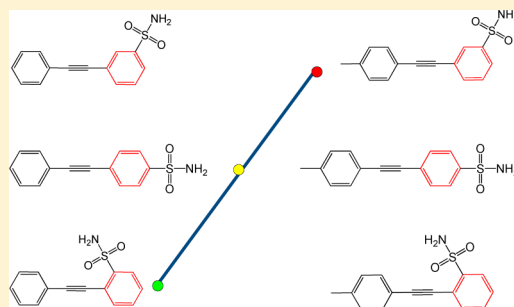


# Systematic Assessment of Compound Series with SAR Transfer Potential

Bijun Zhang,<sup>†</sup> Anne Mai Wassermann,<sup>#,†</sup> Martin Vogt,<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, Dahlmannstrasse 2, D-53113 Bonn, Germany

**ABSTRACT:** Compound series with different core structures that contain pairs of analogs with corresponding substitution patterns and similar activity represent structure–activity relationship (SAR) transfer events. On the basis of the matched molecular pair (MMP) formalism and linear regression analysis of compound potencies, a general approach is introduced for the identification of SAR transfer series (SAR-TS) and SAR-TS with regular potency progression (SAR-TS-RP). We have systematically extracted such series from public domain compound data and analyzed their size distribution and structural characteristics. More than 900 SAR-TS and 500 SAR-TS-RP with high-confidence potency annotations were identified in various compound activity classes. These series provide a substantial knowledge base for the analysis and prediction of SAR transfer and are made publicly available.



## INTRODUCTION

SAR transfer is a frequently considered topic in medicinal chemistry. This process can be conceptualized in at least two different ways. First, one might consider an individual compound series binding to two (related) targets. If similar structure–activity relationships (SARs) are observed, then the situation might be regarded as SAR transfer from one target to another. Second, one might consider two chemically different series with corresponding analogs that display pairwise similar activity against a given target. Then, the target-specific SAR might be transferred from one series to another. SAR transfer from one series to another is particularly relevant for practical hit-to-lead and lead optimization efforts in medicinal chemistry. This is the case because one is often interested in replacing a lead series that displays promising SAR progression but is associated with liabilities (such as toxicity or unfavorable pharmacology) with a structurally distinct series to circumvent these complications. In addition, replacement of existing series is often considered to establish new patent positions. This would ideally be accomplished by generating corresponding analogs on a different structural background that yield comparably attractive SAR trends. However, finding suitable replacement series is in practice often a challenging task because it is generally difficult to predict corresponding series with similar activity and SAR characteristics. Despite its high relevance for medicinal chemistry applications, SAR transfer has thus far only been little investigated from a computational perspective; perhaps surprisingly. SAR transfer can also be rationalized as an advanced form of scaffold hopping where not only the replacement of active core structures is required but also conservation of SAR characteristics.

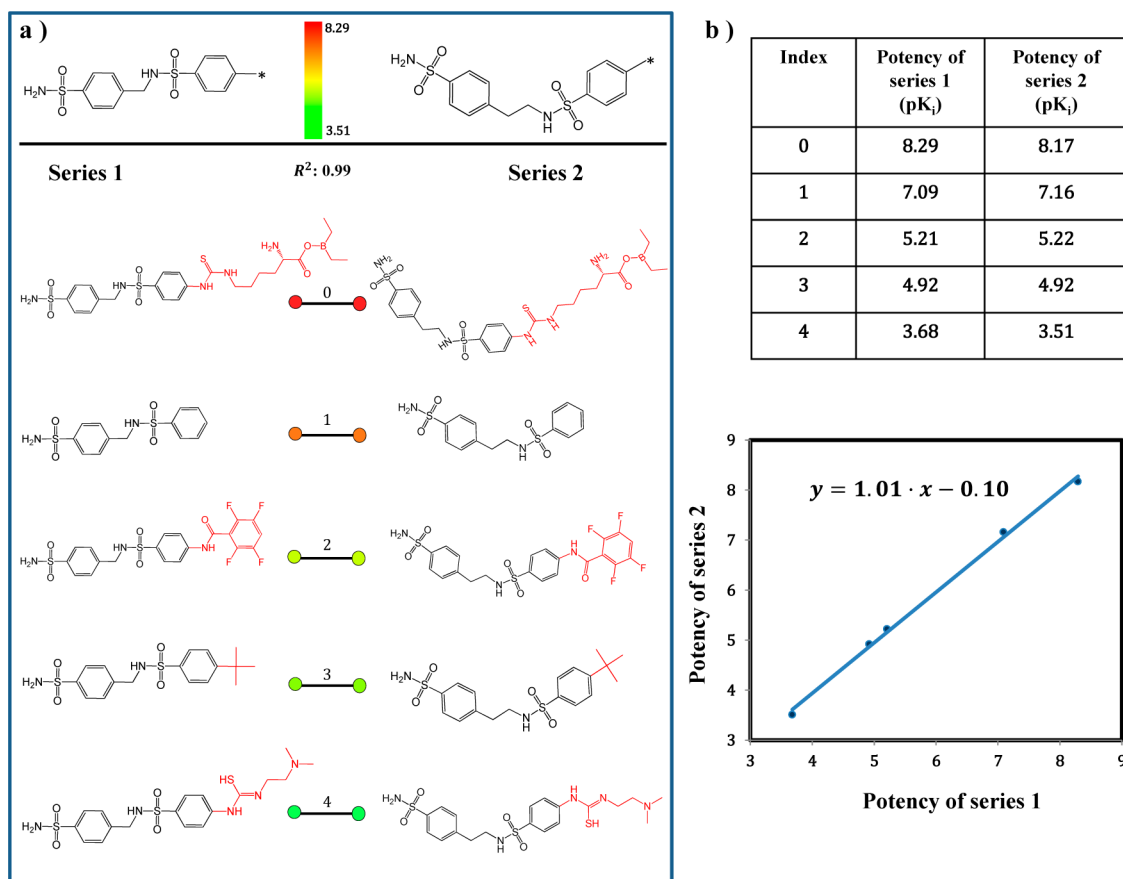
In our exploration of SAR transfer events, we have focused on the identification of different compound series with similar

SAR behavior against individual targets. An opportunity to rationalize SAR transfer on a large scale is provided by mining of currently available compound data. To these ends, we have previously introduced a data mining approach to identify SAR transfer series (SAR-TS).<sup>1</sup> The method was primarily designed to find a matching SAR transfer series for a given series of analogs. Therefore, the compound scaffold<sup>2</sup> shared by an analog series of interest was extracted from the series, and a search was carried out to identify scaffolds in other series that differed from the starting scaffold by the replacement of one contiguous ring system. Template and target scaffolds were compared by generating matched molecular pairs (MMPs).<sup>3</sup> An MMP is defined as a pair of compounds that differ only by a chemical change at a specific site, i.e., by the exchange of a pair of substructures. Analog series with qualifying scaffolds were subjected to R-group decomposition, and compound pairs with corresponding substitutions were identified. A potency-based scoring function was applied to identify SAR transfer series with regular potency progression (SAR-TS-RP).<sup>1</sup> The methodology was also adapted to search for possible scaffold pairings meeting the one-ring-system difference criterion. In BindingDB,<sup>4</sup> a total of 405 possible series alignments were detected, 61 of which qualified as SAR-TS-RP, with an average of close to four analog pairs per series.<sup>1</sup> In addition to this approach, we have also designed a methodology to mine compound network representations of individual data sets for characteristic subgraphs that might indicate the presence of SAR transfer events.<sup>5</sup>

Herein we introduce a generalized approach to identify compound series with SAR transfer potential that combines the

**Received:** October 9, 2012

**Published:** November 27, 2012



**Figure 1.** Illustration of SAR transfer. In (a), an exemplary SAR-TS with regular potency progression consisting of carbonic anhydrase I inhibitors is represented according to ref 1. Compound pairs are arranged in the order of increasing potency as illustrated by pairs of color-coded dots using a continuous color spectrum from green (lowest potency) over yellow to red (highest potency). Common core structures of analogs are drawn in black and substituents (R-groups) in red. This series yields an  $R^2$  value of 0.99. In (b), the potency values for the compound pairs in this series and the corresponding regression model are reported. This model yields a slope  $S$  of 1.01. This SAR-TS-RP represents a nearly perfect SAR transfer event, consistent with its  $R^2$  and  $S$  values.

MMP formalism without restriction to predefined scaffolds with linear regression analysis of compound potency value distribution in different series. The methodology is shown to be a reliable indicator of SAR transfer events. It has been applied to systematically search for SAR-TS and SAR-TS-RP in more than 900 target-directed compound data sets, identifying more than 900 qualifying series consisting of at least three corresponding analog pairs. Hence, as presented, the analysis is general.

## MATERIALS AND METHODS

**Matched Molecular Pairs.** MMPs were systematically generated using an in-house implementation of the algorithm by Hussain and Rea.<sup>6</sup> Compounds were fragmented by deleting one, two, or three nonring single bonds (referred to single-, double-, or triple-cuts, respectively) attached to a ring system. A single-cut generated two fragments of a compound (represented as canonical SMILES strings<sup>7</sup>). The larger fragment was selected as the key fragment and the smaller fragment as the corresponding value fragment. Values and corresponding keys were stored in an index table. Nonidentical fragments of the same size were each stored once as the key fragment and value. For double- or triple-cuts, core fragments with two or three attachment points were obtained, respectively, plus two or three terminal fragments (connectivity information was

retained). In these cases, the larger fragment(s) were considered the key and the smaller the value. Due to the large number of possible fragments consisting of single hydrogen atoms, these fragments were retained in a postprocessing step only if they were involved in MMP formation.

**Qualifying Compound Series.** Analog series were identified by searching the index table for MMPs sharing the same key and different values. Two matching series (MS) were identified by pairwise comparison of values in the index table. If two compounds had different keys but shared the same value, they formed a pair of matching compounds. A minimum of three matching pairs was required to form an MS. In addition, we also required that MS spanned at least two orders of magnitude in potency, regardless of their size, in order to prequalify series for potentially attractive SAR transfer events (different from “flat” SARs).

**Linear Regression.** If compounds forming an MS have pairwise comparable or corresponding potency, the SAR transfer condition is met. The correspondence of logarithmic potency values was assessed by linear regression. For regression modeling, each data point ( $x_i, y_i$ ) represented a matching pair of compounds. Least-squares fitting was applied to minimize the total sum of squares of residuals, derive the best linear model, and calculate the coefficient of determination:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - a - S \cdot x_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Here,  $\bar{y}$  is the mean of the  $y_i$ ,  $n$  refers to the total number of data points, and  $a$  and  $S$  are the intercept and slope of the regression line, respectively.

**SAR Transfer Criteria.** Prequalified MS (as discussed above) were subjected to the assessment of SAR transfer potential. The  $R^2$  value of the regression model was used as a primary measure of SAR transfer (i.e., similarity of potency values of analogs forming a pair). Furthermore, the slope  $S$  of the fitted linear model was used as a measure of potency progression within the MS. A slope of 1 would indicate linearly increasing logarithmic potency of the same magnitude between pairs of corresponding analogs, hence representing an “ideal” SAR transfer event. On the basis of initial test calculations to evaluate SAR transfer potential for exemplary series, the following criteria were applied:

(1). *SAR Transfer Series.* To qualify as SAR-TS, an  $R^2$  value of at least 0.98 and a slope  $S$  between 0.5 and 2.0 were required.

(2). *SAR Transfer Series with Regular Potency Progression.* Furthermore, to qualify as SAR-TS-RP, an  $R^2$  value of at least 0.98 and a slope  $S$  between 0.9 and 1.1 were required.

It should be emphasized that the criteria applied for SAR-TS-RP formation also select for essentially ideal SAR transfer events.

**Implementation.** Steps required for the generation of MMPs, identification of MS, construction of regression models, and the search for extended SAR transfer series were implemented in Java using the OpenEye chemistry toolkit.<sup>8</sup>

**Data Sets.** From the current version of BindingDB,<sup>4</sup> a total of 918 compound data sets (activity classes) for human targets were extracted that contained six or more compounds for which  $K_i$  values were available as potency measurements. If multiple  $K_i$  values were available for a compound, their geometric mean was calculated as the final potency annotation. The 918 qualifying data sets contained 64,633 unique compounds with a total of 110,392 potency annotations.

## RESULTS AND DISCUSSION

**SAR Transfer Events.** Figure 1a shows a nearly ideal example of an SAR transfer event. The SAR-TS consists of five

Table 1. SAR Transfer Series<sup>a</sup>

Class size	# MS	# SAR-TS	# SAR-TS-RP
[6, 10]	0	0	0
[10, 50)	30	8	4
[50, 100)	63	15	6
[100, 200)	131	29	19
[200, ∞)	3582	897	567
total	3806	949	596

<sup>a</sup>Reported are the number of prequalified matching series (MS), SAR transfer series (SAR-TS), and SAR transfer series with regular potency progression (SAR-TS-RP) extracted from BindingDB. “Class size” refers to the size (number of compounds) of the activity classes from which the series originated.

pairs of carbonic anhydrase I inhibitors with different core structures and corresponding substituents. Both analog series display regular (linear) and closely corresponding progression toward higher potency, consistent with an  $R^2$  value of 0.99 and a slope of 1.01 of the regression model (Figure 1b). Following

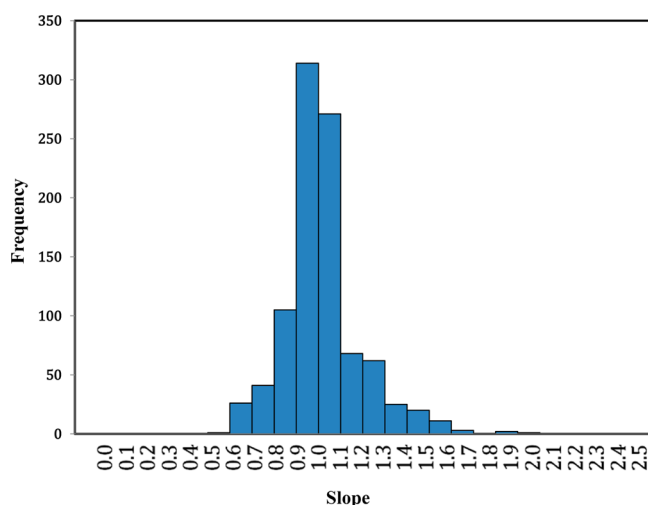
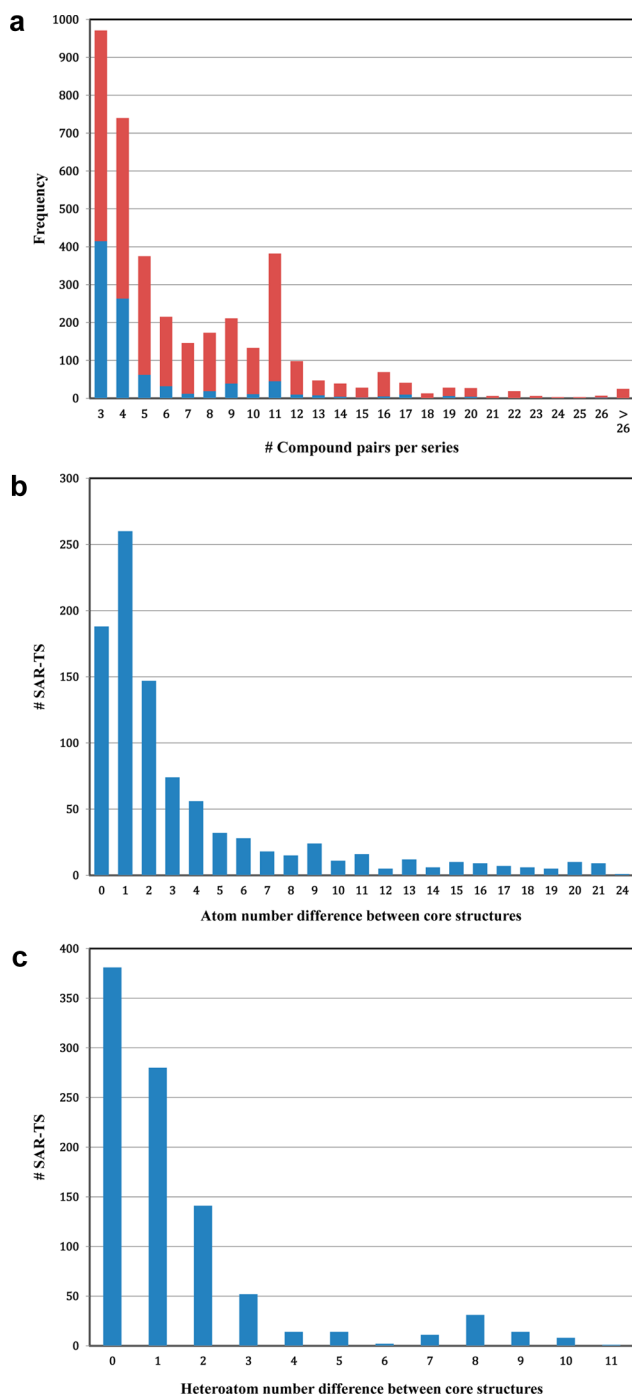


Figure 2. Slope distribution. For all SAR-TS, the distribution of regression model slopes is reported. The apparent asymmetry of the slope distribution is artificial. It depends on which of two compared series is used as independent variable (reference). On a logarithmic scale, the slope distribution is symmetrical.

the terminology used herein, this series represents a prototypic SAR-TS-RP. The participating analog series reveal equivalent SAR information, with corresponding substitutions of different cores leading to steadily increasing compound potency. In this case, core structures (distinguished by linker length between the two sulfonamide moieties) can be readily exchanged without affecting SAR behavior. This is the type of SAR transfer events that are of particular interest in medicinal chemistry. As shown below, there can be much larger differences between core structures of an SAR-TS than in this example.

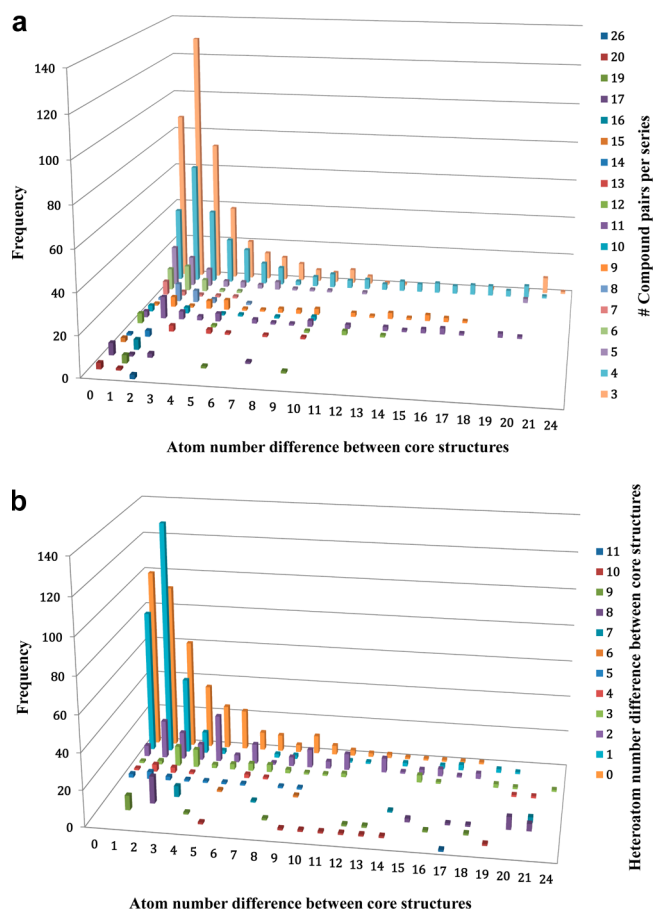
**Methodological Concept.** With the approach introduced herein, we have aimed to generalize the search for SAR transfer series. On the basis of the MMP formalism, matching analog series can be readily identified in compound databases. The additional criterion that MS should span at least two orders of magnitude in potency focuses the analysis on series that have the potential to qualify as transfer series. The pairwise alignment of candidate series is also straightforward on the basis of key and value information contained in the MMP index table. Furthermore, the generation of regression models of corresponding potency values then identifies series that represent SAR transfer events (based on  $R^2$  values reflecting correlation) and, in addition, transfer events with regular potency progression (based on slope values of the regression models). The latter series are of particular interest for medicinal chemistry, given their SAR continuity over a large potency range. In addition, these series are also amenable to comparative QSAR modeling.

**Identification of SAR Transfer Series.** Applying the MMP-based methodology, we have identified a total of 3806 MS in BindingDB with at least three corresponding analog pairs (Table 1). As expected, most MS were found in large activity classes, but MS were also present in relatively small classes (Table 1). On the basis of our regression criteria, the 3806 MS contained 949 SAR-TS, and a subset of 596 of these series qualified as SAR-TS-RP. Thus, ~25% of all MS represented SAR transfer events and ~16% transfer events with regular potency progression. SAR-TS were found in 83 of the 918 activity classes we analyzed and SAR-TS-RP in 61 of them. Thus, SAR transfer was detected in ~9% of all activity classes,



**Figure 3.** Distribution of SAR transfer series. In (a), the length distribution of MS (blue plus red bars) and SAR-TS (blue bars) is reported. MS peaks for larger numbers of compound pairs per series (e.g., 11) result from many different activity classes. (b) shows the distribution of SAR-TS with respect to the size difference (number of non-hydrogen atoms) between the core structures representing the two participating analog series. In (c), the corresponding distribution is shown for differences in the heteroatom composition between core structures.

indicating that SAR transfer was an overall rare event. However, in classes where SAR transfer occurred, multiple SAR-TS series were typically detected with, on average, 11–12 series per class. Figure 2 shows the distribution of slopes of the regression models of all 949 SAR-TS. The distribution reveals a clear peak in the *S* value interval [0.9, 1.1]. Thus, the majority of transfer



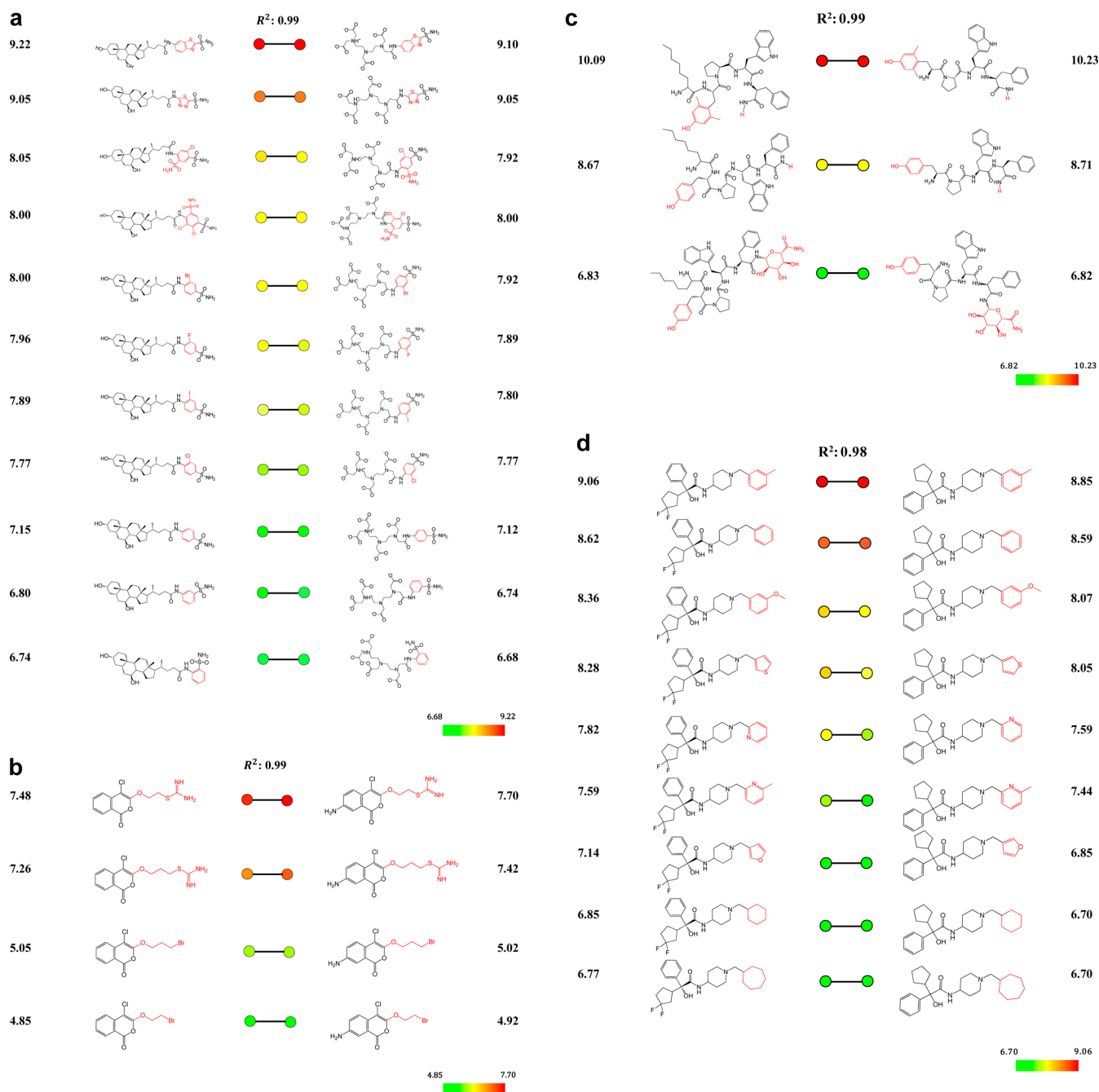
**Figure 4.** Length and composition characteristics. In (a), the distribution of SAR-TS with increasing length (color-coded) and size difference between core structures is reported and in (b), the distribution of SAR-TS with respect to increasing size difference between core structures and their heteroatom difference (color-coded).

series displayed nearly linear potency progression, an unexpected finding.

**Characterization of SAR Transfer Series.** We next analyzed the composition of all SAR-TS. Figure 3a reports the length distribution of transfer series (compared to MS). The majority of SAR-TS consisted of three to four analog pairs. Series with five to 11 pairs were also frequently observed, and SAR-TS with up to 20 analog pairs were detected.

As reported in the Methods section, the size of substitutions was limited to ensure that SAR-TS consisted of series of structurally closely related compounds. However, for core structures in series, no size or composition difference restrictions were applied. In fact, SAR-TS with significantly different core structures are of particular interest for SAR transfer analysis. In principle, the more core structures differ in size and composition, the less likely it is to facilitate SAR transfer. Figure 3b reports the size difference (number of non-hydrogen atoms) between core structures of analog series forming an SAR-TS. The majority of core structures were of the same or similar size, i.e., they differed at most by one or two atoms. In addition, in about 150 SAR-TS, core structures differed by three to five atoms, and large size differences of up to 24 atoms were also observed. Figure 3c reports the difference in heteroatoms between core structures, another measure of chemical diversity. The majority of core structures in transfer





**Figure 5.** Exemplary series. SAR-TS-RP of different length and composition are shown according to Figure 1a: (a) carbonic anhydrase II inhibitors, (b) urokinase-type plasminogen activator inhibitors, (c)  $\mu$ -opioid receptor antagonists, (d) muscarinic acetylcholine receptor M1 antagonists.

series had the same heteroatom content or differed by one or two heteroatoms. However, core structure pairs with differences of up to 11 heteroatoms were also found.

In Figure 4a, SAR-TS are classified according to their size and the difference in atom numbers of core structures. Here, a number of series with 10 or more corresponding analog pairs and large differences in heteroatom content were detected. In addition, in Figure 4b, SAR-TS are classified considering the difference in size and heteroatom content of their core structures. In this case, a number of instances with large difference in both the size and heteroatom content of core structures were observed, hence indicating the presence of chemically distinct cores in these transfer series.

**Exemplary Series.** Figure 5 shows representative SAR-TS-RP with different characteristics. In Figure 5a, a transfer series of carbonic anhydrase II inhibitors is shown that consisted of 11 corresponding analog pairs. In this case, the core structures contained the same number of atoms but were chemically completely distinct. It would essentially be impossible to predict this SAR transfer event, although it can be rationalized. A hallmark of carbonic anhydrase inhibition is the stringent requirement of a sulfonamide group (to complex the catalytic zinc ion in the active site), which represent the major inhibitory interaction. Accordingly, the sulfonamide was strictly conserved in all analogs. This requirement explains the tolerance for completely different scaffolds in carbonic anhydrase inhibitors. This transfer series also illustrates the utility of the MMP

formalism for series identification. The exchanged fragments (colored in red) that structurally distinguish these inhibitors were a part of the sulfonamide-presenting substituent and were automatically detected as differences in values. The sulfonamide group itself and the core structure together represent the conserved (discontinuous) key substructure of each analog series. Figure 5b shows another SAR-TS-RP comprising four analog pairs of urokinase-type plasminogen activator inhibitors, which provides an example of a transfer series with closely related relatively small core structures and large R-groups, yielding nearly linear potency progression. In Figure 5c, an SAR-TS-RP consisting of three pairs of  $\mu$ -opioid receptor antagonists is shown. This series is representative of large and complex core structures with corresponding substituents at two sites resulting from double-cut MMPs. Series of such complexity would be difficult to identify without applying the MMP formalism. Finally, in Figure 5d, a transfer series consisting of nine pairs of muscarinic acetylcholine receptor M1 antagonists is presented that contained chemically distinct scaffolds of similar size with various ring substituents at a single site. This series represents another exemplary SAR transfer event. Taken together, these examples illustrate that SAR transfer series often have different structural characteristics.

## CONCLUSIONS

Herein we have introduced a generally applicable methodology for the identification of SAR transfer events that combines the MMP concept and potency regression analysis. The approach does not utilize predefined molecular scaffolds as core structures of compound series but considers all possible structural modifications captured in MMP index tables to define matching analog series. By adding stereochemistry- or tautomer-sensitive molecular representations to MMP-based transfer alignments, this information can be easily incorporated to further evaluate series alignments. In our systematic search of more than 900 compound activity classes, SAR-TS were detected in ~9% of these classes, consistent with the expectation that SAR transfer events are quite rare in compounds originating from medicinal chemistry programs. However, in 83 activity classes, each of which corresponded to a unique target, a total of 949 SAR-TS were identified. In cases where SAR transfer was detected multiple transfer series were usually identified. Hence, it is conceivable that target proteins of these series are particularly tolerant to SAR transfer, i.e., capable of accommodating various core structures with similar binding modes, which is expected to provide useful information for structure-based compound design and medicinal chemistry projects. Furthermore, the SAR-TS we identified often contained chemically distinct cores. Moreover, many transfer series displayed nearly linear potency progression. Hence, these transfer series are continuous in their SAR character and fall into the applicability domain of QSAR methods. The collection of SAR-TS identified herein and the associated target information represent a substantial knowledge base for the study of SAR transfer events and the selection of core structures as possible replacements for compound series associated with liabilities. Upon publication, this information is freely available via the following URL: <http://www.lifescienceinformatics.uni-bonn.de>.

## AUTHOR INFORMATION

### Corresponding Author

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

### Present Address

<sup>#</sup>In Silico Lead Discovery, Novartis Institutes for BioMedical Research, Inc., 250 Massachusetts Avenue, Cambridge, MA 02139, USA.

### Notes

The authors declare no competing financial interest.

## REFERENCES

- (1) Wassermann, A. M.; Bajorath, J. A Data Mining Method to Facilitate SAR Transfer. *J. Chem. Inf. Model.* **2011**, *51*, 1857–1866.
- (2) Bemis, G. W.; Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *J. Med. Chem.* **1996**, *39*, 2887–2893.
- (3) Kenny, P. W.; Sadowski, J. Structure Modification in Chemical Databases. In *Chemoinformatics in Drug Discovery*; Oprea, T. I., Ed.; Wiley-VCH: Weinheim, Germany, 2005; pp 271–285.
- (4) 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.
- (5) Gupta-Ostermann, D.; Wawer, M.; Wassermann, A. M.; Bajorath, J. Graph Mining for SAR Transfer Series. *J. Chem. Inf. Model.* **2012**, *52*, 935–942.
- (6) 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.
- (7) Weininger, D. SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules. *J. Chem. Inf. Comput. Sci.* **1988**, *28*, 31–36.
- (8) OEChem TK version 1.7.4.3; OpenEye Scientific Software Inc.: Santa Fe, NM, 2010.