

Manuscript Number:	GIGA-D-18-00184R1	
Full Title:	The metagenome of the female upper reproductive tract	
Article Type:	Research	
Funding Information:	Shenzhen Municipal Government of China (JCYJ20160229172757249)	Mrs Chen Chen
	Shenzhen Municipal Government of China (JCYJ20150601090833370)	Mrs Chen Chen
	Macau Technology Development Fund (102/2016/A3)	Mrs Chen Chen
Abstract:	<p>Background</p> <p>The human uterus is traditionally believed to be sterile, while the vaginal microbiota plays important role in fending off pathogens. Emerging evidence demonstrates the presence of bacteria beyond the vagina. However, a microbiome-wide metagenomic analysis identifying the overall microorganism communities has been lacking.</p> <p>Results</p> <p>We performed shotgun-sequencing by Illumina platform of 52 samples from the cervical canal and the peritoneal fluid of Chinese women in reproductive age. Direct annotation of sequencing reads identified the taxonomy of bacteria, archaea, fungi and viruses, confirming and extending the results from our previous study. We replicated the findings in another 24 samples from the vagina, the cervical canal, the uterus and the peritoneal fluid using BGISEQ-500 platform, revealing that microorganisms in the samples from the same individual were largely shared in the whole reproductive tract. Over 99% human sequences were detected in the 20GB raw data. After filtering, vaginal microorganisms were well covered in the generated reproductive tract gene catalogue, while the more diverse upper reproductive tract microbiota might need greater depth of sequencing and more samples to meet the full coverage scale.</p> <p>Conclusions</p> <p>Microbiota in unprecedented data for uncharted body site, female upper reproductive tract, were analyzed in this study. The community results indicated that an intra-individual continuum of all types of microorganisms gradually changed from the vagina to the peritoneal fluid. A framework was also established in this study aiming at understanding the implications of the composition and functional potential of this distinct microbial ecosystem in relation to health and disease.</p>	
Corresponding Author:	Huijue Jia	
	CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Chen Chen	
First Author Secondary Information:		
Order of Authors:	Chen Chen	
	Fei Li	
	Weixia Wei	
	Zirong Wang	
	Juanjuan Dai	
	Lilan Hao	

	Liju Song
	Xiaowei Zhang
	Liping Zeng
	Hui Du
	Huiru Tang
	Na Liu
	Huanming Yang
	Jian Wang
	Lise Madsen
	Susanne Brix
	Karsten Kristiansen
	Xun Xu
	Junhua Li
	Ruifang Wu
	Huijue Jia
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Reviewer reports:</p> <p>Reviewer #1: The manuscript has vastly improved since the first submission, although some additional information and language editing is required.</p> <p>I would like to thank the authors for including more in-depth information about the functional data; i think p-values for the aforementioned differences in gene families have to be provided; maybe marking the significant differences on the figure. I also think that since this is actually the most novel data in the paper, the Supplementary Figure 3 should be moved to the main body of the paper.</p> <p>We thank reviewer for the valuable comment. For the functional data, we evaluated the difference between the CV and PF by comparing all the functional genes from the built catalogues instead of the individual samples to avoid the deviation. Hence, we can't do the statistical analysis, such as p-value calculation. Additionally, we agreed with the suggestion from reviewer and moved Supplementary Figure 3 to the main body to make the paper more fulness and integrated.</p> <p>Authors have enough of samples to claim the gradient in microbiota over the reproductory tract; as well as enough of samples to perform functional comparisons between PF and CU samples, but believe that these data is not enough to address differences in alpha- and beta-diversity between PF/CU. I find this argument a little vague and I strongly believe that the paper would benefit from including this information, but since it initially was just my suggestion to the authors, they are free to ignore it.</p> <p>We would like to thank the reviewer's question. We consider that the sample number of PF is enough to explain the microbiota and the community function since the microbiota and its function were mainly dressed by the dominate species which could be fully explained by the data achieved from these samples. However, we claim that it is not enough to address differences in alpha- and beta-diversity between PF/CU because that the rarefaction curve of PF did not reach the saturation, while CU did. That means some species may still not been detected in the PF samples and this will cause the deviation of the bacterial alpha- and beta-diversity results.</p> <p>Line 161, figure 3, the PF line is not far from reaching the plateau - it does seem that the line is approaching the asymptote, so i believe that 'far from saturation' is an overstatement and should be toned down.</p> <p>We agree with the suggestion raised from the reviewer. After the double checking, we totally agreed with the reviewer and revised this statement to make it toned down (Line 160-162).</p>

	<p>Reviewer #2: Still think that a quick read for English would be useful for this manuscript though the writing has improved. Examples of writing that needs correcting are as follows: Lines 34-35 in the abstract: "the vaginal microbiota plays important roles ... " should be "role"; Line 91 "(the stringent selection rules ..." should be "(for the stringent selection rules ..." We thank reviewer very much for the carefully reading and correction. We revised the manuscript according to the comments and checked the English carefully throughout the manuscript.</p> <p>The clustering process is still not described. Did the authors use hierarchical clustering, kmeans clustering What was the cutoff used to identify the clusters? Based on this and on figure Sup Fig 2 I don't believe that they achieved individual sub-cluster representation. We thank reviewer for pointing out this information which should not be omitted. We applied centroid-linkage method for the hierarchical clustering in this study and the detailed information has been added in the methods (Line 194-195). After hierarchical clustering, we selected the samples with enough DNA amount in each sub-cluster as the representative candidates to do the further analysis.</p> <p>Line 92 Supplementary Figure 2 does not provide details of the stringent selection rules. We would like to thank the reviewer for the question. The reason why we cite Supplementary Fig. 2 here is one of the selection rules is according to the results of the clustering which presented in the Supplementary Fig. 2. However, as the reviewer pointed out, this figure did not provide the details of the stringent selection rules, so we deleted the citation in line 92.</p> <p>Line 136 The authors should provide the reference to the "previous" study when they mention it. We are grateful for the reviewer's suggestion. We agreed with the reviewer and added the reference here, which is our recent study using 16S rRNA amplicon sequencing.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
Resources	Yes

<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

[Click here to view linked References](#)

The metagenome of the female upper reproductive tract

Chen Chen^{2,3,*}, Fei Li^{1,2,3*}, Weixia Wei^{4,5,*}, Zirong Wang², Juanjuan Dai^{4,5}, Lilan Hao^{2,3}, Liju Song^{2,3}, Xiaowei Zhang^{2,3}, Liping Zeng^{4,5}, Hui Du^{4,5}, Huiru Tang^{4,5}, Na Liu⁶, Huanming Yang^{2,9}, Jian Wang^{2,9}, Lise Madsen^{2,7,11}, Susanne Brix¹², Karsten Kristiansen^{2,7}, Xun Xu^{2,3}, Junhua Li^{2,3,8,13}, Ruifang Wu^{4,5†}, Huijue Jia^{2,3,8,10,†}

¹BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, China.

²BGI-Shenzhen, Shenzhen 518083, China.

³China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China.

⁴Peking University Shenzhen Hospital, Shenzhen 518036, China.

⁵Shenzhen Key Laboratory on Technology for Early Diagnosis of Major Gynecological diseases, Shenzhen, PR China

⁶BGI genomics, BGI-Shenzhen, Shenzhen 518083, China

⁷Laboratory of Genomics and Molecular Biomedicine, Department of Biology, University of Copenhagen, Universitetsparken 13, 2100 Copenhagen, Denmark.

⁸Shenzhen Key Laboratory of Human Commensal Microorganisms and Health Research, BGI-Shenzhen, Shenzhen 518083, China.

⁹James D. Watson Institute of Genome Sciences, Hangzhou310000, China.

¹⁰Macau University of Science and Technology, Taipa, Macau 999078, China.

¹¹Institute of Marine Research (IMR), Postboks 1870, Nordnes, N-5817, Bergen, Norway.

¹²Department of Biotechnology and Biomedicine, Technical University of Denmark, Soltofts Plads, 2800 Kongens. Lyngby, Denmark.

¹³School of Bioscience and Biotechnology, South China University of Technology, Guangzhou 510006, China;

*These authors contributed equally to this work.

†Correspondence should be addressed to H. J. (jiahuijue@genomics.cn, ORCID: 0000-0002-3592-126X) or R. W. (wurf100@126.com).

1 31 **Abstract**

2
3 32 *Background*

4
5
6 33 The human uterus is traditionally believed to be sterile, while the vaginal microbiota plays
7
8
9 34 important role in fending off pathogens. Emerging evidence demonstrates the presence of
10
11
12 35 bacteria beyond the vagina. However, a microbiome-wide metagenomic analysis identifying
13
14
15 36 the overall microorganism communities has been lacking.

16
17 37 *Results*

18
19
20 38 We performed shotgun-sequencing using the Illumina platform of 52 samples from the
21
22
23 39 cervical canal and peritoneal fluid of Chinese women in reproductive age. Direct annotation
24
25
26 40 of sequencing reads identified the taxonomy of bacteria, archaea, fungi and viruses,
27
28
29 41 confirming and extending the results from our previous study. We replicated the findings in
30
31
32 42 another 24 samples from the vagina, the cervical canal, the uterus and peritoneal fluid using
33
34
35 43 BGISEQ-500 platform, revealing that microorganisms in the samples from the same
36
37
38 44 individual were largely shared in the whole reproductive tract. Over 99% human sequences
39
40
41 45 were detected in the 20GB raw data. After filtering, vaginal microorganisms were well
42
43
44 46 covered in the generated reproductive tract gene catalogue, while the more diverse upper
45
46
47 47 reproductive tract microbiota might need greater depth of sequencing and more samples to
48
49
50 48 meet the full coverage scale.

51 49 *Conclusions*

52
53 50 Microbiota in unprecedented data for uncharted body site, female upper reproductive tract,
54
55
56 51 were analyzed in this study. The community results indicated that an intra-individual
57
58
59 52 continuum of all types of microorganisms gradually changed from the vagina to the peritoneal
60
61
62
63
64
65

1 53 fluid. A framework was also established in this study aiming at understanding the
2
3 54 implications of the composition and functional potential of this distinct microbial ecosystem
4
5
6 55 in relation to health and disease.
7

8
9 **56 Keywords**

10
11 57 Metagenomics, Microbiota, Female upper reproductive tract
12
13
14
15 58

16
17 **59 Background**

18
19
20 60 As humans evolved, the female reproductive tract has formed complex and unique structures
21
22 61 such as the uterus, cervix and the vagina. The human vagina hosts trillions of bacteria that can
23
24
25 62 significantly impact the health of women and their neonates. The cervix has been regarded to
26
27
28 63 be a perfect barrier between the vagina and uterus leading to the assumption that the upper
29
30
31 64 reproductive tract functions in a sterile environment. However, judging from evidence in
32
33
34 65 insects and other animals, humans are probably no exception with regard to vertical
35
36
37 66 transmission of the mothers' microbiota before birth [1]. Thus, in humans, bacterial DNA has
38
39
40 67 been detected in the placenta [2,3]. Based on our recent analyses using 16S rRNA amplicon
41
42
43 68 sequencing, the upper reproductive tract, including cervix, uterus, fallopian tubes and
44
45
46 69 peritoneal fluid harbor diverse communities of bacteria, though at low abundance [4].

47
48 70 Recently, the studies of female reproductive tract microbiota have mainly focused on the
49
50
51 71 vagina using 16S rRNA amplicon sequencing [5–7]. Studies using 16S rRNA gene amplicon
52
53
54 72 sequencing have limitations in relation to lower taxonomic resolution and the lack of ability to
55
56
57 73 perform species-specific functional inference. Metagenomic shotgun sequencing can address
58
59
60 74 these limitations, but only a few studies have applied metagenomic shotgun sequencing on the
61
62
63
64
65

1 75 vaginal microbiota [8], and no studies have characterized the compositional range of the
2
3 76 upper reproductive tract microbiome using metagenomic analysis. The present study is the
4
5
6 77 first to provide metagenomic data from the female upper reproductive tract.
7
8
9 78

10 11 79 **Data description**

12
13 80 Samples of six locations (CL, lower third of vagina; CU, posterior fornix; CV, cervical mucus
14
15
16 81 drawn from the cervical canal; ET, endometrium; FLL and FRL, left and right fallopian tubes;
17
18
19 82 PF, peritoneal fluid from the pouch of Douglas) throughout the female reproductive tract from
20
21
22 83 137 Chinese women of reproductive age, undergoing surgery for conditions not known to
23
24
25 84 involve infection (**Supplementary Table 1**) were collected for this study. 16S rRNA gene
26
27
28 85 amplicon sequencing was performed on 665 of these samples. The results from 476 of these
29
30
31 86 have been published previously [4], and those from the remaining 189 were presented in this
32
33
34 87 study. Two samples (1 CV and 1 CU) were subjected to shotgun sequencing with or without
35
36
37 88 prior removal of human DNA using a commercial kit to test the experimental effect of host
38
39
40 89 sequencing removing (refer to **Methods** section). Then, 25 PF and 25 CV samples were
41
42
43 90 sequenced on the Illumina HiSeq platform using 100 bp paired-end (PE) sequencing (for the
44
45
46 91 stringent selection rules of samples, see **Methods** for details). For these 52 samples, 20GB of
47
48
49 92 raw data per sample, corresponding to a total of 0.99 TB were generated. Additionally,
50
51
52 93 intra-individual similarity in the vagino-uterine microbiota were also examined basing on 24
53
54
55 94 samples from different sites of the reproductive tract (CL, CU, CV, ET, PF) in 6 women.
56
57
58 95 These samples were sequenced on the BGISEQ-500 sequencer using 100 bp single-end (SE)
59
60
61 96 sequencing and generated 60GB of raw data per sample, totaling 1.40 TB. The dataset after
62
63
64
65

1 97 filtering out low-quality and host reads (refer to **Methods** section) is available via the EBI
2
3 98 database using the accession number PRJEB24147.
4
5
6 99

9 100 **Analyses and Discussion**

10 101 *Metagenomic sequencing*

11
12
13
14 102 According to shotgun-sequencing of vaginal samples by the Human Microbiome Project
15
16
17 103 (HMP) and of placental samples by Aagaard *et al.*, over 90% of the sequences were derived
18
19
20 104 from human host DNA [2,9]. To overcome this problem, we first tested a commercial kit that
21
22
23 105 removes human DNA by binding and precipitating CpG-methylated DNA. Unfortunately,
24
25
26 106 after the kit treatment, a considerable amount (99.9% for CV sample and 79% for CU sample)
27
28
29 107 of host DNA still remained (**Supplementary Fig. 1a**). Furthermore, the bacteria compositions
30
31
32 108 varied by kit treatment when comparing with the control group (**Supplementary Fig. 1b**). We
33
34
35 109 therefore abandoned the strategy of host DNA removal prior to shotgun metagenomics
36
37
38 110 sequencing.

39 111 The sample selection was founded on the data from CV and PF samples [4], which we
40
41
42 112 identified as robust representations of the overall samples. Since higher amounts of DNA is
43
44
45 113 required for shotgun-sequencing results, a more stringent rule was set as the following two
46
47
48 114 criteria: individual sub-clusters representation and sufficient DNA amount (see details in
49
50
51 115 **Methods** section). To follow the former criterion, clustering results based on the relative
52
53
54 116 abundances of OTUs in the PF and CV samples showed that the samples marked with red (all
55
56
57 117 containing DNA > 1 µg) were well distributed amongst all collected samples
58
59 118 (**Supplementary Fig. 2**), so these were selected for shotgun-sequencing in this study. As a
60
61
62
63
64
65

1 119 result, 25 PF and 25 CV samples were selected for sequencing using the Illumina HiSeq 4000
2
3 120 platform. After quality control, high-quality reads were aligned to hg 19 using SOAP and
4
5
6 121 GRCh38 using DeconSeq to remove human reads (see details in **Methods** section). The
7
8
9 122 average host contamination rate of 99.72% for CV and 99.93% for PF (**Supplementary**
10
11 123 **Table 2**), which were lower than that previously reported for placenta samples [2].

12 124 The findings further expanded by inclusion of additional 24 samples subjected to sequencing
13
14 125 on the BGISEQ-500 platform, in which we also examined the intra-individual similarity in
15
16
17 126 the vagino-uterine microbiota based on samples from different sites of the reproductive tract
18
19
20 127 (CL, CU, CV, ET, PF). The average host contamination rate for vagina (CL, CU) samples
21
22
23 128 was 96.55%, and lower than those of the CV, ET and PF samples, which all above 99.5%
24
25
26 129 (**Supplementary Table 2**).

27
28
29 130 *A diverse microbiome in the cervical canal and the peritoneal fluid of reproductive age*
30
31 131 *women*

32
33
34 132 To obtain an overview of the overall composition of the vagino-uterine microbiome, we used
35
36
37 133 Kraken to directly assign sequencing reads to all types of microbial taxa [10]. The dominant
38
39
40 134 *Lactobacillus* spp. in CV and *Pseudomonas* spp. in PF were detected in the present study and
41
42
43 135 in corresponded with the previous study [4]. In addition, the microbiome that comprise
44
45
46 136 methane-producing archaea, yeasts, herpesviruses, papillomaviruses, and bacteriophages were
47
48
49 137 also founded (**Fig. 1a, b**).

50
51
52 138 The abundance of these taxonomic units varied among samples, and those constituting more
53
54
55 139 than 0.1% of the total reads number were identified in the CV and PF samples from the same
56
57
58 140 individual (**Fig. 1c**).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 141 To gain further insight into compositional similarities of the microbiota at different sites of
2
3 142 the reproductive tract in the same individual, we selected taxa at the family level which
4
5
6 143 fulfilled two criteria: they were presented in at least two sites of the same individual and the
7
8
9 144 relative abundance was higher than 0.1%. Taxa fulfilling these criteria made up more than 45%
10
11 145 of the microorganisms presented in the samples across the 6 individuals subjected for this
12
13 146 detailed analysis (**Fig. 2**). Lactobacillaceae or Bifidobacteriaceae dominated in vagina (CL
14
15 and CU), but not in the upper reproductive tract, where microorganisms such as
16
17 147 Pseudomonadaceae, Propionibacteriaceae, Streptococcaceae and Moraxellaceae constituted a
18
19 148 notable fraction of the microbiota. In addition, eukaryotes, viruses and archaea, such as
20
21 149 Saccharomycetaceae, Herpesviridae, Ferroplasmaceae were also found in the female
22
23 150 reproductive tract. The results at the bacterial level are in keeping with our findings in a
24
25 151 recent study [4], and the current data further demonstrates an intra-individual continuum of all
26
27 152 types of microorganisms that gradually changes from the vagina to the peritoneal fluid.
28
29
30
31
32
33
34
35

36 154 *Genes from the vagino-uterine microbiota*

37
38
39 155 Reference gene catalogs have greatly facilitated analyses of the microbiome, especially the
40
41 156 human gut microbiome [11–13]. Here, we established the first gene catalog of the
42
43 157 microbiome of the female upper reproductive tract, which comprises of 60,699 genes.
44
45
46

47 158 Rarefaction analysis based on gene number revealed a curve approaching saturation with
48
49 159 about 23 CV samples (**Fig. 3**). However, rarefaction analysis based on gene numbers in PF
50
51 160 samples revealed a curve that close to saturation but still did not reach the plateau, possibly
52
53 161 due to a more diverse microbiota in PF. Therefore, with 20GB sequences per sample, vaginal
54
55 162 bacteria could be well covered, whereas characterization of bacteria from the upper
56
57
58
59
60
61
62
63
64
65

1 163 reproductive tract would require a higher amount of sequences and more samples.
2
3 164 We annotated the genes in the gene catalog according to the Kyoto Encyclopedia of Genes
4
5
6 165 and Genomes (KEGG) [14]. The matched genes in PF (15,316 genes) were all covered within
7
8
9 166 CV (39,087 genes). Comparing CV and PF in the distribution of KEGG pathways, PF showed
10
11
12 167 a greater proportion of genes in carbohydrate metabolism, replication and repair, membrane
13
14
15 168 transport and drug resistance, whereas the genes involved in translation, energy metabolism
16
17
18 169 and metabolism of cofactors and vitamins were enriched in CV (**Fig. 4**). In KO modules, CV
19
20
21 170 showed enrichment of transport systems for thiamine, cystine, teichoic acid, taurine and
22
23
24 171 putative ABC transport systems compared to PF. Regulatory systems of aerobic and
25
26
27 172 anaerobic respiration, osmotic stress response and multicellular behavior control also enriched
28
29
30 173 in CV (**Supplementary Table 3**).

31 174

33 175 **Methods**

36 176 *Sample description*

37
38
39 177 A total of 137 Chinese women of reproductive age, undergoing surgery for conditions not
40
41
42 178 known to involve infection (hysteromyoma, adenomyosis, endometriosis, and
43
44
45 179 salpingemphraxis) were enrolled in this study (**Supplementary Table 1**). Samples were taken
46
47
48 180 from the CL, CU and CV on the day of the clinical visit without any prior disturbance.
49
50
51 181 Depending on the clinical conditions, laparoscopy or laparotomy were performed, and
52
53
54 182 samples from the ET, FLL, FRL and PF were taken during surgery (**Supplementary Table 1**).
55
56
57 183 The study was approved by the institutional review boards at Peking University Shenzhen
58
59
60 184 Hospital and BGI-Shenzhen, and all women provided written informed consent. The subject
61
62
63
64
65

1 185 exclusion criteria, sampling and DNA extraction methods can be found in [4].
2
3 186 To test the effect of experimental removal of human DNA, one CU sample and one CV
4
5
6 187 sample were used to shotgun sequencing on Illumina HiSeq2000 platform with or without
7
8
9 188 prior removal of human DNA, respectively. The NEBNext Microbiome DNA Enrichment Kit
10
11
12 189 was used here according to the manufacturer's instructions with a total of 10 µg input DNA
13
14
15 190 per sample.
16
17 191 Then we made a prior selection of samples to undergo shotgun-sequencing. The selection was
18
19
20 192 founded on the data from CV and PF samples [4] based on the following two criteria: i)
21
22
23 193 samples should represent individual sub-clusters when subjected to hierarchical
24
25
26 194 (centroid-linkage) clustering based on relative abundances of operational taxonomic units
27
28
29 195 (OTUs) from 16S rRNA gene amplicon sequencing; ii) the amount of DNA should be above
30
31
32 196 1 µg. The samples with good scattering in different clusters based on the relative abundances
33
34
35 197 of OTUs in the PF and CV samples were selected for shotgun-sequencing on Illumina
36
37
38 198 HiSeq4000 platform.
39
40
41 199 We replicated the findings in another 24 samples on the BGISEQ-500 platform, where
42
43
44 200 additional sites (CL, CU, CV, ET and PF) of 6 women were moreover involved. To meet the
45
46
47 201 need of library construction, the amount of DNA in the all 24 samples were above 1 µg. And
48
49
50 202 three qualified samples for each woman were set as a threshold.

51 203 52 53 204 *Metagenomic shotgun sequencing*

54
55
56 205 Library construction and shotgun sequencing using Illumina HiSeq2000/4000 platforms
57
58
59 206 (insert size 350 bp; 100 bp of PE reads; two replicate libraries were constructed for each lane.)
60
61
62
63
64
65

1 207 and BGISEQ-500 (100 bp of SE reads; one library was constructed for each lane) were
2
3 208 performed as previously described [15] (and see protocol in protocols.io[16]). The quality
4
5
6 209 control of sequencing data from the HiSeq and BGISEQ platforms were also followed this
7
8
9 210 study. Then, human sequences were eliminated by alignment to the hg19 reference genome
10
11 211 using SOAP2.22 (SOAPaligner/soap2, RRID:SCR_005503). As the resulting data still
12
13
14 212 contained human sequences, a more stringent procedure using DeconSeq by aligning data to
15
16
17 213 the GRCh38 reference genome was applied [17].
18
19

20 214 ***Taxonomic assignment of sequencing reads***

21
22 215 High-quality, non-human sequences were tentatively assigned to microbial taxa using Kraken
23
24
25 216 with default parameters (Kraken, RRID:SCR_005484)[10]. For pair-end reads Kraken
26
27
28 217 concatenated the pairs together with a single N between the sequences automatically
29
30
31 218 with default parameters and the manual clarified that this software raised the sensitivity
32
33
34 219 by about 3 percentage points over classifying the sequences as single-end reads.
35

36 220 ***Construction of a gene catalog***

37
38
39 221 The high-quality, non-human sequencing reads of 52 samples sequenced by Illumina HiSeq
40
41
42 222 platforms were *de novo* assembled into contigs using IDBA-UD (IDBA-UD
43
44
45 223 (RRID:SCR_011912))[18]. We used the same strategy as previous study [12,13], where genes
46
47
48 224 were predicted from the contigs by MetaGeneMark [19], and highly similar genes (95%
49
50
51 225 identity, 90% overlap) were removed as redundancy using CD-HIT (CD-HIT,
52
53
54 226 RRID:SCR_007105) [20]. Functional annotations were made by BLASTP (v2.2.24) based on
55
56
57 227 KEGG (v76) databases (KEGG , RRID:SCR_012773)[14].
58
59

60
61
62
63
64
65 228

1 229 ***Availability of Supporting Data***

2
3 230 The sequencing data after filtering out low-quality and host reads is available via the EBI
4
5
6 231 database using the accession number PRJEB24147. Additional supporting data is available
7
8
9 232 via the *GigaScience* GigaDB database [21].

10
11
12 233 ***Abbreviations***

13
14
15
16 234 bp: base pair; GB: Gigabase; HMP: Human Microbiome Project; KEGG: Kyoto Encyclopedia
17
18 235 of Genes and Genomes; OTU: operational taxonomic units; PE: paired-end; SE: single-end.

19
20
21
22 236 ***Author's contributions***

23
24
25 237 H.J. and R.W. conceived and directed the project. W.W., J.D., L.Z., H.D., H.T., and R.W.
26
27 238 performed the clinical diagnosis, sample collection. C.C., Z.W., F.L., and L.H. performed the
28
29 239 bioinformatic analyses and prepared display items. C.C., F.L., Z.W., X.Z., J.L. and H.J. wrote
30
31
32 240 the first version of the manuscript. L.M., S.B. and K.K. revised the manuscript. All authors
33
34
35 241 contributed to the final revision of the manuscript.

36
37
38
39 242
40
41 243 ***Acknowledgements***

42
43
44 244 The study was supported by the Shenzhen Municipal Government of China
45
46 245 (JCYJ20160229172757249, JCYJ20150601090833370), the grant from the Macau
47
48
49 246 Technology Development Fund (102/2016/A3). We gratefully acknowledge colleagues at
50
51
52 247 BGI-Shenzhen for DNA quality control, library construction, sequencing, and helpful
53
54
55 248 discussions.

56
57
58 249
59
60
61
62
63
64
65

1 250 **Competing financial interests**

2
3 251 The authors declare no competing financial interests.

4
5
6 252 **References**

- 7
8
9 253 1. Funkhouser LJ, Bordenstein SR. Mom Knows Best: The Universality of Maternal
10 254 Microbial Transmission. *PLoS Biol.* 2013;11:e1001631.
- 11
12 255 2. Aagaard K, Ma J, Antony KM, Ganu R, Petrosino J, Versalovic J. The placenta harbors a
13 256 unique microbiome. *Sci. Transl. Med.* 2014;6:237ra65.
- 14
15
16 257 3. Antony KM, Ma J, Mitchell KB, Racusin DA, Versalovic J, Aagaard K. The Preterm
17 258 Placental Microbiome Varies in Association with Excess Maternal Gestational Weight Gain.
18 259 *Am J Obs. Gynecol.* Elsevier; 2014;212:653.e1-653.e16.
- 19
20
21 260 4. Chen C, Song X, Wei W, Zhong H, Dai J, Lan Z, et al. The microbiota continuum along
22 261 the female reproductive tract and its relation to uterine-related diseases. *Nat Commun.* 2017
23 262 Oct 17;8(1):875. doi: 10.1038/s41467-017-00901-0.
- 24
25
26 263 5. Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SSK, Mcculle SL, et al. Vaginal
27 264 microbiome of reproductive-age women. *Proc. Natl. Acad. Sci.* 2010;108:4680–7.
- 28
29 265 6. Gajer P, Brotman RM, Bai G, Sakamoto J, Schutte UME, Zhong X, et al. Temporal
30 266 Dynamics of the Human Vaginal Microbiota. *Sci. Transl. Med.* 2012;4:132ra52-132ra52.
- 31
32 267 7. Ding T, Schloss PD. Dynamics and associations of microbial community types across the
33 268 human body. *Nature.* 2014 May 15;509(7500):357-60. doi: 10.1038/nature13178.
- 34
35
36 269 8. Lloyd-price J, Mahurkar A, Rahnvard G, Crabtree J, Orvis J, Hall AB, et al. Strains ,
37 270 functions and dynamics in the expanded Human Microbiome Project. *Nature.* 2017 Oct
38 271 5;550(7674):61-66. doi: 10.1038/nature23889.
- 39
40 272 9. Methé BA, Nelson KE, Pop M, Creasy HH, Giglio MG, Huttenhower C, et al. A
41 273 framework for human microbiome research. *Nature.* 2012;486:215–21.
- 42
43
44 274 10. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using
45 275 exact alignments. *Genome Biol.* 2014;15:R46.
- 46
47 276 11. Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. *Nat.*
48 277 *Rev. Microbiol.* 2016;14:508–22.
- 49
50
51 278 12. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of
52 279 reference genes in the human gut microbiome. *Nat. Biotechnol.* 2014;32:834–41.
- 53
54 280 13. Xie H, Guo R, Zhong H, Feng Q, Lan Z, Qin B, et al. Shotgun Metagenomics of 250
55 281 Adult Twins Reveals Genetic and Environmental Impacts on the Gut Microbiome. *Cell Syst.*
56 282 2016;
- 57
58
59 283 14. Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. KEGG as a reference
60 284 resource for gene and protein annotation. 2016;44:457–62.
- 61
62
63
64
65

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 285 15. Fang C, Zhong H, Lin Y, Chen B, Han M, Ren H, et al. Assessment of the cPAS-based
286 BGISEQ-500 platform for metagenomic sequencing. *Gigascience*. 2018 Mar 1;7(3):1-8. doi:
287 10.1093/gigascience/gix133.
- 288 16. Jie Huang, Xinming Liang, Yuankai Xuan, Chunyu Geng, Yuxiang Li, Haorong Lu,
289 Shoufang Qu, Xianglin Mei, Hongbo Chen, Ting Yu, Nan Sun, Junhua Rao, Jiahao Wang,
290 Wenwei Zhang, Ying Chen, Sha Liao, Hui Jiang, Xin Liu, Zhaopeng Yang, Feng Mu,
291 Shangxian Gao (2018). BGISEQ-500 WGS library construction. *protocols.io*
292 [dx.doi.org/10.17504/protocols.io.ps5dng6](https://doi.org/10.17504/protocols.io.ps5dng6)
- 293 17. Schmieder R, Edwards R. Fast identification and removal of sequence contamination from
294 genomic and metagenomic datasets. *PLoS One*. 2011;6.
- 295 18. Peng Y, Leung HCM, Yiu SM, Chin FYL. IDBA-UD: a de novo assembler for single-cell
296 and metagenomic sequencing data with highly uneven depth. *Bioinformatics*. 2012;28:1420–
297 8.
- 298 19. Hyatt D, Locascio PF, Hauser LJ, Uberbacher EC. Gene and translation initiation site
299 prediction in metagenomic sequences. *Bioinformatics*. 2012;28:2223–30.
- 300 20. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein
301 or nucleotide sequences. *Bioinformatics*. 2006;22:1658–9.
- 302 21. Chen C; Li F; Wei W; Wang Z; Dai J; Hao L; Song L; Zhang X; Zeng L; Du H; Tang H; Liu
303 N; Yang H; Wang J; Madsen L; Brix S; Kristiansen K; Xu X; Li J; Wu R; Jia H (2018):
304 Supporting data for "The metagenome of the female upper reproductive tract" *GigaScience*
305 Database. <http://dx.doi.org/10.5524/100491>

1 307 **Figure legends**

2
3
4 308 **Figure 1: The overall microbiome composition of the cervical canal and the peritoneal**
5
6 309 **fluid of reproductive age women.** Cumulative bar charts of the main taxa at domain (a) and
7
8
9 310 family (b) levels in CV and PF samples. (c) Compositional overlap at family level of CV and
10
11
12 311 PF samples from the same individuals. Relative number of reads was calculated as $N_p =$
13
14 312 $\frac{a_p}{a_t} \times m$, where a_p is the number of reads within p taxa in a sample. a_t is the total number
15
16
17 313 of reads within a sample, m is median number of reads within all 50 samples. When p taxa
18
19
20 314 is shared by CV and PF samples from the same individuals, and at the same time both
21
22
23 315 N_p values are higher than $0.1\% \times m$, the p taxa is included in the cumulative bar charts.
24
25
26 316 Taxa names (b, c) in black, purple, and blue denote bacteria, eukaryote and viruses,
27
28
29 317 respectively.

30
31
32 318 **Figure 2: Composition of the vagino-uterine microbiota.** (a, c, e, g, i, k) Venn diagram
33
34
35 319 depicting shared taxa at the family levels in samples collected at different sites in the same
36
37
38 320 individual. (b, d, f, h, k, l) Cumulative bar charts of the taxa with relative abundance higher
39
40
41 321 than 0.1% and present in at least two sites of the same individual. Taxa names (b, c) in black,
42
43
44 322 purple, blue, and grey denote bacteria, eukaryote, viruses and archaea, respectively.

45
46
47 323 **Figure 3: Rarefaction of microbial gene content in CV (a) and PF (b) samples.** The
48
49
50 324 number of genes in each group was calculated after 100 random samplings with replacement.
51
52
53 325 Boxes denote the interquartile range (IQR) between the first and third quartiles (25th and 75th
54
55
56 326 percentiles, respectively) and the line inside denotes the median. Whiskers denote the lowest
57
58
59 327 and highest values within 1.5 times IQR from the first and third quartiles, respectively.
60
61
62
63
64
65

1 328 Circles denote outliers beyond the whiskers.

2
3
4 329 **Figure 4: KEGG pathway classification of the vagino-uterine microbiome.** Comparison
5
6
7 330 of CV (red) and PF (blue) data based on KEGG annotation, which emphasizes functional
8
9
10 331 similarity of the CV and PF microbiota.

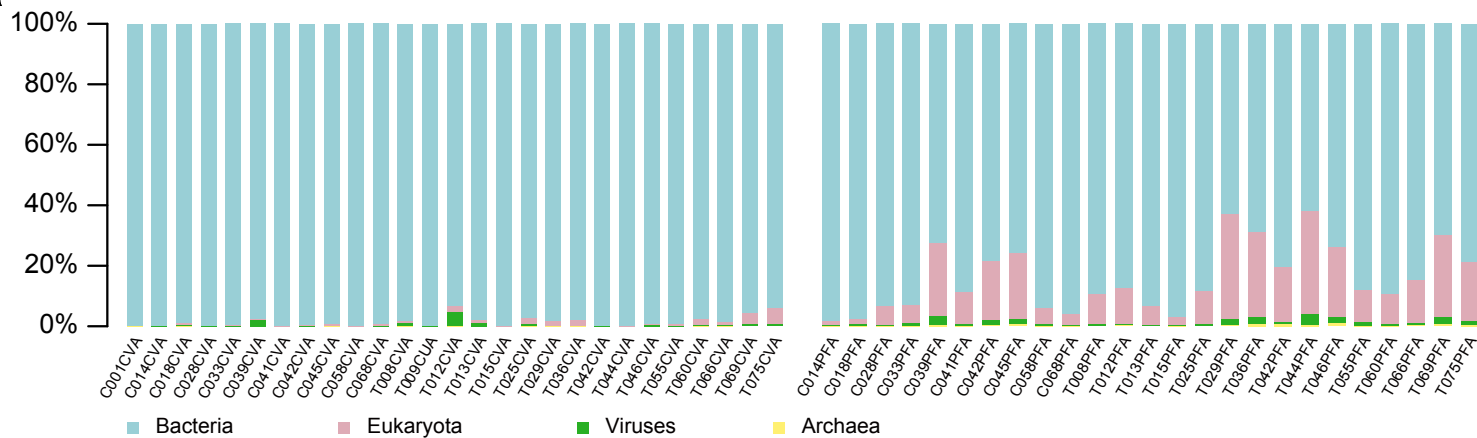
11
12
13
14 332

15
16
17 333 **Supplementary Figure legends**

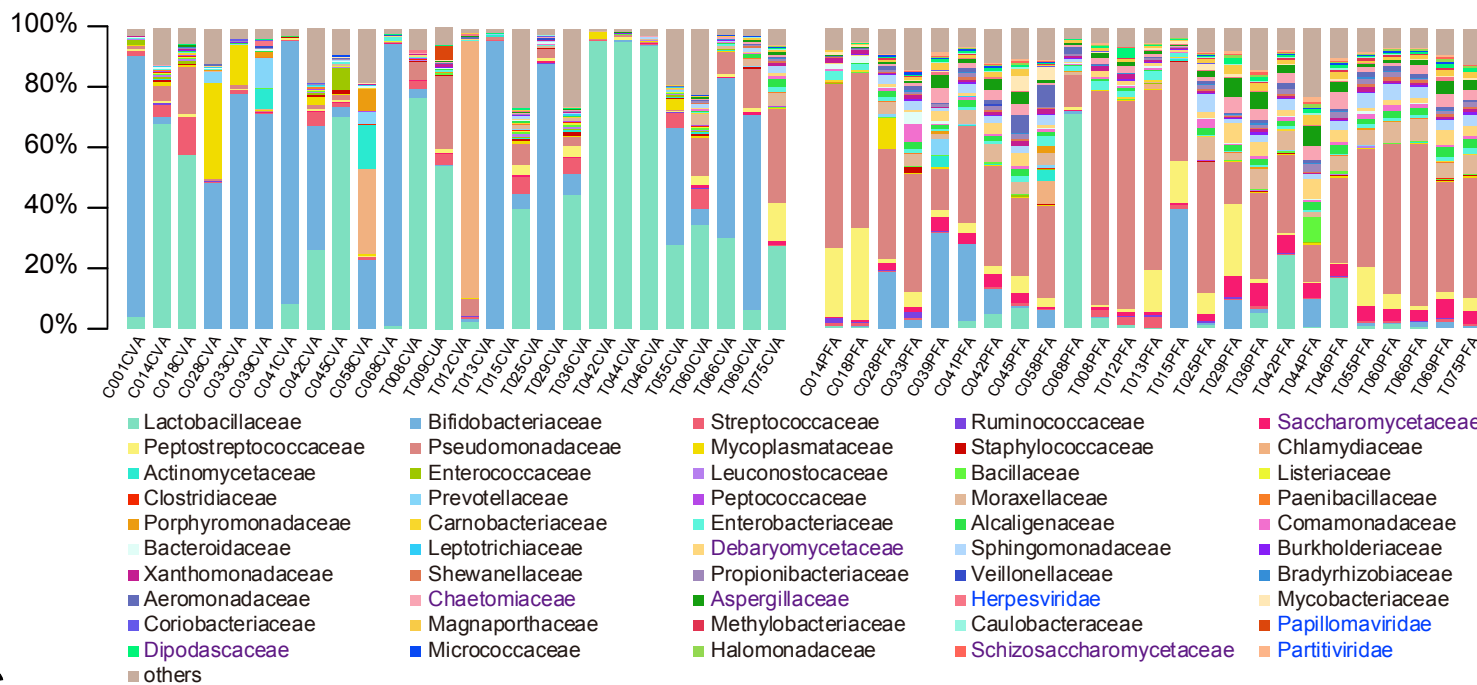
18
19
20
21 334 **Supplementary Figure 1: Evaluation of the NEBNext Microbiome DNA Enrichment Kit**
22
23
24 335 **by two comparative strategies.** Sample names suffixed by “-HR” represent DNA samples
25
26
27 336 that were treated with the kit for removal of host DNA before shotgun sequencing, while
28
29
30 337 sample names suffixed by A represent DNA samples that were subjected to shotgun
31
32
33 338 sequencing directly (a). The table data shows the obtained read number, and remaining reads
34
35
36 339 after removal of host DNA reads in the two samples. b) Influence of host DNA presence on
37
38
39 340 bacterial DNA identification during shotgun sequencing. The plots display the compositional
40
41
42 341 difference amongst major bacteria genera in samples with and without (-HR) host DNA
43
44
45 342 presence. Data were analyzed by mapping reads to the ICG bacterial reference gene catalog
46
47 343 [12].

48
49
50 344 **Supplementary Figure 2: Samples selected for metagenomic sequencing.** Hierarchical
51
52
53 345 clustering of CV (a) and PF (b) samples based on the relative abundances of OTUs. Samples
54
55
56 346 which represent individual sub-clusters and hold DNA amounts above 1 µg were selected for
57
58
59 347 shotgun-sequencing (red).

a



b



c

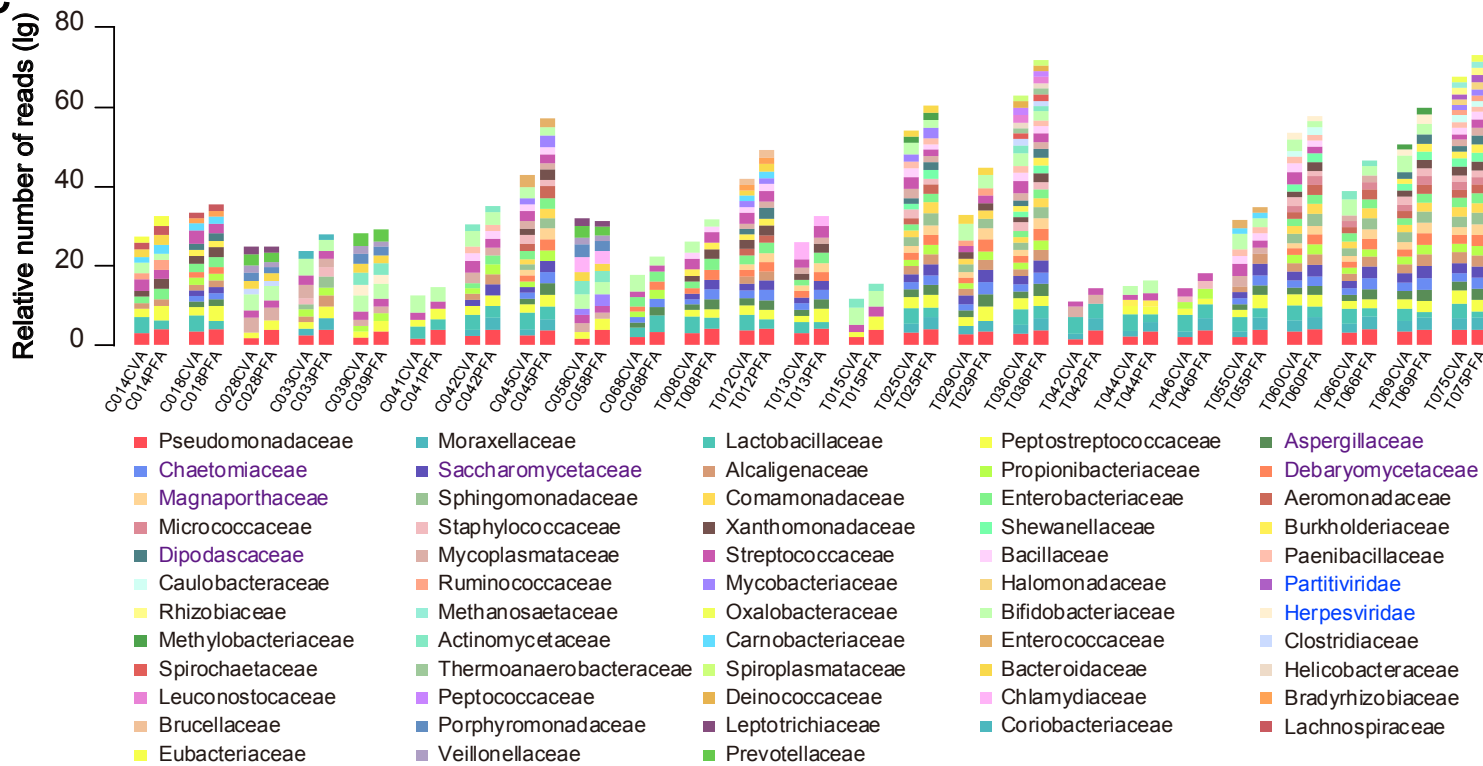
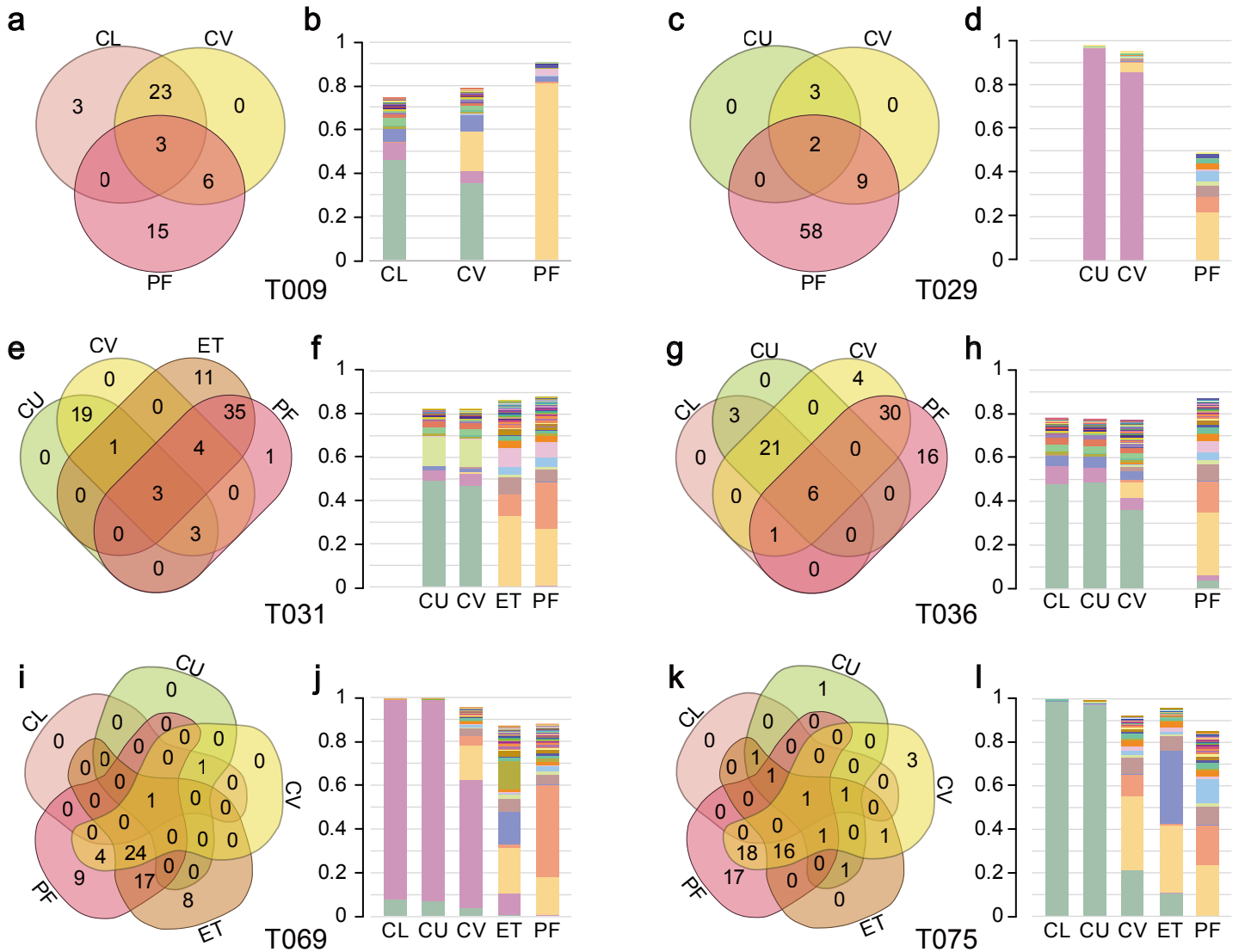


Figure 2

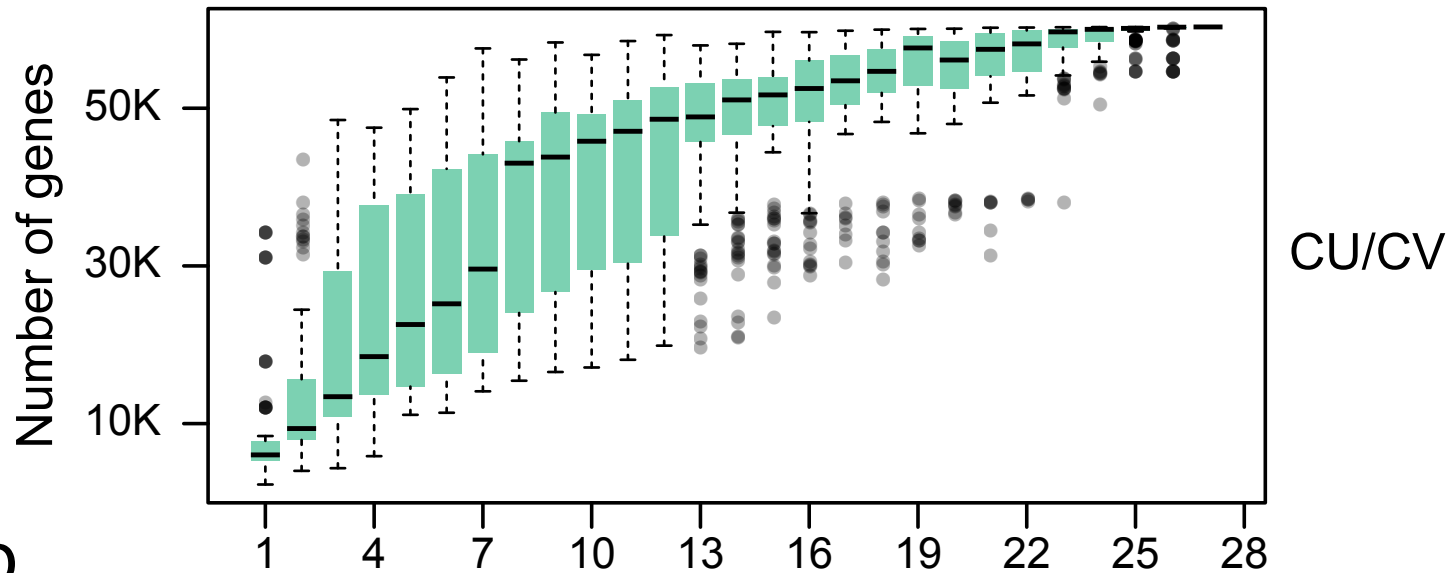
[Click here to access/download;Figure;FIG2-20180126.pdf](#)



- Lactobacillaceae
- Moraxellaceae
- Staphylococcaceae
- Cordycipitaceae
- Ajellomycetaceae
- Chaetomiaceae
- Aspergillaceae
- Herpotrichiellaceae
- Clostridiaceae
- Arthrodermataceae
- Carnobacteriaceae
- Aureobasidiaceae
- Erysipelotrichaceae
- Bradyrhizobiaceae
- Herpesviridae
- Leptotrichiaceae
- Bifidobacteriaceae
- Saccharomycetaceae
- Alcaligenaceae
- Peptoniphilaceae
- Glomerellaceae
- Bacillaceae
- Dermateaceae
- Pleosporaceae
- Neisseriaceae
- Corynebacteriaceae
- Plectosphaerellaceae
- Atopobiaceae
- Dictyoglomaceae
- Brucellaceae
- Leptosphaeriaceae
- Tremellaceae
- Pseudomonadaceae
- Comamonadaceae
- Ruminococcaceae
- Mycobacteriaceae
- Debariomycetaceae
- Burkholderiaceae
- Methylobacteriaceae
- Listeriaceae
- Paenibacillaceae
- Magnaporthaceae
- Campylobacteraceae
- Helicobacteraceae
- Ferroplassmaceae
- Spiroplasmataceae
- Pseudeurotiaceae
- Pucciniaceae
- Propionibacteriaceae
- Enterobacteriaceae
- Peptostreptococcaceae
- Micrococcaceae
- Enterococcaceae
- Hypocreaceae
- Aeromonadaceae
- Peptococcaceae
- Deinococcaceae
- Rhizobiaceae
- Spirochaetaceae
- Rhodobacteraceae
- Mycosphaerellaceae
- Planctomycetaceae
- Thermoanaerobacteriales_Family_III_Incertae_Sedis
- Streptococcaceae
- Sphingomonadaceae
- Xanthomonadaceae
- Sordariaceae
- Leuconostocaceae
- Clavicipitaceae
- Mycoplasmataceae
- Actinomycetaceae
- Shewanellaceae
- Ustilaginaceae
- Orbiliaceae
- Sclerotiniaceae
- Rhodospirillaceae
- Nostocaceae
- Agaricaceae

a

Rarefaction of genes



b

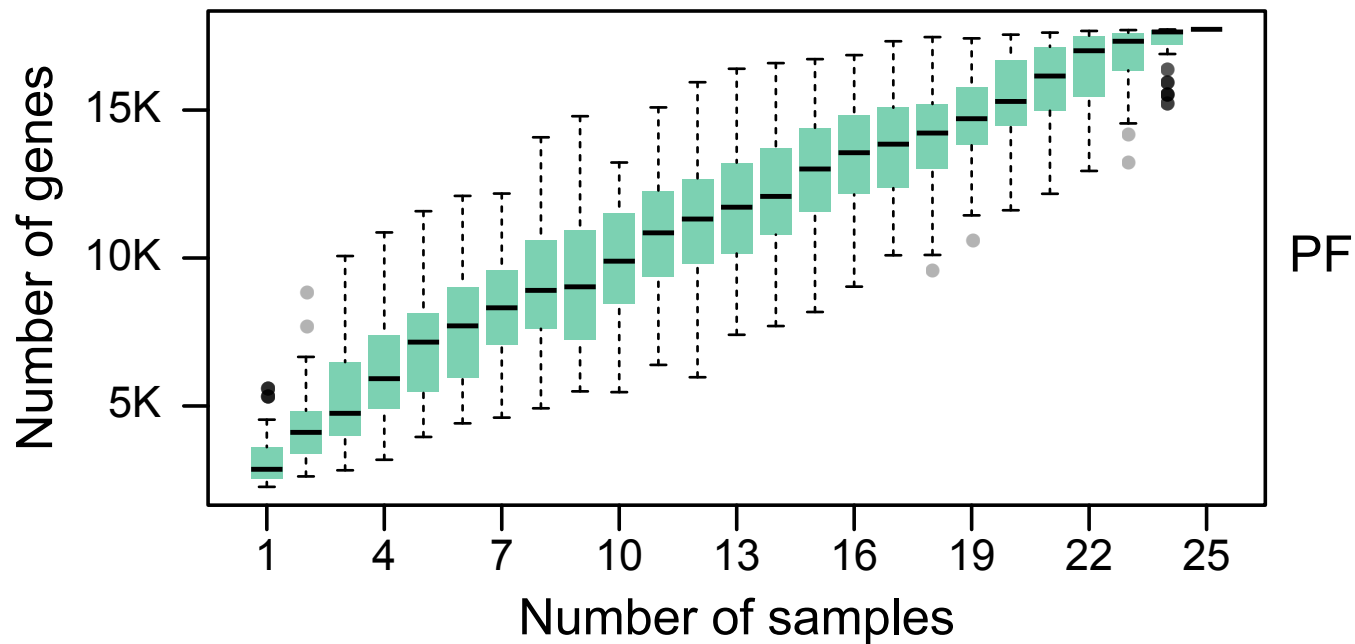
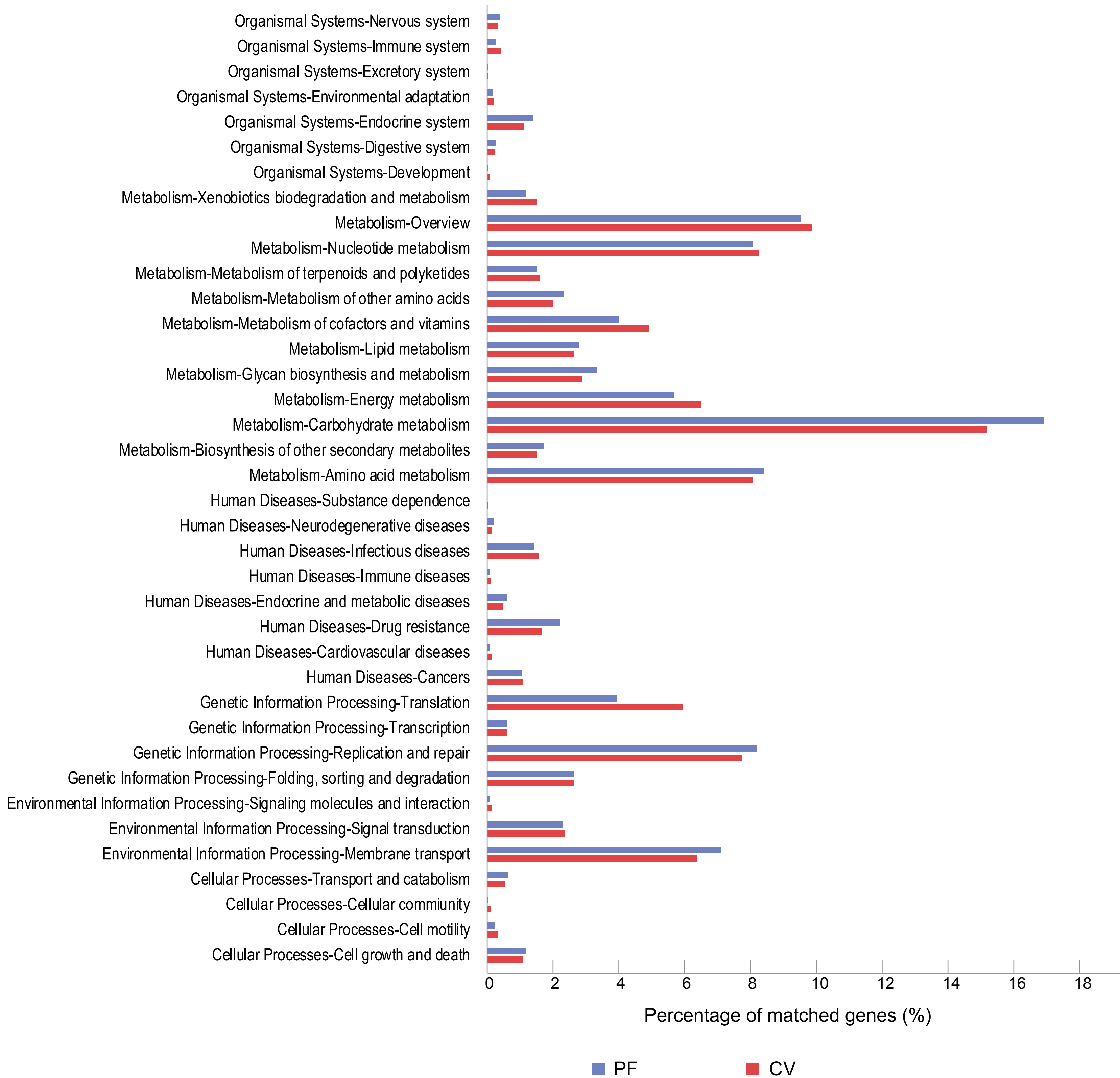


Figure 4

KEGG pathway classification





Click here to access/download
Supplementary Material
SFIG1-20180126.pdf





Click here to access/download
Supplementary Material
SFIG2-20180126.pdf





[Click here to access/download](#)

Supplementary Material

Supplementary Table 1-20180126.xlsx





[Click here to access/download](#)

Supplementary Material

Supplementary Table 2-20180126.xlsx





Click here to access/download
Supplementary Material
Supplementary Table 3.xlsx

