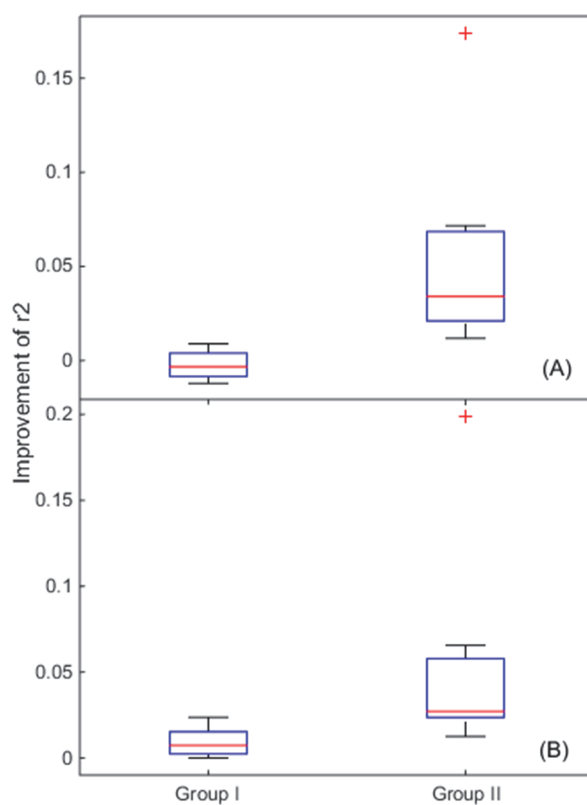# Supporting Information

## Supporting Figures



**Fig. S1.** Influence of the number of ligand samples on the model performance. (A) The improvement in $r^2$ of the model based on the top-300 features selected from 1,024 bits against the baseline model. (B) The improvement in $r^2$ of the optimal model (highlighted in boldface in Table 1) against the model based on the top-300 features selected from 1024 bits. Group I: GPCR datasets with more than 600 ligands, i.e., P08908, Q9Y5N1, P28335, P35372, Q99705, P0DMS8, Q16602, P51677, P48039; Group II: GPCR datasets with fewer than 600 ligands, i.e., Q9H228, Q8TDU6, Q8TDS4, Q9HC97, P41180, Q14833, Q99835.

# Supporting Tables

## Table S1. Description of datasets used in this study

| UniProt ID | Gene Name | Protein Name | Class | Subfamily | # of Ligands | # of Controls | Clinical Significance |
|---|---|---|---|---|---|---|---|
| P08908 | HTR1A | 5-hydroxytryptamine receptor 1A | A | Aminergic receptors | 4322 | 850 | Blood pressure, heart rate, antidepressant, anxiolytic, schizophrenia and Parkinson (H Ito, 1999) |
| Q9Y5N1 | HRH3 | Histamine H3 receptor | A | Aminergic receptors | 3644 | 700 | Cognitive disorders (Esbenshade, et al., 2008) |
| P28335 | HTR2C | 5-hydroxytryptamine receptor 2C | A | Aminergic receptors | 3286 | 650 | mood, anxiety, feeding, and reproductive behavior(Heisler, et al., 2007) |
| P35372 | OPRM1 | Mu-type opioid receptor | A | Peptide receptors | 4591 | 900 | Morphine-induced analgesia and itching (Liu, et al., 2011) |
| Q99705 | MCHR1 | Melanin-concentrating hormone receptors 1 | A | Peptide receptors | 3663 | 700 | Appetite, anxiety and depression (Rivera, et al., 2008) |
| P0DMS8 | ADORA3 | Adenosine receptor A3 | A | Nucleotide receptors | 3664 | 700 | Bronchial asthma(Jacobson, et al., 2008)and rheumatoid arthritis(Silverman, et al., 2008) |
| Q9H228 | S1PR5 | Sphingosine 1-phosphate receptor 5 | A | Lipid receptors | 320 | 60 | Huntington's disease(Buttari, 2018) |
| P51677 | CCR3 | C-C chemokine receptor type 3 | A | Protein receptors | 1147 | 200 | Binds and responds to a variety of chemokines (H, et al., 1996) |
| P48039 | MTNR1A | Melatonin receptor type 1A | A | Melatonin receptors | 946 | 190 | Circadian rhythm(Slaugenhaupt, et al., 1995) |
| Q8TDU6 | GPBAR1 | G-protein coupled bile acid receptor 1 | A | Steroid receptors | 464 | 90 | Suppression of macrophage functions and regulation of energy homeostasis by bile acids(Wang, et al., 2011) |
| Q8TDS4 | HCAR2 | Hydroxycarboxylic acid receptor 2 | A | Alicarboxylic acid receptors | 500 | 100 | Dyslipidemia(Hu, et al., 2015) |
| Q9HC97 | GPR35 | G-protein coupled receptor 35 | A | Orphan receptors | 330 | 60 | Brachydactyly mental retardation syndrome(Shrimpton, et al., 2004) |
| Q16602 | CALCRL | Calcitonin gene-related peptide type 1 receptor | B | Peptide receptors | 691 | 140 | migraine(Edvinsson, 2008) |
| P41180 | CASR | Extracellular calcium-sensing receptor | C | Ion receptors | 423 | 80 | Alzheimer's disease, asthma(Kim, et al., 2014) |
| Q14833 | GRM4 | Metabotropic glutamate receptor 4 | C | Amino acid receptors | 548 | 110 | Hallucinogenesis |
| Q99835 | SMO | Smoothened homolog | F | Protein receptors | 591 | 120 | Developmental disorders |

**Table S2.** Influence of regression models on the SED performance.

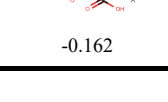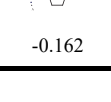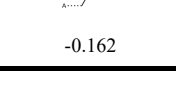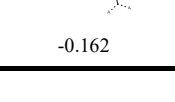| Group[a] | GPCRs | EC[b] | RF | GBDT | SVR | DNN |
|---|---|---|---|---|---|---|
| I | P08908 | $r^2$ (↑) | 0.9180 | 0.9013 | 0.9027 | **0.9314** |
| | | RMSE(↓) | 1.0450 | 1.1556 | **0.9015** | 0.9879 |
| | Q9Y5N1 | $r^2$ (↑) | 0.9380 | 0.9167 | 0.9250 | **0.9598** |
| | | RMSE(↓) | **0.9470** | 1.0886 | 1.0421 | 0.9486 |
| | P28335 | $r^2$ (↑) | 0.8911 | 0.8781 | 0.8895 | **0.9095** |
| | | RMSE(↓) | 1.1482 | 1.2048 | 1.1503 | **1.1184** |
| | P35372 | $r^2$ (↑) | 0.8908 | 0.8697 | 0.8642 | **0.8954** |
| | | RMSE(↓) | 1.1660 | 1.2555 | 1.2807 | **1.1547** |
| | Q99705 | $r^2$ (↑) | 0.9214 | 0.9142 | 0.9115 | **0.9436** |
| | | RMSE(↓) | 0.9959 | 1.0137 | 1.0424 | **0.8928** |
| | P0DMS8 | $r^2$ (↑) | 0.8880 | 0.8649 | 0.8616 | **0.8938** |
| | | RMSE(↓) | **1.1650** | 1.2719 | 1.2857 | 1.1907 |
| | Q16602 | $r^2$ (↑) | 0.9489 | 0.8722 | 0.7980 | **0.9533** |
| | | RMSE(↓) | **0.8795** | 1.3374 | 1.7244 | 1.4746 |
| | P51677 | $r^2$ (↑) | 0.9175 | 0.8904 | 0.8860 | **0.9405** |
| | | RMSE(↓) | **1.0223** | 1.1356 | 1.1887 | 1.0280 |
| | P48039 | $r^2$ (↑) | 0.9027 | 0.8571 | 0.8922 | **0.9147** |
| | | RMSE(↓) | **1.2811** | 1.4368 | 1.3275 | 1.3635 |
| II | Q9H228 | $r^2$ (↑) | 0.8725 | 0.8889 | 0.8826 | **0.9100** |
| | | RMSE(↓) | 1.4275 | 1.2885 | **1.2426** | 1.3231 |
| | Q8TDU6 | $r^2$ (↑) | 0.9275 | 0.8961 | 0.8921 | **0.9329** |
| | | RMSE(↓) | **0.9789** | 1.0694 | 1.1623 | 1.0253 |
| | Q8TDS4 | $r^2$ (↑) | 0.8912 | 0.9098 | 0.9006 | **0.9378** |
| | | RMSE(↓) | 1.0791 | **0.8734** | 0.9731 | 0.9567 |
| | Q9HC97 | $r^2$ (↑) | 0.7272 | 0.5998 | 0.8371 | **0.8508** |
| | | RMSE(↓) | 1.7346 | 1.7818 | 1.4090 | **1.3631** |
| | P41180 | $r^2$ (↑) | 0.7384 | 0.6673 | 0.7374 | **0.8435** |
| | | RMSE(↓) | 1.9836 | 2.0982 | 1.8910 | **1.5410** |
| | Q14833 | $r^2$ (↑) | 0.7559 | 0.6268 | 0.7719 | **0.7947** |
| | | RMSE(↓) | 1.5008 | 1.8198 | **1.4169** | 1.4635 |
| | Q99835 | $r^2$ (↑) | 0.8627 | 0.7840 | 0.8235 | **0.9028** |
| | | RMSE(↓) | 1.2819 | 1.5412 | 1.3882 | **1.1239** |

RF: Random Forest, GBDT: Gradient Boosting Decision Tree, SVR: Support Vector Regression, DNN: Deep Neural Network. [a]Group I: original number of ligands ＞600; II: original number of ligands ≤ 600. [b]Evaluation Criterion: ↑ (↓) indicates that larger (smaller) values are better; the best results for each evaluation criterion are highlighted in boldface. The input of each regression model was the top-300 features selected from the optimal bits of the ECFPs. For each GPCR dataset, the optimal bit is the ECFP length corresponding to the optimal result (highlighted in boldface in Table 1).

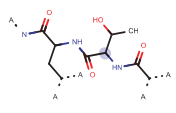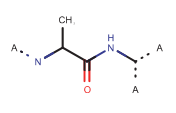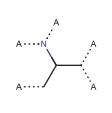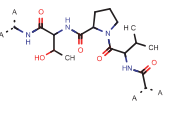# Table S3. Comparison of SED with the WDL-RF method

| Group[a] | GPCRs | EC[b] | WDL-RF | SED[c] |
|---|---|---|---|---|
| I | P08908 | $r^2$(↑) | 0.7062 | 0.9314* |
| | | RMSE(↓) | 1.5038 | 0.9879* |
| | Q9Y5N1 | $r^2$(↑) | 0.8242 | 0.9598* |
| | | RMSE(↓) | 1.6293 | 0.9486* |
| | P28335 | $r^2$(↑) | 0.5764 | 0.9095* |
| | | RMSE(↓) | 2.2714 | 1.1184* |
| | P35372 | $r^2$(↑) | 0.7428 | 0.8954* |
| | | RMSE(↓) | 1.5127 | 1.1547* |
| | Q99705 | $r^2$(↑) | 0.8236 | 0.9436* |
| | | RMSE(↓) | 1.3342 | 0.8928* |
| | P0DMS8 | $r^2$(↑) | 0.7580 | 0.8938* |
| | | RMSE(↓) | 1.6753 | 1.1907* |
| | Q16602 | $r^2$(↑) | 0.7482 | 0.9533* |
| | | RMSE(↓) | 0.8153 | 1.4746* |
| | P51677 | $r^2$(↑) | 0.8491 | 0.9405* |
| | | RMSE(↓) | 1.3391 | 1.0280* |
| | P48039 | $r^2$(↑) | 0.7874 | 0.9147* |
| | | RMSE(↓) | 1.8357 | 1.3635* |
| II | Q9H228 | $r^2$(↑) | 0.5839 | 0.9100* |
| | | RMSE(↓) | 0.6165 | 1.3231* |
| | Q8TDU6 | $r^2$(↑) | 0.7303 | 0.9329* |
| | | RMSE(↓) | 1.6806 | 1.0253* |
| | Q8TDS4 | $r^2$(↑) | 0.4556 | 0.9378* |
| | | RMSE(↓) | 0.7637 | 0.9567* |
| | Q9HC97 | $r^2$(↑) | 0.4797 | 0.8508* |
| | | RMSE(↓) | 0.7764 | 1.3631* |
| | P41180 | $r^2$(↑) | 0.6619 | 0.8435* |
| | | RMSE(↓) | 2.0476 | 1.5410* |
| | Q14833 | $r^2$(↑) | 0.4169 | 0.7947* |
| | | RMSE(↓) | 0.8373 | 1.4635* |
| | Q99835 | $r^2$(↑) | 0.4834 | 0.9028* |
| | | RMSE(↓) | 2.1368 | 1.1239* |

[a]Group I: original number of ligands ＞600; II: original number of ligands ≤ 600. [b]Evaluation Criterion: ↑ (↓) indicates that larger (smaller) values are better. [c]SED: For each GPCR dataset, the optimal result was selected for comparison (highlighted in boldface in Table 1). * indicates that the performance of the SED method is significantly better than that of the WDL-RF methods (Wu, et al., 2018) based on Wilcoxon signed-rank test.
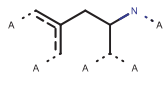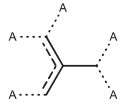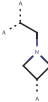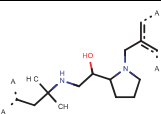
**Table S4.** Top 51 to 300 Substructures Identified by SED Accompanying with the Associated Pearson Correlation Coefficients.

| Top 51-60 (Pearson correlation) | Top 61-70 (Pearson correlation) | Top 71-80 (Pearson correlation) | Top 81-90 (Pearson correlation) | Top 91-100 (Pearson correlation) |
|---|---|---|---|---|
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |
|  -0.162 |  -0.162 |  -0.162 |  -0.162 |  -0.162 |

**(Table S4 Continued)**

| Top 101-110 | Top 111-120 | Top 121-130 | Top 131-140 | Top 141-150 |
|---|---|---|---|---|
| (Pearson correlation) | (Pearson correlation) | (Pearson correlation) | (Pearson correlation) | (Pearson correlation) |
| -0.162 | -0.162 | -0.162 | -0.162 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.162 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.162 | -0.199 |
| -0.162 | -0.162 | -0.162 | -0.162 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.115 |

(Table S4 Continued)

| Top 151-160 | Top 161-170 | Top 171-180 | Top 181-190 | Top 191-200 |
|---|---|---|---|---|
| (Pearson correlation) | (Pearson correlation) | (Pearson correlation) | (Pearson correlation) | (Pearson correlation) |
| -0.23 | -0.214 | 0.28 | -0.162 | -0.199 |
| -0.199 | 0.288 | 0.191 | 0.223 | -0.347 |
| -0.484 | -0.134 | -0.199 | -0.159 | -0.23 |
| -0.162 | -0.162 | -0.326 | -0.199 | -0.122 |
| -0.199 | 0.237 | -0.162 | -0.162 | 0.291 |
| 0.288 | 0.242 | 0.04 | -0.199 | -0.199 |
| 0.255 | -0.133 | -0.199 | -0.267 | 0.033 |
| -0.367 | -0.196 | 0.063 | -0.162 | -0.23 |
| 0.276 | -0.199 | 0.299 | -0.347 | -0.306 |
| 0.276 | -0.199 | -0.257 | -0.258 | 0.022 |

(Table S4 continued)

| Top 201-210 (Pearson correlation) | Top 211-220 (Pearson correlation) | Top 221-230 (Pearson correlation) | Top 231-240 (Pearson correlation) | Top 241-250 (Pearson correlation) |
|---|---|---|---|---|
| -0.258 | -0.115 | -0.115 | -0.115 | -0.115 |
| 0.211 | -0.115 | -0.115 | -0.162 | -0.162 |
| 0.17 | -0.162 | -0.162 | 0.162 | -0.162 |
| -0.162 | -0.162 | 0.036 | -0.162 | -0.162 |
| -0.115 | -0.162 | -0.162 | -0.115 | -0.115 |
| -0.162 | -0.162 | -0.162 | -0.115 | -0.199 |
| -0.199 | -0.162 | -0.258 | -0.162 | -0.162 |
| -0.115 | -0.115 | -0.115 | -0.162 | -0.115 |
| -0.115 | -0.162 | -0.115 | -0.115 | -0.162 |
| -0.162 | -0.162 | -0.115 | -0.115 | -0.162 |

(Table S4 continued)

| Top 251-260 (Pearson correlation) | Top 261-270 (Pearson correlation) | Top 271-280 (Pearson correlation) | Top 281-290 (Pearson correlation) | Top 291-300 (Pearson correlation) |
|---|---|---|---|---|
| -0.115 | -0.162 | -0.23 | 0.021 | 0.021 |
| -0.162 | 0.021 | 0.021 | 0.021 | -0.115 |
| -0.162 | -0.23 | 0.021 | 0.021 | -0.115 |
| -0.115 | -0.23 | -0.23 | -0.23 | -0.115 |
| -0.115 | 0.021 | -0.199 | 0.021 | -0.115 |
| -0.162 | 0.021 | 0.021 | 0.021 | -0.115 |
| -0.23 | -0.23 | 0.021 | 0.021 | -0.115 |
| -0.115 | 0.021 | 0.021 | 0.021 | -0.115 |
| -0.23 | -0.23 | -0.162 | 0.021 | -0.115 |
| -0.23 | 0.021 | -0.162 | 0.021 | -0.162 |

# Supporting Texts

## Text S1. Code usage of SED

We have developed a demonstration program for the sparse screening of ECFPs and ligand-based virtual screening; the source codes and datasets are available through https://zhanglab.ccmb.med.umich.edu/SED/. The code for SED was developed in Matlab2014, Python 2.7. This provides a general framework for the screening for Lasso of ECFPs and ligand-based virtual screening, which allows users to develop their own key substructure recognition and virtual screening tools for the drug targets of their choice on the basis of our code. The program GenerateMD for the generation of ECFPs is authorized by the ChemAxon Ltd. with the free license for academic research. Due to the copyright issues, we provided the open-source cheminformatics software RDKit for generating the Morgan fingerprints which are a suitable alternative of ECFPs here.

<u>Input</u>: Compounds in the format of canonical SMILES and their bioactivity values.

<u>Output</u>: Model performance ($RMSE, r^2$).

The procedure of the pipeline is as follows: Input compounds in the format of canonical SMILES and their bioactivity values → Generate fingerprints by RDKit → Obtain top features selected by Lasso screening of fingerprints → Rebuild the ligand dataset using the selected top features → Construct DNN regression models → Obtain the model performance. The details of implementing the SED package can see the README.txt file through https://zhanglab.ccmb.med.umich.edu/SED/README.txt.

# References

Alba D.P.*, et al.* (2018) Stimulation of S1PR5 with A-971432, a selective agonist, preserves blood–brain barrier integrity and exerts therapeutic effect in an animal model of Huntington's disease, *Human Molecular Genetics*, **27**:2490-2501.

Edvinsson, L. (2008) CGRP-receptor antagonism in migraine treatment, *Lancet*, **372**, 2089-2090.

Esbenshade, T.A.*, et al.* (2008) The histamine H3 receptor: an attractive target for the treatment of cognitive disorders, *British Journal of Pharmacology*, **154**, 1166-1181.

H, C.*, et al.* (1996) The beta-chemokine receptors CCR3 and CCR5 facilitate infection by primary HIV-1 isolates, *Cell*, **85**, 1135-1148.

H Ito, C.H., L Farde (1999) Localization of 5-HT1A receptors in the living human brain using [carbonyl-11C]WAY-100635: PET with anatomic standardization technique, *Journal of Nuclear Medicine*, **40**, 102-109.

Heisler, L.K.*, et al.* (2007) Serotonin 5-HT(2C) receptors regulate anxiety-like behavior, *Genes Brain & Behavior*, **6**, 491-496.

Hu, M.*, et al.* (2015) Pharmacogenetics of cutaneous flushing response to niacin/laropiprant combination in Hong Kong Chinese patients with dyslipidemia, *Pharmacogenomics*, **16**, 1387-1397.
Jacobson, K.A.*, et al.* (2008) Flexible modulation of agonist efficacy at the human A 3 adenosine receptor by the imidazoquinoline allosteric enhancer LUF6000, *Bmc Pharmacology*, **8**, 20.

Kim, J.Y.*, et al.* (2014) Calcium-sensing receptor (CaSR) as a novel target for ischemic neuroprotection, *Annals of Clinical & Translational Neurology*, **1**, 851-866.

Liu, X.Y.*, et al.* (2011) Unidirectional cross-activation of GRPR by MOR1D uncouples itch and analgesia induced by opioids, *Cell*, **147**, 447-458.

Rivera, G.*, et al.* (2008) Melanin-concentrating hormone receptor 1 antagonists: a new perspective for the pharmacologic treatment of obesity, *Current Medicinal Chemistry*, **15**, 1025-1043.

Shrimpton, A.E.*, et al.* (2004) Molecular delineation of deletions on 2q37.3 in three cases with an Albright hereditary osteodystrophy-like phenotype, *Clinical Genetics*, **66**, 537-544.

Silverman, M.H.*, et al.* (2008) Clinical evidence for utilization of the A3 adenosine receptor as a target to treat rheumatoid arthritis: data from a phase II clinical trial, *Journal of Rheumatology*, **35**, 41-48.

Slaugenhaupt, S.A.*, et al.* (1995) Mapping of the gene for the Mel 1a-melatonin receptor to human chromosome 4 (MTNR1A) and mouse chromosome 8 (Mtnr1a), *Genomics*, **27**, 355-357.

Wang, Y.D.*, et al.* (2011) The G-Protein-coupled bile acid receptor, Gpbar1 (TGR5), negatively regulates hepatic inflammatory response through antagonizing nuclear factor kappa light-chain enhancer of activated B cells (NF-κB) in mice, *Hepatology*, **54**, 1421-1432.

Wu, J.*, et al.* (2018) WDL-RF: Predicting Bioactivities of Ligand Molecules Acting with G Protein-coupled Receptors by Combining Weighted Deep Learning and Random Forest, *Bioinformatics*,**34**, 2271-2282.