

MedChemComm

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1 Identifying Farnesoid X Receptor Agonists with Naïve Bayesian and Recursive 2 Partitioning Approaches

3
4 Qianzhi Ding^a, Chanjuan Li^a, Ling Wang^b, Yali Li^a, Huihao Zhou^a, Qiong Gu^{a*}, and Jun Xu^{a*}

5 6 Abstract

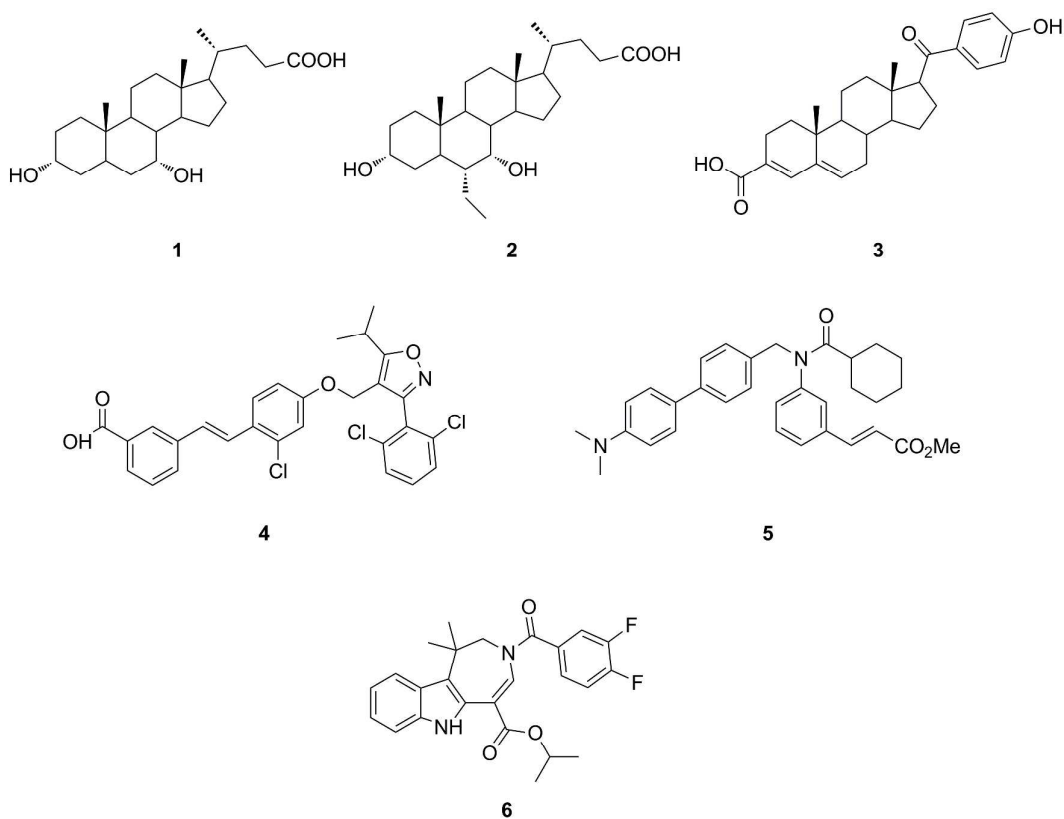
7 The farnesoid X receptor (FXR), a ligand-modulated transcription factor, is a multiple
8 functional hepatic cell protector. Therefore, FXR agonists represent promising
9 dyslipidemia and anti-diabetes agents. To identify novel FXR agonists, models were
10 created from 144 known FXR agonists with naïve Bayesian (NB) and recursive
11 partitioning (RP) approaches. The predictive and reliable models were selected with
12 Matthews correlation coefficient (MCC) criterion (>0.900 with 117 testing
13 compounds). The top 4 models were validated with the external data (282 compounds
14 having cell-free activities and 500 decoys). Two optimal FXR agonist models (one
15 from the NB method and the other from the RP method) were obtained from the top
16 models by further validations. A virtual screening campaign was conducted against
17 our in-house compound library with the optimal models and produced 15 virtual hits,
18 which were further confirmed with cell-based luciferase assays. Finally, we
19 discovered two new FXR agonists. Molecular docking studies indicated that the two
20 new FXR agonists have similar binding modes to the known FXR agonists. This work
21 demonstrated that a machine learning approach with combined NB and RP methods
22 was able to identify novel FXR agonists and that the approach could be applied in
23 other lead identification processes.

24 25 1 Introduction

26 The farnesoid X receptor (FXR), a ligand-activated transcriptional factor and multiple
27 functional hepatic cell protector, is mainly expressed in liver and small intestine¹.
28 After activation by bile acids (BAs) or other agonists, FXR binds to specific DNA
29 response elements as a heterodimer with the retinoid X receptor (RXR). Subsequently,
30 with regulating the expression of genes, such as SREPB1c (sterol regulatory element

1 binding protein 1c), PEPCK1 (Phosphoenolpyruvate carboxykinase 1), BSEP (bile
2 salt export pump), FXR is involved in bile acids, hepatic triglyceride, lipid and
3 glucose homeostasis, liver regeneration/repair, hepatocyte survival, and tumor
4 suppression². FXR also down regulates hepatic inflammation and oncogenes^{3, 4}.
5 Hence, FXR is a potential drug target for the therapy of metabolic diseases^{5, 6}.

6 Known FXR agonists include steroidal and non-steroidal agents⁷. The steroidal
7 agonists, such as BA derivatives⁸ (e.g., chenodeoxycholic acids⁹, CDCA, **1**), are
8 endogenous ligands; steroidal agonists, such as 6 α -ethyl-CDCA (6-ECDCA, **2**)¹⁰ and
9 MFA-1 (**3**)¹¹, are non-endogenous FXR agonists. Among these steroidal agonists,
10 6-ECDCA (**2**) has entered into the clinical research for treating primary biliary
11 cirrhosis (PBC) or non-alcoholic fatty liver disease (NAFLD)¹²⁻¹⁵. Moreover, a
12 number of non-steroidal FXR agonists have been found, such as GW4064 (**4**)¹⁶,
13 fexaramine (**5**)¹⁷, XL335 (**6**)¹⁸ and the others with new scaffolds¹⁹. Some of the
14 agonists have been tested in a reporter gene assay that is widely used in discovering
15 FXR agonists and exhibited potent activities²⁰. However, many of the known agonists,
16 such as GW4064 and XL335, were not druggable due to problems of intrinsic toxicity,
17 absorption, or metabolism^{7, 18, 21, 22}. More new FXR agonists with acceptable
18 pharmaceutical properties are in high demand.



1
2 **Fig 1.** The known steroidal and non-steroidal FXR agonists.

3 FXR binding pocket is flexible²³ and allows structurally diverse agonists (Fig 1)
4 with a number of different binding modes^{11, 17, 24-26}. A number of structure-based
5 pharmacophore models have been created since 2011^{21,27,28} to predict FXR agonists.
6 However, these models were ligand-scaffold dependent and referenced with a smaller
7 number of ligands (<10 ligands). Therefore, it would be difficult for them to predict
8 novel scaffold for FXR agonists.

9 To avoid the limitations of structure-based pharmacophore modeling,
10 ligand-based machine learning approaches have been successfully applied for
11 searching anti-MRSA (anti-methicillin-resistant *Staphylococcus aureus*) compounds²⁹
12 or identifying agents regulating PPAR (peroxisome proliferators-activated receptor)³⁰,
13 LXR (liver X receptor)³¹ and other proteins³². Therefore, in this paper, we built FXR
14 agonist models using multiple machine learning approaches (NB and RP) based upon
15 a larger training data set with diverse scaffolds.

16 To achieve optimal models, we generated many ligand-based models, which were
17 evaluated by cross validations and external validations. The optimal models were used

1 in a virtual screening campaign against our in-house compound library for FXR
2 agonists. The virtual screening hits were then tested with *in vitro* cell-based luciferase
3 assays.

4 **2 Materials and Methods**

5 **2.1 Data set**

6 The human FXR agonist cell-based assay data used in this work were derived from
7 the ChEMBL database (version 19)³³, and the data were selected by the following
8 criteria: (1) the data of the human FXR agonist assay were selected; (2) the data of the
9 cell-based assay were selected; and (3) duplicated compounds were removed. This
10 approach resulted in 170 human FXR agonists with EC₅₀ values ranging from 2 to
11 over 100,000 nM (that is, five-order of magnitude). One hundred forty-four of the 170
12 human FXR agonists were marked as “active” (EC₅₀ values were under or equal to 5
13 μM); the remaining 26 compounds were marked as “inactive” (EC₅₀ values were
14 greater than 5 μM). The activity threshold was set at 5 μM (see Fig S1 in the
15 Supplementary information for more details).

16 **2.2 Decoy generation**

17 Decoys data were generated from DUD-E (a database of useful decoys: enhanced)³⁴
18 (<http://dude.docking.org/>), and added to the training data to keep it balanced. Ten
19 diverse structures were selected from the “active” part of the database using the
20 diverse molecules module in Pipeline Pilot 7.5 (Accelrys, Inc., San Diego, CA.).
21 Subsequently, these 300 decoy structures were generated by calculating their
22 molecular properties based upon the 10 reference structures in the DUD-E server.
23 Three hundred decoys marked “inactive” were added into the database, and the whole
24 data set was optimized using MOE 2013.08 (Chemical Computing Group Inc.) based
25 on the MMFF94 force field³⁵. All structures were saved as MACCS (Molecular
26 ACCess System) sdf files and SMILES (Simplified molecular input line entry
27 specification) files. Finally, the whole database was divided into two parts, a training
28 set (353) and a test set (117), based on the random algorithm in Discovery Studio
29 2.5.5 (DS2.5.5, Accelrys, Inc., San Diego, CA.). The number of molecules in the
30 training set was three times as many as that in the test set. This proportion was

1 employed in reference³⁶.

2 **2.3 Calculation of molecular properties**

3 The computed molecular properties (MP) were molecular weight (MW), the
4 octanol/water partitioning coefficient (ALogP) based on the Ghose and Crippen's
5 method, the molecular solubility (Molecular_Solubility), the apparent partition
6 coefficient at pH = 7.4 (LogD) based on the Csizmadia's method, the molecular
7 surface area (MSA), the molecular polar surface area (MPSA), the molecular
8 fractional polar surface area (MFPSA), the number of rings (nR), the number of
9 aromatic rings (nAR), the number of hydrogen bond donors (nHBDdon), the number of
10 hydrogen bond acceptors (nHBAcc), the count of oxygen and nitrogen (NPlusO), and
11 the number of rotatable bonds (nRB). These values were all calculated with DS 2.5.5.

12 **2.4 Calculation of molecular fingerprints**

13 Two sets of fingerprints, SciTegic extended-connectivity fingerprints (ECFP, FCFP
14 and LCFP) and Daylight-style path-based fingerprints (EPFP, FFPF and LPFP), were
15 calculated using DS 2.5.5. Each type of fingerprint was used in four diameters: 4, 6, 8,
16 and 10. All of these fingerprints are frequently applied in ADME, QSAR (quantitative
17 structure–activity relationship), and QSPR (quantitative structure-property
18 relationship) models^{36, 37}.

19 **2.5 Naïve Bayesian**

20 Naïve Bayesian is a simple probabilistic classification approach based on Bayes'
21 theorem. Naïve Bayesian is highly scalable and unsupervised in a learning problem.

22 The core function is eq. 1.

$$23 \quad P(H|E) = \frac{P(E|H)P(H)}{P(E)} \quad (1),$$

24 where, H is the hypothesis or model, E is the observed data, P(H) is the probability of
25 hypothesis H before observing any data, P(E) is the marginal probability of the data,
26 and P(H|E) is the probability that the hypothesis H is correct for the observed data.
27 P(E|H) is the likelihood that the probability of data E if hypothesis H is true. More
28 details can be found in reference³⁸. In this work, a Laplacian-corrected Bayesian
29 classifier algorithm (implemented in DS2.5.5) was applied for building Bayesian

1 models. In our case, the models were trained with both the agonist (“active”) and
2 non-agonist (“inactive”) data and considered each of the MPs and molecular
3 fingerprint as the features to gain the knowledge to distinguish active from inactive.
4 This building process is unbiased and takes the complexity of the model into
5 consideration, which can avoid the over-fitting problem.

6 **2.6 Recursive partitioning**

7 Recursive partitioning (RP) is a type of accurate and comprehensible classification
8 method that is used to discover the relationship between a dependent property (Y
9 variable) and a number of independent properties (X variables). A decision tree will
10 be created to classify the data points in the training set when RP proceeds. RP is a
11 dichotomous process that divides independent variables (fingerprints and MPs). All
12 RP models were built based on 12 fingerprints and 13 molecular descriptions in this
13 study. Subsequently, 5-fold cross-validation was employed to determine the degree of
14 pruning, which was required for the best predictive model. More details can be found
15 in reference³⁹.

16 **2.7 Evaluation of the model performance**

17 To evaluate the performance of Bayesian and RP classifiers, 5-fold cross-validation
18 was used in this study. A set of evaluation indexes, including true positives (TP), true
19 negatives (TN), false positives (FP), false negatives (FN), sensitivity (SE), specificity
20 (SP), the prediction accuracy for agonists (Q_a), the prediction accuracy for
21 non-agonists (Q_{na}), overall predictive accuracy (Q), and the Matthews correlation
22 coefficient (MCC), were calculated with the formulas (2) to (7), and the receiver
23 operating characteristic (ROC) curve was plotted. The area under the curves (AUC),
24 which represents the classification ability of a binary classifier, was calculated
25 through iteratively seeking the proper classifier threshold⁴⁰.

$$26 \quad SE = \frac{TP}{TP+FN} \quad (2)$$

$$27 \quad SP = \frac{TN}{TN+FP} \quad (3)$$

$$28 \quad Qa = \frac{TP}{TN+FP} \quad (4)$$

$$Qna = \frac{TN}{TN+FN} \quad (5)$$

$$Q = \frac{TP+TN}{TP+FN+TN+FP} \quad (6)$$

$$MCC = \frac{TP \times TN - FN \times FP}{\sqrt{(TP+FN)(TP+FP)(TN+FN)(TN+FP)}} \quad (7)$$

The MCC values are the measures for the classification accuracies of the models.

2.8 Cell culture

HEK293T cells were cultured in Dulbecco's modified Eagle's medium (DMEM) supplemented with 10% fetal bovine serum (FBS), 100 U/ml penicillin and 100 μ g ml⁻¹ streptomycin at 37°C in 5% CO₂ (V/V). The tested compounds were dissolved in DMSO and supplemented at indicated concentrations.

2.9 Transfection and luciferase assay

Activation studies on FXR were performed according to the method of Andrea A. Cronican⁴¹ with a few modifications. To be brief, the HEK-293T cells were seeded into 96-well plates at 3×10^4 cells per well and allowed to attach overnight at 37°C. Plasmids pSG5/hFXR and pSG5/hRXR α , reporter plasmid pGL3/(DR-4)-c-fos-FF-luc, and the internal control plasmid pCMV/Renilla-luc were kindly gifts from Prof. Qing Song (University of Science and Technology, Beijing, China)⁴². These plasmids were co-transfected into cells using LipofectamineTM 2000 (Invitrogen, USA) in accordance with the manufacturer's instructions. After 10 hours, cells were treated with tested compounds. FXR agonist GW4064 was used as a positive control; 0.1% DMSO was taken as vehicle. Luminescence measurements were processed 20 hours later. The results are expressed as relative firefly luciferase activity normalized to the renilla luciferase activity (fold change compared to vehicle control).

2.10 Molecular docking

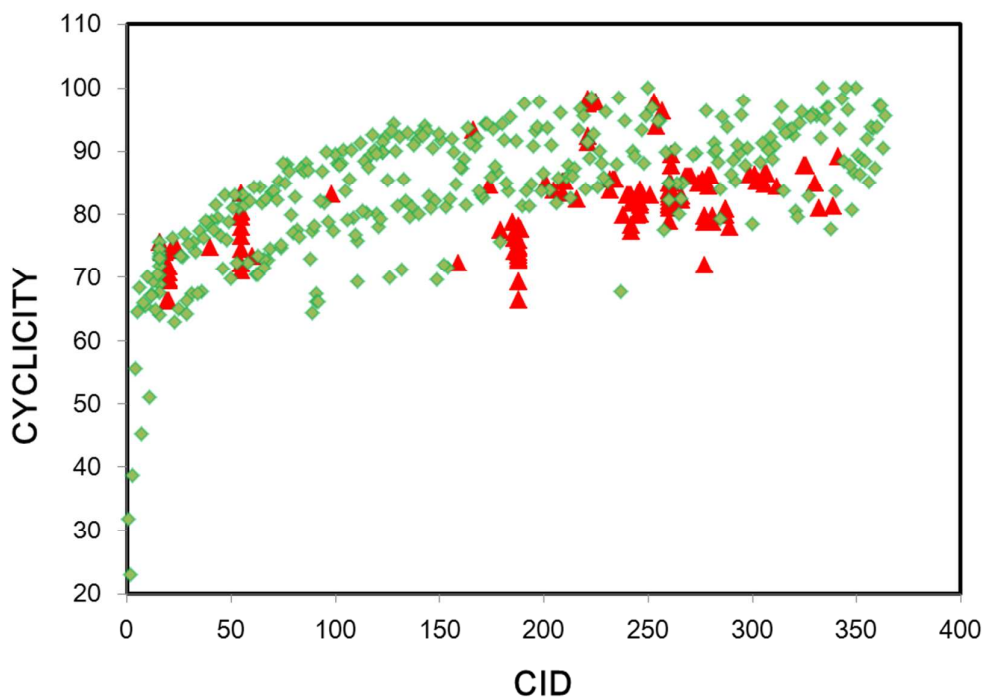
Molecular docking was employed to gain an insight into the binding modes of two active compounds and FXR. The FXR-GW4064 crystal complex (PDB code: 3DCT) was downloaded from Protein Data Bank (PDB, <http://www.rcsb.org/pdb/home/home.do>) and prepared using the protein preparation protocol in the Schrödinger 2013.01. The extra precision (XP) mode in the Glide 5.9⁴³⁻⁴⁵ of the Schrödinger software suite was employed to study the binding modes.

1 The docking parameters were all validated using re-docking methods. Two active
2 compounds were prepared by Ligprep module in the Schrödinger software suite and
3 docked into the FXR crystal structure in XP mode.

4 **3 Results and discussion**

5 **3.1 Chemical space and structural diversity analysis**

6 The structural diversity of the training and testing sets has a significant influence on
7 the reliability and predictive ability of the models. In this study, an S-cluster approach
8 (SCA)⁴⁶ (*in-house* software) was employed to measure the structural diversity (Fig 2).
9 The cluster ID (CID), the serial number to each compound cluster, is proportional to
10 the chemical structure complexity. More CIDs indicate higher structural diversity. The
11 cyclicity is the metric of the cyclic degree of a molecule. The higher cyclicity value
12 indicates the molecule has fewer/shorter substituents. Fig 2 demonstrates that the data
13 points from both FXR agonists and non-agonists are widely spread, indicating that the
14 structures of the 470 molecules are diverse.



15
16 **Fig 2.** The SCA-plot of 470 compounds for model building. Red: FXR agonists.
17 Green: FXR non-agonists.

18 **3.2 Correlation analyses of molecular properties and FXR agonist activity**

19 Thirteen physicochemical properties, including MW, ALogP, Molecular_Solubility,

1 LogD, MSA, MPSA, MFPSA, nR, nAR, nHBDOn, nHBAcc, NPlusO and nRB, were
 2 calculated. The correlations between the 13 properties and FXR agonist activity were
 3 measured using correlation coefficients (R), and the significances of the difference
 4 between paired samples were evaluated with the student's *t* test (*p*-value) as listed in
 5 Table 1. LogD, MW, Molecular_Solubility, nAR, MSA, and MFPSA were
 6 significantly different from others. Molecular_Solubility was identified as the best
 7 property to discriminate FXR agonists and non-agonists, although it was not a strong
 8 predictor. ALogP, LogD, nR, nHBDOn and nHBAcc exhibited relatively higher
 9 correlations with FXR agonist activity. Because no molecular property had a
 10 significantly higher value for FXR agonist activity, multiple molecular properties had
 11 to be used in building machine learning models.
 12 **Table 1.** The MPs and their relationships with FXR agonist activity (R) and *p*-values
 13 (significances).

MP	<i>p</i> -value ^a	R ^b
ALogP	5.64×10^{-7}	0.154
LogD	7.98×10^{-11}	0.139
Molecular_Solubility	3.22×10^{-20}	0.080
MW	4.75×10^{-11}	0.045
NPlusO	2.36×10^{-3}	0.078
nRB	1.77×10^{-6}	0.076
nR	6.16×10^{-3}	0.100
nAR	9.30×10^{-10}	0.085
nHBAcc	9.22×10^{-3}	0.131
nHBDOn	5.39×10^{-2}	0.153
MSA	5.31×10^{-14}	0.074
MPSA	1.08×10^{-3}	0.077
MFPSA	1.31×10^{-14}	0.085

14 ^a*p*-value: the statistical significance between FXR agonists and non-agonists.

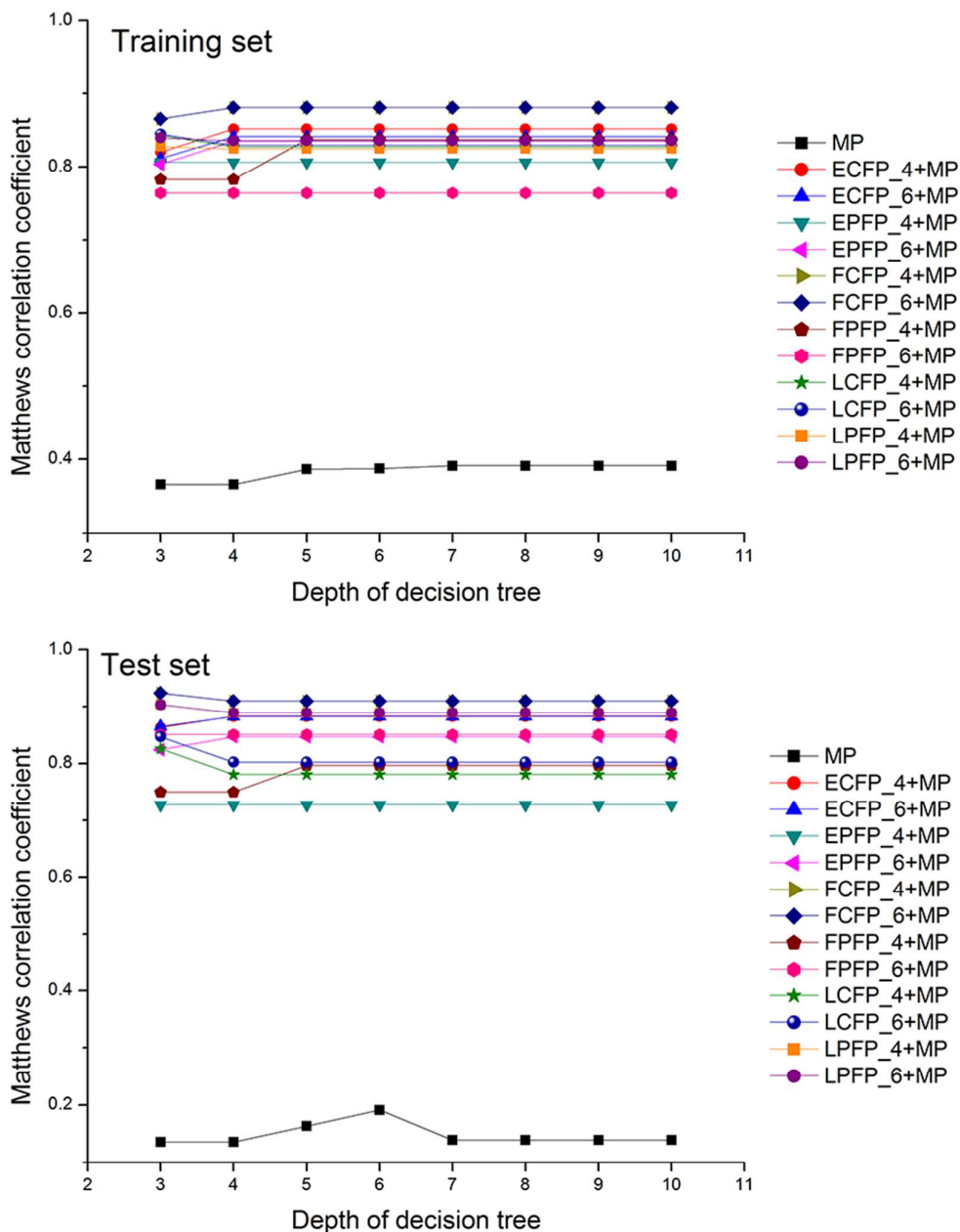
15 ^bR: the correlation coefficient between a descriptor and EC₅₀ (FXR).

1 3.3 Recursive partitioning models

2 The decision tree generated with RP is more intuitive compared with those “blind
3 modeling” approaches, such as ANN (Artificial Neural Network) and SVM (Support
4 Vector Machine). A deeper decision tree is more accurate, but it may cause
5 over-fitting problems. A shorter decision tree may increase the possibility of applying
6 the tree to new data sets, but it may reduce the accuracy of the prediction³⁶. To
7 optimize the depth of a decision tree for the best prediction performance, a number of
8 experiments with depth thresholds ranging from 3 to 10 were tried. A total of 104 RP
9 models were built and evaluated with evaluation indexes. The 5-fold cross-validation
10 method was used to measure the robustness of those models.

11 With increasing depth thresholds from 3 to 10, 8 decision trees were built using
12 13 molecular properties. The MCC values of a test set indicated that the tree with the
13 depth of 6 reached the best performance. The evaluation indexes of the best RP model
14 based on MPs are listed in Table 2. In Table 2, MP represents 13 descriptors
15 calculated by DS 2.5.5; Depth* represents the best tree depth for the corresponding
16 RP model. For the training set, the best model with depth 6 achieves a sensitivity of
17 81.4%, specificity of 61.4%, MCC value of 0.387, and an AUC value of 0.716. For
18 the test set, the performance was poor according to those evaluation indexes
19 ($SE_{\text{test}}=68.4\%$, $SP_{\text{test}}=51.9\%$, $MCC=0.191$, and $AUC=0.607$, Table 2). All results of
20 the test set suggested that the best RP model based on molecular properties is limited
21 in distinguishing agonist from non-agonist because MPs represent whole molecular
22 structure contributions, not sub-structural contributions. To take sub-structural
23 contributions into account, molecular fingerprints must be taken into consideration.
24 Therefore, 96 RP models were generated using MPs and 12 sets of molecular
25 fingerprints.

26 As shown in Fig 3, the RP models derived from MPs and molecular fingerprints
27 have much better MCC values than those derived only from MPs. The differences of
28 the MCC values are more significant when the models were validated with the test
29 data set.



1
 2 **Fig 3.** The relationships between the Matthews correlation coefficient values and the
 3 decision tree depths, and descriptors. MP: the descriptors consist of only MPs
 4 (calculated with DS 2.5.5.). *+MP: the descriptors consist of different fingerprints
 5 plus MPs.

6 Table 2 lists the validation parameters for all of the best RP models tested with
 7 the test set and training set. In this table, RP models using FCFP_4 and FCFP_6
 8 fingerprints have the highest MCC values (0.924) for the test set. The best decision
 9 tree depth of both models is 3. The two models have the same sensitivities (97.4%),

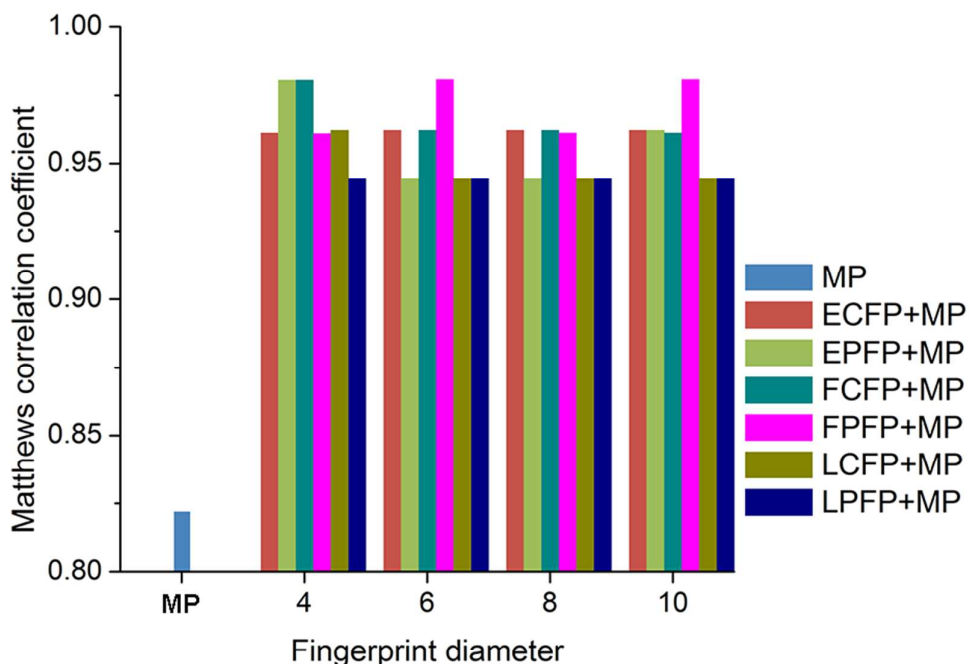
- 1 specificities (96.2%) and AUC values (0.975) for the test set.
- 2 **Table 2.** Performance of the best RP models with the combination of different
- 3 fingerprints and MPs.

Models	Training set								Test set							
	TP	FN	TN	FP	SE	SP	MCC	AUC	TP	FN	TN	FP	SE	SP	MCC	AUC
MP ^a _depth ^b	83	19	154	97	0.814	0.614	0.387	0.716	26	12	41	38	0.684	0.519	0.191	0.607
ECFP_4_depth4+MP	94	8	237	14	0.922	0.944	0.852	0.939	33	5	78	1	0.868	0.987	0.883	0.956
ECFP_6_depth4+MP	95	7	234	17	0.931	0.932	0.841	0.934	34	4	77	2	0.895	0.975	0.882	0.98
EPFP_4_depth3+MP	97	5	225	26	0.951	0.896	0.806	0.922	34	4	68	11	0.895	0.861	0.727	0.894
EPFP_6_depth4+MP	95	7	233	18	0.931	0.928	0.835	0.947	35	3	74	5	0.921	0.937	0.847	0.968
FCFP_4_depth3+MP	95	7	238	13	0.931	0.948	0.865	0.93	37	1	76	3	0.974	0.962	0.924	0.975
FCFP_6_depth3+MP	95	7	238	13	0.931	0.948	0.865	0.93	37	1	76	3	0.974	0.962	0.924	0.975
FPPF_4_depth5+MP	96	6	232	19	0.941	0.924	0.837	0.921	35	3	71	8	0.921	0.899	0.796	0.945
FPPF_6_depth3+MP	93	9	223	28	0.912	0.888	0.764	0.908	36	2	73	6	0.947	0.924	0.851	0.946
LCFP_4_depth3+MP	96	6	233	18	0.941	0.928	0.843	0.933	34	4	74	5	0.895	0.937	0.826	0.917
LCFP_6_depth3+MP	97	5	232	19	0.951	0.924	0.845	0.941	35	3	74	5	0.921	0.937	0.847	0.93
LPFP_4_depth3+MP	94	8	233	18	0.922	0.928	0.828	0.936	36	2	76	3	0.947	0.962	0.903	0.965
LPFP_6_depth3+MP	94	8	235	16	0.922	0.936	0.840	0.936	36	2	76	3	0.947	0.962	0.903	0.965

- 4 ^aMP: the 13 descriptors calculated with DS 2.5.5. ^bDepth*: the best tree depth for
- 5 the corresponding model.

6 **3.4 Naïve Bayesian models**

- 7 One NB model was derived from the MPs calculated with DS 2.5.5 software.
- 8 Twenty-four models were derived from 13 MPs combined with different types of
- 9 molecular fingerprints (four diameters and six types). The MCC values for all NB
- 10 models are depicted in Fig 4. The NB MP-only model has much lower MCC value
- 11 than the MCC values of NB models derived from the combination of MPs and
- 12 fingerprints. This is consistent with the cases of RP models (Fig 3).



1

2 **Fig 4.** The Matthews correlation coefficient values for the NB MP model and the NB
3 models derived from MP+fingerprints with different diameters.

4 The performance parameters for the top NB models are listed in Table 3. For the
5 test data set, the NB models achieved the same performance (MCC = 0.981). Because
6 greater diameter fingerprint requires higher computation resource, the best NB model
7 was determined to be NB_FPFP_6+MP, which achieves a sensitivity of 100.0%,
8 specificity of 98.7%, prediction accuracies for FXR agonist class of 97.4% and, an
9 AUC value of 0.999. For the training set, the NB_FPFP_6+MP model achieves
10 $SE_{\text{training}}=98.0\%$, $SP_{\text{training}}=95.2\%$, $MCC_{\text{training}}=0.908$, and $AUC = 0.987$.

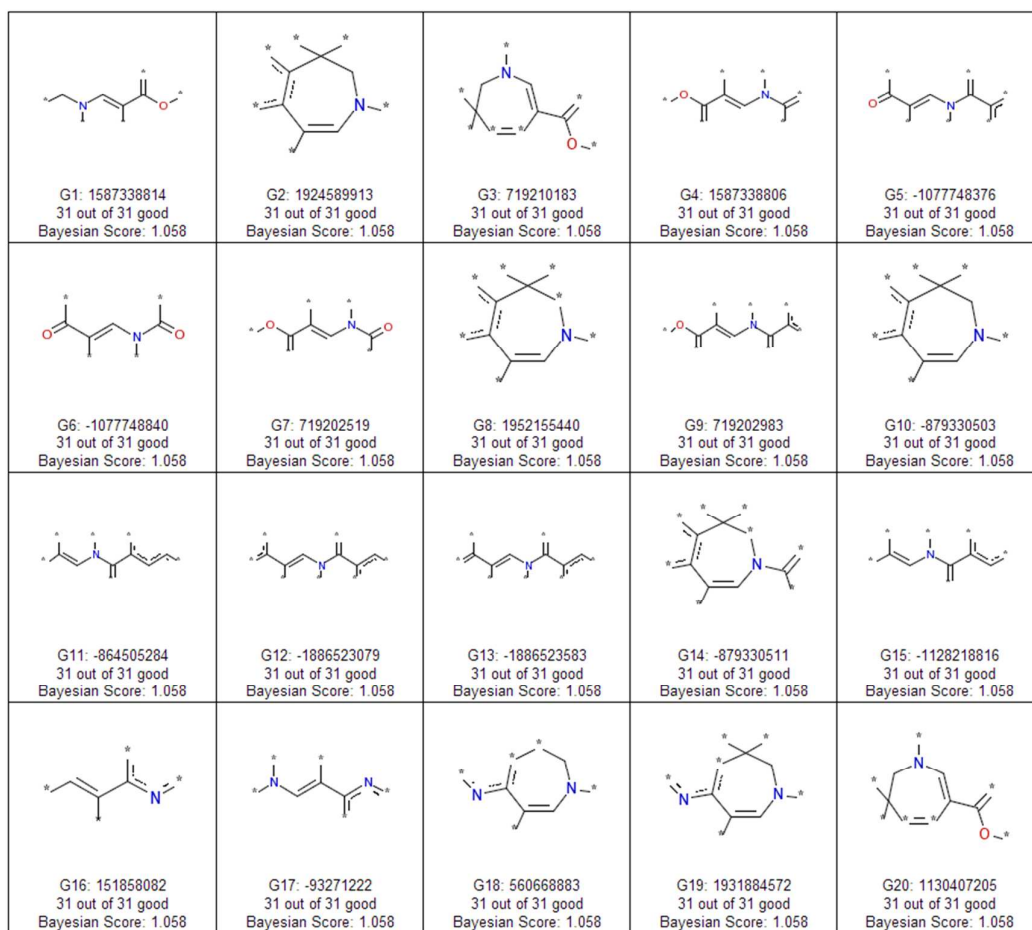
11 **Table 3.** Performance parameters for the best RP and NB models

Models	Training set								Test set							
	TP	FN	TN	FP	SE	SP	MCC	AUC	TP	FN	TN	FP	SE	SP	MCC	AUC
RP_FCFP_4_depth3 ^b +MP ^a	95	7	238	13	0.931	0.948	0.865	0.930	37	1	76	3	0.974	0.962	0.924	0.975
RP_FCFP_6_depth3+MP	95	7	238	13	0.931	0.948	0.865	0.930	37	1	76	3	0.974	0.962	0.924	0.975
NB_FPFP_6+MP	100	2	239	12	0.980	0.952	0.908	0.987	38	0	78	1	1.000	0.987	0.981	0.999
NB_FPFP_10+MP	100	2	236	15	0.980	0.940	0.890	0.984	38	0	78	1	1.000	0.987	0.981	1.000

12 ^aMP: the 13 descriptors calculated with DS 2.5.5. ^bDepth*: the best tree depth for the
13 corresponding model.

1 3.5 Interpreting fingerprint modeling results

2 Structural fragments that make positive contributions to FXR agonists can be derived
 3 from the best NB model. These fragments (or privileged fragments) were exported
 4 from the top- n ($n > 0$) fragments that have $P(H|E)$ values greater than zero. The top-20
 5 privileged fragments are listed in Fig 5 and represent the guidelines for FXR agonist
 6 design, virtual screening or lead optimizations. Some of the privileged fragments can
 7 be merged, such as, G2, G8 and G10 belong to the same fragment family; G3 and
 8 G20 belong to another family; etc. Many privileged fragments are alkaloids with
 9 conjugated double bonds system. Unsaturated seven-membered ring alkaloids are
 10 privileged scaffolds for FXR agonists.



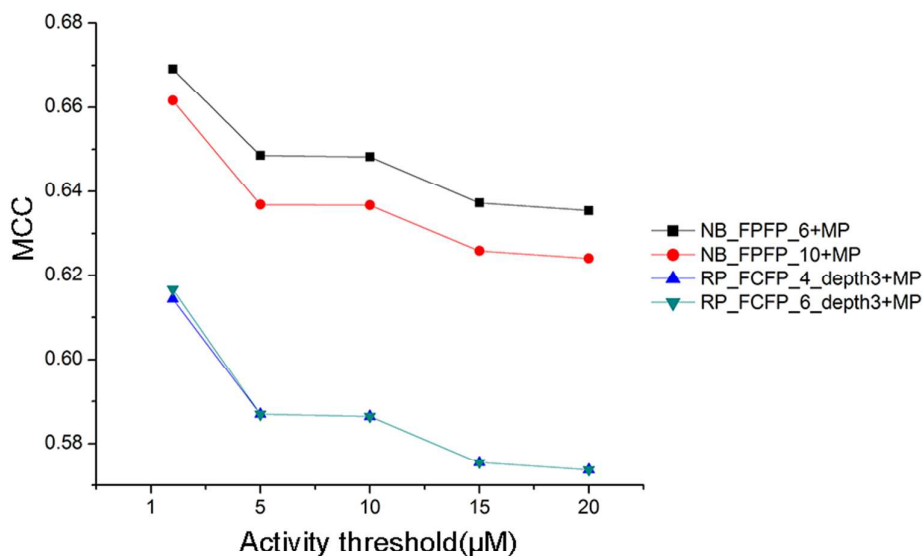
11
 12 **Fig 5.** Privileged fragments exported from the best NB model, which was created
 13 from FFPF_6 fingerprints in DS 2.5.5.

14 3.6 Validating the best models with external cell-free data

15 In FXR agonist assay experiments, cell-free assays are more confirming than the

1 cell-based assays. Therefore, we use the FXR cell-free assay data to validate the
2 models and determine the final FXR agonist predictive model.

3 The top-4 models, two RP models (RP_FCFP_4+MP and RP_FCFP_6+MP) and
4 two NB models (NB_FFPF_6+MP and NB_FFPF_10+MP), were validated with an
5 external dataset including 282 cell-free activity data and 500 decoy compounds. The
6 5-fold cross-validation was employed in this test. The external dataset was divided
7 into five sub-datasets with five activity thresholds (1 μ M, 5 μ M, 10 μ M, 15 μ M, and
8 20 μ M). Twenty MCC values for the twenty test cases (four models by five
9 sub-datasets) are depicted in Fig 6. One micromolar was determined to be the best
10 activity threshold to define active compounds. This is also consistent with the active
11 threshold definition (5 μ M) we used when we built the models. Hence, the NB model,
12 NB_FFPF_6+MP, is the best FXR agonist predictive model. The performance data of
13 NB_FFPF_6+MP can be found in Table 3.

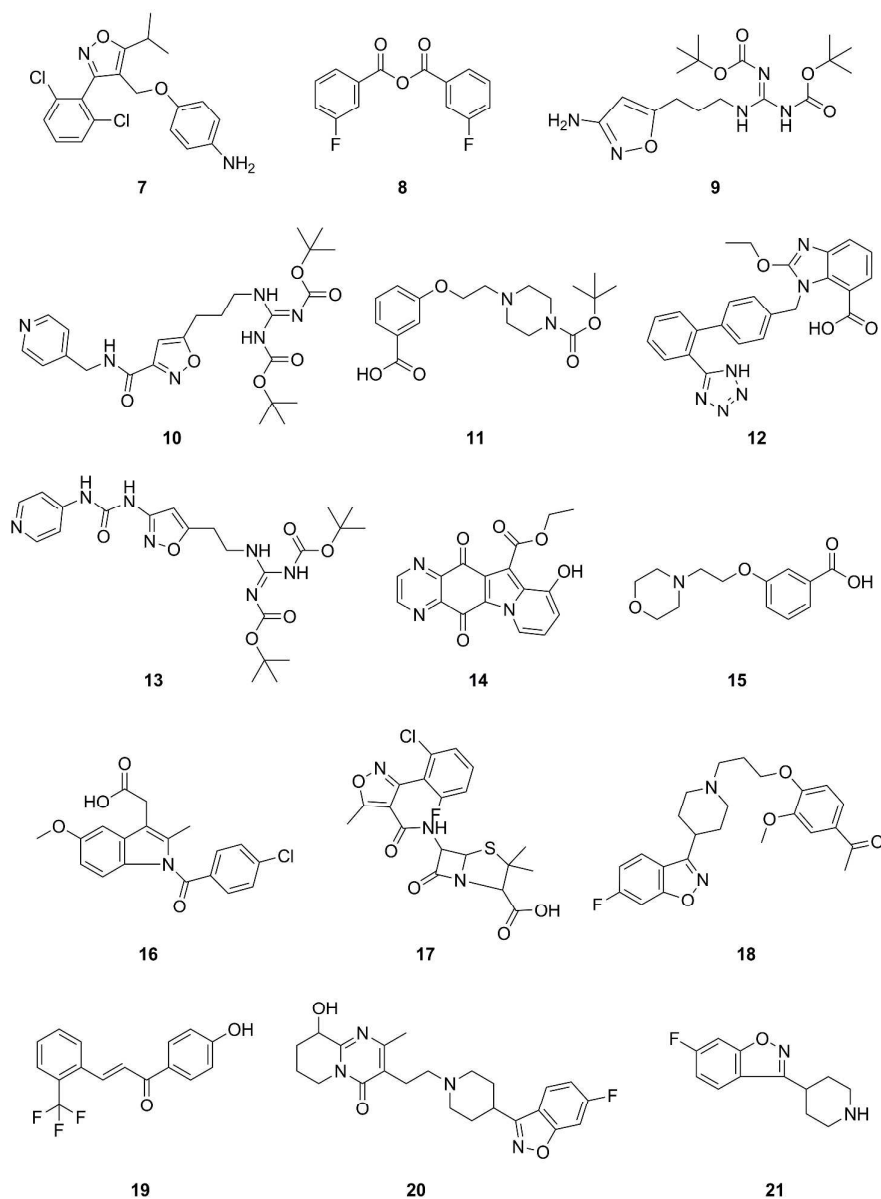


14
15 **Fig 6.** External validation: The relationship between MCC and activity threshold. The
16 top line is for the best predictive model.

17 3.7 Virtual screening FXR agonists with the best models

18 Although NB_FFPF_6+MP was the most recommended FXR agonist predictive
19 model, the virtual screening campaign still combined the results from the best model

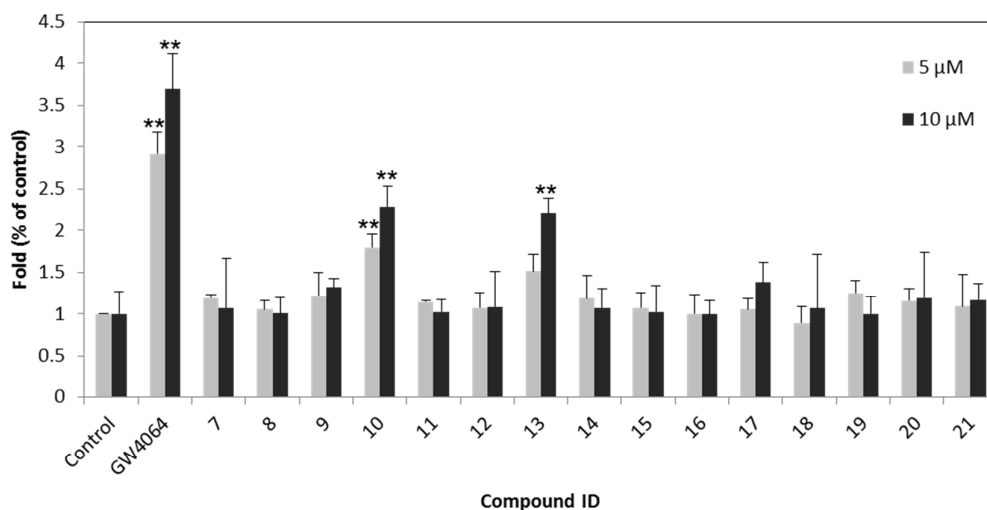
1 of the RP approach (RP_FCFP_4_depth3+MP) to avoid potential false negatives. The
2 virtual library is our in-house library, the Guangdong small molecule tangible library
3 (GSMTL)⁴⁷, which has more than 7,500 chemical compounds with average
4 purities >95%. The virtual screening resulted in 195 virtual hits (162 hits from NB,
5 and 33 hits from RP). According to a previously study²⁹, 57 compounds with simple
6 scaffolds and low molecular weight (<200) were abandoned. Finally, 15 compounds
7 (Fig 7) were picked for an *in vitro* cell-based luciferase assay considered the diversity
8 of the scaffolds and their availability.



9
10 **Fig 7.** The 15 compounds confirmed with cell-based luciferase assays.

1 3.8 *In vitro* cell-based bioassay results

2 GW4064 was used as a positive control¹⁶ for assaying the 15 compounds. The results
3 are represented in relative firefly luciferase activities normalized to the renilla
4 luciferase activities. The computational formula of the fold-activation is (firefly
5 luciferase activities / renilla luciferase activities of test compound) / (firefly luciferase
6 activities / renilla luciferase activities of control). The agonist activities of the 15
7 compounds are depicted in Fig 8 as a bar chart. Compounds **10** and **13** significantly
8 activated FXR relative to the blank control. The EC₅₀ values of compounds **10** and **13**
9 are 15.39 and 29.94 μM. The activation curves of compounds **10** and **13** are shown in
10 Fig S3, which exhibits a clear dose-dependent effect. The EC₅₀ values are not very
11 strong. This may be due to the poor bioavailability, which can be improved in the lead
12 optimization process.



13

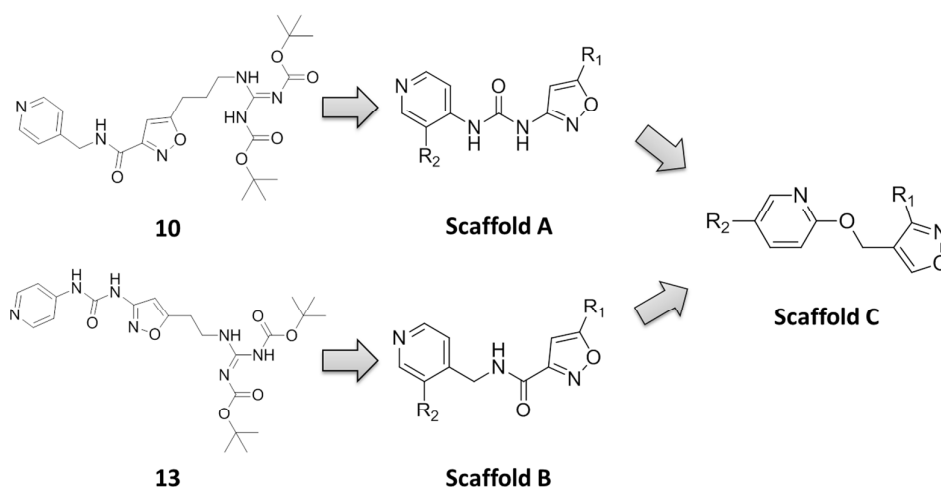
14 **Fig 8.** Bar chart for cell-based assay against FXR. The data are presented as the mean
15 ± SE. Fold = (firefly luciferase activities / renilla luciferase activities of test
16 compound) / (firefly luciferase activities / renilla luciferase activities of control).

17 Compared with the control: ** $p < 0.01$; (n = 3).

18 3.9 Scaffold analyses

19 By inspecting the structures of compounds **10** and **13**, we suggest new scaffolds A and
20 B (Fig 9) for FXR agonists. It is worth noting that scaffolds A and B are topologically

1 different, but structurally similar to each other. Both scaffolds A and B can be traced
 2 back to a known scaffold C that was derived from the training data set by means of
 3 the SCA⁴⁶ method. The old and new scaffolds all have the isoxazole ring and pyridine
 4 ring. However, these rings are connected by different links in different substitute
 5 positions. Furthermore, the side chains of compounds **10** and **13** have a similar
 6 structure to the link-like privileged fragments derived from the best NB models.
 7 These results demonstrate that established models are powerful in discovering FXR
 8 agonists with new scaffolds, and that the privileged fragments can guide the virtual
 9 screening of FXR agonists.



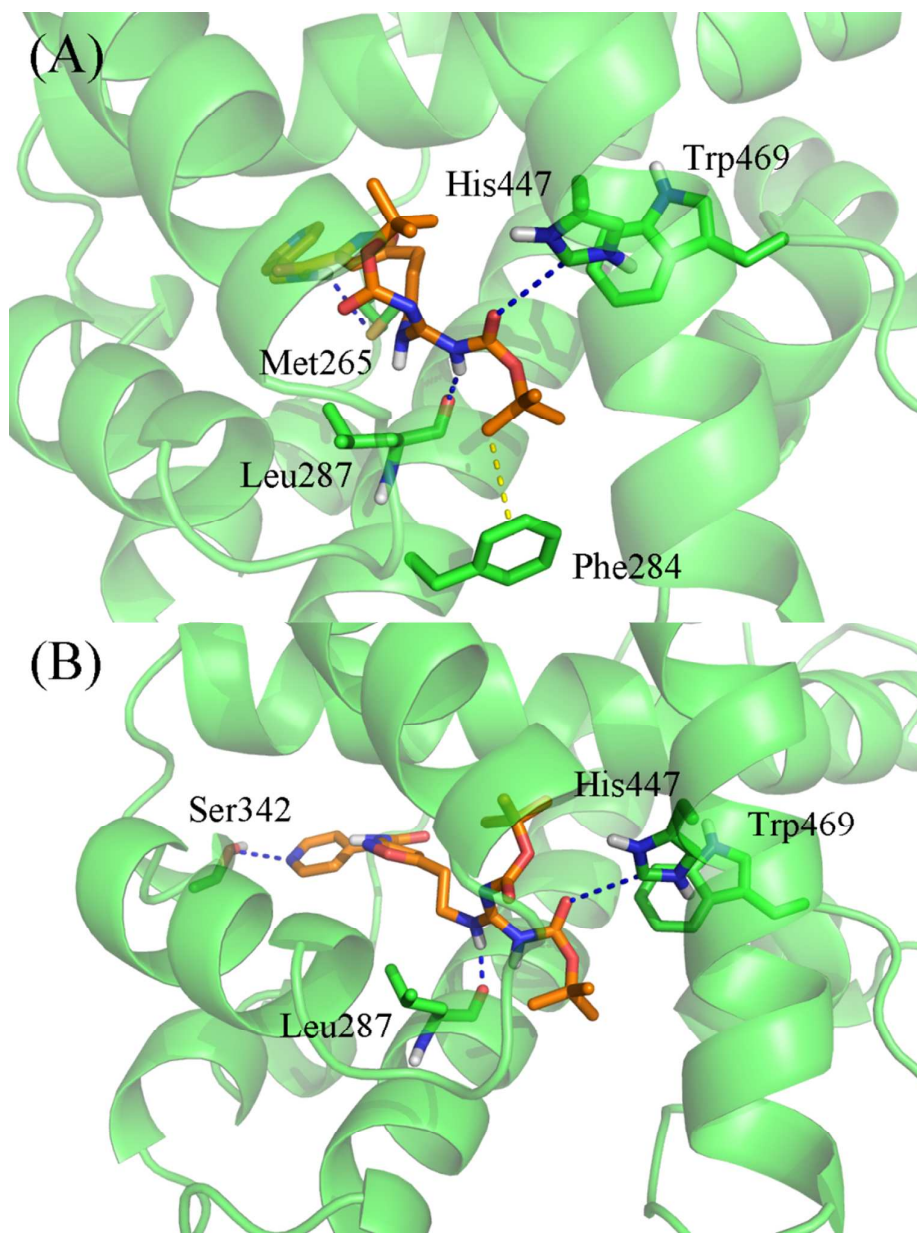
10
 11 **Fig 9.** The scaffold analyses for compounds **10** and **13**.

12 **3.9 Binding analyses for compounds 10 and 13**

13 To ensure the XP mode of Glide 5.9 (Schrödinger, Inc.) is good for docking a ligand
 14 to FXR, the crystallized ligand was extracted from an experimental FXR co-crystal
 15 structure (PDB code: 3DCT), and the ligand was docked back to the FXR structure.
 16 This resulted in a number of ligand-FXR complexes. The average RMSD of the
 17 top-10 docking poses was 1.06Å, indicating that the XP mode of Glide 5.9 is suitable
 18 for FXR system.

19 Compounds **10** and **13** were prepared with Ligprep module, and docked to the
 20 FXR structure by means of the XP mode of Glide 5.9. The proposed binding modes
 21 are depicted in Fig 10. Both compounds interact with His447, which is consistent with
 22 the FXR-GW4064 complex²⁵ (PDB code: 3DCT) and FXR-MFA-1 complex¹¹ (PDB

1 code: 3BEJ). This interaction stabilizes the activation conformation of helix 12.
2 Furthermore, the interaction of compound **10** with Met265 is consistent with the
3 interaction of the existing FXR-GSK-8062 complex²⁵ (PDB code: 3DCU). The
4 binding interactions further support the experimental data.



5
6 **Fig 10.** Binding analyses. A: Binding mode for compound **10**; B: Binding mode for
7 compound **13**. The blue dashed line represents the H-bond interaction; the yellow
8 dashed line represents the CH- π interaction.

9 **4 Conclusions**

1 It is a significant challenge to predict new FXR agonists because the FXR binding site
2 is highly flexible. RP and NB approaches can be employed in building FXR agonist
3 predictive models to avoid the problem caused by the flexibility. The keys for these
4 ligand-based approaches are to identify proper descriptors. This study demonstrates
5 the following: (1) the descriptors composed from the combinations of MPs and
6 fingerprints are better than MPs alone; (2) privileged structural fragments can be
7 derived from the best models using structural fingerprints that can serve as guidelines
8 for FXR agonist design; and (3) the naïve Bayesian approach seems capable of
9 producing better models. However, to avoid potential false negatives, we suggest that
10 the best models from both NB and RP approaches are used, as they may generate
11 similar hits with very different topological scaffolds.

12

13 **Corresponding Author**

14 *Phone: +86-20-39943077. Fax: +86-20-39943077. E-mail:

15 guqiong@mail.sysu.edu.cn (Q.G.).

16 *Phone: +86-20-39943023. Fax: +86-20-39943023. Email:

17 junxu@biochemomes.com. (J.X.).

18

19 **Acknowledgement**

20 This work was supported by the National Science Foundation of China (81173470,
21 81473138), National High Technology Research and Development Program of China
22 (863 Program, 2012AA020307), National Supercomputer Center in Guangzhou
23 (2012Y2-00048/2013Y2-00045, 201200000037), the introduction of innovative R&D
24 team program of Guangdong Province (2009010058), Guangdong Provincial Key
25 Laboratory of Construction Foundation (2011A060901014). We thank Prof. Qing
26 Song (University of Science and Technology, Beijing, China) for kindly providing
27 plasmids.

28

29 **Notes**

30 *^aSchool of Pharmaceutical Sciences & Institute of Human Virology, Sun Yat-Sen*

- 1 *University, 132 East Circle Road at University City, Guangzhou, 510006, China.*
2 ^b*Pre-Incubator for Innovative Drugs & Medicine, School of Bioscience and*
3 *Bioengineering, South China University of Technology, Guangzhou 510006, China.*
4 † Electronic Supplementary Information (ESI) available.

5

6 **References**

- 7 1. B. M. Forman, E. Goode, J. Chen, A. E. Oro, D. J. Bradley, T. Perlmann, D. J. Noonan, L. T. Burka, T.
8 McMorris, W. W. Lamph, R. M. Evans and C. Weinberger, *Cell*, 1995, 81, 687-693.
9 2. Y. Jiao, Y. Lu and X. Y. Li, *Acta Pharmacol Sin*, 2015, 36, 44-50.
10 3. M. Makishima, A. Y. Okamoto, J. J. Repa, H. Tu, R. M. Learned, A. Luk, M. V. Hull, K. D. Lustig, D. J.
11 Mangelsdorf and B. Shan, *Science*, 1999, 284, 1362-1365.
12 4. X. F. Huang, W. Y. Zhao and W. D. Huang, *Acta Pharmacol Sin*, 2015, 36, 37-43.
13 5. S. Fiorucci, A. Mencarelli, E. Distrutti, G. Palladino and S. Cipriani, *Current medicinal chemistry*,
14 2010, 17, 139-159.
15 6. S. Fiorucci, A. Mencarelli, E. Distrutti and A. Zampella, *Future medicinal chemistry*, 2012, 4,
16 877-891.
17 7. D. Merk, D. Steinhilber and M. Schubert-Zsilavecz, *Future medicinal chemistry*, 2012, 4,
18 1015-1036.
19 8. D. J. Parks, S. G. Blanchard, R. K. Bledsoe, G. Chandra, T. G. Consler, S. A. Kliewer, J. B. Stimmel, T.
20 M. Willson, A. M. Zavacki, D. D. Moore and J. M. Lehmann, *Science*, 1999, 284, 1365-1368.
21 9. D. W. Russell, *Annual review of biochemistry*, 2003, 72, 137-174.
22 10. R. Pellicciari, S. Fiorucci, E. Camaioni, C. Clerici, G. Costantino, P. R. Maloney, A. Morelli, D. J.
23 Parks and T. M. Willson, *Journal of medicinal chemistry*, 2002, 45, 3569-3572.
24 11. S. M. Soisson, G. Parthasarathy, A. D. Adams, S. Sahoo, A. Sitlani, C. Sparrow, J. Cui and J. W.
25 Becker, *Proceedings of the National Academy of Sciences of the United States of America*, 2008,
26 105, 5337-5342.
27 12. A. H. Ali, E. J. Carey and K. D. Lindor, *Annals of translational medicine*, 2015, 3, 5.
28 13. L. Adorini, M. Pruzanski and D. Shapiro, *Drug Discov Today*, 2012, 17, 988-997.
29 14. M. G. Silveira and K. D. Lindor, *Expert opinion on pharmacotherapy*, 2014, 15, 365-372.
30 15. S. Mudaliar, R. R. Henry, A. J. Sanyal, L. Morrow, H. U. Marschall, M. Kipnes, L. Adorini, C. I.
31 Sciacca, P. Clopton, E. Castelloe, P. Dillon, M. Pruzanski and D. Shapiro, *Gastroenterology*, 2013,
32 145, 574-582 e571.
33 16. P. R. Maloney, D. J. Parks, C. D. Haffner, A. M. Fivush, G. Chandra, K. D. Plunket, K. L. Creech, L. B.
34 Moore, J. G. Wilson, M. C. Lewis, S. A. Jones and T. M. Willson, *Journal of medicinal chemistry*,
35 2000, 43, 2971-2974.
36 17. M. Downes, M. A. Verdecia, A. J. Roecker, R. Hughes, J. B. Hogenesch, H. R. Kast-Woelbern, M. E.
37 Bowman, J. L. Ferrer, A. M. Anisfeld, P. A. Edwards, J. M. Rosenfeld, J. G. Alvarez, J. P. Noel, K. C.
38 Nicolaou and R. M. Evans, *Molecular cell*, 2003, 11, 1079-1092.
39 18. B. Flatt, R. Martin, T. L. Wang, P. Mahaney, B. Murphy, X. H. Gu, P. Foster, J. Li, P. Pircher, M.
40 Petrowski, I. Schulman, S. Westin, J. Wrobel, G. Yan, E. Bischoff, C. Daige and R. Mohan, *Journal*
41 *of medicinal chemistry*, 2009, 52, 904-907.
42 19. D. Merk, C. Lamers, K. Ahmad, R. Carrasco Gomez, G. Schneider, D. Steinhilber and M.

- 1 Schubert-Zsilavecz, *Journal of medicinal chemistry*, 2014, 57, 8035-8055.
- 2 20. D. Merk, D. Steinhilber and M. Schubert-Zsilavecz, *Expert Opin Drug Discov*, 2014, 9, 27-37.
- 3 21. J. Fu, P. Si, M. Zheng, L. Chen, X. Shen, Y. Tang and W. Li, *Bioorganic & medicinal chemistry letters*,
- 4 2012, 22, 6848-6853.
- 5 22. D. L. Howarth, S. H. W. Law, J. M. Law, J. A. Mondon, S. W. Kullman and D. E. Hinton, *Toxicology*
- 6 *and applied pharmacology*, 2010, 243, 111-121.
- 7 23. D. Schuster, P. Markt, U. Grienke, J. Mihaly-Bison, M. Binder, S. M. Noha, J. M. Rollinger, H.
- 8 Stuppner, V. N. Bochkov and G. Wolber, *Bioorganic & medicinal chemistry*, 2011, 19, 7168-7180.
- 9 24. L. Z. Mi, S. Devarakonda, J. M. Harp, Q. Han, R. Pellicciari, T. M. Willson, S. Khorasanizadeh and F.
- 10 Rastinejad, *Molecular cell*, 2003, 11, 1093-1100.
- 11 25. A. Akwabi-Ameyaw, J. Y. Bass, R. D. Caldwell, J. A. Caravella, L. Chen, K. L. Creech, D. N. Deaton, S.
- 12 A. Jones, I. Kaldor, Y. Liu, K. P. Madauss, H. B. Marr, R. B. McFadyen, A. B. Miller, F. N. Iii, D. J.
- 13 Parks, P. K. Spearing, D. Todd, S. P. Williams and G. B. Wisely, *Bioorganic & medicinal chemistry*
- 14 *letters*, 2008, 18, 4339-4343.
- 15 26. J. Achenbach, M. Gabler, R. Steri, M. Schubert-Zsilavecz and E. Proschak, *Medchemcomm*, 2013,
- 16 4, 920-924.
- 17 27. D. Schuster, P. Markt, U. Grienke, J. Mihaly-Bison, M. Binder, S. M. Noha, J. M. Rollinger, H.
- 18 Stuppner, V. N. Bochkov and G. Wolber, *Bioorganic & medicinal chemistry*, 2011, 19, 7168-7180.
- 19 28. M. Marinozzi, A. Carotti, E. Sansone, A. Macchiarulo, E. Rosatelli, R. Sardella, B. Natalini, G. Rizzo,
- 20 L. Adorini, D. Passeri, F. De Franco, M. Pruzanski and R. Pellicciari, *Bioorganic & medicinal*
- 21 *chemistry*, 2012, 20, 3429-3445.
- 22 29. L. Wang, X. Le, L. Li, Y. Ju, Z. Lin, Q. Gu and J. Xu, *Journal of chemical information and modeling*,
- 23 2014, 54, 3186-3197.
- 24 30. S. Derksen, O. Rau, P. Schneider, M. Schubert-Zsilavecz and G. Schneider, *Chemmedchem*, 2006,
- 25 1, 1346-+.
- 26 31. Y. Li, L. Wang, Z. Liu, C. Li, J. Xu and Q. Gu, *Molecular bioSystems*, 2015, DOI:
- 27 10.1039/c4mb00718b.
- 28 32. N. Ai, M. D. Krasowski, W. J. Welsh and S. Ekins, *Drug Discov Today*, 2009, 14, 486-494.
- 29 33. A. Gaulton, L. J. Bellis, A. P. Bento, J. Chambers, M. Davies, A. Hersey, Y. Light, S. McGlinchey, D.
- 30 Michalovich, B. Al-Lazikani and J. P. Overington, *Nucleic acids research*, 2012, 40, D1100-1107.
- 31 34. M. M. Mysinger, M. Carchia, J. J. Irwin and B. K. Shoichet, *Journal of medicinal chemistry*, 2012,
- 32 55, 6582-6594.
- 33 35. T. A. Halgren, *Journal of computational chemistry*, 1996, 17, 490-519.
- 34 36. L. Chen, Y. Y. Li, Q. Zhao, H. Peng and T. J. Hou, *Mol Pharmaceut*, 2011, 8, 889-900.
- 35 37. S. Tian, Y. Y. Li, J. M. Wang, J. Zhang and T. J. Hou, *Mol Pharmaceut*, 2011, 8, 841-851.
- 36 38. P. Watson, *Journal of chemical information and modeling*, 2008, 48, 166-178.
- 37 39. G. De'ath and K. E. Fabricius, *Ecology*, 2000, 81, 3178-3192.
- 38 40. P. Baldi, S. Brunak, Y. Chauvin, C. A. F. Andersen and H. Nielsen, *Bioinformatics*, 2000, 16,
- 39 412-424.
- 40 41. A. A. Cronican, N. F. Fitz, T. Pham, A. Fogg, B. Kifer, R. Koldamova and I. Lefterov, *Biochem*
- 41 *Pharmacol*, 2010, 79, 1310-1316.
- 42 42. J. Fukuchi, C. Song, Q. Dai, R. A. Hiipakka and S. Liao, *J Steroid Biochem*, 2005, 94, 311-318.
- 43 43. R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E. H.
- 44 Knoll, M. Shelley, J. K. Perry, D. E. Shaw, P. Francis and P. S. Shenkin, *Journal of medicinal*

- 1 *chemistry*, 2004, 47, 1739-1749.
- 2 44. R. A. Friesner, R. B. Murphy, M. P. Repasky, L. L. Frye, J. R. Greenwood, T. A. Halgren, P. C.
- 3 Sanschagrin and D. T. Mainz, *Journal of medicinal chemistry*, 2006, 49, 6177-6196.
- 4 45. T. A. Halgren, R. B. Murphy, R. A. Friesner, H. S. Beard, L. L. Frye, W. T. Pollard and J. L. Banks,
- 5 *Journal of medicinal chemistry*, 2004, 47, 1750-1759.
- 6 46. J. Xu, *Journal of medicinal chemistry*, 2002, 45, 5311-5320.
- 7 47. Q. O. Gu, J. Xu and L. Q. Gu, *Molecules*, 2010, 15, 5031-5044.

8

9