

RSC Advances



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. This *Accepted Manuscript* will be replaced by the edited, formatted and paginated article as soon as this 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 Discovering New DNA Gyrase Inhibitors Using Machine Learning Approaches

2
3 Long Li¹, Xiu Le¹, Ling Wang^{2,3}, Qiong Gu¹, Huihao Zhou¹ and Jun Xu*¹

4 ¹*Research Center for Drug Discovery, School of Pharmaceutical Sciences, Sun*
5 *Yat-Sen University, Guangzhou 510006, China*

6 ²*Pre-Incubator for Innovative Drugs & Medicine, School of Bioscience and*
7 *Bioengineering, South China University of Technology, Guangzhou 510006, China*

8 ³*Guangdong Provincial Key Laboratory of Fermentation and Enzyme Engineering,*
9 *School of Bioscience and Bioengineering, South China University of Technology,*
10 *Guangzhou 510006, China*

11
12 *Correspondent Author: Jun Xu, junxu@biochemomes.com.

13 14 Abstract

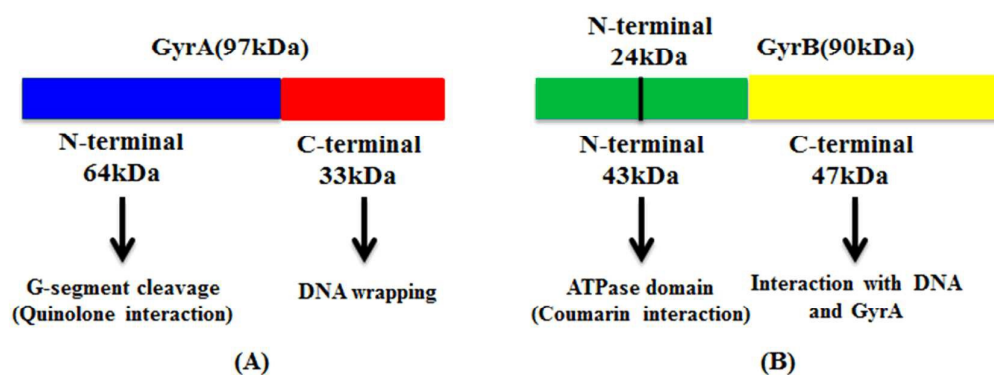
15 The bacterial DNA gyrase is not expressed in eukaryotes. It is a promising target for
16 broad-spectrum antibiotics. This paper reports new DNA gyrase inhibitors as
17 broad-spectrum antibacterial agents discovered by means of target-based *in silico* and
18 *in vitro* models. Two machine learning methods (naïve Bayesian and recursive
19 partitioning) were employed to build the *in silico* models based on physicochemical
20 descriptors and structural fingerprints. For both training and testing sets, the overall
21 predictive accuracies of the best *in silico* models were greater than 80%. The best 11
22 models were used to virtually screen a molecular database to identify DNA gyrase
23 inhibitors. The *in vitro* models were used to verify the virtual hits activities against
24 *Escherichia coli*, methicillin-resistant *staphylococcus aureus* and other bacteria, and
25 DNA gyrase. The MIC values of the confirmed DNA gyrase inhibitors range 1~32
26 µg/mL and, the relatively inhibition rates of the inhibitors range 42%~75% at 1 µM.
27 Cell-based cytotoxicity assays demonstrated that the confirmed DNA gyrase
28 inhibitors were not toxic. *In silico* studies indicated that the new DNA gyrase
29 inhibitors have the similar binding modes of the reported inhibitors.

30
31 **Keywords:** Antibiotic, DNA gyrase inhibitor, machine learning, virtual screening.

32 33 1. Introduction

34 Growing multidrug-resistant bacteria and declining available antibacterial agents
 35 are threatening public health.¹⁻³ New agents against drug-resistant bacteria are
 36 demanded.^{4,5} DNA gyrase is a promising antibacterial drug target because it is
 37 required for all bacteria, and absent in eukaryotes. DNA gyrase is a type II
 38 topoisomerase that mediates negative supercoiling to the relaxed closed circular
 39 DNA^{6,7} and well-studied as an anti-bacterial target.^{8,9} However, only one compound
 40 (ETX0914) is in clinical trials. Others DNA gyrase inhibitors were failure due to side
 41 effects or poor bioavailability.

42 DNA gyrase is a hetero tetramer made up of two GyrA and two GyrB subunits.⁸
 43 GyrA consists of two stable fragments GyrA33 and GyrA64.¹⁰ GyrA64 catalyzes
 44 supercoiling reaction while the GyrB exists and, associates with DNA cleavage and
 45 ligation under the condition of holoenzyme. GyrA33 directly effects on DNA and
 46 forms DNA-enzyme complex that catalyzes supercoiling reaction together with
 47 GyrA64 and GyrB.^{11,12} In the same way, GyrB consists of fragments GyrB43 and
 48 GyrB47. The N-terminal of GyrB43 hydrolyses ATP. As a part of GyrB43, GyrB24
 49 binds DNA gyrase inhibitors such as novobiocin,¹³ aminocoumarin^{13,14} and
 50 cyclothialidine;^{15,16} The C-terminal GyrB47 catalyzes supercoiling DNA to relaxed
 51 DNA in the presence of GyrA (Figure 1).^{17,18}



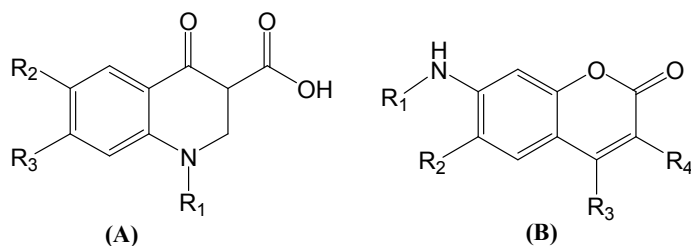
52

53

Figure 1. The hetero components of DNA gyrase.

54 DNA gyrase inhibitors (such as, GSK299423, NXL101 and gyramide) contain
 55 either have quinolone scaffold (A)^{19,20} or aminocoumarins scaffold (B)²¹ (Figure 2).
 56 Quinolones may inhibit supercoiling activity or, induce DNA double-strand breaking.
 57 As examples of scaffold A, fluoroquinolones (FQ) are bacterial topoisomerase
 58 inhibitors.²² The aminocoumarins (such as aminopyrazinamides, hiazolopyridine
 59 ureas, and pyrrolamides) are the competitive inhibitors of ATP hydrolysis, and inhibit

60 DNA supercoiling activities.²³⁻²⁵



61

62

Figure 2. Chemical scaffolds of quinolones (A) and aminocoumarins (B)

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

So far, only three anti drug-resistant bacteria agents (daptomycin, linezolid and bedaquiline) were reported since 1960. DNA gyrase (an anti drug-resistant bacteria drug target) has only one compound (ETX0914), which is under Phase II clinical trials.²⁵ It is demanded for new DNA gyrase inhibitors. Known DNA gyrase inhibitors have diverse scaffolds (Figure 3), which mean that the active sites of the target can adopt diversified ligand shapes. The relations between structures and DNA gyrase inhibitory activities cannot be assumed as being linear or other continuous functional. Hence, we employ two machine learning approaches, naïve Bayesian (NB) learning and recursive partitioning (RP) approaches to generate virtual screening models from target-based DNA gyrase inhibitory data.²⁶ To assure the robustness of the models, we evaluated the models by means of 5-fold cross validations. An external testing data set was also used to test the models. Then, the models were used to virtually screen an in-house compound library, which consisted of 488 tangible compounds.^{27, 28} The virtual hits were validated with cell-based and target-based microbial assays, and following with cytotoxicity assays. The binding modes of confirmed DNA gyrase inhibitors were investigated.

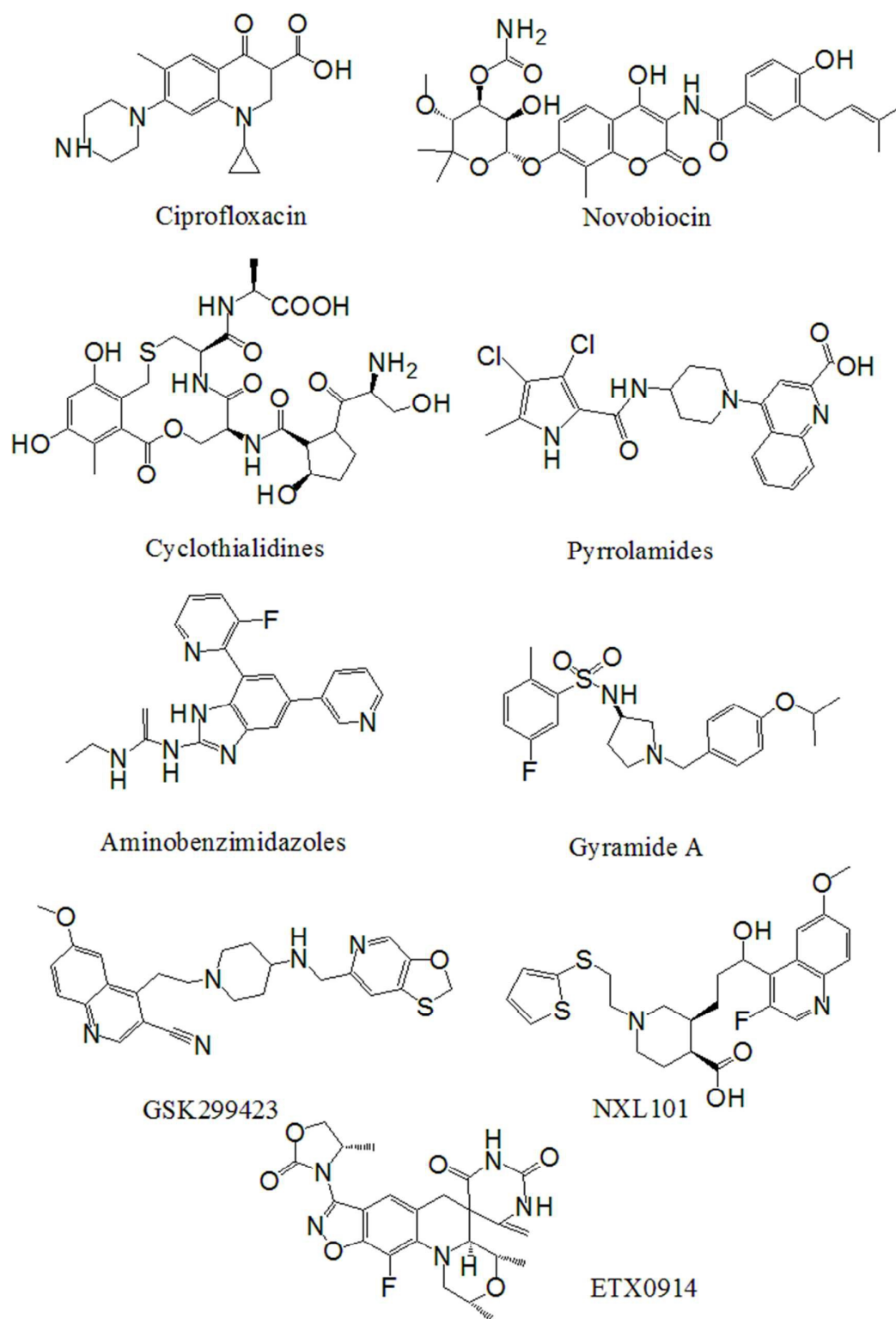


Figure 3.The structures of known DNA gyrase inhibitors.

79

80

81

82 2. Materials and Methods

83 2.1 Data for generating virtual screening models

84 The DNA gyrase inhibitor bioassay data were extracted from the ChEMBL²⁹ and
 85 BindingDB databases by taking target-based *Escherichia coli* strain bioassay data.
 86 Duplicated records or records without IC₅₀ values were removed. This resulted in 137
 87 DNA gyrase inhibitors with IC₅₀ values ranging from 0.9 to 1,000,000 nM. These
 88 compounds were categorized into positives and negatives based upon their IC₅₀ values
 89 (the compounds with IC₅₀ values less than or equal to 5 μM were marked with “1” for
 90 positives. Others were marked with “0” for negatives). The entire data set was
 91 randomly divided into four portions. A training set was made of the three portions
 92 containing 103 compounds. The remaining portion was used as a testing set
 93 containing 34 compounds.³⁰ This process was done with DS (Discovery Studio 3.5,
 94 Accelrys, San Diego, USA). The detailed process can be examined in Supporting
 95 Information.

96 2.2 Molecular descriptor calculation and selection

97 The molecular descriptors of the data set were computed with MOE 2013.08 (CCG,
 98 Montreal, Canada) and DS, resulting in 192 MOE molecular descriptors and 252 DS
 99 molecular descriptors for each compound in the data set.

100 With Pearson correlation analyses, the redundant molecular descriptors (selective
 101 ratio > 0.9) were removed, the molecular descriptors (selective ratio < 0.1), which
 102 were unrelated to the DNA gyrase inhibitory activities, were excluded.^{27,28,31} This
 103 resulted in 36 MOE descriptors and 15 DS descriptors (Table 1).

104

105

Table 1. Selected molecular descriptors

Class	Number	Descriptor
MOE	36	GCUT_PEOE_0, GCUT_SLOGP_0, GCUT_SLOGP_1, GCUT_SLOGP_2, GCUT_SLOGP_3, GCUT_SMR_0, GCUT_SMR_1, GCUT_SMR_2, PEOE_VSA+0, PEOE_VSA+1, PEOE_VSA+2, PEOE_VSA+3, PEOE_VSA+4, PEOE_VSA+5, PEOE_VSA+6, PEOE_VSA-0, PEOE_VSA-2, PEOE_VSA-3, PEOE_VSA_FPOS, SMR_VSA2, SMR_VSA3, SMR_VSA6, SMR_VSA7, SlogP, SlogP_VSA4, SlogP_VSA8, a_ICM, ast_violation_ext, b_maxl1en, b_rotR, mutagenic, petitjeanSC, reactive, rsynth, vsa_acc, vsa_hyd

DS	15	E_DIST_equ, SIC, CHI_V_3_P, JX, HBA_Count, HBD_Count, NPlusO_Count, Num_Hydrogens, Num_RingBonds, Num_AromaticBonds, Num_RingAssemblies, Num_Rings6, Num_AliphaticDoubleBonds, Num_TerminalRotomers, Num_TrueStereoAtoms
----	----	--

106

107 **2.3 Structural fingerprints calculation**

108 Structural fingerprints were calculated using DS software. The fingerprints consist
109 of Daylight-style path-based fingerprints and SciTegic extended-connectivity
110 fingerprints.

111 **2.4 Machine learning approaches**

112 Two machine learning methods, NB and RP, were applied through DS software.

113 **2.4.1 NB method**

114 NB method is a supervised learning approach, and directly calculates the overall
115 distribution based on the prior distribution of parameters and the posterior distribution
116 of parameters obtained from the sample data. The method is based on the Bayes'
117 theorem and the maximum posteriori hypothesis,³² requires the training objects are
118 marked with positives or negatives.³³

119 **2.4.2 RP method**

120 RP (or decision tree) is a statistical method for multivariable analysis and, based on
121 hierarchical rules. It creates a decision tree to describe the relationship between an
122 active and a set of properties/descriptors of objects.^{34, 35}

123 **2.5 Decoys generation**

124 The decoy data were generated from DUD-E³⁶ (<http://dude.docking.org/>) through
125 the Pipeline Pilot 7.5 module of DiscoveryStudio . 10 diverse compounds were used
126 as reference compounds, which were randomly selected from the positives in the
127 input data set. The decoys were selected from DUD-E based upon the dissimilarity to
128 the reference compounds. 80 decoys, which were regarded as negatives, were selected
129 for external tests.

130 **2.6 Method for model performance evaluation**

131 A 5-fold cross validation was used to evaluate the performances of NB and RP
132 models. True positives (TP), true negatives (TN), false positives (FP), false negatives
133 (FN), sensitivity (SE), specificity (SP), overall predictive accuracy (Q), the Matthews

134 correlation coefficient (C) and the receiver operating characteristic (ROC) curve were
135 defined as follows to measure the performance:³⁷

$$136 \quad SE = \frac{TP}{TP+FN}$$

$$137 \quad SP = \frac{TN}{TN+FP}$$

$$138 \quad Q = \frac{TP+TN}{TP+FN+TN+FP}$$

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

139 **2.7 Compound library for virtual screening campaigns**

140 The in-house tangible compound library, which contains 488 natural products or
141 chemically modified natural products, were virtually screened with the best machine
142 learning models.

143 **2.8 *In vitro* antimicrobial assay**

144 **Minimum inhibitory concentration testing.** The test was performed to determine
145 the minimum concentration of the indicated agent necessary to inhibit visible growth
146 of bacteria. In this study, our compounds were tested against bacteria including
147 MRSA ST239, MRSA ST5, MRSA 252, *Staphylococcus aureus*, *Fecal bacteria*,
148 *Staphylococcus epidermidis*, *Pneumonia*, ATCC 25922 and *Shigella flexneri*.
149 Ampicillin and vancomycin sodium were used as positive control agents. The MIC
150 values were determined using Mueller-Hinton broth method based on national
151 committee for clinical laboratory standard^{38,39}. Each compound was tested for 11
152 concentrations (256, 128, 64, 32, 16, 8, 4, 2, 1, 0.5, 0.25 µg/mL). 90 µL bacterial
153 culture medium was added into the first column of wells of flat bottomed 96-well
154 tissue culture plates, and other wells were added with 50 µL same medium, and then
155 10 µL solution of compound was added into the first column of wells. Then, 50 µL
156 mixture extracted from the first column wells were transferred to the second column
157 of wells, and repeated this operation column by column till the second last column of
158 wells. After this step, the 50 µL bacterial culture solutions in last column of wells
159 were discarded. Finally, 50 µL bacterial solution was diluted by culture medium, and
160 added into all wells in the 96-well plate. The last row wells were for positive controls,
161 and the last column wells were for negative controls. The plates were incubated at
162 37°C overnight in electro-heating standing-temperature cultivator before the
163 measurement of the absorbance value. We used a multifunction microplate reader to
164 measure the optical density values at 600 nm. Each antimicrobial assay was replicated

165 four times.

166 **2.9 DNA gyrase expression and purification**

167 The recombinant protein was expressed with plasmids pET-15-GyrA and
168 pET-15b-GyrB in *E.coli*, and purified using Ni-NTA column. After the SDS-PAGE
169 verification, we mixed GyrA and Gry B at 1:1 molar ratio, and incubated on ice for 30
170 min before DNA supercoiling assay.³⁹

171 **2.10 DNA gyrase-mediated pHOT-1 supercoiling assay**

172 The DNA supercoiling assay was conducted to test the inhibitory activity on the
173 enzyme reaction. Firstly, 4 μ L 5 \times DNA gyrase assay buffer, 0.1U relaxed pHOT-1
174 DNA and 12.9 μ L ddH₂O were mixed.³⁸ Then, 17 μ L mixture mentioned above, 2 μ L
175 compounds and 1 μ L reconstituted DNA gyrase were mixed, and incubated at 37 $^{\circ}$ C.
176 After 1 h, 4 μ L 5 \times stop buffer was added to stop the reaction. Novobiocin was used as
177 positive control, and 1% DMSO was employed as blank control. To separate the DNA
178 products, electrophoresis on a 1% agarose gel run used. The gel was stained for 20
179 min in ethidium bromide, decolorized for 15 min in water and visualized with UV light.
180 The optical density of the bands for supercoiling and relaxed DNA was quantified
181 using the Quantity One software. The inhibition rates were used to calculate the IC₅₀
182 values with GraphPad Prism 5. The IC₅₀ values were measured with 7 concentration
183 points, and repeated for three times.

184 **2.11 Cytotoxicity Assay**

185 HEK-293, a human embryonic kidney normal cell line, was used to evaluate the
186 cytotoxicity of the compounds. HEK-293 cells were inoculated in 96-well plates with
187 DMEM medium containing 10% fetal bovine serum at 37 $^{\circ}$ C in 5% CO₂ incubator.
188 Then, the cells were intervened with different compounds at 20 μ M for 24h after cells
189 were adherent and each compound was added into three parallel double wells. Blank
190 control group and empty wells were prepared. Then 20 μ L 2.5 mg/mL MTT was
191 added to each well and incubated for 4h, and 100 μ L DMSO was added every well
192 lastly. Absorption values were measured at 492 nm after 20 minutes' oscillation. The
193 inhibition rate of each compound against 293T cell lines was calculated with the
194 following formula: Inhibition of cell (%) = 1- (A_{experimental group} - A_{blank}) / (A_{control group} -
195 A_{blank}) \times 100%.⁴⁰

196 **2.12 Molecular docking**

197 The intact DNA gyrase (PDB code: 3G7E)²⁵ was used as the template to explore

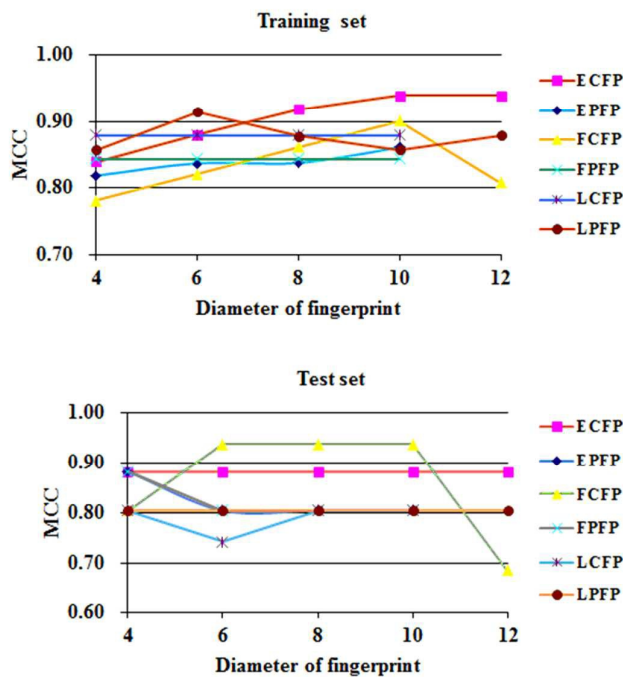
198 the binding modes of the confirmed DNA gyrase inhibitors. The structure data was
 199 processed using a protocol from Schrödinger software 2013.01(Schrödinger Inc., New
 200 York, USA). The active compounds were prepared by Ligprep module in the
 201 Schrödinger software. The extra precision Glide 5.9⁴¹⁻⁴³ of Schrödinger software was
 202 used to dock the active compound structures into the binding pocket of the DNA
 203 gyrase. The active compounds were also superimposed with the native ligand using
 204 WEGA algorithm⁴⁴ to ensure the correct docking pose.

205

206 3. Results

207 3.1 Classifiers derived from molecular descriptors or structural fingerprints

208 Figure 4 indicates that the size (the diameter of a fingerprint) of a structural
 209 fingerprint or the type (ECFP, etc.) of a structural fingerprint can change the model
 210 performance (MCC value). But, there is no general trend. SciTeGic
 211 extended-connectivity fingerprints resulted in better performance in general.



212

213 **Figure 4.** The relations among MCCs and fingerprint sizes or types.

214 Table 2 lists the performance parameters of top-10 machine learning models
 215 running on training set and testing set. The top-10 models were all generated from NB
 216 method with overall predictive accuracies greater than 94.1% for both training set and
 217 test set. For the testing set, the models using FCFP_6, FCFP_8 and FCFP_10
 218 fingerprints achieved better performances with the sensitivity of 95.4%, the specificity

219 of 100.0%, overall prediction accuracies of 97.1%, and the AUC value of 0.992.

220

221 **Table 2.** Performances of top-10 models using descriptors* or fingerprints

Models	Training set								
	TP	FN	TN	FP	SE	SP	C	AUC	Q
FCFP_6	56	6	38	3	0.903	0.927	0.821	0.918	0.913
FCFP_8	57	5	39	2	0.919	0.951	0.861	0.914	0.932
FCFP_10	58	4	40	1	0.935	0.976	0.902	0.911	0.951
ECFP_4	57	5	38	3	0.919	0.927	0.840	0.926	0.922
ECFP_6	58	4	39	2	0.935	0.951	0.880	0.923	0.942
ECFP_8	60	2	39	2	0.968	0.951	0.919	0.92	0.961
ECFP_10	60	2	40	1	0.968	0.976	0.940	0.919	0.971
ECFP_12	60	2	40	1	0.968	0.976	0.940	0.919	0.971
EPFP_4	57	5	37	4	0.919	0.902	0.819	0.893	0.913
FPPF_4	56	6	39	2	0.903	0.951	0.843	0.889	0.922
Models	Testing set								
	TP	FN	TN	FP	SE	SP	C	AUC	Q
FCFP_6	21	1	12	0	0.954	1.000	0.939	0.992	0.971
FCFP_8	21	1	12	0	0.954	1.000	0.939	0.992	0.971
FCFP_10	21	1	12	0	0.954	1.000	0.939	0.992	0.971
ECFP_4	20	2	12	0	0.909	1.000	0.883	0.992	0.941
ECFP_6	20	2	12	0	0.909	1.000	0.883	0.989	0.941
ECFP_8	20	2	12	0	0.909	1.000	0.883	0.989	0.941
ECFP_10	20	2	12	0	0.909	1.000	0.883	0.989	0.941
ECFP_12	20	2	12	0	0.909	1.000	0.883	0.989	0.941
EPFP_4	20	2	12	0	0.909	1.000	0.883	0.989	0.941
FPPF_4	20	2	12	0	0.909	1.000	0.883	0.973	0.941

222 * The models using descriptors are not listed in this table because they are not

223 ranked in the top-10 models.

224

225 3.2 Performance of models using combined molecular descriptors and structural 226 fingerprints

227 Descriptors (physiochemical properties) and fingerprints (substructures) represent

228 different attributions of compound structures. We thought the models using both
 229 might result in better performances. 54 NB models and 324 RP models generated
 230 from the combinations of descriptors and fingerprints (detailed modeling data can be
 231 found in Supporting Information Table S4/S5 and Figure S4/S5). The top-10 models
 232 are listed in Table 3.

233

234 **Table 3.** Top-10 models using combined descriptors and fingerprints

Models	Training set								
	TP	FN	TN	FP	SE	SP	C	AUC	Q
MOE ^a +EFCP_4-4*	28	13	56	6	0.683	0.903	0.610	0.793	0.816
MOE+FPFP_4-4*	28	13	56	6	0.683	0.903	0.610	0.793	0.816
MOE + EPFP_8	58	4	36	5	0.935	0.878	0.817	0.915	0.913
FCFP_6	56	6	38	3	0.903	0.927	0.821	0.918	0.913
FCFP_8	57	5	39	2	0.919	0.951	0.861	0.914	0.932
FCFP_10	58	4	40	1	0.935	0.976	0.902	0.911	0.951
DS ^b +EPFP_4-5*	27	14	61	1	0.659	0.984	0.707	0.8226	0.854
DS + EPFP_4	57	5	38	3	0.919	0.927	0.840	0.894	0.922
DS + FPFP_4	54	8	39	2	0.871	0.951	0.808	0.892	0.903
MOE + EPFP_4	59	3	37	4	0.952	0.902	0.858	0.894	0.932
Models	Test set								
	TP	FN	TN	FP	SE	SP	C	AUC	Q
MOE+EFCP_4-4*	22	0	1	0	2.000	1.000	1.000	0.800	0.909
MOE+FPFP_4-4*	22	0	12	0	1.000	1.000	1.000	1.000	1.000
MOE + EPFP_8	21	1	12	0	0.955	1.000	0.939	0.992	0.971
FCFP_6	21	1	12	0	0.955	1.000	0.939	0.992	0.971
FCFP_8	21	1	12	0	0.955	1.000	0.939	0.992	0.971
FCFP_10	21	1	12	0	0.955	1.000	0.939	0.992	0.971
DS+EPFP_4-5*	22	0	11	1	1.000	0.917	0.936	0.958	0.971
DS + EPFP_4	20	2	12	0	0.909	1.000	0.883	0.970	0.941
DS + FPFP_4	20	2	12	0	0.909	1.000	0.883	0.973	0.941
MOE + EPFP_4	20	2	12	0	0.909	1.000	0.883	0.985	0.941

235 * RP models.

236 ^a MOE: descriptors calculated from MOE software.

237 ^b DS: descriptors calculated from DS software.

238 Comparing tables 2 and 3, we find that NB models using combined molecular
239 descriptors and structural fingerprints are actually worse than the NB models using
240 molecular descriptors or structural fingerprints. However, the RP models using
241 combined molecular descriptors and structural fingerprints can result better
242 performance than the ones of using non-combined descriptors or fingerprints.

243

244 3.3 Determining and external testing final models

245 Combining tables 2 and 3, we get top-11 models after removed duplicated models.
246 The 11 final models were tested with the external testing data set. Table 4 lists the
247 results.

248 **Table 4.** The external testing results for the top-11 final models

Models	Test set								
	TP	FN	TN	FP	SE	SP	C	AUC	Q
FCFP_10	9	2	95	14	0.818	0.872	0.506	0.927	0.867
DS + FFPF_4	91	18	11	0	0.835	1.000	0.563	0.583	0.850
DS+EPFP_4-5*	91	18	11	0	0.835	1.000	0.563	0.583	0.850
FCFP_8	13	11	93	9	0.542	0.912	0.469	0.505	0.841
ECFP_6	13	11	93	9	0.542	0.912	0.469	0.505	0.841
MOE + FFPF_4	13	11	93	9	0.54	0.912	0.469	0.505	0.841
MOE+ECFP_4-4*	13	11	93	9	0.54	0.912	0.469	0.505	0.841
FFPF_4	2	17	90	4	0.105	0.957	0.105	0.764	0.814
MOE + EPFP_8	2	17	89	5	0.105	0.947	0.081	0.501	0.805
FCFP_6	12	6	66	29	0.667	0.695	0.275	0.505	0.690
MOE + EPFP_4	16	3	52	44	0.842	0.553	0.296	0.666	0.602

249

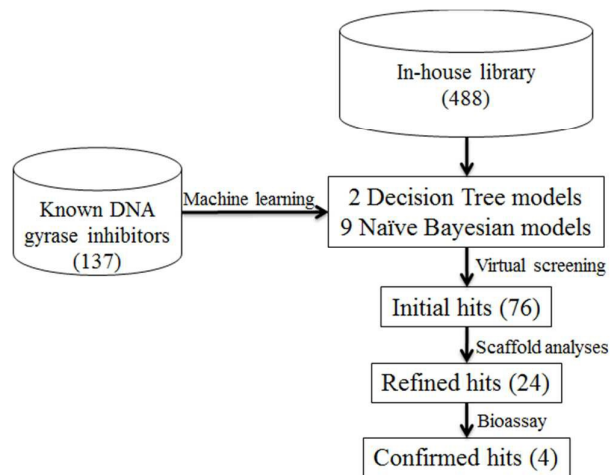
250 The overall prediction accuracies of the final models are greater than 80% (except
251 models FCFP_6, MOE + EPFP_4). The top model (FCFP_10) was generated from
252 NB method (see the first row in Table 4).

253

254 3.4 Virtual screening DNA gyrase inhibitors with the final models

255 Our in-house library, which has 488 tangible compounds, was virtually screened
256 with the top-11 predictive models (Table 4), which consist of nine NB models, and

257 two RP models. The NB models resulted in 67 hits, and the RP models resulted in 19
 258 hits. By combining the two hit sets, we got 76 initial hits without duplicates. The
 259 initial hits were further refined by scaffold analyzing processes, which removed
 260 known antibacterial scaffolds (such as, flavone derivatives), and resulted in 24 refined
 261 hits.. These refined hits were tested with cell-based microbiological assays. The
 262 flow-chart of discovering new DNA gyrase inhibitors using machine learning
 263 approaches is depicted in Figure 5.



264

265 **Figure 5.** The flow-chart of discovering new DNA gyrase inhibitors using machine
 266 learning approaches

267 3.5 Cell-based microbiological assay results

268 Both G+ and G- strains were tested in the cell-based microbiological assays.
 269 Ampicillin sodium and vancomycin sodium were used as positive controls. 4
 270 compounds actively inhibited *E. coli* and MRSA strains (XGS00156, XGS00157,
 271 XGS00158 and XGS00159). As shown in Table 5, the 4 active compounds have MIC
 272 values < 10 μ M. The advantages of the 4 compounds are that these compounds
 273 exhibited broader spectrum of antibacterial activities than ampicillin or vancomycin.
 274 The activities of compound XGS00159 are comparable with the ones of ampicillin or
 275 vancomycin. All active compounds share the same scaffold. Their initial SAR is
 276 established (Figure 6).

277

Table 5. Cell-based microbiological study results (μ M)

	G ⁺ ^a							G ⁻ ^a	
ID	MRSA ST239	MRSA ST5	MRSA 252	ATCC 29213	ATCC 29212	ATCC 12228	<i>Pneum</i> <i>onia</i>	ATCC 25922	CMCC 51572

XGS00156	16.42	15.42	15.42	15.42	-	15.42	7.71	3.85	61.67
XGS00157	32.01	65.67	131.34	16.42	-	32.84	131.34	8.21	-
XGS00158	10.17	20.35	10.17	10.17	-	10.50	5.09	5.09	81.4
XGS00159	5.04	5.04	5.04	5.04	-	5.04	2.52	5.04	-
Amp ^b	-	-	-	-	-	2	2	2	2
WG ^c	2	2	2	2	2	2	2	-	-

278 ^a MRSA: methicillin resistant *staphylococcus aureus*,

279 ATCC 29213: *Staphylococcus aureus*, ATCC 29212: *Fecal bacteria*,

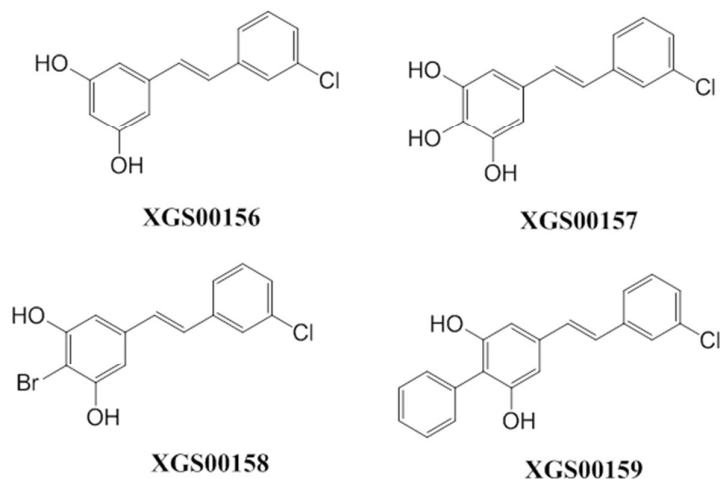
280 ATCC 12228: *Staphylococcus epidermidis*, ATCC 25922: *Escherichia coli*,

281 CMCC 51572: *Shigella flexneri*;

282 ^bAmp: ampicillin sodium;

283 ^cWG: vancomycin sodium, positive control.

284



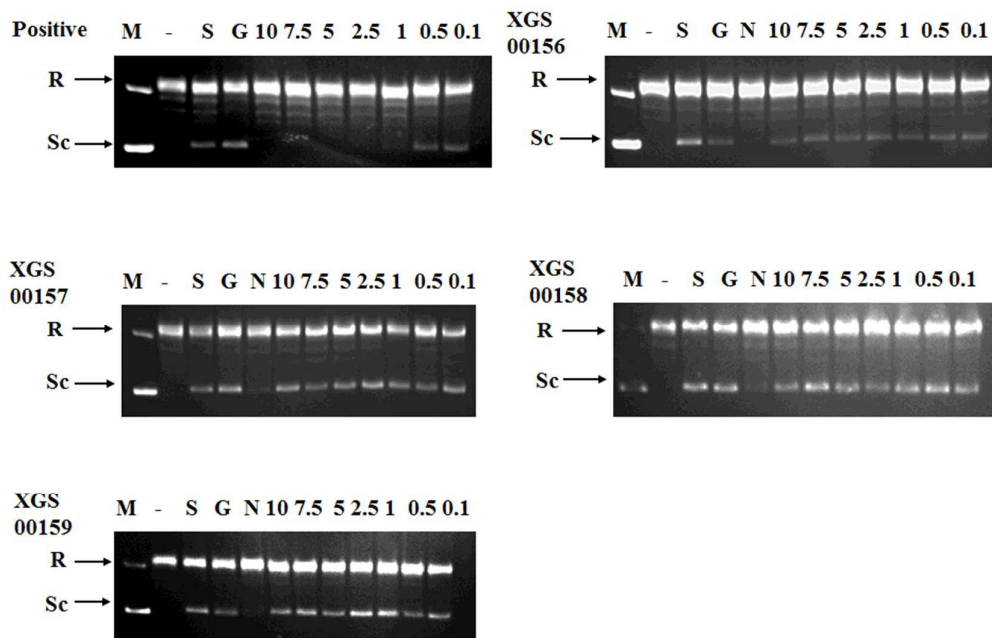
285

286

Figure 6. Initial SAR of 4 confirmed hits.

287 3.6 DNA supercoiling assay results

288 The 4 compounds were tested with DNA supercoiling assays. Novobiocin was used
 289 as a positive control. The results were depicted in Figure 7, and indicated that the 4
 290 compounds dose-dependently inhibited DNA supercoiling. Thus, the 4 compounds
 291 have been proved that they are DNA gyrase inhibitors.



292

293 **Figure 7.** DNA supercoiling assay results for the 4 compounds (XGS00156,

294 XGS00157, XGS00158 and XGS00159). R: relaxed DNA; Sc: supercoiled DNA; M:

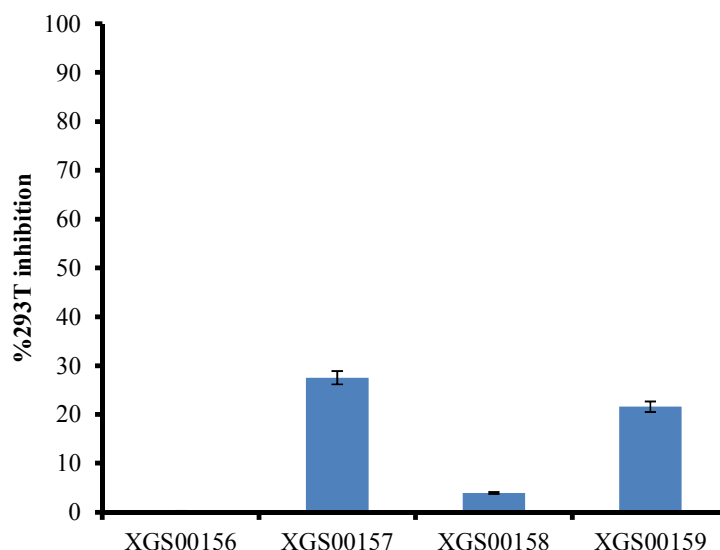
295 Marker; -: negative control; S: activity of enzyme; G: 1% DMSO; N: Positive control,

296 novobiocin.

297

298 3.7 Cytotoxicity assay results

299 Figure 8 depicts the cytotoxicity assay results for the 4 active compounds. At 20
300 μM , two active compounds inhibited $< 5\%$ of 293T cell lines, other two compounds
301 inhibited $< 30\%$ of 293T cell lines. Thus, the 4 active compounds are considered as
302 promising drug leads, and worth further lead optimization processing.⁴⁰



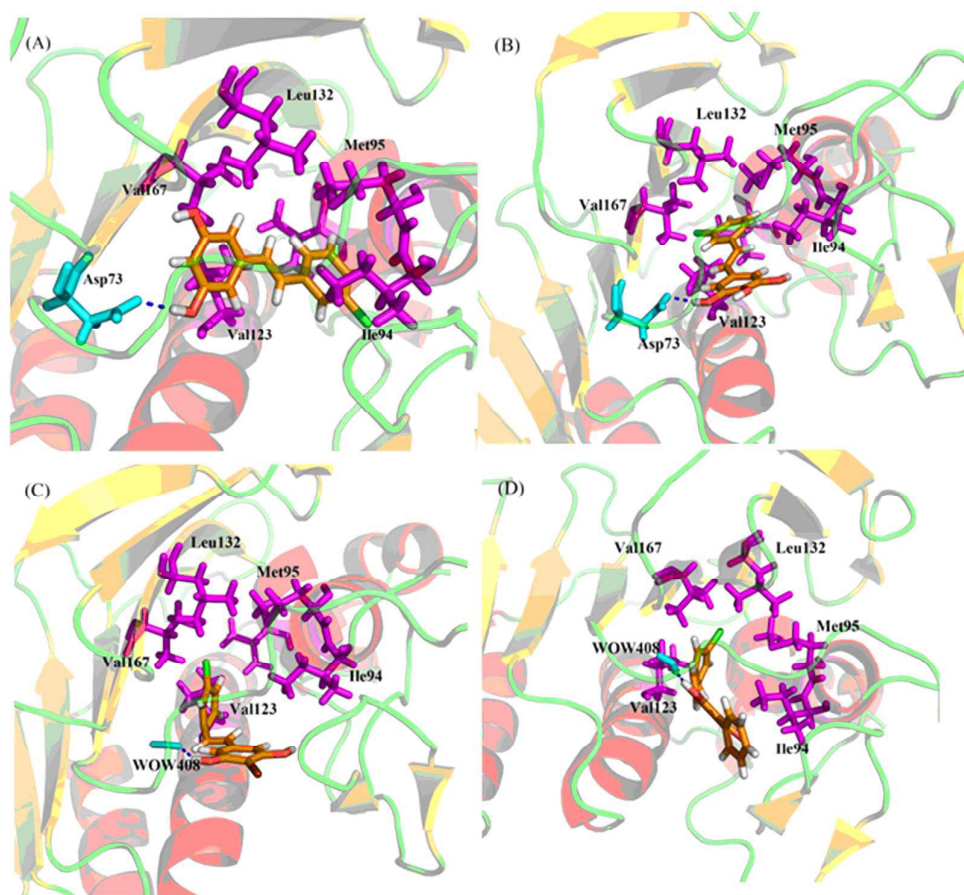
303

304

Figure 8. Cytotoxicity assay results

305 3.8 Molecular docking study results

306 The 4 compounds were docked to the crystal structure (3G7E), in which the native
307 ligand was removed. The docking processes were executed with both the extra
308 precision Glide and WEGA algorithm. The docking poses of the compounds were
309 consistent. This demonstrated that the docking processes were reliable. Figure 7
310 depicts the binding modes of the 4 active compounds. All 4 compounds have the
311 similar interactions with the known key residues, such as WOW 408 or Asp 73, which
312 is a hydrogen bond donor. The hydrophobic groups of the compounds interact with
313 the receptor hydrophobic pocket (Val 43, Met 95, Ile 94, Val 123, Leu 132 and Val
314 167).^{45,46} Thus, these active compounds binding modes support the observations of
315 the *in vitro* results.



316

317 **Figure 9.** Molecular docking study results. (A): Binding mode of XGS00156; (B):
318 Binding mode of XGS00157; (C): Binding mode of XGS00158; (D): Binding mode
319 of XGS00159. The molecules in orange are the active compounds; the residues in
320 blue donor hydrogen bonds; the residues in red provide hydrophobic interactions. The
321 deep blue dashed lines represent hydrogen bonds.

322

323 4. Conclusions

324 DNA gyrase is a promising drug target, but, there are not many DNA gyrase
325 inhibitors under clinic trials. Existing DNA gyrase inhibitors are structurally diverse,
326 it would be difficult to discover novel DNA gyrase inhibitors through structure-based
327 molecular design, or individual ligand-based modeling technology, or traditional
328 QSAR techniques. This work demonstrates that we can discover a novel scaffold of
329 DNA gyrase inhibitors by combining multiple machine learning methods and
330 target-based approaches. There are many ways to build virtual screening models due
331 to many types of structural descriptors or fingerprints. Since we did not discover
332 specific descriptors or fingerprints were particular superior to the others for the virtual

333 screening campaign. To do our best to include excellent virtual screening models, we
334 have explored 424 machine learning models derived from the combinations of the
335 descriptors or fingerprints. The confirmed hits were generated from the top-11
336 models.

337 **Acknowledgements**

338 This work was supported by the National Science Foundation of China (81173470),
339 National High Technology Research and Development Program of China (863
340 Program, 2012AA020307), National Supercomputer Center in Guangzhou
341 (2012Y2-00048/2013Y2-00045, 201200000037), the introduction of innovative R&D
342 team program of Guangdong Province (2009010058), Guangdong Provincial Key
343 Laboratory of Construction Foundation (2011A060901014), and the Fundamental
344 Research Funds for the Central Universities (2013HGCH0015).

345

346 **Transparency declarations**

347 The authors declare no competing financial interest.

348 **Notes and references**

349 LL E-mail: luckylilong1012@163.com.

350 XL E-mail: lexiu2012@163.com

351 LW E-mail: lingwang@scut.edu.cn

352 HZ E-mail: zhuihao@mail.sysu.edu.cn

353 QG Email: guqiong@mail.sysu.edu.cn.

354 JX Email: junxu@biochemomes.com.

355 † The experiment design JX, LL, XL, LW. Implementation: LL, HZ, QG. Manuscript
356 revision and submission: LL and JX.

357

358

- 359 1. R. Janupally, B. Medepi, P. Brindha Devi, P. Suryadevara, V. U. Jeankumar, P.
360 Kulkarni, P. Yogeewari and D. Sriram, *Chemical biology & drug design*, 2015, DOI:
361 10.1111/cbdd.12529.
- 362 2. S. B. Singh, D. E. Kaelin, J. Wu, L. Miesel, C. M. Tan, P. T. Meinke, D. Olsen, A.
363 Lagrutta, P. Bradley, J. Lu, S. Patel, K. W. Rickert, R. F. Smith, S. Soisson, C. Wei, H.
364 Fukuda, R. Kishii, M. Takei and Y. Fukuda, *ACS medicinal chemistry letters*, 2014,
365 5, 609-614.
- 366 3. G. S. Basarab, J. I. Manchester, S. Bist, P. A. Boriack-Sjodin, B. Dangel, R.
367 Illingworth, B. A. Sherer, S. Sriram, M. Uria-Nickelsen and A. E. Eakin, *Journal of*
368 *medicinal chemistry*, 2013, 56, 8712-8735.
- 369 4. L. Feng, M. M. Maddox, M. Z. Alam, L. S. Tsutsumi, G. Narula, D. F. Bruhn, X. Wu, S.
370 Sandhaus, R. B. Lee, C. J. Simmons, Y. C. Tse-Dinh, J. G. Hurdle, R. E. Lee and D. Sun,
371 *J Med Chem*, 2014, 57, 8398-8420.
- 372 5. H. Nimesh, S. Sur, D. Sinha, P. Yadav, P. Anand, P. Bajaj, J. S. Viridi and V. Tandon,
373 *Journal of medicinal chemistry*, 2014, 57, 5238-5257.
- 374 6. J. J. Champoux, *Annual review of biochemistry*, 2001, 70, 369-413.
- 375 7. V. U. Jeankumar, R. S. Reshma, R. Janupally, S. Saxena, J. P. Sridevi, B. Medapi, P.
376 Kulkarni, P. Yogeewari and D. Sriram, *Organic & biomolecular chemistry*, 2015,
377 13, 2423-2431.
- 378 8. K. M. Martin Gellert, Mary H.O'Dea, *Biochemistry-U.S.*, 1976, 73, 3872-3876.
- 379 9. M. Nollmann, N. J. Crisona and P. B. Arimondo, *Biochimie*, 2007, 89, 490-499.
- 380 10. J. GR and R. JP., *Biochemistry-U.S.*, 1976, 15, 5105-5110.
- 381 11. A. M. Richard J. Reece, *The Journal of Biological chemistry*, 1991, 266, 3540-3546.
- 382 12. A. M. Richard J. Reece, *Nucleic acids research*, 1991, 19, 1399.
- 383 13. F. Collin, S. Karkare and A. Maxwell, *Applied microbiology and biotechnology*,
384 2011, 92, 479-497.
- 385 14. C. Anderle, M. Stieger, M. Burrell, S. Reinelt, A. Maxwell, M. Page and L. Heide,
386 *Antimicrob Agents Ch*, 2008, 52, 1982-1990.
- 387 15. N. Nakada., H. Shimada., T. Hirata. and Y. Aokl., *Antimicrobial agents and*
388 *chemotherapy*, 1993, 37, 2656-2661.
- 389 16. L. M. Riley, M. Veses-Garcia, J. D. Hillman, M. Handfield, A. J. McCarthy and H. E.
390 Allison, *BMC microbiology*, 2012, 12, 42.
- 391 17. M. M. Toshiro Adachil, Elizabeth A. Robinson, *Nucleic acids research*, 1987, 15.
- 392 18. L. W. Tari, X. Li, M. Trzoss, D. C. Bensen, Z. Chen, T. Lam, J. Zhang, S. J. Lee, G.
393 Hough, D. Phillipson, S. Akers-Rodriguez, M. L. Cunningham, B. P. Kwan, K. J.
394 Nelson, A. Castellano, J. B. Locke, V. Brown-Driver, T. M. Murphy, V. S. Ong, C. M.

- 395 Pillar, D. L. Shinabarger, J. Nix, F. C. Lightstone, S. E. Wong, T. B. Nguyen, K. J. Shaw
396 and J. Finn, *PloS one*, 2013, 8, e84409.
- 397 19. G. A. Jacoby, *L. Clinic., Burlington. and Massachusetts.*, Supplement article, 2015.
- 398 20. J. Ruiz, *The Journal of antimicrobial chemotherapy*, 2003, 51, 1109-1117.
- 399 21. M. Gellert, M. H. O'Dea and T. Itoh, *Biochemistry-Us*, 1976, 73, 4474-4478.
- 400 22. P. S. Hameed, A. Raichurkar, P. Madhavapeddi, S. Menasinakai, S. Sharma, P. Kaur,
401 R. Nandishaiah, V. Panduga, J. Reddy, V. K. Sambandamurthy and D. Sriram, *ACS*
402 *medicinal chemistry letters*, 2014, 5, 820-825.
- 403 23. L. A. M. Silke Alt, Anthony Maxwell, Lutz Heide, *The Journal of antimicrobial*
404 *chemotherapy*, 2011, 66, 2061-2069.
- 405 24. M. Rajendram, K. A. Hurley, M. H. Foss, K. M. Thornton, J. T. Moore, J. T. Shaw and
406 D. B. Weibel, *ACS chemical biology*, 2014, 9, 1312-1319.
- 407 25. Gregory S. Basarab¹, Gunther H. Kern², John McNulty³, John P. Mueller⁴, Kenneth
408 Lawrence⁴, Karthick Vishwanathan², Richard A. Alm², Kevin Barvian⁵, Peter
409 Doig², Vincent Galullo², Humphrey Gardner², Madhusudhan Gowravaram⁶,
410 Michael Huband⁷, *Scientific Reports*, 2015, 5, 11827.
- 411 26. L. Wang, X. Le, L. Li, Y. Ju, Z. Lin, Q. Gu and J. Xu, *Journal of chemical information*
412 *and modeling*, 2014, 54, 3186-3197.
- 413 27. H. Cui, B. Xu, T. Wu, J. Xu, Y. Yuan, Q. Gu, *Journal of Natural Products*, 2014, 77,
414 100-110.
- 415 28. T. Wu, Q. Wang, C. Jiang, H. Cui, Y. Wang, J. Xu, Q. Gu, *Journal of Natural Products*,
416 2015, 78, 500-509.
- 417 29. A. Gaulton, L. J. Bellis, A. P. Bento, J. Chambers, M. Davies, A. Hersey, Y. Light, S.
418 McGlinchey, D. Michalovich, B. Al-Lazikani and J. P. Overington, *Nucleic acids*
419 *research*, 2012, 40, D1100-1107.
- 420 30. Y. Y. L. L. Chen, Q. Zhao, H. Peng and T. J. Hou, *Mol Pharmaceut*, 2011, 8, 889-900.
- 421 31. J. Fang, R. Yang, L. Gao, D. Zhou, S. Yang, A. L. Liu and G. H. Du, *Journal of chemical*
422 *information and modeling*, 2013, 53, 3009-3020.
- 423 32. S. Y. Yang, *Drug Discov Today*, 2010, 15, 444-450.
- 424 33. L. C. D. Li, Y. Li, S. Tian, H. Sun and T. Hou, *Mol Pharmaceut*, 2014, 11, 716-726.
- 425 34. G. De'ath and K. E. Fabricius, *the Ecological Society of America*, 2000, 81,
426 3178-3192.
- 427 35. L. Chen, Y. Li, Q. Zhao, H. Peng and T. Hou, *Mol Pharm*, 2011, 8, 889-900.
- 428 36. M. C. M. M. Mysinger, J. J. Irwin and B. K. Shoichet, *Journal of medicinal chemistry*,
429 2012, 55, 6582-6594.
- 430 37. Pierre Baldi and Y. Chauvin, *Bioinformatics review*, 2000, 16, 412-424.

- 431 38. C. a. L. S. Institue, 2012.
- 432 39. E. Cf. A. S. Eot. E. So. C. Ma. Infectious and D. ESCMID), Clinical Microbiology and
433 Infection, 2000, 6, 509-515.
- 434 40. Ding, QZ. Li, CJ. Wang, L. Li, YL. Xu, J, Med. Med. Chem. Commun.2015,6,
435 1393-1403
- 436 41. J. L. B. R. A. Friesner, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P.
437 Repasky, E. H. Knoll, M. Shelley, J. K. Perry, D. E. Shaw, P. Francis and P. S. Shenkin,
438 Journal of medicinal chemistry, 2004, 47, 1739-1749.
- 439 42. R. B. M. T. A. Halgren, R. A. Friesner, H. S. Beard, L. L. Frye, W. T. Pollard and J. L.
440 Banks, Journal of medicinal chemistry, 2004, 47, 1750-1759.
- 441 43. R. B. M. R. A. Friesner, M. P. Repasky, L. L. Frye, J. R. Greenwood, T. A. Halgren, P. C.
442 Sanschagrín and D. T. Mainz, Journal of medicinal chemistry, 2006, 49,
443 6177-6196.
- 444 44. Yan, X. Li, J. Liu, Z. Zheng, M. Ge, H. Xu, J, Journal of Chemical Information and
445 Modeling, 2013, 53, (8), 1967-78.
- 446 45. J. Sun, P.-C. Lv, Y. Yin and R.-J. Yuan, PloS one, 2013, 8, e9751.
- 447 46. M. Brvar, A. Perdih, M. Renko, G. Anderluh, D. Turk and T. Solmajer, J Med Chem,
448 2012, 55, 6413-6426.