, approaches for activity cliff prediction are currently not available. We have derived support vector machine (SVM) models to successfully predict activity cliffs. A key aspect of the approach has been the design of new kernels to enable SVM classification on the basis of molecule pairs, rather than individual compounds. In test calculations on different data sets, activity cliffs have been accurately predicted using specifically designed structural representations and kernel functions.



## 1. INTRODUCTION

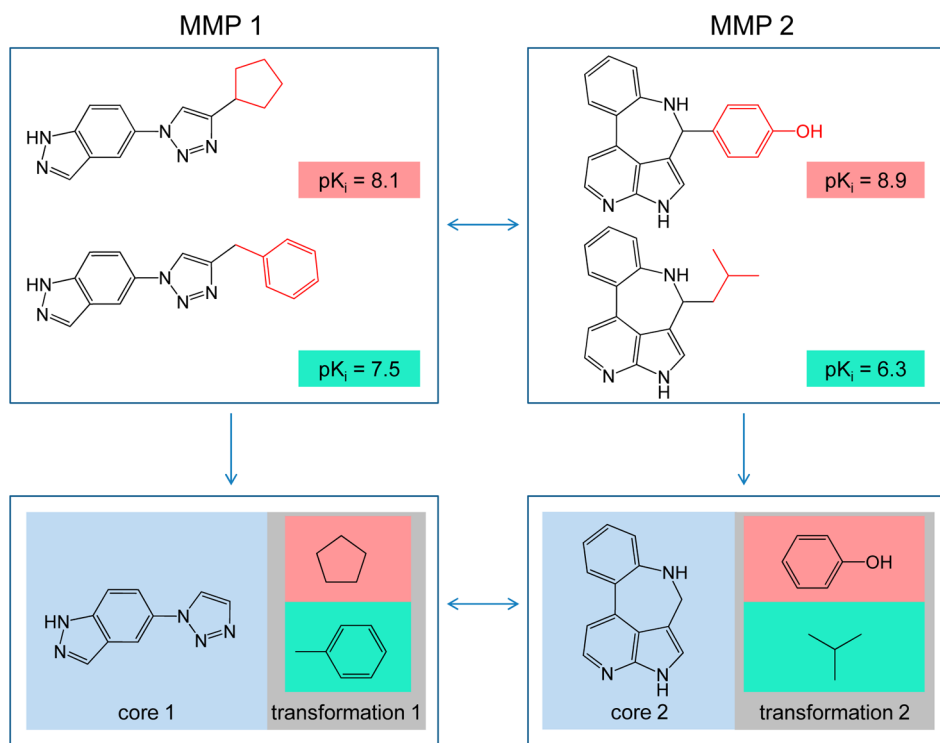
In medicinal chemistry and chemoinformatics, the study of activity cliffs has experienced increasing interest in recent years.<sup>1</sup> Activity cliffs are generally defined as pairs or groups of chemically similar compounds with large potency differences (i.e., usually at least 2 orders of magnitude).<sup>1,2</sup> In chemoinformatics, the exploration of activity cliffs is a topic of interest because qualifying compound pairs can be identified through mining of compound data sets, hence enabling large-scale SAR analysis.<sup>3</sup> In medicinal chemistry, activity cliffs and their structural neighborhoods are considered a prime source of SAR information, given that small chemical differences lead to large bioactivity effects.<sup>1</sup> In traditional medicinal chemistry, activity cliffs are often analyzed in individual compound series. However, they are also systematically explored. For example, in independent studies, activity cliffs were systematically identified and characterized in different data sets.<sup>3–6</sup> In these investigations, cliffs were often defined and represented in rather different ways. Furthermore, activity cliff distributions in current bioactive compounds have been determined through systematic data mining.<sup>7</sup> Moreover, many compound data sets have also been searched for higher-order activity cliff arrangements such as activity ridges<sup>8</sup> and coordinated cliffs.<sup>9</sup> Taken together, these investigations were primarily focused on compound data mining for and visualization of activity cliffs, as mentioned above. Clearly, activity cliff analyses available at present are descriptive in nature, as an integral part of large-scale SAR exploration.<sup>3</sup> By contrast, no attempts have thus far been reported to develop computational models for activity cliff prediction. Here, we present a first step in this direction. Support vector machine (SVM)<sup>10</sup> models have been developed

to screen compound data sets and predict activity cliffs. In the following, the derivation and evaluation of these models and the underlying methodology are described in detail.

## 2. ACTIVITY CLIFF REPRESENTATION AND DATA SETS

**2.1. Compound Pairs.** For any study of activity cliffs, a molecular representation must be selected that provides a basis for pairwise compound similarity assessment. For our analysis, we have applied the concept of matched molecular pairs (MMPs)<sup>11</sup> to represent activity cliffs, following the recent introduction of MMP-based cliffs.<sup>12</sup> An MMP is defined as a pair of compounds that differ only at a single site, i.e., a substructure such as a ring or an R group. Hence, two compounds forming an MMP share a common core and are distinguished by a molecular transformation, i.e., the exchange of a pair of substructures, which converts one compound into the other. The exchange of a substructure can also induce changes in physicochemical properties such as, for example, lipophilicity, charge, or hydrogen bond potential. Compared to other similarity measures, an advantage of the MMP formalism in the context of activity cliff analysis is that the structural difference between compounds in a pair is well-defined and limited to a single substructure. This represents a clearly defined and chemically intuitive criterion for cliff formation that does not rely on calculated similarity values. Furthermore, this approach is consistent with the basic idea of the activity cliff

Received: July 2, 2012



**Figure 1.** MMP comparison. Two MMPs are compared. On the basis of compound potency differences, MMP 1 is an MMP-nonCliff, whereas MMP 2 represents an MMP-cliff. In the upper panels, transformation substructures are shown in red. In the lower panels, the MMPs are divided into the common core (blue background) and the molecular transformation (gray background). Substructures originating from the compounds with higher and lower potency are highlighted (red and green background, respectively).

concept that compounds must be similar; i.e., structural differences must be limited.

MMPs were derived using an in-house Java implementation of the Hussain and Rea algorithm.<sup>13</sup> MMP generation was restricted to molecular transformations of terminal groups; i.e., only single bond cuts were considered. Furthermore, the maximal size of exchanged substructures was restricted to 13 heavy atoms, and the maximal size difference was limited to eight heavy atoms.<sup>12</sup> Furthermore, we concentrated on the smallest of all possible transformations to define a given MMP. Consequently, MMP core structures consisted of coherent fragments, for which other molecular representations could be calculated, and typically small substituents. For model building, as described below, MMPs were either represented as pairs of complete compounds or, alternatively, only by the transformations defining them.

In the following, we use the term ‘substructures’ to refer to fragments exchanged during a transformation and ‘core structure’ to refer to the common core of MMPs.

**2.2. Compound Data Sets.** Nine compound data sets were extracted from BindingDB.<sup>14,15</sup> The data sets were selected because they yielded large numbers of MMP-cliffs that were exclusively formed by compounds with at least 10  $\mu$ M potency on the basis of  $K_i$  measurements. If several  $K_i$  values were available for a compound, the geometric mean was calculated as the final potency annotation. For fingerprint calculations, only compounds in which all atoms were assigned to Sybyl atom types were considered.<sup>16</sup> These atom types were used to enable calculations with a combinatorial feature fingerprint, as described below. Additionally, an MMP was omitted from the calculations if a chosen molecular representation (see below) did not unambiguously specify the underlying transformation.

For each data set, the resulting MMPs were divided into MMPs forming activity cliffs (MMP-cliffs), MMP-nonCliffs, and other MMPs based on the following potency difference criteria: To qualify as an MMP-cliff, compounds forming the pair were required to have a potency difference of at least 2 orders of magnitude. To control the potential influence of potency boundary effects on activity cliff prediction, the potency difference of compounds forming an MMP-nonCliff was limited to at most 1 order of magnitude. Accordingly, MMPs with compounds having a potency difference between 1 and 2 orders of magnitude were not further considered for SVM modeling and were assigned to the class of ‘other MMPs’.

The partition of an MMP into its common core and transformation is illustrated in Figure 1 for two exemplary MMPs forming an MMP-cliff and MMP-nonCliff, respectively. Compound sets and MMP statistics are reported in Table 1. The top five data sets in Table 1 contained the largest number of MMP-cliffs. These data sets were relatively unbalanced because the ratio of MMP-nonCliffs to MMP-cliffs varied between 6 and 21. Because data sets of unbalanced composition typically present a difficult scenario of SVM modeling,<sup>17,18</sup> we also selected four more balanced data sets (ranks six to nine in Table 1). In these cases, the MMP-nonCliff/MMP-cliff ratio was less than 4.

### 3. SUPPORT VECTOR MACHINE MODELING

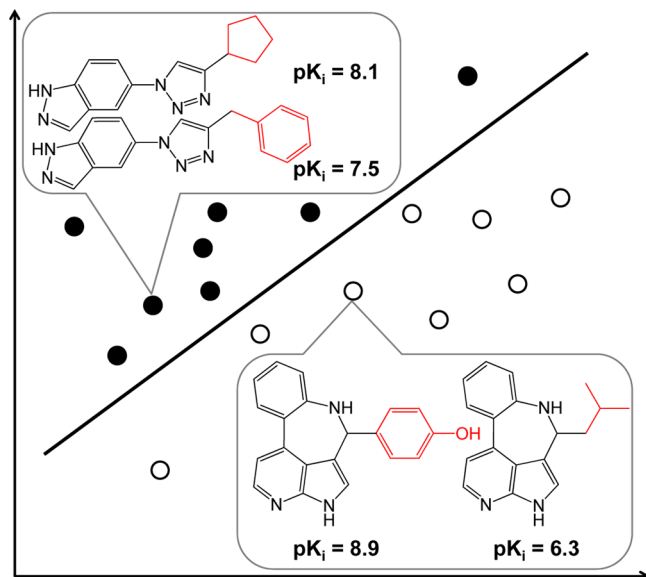
**3.1. Motivation and Strategy.** SVMs<sup>10</sup> are supervised machine learning algorithms for binary object classification and ranking. The prediction of activity cliffs requires the comparison of compound pairs instead of individual compounds, which presents an off-the-beaten path scenario for machine learning and classification methods. For SVM

Table 1. Data Sets<sup>a</sup>

target name	target code	no. of cpds	no. of MMPs	no. of MMP-cliffs	no. of MMP-nonCliffs	no. of other MMPs
factor Xa	fxa	2202	14493	1161	10108	3224
melanocortin receptor 4	mcr4	1159	13053	449	9618	2986
kappa opioid receptor	kor	1645	10104	649	7190	2265
thrombin	thr	2037	9585	1103	6390	2092
adenosine a3 receptor	aa3	1862	9575	681	6752	2142
calpain 2	cal2	121	1206	387	718	101
cathepsin b	catb	150	681	120	451	110
dipeptidyl peptidase 8	dpp8	44	602	141	421	40
janus kinase 2	jak2	58	366	109	186	71

<sup>a</sup>For each of the nine compound sets, the target name, a target code (abbreviation), the number of compounds (cpds), and the number of MMPs are reported. MMPs are divided into the number of compound pairs forming activity cliffs (MMP-cliffs), no activity cliffs (MMP-nonCliffs), and other MMPs, following the potency difference-based definition detailed in the Methods section. The data sets are sorted by decreasing numbers of MMPs.

modeling, kernel functions can be designed to account for specific relationships between objects and facilitate classification on the basis of these relationships. Our focus on the SVM approach for activity cliff prediction was largely motivated by the design of new kernel functions to facilitate comparisons of compound pairs, as illustrated in Figure 2. Our approach to facilitate compound pair-based predictions included, as a basis, the generation of training sets of MMP-cliffs and MMP-nonCliffs. In addition, an integral part of our approach was to attempt a systematic analysis of structural differences between



**Figure 2.** Activity cliff prediction using SVMs. The schematic figure illustrates the principal idea of SVM-based activity cliff prediction. In this case, the basic classification unit is a compound pair, different from standard compound classification tasks. Compound pairs forming MMP-cliffs (nonfilled circles) and MMP-nonCliffs (black circles) are separated by a hyperplane. In molecular graphs, transformation substructures are colored red. For each compound, its  $pK_i$  value is reported.

compounds in cliff and noncliff pairs. The underlying hypothesis was that there should be structural features among compounds sharing a specific activity that are responsible for high and low potency and thus, ultimately, for the formation of activity cliffs. Although this hypothesis was intuitive, its potential utility for activity cliff prediction remained to be evaluated. Methodologically, this was not a trivial task because it required, first, relating features of compounds forming pairs to each other and, second, comparing feature differences across pairs.

**3.2. SVM Theory in Brief.** SVMs make use of labeled training data that are mapped into a feature space to build a linear classification model. A set of  $n$  training objects  $\{\mathbf{x}_i, y_i\}$  ( $i = 1, \dots, n$ ) are represented by a feature vector  $\mathbf{x}_i \in \chi$  (e.g.,  $\mathbb{R}^d$ ) and an associated class label  $y_i \in \{-1, 1\}$  corresponding to the ‘negative’ and ‘positive’ classes, respectively. By solving a convex quadratic optimization problem, a hyperplane  $H$  is derived that best separates positive from negative training data (Figure 2). During training, the cost parameter  $C$  penalizes the misclassification of training data and achieves a balance between minimizing the training error and maximizing the generalization of the classification.

The hyperplane  $H$  is defined by the normal weight vector  $\mathbf{w}$  and the bias  $b$ , so that  $H = \{\mathbf{x} | \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}$ , where  $\langle \cdot, \cdot \rangle$  is a scalar product. Test data are mapped into the same feature space  $\chi$  and classified by the linear decision function  $f(\mathbf{x}) = \text{sgn}(\langle \mathbf{x}, \mathbf{w} \rangle + b)$ , i.e., depending on which side of the hyperplane they fall. In our calculations, the positive class consisted of the MMP-cliffs and the negative class of MMP-nonCliffs.

If the training data are not linearly separable in the feature space  $\chi$ , the so-called *Kernel trick*<sup>19</sup> can be applied to replace the scalar product  $\langle \cdot, \cdot \rangle$  by a kernel function  $K(\cdot, \cdot)$ . Kernel functions are used to calculate the scalar product of two feature vectors in a higher dimensional space  $\mathcal{H}$  without explicitly calculating the mapping  $\Phi: \chi \rightarrow \mathcal{H}$ . In the higher dimensional space  $\mathcal{H}$ , a linear separation of the training data might be feasible. Kernel functions are of the form  $K(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$ , where  $\mathbf{u}$  and  $\mathbf{v}$  are feature vector representations.

**3.3. Standard Kernel Functions.** The following four popular kernels are often used in SVM calculations:

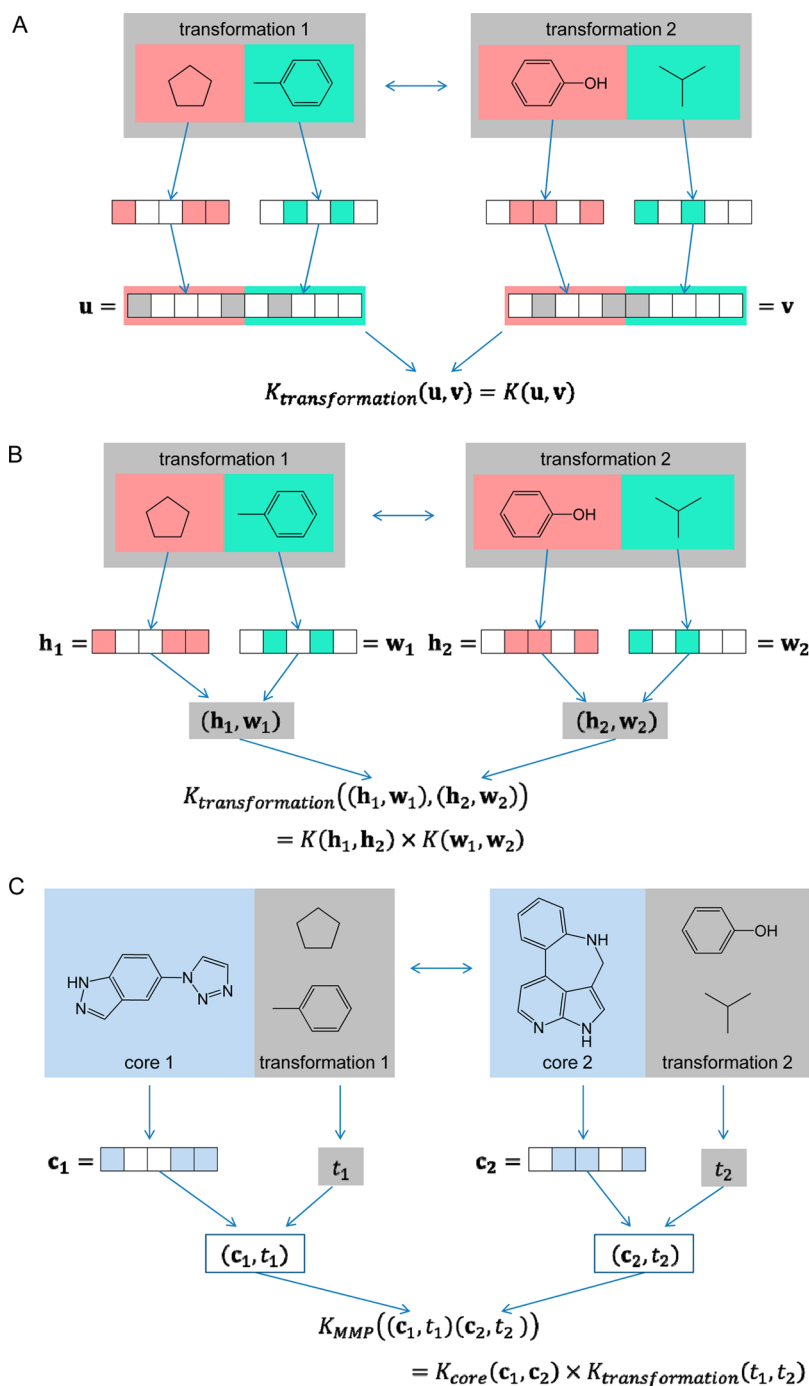
$$K_{\text{linear}}(\mathbf{u}, \mathbf{v}) = \langle \mathbf{u}, \mathbf{v} \rangle$$

$$K_{\text{Gaussian}}(\mathbf{u}, \mathbf{v}) = \exp(-\gamma \|\mathbf{u} - \mathbf{v}\|^2)$$

$$K_{\text{polynomial}}(\mathbf{u}, \mathbf{v}) = (\langle \mathbf{u}, \mathbf{v} \rangle + 1)^d$$

$$K_{\text{Tanimoto}}(\mathbf{u}, \mathbf{v}) = \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle + \langle \mathbf{v}, \mathbf{v} \rangle - \langle \mathbf{u}, \mathbf{v} \rangle}$$

The linear kernel corresponds to the standard scalar product. The Gaussian kernel is also known as the radial basis function kernel and depends on an adjustable parameter  $\gamma$ . In the polynomial kernel, the parameter  $d$  determines the degree of the polynomial function. The Tanimoto kernel<sup>20</sup> was introduced given the popularity of the Tanimoto coefficient for quantifying compound similarity. On the basis of these kernels, new kernel functions were designed specifically for activity cliff prediction, as described in the following.



**Figure 3.** Kernel functions. The design of new kernel functions for SVM-based activity cliff prediction is illustrated. (A) Substructure-difference kernel. A fingerprint representation is generated for each substructure representing a given transformation. Then, a difference vector is calculated for the two substructure fingerprints. In kernel calculations, difference vectors for different transformations are compared. (B) Substructure-pair kernel. Fingerprint representations of substructures representing a transformation are combined to yield substructure pairs. Kernel calculations then compare the substructure pairs of different transformations. (C) MMP kernel. A fingerprint representation is calculated for the common core of each MMP. The corresponding transformations are represented by a transformation object that is either the substructure-difference vector or the substructure-pair representation (according to A and B, respectively). The core structure vector and the transformation object are then combined for kernel calculations.

#### 4. DESIGN OF KERNEL FUNCTIONS FOR ACTIVITY CLIFF PREDICTION

**4.1. Substructure-Difference Kernel.** In order to use MMPs for SVM calculations, a feature vector representation of MMPs must be generated. As discussed above, MMPs consist of a common core structure and two differentiating

substructures (substituents) that constitute the molecular transformation. We first designed a kernel that utilized only the transformation to create a single feature vector. The substituents were classified according to the highly potent partner in the MMP, termed ‘highly potent substructure’, and the weakly potent MMP compound, referred to as ‘weakly potent substructure’ (for MMP-nonCliffs, these potency



differences were within an order of magnitude). For both substructures, a keyed fingerprint of size  $n$  was calculated. Then, a difference fingerprint of size  $2n$  was created that contained as the first  $n$  positions only those features present in the highly but not weakly potent substructure. If a feature was present in both substructures, the corresponding bit in the difference fingerprint was set off. The last  $n$  positions in the difference fingerprint contained features present only in the weakly but not the highly potent substructure. Accordingly, this difference vector uniquely described the transformation defined by the substructures on the basis of fingerprint features. The design of the difference fingerprint and substructure-difference kernel is illustrated in Figure 3A. Because this compound pair representation only comprises a single vector, kernel calculations can be performed as follows:

$$K_{\text{transformation}}(\mathbf{u}, \mathbf{v}) = K(\mathbf{u}, \mathbf{v})$$

where  $\mathbf{u}$  and  $\mathbf{v}$  are the substructure-difference vectors of two MMPs and  $K(\mathbf{u}, \mathbf{v})$  might, for example, be the Tanimoto or the polynomial kernel.

Because substructures were sorted on the basis of potency relationships between MMP compounds, this process is referred to as potency-based ordering. This information was taken into account during the learning and test phase. In addition, ordering of substructures by size (without considering potency relationship) was also investigated.

**4.2. Substructure-Pair Kernel.** Another new kernel represented a transformation as a pair of substructures. Again, substructures were classified according to the potency of the compounds from which they originate, and fingerprint representations were calculated. However, in this case, a transformation was represented as the pair  $(\mathbf{h}, \mathbf{w})$ , where  $\mathbf{h}$  is the fingerprint vector of the highly potent substructure and  $\mathbf{w}$  the feature vector of the weakly potent substructure. Given two substructure pairs  $(\mathbf{h}_i, \mathbf{w}_i)$  and  $(\mathbf{h}_j, \mathbf{w}_j)$ , a 'transformation kernel' was defined as the product of two separate kernels for the highly potent and weak potent substructures:

$$K_{\text{transformation}}((\mathbf{h}_i, \mathbf{w}_i), (\mathbf{h}_j, \mathbf{w}_j)) = K(\mathbf{h}_i, \mathbf{h}_j) \times K(\mathbf{w}_i, \mathbf{w}_j)$$

Thus, two independent kernels for highly potent and weakly potent substructures were combined to account for pairwise transformation similarities. The two kernels could again be implemented using standard kernel functions. The design of the substructure-pair kernel is illustrated in Figure 3B.

**4.3. MMP Kernel.** So far, we only considered molecular transformations to represent structural changes in MMPs that potentially lead to the formation of activity cliffs. However, the common core of an MMP might add further information for the classification of MMP-cliffs and MMP-nonCliffs because it defines the structural environment of a transformation. A potential caveat associated with considering the common core was that a given core structure might appear in both the positive and the negative class. This might be the case if a compound formed an MMP-cliff and an MMP-nonCliff with different partners. Hence, it was difficult to predict how the inclusion of the core might influence the classification calculations.

In order to generate a kernel that contains core information, an MMP was represented by combining the common core and the transformation, i.e.,  $(\mathbf{c}, t)$ , where  $\mathbf{c}$  is the feature vector representation from the common core and  $t$  is a transformation object that can either be described by the substructure-

difference vector or the substructure pair. Thus, the MMP kernel is defined by

$$K_{\text{MMP}}((\mathbf{c}_i, t_i)(\mathbf{c}_j, t_j)) = K_{\text{core}}(\mathbf{c}_i, \mathbf{c}_j) \times K_{\text{transformation}}(t_i, t_j)$$

The kernel function for pairs is again separated into independent kernels for each data type. The design of the MMP kernel is illustrated in Figure 3C. Because the common core was represented by a single feature vector, standard kernel functions could replace the core kernel (see above).

## 5. CALCULATION SETUP

**5.1. Cost Factor.** All SVM calculations were carried out using SVM<sup>light</sup>,<sup>21</sup> a freely available SVM implementation. With two exceptions, suggested default parameters of SVM<sup>light</sup> were used to render the calculations reproducible. Apart from adjustable parameters in kernel functions, as specified above, we only modified the cost factor for the treatment of unbalanced data sets. SVM calculations on significantly unbalanced data sets often result in the generation of a hyperplane that is proximal to under-represented training examples,<sup>17,18</sup> here the positive examples (MMP-cliffs). As a consequence, positive instances are often predicted at only low rates. The cost factor defines the ratio of training error costs on the positive class ( $C^+$ ) to penalties on the negative class ( $C^-$ ):<sup>22</sup>

$$\text{cost factor} = \frac{C^+}{C^-}$$

The default value of the cost factor is 1; i.e., the same penalty is applied to positive and negative examples. However, increasing the error cost  $C^+$ , i.e., the penalty to predict a false-negative, repositions the hyperplane farther away from the positive examples. An often recommended cost-factor adjustment<sup>17,21</sup> can be expressed as

$$\text{cost factor} = \frac{\text{NTE}}{\text{PTE}}$$

where NTE and PTE are the number of negative and positive training examples, respectively. Thus, the potential total cost of false negative errors and the potential total cost of false positive errors are the same.<sup>22</sup>

**5.2. Statistics.** We performed 10-fold cross-validation as a reasonable compromise between data perturbation and training data size.<sup>23</sup> The MMP-cliff and MMP-nonCliff classes were randomly partitioned into 10 samples such that the global ratio between positive and negative training examples was constant. Nine of 10 samples were utilized for SVM learning and model building including all positive and negative training examples, and the remaining sample was used as a test set for prediction. In systematic classification calculations, each sample was used once as a test set. Average statistics were calculated over all 10 trials and used for performance evaluation. The following statistics were calculated:

$$\text{accuracy} = \text{AC} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FN} + \text{TN} + \text{FP}}$$

$$\text{recall} = \text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{specificity} = \text{TNR} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

Table 2. Cost-Factor Settings<sup>a</sup>

target	cost factor = 1					cost factor = NTE/PTE				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	92.47	32.38	99.38	85.82	46.89	89.82	72.69	91.79	50.45	59.51
mcr4	97.01	35.62	99.88	94.17	51.34	95.90	78.15	96.72	53.02	63.01
kor	92.69	13.56	99.83	86.75	23.08	87.09	66.87	88.92	35.44	46.26
thr	91.69	52.07	98.53	85.93	64.64	88.90	81.15	90.23	58.99	68.29
aa3	92.47	25.84	99.19	76.81	38.51	86.95	71.95	88.46	38.63	50.19
cal2	95.65	95.11	95.95	92.83	93.86	96.10	97.44	95.40	92.14	94.65
catb	95.80	81.67	99.56	98.57	88.67	95.62	88.33	97.56	91.10	89.39
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.33	88.18	91.58	86.24	87.01	90.00	92.73	88.42	82.82	87.30

<sup>a</sup>For each data set, the average accuracy (AC), true positive rate (TPR), true negative rate (TNR), precision (P), and F score are reported (in %) for SVM calculations with different cost-factor settings. In these calculations, transformations were represented using the substructure-difference vector generated with MACCS, and the Tanimoto kernel was used as part of the substructure-difference kernel.

$$\text{precision} = P = \frac{TP}{TP + FP}$$

where TP, FN, TN, and FP define the number of predicted true positives, false negatives, true negatives, and false positives, respectively. TPR and TNR denote the true positive rate and true negative rate, respectively, and are used in the following to account for recall and specificity. Because so assessed prediction accuracy is not a very informative measure when the number of negative examples is much larger than the number of positive examples,<sup>24</sup> we also calculated the (balanced) F score that accounts for both precision and recall (and ranges from 0% to 100%):<sup>24</sup>

$$\text{F score} = 2 \times \frac{P \times TPR}{P + TPR}$$

**5.3. Fingerprints.** For representing substructures, two fragment-type fingerprints were used. The bonded-atom pair fingerprint (BAP)<sup>25</sup> encodes 117 different atom pairs with a focus on short-range connectivity information. In order to account for substructures comprising single atoms, we added three features describing carbon atoms, heteroatoms, and hydrogens resulting in a final fingerprint consisting of 120 structural descriptors. In addition, the MACCS<sup>26</sup> fingerprint was used that consists of 166 structural keys encoding substructures with one to 10 non-hydrogen atoms. An additional feature for a single hydrogen atom was also added in this case (because it might participate in a transformation). Furthermore, we evaluated the combination of both fingerprints (MACCS+BAP), resulting in a descriptor with 284 (117 + 166 + 1) features (two features added to the BAP fingerprint correspond to MACCS structural keys).

As a molecular representation of the common core, we used MACCS and Molprint2D.<sup>27</sup> The Molprint2D fingerprint requires the use of Sybyl atom types and encodes circular atom environments by fusing each atom in the structure with its neighboring atoms until a specific bond radius is reached. Here, we used features with a maximal bond radius of 2.

## 6. INITIAL TRIALS, COST-FACTOR ADJUSTMENT, AND SUBSTRUCTURE REPRESENTATION

**6.1. Basic Classification Performance.** To evaluate the potential of our SVM-based approach, we first determined the classification performance in calculations in which substructures were represented using the substructure-difference vector generated with MACCS, and the Tanimoto kernel was used

as part of the substructure-difference kernel. In addition, a constant cost factor of 1 was applied. The results of these calculations are reported in Table 2. To present comprehensive statistics for performance evaluation, we report for all cross-validated calculations the average accuracy (AC), true positive (TPR) and true negative (TNR) rates, the precision (P), and the F score. In the following discussion, most emphasis is put on TPR, P, and F score values. The results in Table 2 for a cost factor of 1 mirror overall successful activity cliff predictions, with notable compound class dependence. In particular, the (un)balance of positive and negative training examples affected the calculations. For the first five data sets in Table 2, which contained many more MMP-nonCliffs than MMP-cliffs, prediction accuracy was lower than for the remaining more balanced sets (i.e., cal2, catb, dpp8, and jak2), as to be expected (see above). MMP-nonCliffs were generally predicted with very high accuracy, leading to TNRs of nearly 100% in all but one (jak2; 91.6%) case. This also led to an overall accuracy of 90–100% of the calculations and to a precision of 76–100%. Significant differences were observed between the rates with which MMP-cliffs were correctly predicted. Here, TPRs ranged from 13.6% to 99.3%, leading to F scores between 23.1% and 99.6%. For the unbalanced data sets, TPRs and F scores ranged from 13.6% to 52.1% and 23.1% to 64.6%, respectively. By contrast, for the balanced data sets, TPRs and F scores of 81.7–99.3% and 87.0–99.6% were observed, respectively. Thus, the results of initial activity cliff predictions were considered encouraging, at least for balanced data sets, and we thus further refined the approach, as discussed in the following.

**6.2. Cost Factor.** We first attempted to address the low TPRs and resulting F scores observed for unbalanced data sets in our initial calculations. Therefore, the default cost factor of 1 was replaced by the adjusted cost factor = NTE/PTE, which introduced a higher penalty on misclassification of positive training instances, i.e., activity cliffs. We repeated cross-validated SVM calculations under these conditions and observed a significant increase in TPRs for unbalanced data sets, as reported in Table 2. For balanced data sets, classification performance remained essentially unchanged, but for unbalanced sets, TPRs and F scores further increased to 66.9–81.2% and 46.3–68.3%, respectively. A trade-off has been a reduction in precision because the TNRs were reduced from on average 99.4% to 91.2%, due to the adjusted cost factor. However, this relatively small reduction in TNRs was clearly overcompensated for by an average TPR increase of 42.3% for unbalanced sets, yielding reasonably to highly accurate activity cliff

Table 3. Comparison of Fingerprints for Substructure Representation<sup>a</sup>

target	BAP					MACCS					MACCS+BAP				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	83.72	72.18	85.04	35.66	47.73	89.82	72.69	91.79	50.45	59.51	90.11	73.21	92.06	51.50	60.39
mcr4	89.58	77.28	90.15	26.90	39.85	95.90	78.15	96.72	53.02	63.01	95.91	77.25	96.78	53.26	62.84
kor	79.36	65.33	80.63	23.37	34.39	87.09	66.87	88.92	35.44	46.26	87.31	67.95	89.05	36.19	47.15
thr	85.11	78.61	86.23	49.66	60.84	88.90	81.15	90.23	58.99	68.29	89.15	81.43	90.49	59.71	68.85
aa3	81.21	68.15	82.52	28.18	39.83	86.95	71.95	88.46	38.63	50.19	86.98	71.37	88.55	38.63	50.04
cal2	92.67	94.87	91.50	86.08	90.16	96.10	97.44	95.40	92.14	94.65	96.11	97.44	95.40	92.23	94.69
catb	93.17	90.83	93.79	80.24	84.91	95.62	88.33	97.56	91.10	89.39	96.15	88.33	98.22	93.38	90.47
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	82.20	80.82	83.02	74.53	76.97	90.00	92.73	88.42	82.82	87.30	89.67	90.00	89.47	83.74	86.56

<sup>a</sup>The performance of the BAP, MACCS, and MACCS+BAP fingerprints for substructure representation is compared. Calculation statistics are reported according to Table 2. The transformations were represented using the substructure-difference vector. The Tanimoto kernel was used as part of the substructure-difference kernel, and the adjusted cost factor was applied.

Table 4. Comparison of Standard Kernels<sup>a</sup>

target	Tanimoto					linear				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	89.82	72.69	91.79	50.45	59.51	80.70	71.66	81.74	31.14	43.38
mcr4	95.90	78.15	96.72	53.02	63.01	87.30	80.82	87.61	23.44	36.28
kor	87.09	66.87	88.92	35.44	46.26	79.93	67.49	81.06	24.39	35.82
thr	88.90	81.15	90.23	58.99	68.29	84.20	76.97	85.45	47.80	58.93
aa3	86.95	71.95	88.46	38.63	50.19	77.88	75.19	78.16	25.90	38.50
cal2	96.10	97.44	95.40	92.14	94.65	95.30	96.41	94.71	91.06	93.55
catb	95.62	88.33	97.56	91.10	89.39	94.92	85.83	97.33	90.46	87.63
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.00	92.73	88.42	82.82	87.30	92.33	95.45	90.53	85.94	90.24
target	Gaussian ( $\gamma = 1/\text{numFeatures}$ )					Gaussian ( $\gamma = 0.1$ )				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	84.83	71.40	86.37	37.65	49.25	91.19	64.77	94.22	56.37	60.20
mcr4	91.76	82.16	92.20	33.10	47.11	96.26	70.80	97.44	56.89	62.86
kor	83.31	69.19	84.59	28.92	40.77	87.86	62.72	90.13	36.66	46.17
thr	86.73	78.87	88.09	53.37	63.63	89.82	78.71	91.74	62.32	69.51
aa3	80.84	77.10	81.22	29.46	42.57	87.50	68.14	89.46	39.50	49.91
cal2	95.48	96.92	94.71	91.09	93.81	95.56	96.67	94.98	91.47	93.90
catb	94.92	85.83	97.33	90.46	87.63	95.98	85.00	98.89	95.92	89.48
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	92.67	95.45	91.05	86.54	90.60	90.33	90.91	90.00	84.50	87.40
target	polynomial ( $d = 2$ )					polynomial ( $d = 3$ )				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	90.67	70.45	93.00	53.74	60.88	92.12	65.20	95.21	61.04	62.97
mcr4	96.03	75.02	97.01	54.04	62.69	97.02	72.13	98.18	65.41	68.39
kor	87.84	65.50	89.86	37.03	47.20	89.35	59.95	92.00	40.60	48.30
thr	90.35	78.97	92.32	64.15	70.74	91.03	74.89	93.82	67.73	71.11
aa3	87.45	70.93	89.11	39.70	50.84	88.69	66.07	90.97	42.44	51.58
cal2	96.11	95.38	96.52	93.81	94.53	95.66	91.52	97.90	96.02	93.55
catb	95.97	82.50	99.56	98.57	89.19	95.62	80.00	99.78	99.17	88.04
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.27	88.09	91.58	86.46	86.77	88.93	82.64	92.63	87.07	84.42

<sup>a</sup>The performance of different kernels as part of the substructure-difference kernel is compared. Performance statistics are reported. The Gaussian kernel was used with two different  $\gamma$  values ( $\gamma = 1/\text{numFeatures}$  and  $\gamma = 0.1$ ) and the polynomial kernel with two different exponents  $d$  ( $d = 2$  and  $d = 3$ ). The parameter numFeatures describes the number of features present in the substructure-difference vector. The substructures were represented using the substructure-difference vector with MACCS, and the adjusted cost factor was applied.

predictions for all nine different data sets (Table 2). Accordingly, the adjusted cost factor was used in all subsequent calculations.

**6.3. Substructure Representation.** Next, we compared different fingerprint representations of transformation sub-

structures. Table 3 reports search results for the comparison of the BAP, MACCS, and (MACCS+BAP) fingerprints used for the generation of the substructure-difference vector. Calculations with the BAP substructure-difference vector resulted in consistently high TPRs but low precision for unbalanced data

Table 5. Comparison of Transformation Kernels<sup>a</sup>

target	substructure-difference vector					substructure pairs				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	89.82	72.69	91.79	50.45	59.51	91.01	74.85	92.87	54.68	63.17
mcr4	95.90	78.15	96.72	53.02	63.01	96.30	74.14	97.34	56.94	64.16
kor	87.09	66.87	88.92	35.44	46.26	88.00	60.55	90.47	36.73	45.61
thr	88.90	81.15	90.23	58.99	68.29	90.03	82.07	91.41	62.34	70.80
aa3	86.95	71.95	88.46	38.63	50.19	88.83	67.11	91.02	43.23	52.45
cal2	96.10	97.44	95.40	92.14	94.65	96.29	96.41	96.24	93.39	94.81
catb	95.62	88.33	97.56	91.10	89.39	95.44	86.67	97.78	91.63	88.85
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.00	92.73	88.42	82.82	87.30	91.20	94.45	89.33	84.48	88.91

<sup>a</sup>The performance of the substructure-difference vector is compared to the substructure-pair representation. The substructures were encoded using MACCS. The Tanimoto kernel was used as part of the two transformation kernels, and the adjusted cost factor was applied.

Table 6. Comparison of Transformation and MMP Kernels<sup>a</sup>

target	transformation kernel					MMP kernel (core structure: MACCS)					MMP kernel (core structure: Molprint2D)				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	89.82	72.69	91.79	50.45	59.51	94.11	81.30	95.58	67.91	73.95	94.68	82.17	96.11	70.92	76.03
mcr4	95.90	78.15	96.72	53.02	63.01	98.02	81.04	98.81	76.50	78.57	98.35	83.05	99.06	80.82	81.82
kor	87.09	66.87	88.92	35.44	46.26	93.21	72.57	95.08	57.38	63.99	94.86	72.58	96.87	67.88	70.04
thr	88.90	81.15	90.23	58.99	68.29	93.07	84.41	94.57	72.93	78.21	93.75	84.05	95.43	76.20	79.85
aa3	86.95	71.95	88.46	38.63	50.19	93.52	74.74	95.41	62.60	67.91	95.12	74.45	97.20	73.23	73.57
cal2	96.10	97.44	95.40	92.14	94.65	97.55	97.69	97.49	95.54	96.57	97.64	97.69	97.63	95.79	96.70
catb	95.62	88.33	97.56	91.10	89.39	96.85	90.00	98.67	95.24	92.30	97.02	90.83	98.67	95.43	92.76
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.00	92.73	88.42	82.82	87.30	90.67	91.82	90.00	84.74	87.89	91.00	91.82	90.53	85.55	88.28

<sup>a</sup>The performance of the transformation kernel is compared to the MMP kernel. The substructure-difference kernel was used as the transformation kernel. Substructures were represented using the substructure-difference vector with MACCS. The MACCS and Molprint2D fingerprints were compared as core structure representations in the MMP kernel. The Tanimoto kernel was used as part of the substructure-difference kernel as well as the MMP kernel, and the adjusted cost factor was applied.

sets, due to a reduction in TNRs, leading to low F scores. MACCS-based calculations yielded comparably high TPRs but higher precision and F scores. No further increases in these rates and scores were observed when the (MACCS+BAP) combination was used. Consequently, MACCS was used in subsequent calculations.

## 7. KERNEL COMPARISON

Kernel design for the treatment of compound pairs as a basic classification object has been a key aspect of our approach to activity cliff prediction. We first compared the performance of standard kernels that provided a basis for the generation of substructure and MMP kernels.

**7.1. Standard kernels.** Table 4 summarizes the results of SVM calculations using different kernel functions as a component of the substructure-difference kernel. The use of the Tanimoto kernel resulted in TPRs that were consistently above 66% for all targets. In these calculations, the precision was low for two data sets (kor and aa3). For unbalanced sets, F scores varied from 46.3% to 68.3%. By contrast, for balanced sets, F scores were consistently higher than 87%. The linear kernel essentially paralleled the results of the Tanimoto kernel for balanced data sets but displayed consistently lower precision for unbalanced sets. The use of the Gaussian kernel with small  $\gamma$  parameter values (0.0034–0.0065), depending on the number of features ( $\gamma = 1/\text{numFeatures}$ ) used in the calculations, also resulted in comparable TPRs but lower precision for unbalanced sets. For a larger  $\gamma$  value of 0.1, increased precision was observed but TPRs were reduced, yielding F scores

comparable to the Tanimoto kernel. Furthermore, the polynomial kernel (with  $d = 2$  and  $d = 3$ ) also produced rates and scores that were similar to those obtained for the Tanimoto kernel. Thus, taken together, differences in prediction performance for different standard kernels were by and large insignificant. Since the Tanimoto kernel was parameter-free, it was selected for further calculations.

**7.2. Transformation and MMP Kernels.** An interesting initial finding was that promising classification results were obtained using the Tanimoto substructure-difference kernel (see section 6.1). This kernel only accounted for differences between transformation substructures, rather than entire MMPs. We then compared the substructure-difference and substructure-pair kernels (based on the Tanimoto kernel). Calculation requirements for these kernels differed. The substructure-difference kernel only required one kernel calculation, but the difference vector must be precalculated. By contrast, for the substructure-pair kernel, no precalculations were required, but the kernel calculation must be carried out for two functions. The results of search calculations using these alternative transformation-only kernels are reported in Table 5. No clear preference for one or the other kernel was detectable. Overall, the substructure-pair kernel produced slightly lower TPRs but slightly higher precision than the substructure-difference kernel (except for the dpp8 set yielding  $P = 100\%$  in both instances), which resulted in similar F scores.

We then included the MMP kernel in the comparison, which was designed to combine the core structure representation of an MMP with its substructure-difference vector. For this



Table 7. Comparison of Potency- and Size-Based Substructure Ordering<sup>a</sup>

target	potency-based ordering					size-based ordering				
	AC	TPR	TNR	P	F score	AC	TPR	TNR	P	F score
fxa	89.82	72.69	91.79	50.45	59.51	87.89	68.38	90.13	44.30	53.73
mcr4	95.90	78.15	96.72	53.02	63.01	95.21	70.15	96.38	47.80	56.71
kor	87.09	66.87	88.92	35.44	46.26	84.17	62.10	86.16	28.81	39.33
thr	88.90	81.15	90.23	58.99	68.29	87.32	77.81	88.97	54.95	64.36
aa3	86.95	71.95	88.46	38.63	50.19	83.79	68.13	85.37	32.06	43.50
cal2	96.10	97.44	95.40	92.14	94.65	91.94	93.03	91.36	85.38	88.89
catb	95.62	88.33	97.56	91.10	89.39	93.17	85.00	95.33	84.27	83.77
dpp8	99.82	99.29	100.00	100.00	99.63	99.82	99.29	100.00	100.00	99.63
jak2	90.00	92.73	88.42	82.82	87.30	86.27	87.18	85.79	79.10	82.44

<sup>a</sup>The performance of potency-based ordering of transformation substructures is compared to size-based ordering. The transformations were represented using the substructure-difference vector. The Tanimoto kernel was used as part of the substructure-difference kernel, and the adjusted cost factor was applied.

purpose, the core structure can be represented using different fingerprints. For substructure representations, fragment fingerprints such as BAP or MACCS are in principle a preferred choice, but for core structure representation, other types of fingerprints might also be used. In Table 6, search results for the substructure-difference kernel are compared to those obtained for two versions of the MMP kernel including one in which the common core was represented using MACCS and another that utilized Molprint2D instead (i.e., a topological atom environment fingerprint). We found that application of the MMP kernel further improved classification performance. For both versions of the MMP kernel, an increase in TPRs and precision was observed compared to the transformation kernel, leading to higher F scores. For the MMP kernel, TPRs were very similar for MACCS and Molprint2D, but F scores were slightly higher for Molprint2D, due to a minor increase in TNRs. On average, F scores were 82.1% for the MACCS- and 84.3% for the Molprint2D-based MMP kernel. Compared to the substructure-difference kernel, which yielded TPRs and F scores of 66.9–99.3% and 46.3–99.6%, respectively, the Molprint2D-based MMP kernel produced TPRs and F scores of 72.6–99.3% and 70.0–99.6%, respectively. Improvements were observed for balanced and unbalanced data sets but were of larger magnitude for the latter. On average, TPRs slightly increased from 83.2% (transformation kernel) to 86.2% (MMP kernel) and F scores (reflecting both recall and precision) from 73.1% to 84.3%. Thus, the incorporation of core structure contributions of MMP-cliffs and MMP-nonCliffs into the kernel function further increased the accuracy of activity cliff predictions.

We also investigated the influence of substructure ordering on the calculations. In the classification scheme underlying our analysis, the ordering of transformation substructures in MMPs was potency-based. As a control, we also evaluated size-based ordering of substructures. The results for the substructure difference are presented in Table 7. With the exception of one set (dpp8; with consistently 100% precision), both TPRs and F scores decreased for size-based ordering. Comparable trends were observed when the MMP kernel was used (data not shown). Hence, potency-based ordering of substructures was generally preferred, but size-based ordering also yielded accurate predictions.

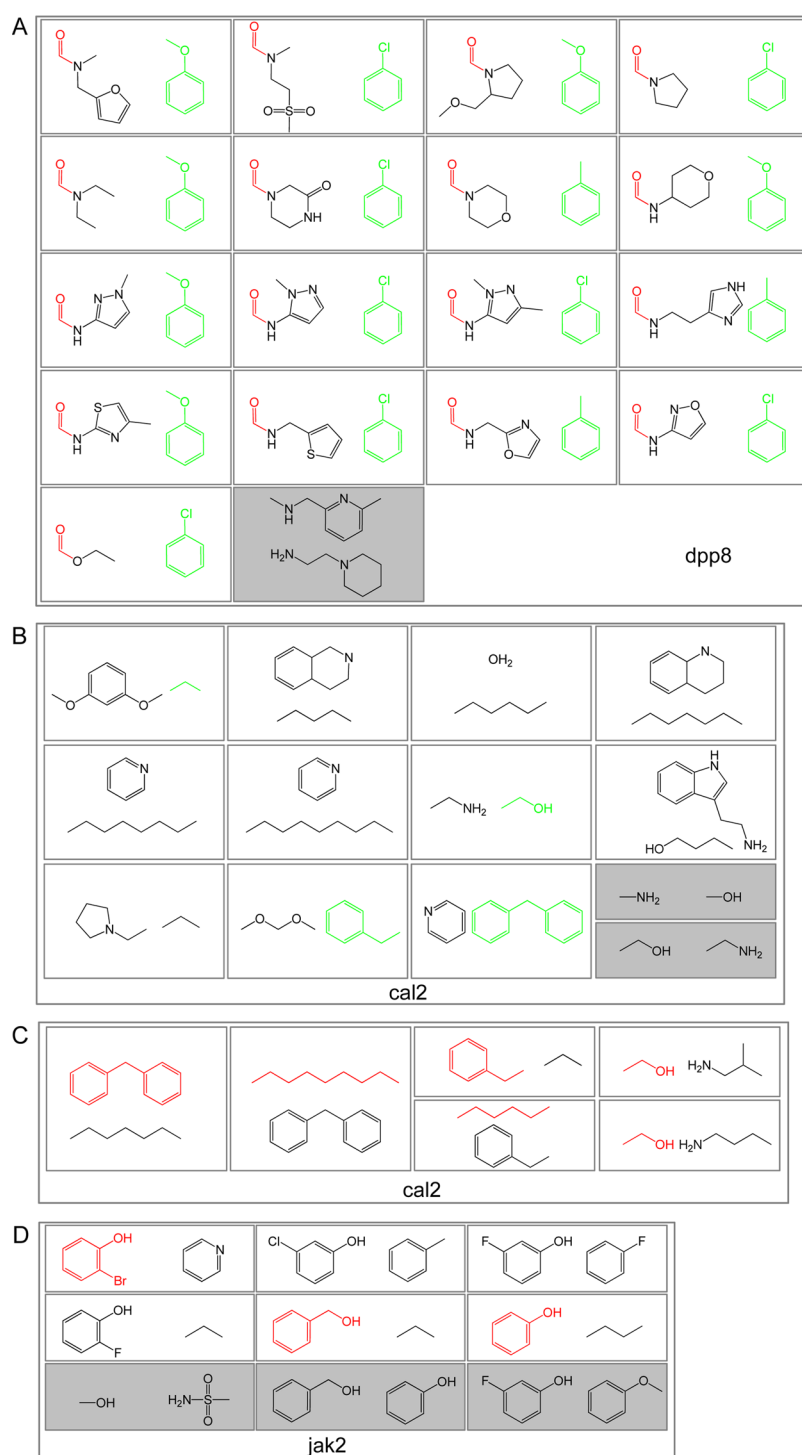
## 8. STRUCTURAL PATTERNS

Given the results of our calculations, we also investigated whether successful predictions of MMP-cliffs might be

rationalized in structural terms. Therefore, we analyzed correctly identified MMP-cliffs and MMP-nonCliffs for the presence of characteristic transformations and structural features. In a number of instances, structural patterns were identified that could be attributed to activity cliff formation. In the following, representative examples are discussed.

Figure 4A shows a number of transformations leading to the formation of MMP-cliffs in the dpp8 set. Most MMP-cliffs were characterized by transformations in which a substructure containing a carbonyl group was replaced by a substituted phenyl group. The carbonyl group was predominantly involved in the formation of amide bonds, but there were also ketone and ether linkages proximal to the carbonyl group. With one exception, all of these MMP-cliffs were correctly classified. The only exception was a structurally very different transformation observed in an MMP-cliff (shown on a gray background), which might present an interesting test case for further analysis of activity cliffs among dpp8 inhibitors. Apart from this exception, the typical MMP-cliff transformation patterns observed in the dpp8 set resulted in perfect classifications, independent of chosen kernel functions and SVM calculations settings.

Figure 4B shows transformations of MMP-cliffs from the cal2 set that were correctly identified or misclassified as MMP-nonCliffs (gray background). The weakly potent substructures of correctly classified MMP-cliffs were linear alkyl chains, oxygen containing (alkyl) substituents, or groups containing phenyl rings. Replacement of these substructures with nitrogen- or oxygen-containing substructures or substituted ring systems caused a strong increase in potency, leading to the formation of activity cliffs. The transformation of misclassified pairs exclusively consisted of small nitrogen- and oxygen-containing substructures. These examples illustrate that clearly defined structural signatures of activity cliffs were not always obvious in the cal2 set, making this case a difficult classification problem. Despite this structural variability, 64.2% of all MMP-cliffs in the cal2 set displayed similar structural patterns and were correctly classified. Figure 4C shows examples of MMP-nonCliff transformations, which further illustrate the presence of complex transformation–potency relationships. In these cases, more potent compounds in pairs contained substructures that were found in weakly potent MMP-cliff compounds shown in Figure 4B. Consequently, these pairs were correctly classified as MMP-nonCliffs. Nevertheless, calculations on the cal2 set yielded accurate predictions using our SVM models, with a TPR and F score of 97.4% and 94.7%.



**Figure 4.** Exemplary MMP transformations. Exemplary MMP-cliff and MMP-nonCliff transformations are shown for different compound data sets (according to Table 1). Highly potent substructures are positioned on the left or at the top of cells and weakly potent ones on the right or at the bottom. Correctly classified MMP-cliff transformations are shown on a white background and misclassified transformations on a gray background. Selected structural features/patterns are colored red (in highly potent substructures) and green (in weakly potent substructures). (A) Data set dpp8/MMP-cliffs, (B) cal2/MMP-cliffs, (C) cal2/MMP-nonCliffs, (D) jak2/MMP-cliffs.

Figure 4D shows MMP-cliff transformations from the jak2 set. Here, replacements of alkyl or phenyl groups with benzyl alcohol or halogen substituted benzyl alcohol groups often led to activity cliffs that were correctly detected. Examples of transformations in misclassified MMP-cliffs included the exchange of methoxy and sulfonamide groups and of pairs of

different benzyl alcohol derivatives that departed from the prevalent substructure patterns among MMP-cliffs.

## 9. CONCLUDING REMARKS

Herein, we have presented a first approach to predict activity cliffs in compound data sets. From a conceptual and methodological point of view, the prediction of activity cliffs

represents a nontrivial task. The underlying assumption is that structural differences in pairs of similar compounds can be directly related to potency differences and then compared across pairs representing cliffs and noncliffs. We have approached the task of activity cliff prediction using SVM modeling because the SVM formalism provides the opportunity to design kernel functions specifically tailored towards this task. To represent activity cliffs and noncliffs, the concept of MMP-cliffs is applied that yields a structurally well-defined representation of activity cliffs on the basis of common core structures and distinguishing substructure transformations. Another general difficulty in activity cliff prediction is the assembly of compound data sets with a balanced composition of positive (cliffs) and negative (noncliff) training examples, which typically is an important prerequisite for effective machine learning. Because activity cliffs are relatively rare among bioactive compounds, data sets are generally unbalanced. We have systematically searched for compound data sets that contained MMP-cliffs at a relatively high frequency and determined all MMP-cliffs and MMP-nonCliffs in these sets. A total of nine data sets were obtained for our analysis that contained significant numbers of MMP-cliffs. However, all of these data sets contained many more MMP-nonCliffs than MMP-cliffs, as expected. For the purpose of our analysis, we considered data sets balanced if the MMP-nonCliff/MMP-cliff ratio was not larger than 4, which was the case for four of our sets. The remaining five sets were characterized by much larger ratios and hence considered unbalanced. However, using newly introduced kernel functions, activity cliffs were predicted with reasonable to high accuracy on the basis of SVM learning and classification. During learning, unbalanced training example distributions were effectively handled by adjusting the cost factor of the SVM calculations. We designed alternative kernel functions that only took transformation substructure differences or transformation and core structure features into account. Interestingly, overall accurate predictions were already obtained when transformation kernels were applied, but prediction accuracy was further improved through the use of MMP kernel functions that considered transformation and core differences. However, much structural information relevant for the formation of activity cliffs was often encoded by transformations, without a critical influence of their specific structural environment. In our analysis, best predictions were obtained when cross-validated SVM calculations with adjusted cost factors were carried out using the Tanimoto kernel-based MMP kernel with a MACCS substructure-difference vector (on the basis of potency-based ordering of substructures) and a Molprint2D representation of common MMP cores. Under these conditions, average true positive rates and F scores of 86.2% and 84.3%, respectively, were achieved in activity cliff predictions, with an average precision of 82.9% and accuracy of 95.8% of the calculations. In many instances, it was possible to rationalize successful predictions of activity cliffs on the basis of structural features of corresponding transformations. Taken together, given the results presented herein, we anticipate that the SVM-based approach to activity cliff prediction should be of considerable interest in the search for cliffs in large compound data sets.

## AUTHOR INFORMATION

### Corresponding Author

\*Tel.: +49-228-2699-306. Fax: +49-228-2699-341. E-mail: [bajorath@bit.uni-bonn.de](mailto:bajorath@bit.uni-bonn.de).

## Author Contributions

<sup>§</sup>The contributions of these authors should be considered equal.

## Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

X.H. is supported by the *China Scholarship Council*. The authors thank Martin Vogt for helpful discussions.

## REFERENCES

- (1) Stumpfe, D.; Bajorath, J. Exploring Activity Cliffs in Medicinal Chemistry. *J. Med. Chem.* **2012**, *55*, 2932–2942.
- (2) Maggiora, G. M. On Outliers and Activity Cliffs – Why QSAR often Disappoints. *J. Chem. Inf. Model.* **2006**, *46*, 1535–1535.
- (3) Wassermann, A. M.; Wawer, M.; Bajorath, J. Activity Landscape Representations for Structure-Activity Relationship Analysis. *J. Med. Chem.* **2010**, *53*, 8209–8223.
- (4) Guha, R.; Van Drie, J. H. Structure-Activity Landscape Index: Identifying and Quantifying Activity Cliffs. *J. Chem. Inf. Model.* **2008**, *48*, 646–658.
- (5) Medina-Franco, J. L.; Martínez-Mayorga, K.; Bender, A.; Marín, R. M.; Giulianotti, M. A.; Pinilla, C.; Houghten, R. A. Characterization of Activity Landscapes Using 2D and 3D Similarity Methods: Consensus Activity Cliffs. *J. Chem. Inf. Model.* **2009**, *49*, 477–491.
- (6) Bajorath, J.; Peltason, L.; Wawer, M.; Guha, R.; Lajiness, M. S.; Van Drie, J. H. Navigating Structure-Activity Landscapes. *Drug Discovery Today* **2009**, *14*, 698–705.
- (7) Wassermann, A. M.; Dimova, D.; Bajorath, J. Comprehensive Analysis of Single- and Multi-Target Activity Cliffs Formed by Currently Available Bioactive Compounds. *Chem. Biol. Drug Des.* **2011**, *78*, 224–228.
- (8) Vogt, M.; Huang, Y.; Bajorath, J. From Activity Cliffs to Activity Ridges: Informative Data Structures for SAR Analysis. *J. Chem. Inf. Model.* **2011**, *51*, 1848–1856.
- (9) Namasivayam, V.; Bajorath, J. Searching for Coordinated Activity Cliffs Using Particle Swarm Optimization. *J. Chem. Inf. Model.* **2012**, *52*, 927–934.
- (10) Vapnik, V. N. *The Nature of Statistical Learning Theory*, 2nd ed.; Springer: New York, 2000.
- (11) Kenny, P. W.; Sadowski, J. Structure Modification in Chemical Databases. In *Chemoinformatics in Drug Discovery*; Oprea, T. I., Ed.; Wiley-VCH: Weinheim, Germany, 2004; pp 271–285.
- (12) Hu, X.; Hu, Y.; Vogt, M.; Stumpfe, D.; Bajorath, J. MMP-Cliffs: Systematic Identification of Activity Cliffs on the Basis of Matched Molecular Pairs. *J. Chem. Inf. Model.* **2012**, *52*, 1138–1145.
- (13) Hussain, J.; Rea, C. Computationally Efficient Algorithm to Identify Matched Molecular Pairs (MMPs) in Large Data Sets. *J. Chem. Inf. Model.* **2010**, *50*, 339–348.
- (14) Liu, T.; Lin, Y.; Wen, X.; Jorissen, R. N.; Gilson, M. K. BindingDB: a Web-Accessible Database of Experimentally Determined Protein-Ligand Binding Affinities. *Nucleic Acids Res.* **2007**, *35*, D198–D201.
- (15) BindingDB. <http://www.bindingdb.org/> (accessed February 8, 2012).
- (16) Clark, M.; Cramer, R. D., III; Van Opdenbosch, N. Validation of the General Purpose Tripos 5.2 Force Field. *J. Comput. Chem.* **1989**, *10*, 982–1012.
- (17) Akbani, R.; Kwek, S.; Japkowicz, N. Applying Support Vector Machines to Imbalanced Datasets. In *Machine Learning: ECML 2004*; Boulicaut, J.-F., Esposito, F., Giannotti, F., Pedreschi, D., Eds.; Springer: Berlin/Heidelberg, 2004; pp 39–50.
- (18) Tang, Y.; Zhang, Y.-Q.; Chawla, N. V.; Krasser, S. SVMs Modeling for Highly Imbalanced Classification. *IEEE Trans. Syst. Man. Cybern. B: Cybern.* **2009**, *39*, 281–288.
- (19) Boser, B. E.; Guyon, I. M.; Vapnik, V. N. A Training Algorithm for Optimal Margin Classifiers. In *Proceedings of the 5th Annual*

*Workshop on Computational Learning Theory*; Pittsburgh, PA, 1992; ACM: New York, 1992; pp 144–152.

(20) Ralaivola, L.; Swamidass, S. J.; Saigo, H.; Baldi, P. Graph Kernels for Chemical Informatics. *Neural Networks* **2005**, *18*, 1093–1110.

(21) Joachims, T. Making Large-Scale SVM Learning Practical. In *Advances in Kernel Methods – Support Vector Learning*; Schölkopf, B., Burges, C. J. C., Smola, A. J., Eds.; MIT-Press: Cambridge, MA, 1999; pp 169–184.

(22) Morik, K.; Brockhausen, P.; Joachims, T. Combining Statistical Learning with a Knowledge-based Approach - A Case Study in Intensive Care Monitoring. In *Proceedings of the 16th International Conference on Machine Learning (ICML-99)*; Morgan Kaufmann: Burlington, MA, 1999.

(23) McLachlan, G. J.; Do, K.-A.; Ambrose, C. *Analyzing Microarray Gene Expression Data*; Wiley & Sons: Hoboken, NJ, 2004.

(24) Kubat, M.; Matwin, S. Addressing the Curse of Imbalanced Training Sets: One-Sided Selection. In *Proceedings of the Fourteenth International Conference on Machine Learning*, San Francisco, 1997; Morgan Kaufmann: Burlington, MA, 1997; pp 179–186.

(25) Ahmed, H. E. A.; Vogt, M.; Bajorath, J. Design and Evaluation of Bonded Atom Pair Descriptors. *J. Chem. Inf. Model.* **2010**, *50*, 487–499.

(26) *MACCS Structural Keys*; Symyx Software: San Ramon, CA, 2002.

(27) Bender, A.; Mussa, H. Y.; Glen, R. C. Molecular Similarity Searching Using Atom Environments, Information-Based Feature Selection, and a Naïve Bayesian Classifier. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 170–178.