

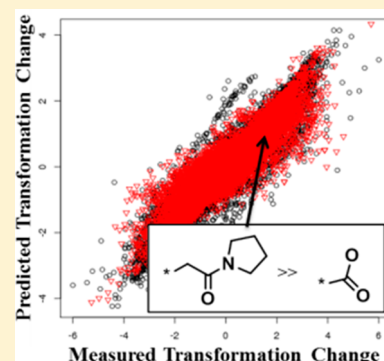
Quantitative Structure–Activity Relationship Models of Chemical Transformations from Matched Pairs Analyses

Jeremy M. Beck and Clayton Springer*

Novartis Institutes for BioMedical Research, 100 Technology Square, Cambridge, Massachusetts, United States

S Supporting Information

ABSTRACT: The concepts of activity cliffs and matched molecular pairs (MMP) are recent paradigms for analysis of data sets to identify structural changes that may be used to modify the potency of lead molecules in drug discovery projects. Analysis of MMPs was recently demonstrated as a feasible technique for quantitative structure–activity relationship (QSAR) modeling of prospective compounds. Although within a small data set, the lack of matched pairs, and the lack of knowledge about specific chemical transformations limit prospective applications. Here we present an alternative technique that determines pairwise descriptors for each matched pair and then uses a QSAR model to estimate the activity change associated with a chemical transformation. The descriptors effectively group similar transformations and incorporate information about the transformation and its local environment. Use of a transformation QSAR model allows one to estimate the activity change for novel transformations and therefore returns predictions for a larger fraction of test set compounds. Application of the proposed methodology to four public data sets results in increased model performance over a benchmark random forest and direct application of chemical transformations using QSAR-by-matched molecular pairs analysis (QSAR-by-MMPA).



INTRODUCTION

In the pharmaceutical field, lead optimization is an iterative process where a project team makes small changes to a molecule to bring about a change in activity, absorption, distribution, metabolism, and excretion (ADME) and/or pharmacokinetic/pharmacodynamic (PK/PD) properties. The identification of activity cliffs in data sets provides information on small structural changes that bring about large changes in activity, as well as regions of a lead where modifications would have the largest effect; because of this, activity cliffs have become a prominent topic in drug design.¹ A recent study found that activity cliff progression through public data sets is relatively low, suggesting that, on the whole, activity cliffs have been underutilized in lead optimization. Systematic methods for the identification of activity cliffs have been reported,^{2,3} such as the structure–activity landscape index (SALI), which compares activity changes between compounds scaled by molecular similarity.² It has been demonstrated that the definition of an activity cliff is dependent on the molecule representation (the similarity term) and recent efforts have utilized 3D representations of the ligand⁴ or crystal structures of bound ligands to rationalize cliffs.⁵ In addition to typical quantitative structure–activity relationship (QSAR) methods, prospective determinations of activity cliffs have recently been demonstrated through QSAR modeling techniques to predict SALI values.⁶ The demonstration by Guha that SALI could be predicted indicates that QSAR models can be trained to predict activity cliffs within an assay and, by association, the activity differences between a pair of compounds.

Another recent topic of interest in the drug design community is that of matched molecular pairs.^{7,8} Matched molecular pairs are generally defined as two compounds that differ by only a small, well-defined chemical transformation, for example an R-group transformation, or substitution of a central ring system. In the framework of SALI, matched molecular pairs can be thought of as the subset of paired compounds that have high similarity scores (i.e., the subset of compounds which could possess meaningful activity cliffs). Initially, these pairs were identified using an MCS routine⁷ although a significantly faster algorithm for matched pairs identification that uses fragmentation was recently developed by Hussain and Rea.⁹ By characterizing all matched molecular pairs within a data set, and then calculating the associated change in an activity or property for each matched pair, the frequency and mean effect of each chemical transformation may be evaluated. Reports of matched pairs analyses have ranged from ADME properties¹⁰ to three-dimensional analysis of binding modes.⁵ Recently, databases of matched pairs have been made publicly available—for example, the SwissBioisostere database¹¹ and VAMMPIRE,¹² which is restricted to 3D matched pairs with corresponding crystallographic binding modes.

The assumption that a chemical transformation should result in a predictable change in property value led to the hypothesis that one could use simple matched molecular pairs analysis (MMPA) to predict the activity of novel compounds. The observed property change for a chemical transformations within

Received: January 8, 2014

a data set can be applied to the prediction of activity values for new compounds that differ from a known compound by only that chemical transformation and determining the resulting activity change from previous observations.¹³ This method, named QSAR-by-MMPA, exhibited improved prediction accuracy of KCNQ1–KCNE1 channel inhibition compared to a nearest neighbors analysis and random forest algorithm.¹⁴ However, the authors note that the ability of the method is tied to the existence of matched pairs between prospective compounds and the training set molecules, and previous knowledge of the activity change associated with those transformations. In many data sets, especially those for active drug discovery projects where exploitable activity cliffs would be highly beneficial, the reliance on knowing the effect of the chemical transformation severely limits the utility of the QSAR-by-MMPA method. In many cases, a transformation is novel, or has been observed too few times to achieve a statistically significant activity change.

In this manuscript we present an adaptation of the QSAR-by-MMPA method¹⁴ that defines pairwise descriptors for each matched pair using the difference between the individual compounds' descriptors. Using those descriptors, we generate a random forest model¹⁵ to predict each matched pair's activity change using a QSAR model on chemical transformations. This methodology of generating pairwise descriptors and constructing QSAR models to predict activity changes was successfully demonstrated for the prediction of activity cliffs.⁶ Predictions for novel compounds could be generated in the same manner as QSAR-by-MMPA, with the only difference being the reliance on a QSAR model. However, the QSAR model would also enable the ability to predict the activity changes associated with novel chemical transformations. The QSAR model is also trained on infrequently observed transformations, which encompass the bulk of the matched pairs. In addition, the use of pairwise physicochemical descriptors to describe the chemical transformations allows for the clustering of transformations based on the changing groups and the chemical context of those transformations. This natural clustering of chemical transformations will allow infrequently observed transformations to be supported by their neighbors in the model, while also subdividing common transformations by the environment on which they occur. In addition, the transformation QSAR model would allow predictions for novel chemical transformations based on their similarity to known transformation. In cases where the SAR is additive, the model should exhibit increased accuracy, while nonadditive SAR would offer no benefit over the standard QSAR model.

METHODOLOGY

Data Sets. For this study, we will use four public data sets obtained from the ChEMBL database. In total, these data sets cover a range of sizes, compound similarity, and QSAR properties. Two end points represent concentration data: the 3D7 data set (Assay ID: 1000386)¹⁶ is an ED₅₀ value, and the Factor Xa data set, which is a collection of assays reporting IC₅₀ values against Factor Xa. (full data set in Supporting Information) The 3D7 data set is relatively small and has a high average similarity between compounds, while the Factor Xa data set is comprised of a larger number of compounds that have a low average similarity. Two additional data sets are also used: volume of distribution at steady-state (Assay ID: 1614670)¹⁷ and logP (octanol/water) of a series of Renin

inhibitors (Assay ID: 1948000) which was previously used for 3D-QSAR methodology testing.¹⁸

Descriptors. To represent each compound, we used three sets of descriptors. A set of molecular descriptors including 51 2-D physicochemical descriptors from MOE (see Supporting Information for a full list), the fraction ionized (sum cation/anion, neutral anion/cation, and fraction ionization) of each compound determined using the Henderson–Hasselbalch equation and pK_as calculated in MoKa, and the chiral Morgan Fingerprint^{19,20} with a radius of 3 generated in RDKit²¹ and hashed to 1024 bits using python's built-in hash function.

All QSAR modeling was carried out using a random forest predictor, as implemented in the randomForest package of R.²² The random forest algorithm is a robust ensemble method originally proposed by Tin Kam Ho²³ and adapted by Breiman.¹⁵ It incorporates feature selection implicitly and is generally considered resistant to overfitting. Unless otherwise noted, a total of 400 decision trees were used to construct each ensemble; this was lowered from the default value as previous studies have indicated that increasing the number of trees generally results in minimal gains for the computational cost.²⁴ In the case of the logP models, the descriptors corresponding to calculated logP and the change in calculated logP are omitted from the model. This choice had a negligible effect on model performance.

Matched pairs calculation: Matched pairs were calculated using the recent algorithm detailed by Hussain and Rea.⁹ To define a pair, the changing portion of the molecule was allowed to contain up to 13 heavy atoms and up to 2 bond cuts to allow pairs corresponding to ring swapping. Two sets of matched pairs were calculated within each data set: "training" pairs and "prospective" pairs. Training pairs are defined as a matched pair between two training set molecules. Prospective pairs are defined as a matched pair between a training set molecule and test set molecule. To evaluate the performance of each MMP-based QSAR method, the same set of training pairs and prospective pairs were used across all methods.

QSAR-by-MMPA. The QSAR-by-MMPA method was used to estimate properties of prospective compounds as described in a previous report.¹⁴ We outline our implementation here and have included a schematic depiction in the Supporting Information. Each data set was split into a training set and test set. Matched pairs between training set molecules were determined, and their chemical transformations were stored with the mean observed property change. For each prospective matched pair, the corresponding chemical transformation was looked up in the training set. If a chemical transformation for a prospective pair was observed within the training set, the associated property change was used for the prospective matched pair. In the event that the chemical transformation was not observed within the training set, the prospective matched pair was discarded. For each prospective matched pair, an estimate of the test set compound was obtained by adding the look-up property change value to the property value of the training set compound. If a test set compound formed multiple pairs with the training set, the mean of the predictions is used.

QSAR-by-MMPA was carried out using three different variations of the QSAR-by-MMPA methodology. The first method, which we will refer to as "global MMPA" throughout the manuscript, is based on the assumption that chemical transformations are globally transferable (i.e., all H >> F transformations are equivalent). This assumption is generally not valid, with the exception of trivial cases (prediction of

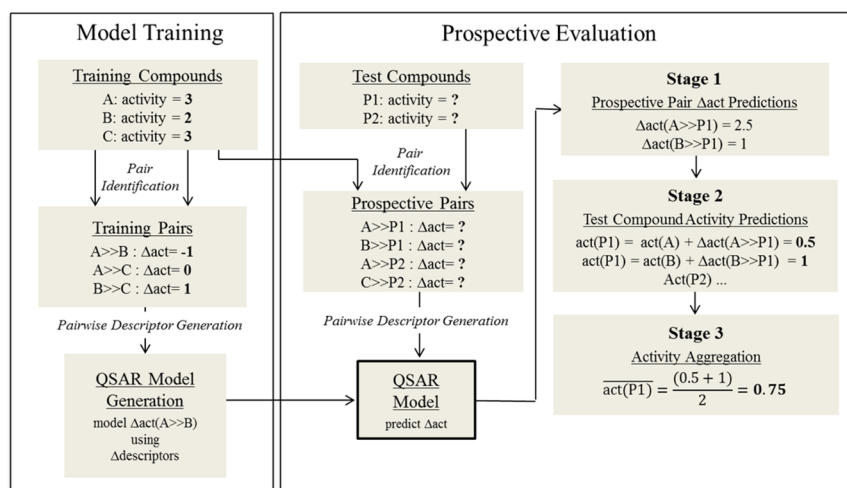


Figure 1. Schematic representation of the transformation QSAR protocol. A, B, and C are molecules in the training set. P1 and P2 are not yet assayed molecules. Analogous figures for QSAR and QSAR-by-MMPA are included in the Supporting Information.

molecular weight changes) or when SAR is additive, for example the prediction of logP. The second and third MMPA methods attempt to distinguish between differing chemical environments for each transformation using two different representations: the radial environment on the MCS by including the portion of the MCS up to three bonds away from the transformation (radial MMPs) or the chemotype the transformation occurs on using similarity of the full MCS (chemotype MMPs). Both of these methods were previously used in the analysis of MMPs to identify structurally significant clusters within each chemical transformation.¹⁰ The radial environment of each pair was stored as a SMILES string, and only transformations with matching environments were considered to estimate prospective pairs. The chemotype filtering of transformations based on the MCS was carried out as follows: for each pair, the MCS was extracted and a Morgan fingerprint (radius = 3) was calculated. For a given prospective matched pair, all examples of that chemical transformation from the training pairs were selected, and the Tanimoto fingerprint similarity between the prospective and training MCS was calculated. If the similarity was greater than 0.7, the two pairs were considered to occur on the same chemotype, and the predicted change associated with the prospective pair was evaluated as the aggregate of the training examples within the same chemotype. Using these three MMPA methods, QSAR-by-MMPA was carried out in the same manner as described previously by Warner.¹⁴ In these smaller data sets, no statistical testing of the property change for each aggregated transformation was possible, as the low number of observations for each pair would eliminate predictions for nearly all compounds.

Transformation QSAR. The protocol used for transformation QSAR is depicted in Figure 1. Descriptors for the training MMPs and prospective MMPs were generated by taking the difference between the individual compounds' descriptors. Using the pairwise descriptors and the known property change for the training pairs, a random forest model was trained to predict the property change for each pair. The random forest model was used to obtain an estimate of the property change for each prospective pair. Using that property change prediction, the property of each test set molecule is calculated in an identical manner to QSAR-by-MMPA. The

only deviation between transformation QSAR and the QSAR-by-MMPA methodology is that predicted activity changes are obtained from a QSAR model built on matched pairs, rather than the mean observed change of matched pairs.

One of the drawbacks of random forest QSAR models is that they lack the interpretability of simpler models; much more so than the simple QSAR-by-MMPA methodology. The randomForest library in R provides several functions to aid the user in interpreting how the QSAR model derives predictions, namely variable importance plots and partial dependence plots. The variable importance plot can be used to visualize the descriptors that contribute strongly to the model predictions. The ten most important descriptors for each activity change model, as well as the important descriptors, are conveyed in the Supporting Information. In each model, the most important descriptor is generally the assay measurement of the training set molecule. Our preliminary work led us to include the property value of the training set compound in the set of descriptors for the model since the measured values may be affected by the dynamic range of each assay. For example, if a compound has extremely low permeability, it is unlikely a transformation will measurably decrease permeability. Likewise, if a compound is at the upper end of the assay's measurable range, it is unlikely a transformation will be identified that measurably increases permeability. This effect appears to be largest for potency end point assays (Factor Xa and Pf-3D7 here) and was measured and confirmed using the partial dependence plot function of the randomForest package. (Supporting Information). The logP model was the exception to this observation, with the training logP value not appearing in the variable importance plot's returned values. In our statistics reported for the logP model, the training set compound's logP value and the ΔSlogP descriptor were excluded during the model training process. An additional model including these descriptors was trained, and exhibits identical performance, although the most important descriptor was, as anticipated, MOE's ΔSlogP descriptor. The partial dependence plot for the logP renin data set was determined using that model. The low magnitude of the logP model response (y-axis) indicates that there is very little dependence on the starting logP value.

Table 1. Data Set Information and Random Forest QSAR Model Performance

	no. comps	SD (exp)	% var. explained	RMSE	R ²	avg. Tan. sim ^a	avg. Manhattan dist ^a
3D7							
training	174	0.85	56	0.56	0.57	0.35	95
test	83	0.75	55	0.5	0.55		
Vdss							
training	445	0.63	58	0.41	0.58	0.04	152
test	219	0.65	56	0.44	0.56		
Factor Xa							
training	1414	1.44	70	0.78	0.71	0.07	168
test	606	1.41	69	0.79	0.7		
logP							
training	200	0.95	92	0.26	0.93	0.32	110
test	135	1.11	92	0.31	0.93		

^aThe Tanimoto similarity and Manhattan distance is the average value of the matrix of Tanimoto similarity values or Manhattan distances.

RESULTS AND DISCUSSION

We have included the standard QSAR models obtained from the random forest algorithm for each of the data sets to serve as a benchmark for the QSAR-by-MMPA and transformation QSAR methods. These QSAR models represent the reference performance for standard QSAR modeling techniques with the descriptors provided, and since we are comparing the random forest algorithm's ability to predict experimental end points to the same algorithm's ability to predict the difference between experimental end points for pairs of compounds, we believe that this metric will give an appropriate baseline of performance to compare to. Performance metrics for each of the standard QSAR models are presented in Table 1 below. The random forest algorithm performs well for the chosen data sets; the models account for between 50 and 60% of the variance for the 3D7 and Vdss data sets, 70% of the variance for the Factor Xa data set, and 95% of the variance for the logP set. As the chosen descriptor set is demonstrated to yield an acceptable random forest model, we would anticipate that the difference between descriptors should be able to encode structural changes for pairs of molecules. For reference, the ten most important variables for each of the random forest models, identified using the internal variable importance measure in the randomForest package, are reported in the Supporting Information. We will utilize these metrics to evaluate the performance of prospective predictions using QSAR-by-MMPA and transformation QSAR.

The performance of QSAR-by-MMPA varies significantly between data sets. (Table 2) The logP model presents an ideal case for transferability of SAR, as the success of additive models for predicting logP is well-known.²⁵ Using global MMPs, predictions returned for the logP data set possess an R² equal to 1.0 and an RMSE of 0.02. (Table 2). In data sets that have nonadditive SAR, such as potency assays where binding interactions with a target are important (e.g., 3D7 and Factor Xa), there will be many matched pairs but the effect of the chemical transformations may not be transferable between scaffolds. For the Factor Xa data set, the accuracy of predicted IC₅₀ values is marginally higher using global application of MMPs (RMSE = 0.69) than for the standard QSAR model (RMSE = 0.71) due to the presence of a larger number of similar transformations. In the 3D7 data set however, prediction accuracy is significantly deteriorated; ED₅₀ values obtained using direct application of MMPs possess an RMSE of 0.68, compared to standard QSAR's RMSE of only 0.40. In these two data sets, the local environment and substitution position should play an important role in changes in the

Table 2. Performance Metrics for the Various QSAR Models on the Subset of Test Set Compounds Covered by the Matched Pairs Protocols

	no. compounds	% var explained ^a	RMSE	R ²
3D7				
std QSAR	26	66	0.40	0.67
global MMPs		25	0.57	0.35
local MMPs		—4	0.68	0.20
chemo. MMPs		23	0.57	0.35
MMP-QSAR		65	0.39	0.65
Vdss				
std QSAR	13	62	0.4	0.68
global MMPs		81	0.68	0.37
local MMPs				
chemo. MMPs				
MMP-QSAR		83	0.46	0.87
Factor Xa				
std QSAR	163	73	0.71	0.74
global MMPs		75	0.69	0.75
local MMPs		74	0.70	0.74
chemo. MMPs		68	0.78	0.69
MMP-QSAR		90	0.44	0.90
logP				
std QSAR	101	96	0.19	0.96
global MMPs		100	0.03	1.0
local MMPs		100	0.004	1.0
chemo. MMPs		100	0.05	1.0
MMP-QSAR		100	0.03	1.0

^aVariance explained as a percentage of the variance in the experimental assay value.

potency measurement. In chemically diverse data sets, especially those corresponding to potency where there is a specific binding pocket, global transferability of SAR should not be assumed. Encoding information of the proper context of each transformation should significantly improve the prediction accuracy.

Inclusion of the radial environment in the matched pairs methodology adds an additional filter to the application of chemical transformations as both the transformation and radial environment must match. This leads to a decrease in the number of test set compounds with associated predictions in the radial MMP method. In the case of the highly homologous series comprising the logP (renin) data set, the decrease is minimal (only two compounds). The 3D7 and Factor Xa data

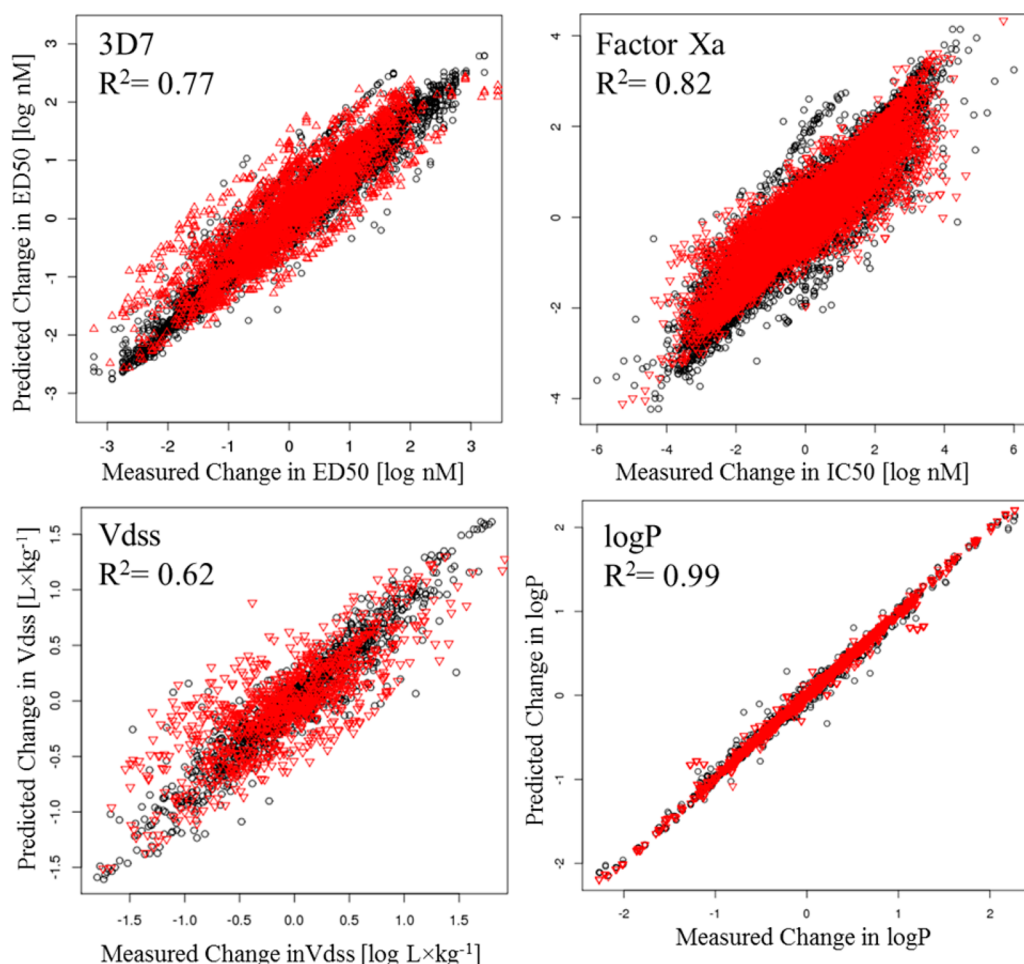


Figure 2. Predicted activity change from training MMPs (black points) and prospective MMPs (red triangles) from transformation QSAR. Correlation coefficients are shown for the prospective predictions in each data set.

sets are significantly impacted by the radial environment filtering of matched pairs: the number of 3D7 predictions drops from 53 to 37, and the number of Factor Xa predictions drops from 424 to 278. Interestingly, incorporation of the radial environment does not result in a notable increase in prediction accuracy in any of the investigated data sets and, in the case of the 3D7 and Factor Xa data sets, actually leads to a decrease in prediction accuracy. The predictions obtained in the 3D7 data set using radial MMPs are worse than random. (Table 2) However, this is likely an artifact of the data sets chosen for this study; the 3D7 and Factor Xa data sets correspond to potency end points, and the chemical environment on which the transformation occurs is likely less important than how that transformation would impact the molecule within the target binding site. On the other hand, logP is an additive property, and therefore, the local environment would have a minimal impact on logP. The RMSE of logP predictions using radial MMPs is slightly lower than global MMP predictions, although the accuracy is already so high that the difference is not significant. It is likely that within larger ADME data sets the effect of including the local environment would be more pronounced. The chemotype-based filtering of matched pairs also does not result in a performance gain relative to the global matched pairs method. (Table 2)

For the 3D7 and Factor Xa data sets, it is likely that the proper answer for the matched pairs calculation is to use 3D matched pairs identified by overlaying the 3D structure and

requiring that the transformations occur in the same region of space. Studies on 3D-based matched pairs have been reported recently,²⁶ and the VAMPPIRE database¹² has been developed to provide a repository for 3D matched pairs. However, the identification of 3D matched pairs is hampered by computational complexity and absence of experimental binding modes to confirm the methods, and its application is less common than the topological-based matched pairs analyses.

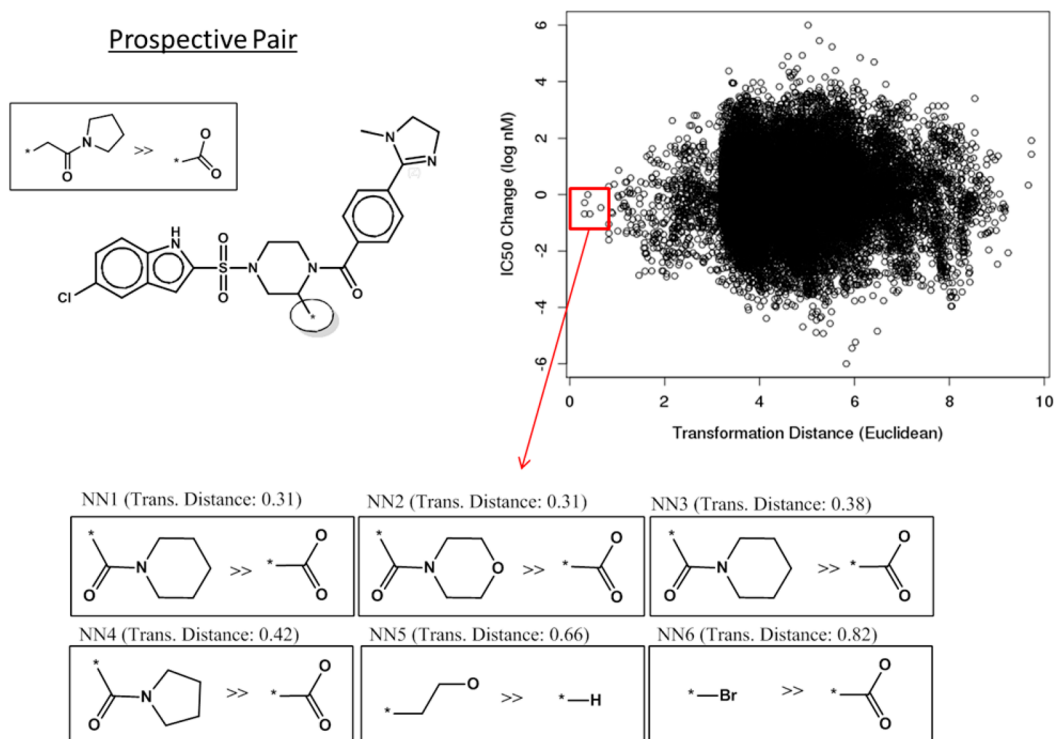
In each of the four data sets treated with QSAR-by-MMPA, predictions are not obtained for a fraction of the test set. These missing predictions are the result of either of two causes: the prospective compound does not form pairs with training set molecules or the prospective compound only forms pairs that correspond to previously unobserved transformations. For the data sets investigated in this study, the second cause is the limiting factor. In the Vdss data set, 101 test set molecules form pairs with the training set, but only 13 compounds form pairs that correspond to known transformations in the global MMP model. In the 3D7 and Factor Xa data sets, roughly one-third of the prospective compounds lack predictions. Incorporation of a transformation's context restricts the number of predictions further. We will note, once again, that no statistical filtering of the transformations was carried out in our QSAR-by-MMPA analysis; if the training set transformations were to be filtered by count/standard deviation, either no predictions would be made for the studied data sets, or predictions for a few prospective compounds would be obtained. The ability to

Table 3. Statistical Parameters for Model Performance of the Model to Predict Property Change, the Calculated Property of the Test Set Molecule for Each Pair, and the Mean Calculated Property for Each Test Set Molecule

	stage 1 (activity change)			stage 2 (per-pair test set activity) ^a			stage 3 (test set activity) ^b		
	R ²	RMSE	% var	R ²	RMSE	% var	R ²	RMSE	% var
3D7	0.77	0.48	77	0.51	0.48	51	0.62	0.48	62
Vdss	0.62	0.38	62	0.5	0.38	50	0.7	0.37	69
logP	0.99	0.05	99	1	0.05	99	1	0.03	100
Factor Xa	0.82	0.59	79	0.8	0.59	79	0.88	0.52	86

^aTest set activity is calculated using the training set molecule's activity plus the predicted activity change for the pair to determine the test set molecule's activity for each pair as depicted in stage two of Figure 1. The calculated activity of each test set compound may be present multiple times.

^bThe aggregated model takes each of the calculated test set activities for a test set molecule and calculates the mean activity for the test set compound as depicted in stage three of Figure 1.

**Figure 3.** Previously unobserved transformation in the Factor Xa test set and its nearest neighbors in the transformation QSAR training set determined by descriptor-based Euclidean distance (distance presented in parentheses).

reinforce the activity change associated with chemical transformations using similar, although not identical, transformations would be a beneficial addition to the matched pairs methodology and may allow for the method to be applied to more diverse data sets such as the Vdss set presented here, where the frequency of each specific chemical transformation is low.

Transformation QSAR was applied to each of the four ChEMBL data sets. Unlike the QSAR-by-MMPA procedure described by Warner, transformation QSAR obtains a prediction of the activity change for a matched pair using a random forest QSAR model trained using chemical transformations from training set matched pairs. Due to this procedural change, a prediction of the activity change is obtained for any prospective pair, regardless of whether or not that pair corresponds to a chemical transformation observed in the training set. The ability of transformation QSAR to make accurate predictions depends on the quality of the QSAR model to predict the activity change associated with each prospective pair, referred to as stage one predictions in Figure

1. Graphic representation of the training and test set predictions from each transformation QSAR model are presented in Figure 2, and the performance metrics are presented in Table 3. In each data set, the accuracy of the activity change predictions is comparable to or better than the standard QSAR model's overall performance on the activity itself. It should be noted that these activity change predictions include chemical transformations that were observed in the training set, as well as chemical transformations that were only seen in the prospective pairs.

We have proposed that the use of pairwise descriptors in the transformation QSAR model allows for clustering of similar transformations, reinforcing the measured activity change associated with infrequently observed transformations, as well as allowing for predictions on novel chemical transformations. This behavior is demonstrated for a prospective chemical transformation from the Factor Xa data set in Figure 3. This prospective MMP, comprised by the compounds ChEMBL61950 and ChEMBL62279, corresponds to a chemical transformation of the pyrrolidine amide to a

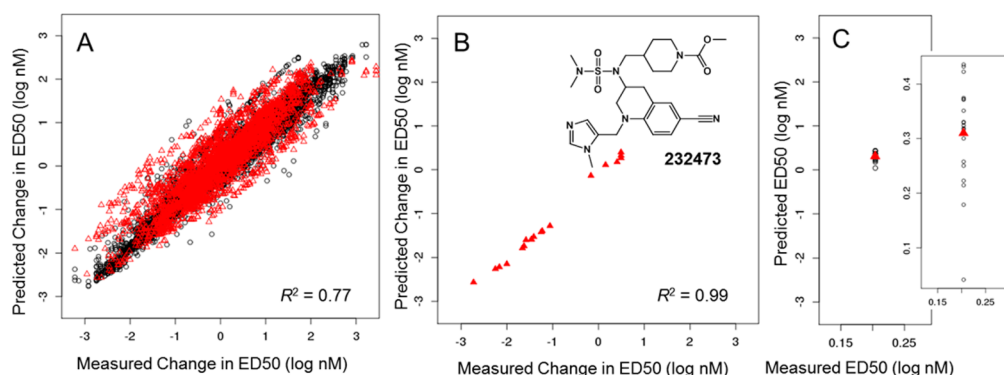


Figure 4. (A) Measured versus predicted change in ED_{50} for the transformation QSAR model of the 3D7 data set with training pairs depicted in black and test pairs in red. (B) Predicted versus measured change in ED_{50} for the 19 pairs involving the test set molecule CHEMBL232473. (C) Calculated ED_{50} of CHEMBL232473 from each of the 19 pairs versus the measured ED_{50} value, with the mean of the 24 log values depicted as a red square.

carboxylate group which was not observed within the training set. By taking the 10 most important descriptors contributing to the random forest model (identified using the variable importance function in R's implementation of randomForest and mean-shifted/variance-scaled), a descriptor-based Euclidean distance may be calculated between the prospective MMP and training MMPs to identify what the model considers to be the most similar chemical transformations from the training set. The six nearest neighbors in the Factor Xa transformation model for the pyrrolidine amide to carboxylate transformation and their Euclidean model-based distance are depicted in Figure 3. The first four transformations correspond to modifications of the cyclic portion of the amide, and all possess similar distances (0.31–0.41). There is then a gap in similarity before the fifth transformation (0.66) and sixth transformation (0.82), which bear little topological similarity to the query. This analysis was carried out across several prospective pairs, including a more common methyl to hydrogen transformation on a carboxylate group, for which 6 of the first 10 training set transformations corresponded to the identical methyl ester to carboxylate transformation, and the remaining pairs corresponded to transformation of a neutral R-group to a carboxylate group. It is clear from these analyses that the model is efficiently grouping similar chemical transformations based on their pairwise descriptors and, therefore, is able to predict the change associated with chemically novel transformations based on similarity to previous occurrences.

The predictions that are obtained using the transformation QSAR model correspond to changes in property/activity. Using simple arithmetic, these activity changes can be combined with the known (training set) value for each prospective pair to obtain an estimate of the novel compound's activity for each MMP (stage two predictions). In stage three, multiple property predictions for each novel compound are aggregated to a mean prediction of that compound's activity. This process is depicted in Figure 4 for a prospective compound in the 3D7 data set, CHEMBL232473. The QSAR model of predicted activity changes in the 3D7 data set performs extremely well ($R^2 = 0.77$) for the 2473 prospective pairs in the 3D7 data set. Out of those 2473 pairs, the prospective compound, CHEMBL232473, forms a total of 19 pairs with the training set; the predicted ΔED_{50} for each of these pairs from the QSAR model is depicted in Figure 4B against the measured ΔED_{50} values. The high correlation ($R^2 = 0.99$) makes it clear that the model is internally consistent for the molecule and is able to accurately

predict the activity change for multiple transformations originating from the prospective compound. This high correlation is observed for the majority of compounds in the data set using a linear regression to the predictions for each test set compound. The second stage of prediction—addition of the training set ED_{50} value to the predicted property change for each pair—results in 19 predictions of CHEMBL232473's ED_{50} (Figure 4C). The mean ED_{50} value from those 19 individual ED_{50} predictions (stage 3) is 0.31 [log nM] depicted as a red triangle in Figure 4C. The measured ED_{50} value for CHEMBL232473 is 0.20 log units.

Predictions for the activity of prospective compounds across each of the four data sets were obtained using the transformation QSAR protocol. The performance of individual stages is presented in Table 3. Transformation QSAR is compared to the previous QSAR-by-MMPA methodology in Table 2 for the smaller subset of compounds returned using the QSAR-by-MMPA methodology, and the performance of transformation QSAR relative to the standard random forest QSAR model is presented in Table 4. Relative to QSAR-by-

Table 4. Performance Metrics of Transformation QSAR and the Standard QSAR (Random Forest) Model

	3D7	Vdss	Factor Xa	logP
Transformation QSAR				
var explained (%)	62	69	100	86
RMSE	0.47	0.37	0.03	0.52
R^2	0.62	0.70	1.0	0.88
no. compounds	80	101	133	580
Standard QSAR RF				
var explained (%)	59	61	95	70
RMSE	0.46	0.41	0.23	0.76
R^2	0.59	0.62	0.96	0.71

MMPA, transformation QSAR exhibits an increase in prediction accuracy for each of the data sets, with the exception of the logP data set where the accuracy is equivalent. (Table 2) This is indicative of the performance relative to the current QSAR-by-MMPA methods. But this performance is only for the relatively small subset of compounds returned by the contextual MMP methods. The true indication of the performance of transformation QSAR relative is relative to the benchmark of standard QSAR modeling techniques. The measures of model performance for the subset of compounds

that share predictions in transformation QSAR and the standard QSAR model are presented in Table 4. In terms of the variance explained in each data set, transformation QSAR is superior to the standard QSAR model; the differences range from only a slight increase in the variance explained in the predicted ED_{50} values (3D7) to a significant improvement (Factor Xa). The RMSE of the predictions also indicates increased accuracy in the transformation QSAR predictions in all data sets except 3D7. However, the 3D7 data set possesses a much larger fraction of qualified data (measured ED_{50} value greater than $5 \mu\text{M}$) than the other data sets studied herein. If the contributions of the out-of-range measurements are excluded, the RMSE of the transformation QSAR predictions (0.42) is slightly lower than the RMSE of the standard QSAR predictions (0.43). One should be careful when considering the activity change between an out-of-range measurement and an in-range measurement, as the associated activity change could, in reality, be significantly larger than the measurement would indicate.

CONCLUSIONS

The utility of matched pairs for prospective design of compounds was recently evaluated using a proprietary data set and was found to outperform conventional QSAR methods.¹⁴ The QSAR-by-MMPA protocol benefits from being extremely simple and interpretable. Our own attempts to apply QSAR-by-MMPA to data sets suggested that the occurrence of matched pairs within a set may be too sparse for widespread utility, a finding discussed in the original publication.¹⁴ This fact led us to evaluate the procedure on several public data sets obtained from ChEMBL. While QSAR-by-MMPA does exhibit excellent performance for several data sets studied herein, we observed that predictions are returned for only a fraction of the test set compounds, and the frequency of transformations obeys a power-law distribution; most chemical transformations are infrequently observed, with the majority of transformations being observed only one or two times and have little to no statistical significance. We have demonstrated that a QSAR model trained to predict the activity change of a matched pair (transformation QSAR) allows enhanced predictions for novel compounds, as well as estimating the activity change of previously unobserved chemical transformations via similarity to known transformations. This is a benefit that matched pairs implementations do not currently incorporate, although care must be used to consider the domain of applicability for the transformation model, just as in standard QSAR modeling. Using four ChEMBL data sets, we observed that transformation QSAR returns predictions of greater than or comparable accuracy to standard QSAR and exhibits superior performance to QSAR-by-MMPA for all cases except logP, for which the performance was equal between methods. We did find that attempts to include the local environment in the QSAR-by-MMPA methodology led to a decrease in performance. It is likely that for the target-binding data sets (3D7 and Factor Xa) the local environments used were not a proper representation of the true transformation context and a 3D representation that correlates with target-binding should be utilized at a significant increase in cost and complexity. Future efforts may focus on evaluating different representations of the chemical environment for both QSAR-by-MMPA and transformation QSAR methods.

The observed drawbacks to the transformation QSAR methodology are increased computational cost and a more complex implementation, and the fact that predictions will not be obtained for some fraction of the data set. Out of the four evaluated data sets shown here, and from our own experience with proprietary data sets, we have not observed an instance where the transformation QSAR performance is lower than standard QSAR modeling using a random forest algorithm. Due to the inherent complexity of the resultant random forest model, the transformation QSAR methodology does lose the interpretability of the straightforward QSAR-by-MMPA. However, we believe that the ability to leverage chemically similar transformations, and return a prediction for novel chemical transformations based on the transformation model are worthwhile gains.

ASSOCIATED CONTENT

Supporting Information

Method protocol schemes, ChEMBL data set retrieval information, list of calculated descriptors, and QSAR model interpretation. This material is available free of charge via the Internet at <http://pubs.acs.org>.

AUTHOR INFORMATION

Corresponding Author

*E-mail: clayton.springer@novartis.com.

Notes

The authors declare no competing financial interest.

REFERENCES

- (1) Bajorath, J. Modeling of activity landscapes for drug discovery. *Exp. Opin. Drug Discov.* **2012**, *7*, 463–473.
- (2) Guha, R.; van Drie, J. H. Structure–activity landscape index: identifying and quantifying activity cliffs. *J. Chem. Inf. Model.* **2008**, *48*, 646–658.
- (3) Heikamp, K.; Hu, X.; Yan, A.; Bajorath, J. Prediction of activity cliffs using support vector machines. *J. Chem. Inf. Model.* **2012**, *52*, 2354–2365.
- (4) Seebeck, B.; Wagener, M.; Rarey, M. From activity cliffs to target-specific scoring models and pharmacophore hypotheses. *ChemMedChem* **2011**, *6*, 1630–1639.
- (5) Hu, Y.; Bajorath, J. Exploration of 3D activity cliffs on the basis of compound binding modes and comparison of 2D and 3D cliffs. *J. Chem. Inf. Model.* **2012**, *52*, 670–677.
- (6) Guha, R. Exploring uncharted territories: predicting activity cliffs in structure-activity landscapes. *J. Chem. Inf. Model.* **2012**, *52*, 2181–2191.
- (7) Sheridan, R. P.; Hunt, P.; Culberson, J. C. Molecular transformations as a way of finding and exploiting consistent local QSAR. *J. Chem. Inf. Model.* **2006**, *46*, 180–192.
- (8) Wassermann, A. M.; Dimova, D.; Iyer, P.; Bajorath, J. Advances in Computational Medicinal Chemistry: Matched Molecular Pair Analysis. *Drug Dev. Res.* **2012**, *73*, 518–527.
- (9) 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.
- (10) Papadatos, G.; Alkarouri, M.; Gillet, V. J.; Willett, P.; Kadirkamanathan, V.; Luscombe, C. N.; Bravi, G.; Richmond, N. J.; Pickett, S. D.; Hussain, J.; Pritchard, J. M.; Cooper, A. W.; Macdonald, S. J. Lead optimization using matched molecular pairs: inclusion of contextual information for enhanced prediction of HERG inhibition, solubility, and lipophilicity. *J. Chem. Inf. Model.* **2010**, *50*, 1872–1886.
- (11) Wirth, M.; Zoete, V.; Michielin, O.; Sauer, W. H. SwissBioisostere: a database of molecular replacements for ligand design. *Nucleic Acids Res.* **2013**, *41*, D1137–D1143.

- (12) Weber, J.; Achenbach, J.; Moser, D.; Proschak, E. VAMMPIRE: A Matched Molecular Pairs Database for Structure-Based Drug Design and Optimization. *J. Med. Chem.* **2013**, *56*, 5203–5207.
- (13) Warner, D. J.; Griffen, E. J.; St-Gallay, S. A. WizePairZ: a novel algorithm to identify, encode, and exploit matched molecular pairs with unspecified cores in medicinal chemistry. *J. Chem. Inf. Model.* **2010**, *50*, 1350–1357.
- (14) Warner, D. J.; Bridgland-Taylor, M. H.; Sefton, C. E.; Wood, D. J. Prospective Prediction of Antitarget Activity by Matched Molecular Pairs Analysis. *Mol. Inform.* **2012**, *31*, 365–368.
- (15) Breiman, L. Random Forests. *Machine Learning* **2001**, *45*, 5–32.
- (16) Coteron, J. M.; Marco, M.; Esquivias, J.; Deng, X.; White, K. L.; White, J.; Koltun, M.; Mazouni, el F.; Kokkonda, S.; Katneni, K.; Bhamidipati, R.; Shackleford, D. M.; Angulo-Barturen, I.; Ferrer, S. B.; Jiménez-Díaz, M. B.; Gamo, F.; Goldsmith, E. J.; Charman, W. N.; Bathurst, I.; Floyd, D.; Matthews, D.; Burrows, J. N.; Rathod, P. K.; Charman, S. A.; Phillips, M. A. Structure-guided lead optimization of triazolopyrimidine-ring substituents identifies potent Plasmodium falciparum dihydroorotate dehydrogenase inhibitors with clinical candidate potential. *J. Med. Chem.* **2011**, *54*, 5540–5561.
- (17) Obach, R. S.; Lombardo, F.; Waters, N. J. Trend Analysis of a Database of Intravenous Pharmacokinetic Parameters in Humans for 670 Drug Compounds. *Drug Metab. Dispos.* **2008**, *36*, 1385–1405.
- (18) Subramanian, G.; Rao, S. N. An integrated computational workflow for efficient and quantitative modeling of renin inhibitors. *Bioorg. Med. Chem.* **2012**, *20*, 851–858.
- (19) Morgan, H. The Generation of a Unique Machine Description for Chemical Structures-A Technique Developed at Chemical Abstracts Service. *J. Chem. Doc.* **1965**, *5*, 107–113.
- (20) Rogers, D.; Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **2010**, *50*, 742–754.
- (21) Landrum, G. *RDKit Documentation*; 2012; pp 1–41; <http://www.rdkit.org/GettingStartedInPython.pdf>.
- (22) Liaw, A.; Wiener, M. Classification and Regression by randomForest. *R News* **2002**, *2*, 18–22.
- (23) Ho, T. K. Random Decision Forests. *Proceedings of the 3rd International Conference on Document Analysis and Recognition*, Montreal, Aug 14–16, 1995; pp 278–282.
- (24) Sheridan, R. P. Three useful dimensions for domain applicability in QSAR models using random forest. *J. Chem. Inf. Model.* **2012**, *52*, 814–823.
- (25) Wang, R.; Fu, Y.; Lai, L. A New Atom-Additive Method for Calculating Partition Coefficients. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 615–621.
- (26) Posy, S. L.; Claus, B. L.; Pokross, M. E.; Johnson, S. R. 3D Matched Pairs: Integrating Ligand- and Structure-Based Knowledge for Ligand Design and Receptor Annotation. *J. Chem. Inf. Model.* **2013**, *53*, 1576–1588.