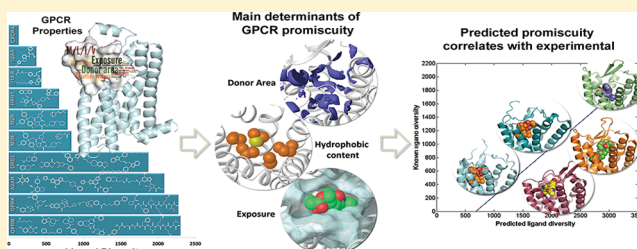


Predicting GPCR Promiscuity Using Binding Site Features

Anat Levit,^{†,‡} Thijs Beuming,[§] Goran Krilov,[§] Woody Sherman,[§] and Masha Y. Niv^{*,†,‡}[†]Institute of Biochemistry, Food Science and Nutrition, Robert H. Smith Faculty of Agriculture Food and Environment, The Hebrew University, Rehovot 76100, Israel[‡]Fritz Haber Center for Molecular Dynamics, The Hebrew University, Jerusalem 91904, Israel[§]Schrodinger Inc., 120 West Forty-Fifth Street, 17th Floor, New York, New York 10036, United States

Supporting Information

ABSTRACT: G protein-coupled receptors (GPCRs) represent a large family of signaling proteins that includes many therapeutic targets. GPCR ligands include odorants, tastants, and neurotransmitters and vary in size and properties. Dramatic chemical diversity may occur even among ligands of the same receptor. Our goal is to unravel the structural and chemical features that determine GPCRs' promiscuity toward their ligands. We perform statistical analysis using more than 30 descriptors related to the sequence, physicochemical, structural, and energetic properties of the GPCR binding sites—we find that the chemical variability of antagonists significantly correlates with the binding site hydrophobicity and anticorrelates with the number of hydrogen bond donors in the binding site. The number of disulfide bridges in the extracellular region of a receptor anticorrelates with the range of molecular weights of its antagonists, highlighting the role of the entrance pathway in determining the size selectivity for GPCR antagonists. The predictive capability of the model is successfully validated using a separate set of GPCRs, using either X-ray structures or homology models.



INTRODUCTION

G protein-coupled receptors (GPCRs) form the largest family of cell surface receptors in the human genome.¹ They are key signaling molecules and are the targets of over 30% of currently approved and marketed drugs.² Endogenous agonists of GPCRs include bioamines, nucleotides, neurotransmitters, peptides, and many other chemical stimuli. Some GPCRs are narrowly tuned toward their agonists, such as the pheromone receptors³ and subsets of the olfactory receptor subfamily,⁴ while others have a broad receptive range. Striking examples of broadly tuned receptors include some olfactory receptors^{5,6} and several bitter taste receptors,^{7,8} in which a single receptor recognizes a broad range of ligands.

A recent systematic analysis of ligand-contacting residues in the transmembrane (TM) ligand-binding pocket of GPCR X-ray structures has revealed that, except for the C-X-C chemokine type 4 receptor (CXCR4) and the neurotensin 1 receptor (NTSR1), all other X-ray structures of Family A GPCRs share similarity in ligand-contacting residues, with topologically equivalent positions in TM3, TM6, and TM7 typically contacting the ligand in nearly all receptors.⁹ Variation in the amino acids occupying these positions accounts for ligand specificity in different receptors. Furthermore, individual GPCRs may accommodate diverse chemical matter by utilizing different subsets of binding site residues,^{8,10} as well as different types of interactions (i.e., polar vs nonpolar).⁸ These differences can be identified by crystallography^{11,12} or via a combination of modeling and mutagenesis studies^{8,10,13} on a

case-by-case basis. Here, we attempt to address whether the receptive range of a receptor can be predicted based on the physicochemical properties of the binding site.

Recently developed methods, such as ligand's eye-view of protein similarity,^{14,15} as well as chemoproteomic¹⁶ or chemogenomic approaches,^{17,18} analyze and predict the relationship between proteins and the ligands they bind based on either ligand similarities or both protein and ligand similarity information. For example, agonist-specific regions were shown to concentrate between TMs 2, 3, and the second extracellular loop (ECL2), while antagonist-specific regions are located at the top of TMs 5 and 6.¹⁸ A ligand-based view of promiscuity can suggest ligand properties that determine its polypharmacology (the number of targets with which the ligand interacts). Studies have found that the most promiscuous drugs tend to be highly hydrophobic (clog P ≥ 3). The relation between ligand size (in terms of molecular weight (MW)) and its promiscuity has also been studied, but no consensus was reached.¹⁹ In some cases, an inverse correlation between mean MW and ligand promiscuity toward targets was found, while another study showed that within a given clog P range, promiscuity decreases with increasing ligand size. A recent study²⁰ on a large set of over 40 000 molecules for which at least three measured affinities (pXC₅₀ ≥ 6) were available in ChEMBL did not find

Received: September 25, 2013

Published: December 16, 2013

any significant correlation between MW and ligand promiscuity.

Our goal is not to study what underlies the promiscuity of a small molecule. Rather, we aim to understand which molecular and structural features of a *receptor* determine the chemical diversity of ligands that it is capable of binding. This ability may also be termed receptive range, multispecificity, or promiscuity of a receptor. Both binding site promiscuity and ligand promiscuity are common in nature. A recent study indicates that more than 1/3 of pockets in a data set of 20 000 ligand binding pockets, interact with multiple, chemically different ligands.²¹ While selectivity or promiscuity of GPCRs in general,^{19,22} and of specific GPCR subtypes in particular,^{23,24} is often mentioned in the literature, a comprehensive metric for quantification of GPCR promiscuity toward their ligands is missing. Here, GPCR promiscuity is assessed by determining the diversity of known antagonists. In turn, antagonist diversity is quantified using several descriptors, including the size of the known ligand set, the number of unique scaffolds within this set and the ranges of physicochemical properties of ligands in the set. The abundance of chemical data for GPCR antagonists²⁵ and the recent increase in the number of X-ray structures of GPCRs,^{9,26,27} provide an unprecedented opportunity for unraveling molecular details of GPCR architecture, which can offer insights into the fundamental features underlying GPCR selectivity. We use this information to develop an approach to predict GPCR selectivity based only on target information.

RESULTS

GPCRs with Determined Structures and Their Corresponding Ligand Sets. Our training set consisted of 10 Family A GPCRs for which X-ray structures had been solved experimentally at the start of this study, and their corresponding ligands sets, retrieved from ChEMBL, as summarized in Table 1.

Table 1. GPCRs in the Training and the Validation Sets, along with X-Ray Structure PDB Codes and ChEMBL Ligand Set Sizes (1 and 10 μ M Activity Cutoff)

receptor	PDB code	ChEMBL ligand set size	
		$\text{IC}_{50}/K_i \leq 1$ μM	$\text{IC}_{50}/K_i \leq 10$ μM
training set			
μ -opioid receptor (OPRM)	4DKL	2300	3103
κ -opioid receptor (OPRK)	4DJH	2273	3220
adenosine 2A receptor (A2AR)	3EML	2084	3107
dopamine D3 receptor (DRD3)	3PBL	1872	2335
muscarinic M2 receptor (M2R)	3UON	786	1024
muscarinic M3 receptor (M3R)	4DAJ	836	1010
histamine H1 receptor (HRH1)	3RZE	669	881
β 1-adrenergic receptor (β 1AR)	2VT4	366	652
β 2-adrenergic receptor (β 2AR)	2RH1	434	661
C–X–C chemokine receptor type 4 (CXCR4)	3ODU	130	184
validation set			
δ -opioid receptor (OPRD)	4EJ4	1941	
N/OFQ opioid receptor (OPRX)	4EA3	1048	
serotonin receptor 1B (SHT1B)	4IAR	762	
serotonin receptor 2B (SHT2B)	4IB4	494	
protease-activated receptor 1 (PAR1)	3VW7	393	
smoothened receptor (SMO)	4JKV	289	

Additional structures of GPCRs have been published during the course of this study, and their structures were used for model validation. We excluded the S1P1 and the NTSR1 receptors from both training and validation sets, as there are only few recently (year 2000 and later) reported ligands in the ChEMBL database, possibly reflecting the fact that their chemical space has not yet been well studied. Family B receptors were not included since preliminary tests showed significant differences from the Family A binding sites.

Variability of GPCRs Ligands Sets. First, we set to assess the diversity of ligands in the different data sets and to explore the relevance of the number of ligands or the number of unique scaffolds in a set as measures of ligand set variability. To this end, we calculated statistics (mean, median, range, standard deviation, and variance) of a default set of physicochemical ligand descriptors for the antagonists sets in our training set (Table 1), all within 1 μ M activity cutoff. The total number of molecules in a set, number of unique scaffolds (NUS), and number of clusters (using either radial or dendritic fingerprints) for each of the 10-ligand sets were also computed ((Supporting Information Figure S1A). Differences in the distribution of MW and ALogP are observed between compounds belonging to different targets, reflecting differences in the nature of the binding sites. The chemical space distribution of the entire 1 μ M activity cutoff set ($n = 11\,750$), as defined by MW and ALogP, is shown in Figure S1B. The compounds included in our data set cover a large chemical space and are not concentrated, for example, only within the boundaries defined by Lipinski rules for drug-like molecules (as indicated by the dotted lines in Figure S1B), suggesting that the data set is diverse and can be used to explore GPCR promiscuity.

Statistical parameters of the ligand physicochemical descriptors were included in a multivariate analysis, to quantify the relationships among them (Figure 1). A correlation matrix including all calculated statistics for ligand properties for the 1 μ M ligand sets is shown in Supporting Information Figure S2A and Table S1 lists the statistically significant pairwise correlations observed. The results are described below and are consistent when repeating the analysis for ligands with up to 10 μ M activity cutoff (data not shown) and when including over 320 additional molecular descriptors of the compounds in the analysis (Figure S2B).

The correlation (r) between the number of molecules in the set and NUS is very high (0.98; p -value < 0.0001). This means, that where there are more ligands, there are also more unique scaffolds and suggests that there is no major overrepresentation of a few scaffolds within a particular ligand set. Furthermore, the *ranges* of almost all calculated physicochemical descriptors (Supporting Information Figure S2A and B) are in striking positive correlation with NUS and with number of molecules. Thus, the larger the number of unique scaffolds or the number of known ligands, the higher the chemical variability, as manifested in properties ranges of the ligands the receptor can interact with.

The medians (Figure 1) and means (Figure S2A) show much lower levels of correlation with NUS, with a few exceptions: the highest positive correlation (r) is observed between NUS and median value of chiral centers in the scaffold (0.65) and with ring count median (0.57). The highest negative correlation is observed between NUS and the rotatable bond median (−0.71), HB donors median (−0.52), and centralization median (−0.515). This means that the more promiscuous receptors tend to have antagonists with more chiral centers and

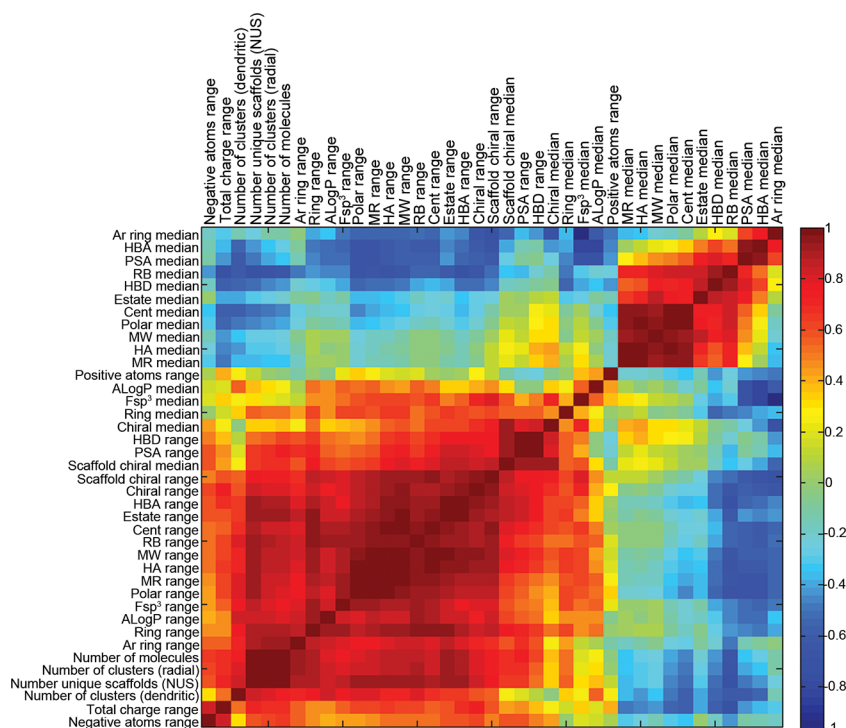


Figure 1. Multivariate correlations between ligand set descriptors obtained from the 1 μM ligand sets. The degree of correlation ranges from red (+1) to blue (−1). Descriptors were clustered together based on correlation values. Median values of total charge, negative atom count, and positive atom count were excluded from the analysis due to lack of variability between the different ligand sets.

rings (being bulky or “fat”) and do not tend to be flexible and polar. This is in agreement with the “fat or flat” idea:²⁸ aromaticity relates to “flatness”, while sp^3 -hybridized (tetrahedral) carbon atom fraction (Fsp^3) relates to fatness or “three-dimensionality” of the molecule.

The anticorrelation between fatness and flatness is also apparent from the opposing positions of these properties along the second principle component obtained in principal component analysis (PCA) of ligand sets molecular properties: The correlations between median values of ligand descriptors of the different targets in our 1 μM training set were subjected to PCA to obtain a global view of the relationship between these descriptors. The 10 GPCR ligand sets were then projected onto the axes of the first two principle components (PCs) (Figure 2). The first two PCs explain 74% of the variability in the data. PC1 is composed mostly of ligand flexibility and polarity descriptors in the positive direction (i.e., number of rotatable bonds and number of HB donors) and of data set size-related descriptors in the negative direction (i.e., NUS, number of clusters, and number of molecules). This means that ligands sets where ligands have more rotatable bonds and hydrogen bond donors typically decompose into fewer unique scaffolds using the Bemis–Murcko protocol. PC2 is composed mostly of the chirality descriptors and flexibility (Fsp^3) in one direction and of aromatic ring count in the opposite direction. The promiscuous opioid receptors map in the region related to high chirality (or fatness). Interestingly, fatness of an individual compound was shown to correlate with its selectivity toward targets.²⁹

The medians and, to a lesser degree, means of polarity descriptors (HB donors, acceptors, etc.) are correlated among themselves and have modest anticorrelation with NUS. Interestingly, the MW median and mean values, which are dominant descriptors of ligand sets variability, are *not* correlated

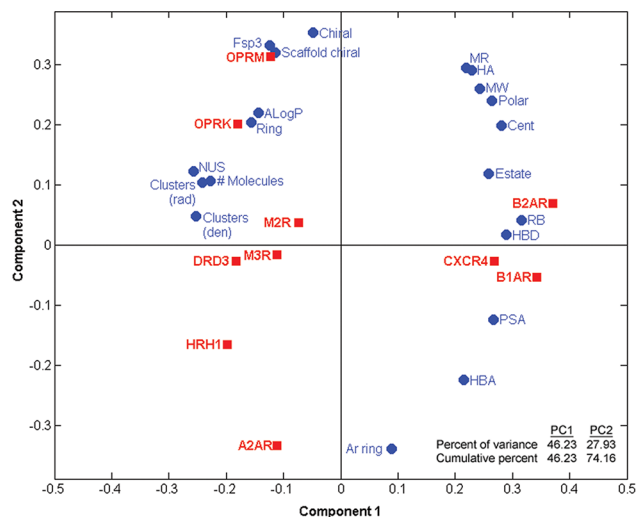


Figure 2. Principal component analysis on median values of 20 ligand descriptors for the 1 μM training set. Biplot of loading and scoring showing the variables (blue) composing the first two principal components and projection of the 10 different ligand sets used (red). The first two principal components account for 74% of the variance in the data. Ar ring—aromatic rings count; HBA—hydrogen bond acceptors; PSA—polar surface area; HBD—hydrogen bond donors; RB—rotatable bonds; Estate—electrotopological states; Cent—centralization; Polar—molecular polarizability; MW—molecular weight; HA—heavy atom count; MR—molecular refractivity; Chiral—chiral center count; Scaffold chiral—chiral center count in scaffolds; Ring—rings count; # molecules—number of molecules; Fsp^3 —fraction of sp^3 -hybridized (tetrahedral) carbon atoms out of total carbon count; NUS—number of unique scaffolds; clusters (rad)—number of clusters using radial FP; clusters (den)—number of clusters using dendritic FP.

with NUS. This finding is not trivial, since it might have been assumed that if large ligands can bind to a receptor, the large binding site may be accessible to many chemically dissimilar ligands.

Ranking the Targets According to Number of Unique Scaffolds (NUS). NUS, which presented high positive correlation with ranges of molecular properties within the examined GPCRs ligand sets, was chosen as a good indication of the diversity of ligands binding to a given GPCR. The ranking of the training set GPCRs based on their NUS, as well as ligand set size is shown in Supporting Information Figure S1A.

It is clear that the GPCRs studied here vary in terms of number and diversity of their ligands, with opioid receptors defined as the most promiscuous using this metric, and the CXCR4 receptor the least promiscuous. To determine receptor features that may be responsible for these differences, binding site descriptors for the GPCRs in the training set (Table 1) are calculated and analyzed next.

Statistical Analysis of Relationships between Ligand Set Diversity and Receptors Descriptors. *Characterizing the Receptor Structures.* The location and features of the orthosteric binding site of GPCRs has long been established from mutagenesis studies^{30–32} and validated in the numerous X-ray structures currently available (analyzed, e.g. in refs 9 and 33). The binding sites of the 10 training set X-ray structures (Table 1) were analyzed in terms of sequence, structure, solvation properties, and flexibility.

For the sequence analysis, the relative percentage of hydrophobic (M/L/I/V), aromatic (W/F/Y), polar (S/T/N/Q/H), basic (K/R), acidic (D/E), and small (A/C/G) binding site residue content was calculated for each of the receptors (binding site positions included in the analysis are listed in Supporting Information Table S2). Next, structural properties of the binding site were analyzed using SiteMap,³⁴ which calculates the size, volume, exposure, and other properties, that relate to the ability of regions on the surface of proteins to bind molecules (see Methods). SiteMap also calculates surfaces representing the hydrophobic, donor, acceptor, and hydrophilic regions of the binding site. The areas of these surfaces were used as descriptors as well.

The ability of ligands to differentially displace and retain specific water molecules solvating the protein binding sites may offer additional insights into the properties of the pocket. The GPCR orthosteric sites were analyzed with the WaterMap algorithm,³⁵ which computes the locations and thermodynamic properties of water molecules in protein binding sites. WaterMap has been used to reveal important aspects of binding sites relating to their ability to bind small molecules, both in globular proteins³⁶ and GPCRs.³⁷ In addition, WaterMap has been successfully used to assess binding selectivity.^{38,39} More recently, WaterMap has been used to understand the role of water networks in the hydrophobic effect and molecular recognition, including validation using isothermal titration calorimetry (ITC) data to assess the entropic and enthalpic binding contributions from the solvent.^{40,41} We hypothesized that thermodynamic parameters of water molecules in the binding site may be related to the binding site promiscuity. Values of enthalpy (ΔH), entropy ($-T\Delta S$), and free energy (ΔG) were estimated using WaterMap for each hydration site in the orthosteric binding site of each GPCR. We used the average values for entropy, enthalpy, and free energy,

as well as the total free energy of all hydration sites as receptor descriptors.

Flexibility of the binding site was estimated by the total number of accessible rotamers from a conformational search within either 5 or 10 kcal/mol of the native state using molecular mechanics with an implicit solvent model (see Methods). Using the number of accessible rotamers is one of many possible ways to account for binding site flexibility and has been used in other applications.^{42,43}

Binding sites were characterized also in terms of intrinsic selectivity ratio (ISR),⁴⁴ a dimensionless quantity that measures the degree to which native-like protein ligand interactions are energetically stabilized relative to other possible protein–ligand binding interactions (see Supporting Information Table S3 for ISR results). A higher value of ISR indicates higher predicted selectivity.

Due to the proximity of the second extracellular loop (ECL2) to the ligand binding site, as evident from GPCR structures, the length and secondary structure of ECL2 may influence the properties of the binding site. Indeed, this loop has been shown to be important for ligand binding in many GPCRs.⁹ The average secondary structure content based on backbone hydrogen bonds, and the total length of the loop were determined. In addition, a conserved Cys in ECL2 forms a disulfide bridge with TM3, which is found in all solved Family A X-ray structures, except for the S1P1 receptor. The total number of disulfide bridges in the extracellular region varies between different receptors and was also included among receptor descriptors.

In total, 32 different structural and sequence-related descriptors were generated for each receptor. Principal component analysis (PCA) suggests that no single descriptor can account for the majority of the observed diversity of the receptors (Supporting Information Figure S3A). This finding is supported by multivariate analysis (Supporting Information Figure S3B and Table S4). Since no single receptor descriptor is highly correlated (either positively or negatively) with all, or a majority of, other descriptors, all 32 receptor descriptors were included in the analysis that is described below.

Next, we aimed to identify those receptor descriptors that are most correlated with, and predictive of, receptor promiscuity (as represented, for example, by NUS). This was carried out in two stages: First, the stepwise regression was applied in an automated process of building a model by successively adding or removing variables based on the *t*-statistics of their estimated coefficients. Second, all variables (in our case, receptor descriptors) that significantly correlate with the analyzed measure of receptor promiscuity serve as input for the generation of a linear regression model for prediction of receptor promiscuity. This model is then applied to an independent set of receptors for validation.

Stage One: Stepwise Regression. Stepwise regression is a method for selecting a subset of predictive variables for a regression model, based on an automatic procedure of searching through different models with combinations of variables, and testing the predictive capabilities of the models. The result is a set of the most important variables (of the 32 receptors' variables used here) to predict the observable (in this case, the chemical diversity of the antagonists that bind to the GPCR, as represented by NUS or other diversity parameters). The outcome of this analysis is summarized in Table 2, showing that the percent of hydrophobic residues (M/L/I/V) in the binding site is the most influential factor on NUS *p*-value

Table 2. Results of Stepwise Regression Analysis Illustrating the Receptor Descriptors Independently Associated with NUS and Various Additional Ligand Descriptors

diversity measure	correlation with NUS (r)	receptor parameter	P value
number of unique scaffolds (NUS)	0.9762	% M/L/I/V	0.0076 ^a
		area donor	0.0277 ^a
		exposure	0.0048 ^a
number of molecules	0.9762	% M/L/I/V	0.0009 ^a
		area donor	0.0474 ^a
		exposure	0.0319 ^a
number of clusters ^b	0.9800	% W/F/Y	0.0683
		% M/L/I/V	0.0016 ^a
		area donor	0.0362 ^a
MW range	0.9312	disulfide bridges	0.0368 ^a
		hydrophilic	0.0143 ^a
		% A/C/G	0.0786
electrotopological state range	0.9363	% M/L/I/V	0.0115 ^a
		disulfide bridges	0.0317 ^a
		balance	0.0230 ^a

^aStatistically significant at the 0.05 level (i.e., 95% confidence). ^bThe number of clusters is based on hierarchical clustering of a Tanimoto similarity matrix derived from radial fingerprints.

0.0076), followed by the SiteMap descriptors for area of the donor region (p -value 0.0277) and exposure (p -value 0.0048).

To test the robustness of this unbiased analysis, we repeated it for four other ligand set descriptors that are strongly

correlated with NUS, namely, the diversity-related descriptors “number of molecules in a set” and “number of clusters” (as calculated using radial fingerprints, which have higher correlation with NUS than dendritic fingerprints—0.98 compared with 0.83) and the molecular size-related descriptors—ranges of “molecular weight” and of “electrotopological states”. The receptor descriptors that were found as most influential and statistically significant differ somewhat, depending on the analyzed ligand descriptor. In general, these include percent of M/L/I/V residues in the binding site, number of extracellular disulfide bridges, and several SiteMap descriptors (area donor, hydrophilic, balance, and exposure). Thus, the structural descriptors found to be correlated with receptor promiscuity are reasonably robust. Somewhat surprisingly, none of the WaterMap, flexibility, and ISR descriptors were found to correlate with any of the ligand-based promiscuity measures. It is possible that indeed there is no correlation between the properties that these descriptors attempt to describe or that the approximations used to develop these descriptors introduce too much noise for a signal to emerge.

Stage Two: Linear Regression Model of GPCR Promiscuity.

On the basis of the stepwise regression results above, the four most influential factors (% M/L/I/V residues in the binding site, number of disulfide bridges, SiteMap area donor, and SiteMap exposure) were entered into a Standard Least Squares fit analysis to develop a predictive model for receptor promiscuity. Either NUS, number of molecules, number of clusters (using radial fingerprints), electrotopological states range, or molecular weight range were used as descriptors of

Table 3. Regression Coefficients of Receptor Descriptors Maximally Affecting NUS and Various Additional Ligand Descriptors, As Determined by the Standard Least Squares Fit Analysis

diversity measure	receptor parameter	estimate (b) ^b	beta ^c	P value
number of unique scaffolds (NUS)	intercept	1365.37		0.0965
	% M/L/I/V	3711.6365	0.641496	0.0008 ^a
	area donor	−2.865893	−0.56276	0.0013 ^a
	exposure	3021.576	0.31258	0.0344 ^a
	disulfide bridges	−143.78	−0.12995	0.2752
number of molecules	intercept	950.23611		0.2511
	% M/L/I/V	4301.9155	0.784294	0.0006 ^a
	area donor	−1.943108	−0.40248	0.0104 ^a
	exposure	1851.2282	0.202012	0.1678
	disulfide bridges	−131.2362	−0.12512	0.3545
number of clusters ^d	intercept	656.68722		0.1170
	% M/L/I/V	1983.1373	0.765061	0.0007 ^a
	area donor	−1.014765	−0.44478	0.0070 ^a
	exposure	740.34228	0.170952	0.2316
	disulfide bridges	−88.90828	−0.17936	0.2045
MW range	intercept	687.34705		0.0238 ^a
	%M/L/I/V	957.79607	0.507617	0.0022 ^a
	area donor	−0.856493	−0.51573	0.0018 ^a
	exposure	1487.5971	0.471897	0.0068 ^a
	disulfide bridges	−106.2566	−0.29449	0.0369 ^a
electrotopological states range	intercept	124.66511		0.0178 ^a
	%M/L/I/V	218.66809	0.663442	0.0005 ^a
	area donor	−0.110723	−0.38167	0.0056 ^a
	exposure	209.79303	0.380985	0.0136 ^a
	disulfide bridges	−18.83734	−0.29887	0.0306 ^a

^aStatistically significant. ^bEstimate (b) is the coefficient of the linear model found by least-squares. ^cBeta is the parameter estimates that would have resulted from the regression had all the variables been standardized to a mean of 0 and a variance of 1. ^dNumber of clusters is based on hierarchical clustering of a Tanimoto similarity matrix derived from radial fingerprints.

the ligands sets, representing the diversity of the ligands and the promiscuity of their receptor. Results indicate that NUS is strongly correlated with the number of hydrophobic residues in the binding site, which has the most dominant and positive contribution, and with SiteMap exposure. NUS is anticorrelated with the SiteMap area donor parameter (Table 3). These results make chemical sense—the more hydrophobic and exposed the binding site is, the more possibilities it has for binding diverse antagonists. Conversely, more HB donors lead to more selectivity (i.e., less promiscuity), in line with previous publications indicating electrostatic interactions as a determinant of selectivity.^{45,46} Hydrophobic interactions are not as sensitive to the geometry of the interaction and thus are less specific than hydrophilic ones, which have a strong orientational dependence. The reason that the donor, but not the acceptor area presents a strong signal is probably related to the fact that the donor areas in GPCRs, as identified with SiteMap, are 25–50% larger than the corresponding acceptor areas.

An illustrative example of promiscuity-relevant descriptors of a narrow (CXCR4) and a broadly tuned (OPRM) receptor binding site is shown in Figure 3. Note the small donor region and large number of binding-site M/L/I/V residues in OPRM, vs the large donor region and single Ile residue in CXCR4.

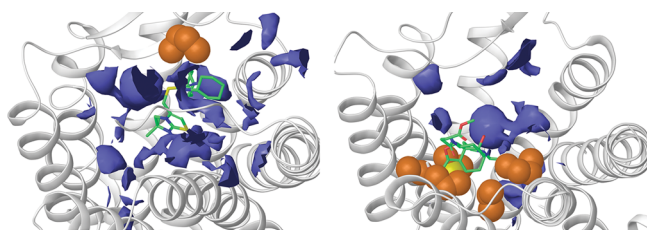


Figure 3. Examples of a relatively specific (CXCR4, left) and a promiscuous (OPRM, right) GPCR. Donor region is shown in blue surface, M/L/I/V residues carbons as orange CPK spheres, and the crystallized ligands carbons as green sticks. Note the larger donor region in the proximity of the ligand in CXCR4 compared to the large number of hydrophobic residues in OPRM.

The MW range of the active ligands is another descriptive estimate of diversity of the ligand sets (and thus of promiscuity of the receptor). When used as the predicted parameter in the least-squares model, MW range is dependent on the same factors as NUS (% M/L/I/V, exposure, and area donor) but is also negatively correlated with the number of disulfide bridges in the extracellular region of the receptor. The number of disulfide bridges is a significant factor when analyzed against the physical parameters of the molecules (MW and electrotopological states), but not when analyzed against parameters related to the size of the ligand set, such as NUS or number of molecules. This might be due to the role of disulfide bridges in enabling entry to the binding site: the less disulfide bridges, the more flexible the access to the binding site and the more space for molecules of variable size.

To examine the sensitivity of the observed trends to the chosen activity cutoff value, the statistical analyses were repeated on a larger set of ligands, obtained using an affinity cutoff of 10 μM (see Methods). The results of the standard least square analysis on the 10 μM set are summarized in Supporting Information Table S5. In general, receptor descriptors that were found to be the most influential in the 1 μM set analysis also had the strongest impact in the expanded 10 μM set.

Model Validation Using an Independent Set. The linear regression model was next applied to GPCRs not included in the training set. The promiscuity values predicted for the validation set receptors (Table 1) agree very well with the experimental NUS and number of molecules ($r = 0.91$ and $r = 0.9$, respectively; see Figure 4), and the fits are statistically significant (p -value = 0.01 and 0.0129, respectively). While the agreement is excellent, there is a consistent overestimation of NUS and number of molecules, which may be explained by the fact that the GPCRs in the validation set are less studied than the GPCRs in the training set. Despite this overestimation, it is clear that the models have strong predictive capability, especially when used in a comparative fashion. Moreover, the validation set receptors include a non-Family A GPCR (smoothened receptor, SMO), suggesting a broad applicability of the linear regression model. Predictions of additional parameters, such as range of MW, are also successful and are shown in Supporting Information Figure S4.

The predictions were repeated using homology models of the six validation set receptors. The results are shown in Figure 4C and D and Supporting Information Figure S4. In most cases, use of homology models results in good and statistically significant predictive ability. Thus, homology models may provide useful information about the ability of GPCRs to accommodate diverse molecules.

DISCUSSION

Ligand/receptor binding promiscuity is intertwined with molecular recognition and establishes the diversity of compounds that can alter the state of living systems. Receptor binding sites can exist along a spectrum, from highly selective to broadly promiscuous, which determines how a system responds to the myriad of natural compounds and xenobiotics that can be presented to a receptor. The evolution of protein function has been speculated to be associated with an initial period of promiscuity followed by iterative tuning of selectivity.^{47,48} Thornton et al. showed that nuclear receptors for steroids evolved according to a principle of “good enough” specificity: these proteins evolved to be specific enough to distinguish among the substances to which they were naturally exposed at that point in time.⁴⁸ Hence, exploration of promiscuity can further guide the study of the processes of evolution.⁴⁹

From a practical standpoint, an improved understanding of promiscuity can facilitate progress in protein engineering and drug design.¹⁹ Recent large-scale computational initiatives to collect and store experimental drug–target interaction data, enhanced by advances in chemogenomics tools, are changing our understanding of the role of drug selectivity. It is now widely recognized that most therapeutically effective molecules are not “magic bullets” that inhibit or activate single protein targets; instead, they are “magic shotguns”⁵⁰ that tend to interact with multiple proteins. Even with the current incomplete drug–target interaction data, only 15% of drugs appear to interact with only a single target and over 50% interact with more than five targets.^{47,51} From the protein point of view, Gao and Skolnick showed that more than 1/3 of about 20 000 known ligand-binding sites interact with multiple chemically different ligands.²¹ We have looked for determinants of receptor promiscuity by analyzing the sequence, physicochemical, structural, and energetic properties of GPCR binding sites and the relationship to the sets of known antagonists.

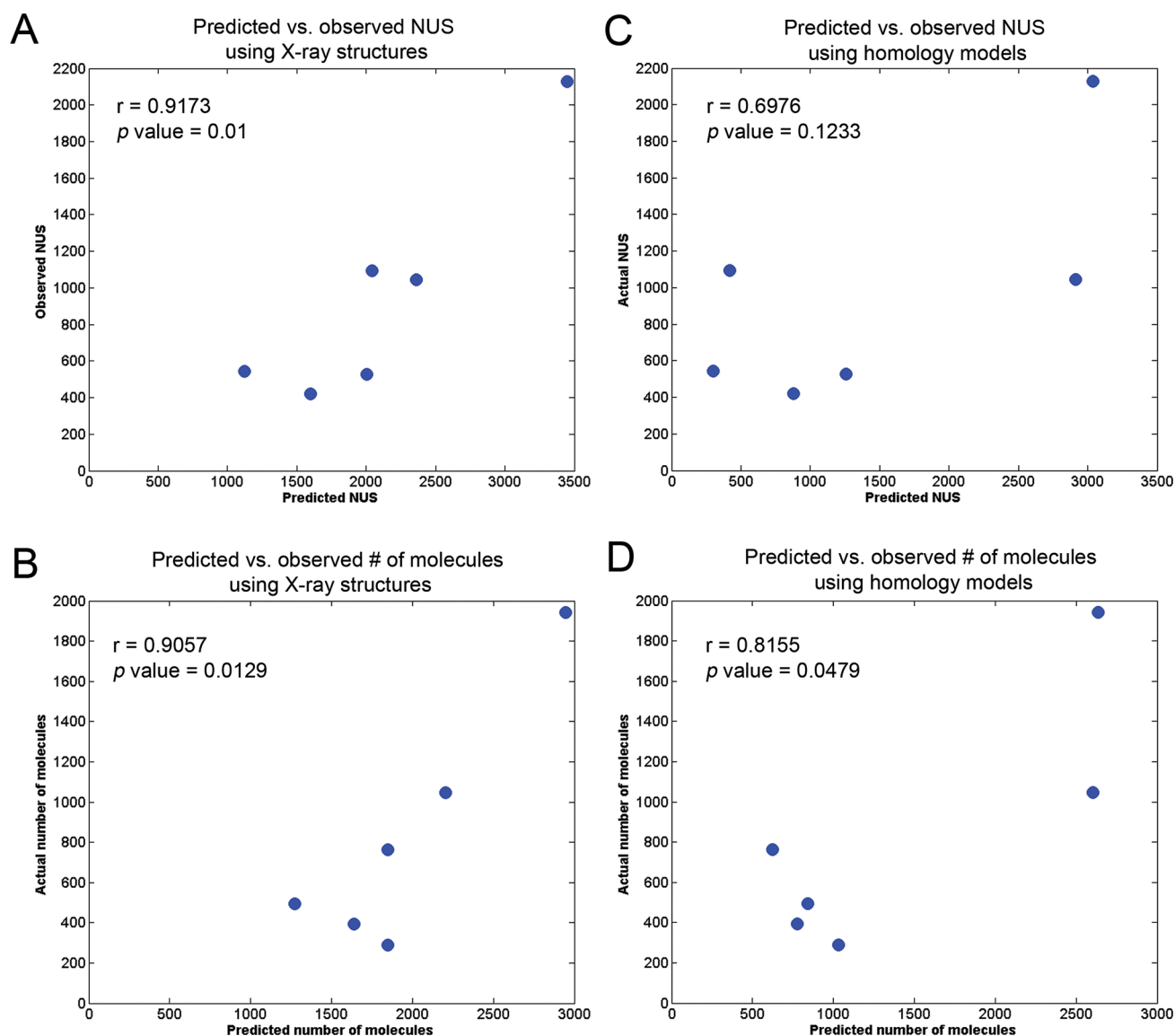


Figure 4. Predictive power of the linear regression models. Predicted vs experimental values of (A and C) NUS and (B and D) number of molecules, as tested on the X-ray (left) and homology model (right) validation sets, respectively. The values of each property were predicted using the linear regression models and compared to the experimental data obtained from ChEMBL. Pairwise correlation (r) and p values are shown for each plot.

To the best of our knowledge, no well-defined comparative measure of receptor promiscuity has previously been established. We therefore evaluated a number of criteria to assess the diversity of ligands that bind to a given GPCR. These included simply the size of the ligands data set for each of the receptors, the number of unique scaffolds in each of these sets, and the number of clusters derived from a Tanimoto similarity matrix. In addition, we measured the *range* of properties typically associated with diversity, including MW and atomic structural information. While no single metric is ideal for the task of classifying receptors based on their promiscuity, we have found that many of these properties are highly correlated, suggesting that together they provide an adequate measure of receptor promiscuity.

The mean hydrophobicity of the antagonists of a receptor is only moderately correlated with the chemical diversity of GPCR antagonists, as represented by the ranges of molecular properties. Thus, the gain obtained through removal of hydrophobic compounds from the aqueous environment and

by forming multiple, weak interactions between a hydrophobic ligand and hydrophobic amino acids, as was proposed for multidrug transporters,^{52,53} cannot solely explain the energetic gain in binding ligands to GPCRs. Nevertheless, hydrophobicity of the binding site of the more promiscuous receptors is higher and the content of HB donors is lower than for the less promiscuous ones.

MW and other size-related properties contribute significantly to the overall variability of the antagonists sets as they represent the main contributions to the principal axis (PC1) obtained in principle component analysis of sets of ligands studied in this work and are well-established dominant descriptors of ligands variability.⁵⁴ However, these do not correlate with the chemical variability of the antagonists (i.e., the number of unique scaffolds; NUS) or with the size of the antagonist sets. Specifically, the mean MW of ligands of the receptors classified as the most promiscuous ones does not differ significantly from the mean MW of ligands of the least promiscuous receptors. It is the *range* of observed MW that correlates with the number of

diverse ligands and which inversely correlates with the number of disulfide bridges in the extracellular region of the corresponding receptor. This may indicate that the higher disulfide bridge content restricts the shape of the binding site by rigidifying the loops. Conversely, the length of ECL2 and secondary structure content did not appear to have a significant impact on receptor promiscuity in our current data set of GPCRs. Neither did we find significant contribution with a crude measure of binding site flexibility, assessed by the estimation of number of accessible binding site side-chain rotamers.

The linear regression model is successful in predicting selectivity of an independent set of GPCRs. This implies that different GPCR targets can be prioritized in terms of the expected diversity of their antagonists using a simplistic model that relies solely on their structure. The successful predictions we obtained using homology models are particularly encouraging, since thus far, X-ray structures are available for only ~20 members of the large GPCR superfamily.

CONCLUSIONS

Our work links structural and sequence characteristics of GPCRs to the putative diversity of antagonists that they can bind. The number of GPCR antagonists correlates with ranges (but typically not with means or medians) of their physicochemical properties. GPCRs with wider ranges of antagonists' properties tend to have more hydrophobic amino acids and less H-bond donors in their orthosteric binding sites. GPCRs with high number of disulfide bridges in the extracellular region typically have a limited range of antagonists' set size. Binding site properties derived from X-rays or homology models of GPCRs that were not included in the original data set successfully predict receptors promiscuity. This analysis can be applied to GPCRs of unknown structure to estimate their expected promiscuity toward ligands, using GPCR homology models, so that the expected range of ligand repertoire can be evaluated even for orphan Family A GPCRs. Such valuable information may guide drug designers in target choice, since targeting promiscuous GPCRs provides room for discovery of multiple chemotypes, and opportunities for multiobjective optimization. On the protein engineering side, narrowing or expansion of the receptive range of a given GPCR by site-directed mutagenesis can be guided by computationally predicting the promiscuity of multiple mutant receptor candidates. The ideas introduced here provide insights into molecular recognition of GPCRs and may be generalized toward assessment and prediction of promiscuity of other protein families.

METHODS

Data Preparation. For each of the receptors in the training and validation sets (see Table 1), chemical structures of its antagonists and corresponding affinity data were retrieved from ChEMBL (version 15).²⁵ For the structures that were solved for species other than human, e.g. the turkey β 1AR, the rat M3R, and the mouse OPRM, the data for the corresponding human homologues were retrieved. Experimentally measured IC_{50} and K_i values were filtered to produce two sets of ligands for each receptor: (1) ligands with affinity of at least 1 μ M; (2) ligands with affinity of at least 10 μ M. In total, 11 750 and 16 177 receptor–ligand pairs were collected for the 1 and 10 μ M data sets, respectively.

Calculation of GPCR Structural Descriptors. *Structure Preparation.* Each X-ray structure in Table 1 was prepared using the Protein Preparation Wizard⁵⁵ implemented in Maestro (version 9.3, Schrödinger, LLC, New York, NY, 2013). For structures with multiple chains in the crystallographic unit chain "A" was selected, except for 2VT4, where we chose chain "B" instead (due to the anomalous kink in TM1 of chain A).

SiteMap Calculations. All GPCR orthosteric TM binding sites were subjected to SiteMap analysis³⁴ (as implemented in Maestro version 9.3), which characterizes binding sites in terms of size (the total number of SiteMap site points); volume, exposure, and enclosure (determined by how buried the site is); hydrophobic and hydrophilic character and balance (measures of the relative hydrophobic/hydrophilic nature of the site); contact (representing the strength of the van der Waals interaction of the site points with the protein); donor/acceptor (the degree to which a ligand might be expected to donate, rather than accept, hydrogen bonds); and SiteScore and DScore (weighted combinations of size, enclosure and hydrophilic character terms, which are found to distinguish drug-binding from nondrug-binding sites). SiteMap also calculates various surfaces representing the hydrophobic, donor, acceptor, and hydrophilic (donor + acceptor) regions of the binding site. The areas of these surfaces were used as descriptors as well (named area hydrophobic, area donor, area acceptor, and area hydrophilic). SiteMap was run in the default mode. All descriptors, including the size of the various surfaces generated by the algorithm, were extracted from the output files. There is a significant degree of redundancy among these parameters, but this can be adequately dealt with by the statistical analyses described below.

WaterMap Calculations. WaterMap (version 1.4, Schrödinger, LLC, New York, NY, 2012) is a method based on molecular dynamics (MD) and statistical thermodynamics to describe solvent energetics around protein surfaces.³⁵ WaterMap calculations were run in the default mode, including a grand canonical Monte Carlo simulation to hydrate the buried GPCR pocket prior to the MD stage. Average thermodynamic properties of all hydration sites in the binding pockets were calculated. In addition, we also calculated the number of distinct clusters of high-energy sites according to a previously published algorithm.³⁶

Flexibility Estimation. The flexibility of binding sites was estimated by performing an analysis of the rotameric states accessible to residues in the protein binding site. All residues within 4 Å of the ligand were considered. Using the rotamer library implemented in Prime (version 3.2, Schrödinger, LLC, New York, NY, 2013), all possible rotameric states for each residue were generated, and their energy was compared to that of the native state. All rotamers within 5 and 10 kcal/mol of the native state were kept. The sum of the total number of rotamers for each binding site was normalized by dividing by the number of residues in the binding site, or by the total number of enumerated states in the rotamer library for the particular residues in the binding site. The rank ordering of GPCR based on predicted flexibility using the different energy cutoffs and normalization schemes were highly similar, suggesting robustness of the predictions with respect to these rather arbitrary choices.

Intrinsic Specificity Ratio. The intrinsic specificity ratio, ISR, is a dimensionless quantity computed as the ratio of (1) the energy gap (δE) between the lowest (i.e., native) binding state

and the average of all other higher energy binding states and (2) the energy variance (ΔE) of the non-native states.⁴⁴ In this work we utilize a variation of the original ISR definition in which we explore the distribution of binding energies using a diverse fragment library rather than an ensemble of poses from the same ligand. This approach provides additional information about native and non-native interactions in the binding site that might not be accessible by a single ligand. The native binding state of a given receptor is assumed to comprise a set of favorable specific interactions necessary for ligand binding. The degree to which these specific interactions are important for stabilizing the protein–ligand complex is probed by computing the ISR for an ensemble of binding states generated by docking a library of diverse fragments spanning a broad range of potential protein–ligand interactions. This fragment library is available on the Schrödinger Web site (<http://schrodinger.com/productpage/14/5/78/>) and consists of 441 unique small fragments (6–37 atoms; molecular weight range 32–226) derived from molecules in the medicinal chemistry literature. The total set, after generating all energetically accessible ionization and tautomeric states, includes 667 fragments. The ISR values computed for the set of 10 GPCRs are given in Supporting Information Table S3. Higher values indicate more selective binding sites.

Ligand Set Descriptors. All chemical structures from ChEMBL were processed to remove inconsistencies and salts, prior to molecular descriptor calculations. The mean, median, range, standard deviation, and variance of the following molecular descriptors of the ligands for each target in our training set, were calculated using Canvas (version 1.5, Schrödinger, LLC, New York, NY, 2012):⁵⁶ molecular weight (MW), lipophilicity (as ALogP, the atomic LogP), hydrogen bond (HB) donor count (HBD), hydrogen bond acceptors count (HBA), rotatable bonds count (RB), polar surface area (PSA), electrotopological states (estate), molecular refractivity (MR), molecular polarizability (polar), centralization (cent), aromatic rings count (Ar ring), rings count (ring), chiral centers count (chiral), heavy atoms count (HA), total charge, negative atoms count, positive atoms count, and fraction of sp^3 -hybridized (tetrahedral) carbon atoms out of total carbon count (Fsp³). There are 18 default ligand descriptors in the *Physicochemical* descriptor set of Canvas, which were chosen based on surveying the literature to determine a small but diverse set of descriptors that are fast to compute and commonly used throughout the field. An extended analysis using a total of 340 ligand descriptors was carried out as well, by calculating all *Topological* and *Ligfilter* descriptor sets.

The number of unique scaffolds (NUS) in each set was computed using the Bemis–Murcko scaffold decomposition protocol⁵⁷ implemented in Canvas. In this process, the largest scaffold in a structure is first obtained by stripping off all terminal side chains with the exception of exocyclic and exolinker double bonds. The resulting scaffold is then split into all possible smaller subscaffolds by breaking bonds and removing linkers between rings, other than those in fused rings. The resulting scaffolds were used to calculate the scaffold chiral centers count. While the Bemis–Murcko strategy provides a good measure of the diversity of molecular frameworks in a given data set, it should be noted that it increases the NUS count for data sets containing fused aromatic rings (the majority of targets studied here) over data sets containing mostly aliphatic compounds (i.e., the lipid S1P1 receptor).

Hierarchical clustering on a Tanimoto similarity matrix using either radial or dendritic fingerprints (FP)⁵⁸ was carried out in Canvas. A clustering level of 0.6 Tanimoto similarity was chosen, based on a preliminary convergence analysis of clusters number. The number of clusters (either for “radial” or “dendritic” FPs) was calculated for each ligand set. Key descriptors found to be significantly related to receptor promiscuity in our subsequent analysis were also calculated for the ligands of the six targets in the validation set.

Statistical Analysis. Principal component analysis (PCA) was performed using Matlab (version R2012b; Mathworks, Inc., MA, USA) on receptor descriptors and on median values of ligand descriptors for the 10 sets in our training set. Prior to analysis, factor scaling was applied to normalize the data. Statistical analyses were carried out in the JMP statistical software package (version 7.0.1; SAS Institute Inc., NC, USA). Correlations between ligand descriptors or between receptor descriptors were computed using the multivariate module. Stepwise regression analysis was used to determine the receptor descriptors that explain most of the variance of a specific ligand diversity measure. These receptor descriptors were then used to generate a regression model for prediction of the diversity measure using standard least squares fit.

Homology Modeling. In order to make models for a subset of GPCRs in the validation set (OPRD, OPRX, SHT1B, SHT2B, and PAR1), the closest template (in terms of sequence identity) available at the time was chosen. This resulted in OPRM as a template for OPRD and OPRX, OPRK for PAR1, and the DRD3 for SHT1B and SHT2B. Models were built using Prime (version 3.1, Schrödinger, LLC, New York, NY, 2012). The alignment was created with the Multiple Sequence Viewer in Maestro using the ClustalW algorithm, ensuring manually the correct alignment of all the conserved motifs in the 7 TM domains and the conserved Cys in ECL2. Gaps and insertions were placed at the center of the loop when needed. No further loop refinement was performed. The model for SMO was obtained by submitting the N-terminus-truncated sequence of smoothened receptor to I-TASSER server,⁵⁹ an iterative threading assembly refinement (I-TASSER) server. Starting from an amino acid sequence, I-TASSER generates three-dimensional atomic models from multiple threading alignments and iterative structural assembly simulations. The X-ray SMO structure (4JKV) was excluded from the templates. Interestingly, the two top templates used in modeling were both Family B GPCRs, 4L6R (glucagon receptor) and 4KSY (CRF1R).

■ ASSOCIATED CONTENT

■ Supporting Information

Figure S1: Chemical characterization of analyzed ligand sets. Figure S2: Multivariate correlations between all statistical properties of the ligand set descriptors. Figure S3: Statistical analysis of receptor descriptors. Figure S4: Prediction power of the linear regression models. Table S1: Statistically significant pairwise correlations between ligand descriptors for the 1 μ M ligand sets. Table S2: Binding site residues for all GPCRs studied. Table S3: ISR results. Table S4: Statistically significant pairwise correlations between receptor descriptors for the 1 μ M ligand sets. Table S5: Regression coefficients of receptor descriptors affecting NUS and various additional ligand descriptors, as determined by standard least-squares fit analysis on the 10 μ M ligand sets. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: masha.niv@mail.huji.ac.il. Tel.: +972-(0)8-9489664.
Fax: +972-(0)8-9476189.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

We thank Dr. Elite Levine for consultation on statistical analysis, Ayana Wiener for helpful discussions, and Prof. Amirav Gordon for his kind encouragement. The Israel Science Foundation (No. 432/12) and the German Research Foundation DFG (ME 1024/8-1) grants to M.Y.N. are gratefully acknowledged. M.Y.N. participates in the European COST Action CM1207 (GLISTEN).

■ REFERENCES

- (1) Fredriksson, R.; Lagerstrom, M. C.; Lundin, L. G.; Schiöth, H. B. The G-protein-coupled receptors in the human genome form five main families. Phylogenetic analysis, paralogon groups, and fingerprints. *Mol. Pharmacol.* **2003**, *63*, 1256–72.
- (2) Rask-Andersen, M.; Almen, M. S.; Schiöth, H. B. Trends in the exploitation of novel drug targets. *Nat. Rev. Drug Discov.* **2011**, *10*, 579–90.
- (3) Leinders-Zufall, T.; Lane, A. P.; Puche, A. C.; Ma, W.; Novotny, M. V.; Shipley, M. T.; Zufall, F. Ultrasensitive pheromone detection by mammalian vomeronasal neurons. *Nature* **2000**, *405*, 792–6.
- (4) Hallem, E. A.; Carlson, J. R. Coding of odors by a receptor repertoire. *Cell* **2006**, *125*, 143–60.
- (5) Li, J.; Haddad, R.; Chen, S.; Santos, V.; Luetje, C. W. A broadly tuned mouse odorant receptor that detects nitrotoluenes. *J. Neurochem.* **2012**, *121*, 881–90.
- (6) Baud, O.; Etter, S.; Spreafico, M.; Bordoli, L.; Schwede, T.; Vogel, H.; Pick, H. The mouse eugenol odorant receptor: structural and functional plasticity of a broadly tuned odorant binding pocket. *Biochemistry* **2011**, *50*, 843–53.
- (7) Meyerhof, W.; Batram, C.; Kuhn, C.; Brockhoff, A.; Chudoba, E.; Bufo, B.; Appendino, G.; Behrens, M. The molecular receptive ranges of human TAS2R bitter taste receptors. *Chem. Senses* **2010**, *35*, 157–70.
- (8) Born, S.; Levit, A.; Niv, M. Y.; Meyerhof, W.; Behrens, M. The Human Bitter Taste Receptor TAS2R10 Is Tailored to Accommodate Numerous Diverse Ligands. *J. Neurosci.* **2013**, *33*, 201–13.
- (9) Venkatakrishnan, A. J.; Deupi, X.; Lebon, G.; Tate, C. G.; Schertler, G. F.; Babu, M. M. Molecular signatures of G-protein-coupled receptors. *Nature* **2013**, *494*, 185–94.
- (10) Brockhoff, A.; Behrens, M.; Niv, M. Y.; Meyerhof, W. Structural requirements of bitter taste receptor activation. *Proc. Natl. Acad. Sci. USA* **2010**, *107*, 11110–5.
- (11) Warne, T.; Moukhametzanov, R.; Baker, J. G.; Nehme, R.; Edwards, P. C.; Leslie, A. G.; Schertler, G. F.; Tate, C. G. The structural basis for agonist and partial agonist action on a beta(1)-adrenergic receptor. *Nature* **2011**, *469*, 241–4.
- (12) Doré, A. S.; Robertson, N.; Errey, J. C.; Ng, I.; Hollenstein, K.; Tehan, B.; Hurrell, E.; Bennett, K.; Congreve, M.; Magnani, F.; Tate, C. G.; Weir, M.; Marshall, F. H. Structure of the Adenosine A2A Receptor in Complex with ZM241385 and the Xanthines XAC and Caffeine. *Structure* **2011**, *19*, 1283–1293.
- (13) Levit, A.; Barak, D.; Behrens, M.; Meyerhof, W.; Niv, M. Homology Model-Assisted Elucidation of Binding Sites in GPCRs. In *Membrane Protein Structure and Dynamics*; Vaidehi, N., Klein-Seetharaman, J., Eds.; Humana Press, 2012; Vol. 914, Chapter 11, pp 179–205.
- (14) Lin, H.; Sassano, M. F.; Roth, B. L.; Shoichet, B. K. A pharmacological organization of G protein-coupled receptors. *Nat. Methods* **2013**, *10*, 140–6.
- (15) van Westen, G. J.; Overington, J. P. A ligand's-eye view of protein similarity. *Nat. Methods* **2013**, *10*, 116–7.
- (16) van Westen, G. J. P.; Wegner, J. K.; IJzerman, A. P.; van Vlijmen, H. W. T.; Bender, A. Proteochemometric modeling as a tool to design selective compounds and for extrapolating to novel targets. *MedChemComm* **2011**, *2*, 16–30.
- (17) Weill, N.; Rognan, D. Development and validation of a novel protein-ligand fingerprint to mine chemogenomic space: application to G protein-coupled receptors and their ligands. *J. Chem. Inf. Model.* **2009**, *49*, 1049–62.
- (18) Wichard, J. D.; Ter Laak, A.; Krause, G.; Heinrich, N.; Kuhne, R.; Kleinau, G. Chemogenomic analysis of G-protein coupled receptors and their ligands deciphers locks and keys governing diverse aspects of signalling. *PLoS One* **2011**, *6*, e16811.
- (19) Nobeli, I.; Favia, A. D.; Thornton, J. M. Protein promiscuity and its implications for biotechnology. *Nat. Biotechnol.* **2009**, *27*, 157–67.
- (20) Gleeson, M. P.; Hersey, A.; Montanari, D.; Overington, J. Probing the links between in vitro potency, ADMET and physicochemical parameters. *Nat. Rev. Drug Discov.* **2011**, *10*, 197–208.
- (21) Gao, M.; Skolnick, J. A comprehensive survey of small-molecule binding pockets in proteins. *PLoS Comput. Biol.* **2013**, *9*, e1003302.
- (22) Kooistra, A. J.; Kuhne, S.; de Esch, I. J. P.; Leurs, R.; de Graaf, C. A structural chemogenomics analysis of aminergic GPCRs: lessons for histamine receptor ligand design. *Br. J. Pharmacol.* **2013**, *170*, 101–126.
- (23) Newman, A. H.; Beuming, T.; Banala, A. K.; Donthamsetti, P.; Pongetti, K.; LaBounty, A.; Levy, B.; Cao, J.; Michino, M.; Luedtke, R. R.; Javitch, J. A.; Shi, L. Molecular determinants of selectivity and efficacy at the dopamine D3 receptor. *J. Med. Chem.* **2012**, *55*, 6689–99.
- (24) Wu, H.; Wacker, D.; Mileni, M.; Katritch, V.; Han, G. W.; Vardy, E.; Liu, W.; Thompson, A. A.; Huang, X. P.; Carroll, F. I.; Mascarella, S. W.; Westkaemper, R. B.; Mosier, P. D.; Roth, B. L.; Cherezov, V.; Stevens, R. C. Structure of the human kappa-opioid receptor in complex with JDTic. *Nature* **2012**, *485*, 327–32.
- (25) 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.* **2012**, *40*, D1100–7.
- (26) Granier, S.; Kobilka, B. A new era of GPCR structural and chemical biology. *Nat. Chem. Biol.* **2012**, *8*, 670–3.
- (27) Wang, C.; Jiang, Y.; Ma, J.; Wu, H.; Wacker, D.; Katritch, V.; Han, G. W.; Liu, W.; Huang, X. P.; Vardy, E.; McCorvy, J. D.; Gao, X.; Zhou, E. X.; Melcher, K.; Zhang, C.; Bai, F.; Yang, H.; Yang, L.; Jiang, H.; Roth, B. L.; Cherezov, V.; Stevens, R. C.; Xu, H. E. Structural Basis for Molecular Recognition at Serotonin Receptors. *Science* **2013**, *340*, 610–614.
- (28) Walters, W. P.; Green, J.; Weiss, J. R.; Murcko, M. A. What do medicinal chemists actually make? A 50-year retrospective. *J. Med. Chem.* **2011**, *54*, 6405–16.
- (29) Leeson, P. D.; St-Gallay, S. A. The influence of the 'organizational factor' on compound quality in drug discovery. *Nat. Rev. Drug Discov.* **2011**, *10*, 749–65.
- (30) de Graaf, C.; Rognan, D. Customizing G Protein-coupled receptor models for structure-based virtual screening. *Curr. Pharm. Des.* **2009**, *15*, 4026–48.
- (31) Levit, A.; Barak, D.; Behrens, M.; Meyerhof, W.; Niv, M. Y. Homology model-assisted elucidation of binding sites in GPCRs. *Methods Mol. Biol.* **2012**, *914*, 179–205.
- (32) Shi, L.; Javitch, J. A. The binding site of aminergic G protein-coupled receptors: the transmembrane segments and second extracellular loop. *Annu. Rev. Pharmacol. Toxicol.* **2002**, *42*, 437–67.
- (33) Katritch, V.; Cherezov, V.; Stevens, R. C. Structure-function of the G protein-coupled receptor superfamily. *Ann. Rev. Pharmacol. Toxicol.* **2013**, *53*, 531–56.
- (34) Halgren, T. New method for fast and accurate binding-site identification and analysis. *Chem. Biol. Drug Des.* **2007**, *69*, 146–8.

- (35) Abel, R.; Young, T.; Farid, R.; Berne, B. J.; Friesner, R. A. Role of the active-site solvent in the thermodynamics of factor Xa ligand binding. *J. Am. Chem. Soc.* **2008**, *130*, 2817–31.
- (36) Beuming, T.; Che, Y.; Abel, R.; Kim, B.; Shanmugasundaram, V.; Sherman, W. Thermodynamic analysis of water molecules at the surface of proteins and applications to binding site prediction and characterization. *Proteins* **2012**, *80*, 871–83.
- (37) Mason, J. S.; Bortolato, A.; Congreve, M.; Marshall, F. H. New insights from structural biology into the druggability of G protein-coupled receptors. *Trends Pharmacol. Sci.* **2012**, *33*, 249–60.
- (38) Beuming, T.; Farid, R.; Sherman, W. High-energy water sites determine peptide binding affinity and specificity of PDZ domains. *Protein Sci.* **2009**, *18*, 1609–19.
- (39) Robinson, D. D.; Sherman, W.; Farid, R. Understanding kinase selectivity through energetic analysis of binding site waters. *ChemMedChem* **2010**, *5*, 618–27.
- (40) Snyder, P. W.; Mecinovic, J.; Moustakas, D. T.; Thomas, S. W., 3rd; Harder, M.; Mack, E. T.; Lockett, M. R.; Heroux, A.; Sherman, W.; Whitesides, G. M. Mechanism of the hydrophobic effect in the biomolecular recognition of arylsulfonamides by carbonic anhydrase. *Proc. Natl. Acad. Sci. U.S.A.* **2011**, *108*, 17889–94.
- (41) Breiten, B.; Lockett, M. R.; Sherman, W.; Fujita, S.; Al-Sayah, M.; Lange, H.; Bowers, C. M.; Heroux, A.; Krilov, G.; Whitesides, G. M. Water networks contribute to enthalpy/entropy compensation in protein-ligand binding. *J. Am. Chem. Soc.* **2013**, *135*, 15579–84.
- (42) Carlson, H. A. Protein flexibility and drug design: how to hit a moving target. *Curr. Opin. Chem. Biol.* **2002**, *6*, 447–52.
- (43) Grove, L. E.; Hall, D. R.; Beglov, D.; Vajda, S.; Kozakov, D. FTFlex: accounting for binding site flexibility to improve fragment-based identification of druggable hot spots. *Bioinformatics* **2013**, *29*, 1218–9.
- (44) Wang, J.; Zheng, X.; Yang, Y.; Drueckhammer, D.; Yang, W.; Verkhivker, G.; Wang, E. Quantifying intrinsic specificity: a potential complement to affinity in drug screening. *Phys. Rev. Lett.* **2007**, *99*, 198101.
- (45) Huggins, D. J.; Sherman, W.; Tidor, B. Rational approaches to improving selectivity in drug design. *J. Med. Chem.* **2012**, *55*, 1424–44.
- (46) Watkins, R. E.; Wisely, G. B.; Moore, L. B.; Collins, J. L.; Lambert, M. H.; Williams, S. P.; Willson, T. M.; Kliewer, S. A.; Redinbo, M. R. The human nuclear xenobiotic receptor PXR: structural determinants of directed promiscuity. *Science* **2001**, *292*, 2329–33.
- (47) Jalencas, X.; Mestres, J. On the origins of drug polypharmacology. *MedChemComm* **2013**, *4*, 80–87.
- (48) Eick, G. N.; Colucci, J. K.; Harms, M. J.; Ortlund, E. A.; Thornton, J. W. Evolution of minimal specificity and promiscuity in steroid hormone receptors. *PLoS Genetics* **2012**, *8*, e1003072.
- (49) Skolnick, J.; Gao, M. Interplay of physics and evolution in the likely origin of protein biochemical function. *Proc. Natl. Acad. Sci. U.S.A.* **2013**, *110*, 9344–9.
- (50) Allen, J. A.; Roth, B. L. Strategies to discover unexpected targets for drugs active at G protein-coupled receptors. *Annu. Rev. Pharmacol. Toxicol.* **2011**, *51*, 117–44.
- (51) Perez-Nueno, V. I.; Ritchie, D. W. Identifying and characterizing promiscuous targets: implications for virtual screening. *Expert Opin. Drug Discov.* **2012**, *7*, 1–17.
- (52) Higgins, C. F. Multiple molecular mechanisms for multidrug resistance transporters. *Nature* **2007**, *446*, 749–57.
- (53) Fluman, N.; Bibi, E. Bacterial multidrug transport through the lens of the major facilitator superfamily. *Biochim. Biophys. Acta* **2009**, *1794*, 738–47.
- (54) Khan, R. M.; Luk, C. H.; Flinker, A.; Aggarwal, A.; Lapid, H.; Haddad, R.; Sobel, N. Predicting odor pleasantness from odorant structure: pleasantness as a reflection of the physical world. *J. Neurosci.* **2007**, *27*, 10015–23.
- (55) Sastry, G. M.; Adzhigirey, M.; Day, T.; Annabhimoju, R.; Sherman, W. Protein and ligand preparation: parameters, protocols, and influence on virtual screening enrichments. *J. Comput.-Aided Mol. Des.* **2013**, *27*, 221–34.
- (56) Sastry, M.; Lowrie, J. F.; Dixon, S. L.; Sherman, W. Large-scale systematic analysis of 2D fingerprint methods and parameters to improve virtual screening enrichments. *J. Chem. Inf. Model.* **2010**, *50*, 771–84.
- (57) Bemis, G. W.; Murcko, M. A. The properties of known drugs. 1. Molecular frameworks. *J. Med. Chem.* **1996**, *39*, 2887–93.
- (58) Duan, J.; Dixon, S. L.; Lowrie, J. F.; Sherman, W. Analysis and comparison of 2D fingerprints: insights into database screening performance using eight fingerprint methods. *J. Mol. Graphics Modell.* **2010**, *29*, 157–70.
- (59) Roy, A.; Kucukural, A.; Zhang, Y. I-TASSER: a unified platform for automated protein structure and function prediction. *Nat. Protoc.* **2010**, *5*, 725–38.