

MetaPrism: A Toolkit for Joint Taxa/Gene Analysis of Metagenomic Sequencing Data

Jiwoong Kim

University of Texas Southwestern Medical Center at Dallas

Shuang Jiang

Southern Methodist University

Guanghua Xiao

University of Texas Southwestern Medical Center at Dallas

Yang Xie

University of Texas Southwestern Medical Center at Dallas

Dajiang Liu

Pennsylvania State University

Qiwei Li

University of Texas at Dallas

Andrew Koh

University of Texas Southwestern Medical Center at Dallas

Xiaowei Zhan (✉ Xiaowei.Zhan@Utsouthwestern.edu)

University of Texas Southwestern Medical Center <https://orcid.org/0000-0002-6249-7193>

Software article

Keywords: metagenomics sequence analysis, joint analysis, microbiome biomarker

Posted Date: February 8th, 2020

DOI: <https://doi.org/10.21203/rs.2.22868/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **MetaPrism: A Toolkit for Joint Taxa/Gene Analysis of**
2 **Metagenomic Sequencing Data**

3 Jiwoong Kim (Jiwoong.Kim@UTSouthwestern.edu)^{1#},

4 Shuang Jiang (Shuangj@mail.smu.edu)^{2#},

5 Guanghua Xiao (Guanghua.Xiao@UTSouthwestern.edu)^{1,3,4},

6 Yang Xie (Yang.Xie@UTSouthwestern.edu)^{1,3,4},

7 Dajiang J. Liu (dxl46@psu.edu)⁵,

8 Qiwei Li (Qiwei.Li@utdallas.edu)⁶,

9 Andrew Koh (Andrew.Koh@UTSouthwestern.edu)^{3,7,8},

10 Xiaowei Zhan (Xiaowei.Zhan@UTSouthwestern.edu)^{1,3,9*}

11 ¹Quantitative Biomedical Research Center, Department of Population and Data Sciences,
12 University of Texas Southwestern Medical Center, Dallas, TX, 75390

13 ²Department of Statistical Science, Southern Methodist University, Dallas, TX 75275

14 ³Harold C. Simmons Cancer Center, University of Texas Southwestern Medical Center, Dallas,
15 Texas, 75390

16 ⁴Department of Bioinformatics, University of Texas Southwestern Medical Center, Dallas,
17 Texas, 75390

18 ⁵Department of Public Health Sciences, Pennsylvania State University, Hershey, Pennsylvania,
19 17033

20 ⁶Department of Mathematical Science, The University of Texas at Dallas, Dallas, Texas, 75080

1 ⁷Department of Microbiology, University of Texas Southwestern Medical Center, Dallas, Texas,
2 75390

3 ⁸Department of Pediatrics, University of Texas Southwestern Medical Center, Dallas, Texas,
4 75390

5 ⁹Center for Genetics of Host Defense, University of Texas Southwestern Medical Center, Dallas,
6 Texas, 75390

7

8 *To whom correspondence should be addressed

9 #These authors contributed equally to this work

10

1 **Abstract**

2 **Background**

3 In microbiome research, metagenomic sequencing generates enormous amounts of data. These
4 data are typically classified into taxa for taxonomy analysis, or into genes for functional analysis.
5 However, a joint analysis where the reads are classified into taxa-specific genes is often
6 overlooked.

7 **Result**

8 To enable the analysis of this biologically meaningful feature, we developed a novel
9 bioinformatic toolkit, MetaPrism, which can analyze sequence reads for a set of joint taxa/gene
10 analyses: 1) classify sequence reads and estimate the abundances for taxa-specific genes; 2)
11 tabularize and visualize taxa-specific gene abundances; 3) compare the abundances between
12 groups, and 4) build prediction models for clinical outcome. We illustrated these functions using
13 a published microbiome metagenomics dataset from patients treated with immune checkpoint
14 inhibitor therapy and showed the joint features can serve as potential biomarkers to predict
15 therapeutic responses.

16 **Conclusions**

17 MetaPrism is a toolkit for joint taxa and gene analysis. It offers biological insights on the taxa-
18 specific genes on top of the taxa-alone or gene-alone analysis.

19 MetaPrism is open source software and freely available at

20 <https://github.com/jiwoongbio/MetaPrism>. The example script to reproduce the manuscript is
21 also provided in the above code repository.

22

23 **Keywords:** metagenomics sequence analysis, joint analysis, microbiome biomarker

1

2 **Background**

3 The human microbiome consists of ~39 trillion bacteria and influences host health. Recently, the
4 use of metagenomic sequencing has become increasingly popular as a more unbiased approach to
5 gut microbiome profiling as compared to 16S rRNA sequencing. A common approach to
6 comparing differences in the gut microbiome between groups (cases and controls) is to identify
7 significant differences in either taxa or microbial genes. Several popular bioinformatic tools have
8 been developed for this purpose, including MetaPhlan2 [1], Kraken [2], HUMAnN2 [3], and
9 FMAP [4] (**Supplementary Table 1**). However, these tools analyze either taxonomic abundances
10 (taxonomic profiling) or gene abundances (function profiling) separately. As each microorganism
11 carries its own genes, taxonomic and functional profiling results are not intrinsically independent.
12 In fact, recent discoveries demonstrated that taxon-specific genes have a causative role in disease
13 progression and treatment responses. For example, Duan et al. found that a specific *Enterococcus*
14 *faecalis* carrying the cytolysin gene promotes alcoholic liver disease[5]. Simms-Waldrip et al. found
15 that the antibiotic resistance genes in the graft-versus-host-disease patients are enriched for
16 *Klebsiella*[6]. Therefore, joint analysis, where taxonomy and functional features are analyzed
17 together, could provide useful biological and clinical insights [7]. However, bioinformatics tools
18 for joint analyses are comparatively lacking.

19 Our innovation in this manuscript is to define and utilize joint taxa/gene features via
20 bioinformatics approach, with the goal of offering biologically interpretable findings. For example,
21 our method characterizes the genes discovered for each species. This allows to quantitative
22 analysis of this species-specific gene, which is usually not readily available. Our approach is

1 initiated from *de novo* assembled contigs which are both taxonomically and functionally annotated.
2 Our simulations showed this method can accurately detect bacterial species and their carried genes.
3 In a recent review article[7], Langille prompted that understanding the gene contents at species
4 level can offer better interpretation than using the taxon or gene content alone, and potentially
5 provide better prediction outcomes. This confirmed that the joint feature is useful for general
6 microbiome studies. Our tool provided these joint features as the first step for a wide range of
7 downstream analysis tasks. For example, we demonstrated that the quantity of taxa-specific gene
8 abundances is a potentially useful biomarker to predict the immunotherapy responses.

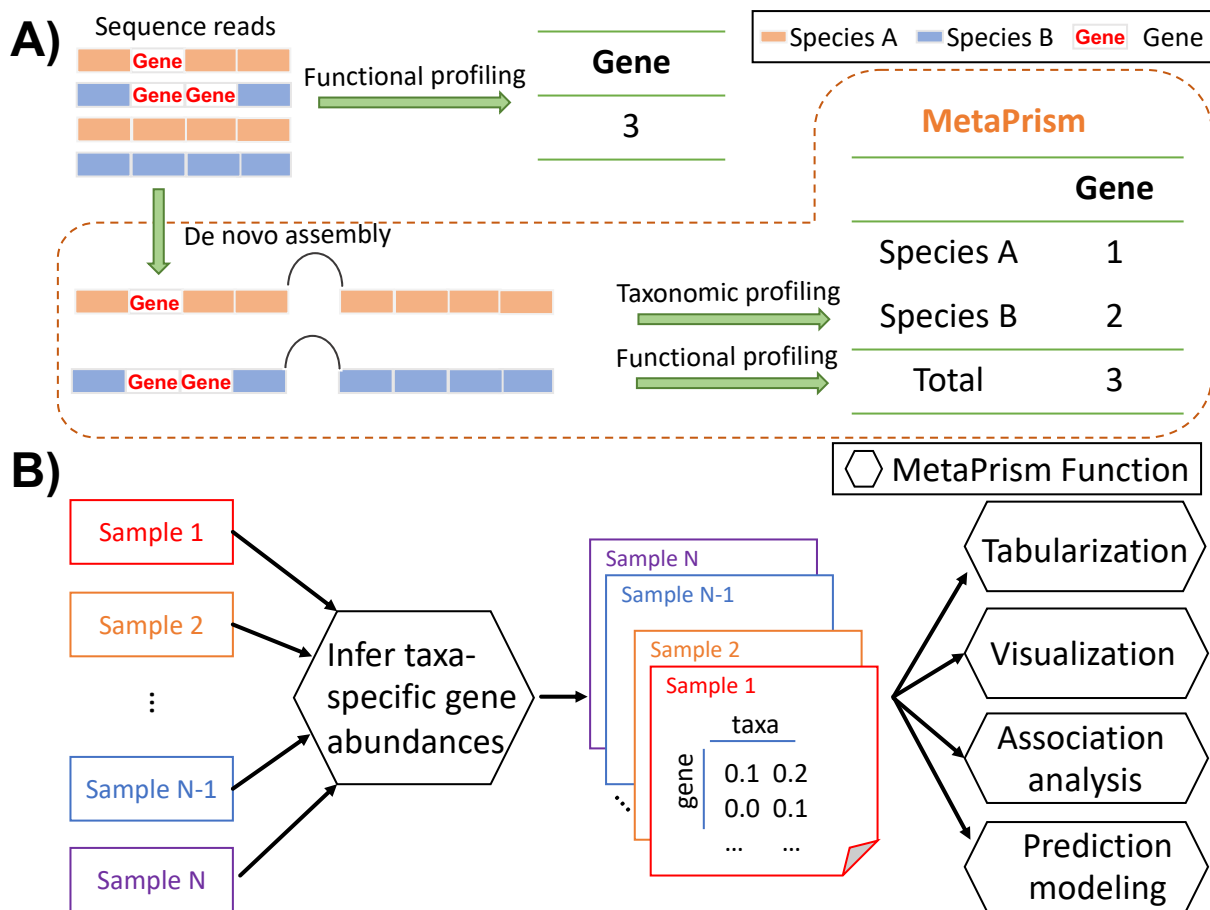
9 To facilitate joint analysis, we developed MetaPrism, a novel bioinformatics tool to (1) classify
10 metagenomic sequence reads into both taxa and gene level, (2) normalize the taxa-specific gene
11 abundances within samples, (3) tabularize or visualize these joint features, (4) perform
12 comparative microbiome studies, and (5) build prediction models for clinical outcomes.
13 MetaPrism is open-sourced and is available at <https://github.com/jiwoongbio/MetaPrism>. Given
14 the advantages of joint analysis, MetaPrism is a useful tool for microbiome metagenomic sequence
15 studies.

16

17 **Implementation**

18 MetaPrism is a toolkit for joint analysis tasks. At its core, MetaPrism will infer the taxa and gene
19 for each metagenome sequence read. One approach is to align each read to bacterial nucleotide
20 reference genomes to obtain its taxonomy and align it to a protein database to obtain its gene
21 functions. However, this approach is technical challenging: due to the short lengths of Illumina
22 sequence reads and the high sequence similarities between bacteria genomes, alignment of short
23 reads is not feasible. We thus developed a novel algorithm (**Figure 1A**) to tackle this challenge.

1 First, we perform *de novo* assembly for each sample using metaSPAdes [8] with all metagenomic
 2 sequence reads to obtain long contigs. As these contigs are much longer than sequence reads, that
 3 allows for accurate taxonomical and functional profiling.
 4



5 **Figure 1. An illustration of the algorithm and the functions in MetaPrism.** A) Illustration of the MetaPrism
 6 algorithm to infer taxa-specific gene abundances. Function profiling alone infers that three reads are mapped to a
 7 gene, but cannot provide further taxonomic information. MetaPrism can estimate two reads are from species A and
 8 one read is from species B; B) An overview of the joint analysis workflow in MetaPrism. Hexagons represent
 9 functions in MetaPrism.

10

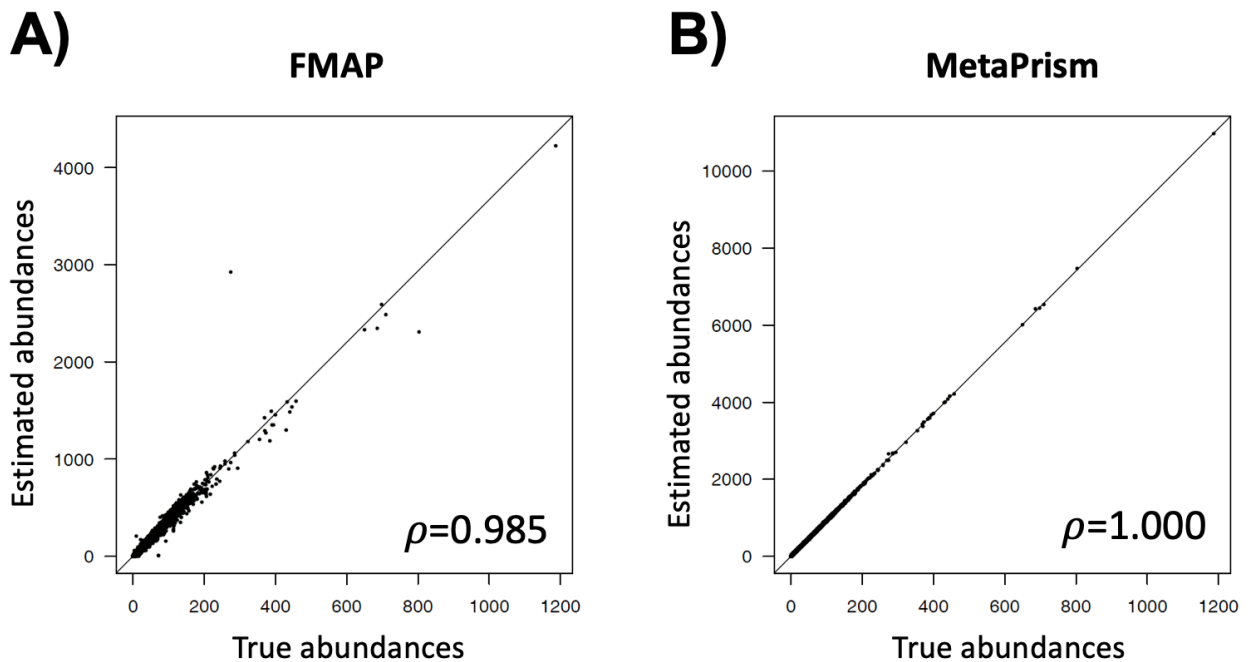
11 Second, we identify the taxonomy of these contigs. All the contigs are aligned to a large
 12 reference database of more than 4,000 bacterial genomes using centrifuge [9]. Ambiguous
 13 alignments will be filtered out from the subsequent analysis.

1 Third, we identify genes and their locations from the contigs. We detect the open reading frames
2 from the contigs, translated the nucleotide bases to amino acids, and aligned them using
3 DIAMOND [10] to a protein database. To comprehensively investigate all bacteria genes, either
4 KEGG protein databases that include protein sequences from KEGG orthologue genes [11] or
5 KFU (KEGG orthology with UniProt protein sequences) [4], can be utilized. By default, we
6 required minimum coverage of 0.8 to ensure good protein alignments.

7 Lastly, we calculate and normalize gene abundance within-sample. We align metagenomic
8 sequence reads to the contigs using BWA [12], and count the number of aligned reads located in
9 the genes of interest. We calculate the read depth normalized by contig length, and this quantity is
10 denoted as mean depth to represent the gene abundances. Larger numbers often indicate higher
11 gene abundance. Other abundance statistics, such as FPKM (Fragment Per Kilobase of transcript
12 per Million reads) or depth per genome (normalized read depth per taxa genome length), are also
13 provided.

14 Through the above steps, the gene abundances are associated with taxonomy information. To
15 assess the accuracy of these estimations, we conducted a simulation study. First, we selected 115
16 bacterial species with complete reference genomes and downloaded their sequences from NCBI
17 FTP (<ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/bacteria>). Then we simulated shotgun
18 metagenomic sequencing reads and generated at 10X coverage to resemble typical read length
19 from the Illumina (100 bp) using ART [13] "art_illumina --in
20 bacteria_complete_reference_genome.fasta --len 100 --fcov 10 --mflen 200 --sdev 50 --noALN --
21 out bacteria_complete_reference_genome.illumina". We re-assembled the simulated metagenome
22 sequences using metaSPAdes [8]. To evaluate the gene abundances calculated by different
23 methods, we compared MetaPrism to our previous program named FMAP using translation

1 alignment (BLASTX), as our previous approach was shown to accurately report gene abundances
2 [4]. The true abundances were determined by the KEGG ortholog (KO) abundance of the genes in
3 the reference genomes by aligning them to the KEGG protein database using DIAMOND [10]. In
4 **Figure 2**, The scatterplot visualized the true abundances (x-axis) and the estimated abundances (y-
5 axis), and MetaPrism showed higher correlation (correlation coefficient = 1.000) compared to
6 FMAP (correlation coefficient = 0.985).



7
8 **Figure 2. Comparison of gene abundances reported by FMAP and MetaPrism.** We used simulations to compare
9 the estimated gene abundances using FMAP and MetaPrism. The Pearson correlation coefficients between true
10 abundances and the software estimated abundances were listed on the bottom right.
11

12 In brief, we simulated metagenomic sequence reads from known species and inferred the gene
13 abundances using FMAP (Kim, et al., 2016b) and MetaPrism. This benchmark showed that gene
14 abundances inferred by MetaPrism were accurate and achieved the highest correlation between
15 inferred abundances and true abundances (**Figure 2**).

16 Based on these joint features, MetaPrism provided the following downstream joint analysis
17 functions (demonstrated in **Figure 1B**): 1) tabularize the abundances of these features

1 (MetaPrism_table.pl); 2) visualize the features in heatmaps (MetaPrism_heatmap.pl); 3) compare
2 the taxa-specific genes abundances across different experimental conditions such as case-control
3 studies (MetaPrism_compare.pl); 4) indicate which features may serve as potential biomarkers in
4 a prediction model (MetaPrism_predict.pl). A list of available functions, command line, and major
5 customization options in MetaPrism are listed in **Supplementary Table 2**.

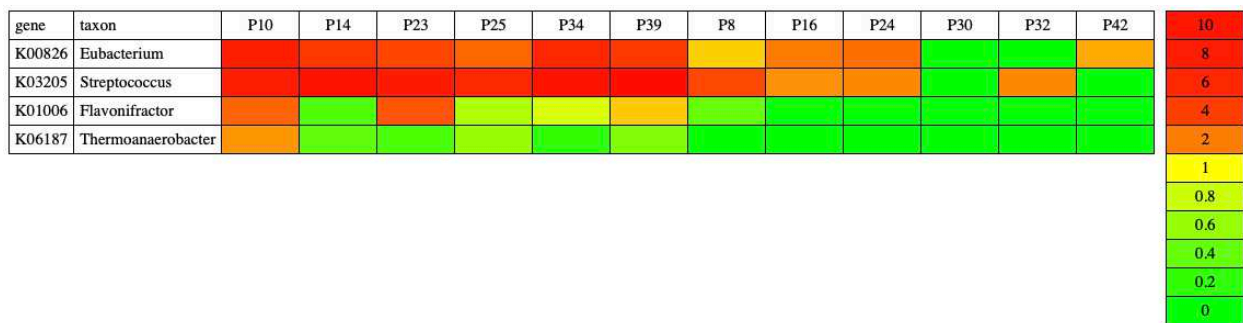
6 **Results**

7 The gut microbiome plays an important role in modulating immune checkpoint therapy [14]. Here
8 we demonstrated a joint analysis using MetaPrism to build a therapy response prediction model.
9 We collected stool samples of 12 melanoma patients before anti-PD1 (pembrolizumab) therapy
10 and performed metagenomic sequencing. 6 patients responded to the therapy and 6 did not.

11 Starting from the metagenomic sequence reads, we performed quality control, including removal
12 of human contamination as previously described [14]. Then all remaining sequence reads were
13 processed in MetaPrism (detailed analysis steps provided in **Supplementary Texts**). On average,
14 MetaPrism inferred the taxonomy and gene features for 1.2 billion reads per sample. Next,
15 MetaPrism normalized the reads within samples by reporting the mean depth per assembled contig.
16 The taxa-specific gene abundances were ranked using a random forest model with leave-one-out
17 cross-validation. This prediction model reached 69% accuracy to predict the immunotherapy
18 responses, which is higher than the accuracy using taxa features alone (54%) or gene features alone
19 (62%). Furthermore, it detected four joint features with variable importance greater than 50%.
20 MetaPrism visualized these abundances with red to green colors representing the depth values
21 (**Figure 3**). The most important feature is the K00826 gene (branched-chain amino acid
22 aminotransferase) from the genus Eubacterium (**Table 1**). It is a novel joint feature that may serve
23 as a potential biomarker for cancer therapy.

1

2 **Figure 3. Heatmap of joint features for predicting immune checkpoint therapy response.** We used
 3 MetaPrism_heatmap.pl to visualize four joint features (taxa-specific gene abundances) in the immune checkpoint
 4 therapy study. The colors from red to green represent the increased gene abundances, the mean depth normalized by
 5 the contig lengths. P10, P14, P23, P25, P34, P39 are patients who respond to the therapy; P8, P16, P24, P30, P32, P42
 6 are patients having progressive outcomes.



7

8 **Table 3 Prediction models and performances for taxonomical analysis, functional analysis, and joint analysis.**

9 We tabularized the details of prediction models used in three types of analyses and their prediction performances.

10

	Taxonomic	Functional	Joint
	profiling	profiling	profiling
Model	Random forest	Random forest	Random forest
Number of trees	500	500	500
Number of features	1,048	5,227	62,086
Top features (if variable importance > 50%) #			
1st feature	Chondromyces (100)	K07705 (100)	K00826 Eubacterium (100)
2nd feature	Roseateles (65)	-	K03205 Streptococcus (89)
3rd feature	-	-	K01006 Flavonifractor (81)
4th feature	-	-	K06187 Thermoanaerobacter (74)
Accuracy*	53.8%	61.5%	69.2%

11

12 #: The variable importance values are listed in parentheses

13 *: Prediction accuracy was evaluated using leave-one-out cross-validations

1 In terms of computation, all the above analyses can be accomplished on a standard computation
2 cluster (e.g., 128GB memory with 2 GB hard drive space per sample).

3 **Conclusion**

4 We present a novel bioinformatics tool, MetaPrism. It implements functions to quantify the joint
5 features (both taxonomic and functional) from metagenomic sequence reads, as well as other
6 functions for downstream data analyses. We demonstrate that the joint features can provide novel
7 insights to understand the microbial role in a cancer immunotherapy study.

8 MetaPrism is flexible and can be customized. For example, to study species-specific antibiotic
9 resistance genes (ARGs), a reference protein database with ARGs, such as ARDB [15] or CARD
10 [16], can be used. MetaPrism can infer taxa-specific ARGs, thus enabling joint resistome profiling.
11 In a GVHD study, with the interests to study the patients' resistome, we performed the analysis
12 using the ARDB in MetaPrism and found increased abundances of antibiotic-resistance genes (e.g.,
13 *mdtG*, *AcrA*, *AcrB*, and *TolC*) in *Klebsiella* and *E. coli* in the GVHD patients compared with the
14 abundances in non-GVHD patients. This finding may hint optimal antibiotic prescription for better
15 management of GVHD.

16 MetaPrism characterizes the joint features based on the contigs that are *de novo* assembled from
17 metagenomic sequence reads. This is a distinct feature compared with other software. For example,
18 HUMAnN2 used a tiered search strategy that relied on a curated reference database for organism-
19 specific genes[3]. As human microbiome contains trillions of microbial genes, reference databases
20 can be inadequate to enumerate the organism-specific genes. Thus we designed the MetaPrism to
21 reduce the dependency on curated reference databases. The tradeoff for this decision is that
22 MetaPrism requires more computational resources for the *de novo* assembling step.

1 In all, MetaPrism is free and useful software to facilitate joint analyses and it is suitable for
2 general microbiome studies. Researchers can expect MetaPrism to quantify species-specific gene
3 abundances and use these interpretable features in association studies and prediction tasks.
4

5 **Availability and requirements**

6 **Project name:** MetaPrism

7 **Project home page:** <https://github.com/jiwoongbio/MetaPrism>

8 **Operating system(s):** Platform independent

9 **Programming language:** Perl, R

10 **Other requirements:** None

11 **License:** GNU GPL

12 **Any restrictions to use by non-academics:** None

13 **Declarations**

14 **Funding**

15 This work has been supported by the following grants: NIH R01 [R01GM115473 (YX), R01GM126479 (DJL,
16 XZ)]; Cancer Center: [P30CA142543 (YX, XZ)]; Specialized Programs of Research Excellence [P50CA070907
17 (YX, XZ)].

18 **Availability of data and materials**

19 The metagenomic shotgun sequence dataset are available from the NCBI BioProject PRJNA397906. The
20 treatment responses for the 12 patients as well as the analysis codes were available in the Supplementary section
21 2. The simulation procedure are available in the Supplementary section 1. The source codes of MetaPrism

1 software is available at: <https://github.com/jiwoongbio/MetaPrism>. That resource contains the software
2 requirements, usage example and documentations for all MetaPrism components (e.g., download bacterial
3 database, quantify species-specific gene abundances, build association models and prediction models, tabularize
4 results and visualize results in heatmap).

5 **Authors' contributions**

6 JK, SJ, and XZ conceived of the project and wrote the first draft of the manuscript. GX, YX, AY, and XZ
7 coordinated and oversaw the study. QL, DL provided critical inputs for the study. JK and SJ developed the
8 software and associated databases. All authors contributed to the review of the manuscript before submission for
9 publication. All authors read and approved the final manuscript.

10 **Ethics approval and consent to participate**

11 Not applicable.

12 **Consent for publication**

13 Not applicable.

14 **Competing interests**

15 The authors declare that they have no competing interests.

16 **Acknowledgement**

17 We thank Jessie Norris for her comments on the manuscript.
18

19 **References**

20

- 21 1. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A,
22 Huttenhower C, Segata N: **MetaPhlAn2 for enhanced metagenomic taxonomic**
23 **profiling**. *Nature methods* 2015, **12**(10):902-903.

- 1 2. Wood DE, Salzberg SL: **Kraken: ultrafast metagenomic sequence classification using**
2 **exact alignments**. *Genome biology* 2014, **15**(3):R46.
- 3 3. Franzosa EA, McIver LJ, Rahnnavard G, Thompson LR, Schirmer M, Weingart G, Lipson
4 KS, Knight R, Caporaso JG, Segata N *et al*: **Species-level functional profiling of**
5 **metagenomes and metatranscriptomes**. *Nature methods* 2018, **15**(11):962-968.
- 6 4. Kim J, Kim MS, Koh AY, Xie Y, Zhan X: **FMAP: Functional Mapping and Analysis**
7 **Pipeline for metagenomics and metatranscriptomics studies**. *BMC Bioinformatics*
8 2016, **17**(1):420.
- 9 5. Duan Y, Llorente C, Lang S, Brandl K, Chu H, Jiang L, White RC, Clarke TH, Nguyen
10 K, Torralba M *et al*: **Bacteriophage targeting of gut bacterium attenuates alcoholic**
11 **liver disease**. *Nature* 2019, **575**(7783):505-511.
- 12 6. Simms-Waldrup TR, Sunkersett G, Coughlin LA, Savani MR, Arana C, Kim J, Kim M,
13 Zhan X, Greenberg DE, Xie Y *et al*: **Antibiotic-Induced Depletion of Anti-**
14 **inflammatory Clostridia Is Associated with the Development of Graft-versus-Host**
15 **Disease in Pediatric Stem Cell Transplantation Patients**. *Biol Blood Marrow*
16 *Transplant* 2017, **23**(5):820-829.
- 17 7. Langille MGI: **Exploring Linkages between Taxonomic and Functional Profiles of**
18 **the Human Microbiome**. *mSystems* 2018, **3**(2).
- 19 8. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA: **metaSPAdes: a new versatile**
20 **metagenomic assembler**. *Genome research* 2017, **27**(5):824-834.
- 21 9. Kim D, Song L, Breitwieser FP, Salzberg SL: **Centrifuge: rapid and sensitive**
22 **classification of metagenomic sequences**. *Genome research* 2016, **26**(12):1721-1729.
- 23 10. Buchfink B, Xie C, Huson DH: **Fast and sensitive protein alignment using**
24 **DIAMOND**. *Nature methods* 2015, **12**(1):59-60.
- 25 11. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M: **KEGG for integration and**
26 **interpretation of large-scale molecular data sets**. *Nucleic acids research* 2012,
27 **40**(Database issue):D109-114.

- 1 12. Li H, Durbin R: **Fast and accurate long-read alignment with Burrows-Wheeler**
2 **transform**. *Bioinformatics*, **26**(5):589-595.
- 3 13. Huang W, Li L, Myers JR, Marth GT: **ART: a next-generation sequencing read**
4 **simulator**. *Bioinformatics* 2012, **28**(4):593-594.
- 5 14. Frankel AE, Coughlin LA, Kim J, Froehlich TW, Xie Y, Frenkel EP, Koh AY:
6 **Metagenomic Shotgun Sequencing and Unbiased Metabolomic Profiling Identify**
7 **Specific Human Gut Microbiota and Metabolites Associated with Immune**
8 **Checkpoint Therapy Efficacy in Melanoma Patients**. *Neoplasia* 2017, **19**(10):848-
9 855.
- 10 15. Liu B, Pop M: **ARDB--Antibiotic Resistance Genes Database**. *Nucleic acids research*
11 2009, **37**(Database issue):D443-447.
- 12 16. McArthur AG, Waglechner N, Nizam F, Yan A, Azad MA, Baylay AJ, Bhullar K,
13 Canova MJ, De Pascale G, Ejim L *et al*: **The comprehensive antibiotic resistance**
14 **database**. *Antimicrob Agents Chemother* 2013, **57**(7):3348-3357.
- 15

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementary.pdf](#)