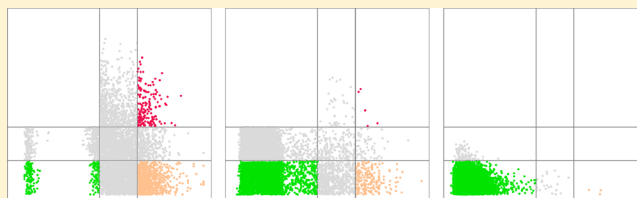


Large-Scale Assessment of Activity Landscape Feature Probabilities of Bioactive Compounds

Shilva Kayastha, Dilyana Dimova, Preeti Iyer, Martin Vogt, and Jürgen Bajorath*

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

ABSTRACT: Activity landscape representations integrate pairwise compound similarity and potency relationships and provide direct access to characteristic structure–activity relationship features in compound data sets. Because pairwise compound comparisons provide the foundation of activity landscape design, the assessment of specific landscape features such as activity cliffs has generally been confined to the level of compound pairs. A conditional probability-based approach has been applied herein to assign most probable activity landscape features to individual compounds. For example, for a given data set compound, it was determined if it would preferentially engage in the formation of activity cliffs or other landscape features. In a large-scale effort, we have determined conditional activity landscape feature probabilities for more than 160,000 compounds with well-defined activity annotations contained in 427 different target-based data sets. These landscape feature probabilities provide a detailed view of how different activity landscape features are distributed over currently available bioactive compounds.



■ INTRODUCTION

Activity landscapes can be rationalized as representations of molecular similarity and potency relationships in sets of active compounds.¹ Modeling of activity landscapes supports SAR analysis of large data sets, helps to visualize SAR patterns, and aids in the identification of compounds that are rich in SAR information.¹ A characteristic feature of activity landscapes is that they generally rely on pairwise compound similarity and potency comparisons. Hence, activity landscape features are typically resolved at the level of compound pairs. Prominent landscape features include *activity cliffs*,² *similarity cliffs*,³ and *smooth pairs*.³ Activity cliffs are pairs of structurally similar compounds with a large difference in potency, similarity cliffs represent pairs of structurally different compounds with similar potency, and smooth pairs are compounds having similar structure and similar potency. Activity cliffs are typically associated with high SAR information content² and similarity cliffs provide a basis for the identification of structurally diverse compounds having similar activity.³ In addition, activity landscapes contain *featureless pairs*, i.e., compounds that are structurally dissimilar and have different potency.

Activity landscape modeling has thus far been limited to compound pair-based analysis. Accordingly, it has not been possible to study and predict activity landscape features for individual compounds. To address this limitation, we have recently introduced a methodology to assign the probability of activity landscape features to single compounds.⁴ This has been accomplished by deriving conditional landscape feature probabilities for individual compounds on the basis of similarity and potency relationships. These probabilities indicate which activity landscape features, or feature combinations, are most likely for each individual compound in data sets of moderate size (100–500 compounds).⁴ For example, a given compound

might have a high probability to engage in activity cliffs or, alternatively, smooth pairs. This methodology represents the first approach to study activity landscape features for single compounds, and comparable or related methods are currently not available. The theoretical framework of the methodology has been described in detail in our previous study, and exemplary applications have been presented.⁴ However, the approach has not yet been applied to derive activity landscape features for compounds on a large scale and analyze their distribution. Therefore, in this study, we have systematically assessed the distribution of activity landscape feature probabilities across currently available bioactive compounds with well-defined activity annotations.

■ MATERIALS AND METHODS

Compound Data Sets. Data sets were assembled from ChEMBL (version 15).⁵ Only compounds with explicit equilibrium constants (K_i values) and half-maximal inhibitory concentrations (IC_{50} values) for human targets at the highest confidence level (confidence score 9)⁵ were considered. Moreover, only exact potency annotations (with relationship type “=”) were further analyzed. For compounds with multiple potency annotations against the same target, the geometric mean of all recorded potency values was taken if they fell within the same order of magnitude. Otherwise, the compound was not further considered. In addition, a qualifying data set was required to contain at least 100 compounds. On the basis of these selection criteria, 127 K_i - and 300 IC_{50} -based data sets were obtained for different targets containing a total of 60,379

Received: November 19, 2013

and 104,434 compounds, respectively. The largest K_i and IC_{50} data sets contained 2307 and 2915 compounds, respectively.

Activity Landscape Feature Categories. On the basis of pairwise similarity and potency relationships, compound pairs can be assigned to four major activity landscape feature categories: *activity cliffs* (AC), i.e., pairs of structurally similar compounds having a significant potency difference; *similarity cliffs* (SC), structurally diverse compounds with similar potency; *smooth pairs* (SP), structurally similar compounds having similar potency; and *featureless pairs* (FL), structurally diverse compounds with a large potency difference.

Landscape features defined on the basis of compound pairs can be intuitively rationalized using structure–activity similarity (SAS) maps,⁶ as illustrated in Figure 1A. SAS maps are simple 2D plots that relate structure and activity similarity of compound pairs to each other.

Structural Similarity and Potency Difference Thresholds. Structural similarity was assessed by calculation of the Tanimoto coefficient (T_c)⁷ using the extended connectivity fingerprint with bond diameter four (ECFP4).⁸ In the SAS map in Figure 1A, compound pairs reaching a $T_c \geq 0.55$ were considered similar (similarity threshold), consistent with widely applied activity cliff definitions.² As a potency difference threshold, a difference in potency between two compounds forming an activity cliff of at least 2 orders of magnitude was required. Activity cliffs of this magnitude are statistically significant.⁹

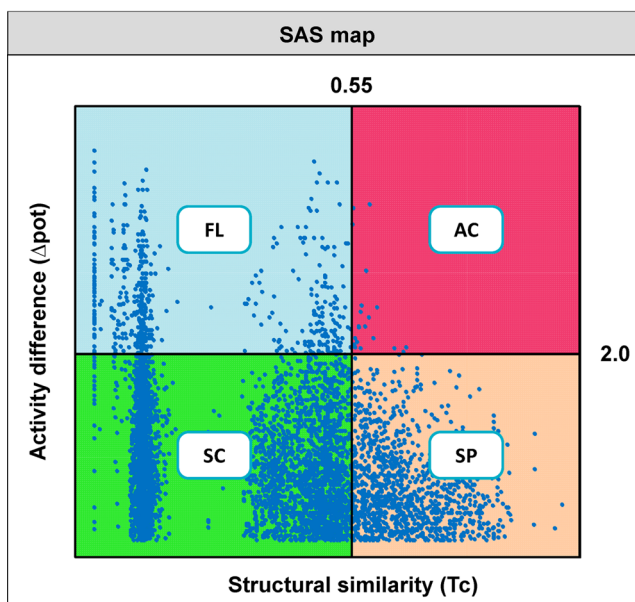
Local SAS Maps. SAS maps can also be calculated for individual compounds, giving rise to local SAS maps,⁴ as illustrated in Figure 1B. A compound-based local SAS map is obtained by only plotting all compound pairs formed by a specific compound. On the basis of these maps, the propensity of an individual compound to form activity cliffs, similarity cliffs, smooth pairs, or featureless pairs can be studied.

Fuzzy Thresholds. To avoid threshold boundary effects in feature assignments, fuzzy thresholds for chemical similarity and potency difference have been introduced.⁴ Fuzzy boundaries for structural similarity were set applying the T_c interval 0.45–0.65 and for activity similarity applying a potency difference interval of 1–2 pK_i or pIC_{50} units, as shown in Figure 2. The three SAR-relevant feature regions in a local SAS map are designated as follows:

- (i) Region \tilde{R}_{00} (lower left): similarity cliffs
- (ii) Region \tilde{R}_{10} (lower right): smooth pairs
- (iii) Region \tilde{R}_{11} (upper right): activity cliffs

Because featureless pairs yield only very little SAR information, we focused on activity cliffs, smooth pairs, and similarity cliffs as the principal landscape features. For compound pairs falling within the boundary intervals, a joint weighted membership function was generated to partially assign them to neighboring landscape regions. Due to its foundation in fuzzy set theory (a detailed discussion is provided in ref 4), this approach makes it possible to use the original partitioning scheme of local SAS maps while softening the thresholds between different regions. For compound pairs m_{ij} with varying potency differences Δpot_{ij} (either ΔpK_i or ΔpIC_{50}), the weighted membership with respect to the interval of 1–2 units is defined by membership functions $\mu_{\Delta potL}$, $\mu_{\Delta potH}$ as follows:

A



B

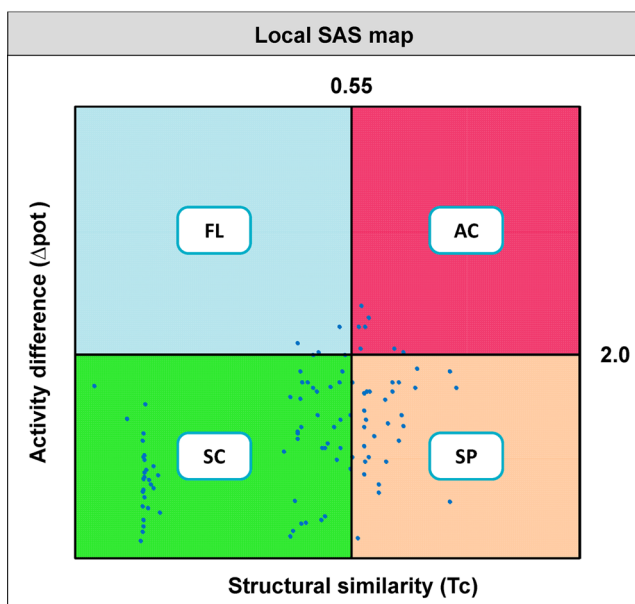


Figure 1. Structure–activity similarity (SAS) maps. (A) An SAS map is shown for an exemplary data set representing all pairwise potency and structural relationships. Structural similarity was assessed by calculation of Tanimoto similarity using the ECFP4 fingerprint. Four different regions (AC: activity cliffs; SP: smooth pairs; SC: similarity cliffs, and FL: SAR noninformative compound pairs) are distinguished on the basis of a similarity threshold ($T_c \geq 0.55$) and potency difference threshold ($\Delta pot \geq 2$). (B) A local SAS map is shown that contains all pairs formed by a given compound. The global and local SAS maps shown here are schematic maps for a model data set.

$$\mu_{\Delta potL}(m_{ij}) = \begin{cases} 0 & \text{if } 2 \leq \Delta pot_{ij} \\ 2 - \Delta pot_{ij} & \text{if } 1 < \Delta pot_{ij} < 2 \\ 1 & \text{if } \Delta pot_{ij} \leq 1 \end{cases}$$

and

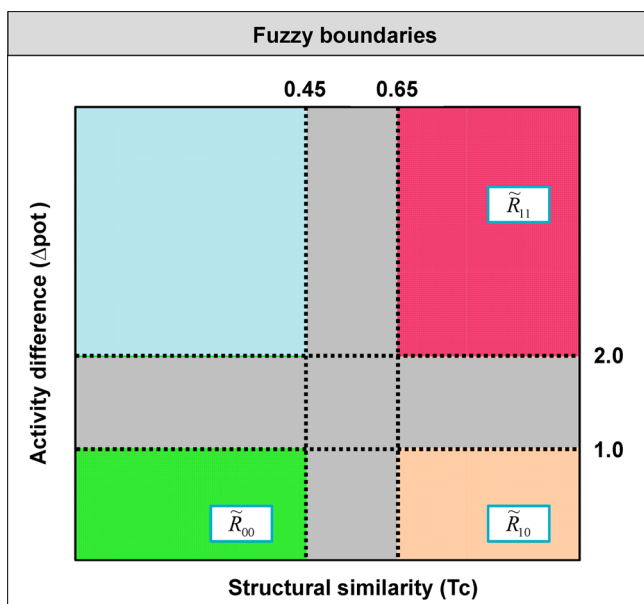


Figure 2. Fuzzy boundaries. The figure illustrates the use of fuzzy threshold value intervals (gray) instead of crisp boundaries.

$$\mu_{\Delta\text{potH}}(m_{ij}) = \begin{cases} 0 & \text{if } \Delta\text{pot}_{ij} \leq 1 \\ \Delta\text{pot}_{ij} - 1 & \text{if } 1 < \Delta\text{pot}_{ij} < 2 \\ 1 & \text{if } 2 \leq \Delta\text{pot}_{ij} \end{cases}$$

Analogously, for compound pairs with varying similarity values Tc_{ij} , the weighted membership with respect to the Tc boundary interval between 0.45 and 0.65 is determined by membership functions μ_{TcL} , μ_{TcH} as follows:

$$\mu_{TcL}(m_{ij}) = \begin{cases} 0 & \text{if } 0.65 \leq Tc_{ij} \\ (0.65 - Tc_{ij})/0.2 & \text{if } 0.45 < Tc_{ij} < 0.65 \\ 1 & \text{if } Tc_{ij} \leq 0.45 \end{cases}$$

and

$$\mu_{TcH}(m_{ij}) = \begin{cases} 0 & \text{if } Tc_{ij} \leq 0.45 \\ (Tc_{ij} - 0.45)/0.2 & \text{if } 0.45 < Tc_{ij} < 0.65 \\ 1 & \text{if } 0.65 \leq Tc_{ij} \end{cases}$$

The weighting scheme for compound pairs within the potency difference or structural similarity boundary intervals was designed to generate corresponding values within the range [0,1] and defines partial memberships to neighboring regions. The combined membership functions for the SAR-relevant regions \tilde{R}_{00} , \tilde{R}_{10} , and \tilde{R}_{11} can be defined as

$$\mu_{\tilde{R}_{00}}(m_{ij}) = \mu_{TcL}(m_{ij})\mu_{\Delta\text{potL}}(m_{ij})$$

$$\mu_{\tilde{R}_{10}}(m_{ij}) = \mu_{TcH}(m_{ij})\mu_{\Delta\text{potL}}(m_{ij})$$

$$\mu_{\tilde{R}_{11}}(m_{ij}) = \mu_{TcH}(m_{ij})\mu_{\Delta\text{potH}}(m_{ij})$$

Feature probabilities. Feature probabilities for individual compounds are calculated using the frequencies with which these compounds participate in the formation of various activity

landscape features. For a given compound i , the fuzzy frequencies $\tilde{R}_{00}(i)$, $\tilde{R}_{10}(i)$, and $\tilde{R}_{11}(i)$ for similarity cliffs, smooth pairs, and activity cliffs, respectively, are estimated from the combined membership functions:

$$|\tilde{R}_{00}(i)| = \sum_{j \neq i} \mu_{\tilde{R}_{00}}(m_{ij}), \quad |\tilde{R}_{10}(i)| = \sum_{j \neq i} \mu_{\tilde{R}_{10}}(m_{ij}), \quad \text{and} \\ |\tilde{R}_{11}(i)| = \sum_{j \neq i} \mu_{\tilde{R}_{11}}(m_{ij})$$

These frequencies can be used to calculate the probabilities of a given compound to form similarity cliffs, smooth pairs, or activity cliffs.

Conditional Feature Probabilities. To determine the feature probabilities for a given compound, its participation in all pairwise relationships is considered. Moreover, it is statistically meaningful to calculate the propensity of a compound to, for example, form activity cliffs with respect to the subset of structurally similar compound pairs, because dissimilar pairs cannot form activity cliffs. Likewise, the probability of a compound to form similarity cliffs is best estimated with respect to the subset of compound pairs that have similar potency. Furthermore, the ability of a compound to participate in the formation of smooth pairs can be calculated either using the subset of structurally similar compound pairs or the subset of compound pairs having similar potency. These conditional relationships can be transformed into conditional probabilities. Thus, for a given compound i , conditional probabilities are calculated using the subset of pairs that satisfy a conditional relationship as follows:

1. If compound i has potency similar to another, then
 - a. the probability of forming a similarity cliff is given by

$$P(TcL|\Delta\text{potL}) = \frac{|\tilde{R}_{00}(i)|}{|\tilde{R}_{00}(i)| + |\tilde{R}_{10}(i)|}$$

- b. the probability of forming a smooth pair is given by

$$P(TcH|\Delta\text{potL}) = \frac{|\tilde{R}_{10}(i)|}{|\tilde{R}_{00}(i)| + |\tilde{R}_{10}(i)|}$$

2. If compound i is structurally similar to another, then
 - a. the probability to form a smooth pair is calculated as

$$P(\Delta\text{potL}|TcH) = \frac{|\tilde{R}_{10}(i)|}{|\tilde{R}_{10}(i)| + |\tilde{R}_{11}(i)|}$$

- b. the probability to form an activity cliff is calculated as

$$P(\Delta\text{potH}|TcH) = \frac{|\tilde{R}_{11}(i)|}{|\tilde{R}_{10}(i)| + |\tilde{R}_{11}(i)|}$$

Accordingly, the frequencies of fuzzy landscape features were calculated and used to determine the conditional probabilities of fuzzy landscape features. Eight SAR-informative feature categories were defined (Cat 1–Cat 8) on the basis of the calculated conditional probabilities and their combinations, as reported in Table 1.

Table 1. Activity Landscape Feature Probability Categories^a

No.	Type	Category
0	single-feature category	SAR noninformative
1		similarity cliffs likely
2		smooth pairs likely/similarity cliffs unlikely
3		smooth pairs likely/activity cliffs unlikely
4	combined-feature category	activity cliffs likely
5		similarity cliffs likely/activity cliffs unlikely
6		similarity cliffs likely/activity cliffs likely
7		similarity cliffs unlikely/activity cliffs unlikely
8		similarity cliffs unlikely/activity cliffs likely

^aThe table summarizes the different activity landscape feature probability feature categories surveyed in this study.

Feature Probability Thresholds and Feature Categories. Previously,⁴ conditional feature probabilities were calculated for 139 target sets assembled from BindingDB.¹⁰

From these, corresponding thresholds at the 90th percentile of the sorted values were determined to categorize compounds.⁴ These thresholds were applied to test the conditional probability-based approach on a few exemplary data sets. For the systematic assessment of activity landscape features across bioactive compounds reported herein, conditional feature probabilities were calculated for all compounds in the 127 K_i and 300 IC₅₀ ChEMBL⁵ target sets and sorted. Conditional probabilities for which the denominator was less than 2 were not further considered. For all target sets, four global conditional feature probability thresholds were determined. The following values were obtained: similarity cliffs likely, 0.998 (K_i) and 1.0 (IC₅₀); smooth pairs likely/similarity cliffs unlikely, 0.109 (K_i) and 0.131 (IC₅₀); smooth pairs likely/activity cliffs unlikely, 1.0 (K_i and IC₅₀); and activity cliffs likely, 0.377 (K_i) and 0.348 (IC₅₀). Compounds were only assigned to a category if their conditional probability met or exceeded the corresponding threshold. As reported in Table 1, these categories correspond to Cat 1–Cat 4 and are *single-feature*

Table 2. Data Sets with the Largest Proportions of Compounds per Feature Category^a

Cat	K _i					IC ₅₀				
	% Cpds (cat)	ChEMBL Target ID	ChEMBL Target Name	% Cpds	# Cpds	% Cpds (cat)	ChEMBL Target ID	ChEMBL Target Name	% Cpds	# Cpds
1	9.5	3729	carbonic anhydrase IV	44.8	290	12.2	1615387	nuclear receptor coactivator 1	80	155
		2885	carbonic anhydrase III	39.5	124		1741221	cysteine protease ATG4B	73.8	325
		4789	carbonic anhydrase VA	37.3	249		4018	neuropeptide Y receptor type 2	73.3	146
2	8.6	2434	interleukin-8 receptor B	89.4	104	8.5	4633	voltage-gated potassium channel subunit Kv1.3	87.7	130
		2001	purinergic receptor P2Y12	85.5	509		3468	caspase-7	79.9	164
		4005	PI3-kinase p110-alpha subunit	82.9	105		5966	IgG receptor FcRn large subunit p51	76.2	172
3	10.2	3969	carbonic anhydrase VB	29.4	177	13.1	3775	dual specificity phosphatase Cdc25A	55.8	104
		332	matrix metalloproteinase-1	26.8	164		4804	dual specificity phosphatase Cdc25B	49	100
		238	dopamine transporter	25.6	870		4191	monoglyceride lipase	47.2	108
4	7.3	333	matrix metalloproteinase-2	21.9	128	6.8	2973	Rho-associated protein kinase 2	40.2	102
		321	matrix metalloproteinase 9	18.4	174		4282	serine/threonine-protein kinase AKT	19.4	505
		332	matrix metalloproteinase-1	16.5	164		204	thrombin	19.4	656
5	0.1	251	adenosine A2a receptor	0.7	2307	-	-	-	-	-
		264	histamine H3 receptor	0.6	1981		-	-	-	-
		226	adenosine A1 receptor	0.6	2030		-	-	-	-
6	0.4	2265	acyl coenzyme A: cholesterol acyltransferase	1.5	137	0.03	4801	caspase-1	1	199
		1790	vasopressin V2 receptor	1.4	146		4625	apoptosis regulator Bcl-X	0.9	107
		205	carbonic anhydrase II	1.3	1519		3746	11-beta-hydroxysteroid dehydrogenase 2	0.8	122
7	0.5	208	progesterone receptor	44.7	114	0.7	278	integrin alpha-4	26.9	108
		2611	furin	40.9	127		4198	inhibitor of apoptosis protein 3	19.6	168
		206	estrogen receptor alpha	17.1	111		4804	dual specificity phosphatase Cdc25B	19	100
8	0.9	2611	furin	25.2	127	0.8	2778	ileal bile acid transporter	21.9	114
		2243	anandamide amidohydrolase	24.8	101		3045	protein kinase C beta	20	110
		4617	phenylethanolamine N-methyltransferase	18.1	144		1921	vasopressin V1b receptor	18.8	133
0	62.6	259	melanocortin receptor 4	85.2	1289	58.0	5071	G protein-coupled receptor 44	84.6	637
		344	melanin-concentrating hormone receptor 1	83.3	872		1985	glucagon receptor	84.1	277
		2014	nociceptin receptor	82.9	642		4015	C–C chemokine receptor type 2	83.4	969

^aFor each activity landscape feature category, the top three ChEMBL target sets are reported with the proportions of compounds (% Cpds) assigned to the category. In addition, the percentage of compounds in all sets (% Cpds (cat)) falling into each category is given.

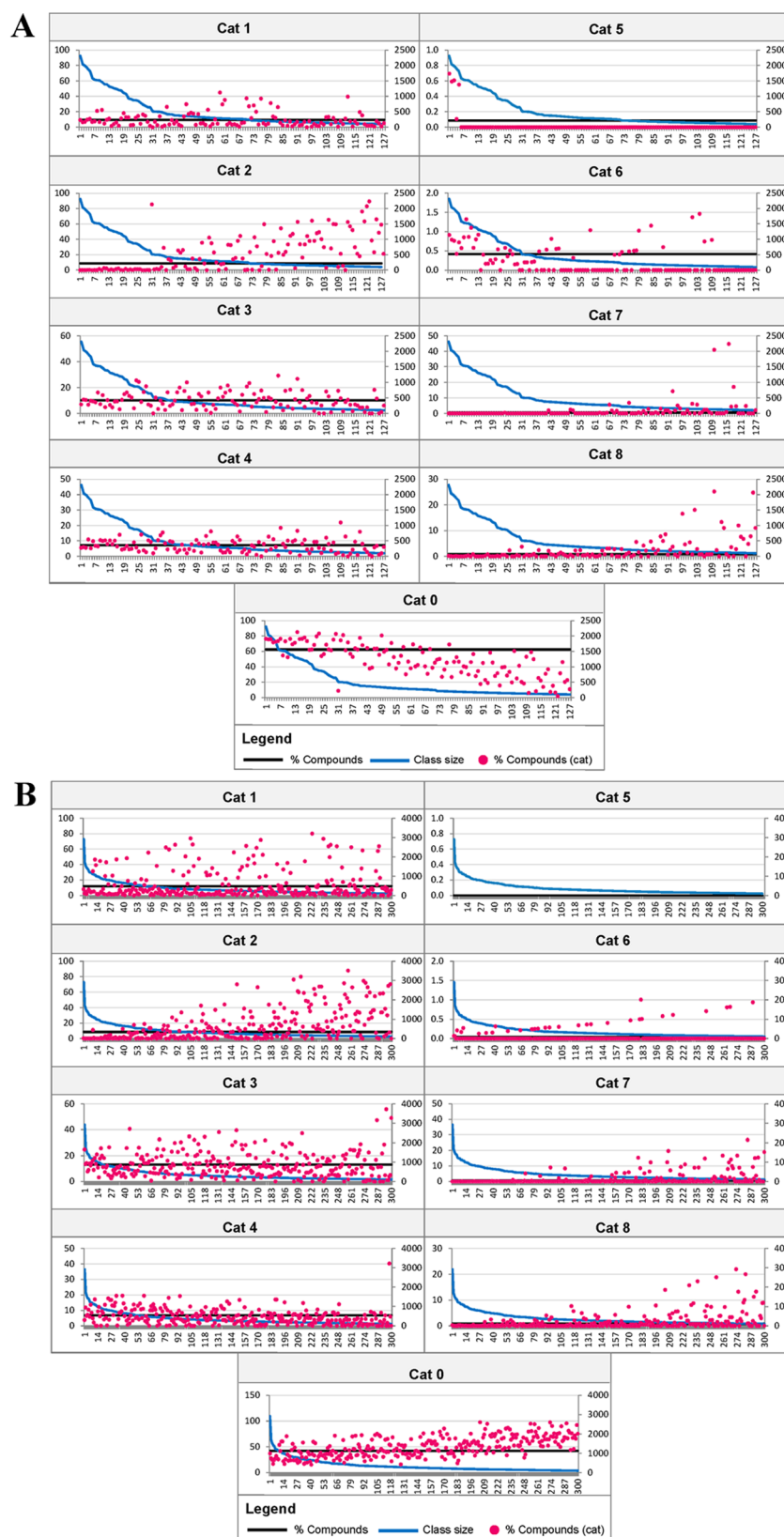


Figure 3. Distribution of activity landscape feature probability categories. For each category, the percentage of compounds is reported for all 127 K, (A) and 300 IC₅₀ (B) data sets. Plots are shown for each feature probability category that report the percentage of compounds (magenta) belonging to this category relative to the data set size (blue). On the x-axis, compound data sets are consecutively numbered. On the y-axis, the percentages (left) and total compound numbers (right) are reported. In addition, the percentage of compounds in all sets assigned to a given category is marked in black.

categories. In addition, four other categories were derived for compounds that exceeded more than one threshold. These categories correspond to Cat 5–Cat 8 in Table 1 and are termed *combined-feature categories*. If compounds met or exceeded more than one threshold, they were assigned to the corresponding combined-feature category. Compounds that did not reach any thresholds were categorized as SAR *non-informative* (Cat 0 in Table 1). A formal derivation of feature thresholds and feature categories is reported in the original publication of the conditional probability approach and its theoretical foundations.⁴ Because the 90th percentile of all sorted values was used as the criterion for the determination of global threshold values, for each conditional probability, 10% of all compounds met or exceeded this value. However, for individual compound sets, this percentage might vary significantly based on the intraset distribution of structurally similar compounds and compounds having similar potency.

RESULTS AND DISCUSSION

Study Goals. In this study, we have systematically assessed activity landscape feature probabilities of bioactive compounds and analyzed the distribution of feature categories across different targets. To these ends, we have selected more than 160,000 compounds with well-defined activity annotations belonging to 427 different target-based data sets and determined the propensity of each compound to engage in the formation of activity cliffs, similarity cliffs, or smooth pairs. Theoretically expected and observed feature probability distributions were compared in order to better understand how activity landscape features were distributed across different target-based compound sets. On the basis of feature probability calculations, we also aimed to distinguish compounds associated with interpretable SAR from SAR noninformative ones. In addition, we have attempted to identify target-based sets that are particularly rich in specific activity landscape features and globally determine landscape features that are formed with high probability by individual compounds.

SAR-Informative and SAR Noninformative Compounds. The assessment of SAR information associated with active compounds is one of the central themes in medicinal chemistry. Yet, the question of what represents valuable SAR information is difficult to address. An SAR-informative compound can be regarded as a compound that is predominantly involved in the formation of one informative activity landscape feature such as activity cliffs, smooth pairs, or similarity cliffs. Consequently, in order to account for SAR information, both the structure and the activity of a given compound must be considered. Clearly, the assessment of SAR information is dependent on the choice of a molecular representation (e.g., different fingerprints) as well as the similarity measure (e.g., the Tanimoto coefficient). In a recent study, the fingerprint descriptor dependence of SAR information has been systematically analyzed and quantified.¹¹ It has been demonstrated that the use of different fingerprint representations might change the nature of an SAR-informative compound and potentially transform such a compound into an SAR noninformative one.¹¹

Expected Distribution of Activity Landscape Feature Categories. On the basis of the derived feature probability thresholds, as discussed above, 10% of all compounds were expected to exceed the threshold of each of the single-feature categories 1–4. Some compounds exceeded two thresholds, which could be expected for 1% of all compounds assuming

threshold independence. As detailed in the Materials and Methods section, compounds exceeding two thresholds were assigned to only one of the combined-feature categories 5 to 8 and not to the corresponding single-feature categories. Therefore, ultimate assignments to single-feature categories 1–4 were expected to be below 10% of the compounds, depending on the probability distribution of combined-feature categories. For example, a compound in category 8 according to Table 1 would not also be assigned to categories 2 and 4 although it exceeded the relevant thresholds for these single-feature categories. These methodological considerations were relevant for rationalizing the observed activity landscape feature probability distributions.

Single-Feature vs Combined-Feature Categories. The percentages of data set compounds assigned to all categories are reported in Table 2 (% cpds (cat)). Values for single-feature categories (Cat 1–Cat 4) ranged from 7.3% (K_i) and 6.8% (IC_{50}) to 10.2% (K_i) and 13.1% (IC_{50}). Thus, these results considerably departed from the theoretically expected values discussed above, for two reasons. First, compounds did not enter conditional probability calculations if there were no structurally similar neighbors in a given data set, which was frequently observed, or if there were only very few data set compounds with similar potency (which was rarely observed). Precisely, for 8642 (14.3%) and 24,763 (23.7%) of the compounds with K_i and IC_{50} annotations, respectively, too small numbers of structural neighbors were available to qualify them for landscape feature probability calculations, whereas for only six (0.009%) and 14 (0.013%) compounds with K_i and IC_{50} annotations, respectively, too small numbers of compounds with similar potency were available to determine conditional feature probabilities. Second, in a number of cases, the 90th percentile yielded thresholds of 1.0. For example, this was the case for Cat 1 (similarity cliffs likely) in the case of IC_{50} data and for Cat 3 (activity cliff unlikely) for both IC_{50} and K_i data. The thresholds for Cat 1 and Cat 3 were 1.0 because more than 10% of all compounds did not form any smooth pairs with compounds having similar potency (Cat 1), and more than 10% of all compounds did not form any activity cliffs with structurally similar ones (Cat 3). For instance, in IC_{50} data sets, 24,763 of 104,434 compounds did not have sufficient numbers of structural neighbors for conditional probability calculations. Of the remaining 79,671 compounds, 14,373 compounds assigned to Cat 3, Cat 5, and Cat 7 did not form any activity cliffs with structurally similar compounds (despite the given conditions of fuzzy threshold boundaries and partial category memberships). In general, larger deviations from the theoretically expected values were observed for IC_{50} data sets compared to K_i data sets.

The observed percentages of compounds falling into the combined-feature categories Cat 5–Cat 8 only ranged from 0.1% (K_i) and 0% (IC_{50}) to 0.9% (K_i) and 0.8% (IC_{50}). Thus, observed values for categories 5–8 departed from the expected values of 1% (representing only an approximate estimate). As further discussed below, exceeding two thresholds was principally difficult for compounds in many instances.

Class-Size Dependence. The observed frequencies for each category and each class are reported in Figure 3. Data sets were sorted in the order of decreasing size. Categories 2, 7, and 8 (for both K_i and IC_{50}) and 6 (for IC_{50}) showed a clear dependence on the size of the data set. As can be seen, data sets with highest frequency of compounds assigned to a given feature category were almost exclusively of relatively small size.

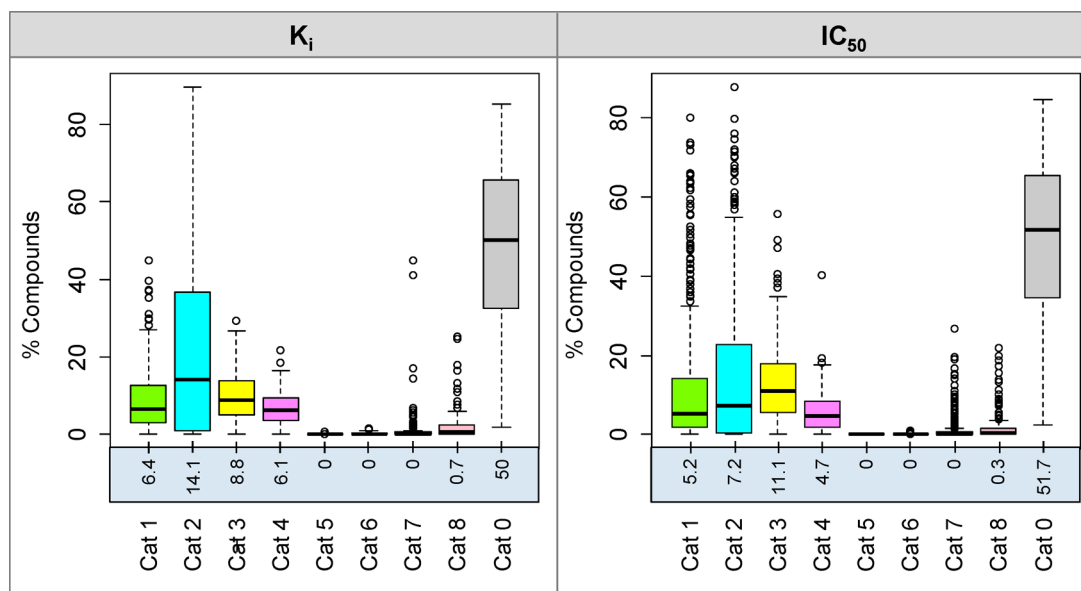


Figure 4. Compound assignment. Box plots are shown that monitor the percentage of compounds assigned to each activity landscape feature category. Furthermore, the median values for each category are reported. Whiskers extend to the most extreme data set within the 1.5 interquartile range of the box limits. Small circles outside the whiskers represent individual data sets.

For example, the interleukin-8 receptor B antagonist K_i data set contained 104 compounds and the percentage of compounds falling into category 2 was 89.4%, which was more than 10 times higher than the average for this category (8.6% over K_i sets). Also, the average percentage of compounds assigned to category 6 in IC_{50} sets was only 0.03%. However, 1% of the 199 compounds in the caspase-1 inhibitor IC_{50} data set fell into this category. It should also be noted that compounds falling into category 2, 7, and 8 were unlikely to form similarity cliffs. Data sets containing such compounds at a high frequency were often found to consist of only a few series of similar analogs.

Criteria for combined-feature categories were in part difficult to meet. For example, compounds assigned to Cat 5 were required to exceed the thresholds for Cat 1 (similarity cliffs likely) and for Cat 3 (smooth pairs likely/activity cliffs unlikely). Our results showed that Cat 5 is extremely rare because only 0.1% of all K_i and no IC_{50} compounds fell into this category. Principal reasons for these very low rates can be rationalized on the basis of a local SAS map (Figure 1B). For a compound to belong to Cat 5 there need to be only few compounds in the top right region (activity cliffs) compared to the bottom right (smooth pairs) and only few compounds in the bottom right (smooth pairs) compared to the bottom left (similarity cliffs). However, for IC_{50} calculations, the determined threshold for similarity cliffs was 1.0. Consequently, for a compound to be assigned to a “similarity cliff likely” category (Cat 1), it was not permitted to form any smooth pair. However, in this case, the compound could not possibly belong to the “activity cliff unlikely” category (Cat 3) and, hence, could principally not be assigned to Cat 5. As shown in Figure 3, compounds assigned to Cat 5 were almost exclusively observed in large K_i data sets.

Hence, data set sizes influenced the frequency of activity landscape feature probabilities. In small data sets, specifically, distributed compound similarity and potency relationships led to prevalence of individual single-feature categories. By contrast, the overall rarely observed combined-feature categories

were mostly detected in large data sets that contained at least a few qualifying compounds.

Distribution of Activity Landscape Feature Categories. Central questions underlying our analysis have been how activity landscape feature categories are distributed over target-based sets of bioactive compounds and which proportions of compounds are assigned to specific categories on the basis of conditional feature probabilities. In Figure 4, the observed frequency of all activity landscape feature categories is reported in box plot representations. For single-feature categories, the largest variation across different data sets was observed for Cat 2 (smooth pairs likely/similarity cliffs unlikely) and to a lesser extent for Cat 1 (similarity cliffs likely). The frequencies of Cat 1 and Cat 2 compounds in a data set were related to its composition. Structurally homogeneous and structurally diverse data sets were expected to contain most Cat 2 and Cat 1 compounds, respectively. In these cases, variations were generally high in different data sets. The smallest variation was observed for Cat 4 (activity cliffs likely). In this case, percentages varied around 10% in almost all data sets and there was no notable dependence on the size of the data set.

Cat 5 and Cat 6 were extremely rare due to the high threshold for similarity cliffs. Interestingly, a number of data sets for Cat 7 (similarity cliff unlikely/activity cliff unlikely) and Cat 8 (similarity cliff unlikely/activity cliff likely) exceeded the upper bound of the interquartile ranges for Cat 3 (activity cliff unlikely) and Cat 4 (activity cliff likely). These sets were structurally homogeneous and showed either small (Cat 7) or large potency variations (for Cat 8).

Data Sets with Prominent Activity Landscape Features. In Table 2, data sets with highest proportions of compounds per landscape feature category are reported. Matrix metalloproteinase inhibitor sets belonging to the peptidase M10A family represented the K_i data sets with largest percentages (16.5%–21.9%) of compounds likely to form activity cliffs (Cat 4). The size of these relatively small sets ranged from 128 to 174 compounds. On the other hand, IC_{50} data sets rich in compounds likely to form activity cliffs

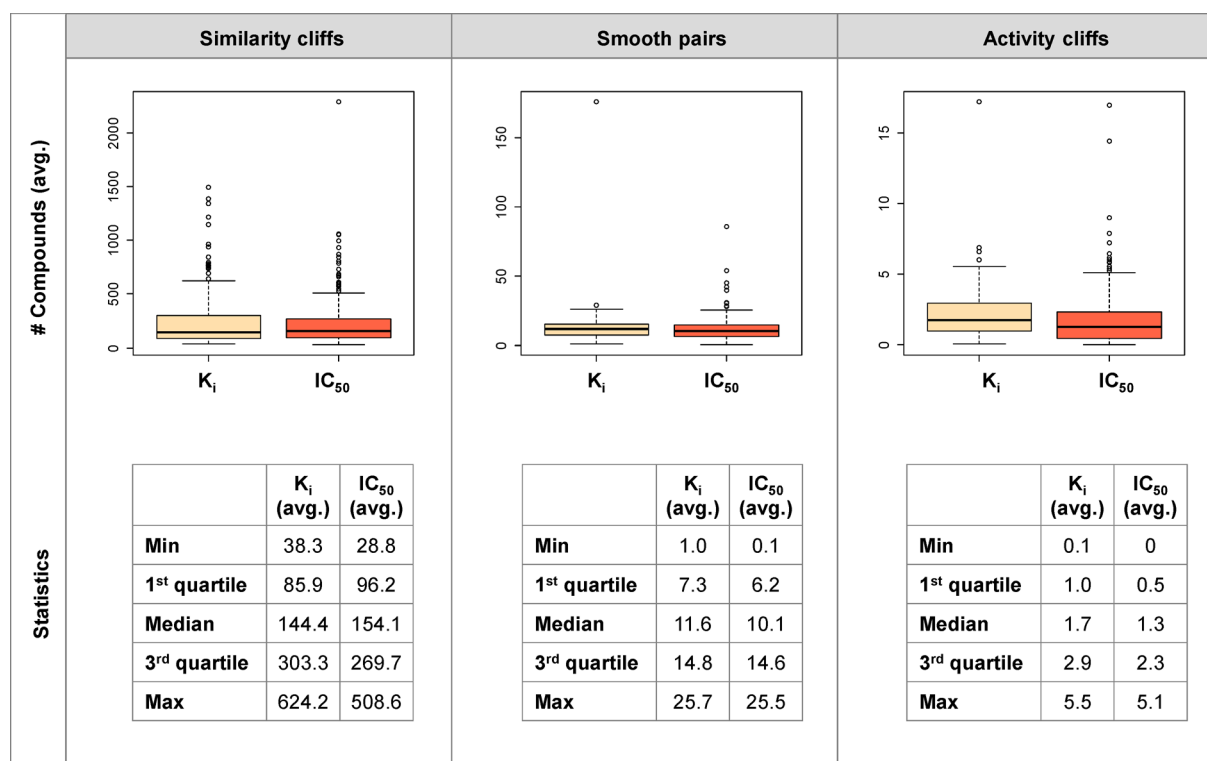


Figure 5. Distribution of feature frequencies. For each compound in a data set, the number of partners forming similarity cliffs, smooth pairs, and activity cliffs with this compound was determined. Box plots show the distribution of target set averages (# compounds (avg.)). Whiskers extend to the most extreme data set within the 1.5 interquartile range of the box limits. Small circles outside the whiskers represent individual data sets. In addition, the statistics (Min, 1st quartile, Median, 3rd quartile, and Max) for each box plot are reported in a table inset.

(percentages ranged here from 19.4% to 40.2%) were members of the kinase and peptidase S1 protein families and were generally larger in size (the number of compounds ranged from 102 to 656). Furthermore, members of the carbonic anhydrase family were among the K_i data sets with the largest proportions of compounds likely to form similarity cliffs (Cat 1). In these data sets, the percentages of compounds falling into Cat 1 ranged from 37.3% to 44.8%, thereby exceeding the average for this category (9.5% over all K_i sets) by a factor of ~ 4.3 . In IC_{50} data sets, the average for Cat 1 was 12.2%, which was exceeded by top-ranked sets by a factor of ~ 6 . In general, data sets with highest percentages of compounds belonging to different activity landscape feature categories were relatively small in size. One notable exception was combined-feature category 5 (similarity cliffs likely/activity cliffs unlikely). As detailed above, no IC_{50} compound was assigned to this category. However, K_i set compounds falling into this category belonged to the largest K_i data sets (with 1981–2307 compounds). It should also be emphasized that approximately 50% of all data set compounds did not reach any threshold value and were thus assigned to the SAR noninformative compound category (Cat 0). Hence, about half of the data set compounds were not likely to yield significant SAR information.

Global Distribution of Activity Landscape Features.

The propensities of individual compounds to form different activity landscape features, i.e. activity cliffs, smooth pairs, or similarity cliffs, have been systematically determined as conditional probabilities. The ability of a given compound to form similarity cliffs or smooth pairs was determined using the subset of partners with similar potencies. Likewise, the participation of this compound in the formation of activity

cliffs or smooth pairs was calculated by exclusively focusing on the subset of its structural neighbors. As explained above, the frequencies of each category were mainly determined by the selection of the 90th percentile as threshold values for the conditional probabilities. On this basis, significantly increased probabilities to form specific landscape features relative to expected probabilities were identified for compounds in K_i and IC_{50} data sets. This normalization procedure compensated for the large differences in the absolute values of landscape features observed previously.³ We also determined the average number of characteristic landscape features formed on a per-compound basis. In Figure 5, the average numbers of similarity cliffs, smooth pairs, and activity cliffs formed by each compound is reported in a box plot format for the K_i and IC_{50} data sets, respectively. The averages depended to a large extent on the size of the data set and varied accordingly. The median average number of similarity cliffs per compound observed was 144.4 for K_i and 154.1 for IC_{50} data sets. For smooth pairs, the median was approximately an order of magnitude smaller, with 11.6 for K_i and 10.1 for IC_{50} data sets. For activity cliffs, the median was approximately 2 orders of magnitude smaller, with 1.7 for K_i and 1.3 for IC_{50} data sets. The interquartile ranges varied accordingly by approximately a factor of 2 around the median. Taken together, these observations were consistent with a previously observed distribution of activity landscape features at the level of compound pairs.³ When assessed at the level of individual compounds, similarity cliffs were the most probable SAR-informative activity landscape feature and activity cliffs the most improbable.

Potential Utility of Activity Landscape Features. The calculated landscape probabilities reflect the propensities of

compounds to be involved in activity landscape features based on the available data. As such, they can aid in prioritizing and identifying key compounds for further optimization in a given data set. However, the transfer of SAR-informative features between different data sets is principally limited. Because the probabilities only reflect the SAR information content of available data, the prediction of new compounds is difficult. However, key compounds for further chemical exploration can be rationally selected from a given data set on the basis of landscape feature probabilities. For the optimization of such compounds, detailed assessment and visual inspection by experts remains a critically important step.

CONCLUSIONS

Herein, we have reported a large-scale analysis of conditional activity landscape feature probabilities across bioactive compounds. Approx. 160,000 publicly available bioactive compounds with high-confidence activity annotations belonging to 427 target-based data sets were profiled. Pairwise structural and potency similarities were systematically determined for all compounds and conditional feature probabilities were derived on the basis of structural and potency comparisons. According to the methodology, 10% and 1% of the compounds were expected to exceed the thresholds for the single-feature and combined-feature categories, respectively. These values provided reference points for the assessment of data set-dependent feature probabilities. Significant deviations from the expected values were observed, thus indicating that activity landscape features were often unevenly distributed across different target-based sets. In many instances, the frequency distributions displayed large variations, more so for IC_{50} than K_i data sets. Furthermore, four feature categories including Cat 2, 7, and 8 (all sets) and 6 (IC_{50} sets) displayed a clear dependence on the data set size. Approximately half of all data set compounds did not reach any feature probability threshold and were hence classified as SAR noninformative. However, individual data sets were identified that were particularly rich in compounds forming specific landscape features. Moreover, when determined at the level of individual compounds, similarity cliffs represented the most prevalent SAR-informative activity landscape feature. On a per-compound basis, the formation of similarity cliffs was ~ 100 -fold more likely than the formation of activity cliffs (in activity landscapes, activity cliffs are rare^{1,2}), and there was a generally high probability for compounds to form similarity cliffs. It should be noted that the compound-based activity landscape feature approach applied herein is descriptive in nature and not used to predict new compounds. Furthermore, it should also be noted that activity landscape feature probabilities have been derived on the basis of 2D molecular representations and similarity assessment. Given the uncertainties in predicting bioactive compound conformations, the accurate assessment of which would be a prerequisite for meaningful activity landscape analysis, the use of conformation-dependent 3D descriptors on the basis of modeled compound conformations is not recommended for the derivation of activity landscape feature probabilities. Recent progress has been made in rationalizing activity cliffs in three dimensions on the basis of systematic X-ray data analysis,¹² and similar structure-based analyses might be carried out, for example, for similarity cliffs. However, the current knowledge base of 3D activity cliffs is still rather limited,¹² which precludes large-scale analysis comparable to the study reported herein. As a final note, the overall high

probability of compounds to form similarity cliffs revealed in our analysis also implies that structurally diverse compounds with similar activity (having a high probability to form similarity cliffs) can be identified for many different targets.

AUTHOR INFORMATION

Corresponding Author

*Tel.: +49-228-2699-306. Fax: +49-228-2699-341. E-mail: bajorath@bit.uni-bonn.de.

Notes

The authors declare no competing financial interest.

REFERENCES

- (1) Wassermann, A. M.; Wawer, M.; Bajorath, J. Activity Landscape Representations for Structure–Activity Relationship Analysis. *J. Med. Chem.* **2010**, *53*, 8209–8223.
- (2) Stumpfe, D.; Bajorath, J. Exploring Activity Cliffs in Medicinal Chemistry. *J. Med. Chem.* **2012**, *55*, 2932–2942.
- (3) Iyer, P.; Stumpfe, D.; Vogt, M.; Bajorath, J.; Maggiora, G. M. Activity Landscapes, Information Theory, and Structure–Activity Relationships. *Mol. Inf.* **2013**, *32*, 421–430.
- (4) Vogt, M.; Iyer, P.; Maggiora, G. M.; Bajorath, J. Conditional Probabilities of Activity Landscape Features for Individual Compounds. *J. Chem. Inf. Model.* **2013**, *53*, 1602–1612.
- (5) Gaulton, A.; Bellis, L. J.; Bento, A. P.; Chambers, J.; Davies, M.; Hersey, A.; Light, Y.; McGlinchey, S.; Michalovich, D.; Al-Lazikani, B.; Overington, J. P. ChEMBL: A Large-scale Bioactivity Database for Drug Discovery. *Nucleic Acids Res.* **2011**, *40*, D1100–D1107.
- (6) Shanmugasundaram, V.; Maggiora, G. M. Characterizing Property and Activity Landscapes Using an Information-Theoretic Approach. Presented at the 222nd American Chemical Society National Meeting, Division of Chemical Information, Chicago, IL, August 26–30, 2001; American Chemical Society: Washington, DC, 2001; abstract no. 77.
- (7) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.
- (8) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**, *50*, 742–754.
- (9) Stumpfe, D.; Bajorath, J. Frequency of Occurrence and Potency Range Distribution of Activity Cliffs in Bioactive Compounds. *J. Chem. Inf. Model.* **2012**, *52*, 2348–2353.
- (10) Liu, T.; Lin, Y.; Wen, X.; Jorissen, R. N.; Gilson, M. K. Binding-DB: A Web-Accessible Database of Experimentally Determined Protein–Ligand Binding Affinities. *Nucleic Acids Res.* **2007**, *35*, D198–D201.
- (11) Dimova, D.; Stumpfe, D.; Bajorath, J. Quantifying the fingerprint descriptor dependence of structure–activity relationship information on a large scale. *J. Chem. Inf. Model.* **2013**, *53*, 2275–2281.
- (12) Stumpfe, D.; Hu, Y.; Dimova, D.; Bajorath, J. Recent Progress in Understanding Activity Cliffs and Their Utility in Medicinal Chemistry. *J. Med. Chem.* **2014**, *57*, 18–28.