

# 1 ***MetaMed: Linking microbiota functions with medicine***

## 2 **therapeutics**

3

### 4 **Supplementary Methods**

5

#### 6 **1 BGC and drug annotation data collection**

7 Host microbes information like biosynthetic gene clusters (BGCs) and their molecular  
8 products are downloaded from the *Minimum Information about a Biosynthetic Gene cluster*  
9 (*MIBiG*) specification (1) (version 1.3), which includes 827 microbes, 1,157 BGCs and 1,157  
10 metabolites. We categorized these microbes into their phylum levels. The final data set  
11 included 1,157 BGCs and 1,157 metabolites for 827 microbes.

12

13 The drug data annotation information is obtained from *DrugBank* datasets (2) (version 5.0.1).

14 We annotated 8,226 drugs from three aspects: (i) The indications of 8,226 drugs obtained  
15 from *DrugBank*, (ii) the common side effects of 332 drugs selected from *SIDER* (3) (version  
16 4.1) and (iii) the immune status transitions of 1310 drug obtained from one recent study (4).

17 We further labeled drugs with ATC classification system level 1 description. These  
18 therapeutic classes include: H, systemic hormonal preparations, excluding sex hormones and

19 insulins; V, various; B, blood and blood-forming organs; P, antiparasitic products; M,  
20 musculoskeletal system; L, antineoplastic and immunomodulating agents; G, genitourinary  
21 system and sex hormones; R, respiratory system; A, alimentary tract and metabolism; D,  
22 dermatologicals; J, anti-infectives for systemic use; S, sensory organs; N, nervous system;

23 and C, cardiovascular system. Finally, we curated the drug data set with 8,226 drugs, and

24 annotated the drugs with indications, side effects, immune transitions and therapeutic classes.

## 25 **2 Calculation of *MetaMed* entity relationships**

26 We constructed a data set containing 1,157 microbe biosynthetic gene clusters (BGCs) as  
27 well as their metabolites from *MIBiG* (1) , and 8,226 drugs obtained from *DrugBank* (2). We  
28 defined a solid similarity score by considering both the molecular structure (5) and  
29 perturbation transcriptional expression profiles (6) to connect microbe functions with  
30 available drug annotation information.

31

32 To integrate the structural similarity and transcriptional expression similarity, we obtained the  
33 final similarity score of microbe-drug pairs by using the mean of structure similarity and  
34 transcriptional expression similarity, or using the structure similarity if the transcriptional  
35 expression is not available. 1,193,324 linkages were collected based on the structure  
36 similarity and finally 18,090 linkages were curated by considering both structure similarity  
37 and transcriptional perturbation profiles similarity.

38

39 Based on the microbe-drug pairs and drug annotations, the scores of microbes with disease  
40 treatment effects or side effects are defined directly as the similarity score of microbe-drug  
41 pairs. The score of microbes with the impacts on human immune transition is defined as:

42  $Score(\text{microbe-immune pair}) = Score(\text{microbe-drug pair}) \times Score(\text{drug-immune pair})$  [1]

43

### 44 **2.1 Calculation of microbe-drug pairs structure similarity**

45 To estimate the structural similarity between microbe metabolite and drug, we used the  
46 simplified molecular-input line-entry system (SMILES) (7) notation to describe these  
47 compounds, and calculated the similarity scores from SMILES notations. The SMILES  
48 notations of drugs and microbe metabolites are translated from their common name by  
49 *Chemspipy* (8). The structural similarity between two compounds is calculated as the

50 Tanimoto coefficient (9) of SMILES structural fingerprint by *Open Babel* (10). Given a  
51 compound *A* and a compound *B*, the Tanimoto coefficient for binary vectors is defined as:

$$52 \frac{N(A,B)}{N(A)+N(B)-N(A,B)} \quad [2]$$

53 where,  $N(A)$  and  $N(B)$  are the number of bits set on ('1' bits) in molecular fingerprints *A* and  
54 *B* respectively, and  $N(A,B)$  is the number of bits shared by *A* and *B*.

55

## 56 **2. 2 Calculation of microbe-drug pairs transcriptional expression similarity**

57 To predict transcriptional expression similarity of microbe-drug pairs, we obtained Drug-  
58 induced transcriptional profile changes determined from human cells lines from NIH LINCS  
59 program (6). We processed the whole data sets, which include more than 110,000  
60 experiments tested over 8,000 compounds. We collected microbe transcriptional perturbation  
61 data by linking microbial metabolites to the LINCS compounds with the same compound  
62 names or same compound structures (SMILES notations). Then, we selected the up-regulated  
63 genes and down-regulated genes of the corresponding microbe metabolite from level 5 data  
64 of experiments in the same cell line at the same time point. We searched each metabolite for  
65 mimic drugs and their transcriptional expression similarity scores (cosine distance between  
66 metabolite and drug gene signatures) by LINCS API. Finally we filtered the mimic drugs by  
67 the same cell line and the same time point of input metabolite, and obtained the microbe-drug  
68 pairs and their similarity scores.

69

## 70 **3 Global landscape of *MetaMed***

71 To examine the global landscape of the *MetaMed*, we first filtered the data by microbe-drug  
72 pairs similarity score. The low similarity score indicates a low correlation between drugs and  
73 microbes, and holds limit clues for identifying the real relationship between drugs and

74 microbes. We set the similarity cutoff of 0.6 for drug and microbe relationship for global  
75 analysis. Then we annotated the drug with ATC classification system level 1 description. We  
76 also categorized the microbe with phylum level. We obtained pairs from 366 microbes and 741  
77 drugs in total, including 14 classes of drugs and 9 classes of microbes. By applying our  
78 algorithm QUalitative BIClustering algorithm (11) using the *QUBIC R* packages (version  
79 1.6.4 <https://bioconductor.org/packages/release/bioc/html/QUBIC.html>), we identified three  
80 potential biclusters, and they were visualized with heatmap.

81

#### 82 **4 Validation of the entity relationships in *MetaMed***

83 According to drug annotation information, we demonstrated the utility of the *MetaMed* from  
84 the following five aspects: meaningful microbe-drug linkings, microbes with disease  
85 treatment effects, microbes with side effects, microbes with the impacts on immune transition  
86 and finally identification of combination drug usage for disease treatment. We selected the  
87 pairs with similarity cutoff over 0.9 for validations.

88

89 We first validated the predicted metabolite-drug links. Two types of pairs are considered. If  
90 the similarity scores of microbe-drug pairs are 1.0, it indicates that these microbes can  
91 generate exactly the same drugs as the secondary metabolites. If the similarity score is lower  
92 than 1.0 while still maintains high similarity, it indicates that these microbes should have  
93 similar therapeutic indications as those of the corresponding drugs. For the first type, we  
94 found these microbes are already annotated to produce the corresponding drugs as indicated  
95 in *DrugBank*. For the second type of microbe-drug pairs, we validated them by published  
96 literature evidence.

97

98 For other identifications including microbes with disease treatment effects, microbes with  
99 side effects and microbes with the impacts on immune transition, we predicted the microbe  
100 impact on human directly. If some of the microbes exist in environment and they cannot  
101 survive in gut, we took their metabolites as the compound which may have impacts on human  
102 health. Then we validated them by available databases or published literature evidence.

103

104 To validate the combination drug usage, we validated the applications in disease treatment by  
105 data obtained from immune checkpoint therapy and Metagenome-wide association studies  
106 (MWAS) (12, 13). First, we combined *MetaMed* with two recent studies of the association of  
107 gut microbes in regulating the efficacy of anti-CTLA-4 and anti-PDL1 cancer therapy (14,  
108 15). By investigation of the outgrowth microbes, we identified drugs with the similar  
109 functions as the secondary metabolites, and these drugs may be taken as the the potential drug  
110 combination for treatments. Next we combined *MetaMed* with MWAS and obtained raw  
111 sequence data (ERP002469) of healthy and T2D samples from MWAS. We processed the  
112 sequence data by *Metapipe* and obtained the abundance of OTUs. The differential OTU  
113 analysis was performed by *edgeR*(16). The outgrowth OTUs with  $p\text{-value} < 0.05$  and  
114  $\log\text{FC} > 1.5$  were selected as differential expressed OTUs. By such analysis, we predicted  
115 potential combination drug usage for T2D treatments.

116

## 117 **5 Visualization of *MetaMed* entity relationships**

118 *Circos* plots (17) are created using the *circlize R* package (version 0.0.7  
119 <https://github.com/jokergoo/circlize>). Network diagrams were produced using *Cytoscape*  
120 (18).

121 All other plots are created using the *R* statistical package.

## 122 **6 Building *MetaMed* platform**

123 We built the web system *MetaMed* V1.1 (<http://metamed.rwebox.com/index>) for browsing  
124 various microbe functions and medicine annotations linkage information, as well as  
125 calculating such linkages directly from personalized metagenomics sequencing data  
126 (*Metapipe*, <https://github.com/adamtongji/metapipe>).  
127

128 *Metapipe* can analyze two types of human sequence data, i.e., WGS data and 16S rRNA data.  
129 For the WGS data, we set the reference sequences from biosynthetic gene clusters (BGCs)  
130 sequences of microbes. The mapping step was performed by *Bowtie 2* (19) and *Samtools* (20).  
131 The output includes microbes and their expression estimated by read counts. Users can filter  
132 the microbe results from *Metapipe* by setting the cutoff of minimum mapped read counts. For  
133 the 16S rRNA data, *Metapipe* uses *QIIME* (21) to estimate the microbe abundance and  
134 obtains the result table including microbes and their corresponding estimated abundance  
135 information. Users can filter the microbe results from *Metapipe* by setting the abundance  
136 cutoff. The output of *Metapipe* can be directly used in the *MetaMed* online platform.  
137

## 138 References

- 139 1. **Medema MH, Kottmann R, Yilmaz P, Cummings M, Biggins JB, Blin K, de**  
140 **Bruijn I, Chooi YH, Claesen J, Coates RC, Cruz-Morales P, Duddela S,**  
141 **Düsterhus S, Edwards DJ, Fewer DP, Garg N, Geiger C, Gomez-Escribano JP,**  
142 **Greule A, Hadjithomas M, Haines AS, Helfrich EJN, Hillwig ML, Ishida K, Jones**  
143 **AC, Jones CS, Jungmann K, Kegler C, Kim HU, Kötter P, Krug D, Masschelein J,**  
144 **Melnik AV, Mantovani SM, Monroe EA, Moore M, Moss N, Nützmänn H-W, Pan**  
145 **G, Pati A, Petras D, Reen FJ, Rosconi F, Rui Z, Tian Z, Tobias NJ, Tsunematsu Y,**  
146 **Wiemann P, Wyckoff E, Yan X, Yim G, Yu F, Xie Y, Aigle B, Apel AK, Balibar**

147 **CJ, Balskus EP, Barona-Gómez F, Bechthold A, Bode HB, Borriss R, Brady SF,**  
148 **Brakhage AA, Caffrey P, Cheng Y-Q, Clardy J, Cox RJ, De Mot R, Donadio S,**  
149 **Donia MS, van der Donk WA, Dorrestein PC, Doyle S, Driessen AJM, Ehling-**  
150 **Schulz M, Entian K-D, Fischbach MA, Gerwick L, Gerwick WH, Gross H, Gust**  
151 **B, Hertweck C, Höfte M, Jensen SE, Ju J, Katz L, Kaysser L, Klassen JL, Keller**  
152 **NP, Kormanec J, Kuipers OP, Kuzuyama T, Kyrpides NC, Kwon H-J, Lautru S,**  
153 **Lavigne R, Lee CY, Linquan B, Liu X, Liu W, Luzhetskyy A, Mahmud T, Mast Y,**  
154 **Méndez C, Metsä-Ketelä M, Micklefield J, Mitchell DA, Moore BS, Moreira LM,**  
155 **Müller R, Neilan BA, Nett M, Nielsen J, O'Gara F, Oikawa H, Osbourn A,**  
156 **Osburne MS, Ostash B, Payne SM, Pernodet J-L, Petricek M, Piel J, Ploux O,**  
157 **Raaijmakers JM, Salas JA, Schmitt EK, Scott B, Seipke RF, Shen B, Sherman**  
158 **DH, Sivonen K, Smanski MJ, Sosio M, Stegmann E, Süßmuth RD, Tahlan K,**  
159 **Thomas CM, Tang Y, Truman AW, Viaud M, Walton JD, Walsh CT, Weber T,**  
160 **van Wezel GP, Wilkinson B, Willey JM, Wohlleben W, Wright GD, Ziemert N,**  
161 **Zhang C, Zotchev SB, Breitling R, Takano E, Glöckner FO. 2015. Minimum**  
162 **Information about a Biosynthetic Gene cluster. Nat Chem Biol 11:625–631.**

163 2. **Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, Sajed T, Johnson**  
164 **D, Li C, Sayeeda Z, Assempour N, Iynkkaran I, Liu Y, Maciejewski A, Gale N,**  
165 **Wilson A, Chin L, Cummings R, Le D, Pon A, Knox C, Wilson M. 2018.**  
166 **DrugBank 5.0: a major update to the DrugBank database for 2018. Nucleic Acids Res**  
167 **46:D1074–D1082.**

168 3. **Kuhn M, Letunic I, Jensen LJ, Bork P. 2016. The SIDER database of drugs and side**  
169 **effects. Nucleic Acids Res 44:D1075–9.**

170 4. **Kidd BA, Wroblewska A, Boland MR, Agudo J, Merad M, Tatonetti NP, Brown**

- 171 **BD, Dudley JT.** 2016. Mapping the effects of drugs on the immune system. *Nat*  
172 *Biotechnol* **34**:47–54.
- 173 5. **Boström J, Hogner A, Schmitt S.** 2006. Do structurally similar ligands bind in a  
174 similar fashion? *J Med Chem* **49**:6716–6725.
- 175 6. **Duan Q, Reid SP, Clark NR, Wang Z, Fernandez NF, Rouillard AD, Readhead B,**  
176 **Tritsch SR, Hodos R, Hafner M, Niepel M, Sorger PK, Dudley JT, Bavari S,**  
177 **Panchal RG, Ma'ayan A.** 2016. L1000CDS2: LINCS L1000 characteristic direction  
178 signatures search engine. *NPJ Syst Biol Appl* **2**:257.
- 179 7. **Pearson WR, Lipman DJ.** 1988. Improved tools for biological sequence comparison.  
180 *Proc Natl Acad Sci USA* **85**:2444–2448.
- 181 8. **Pence HE, Williams A.** 2010. ChemSpider: An Online Chemical Information  
182 Resource. *Journal of Chemical Education* **87**:1123–1124.
- 183 9. **Willett P, Barnard JM, Downs GM.** 1998. Chemical Similarity Searching. *Journal of*  
184 *Chemical Information and Computer Sciences* **38**:983–996.
- 185 10. **O'Boyle NM, Morley C, Hutchison GR.** 2008. Pybel: a Python wrapper for the  
186 OpenBabel cheminformatics toolkit. *Chemistry Central Journal* **2**:5.
- 187 11. **Li G, Ma Q, Tang H, Paterson AH, Xu Y.** 2009. QUBIC: a qualitative biclustering  
188 algorithm for analyses of gene expression data. *Nucleic Acids Res* **37**:e101–e101.
- 189 12. **Wang J, Jia H.** 2016. Metagenome-wide association studies: fine-mining the  
190 microbiome. *Nat Rev Microbiol* **14**:508–522.
- 191 13. **Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B,**



- 192 **Nielsen J, Bäckhed F.** 2013. Gut metagenome in European women with normal,  
193 impaired and diabetic glucose control. *Nature* **498**:99–103.
- 194 14. **Vetizou M, Pitt JM, Daillere R, Lepage P, Waldschmitt N, Flament C,**  
195 **Rusakiewicz S, Routy B, Roberti MP, Duong CPM, Poirier-Colame V, Roux A,**  
196 **Becharef S, Formenti S, Golden E, Cording S, Eberl G, Schlitzer A, Ginhoux F,**  
197 **Mani S, Yamazaki T, Jacquelot N, Enot DP, Berard M, Nigou J, Opolon P,**  
198 **Eggermont A, Woerther PL, Chachaty E, Chaput N, Robert C, Mateus C,**  
199 **Kroemer G, Raoult D, Boneca IG, Carbonnel F, Chamillard M, Zitvogel L.** 2015.  
200 Anticancer immunotherapy by CTLA-4 blockade relies on the gut microbiota. *Science*  
201 **350**:1079–1084.
- 202 15. **Sivan A, Corrales L, Hubert N, Williams JB, Aquino-Michaels K, Earley ZM,**  
203 **Benyamin FW, Lei YM, Jabri B, Alegre M-L, Chang EB, Gajewski TF.** 2015.  
204 Commensal Bifidobacterium promotes antitumor immunity and facilitates anti-PD-L1  
205 efficacy. *Science* **350**:1084–1089.
- 206 16. **Robinson MD, McCarthy DJ, Smyth GK.** 2010. edgeR: a Bioconductor package for  
207 differential expression analysis of digital gene expression data. *Bioinformatics* **26**:139–  
208 140.
- 209 17. **Gu Z, Gu L, Eils R, Schlesner M, Brors B.** 2014. circlize implements and enhances  
210 circular visualization in R. *Bioinformatics* **30**:2811–2812.
- 211 18. **Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N,**  
212 **Schwikowski B, Ideker T.** 2003. Cytoscape: a software environment for integrated  
213 models of biomolecular interaction networks. *Genome Res* **13**:2498–2504.
- 214 19. **Langmead B, Salzberg SL.** 2012. Fast gapped-read alignment with Bowtie 2. *Nat*

215 Methods **9**:357–359.

216 20. **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis**  
217 **G, Durbin R, 1000 Genome Project Data Processing Subgroup.** 2009. The  
218 Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**:2078–2079.

219 21. **Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK,**  
220 **Fierer N, Peña AG, Goodrich JK, Gordon JI, Huttley GA, Kelley ST, Knights D,**  
221 **Koenig JE, Ley RE, Lozupone CA, McDonald D, Muegge BD, Pirrung M, Reeder**  
222 **J, Sevinsky JR, Turnbaugh PJ, Walters WA, Widmann J, Yatsunenko T,**  
223 **Zaneveld J, Knight R.** 2010. QIIME allows analysis of high-throughput community  
224 sequencing data. *Nat Methods* **7**:335–336.

225