

1 **Evaluating bacterial and functional diversity of human gut microbiota by complementary**  
2 **metagenomics and metatranscriptomics**

3

4 Ravi Ranjan<sup>1# $\$$ a</sup>, Asha Rani<sup>1# $\$$</sup> , Patricia W. Finn<sup>1@</sup> and David L. Perkins<sup>1,3 $\$$ @</sup>

5

6 <sup>1</sup>Department of Medicine, <sup>2</sup>Department of Bioengineering, <sup>3</sup>Department of Surgery, University of Illinois,  
7 Chicago, IL 60612 USA

8

9  **$\$$ Correspondence:**

10 David Perkins, MD, PhD

11 Email: [perkinsd@uic.edu](mailto:perkinsd@uic.edu), Phone: 312-413-3382, Fax: 312-355-0499

12

13 Ravi Ranjan, PhD

14 Email: [ranjan@uic.edu](mailto:ranjan@uic.edu)

15

16 Asha Rani, PhD

17 Email: [asharani@uic.edu](mailto:asharani@uic.edu)

18

19 Department of Medicine

20 University of Illinois at Chicago

21 Chicago IL 60612 USA

22

23 # These authors contributed equally and considered as co-first authors.

24 @These authors contributed equally and considered as co last authors.

25

26

27

28

29

30

31

32

33

34

35

36

37

## 1 **ABSTRACT**

2 It is well accepted that dysbiosis of microbiota is associated with disease; however, the biological  
3 mechanisms that promote susceptibility or resilience to disease remain elusive. One of the major limitations  
4 of previous microbiome studies has been the lack of complementary metatranscriptomic (functional) data  
5 to complement the interpretation of metagenomics (bacterial abundance). The purpose of the study was  
6 twofold, first to evaluate the bacterial diversity and differential gene expression of gut microbiota using  
7 complementary shotgun metagenomics (MG) and metatranscriptomics (MT) from same fecal sample.  
8 Second, to compare sequence data using different Illumina platforms and with different sequencing  
9 parameters as new sequencers are introduced and determine if the data are comparable on different  
10 platforms. In this study, we perform ultra-deep metatranscriptomic shotgun sequencing for a sample that  
11 we previously analyzed with metagenomics shotgun sequencing. We validated the sequencing and  
12 analysis methods using different Illumina platform, and with different sequencing and analysis parameters.  
13 Our results suggest that use of different Illumina platform did not lead to detectable bias in the sequencing  
14 data. The analysis of the sample using MG and MT approach shows that some species genes are more  
15 highly represented in the MT than in the MG, indicating that some species are highly metabolically active.  
16 Our analysis also shows that ~52% of the genes in the metagenome are in the metatranscriptome, and  
17 therefore are robustly expressed. The functions of the low and rare abundance bacterial species remain  
18 poorly understood. Our observations indicate that among the low abundant species analyzed in this study  
19 some were found to be more metabolically active compared to others and can contribute distinct profiles of  
20 biological functions that may modulate the host-microbiota and bacteria-bacteria interactions.

21

## 22 **KEYWORDS**

23 Metatranscriptomics; Metagenomics; Microbiome; Microbiota; Next-generation sequencing; RNAseq;  
24 Shotgun sequencing; 16S rRNA; Targeted amplicon sequencing

25

## 26 **INTRODUCTION**

27 The human microbiota represents a complex community of numerous and diverse microbes that is linked  
28 with our development, metabolism, physiology, health, and is considered functionally comparable to an  
29 organ of the human body (Cho & Blaser 2012, Human Microbiome Project 2012). Previous studies have  
30 established that a healthy human microbiota is associated with maintaining health, whereas dysbiosis has  
31 been associated with various pathologies and diseases such as obesity, inflammatory bowel disease,  
32 pulmonary diseases, urinary tract infection etc., (Iebba et al 2016, Pflughoeft & Versalovic 2012).  
33 Traditionally, identifying microbes relied on culture based techniques, however the majority (>90 – 95 %)   
34 of microbial species cannot be readily cultured using current laboratory techniques (Sharma et al 2005) .  
35 Advancements in culture- and cloning-independent molecular methods, coupled with high-throughput next-  
36 generation DNA sequencing technologies have rapidly advanced our understanding of the microbiota.  
37 Additionally, with the rate of recent technological advancements, the DNA sequencing ventures have been

1 introducing new DNA sequencers with versatile sequencing parameters. This has also complicated the  
2 comparison of data within and among the samples. Thus, there is a need to compare the sequencing data  
3 from same samples using different platforms. Many previous studies employed targeted amplicon  
4 sequencing of the conserved prokaryotic 16S ribosomal RNA (16S rRNA) gene (Human Microbiome Project  
5 2012, Huse et al 2012, Stulberg et al 2016). This method identifies operational taxonomic units (OTUs) and  
6 are correlated with bacterial taxa; however, assignment of taxa defined by OTUs is commonly limited to the  
7 genus level due to low accuracy at the species level. In contrast, metagenomics shotgun sequencing  
8 (MGS), which is employed in our study, can determine taxonomic annotations at the species level.

9         Although the association of multiple diseases with dysbiosis of the microbiome has been  
10 established, the elucidation of the underlying biologic mechanisms that promote pathological phenotypes  
11 has been elusive in most cases. A major limitation of both targeted amplicon and metagenome shotgun  
12 sequencing is that bacterial functions are predicted based on the genome sequence of the associated taxa.  
13 However, it is well established that there is differential bacterial gene expression at the transcriptional level  
14 in response to environmental and dietary exposures. For example, it has been reported that there is a set  
15 of constitutively expressed core genes that mediate core microbial functions as well as a highly regulated  
16 subset of genes that respond to unique environmental influences (Booijink et al 2010, Ursell & Knight 2013).  
17 In addition, some bacteria may exist in an inert state or spore form and thus not contribute to the biological  
18 response (Franzosa et al 2014). Thus, an analysis of bacterial gene expression with metatranscriptomics  
19 approach could provide additional insight into the biological functions of specific microbiomes.

20         The gut microbiota is composed of highly abundant few species and less abundant many rare  
21 bacterial species, thus to understand the complex functions of the microbiota it is essential to understand  
22 the functions of both the high- and low-abundant bacterial species. Analyses of MG and MT data are often  
23 challenged by the sequencing depth, parameters, and sequencing platforms, which limits the power of  
24 functional classification and abundance estimation, this in turn hampers the downstream data analyses of  
25 differentially expressed genes. The unique feature of our study is that we are comparing the sequencing  
26 reads at different depths, platform, read length, read and contig based comparison for MG and MT for the  
27 same sample. To develop a comprehensive understanding of the ecological functions of a microbiome, it  
28 is essential to determine not only the metatranscriptome, but also to ascertain the functional contributions  
29 of both the abundant and the rare species in a microbiome. To investigate these questions, we analyzed  
30 both the metagenome and the metatranscriptome using shotgun sequencing which can determine the  
31 abundance of gene transcripts relative to the abundance of the genome. This allowed us to identify both  
32 over- and under-expressed transcripts. In this study, we identified biological functions in both rare and  
33 abundant bacterial species using metagenomic and metatranscriptomic methods optimized and validated  
34 in our laboratory.

35

36

37

## 1 MATERIAL AND METHODS

2 **Subject recruitment and sample collection:** The study was approved by the Institutional Review Board  
3 of the University of Illinois at Chicago, and the experimental methods were performed in accordance with  
4 the approved guidelines. A 33 year, male subject without known medical conditions provided the signed  
5 informed consent and self-collected stool in a EasySampler Stool Collection kit (Alpco Diagnostics). The  
6 fecal sample was immediately aliquoted into sterile 1.5 ml Eppendorf safe-lock tubes and stored at -80°C  
7 till further DNA and RNA isolation was carried out.

8  
9 **RNA isolation from fecal sample and mRNA enrichment:** The objective of the study was to perform  
10 matched metagenome and metatranscriptome studies of the same fecal sample. We investigated the same  
11 fecal sample we previously analyzed by metagenomics sequencing. Total RNA was isolated using the  
12 PowerMicrobiome RNA Isolation Kit (Catalog # 26000-50, MO BIO Laboratories, Inc) from a fecal sample.  
13 For efficient lysis of the microbes in the sample, 200 µL of Phenol/Chloroform/Isoamyl alcohol (25:24:1)  
14 (Catalogue #327115000, Acros Organics) was added to the reagents provided with the kit. The contents  
15 were vortexed for 1-2 min with a table top vortexer and homogenized twice at speed 10 for 5 min with air-  
16 cooling using the Bullet Blender Storm Homogenizer (Catalogue # BBY24M, Next Advance Inc). Total RNA  
17 was isolated with the manufacturer's recommended procedure including the on-column DNase treatment  
18 (to remove the potentially co-isolated DNA). The RNA was eluted with 1×TE, pH 8.0, and stored at -80°C.  
19 The quality and quantity of the RNA was assessed using a spectrophotometer (NanoPhotometer Pearl,  
20 Denville Scientific, Inc), agarose gel electrophoresis, fluorometer (Qubit® RNA Broad Range assay, Life  
21 Technologies Corporation), and Agilent RNA 6000 Nano Kit on 2100 Bioanalyzer instrument (Agilent  
22 Technologies, Inc.). Total RNA was enriched for mRNA by subtractive hybridization using the  
23 MICROBExpress™ Bacterial mRNA Enrichment Kit following manufacturers recommended protocol  
24 (Ambion, Life Technologies). The mRNA enrichment and rRNA depletion was analyzed using an Agilent  
25 RNA 6000 Nano Kit on 2100 Bioanalyzer instrument (Agilent Technologies, Inc.).

26  
27 **Fecal metatranscriptome library preparation and shotgun sequencing:** The enriched mRNA was  
28 mechanically fragmented to a size range of ~200 bp with an ultrasonicator using the adaptive focused  
29 acoustics with the following manufacturer recommended protocols (Covaris S220 instrument, Covaris Inc).  
30 The fragmentation of mRNA was assessed using Agilent RNA 6000 Pico Kit on 2100 Bioanalyzer  
31 instrument (Agilent Technologies, Inc). The metatranscriptome libraries were prepared using NEBNext  
32 Ultra RNA Library Prep Kit for Illumina (New England BioLabs Inc). The quality and quantity of all the final  
33 libraries were analyzed with an Agilent DNA 1000 Kit on the 2100 Bioanalyzer Instrument and Qubit. The  
34 final libraries were quantitated and validated by qPCR assay using the PerfeCTa NGS Library Quantification  
35 Kit for Illumina (Quanta Biosciences, Inc.) using the CFX Connect Real-Time PCR Detection System (Bio-  
36 Rad Laboratories, Inc). Sequencing of one of the MT library was performed on a Illumina HiSeq 2000 using  
37 the TruSeq SBS v3 reagent for paired-end 100 read length (BGI Americas) (labeled as HS100), and on

1 Illumina MiSeq using v3-600 cycle kit for paired-end 301 bases (labeled as MS301). Another set of twelve  
2 libraries was sequenced on Illumina MiSeq using 151 paired end chemistry (labeled as MS151).  
3 Manufacturer's recommended protocol was used for performing the sequencing reaction on both the HiSeq  
4 and MiSeq platforms.

5

6 **Data analysis.** The individual twelve libraries were analyzed for taxonomic and functional annotation, also  
7 all of the 12 sequence files were combined in silico and were labeled as (MS151)-Lib-All. The sequence  
8 files (HS100, 12(x) MS151, and MS301) were combined in silico and labelled as HS100+MS151+MS301.:  
9 The sequence reads were processed and analyzed using the CLC Genomics workbench version 7.5  
10 (Qiagen, Aarhus, Denmark). Raw reads were trimmed to a minimum Phred quality score of 20. Raw reads  
11 were filtered by mapping against human reference genome to remove human sequences. The non-human  
12 reads were de novo assembled using the CLC assembler using a word size (k-mer) of 50, minimum contig  
13 length 200bp, to construct the de bruijn graphs. De novo assembly was used to map reads back to the  
14 contigs (mismatch cost 2, insertion cost 3, deletion cost 2, length fraction 0.8, similarity fraction 0.8).  
15 Taxonomic and functional annotations of the reads and contigs were obtained using the automated  
16 annotation pipeline at MG-RAST web server using the default parameters (best hit classification, maximum  
17 *e*-value 1e-5 cutoff, and minimum 60% identity cutoff) using M5NR and KEGG databases (Meyer et al 2008,  
18 Mitra et al 2011). The limma analysis was used to identify species and KEGG functional pathways that were  
19 differentially abundant between metagenome (MG) and metatranscriptome (MT) (Praveen et al 2015).  
20 Limma uses an empirical Bayes method to test the differential expression based on the fitting of each  
21 species/gene to a linear model (Smyth 2004). This provides the rich features for complex experimental  
22 designs and overcomes the small sample size problem, in addition to providing enhanced biological  
23 interpretation for co-regulated sets of genes (Ritchie et al 2015). A *p* value cutoff of 0.05 after multiple  
24 testing correction based on Benjamini-Hochberg method (Benjamini & Hochberg 1995), and a log<sub>2</sub> fold  
25 change ≥1 were used to select the differentially abundant species and pathways. The data files were  
26 visualized in MeV v 4.9.0 (TM4, Boston, MA, USA) (Saeed et al 2003). The metatranscriptome data was  
27 used to compare with the previously reported metagenome data of the same sample from our group  
28 (Ranjan et al 2016).

29

## 30 **RESULTS**

### 31 **Ultra-deep metatranscriptomic shotgun sequencing (MTS)**

32 In our previous study of ultra-deep metagenome shotgun sequencing (MGS) we demonstrated effective  
33 identification of abundant species (defined as >1% relative abundance) with as few as 500 reads; however,  
34 the detection of low abundance or rare species required high numbers of sequence reads. For example,  
35 with a total of 163.7 million sequence reads generated by metagenome shotgun sequencing (MGS), the  
36 rarefaction curve did not show saturation for the identification of additional species (Ranjan et al 2016) .  
37 Based on these data, in the current study of the metatranscriptome we performed ultra-deep MTS

1 sequencing. We performed optimization and validation of our sequencing protocol using multiple  
2 sequencing platforms and analytic strategies (Fig. 1). High quality total RNA was isolated (Supplementary  
3 Fig. 1A), and the bacterial mRNA was enriched from the total RNA using subtractive hybridization, which  
4 depleted most of the rRNA (Supplementary Fig. 1B). The enriched mRNA was mechanically fragmented  
5 and libraries were constructed (Supplementary Figs. 1C and 1D). To evaluate technical reproducibility, we  
6 constructed 12 unique indexed metatranscriptome libraries from a single fecal sample. High quality libraries  
7 were prepared for sequencing on Illumina's MiSeq and HiSeq 2000 platforms (Supplementary Fig. 1E). We  
8 obtained from 3.6 to 5.4 million high quality sequence reads for the 12 replicate libraries sequenced on  
9 MiSeq for 151 PE and 32.7 to 56.5 million reads on a HiSeq 2000 platform using 100 and 151 PE  
10 sequencing parameters. In total, we obtained a total of 139.6 million sequence reads by combining the  
11 HiSeq and MiSeq sequence data in silico (HS100+MS151+MS301) (Table 1).

12

### 13 **Comparison of analytic strategies**

14 In our previous analysis of ultra-deep MGS data, we observed a substantial increase in the average length  
15 of the assembled contigs (904 bp) compared with the average read length 170 bp., and the average N50  
16 length of the contigs was 6,262 bp (Ranjan et al 2016). Therefore, we compared the effect of analyzing the  
17 reads versus assembled contigs in the metatranscriptome (MT) data. In the MT data, the average contig  
18 length was 268 bp which was modestly longer than the average read length of 136 bp (Table 1). The short  
19 length of the assembled MT contigs compared to the metagenomic (MG) contigs is likely due to the smaller  
20 size of the microbial transcripts compared to the larger size of the genomes. In terms of reproducibility, we  
21 did not detect significant differences between the number of reads or assembled contigs among the 12  
22 replicate libraries as analyzed by *Shapiro-Wilk* normality test (data not shown). Thus, the assembly of the  
23 contigs generated a modest increase in length compared with average read length of the MT reads.

24 Next, we compared the bacterial taxonomic assignments based on read and contig analyses.  
25 Analysis at the phyla, genera and species levels all demonstrated the reproducibility of the replicate  
26 libraries, respectively (Supplementary Figs. 2-4). However, we detected differences in the relative  
27 abundance of specific taxa in the read and contig based analyses. Thus, the taxonomic identification was  
28 inconsistent between read and contig based analysis at both phylum and genus level. For example, we  
29 observed an increase in the Bacteroidetes and decrease in Firmicutes with the contig analysis. Differences  
30 in relative abundance in the MT data were also observed at the genus and species levels. There were 21  
31 and 11 genera, and 22 and 19 species in the read and contig based analysis that were above 1%  
32 abundance, respectively (Supplementary Figs. 3 and 4, Supplementary Tables 1 and 2). We further  
33 analyzed the bacterial diversity of combined MT datasets (HS100, MS151, MS301 and HS100-MS151-  
34 MS301) to increase the sequencing depth and coverage. We find similar observations in the distribution of  
35 bacterial phyla (Supplementary Fig. 5A). We observe that the increase in number of reads resulted in  
36 increase of depth of coverage, whereas no significant increase in contig length was detected. In summary,  
37 we previously showed that a contig based analysis is more specific for species identification (Ranjan et al

1 2016) in the MGS dataset; however, these data suggest that a read based analysis is more comprehensive  
2 for identification of both genera and species in metatranscriptome data.

3 To determine if different numbers of reads were skewing the analyses, we generated datasets that  
4 contained an equal number of reads. We randomly sampled 30 million reads from the HiSeq 100 PE, MiSeq  
5 151 PE (MS151) and MiSeq 301 PE (MS301) data, and the reads were assembled into contigs. More  
6 contigs were generated in MS301 (97,631) compared to HS100 (8,253) and MS151 (42,153), most likely  
7 because of a longer sequencing read length. However, there was no substantial increase in the average  
8 length of contigs most likely due to the limitation based on transcript length (Supplementary Table 3). We  
9 observed a similar abundance profile of bacterial phyla, genera and species as in the complete datasets  
10 indicating that differences in read number were not skewing the assignment of taxa in the contig analyses  
11 (Supplementary Fig. 5B, and Supplementary Table 4).

### 13 **Comparison of the metatranscriptome with the metagenome**

14 In total, we identified 1,888 and 1291 bacterial species in the metagenome (MG) [MG-HS100-MS151-  
15 MS301, (Ranjan et al 2016)], and the metatranscriptome (MT) (MT-HS100-MS151-MS301) data,  
16 respectively (Fig. 2A). 1245 bacterial species were shared among the MG and MT (Fig. 2A), representing  
17 the metabolically active species, in the sample at this particular time point. In the phylum Firmicutes,  
18 Bacteroidetes, Actinobacteria, Proteobacteria, Fusobacteria, and Verrucomicrobia 356, 117, 138, 439, 23,  
19 and 6 species were shared, respectively. This accounted for 60% to 92% of the species shared between  
20 the MG and MT defined phyla (Fig. 2B). The detection of MG sequences lacking corresponding MT reads  
21 suggests unexpressed genes or even dormant bacteria. As expected, very few sequences were unique to  
22 the MT, and they were present in extremely low abundance (< 0.001%) presumably because transcripts  
23 are not expressed in the absence of the genome, and likely these sequences were not identified in MG  
24 because of relatively low abundance (Supplementary Table 6). Most (50%) of the sequences identified in  
25 the phylum proteobacteria were closely related to uncultured bacterial sequences. To determine the relative  
26 transcriptional activity of individual taxa and individual genes, we compared the relative abundance in the  
27 combined MT data (HS100-MS151-MS301) to our previously reported MG data for the same sample  
28 (Ranjan et al 2016). In an analysis of the MT at the phyla level, we observed that the abundance of  
29 Bacteroidetes transcripts was high, whereas the abundance of transcripts representing Firmicutes,  
30 Actinobacteria, Fusobacteria, and Verrucomicrobia was low. This was observed across all the sequencing  
31 platforms and read lengths (Fig. 2C). The abundance of the Fusobacteria and Verrucomicrobia was  
32 approximately 100-fold lower than the other Phyla (note Y-axis scale).

### 34 **Analysis of predicted biological functions**

35 We analyzed the functional profiles based on gene expression in the metatranscriptome using the MG-  
36 RAST KEGG annotation suite. KEGG annotates functions from level 1 through 4 with level 1 containing the  
37 most general categories and level 4 the most specific (Mitra et al 2011). We analyzed the data for biological

1 functions at all four levels. Of note, a similar relative abundance of the functions was detected at levels 1  
2 to 4 among the both read and contig based analysis, respectively (Supplementary Figs. 6-9), although  
3 minor differences were detected in the abundances of some functions at levels 3 and 4. We observe the  
4 similar distribution trends in the 30 million randomized MT reads and the assembled contigs  
5 (Supplementary Fig. 10). This implies that the identified functions are similar in either the read or contig  
6 based analysis of the MT data with slight variations.

7 We investigated the MG and MT data at the species level. Interestingly, we observed that few of  
8 the species (for example, *Faecalibacterium prausnitzii*, *Bacteroides spp.*, *B. thetaiotaomicron*, *B. vulgatus*,  
9 *B. ovatus* among others) had a higher relative representation in MT than MG, indicating that these species  
10 are highly transcriptionally active (Fig. 3 and Supplementary Fig. 11). However, the species *B. fragilis* did  
11 not have increased transcriptional activity as compared to other *Bacteroides spp.* As shown in a scatter  
12 plot, *F. prausnitzii*, *Bacteroides spp.*, and *Alistipes putredinis* were highly transcriptionally active at a  
13 significant level ( $\log_2$  fold difference  $\geq 3$ ,  $p$  adj.  $< 0.05$ ) whereas *Clostridium saccharolyticum*, *Eubacterium*  
14 *rectale* and *Ruminococcus obeum* ( $\log_2$  fold difference  $\geq -1$ ,  $p$  adj.  $< 0.05$ ) were low in transcriptional activity  
15 (Fig. 4A).

16 We compared the abundance of KEGG functions detected in the MT data to the predicted functions  
17 in the MG data. The analysis revealed that genes involved in translation, carbohydrate metabolism, and  
18 transcription were highly abundant in MT ( $\log_2$  fold change  $> 3$ ,  $p < 0.05$ ), compared to low abundance of  
19 glycan biosynthesis and metabolism, metabolism of cofactors and vitamins, replication and repair,  
20 membrane transport and amino acid metabolism ( $\log_2$  fold change  $> -2$ ,  $p$  adj.  $< 0.05$ ) (Fig. 4B). Translation  
21 and amino acid metabolism showed the largest differential expression with a fold change of  $> \pm 5$  ( $p$  adj.  
22  $< 0.05$ ), respectively. We observed similar patterns at the more specific levels 2, 3 and 4 (Supplementary  
23 Fig. 12-15). In this fecal sample, in total we detected 1916 functions at KEGG level 4 assignments in MG,  
24 compared to 1067 in MT. The MG and MT data shared 52% (1014) of the total functions, revealing the  
25 shared functional genes involved in active physiological functions of the gut microbiota which can be  
26 detected in MG and MT in a given time point (Fig. 5). Our analysis indicated that MG and MT overlapping  
27 genes are metabolically active genes. Genes which are only detected in the MT are even more  
28 metabolically active. On the other hand, if genes were detected only in MG and not in the MT, this may also  
29 suggest that genes may be present but not active in a given time.

30

### 31 **Contribution of functions in the metatranscriptome by individual bacterial phylum**

32 We further explored the functional contribution of the gut microbiota at the individual phylum level  
33 comprising of Firmicutes, Bacteroidetes, Actinobacteria, Proteobacteria, Fusobacteria and  
34 Verrucomicrobia, as these are abundant in the gut. There were differences in the expression of the genes  
35 in each phylum (Supplementary Figs. 16-18). At the KEGG Level 1 functional category, 50% of the functions  
36 were related to metabolism in each phylum (Firmicutes, Bacteroidetes, Actinobacteria, Proteobacteria,  
37 Fusobacteria and Verrucomicrobia), followed by genetic and environmental information processing



1 functional categories. Of note few functional categories related to the phylum Fusobacteria and  
2 Verrucomicrobia were detected (Supplementary Fig. 16). We further focused our analysis on Fusobacteria  
3 and Verrucomicrobia, as these phyla are present in low abundance (<1% and <0.1% abundance,  
4 respectively) and not well characterized in the gut microbiota (Fig. 2C).

5 In phyla - Firmicutes, Bacteroidetes, Actinobacteria, and Proteobacteria, the genes involved in  
6 carbohydrate metabolism were abundant, followed by amino acid metabolism and translation. There were  
7 no translation and/or transcription functions detected in Fusobacteria and Verrucomicrobia (Supplementary  
8 Fig. 17). However, Fusobacteria and Verrucomicrobia contributed towards the expression of specific genes  
9 involved in carbohydrate and amino acid metabolism pathways compared to other phyla (Figs. 8 and 9,  
10 Supplementary Fig. 18). For example, the genes *glgB* (1,4-alpha-glucan branching enzyme), *pgi* (glucose-  
11 6-phosphate isomerase) involved in starch, and sucrose metabolism and glycolysis/gluconeogenesis were  
12 highly expressed by Fusobacteria (Supplementary Fig. 18). Also, the genes involved in oxidative  
13 phosphorylation such as *atpD* (F-type H<sup>+</sup>-transporting ATPase subunit beta), *ppa* (inorganic  
14 pyrophosphatase) and *nuoE* (NADH-quinone oxidoreductase subunit E) were also enriched in Fusobacteria  
15 (Figs. 6, and Supplementary Fig. 18). On the other hand, the phylum Verrucomicrobia was enriched for  
16 genes involved in alanine, aspartate and glutamate metabolism [*gdhA*: glutamate dehydrogenase (NADP<sup>+</sup>),  
17 *purB*: adenylosuccinate lyase], ABC transporters [*msmX*: maltose/maltodextrin transport system ATP-  
18 binding protein] and amino sugar and nucleotide sugar metabolism [*npdA*: NAD-dependent deacetylase]  
19 (Fig. 7 and Supplementary Fig. 18). These results show the high abundance of transcripts contributed by  
20 the rare abundant bacterial species in the community may contribute unique biological functions to the  
21 microbiome that have the potential to affect the host physiology.

### 22 23 **Diversity analysis of bacterial species and functions**

24 The Shannon diversity index for estimating the bacterial diversity in MG ( $5.4 \pm 0.1$ ) and MT ( $4.9 \pm$   
25  $0.1$ ) was significantly different ( $p < 0.05$ ), however no significant difference was observed in species  
26 evenness ( $0.7 \pm 0.0$ ). Similarly, the index for diversity of functional genes in MG ( $6.7 \pm 0.0$ ) and MT ( $6.0 \pm$   
27  $0.3$ ) was significantly different ( $p < 0.05$ ), also a significant difference was observed in functional evenness  
28 in MG ( $0.89 \pm 0.01$ ) and MT ( $0.93 \pm 0.01$ ). The Shannon diversity index analysis at both taxonomic and  
29 functional level indicated that the MG was more diverse than the MT, most likely due to unexpressed genes  
30 or dormant bacteria (Supplementary Fig. 19).

### 31 32 **Mapping the genomic and transcriptomic KEGG pathways**

33 We mapped the predicted (MG) and expressed (MT) functions onto pathways using KEGG Mapper suite.  
34 Almost all (more than 99%) of the functions identified by MT were also identified in MG (Fig. 8 and  
35 Supplementary Fig. 20). However, some functions were identified only in the MG dataset suggesting that  
36 not all of the predicted functions in the metagenome are expressed, which supports the notion that the  
37 metagenome may not be an accurate proxy of microbiota function. The genes are in the (meta)genomes;

1 they could be expressed under different conditions; therefore, they define the functional potential of the  
2 organisms. Linear regression analysis was applied to the MT and MG data examined from the perspective  
3 of species and function. The linear regression analysis at the species level was correlated among the MG  
4 and MT and 58% of the variation in the MT can be explained by the species composition of the MG  
5 (Spearman's  $r = 0.83$ ;  $r^2=0.58=58\%$ ) (Fig. 9A). A similar correlation was observed at functional level 4 in  
6 MG and MT (Spearman's  $r = 0.76$ ;  $r^2=0.53=53\%$ ) (Fig. 9B). In other words, more than 50% of the variation  
7 in the microbial community MT can be explained by MG composition at species level, or conversely,  
8 approximately 50% of transcriptional activity is regulated and presumably dependent on host or  
9 environmental factors.

## 11 DISCUSSION

12 Dysbiosis of the microbiome has been associated with multiple disease states including obesity,  
13 inflammatory bowel disease, asthma, urinary tract infection, cardiovascular disease and cancer (Pflughoeft  
14 & Versalovic 2012, Rani et al 2016a, Rani et al 2016b). However, the biological mechanisms that link the  
15 complex community of a microbiota with the pathogenesis of most diseases remains elusive. One limitation  
16 of many studies has been the use of targeted 16S rRNA amplicon sequencing which is generally limited to  
17 the genus and or OTU level of classification, thus, a more specific classification at the species level is not  
18 available (Metwally et al 2016, Metwally et al 2018) . In contrast, MGS deep sequencing can accurately  
19 classify bacteria at the species level and also facilitates the annotation and identification of genes which  
20 predict putative biological functions. Further, due to the transcriptional regulation of many genes, MGS  
21 sequencing does not reveal gene expression levels. To address both the challenges, in this project we have  
22 optimized and evaluated the combination of metagenomic and metatranscriptomic shotgun sequencing  
23 data to evaluate methods to analyze the functional roles of both abundant and rare species in the  
24 microbiota. We generated 139.6 million metatranscriptomic reads which we compared to our previously  
25 reported metagenome shotgun sequencing data on the same sample that included 163.7 million reads  
26 (Ranjan et al 2016). One of the limitation of this study is sample size, as it is focused on n-of-1, and these  
27 findings may not be observed in different biological samples. However, with the advent of personalized  
28 medicine and clinical translational studies, there has been surge of n-of-1 studies. Many of clinical cases  
29 possess unique features that may not be identified by classical studies involving large number of samples  
30 (Nikles et al., 2010;Lillie et al., 2011;Schork, 2015).

31 First, our study shows that the different Illumina platforms do not contribute detectable bias in our  
32 analyses (Fig. 2). To validate the technical reproducibility of the sequencing and data analysis methods,  
33 we generated 12 replicates of a single sample that generated a similar number of reads, total bases and  
34 assembled contigs (Table 1). In addition, our analysis identified a reproducible number of both phyla and  
35 species (Supplementary Figs. 2 and 4, respectively). Furthermore, the functional analysis identified similar  
36 abundance of KEGG annotations at all functional levels from 1-4 (see Supplementary Figs. 6-9). Our  
37 investigation of the effect of contig assembly showed that assembly only modestly increased length,

1 presumably due to the short length of the mRNA transcripts. This similar observation has also been reported  
2 in a forest floor community metatranscriptomics (Hesse et al 2015). This suggests that emerging  
3 technologies that produce longer read lengths, particularly in view of their increased error rates, although  
4 useful for metagenomics studies, may not be preferable for metatranscriptomic studies.

5 Our investigation of the effects of contig assembly showed that the relative abundance of some  
6 taxa was modified by assembly. For example, analysis of assembled reads resulted in greater abundance  
7 of Bacteroidetes and lesser abundance of Firmicutes, Actinobacteria and Proteobacteria (Supplementary  
8 Fig. 5). Similar differences were also observed at the level of genus and species. Interestingly, we also  
9 observed similar changes in relative abundance of Bacteroidetes and Firmicutes in our previous analysis  
10 of taxa assignment in our metagenomics data (Ranjan et al 2016). Our results also show that the assembly  
11 of reads into contigs can decrease the detection of taxa. Overall, the results suggest that reads are the  
12 most comprehensive, and contigs are more specific, method to annotate taxa.

13 Most previous microbiota studies have not been performed with matched metagenome and  
14 metatranscriptome datasets of the same sample, thus there is huge knowledge gap in understanding the  
15 role of gene expression of the microbiota in human health and diseases. Our comparison of the predicted  
16 functions in the metagenome in this sample, with the expressed functions in metatranscriptome, identified  
17 more than 1000 functions, which included carbohydrate metabolism, nucleotide metabolism, amino acid  
18 metabolism, translation etc., (Fig. 5). The diversity analysis also suggest that the actual metabolically active  
19 bacterial species and functions are in fact less diverse compared to predicted metagenome diversity (both  
20 taxonomic and functional) (Supplementary Fig. 19).

21 It is well established that the diverse community of bacteria in a microbiome is composed of a small  
22 number of abundant species plus a large number of low or rare abundance species (Ranjan et al 2016);  
23 however, the functional role of the abundant versus rare species is not well understood. Our comparison of  
24 the metatranscriptome with the metagenome data suggests that both the abundant and rare bacteria may  
25 be actively engaged in the gut ecosystem. For instance, bacterial transcripts representing phyla Firmicutes  
26 (*F. prausnitzii*), and Bacteroidetes (*Bacteroides spp.*, and *B. uniformis*) were highly abundant in MT (Fig.  
27 3). Bacterial phyla - Fusobacteria and Verrucomicrobia are relatively less abundant in human gut, but are  
28 known to play an important role gut physiology (Everard et al 2013, Tremaroli & Backhed 2012). For  
29 instance, in our sample, both these phyla actively contributed in expression of specific genes involved in  
30 carbohydrate and amino acid metabolism pathways (Figs. 6 and 7). For example, genes such as *glgB* (1,4-  
31 alpha-glucan branching enzyme) and *pgi* (glucose-6-phosphate isomerase) involved in starch and sucrose  
32 metabolism and gluconeogenesis/glycolysis were highly expressed by Fusobacteria. These data suggest  
33 that the low abundant bacterial species are not just mere bystanders but actively contribute to the gut  
34 ecology. A similar study using the matched metagenomics and metatranscriptomics of the same sample  
35 have observed comparable findings that microbial and metabolic potential vary and are not concordant with  
36 their taxonomic abundance (Franzosa et al 2014). The functional potential of the more and less abundant  
37 bacterial species remain poorly understood. However, our observations indicate that the less abundant

1 species are also metabolically active and may play unique roles in host-bacteria and bacteria-bacteria  
2 interactions and may actively contribute to the gut microbiota and physiology.

### 4 **ACKNOWLEDGEMENTS**

5 This work was supported in part by NIH RO1 HL081663 and NIH RO1 AI053878 to DLP and PWF. The  
6 authors acknowledge Mr. Samer Sabbagh for help with preparing the libraries.

### 8 **AUTHOR CONTRIBUTIONS**

9 DLP, PWF, RR and AR designed the study: RR prepared libraries and performed sequencing, AR and RR  
10 performed data analysis, RR, AR, PWF and DLP wrote the manuscript.

### 12 **COMPETING FINANCIAL INTERESTS**

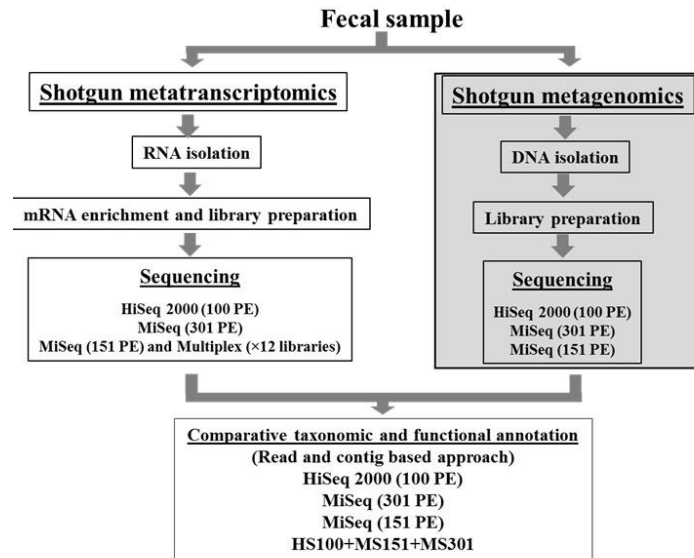
13 The authors have declared that there is no conflict of interest. The funders had no role in study design, data  
14 collection and analysis, decision to publish, or preparation of the manuscript.

16 **SEQUENCE DATASETS:** The sequence data files have been submitted to MG-RAST and the accession  
17 numbers are mentioned in Supplementary Table 5.

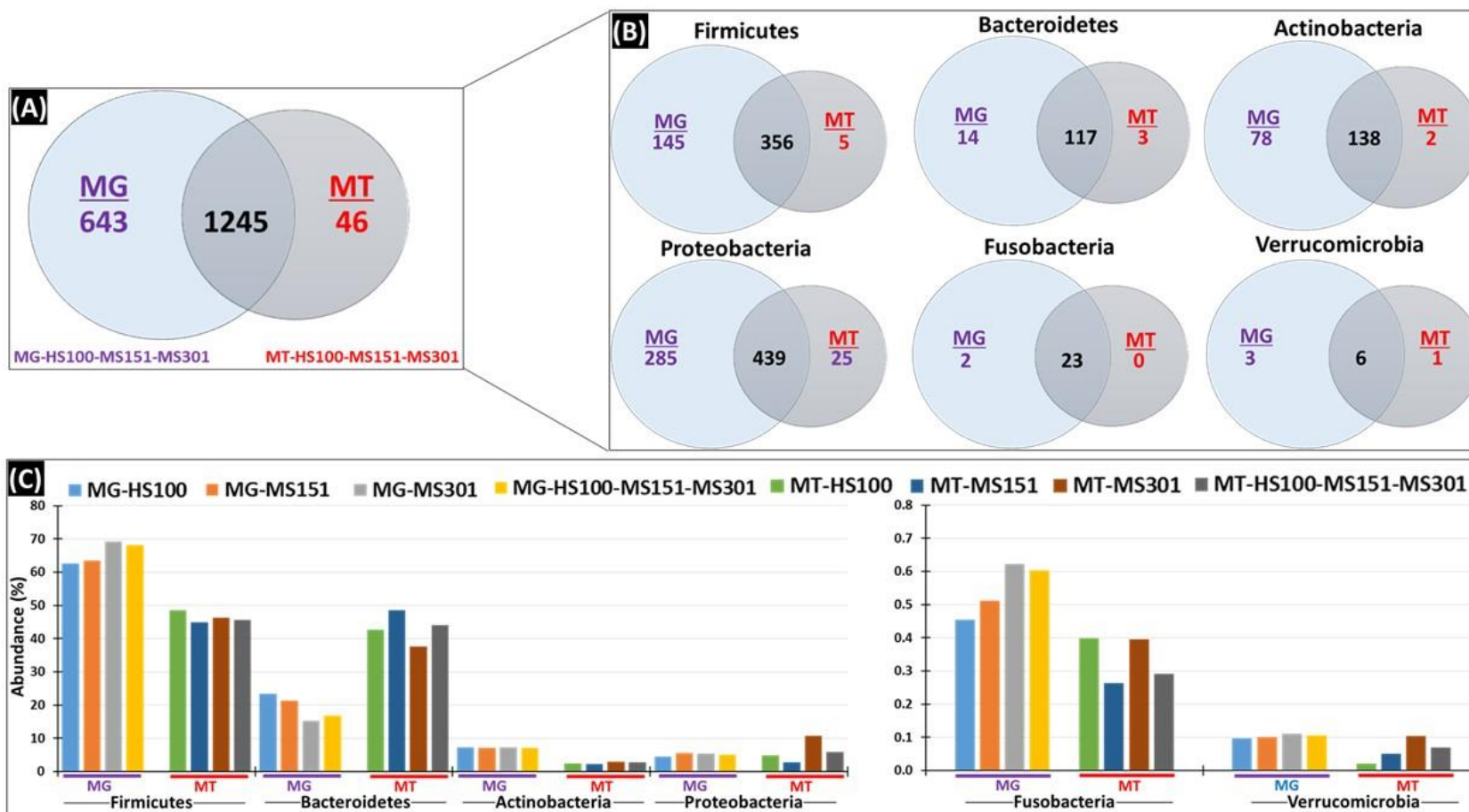
### 19 **REFERENCES**

- 20 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat*  
21 *Soc* 57: 289 - 300
- 22 Boojink CC, Boekhorst J, Zoetendal EG, Smidt H, Kleerebezem M, de Vos WM. 2010. Metatranscriptome analysis of the human  
23 fecal microbiota reveals subject-specific expression profiles, with genes encoding proteins involved in carbohydrate  
24 metabolism being dominantly expressed. *Appl. Environ. Microbiol.* 76: 5533-40
- 25 Cho I, Blaser MJ. 2012. The human microbiome: at the interface of health and disease. *Nat. Rev. Genet.* 13: 260-70
- 26 Everard A, Belzer C, Geurts L, Ouwerkerk JP, Druart C, et al. 2013. Cross-talk between *Akkermansia muciniphila* and intestinal  
27 epithelium controls diet-induced obesity. *Proc. Natl. Acad. Sci. U. S. A.* 110: 9066-71
- 28 Franzosa EA, Morgan XC, Segata N, Waldron L, Reyes J, et al. 2014. Relating the metatranscriptome and metagenome of the  
29 human gut. *Proc. Natl. Acad. Sci. U. S. A.* 111: E2329-38
- 30 Hesse CN, Mueller RC, Vuyisich M, Gallegos-Graves LV, Gleasner CD, et al. 2015. Forest floor community metatranscriptomes  
31 identify fungal and bacterial responses to N deposition in two maple forests. *Front. Microbiol.* 6
- 32 Human Microbiome Project C. 2012. Structure, function and diversity of the healthy human microbiome. *Nature* 486: 207-14
- 33 Huse SM, Ye Y, Zhou Y, Fodor AA. 2012. A Core Human Microbiome as Viewed through 16S rRNA Sequence Clusters. *PLoS One*  
34 7: e34242
- 35 Iebba V, Totino V, Gagliardi A, Santangelo F, Cacciotti F, et al. 2016. Eubiosis and dysbiosis: the two sides of the microbiota. *New*  
36 *Microbiol.* 39: 1-12
- 37 Metwally AA, Dai Y, Finn PW, Perkins DL. 2016. WEVOTE: Weighted Voting Taxonomic Identification Method of Microbial  
38 Sequences. *PLoS One* 11: e0163527
- 39 Metwally AA, Yang J, Ascoli C, Dai Y, Finn PW, Perkins DL. 2018. MetaLonDA: a flexible R package for identifying time intervals of  
40 differentially abundant features in metagenomic longitudinal studies. *Microbiome* 6: 32
- 41 Meyer F, Paarmann D, D'Souza M, Olson R, Glass EM, et al. 2008. The metagenomics RAST server - a public resource for the  
42 automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* 9: 386

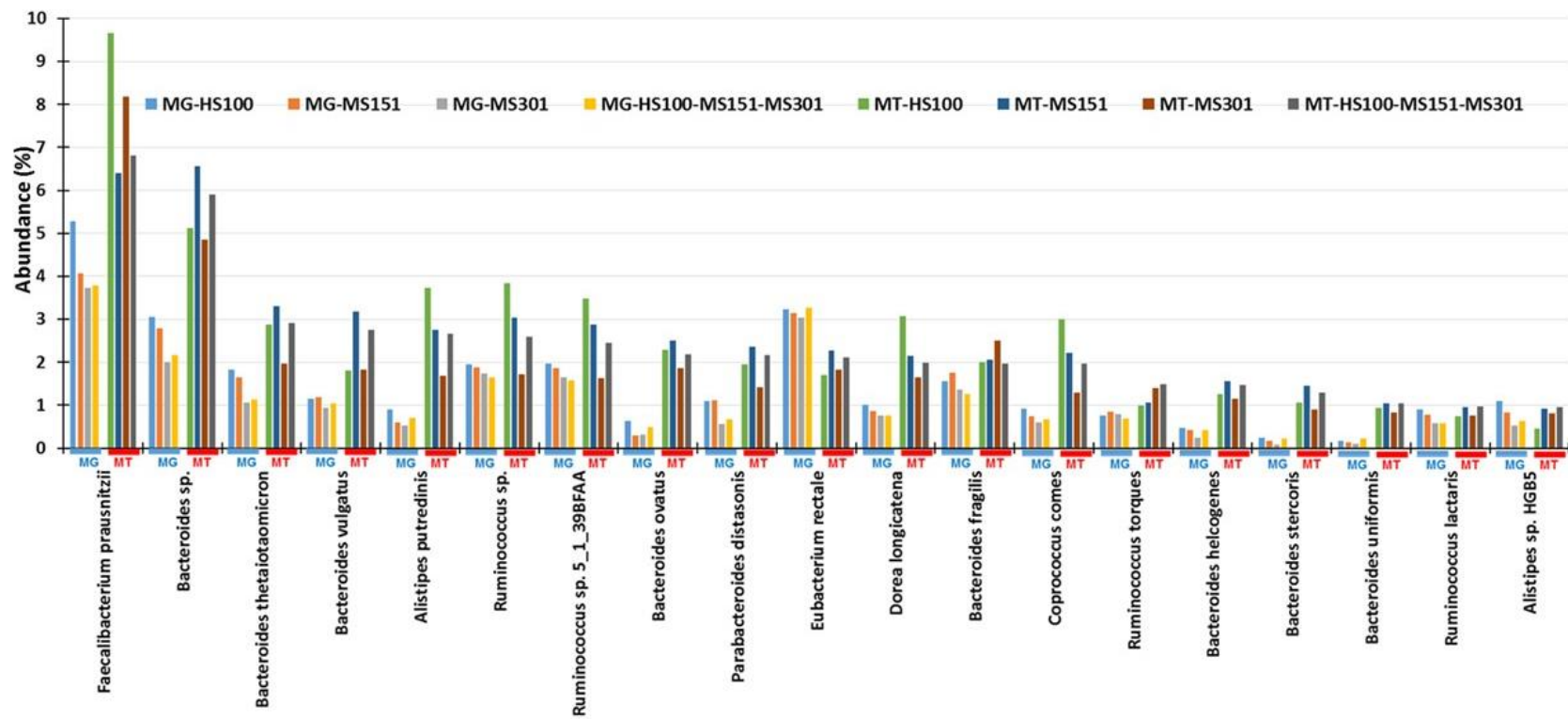
- 1 Mitra S, Rupek P, Richter DC, Urich T, Gilbert JA, et al. 2011. Functional analysis of metagenomes and metatranscriptomes using  
2 SEED and KEGG. *BMC Bioinformatics* 12 Suppl 1: S21
- 3 Pflughoeft KJ, Versalovic J. 2012. Human microbiome in health and disease. *Annu. Rev. Pathol.* 7: 99-122
- 4 Praveen P, Jordan F, Priami C, Morine MJ. 2015. The role of breast-feeding in infant immune system: a systems perspective on the  
5 intestinal microbiome. *Microbiome* 3: 41
- 6 Rani A, Ranjan R, McGee HS, Andropolis KE, Panchal DV, et al. 2016a. Urinary microbiome of kidney transplant patients reveals  
7 dysbiosis with potential for antibiotic resistance. *Transl. Res.*
- 8 Rani A, Ranjan R, McGee HS, Metwally A, Hajjiri Z, et al. 2016b. A diverse virome in kidney transplant patients contains multiple  
9 viral subtypes with distinct polymorphisms. *Sci. Rep.* 6: 33327
- 10 Ranjan R, Rani A, Metwally A, McGee HS, Perkins DL. 2016. Analysis of the microbiome: Advantages of whole genome shotgun  
11 versus 16S amplicon sequencing. *Biochem. Biophys. Res. Commun.* 469: 967-77
- 12 Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, et al. 2015. limma powers differential expression analyses for RNA-sequencing and  
13 microarray studies. *Nucleic Acids Res.* 43: e47
- 14 Saeed AI, Sharov V, White J, Li J, Liang W, et al. 2003. TM4: a free, open-source system for microarray data management and  
15 analysis. *BioTechniques* 34: 374-8
- 16 Sharma R, Ranjan R, Kapardar RK, Grover A. 2005. 'Unculturable' bacterial diversity: An untapped resource. *Curr. Sci.* 89: 72-77
- 17 Smyth GK. 2004. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat.*  
18 *Appl. Genet. Mol. Biol.* 3: Article3
- 19 Stulberg E, Fravel D, Proctor LM, Murray DM, LoTempio J, et al. 2016. An assessment of US microbiome research. *Nature*  
20 *Microbiology* 1: 15015
- 21 Tremaroli V, Backhed F. 2012. Functional interactions between the gut microbiota and host metabolism. *Nature* 489: 242-9
- 22 Ursell LK, Knight R. 2013. Xenobiotics and the human gut microbiome: metatranscriptomics reveal the active players. *Cell Metab.*  
23 17: 317-8
- 24



**Figure 1. Experimental strategy to compare the metatranscriptome and metagenome using multiple Illumina sequencing platforms and data analysis.** Schematic for metagenome and metatranscriptome sequence analysis by shotgun sequencing approach. The shotgun sequencing was performed using Illumina HiSeq 2000 (100 paired-end), and Illumina MiSeq (151 and 301 paired-end). The data was analyzed by read and contig based approach using the MG-RAST. Note that the metagenome data has been published (Ranjan et al., 2016), represented in shaded box.

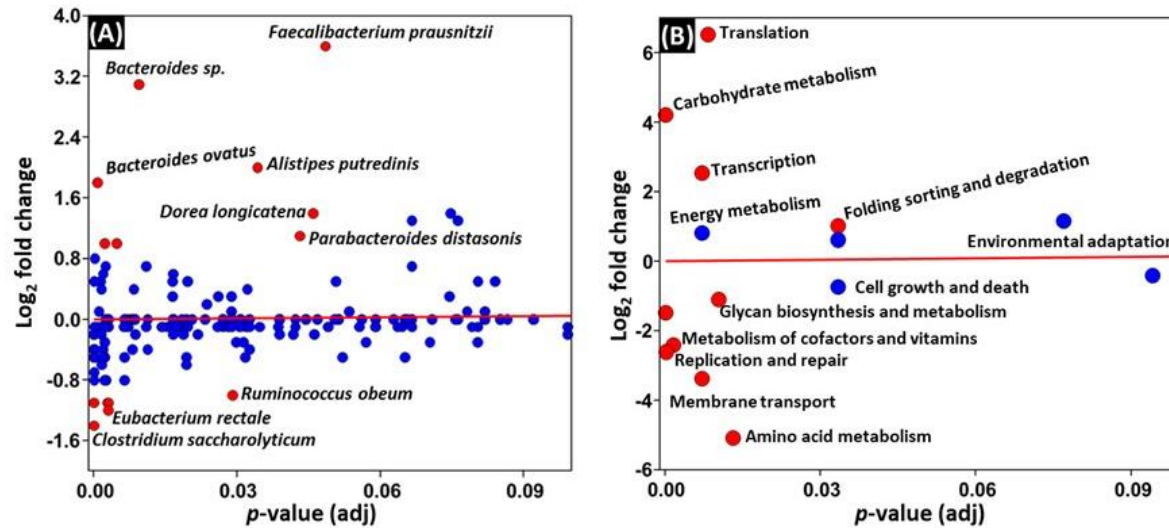


**Figure 2. Taxonomic analysis: Comparison of metagenome (MG) and Metatranscriptome (MT).** The MG and MT sequence obtained after sequencing using platforms (HS100, MS151 and MS301) were assembled into contig and were analyzed for taxonomic annotation. (A) The total bacterial species in MG-HS100-MS151-MS301 and MT-HS100-MS151-MS301 data. (B) Bacterial species in MG-HS100-MS151-MS301 and MT-HS100-MS151-MS301 in different phyla. (C) The abundance of bacterial phyla in MG and MT with different sequencing parameters - Firmicutes, Bacteroidetes, Actinobacteria, Proteobacteria, Fusobacteria and Verrucomicrobia.

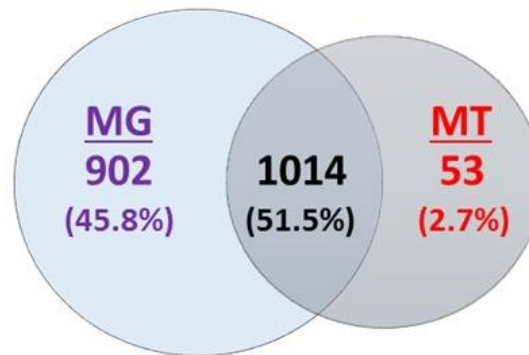


**Figure 3. Abundance of bacterial species in metagenome and metatranscriptome.** Bacterial species above 1% (sorted high to low) are shown in MT-HS100-MS151-MS301.

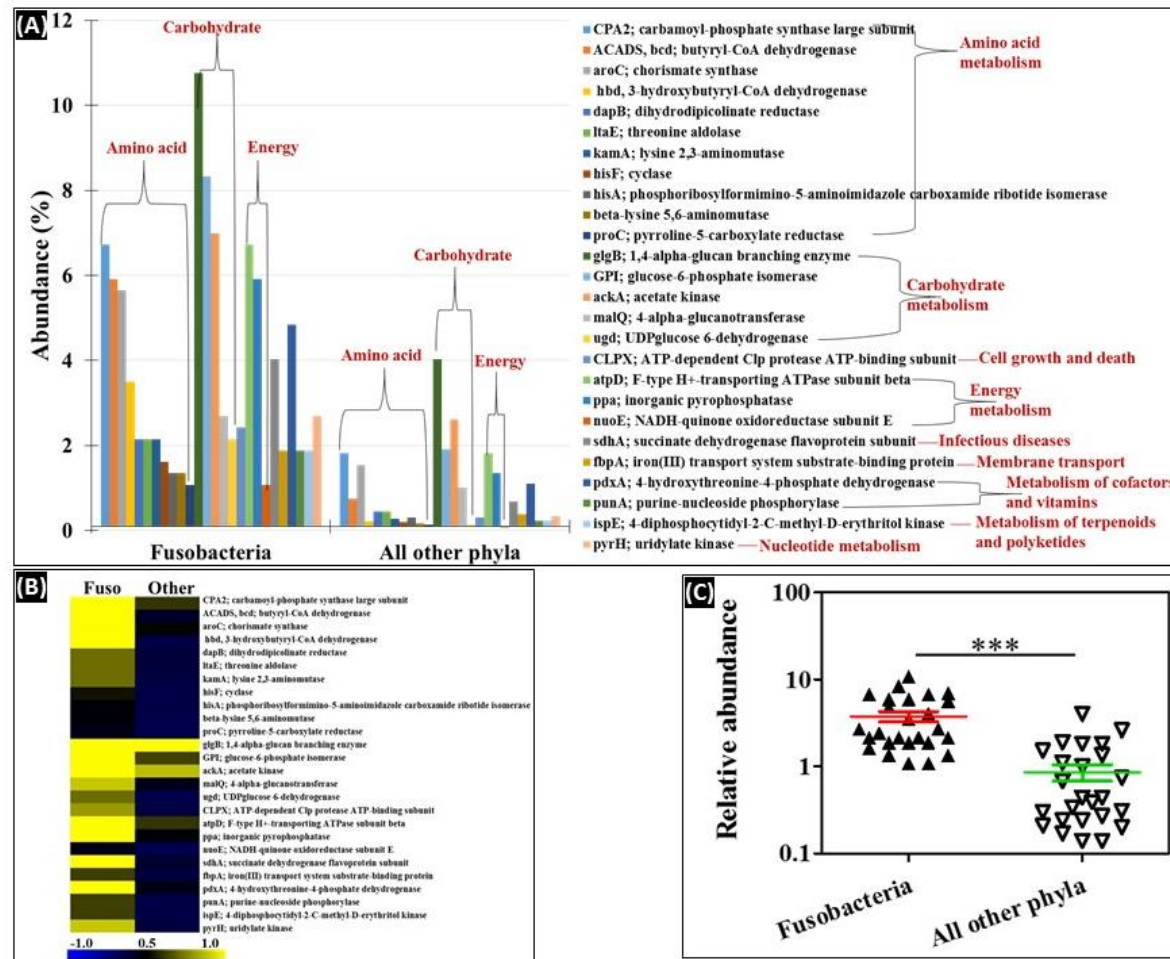




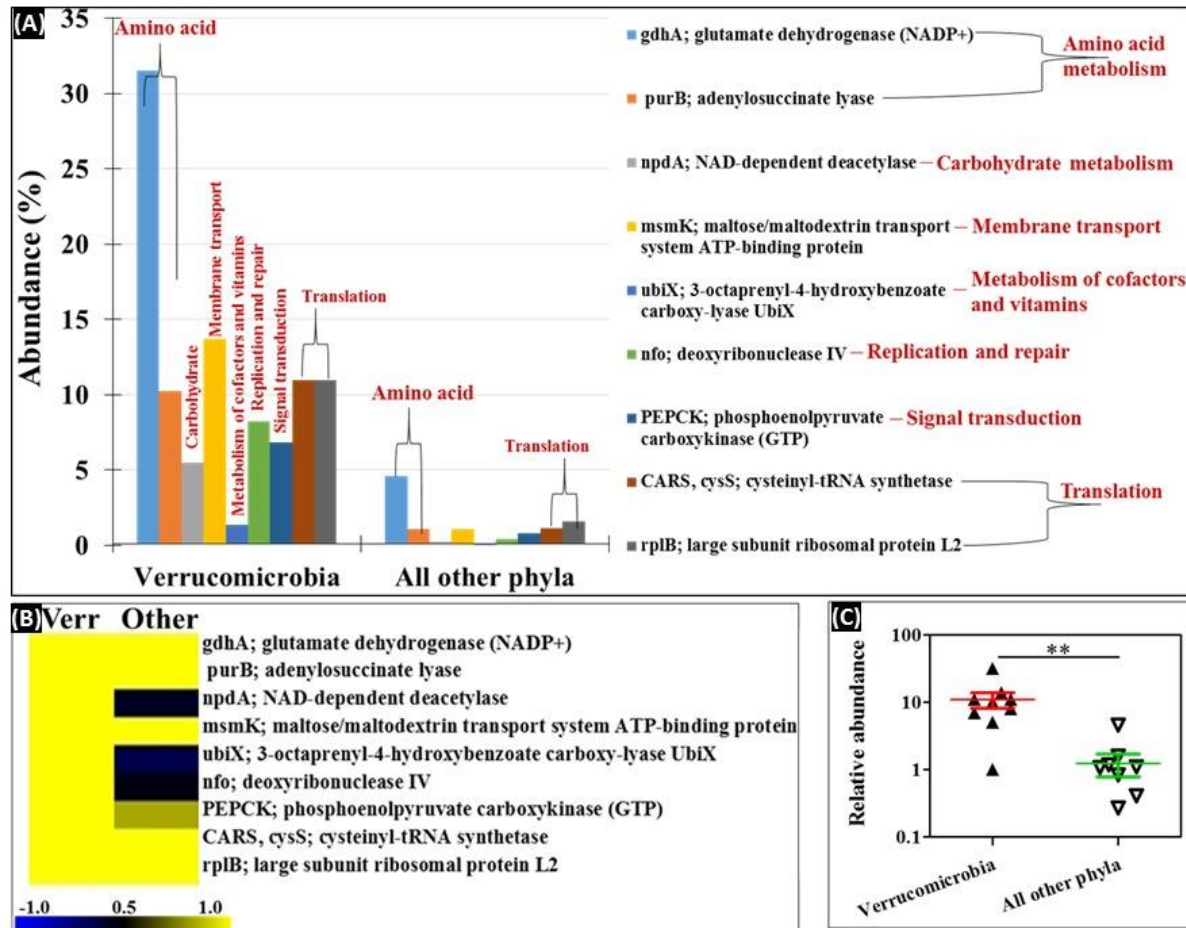
**Figure 4. Differential abundant species and KEGG functional categories.** The scatter plot for differential abundant bacterial species (A) and differentially predicted and expressed KEGG functional categories (B) in the metagenome and metatranscriptome. A  $p$  value cutoff of 0.05 (after FDR correction based on Benjamini-Hochberg method) and a log fold change  $\geq 1$  were used to select the differentially abundant species and functional categories. Significant values for different species and pathways are shown in red and non-significant values are shown with blue circles.



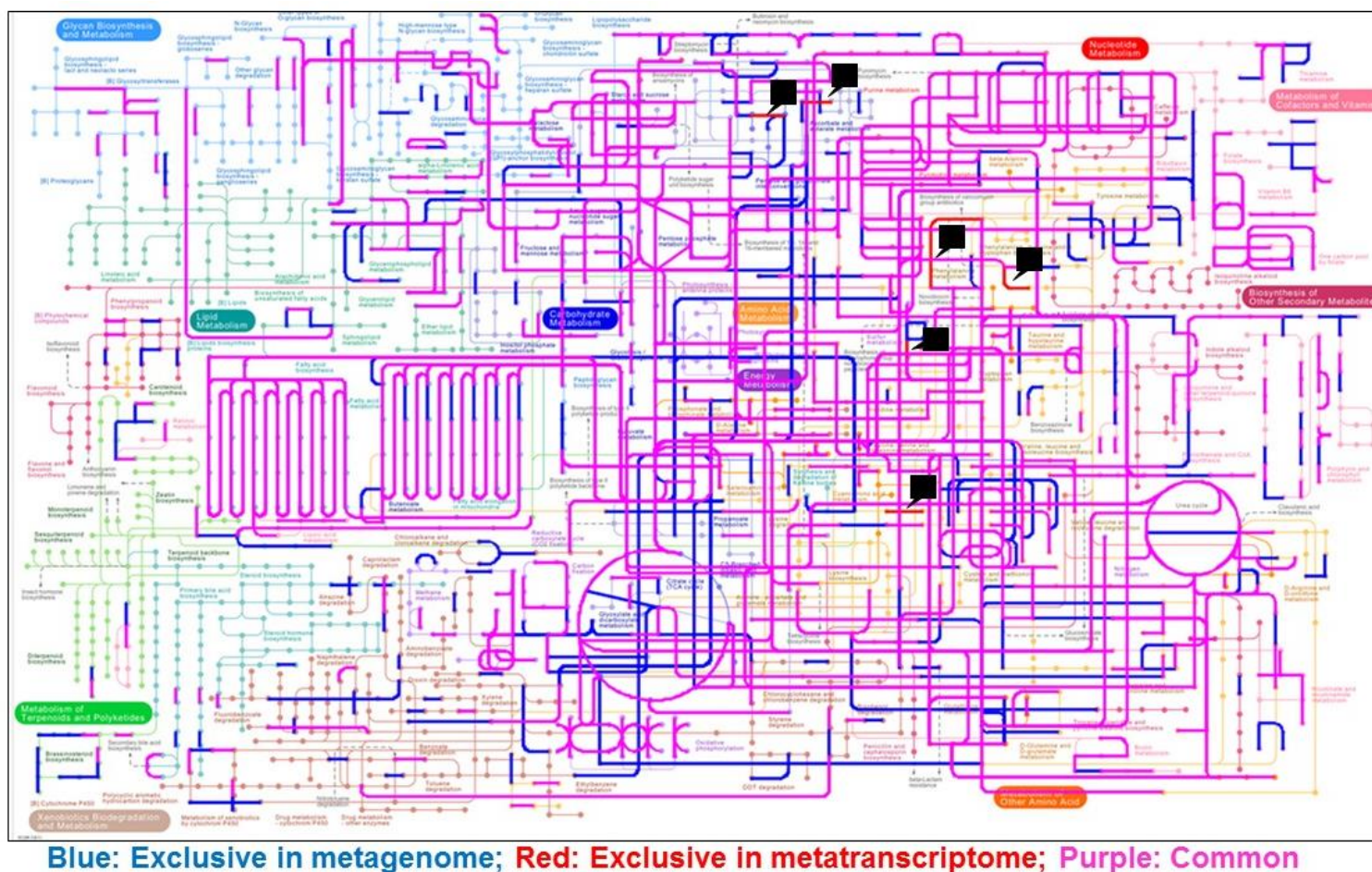
**Figure 5. Comparison of the metabolic functional of metagenome (MG) and metatranscriptome (MT).** Venn diagram for unique and shared metabolic functions identified by KEGG at functional level 4 in the MG (MG-HS100-MS151-MS301) and MT (MT-HS100-MS151-MS301).



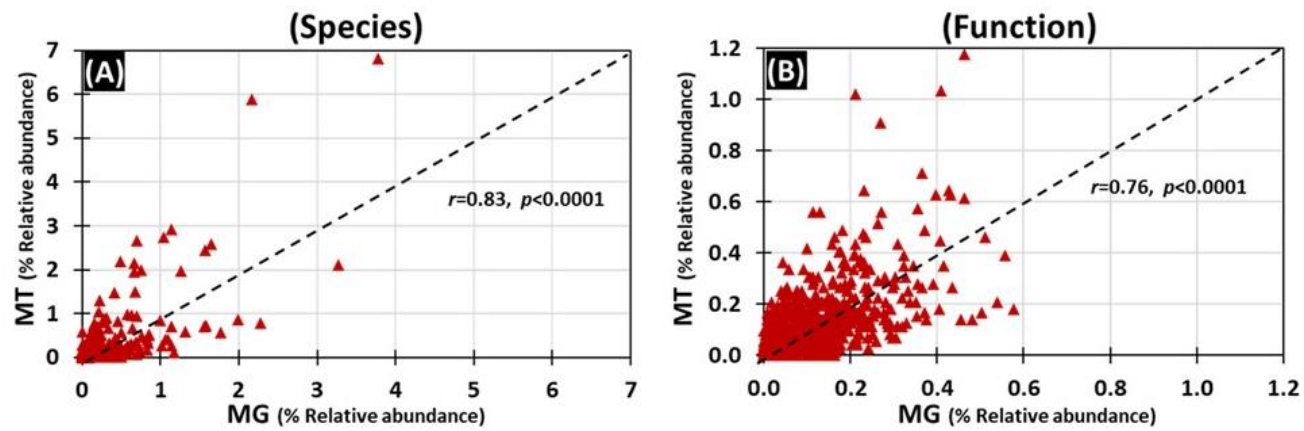
**Figure 6. Metatranscriptome analysis of phylum Fusobacteria.** (A) Relative abundance of Fusobacteria genes compared to all other phyla. (B) Heat-map representation of the genes. The color scheme represents the range of gene abundance values based on Spearman Rank correlation. (C) Significant difference in log abundance of genes highly abundant in Fusobacteria compared to all other phyla.  $p < 0.05$ , Mann-Whitney  $U$  test. Other phyla include Firmicutes, Bacteroidetes, Proteobacteria, Actinobacteria and Verrucomicrobia.



**Figure 7. Metatranscriptome analysis of phylum Verrucomicrobia.** (A) Relative abundance of Verrucomicrobia genes compared to all other phyla. (B) Heat-map representation of the genes. The color scheme represents the range of gene abundance values based on Spearman Rank correlation. (C) Significant difference in log abundance of genes highly abundant in Verrucomicrobia compared to all other phyla.  $p < 0.05$ , Mann-Whitney  $U$  test. Other phyla include Firmicutes, Bacteroidetes, Proteobacteria, Actinobacteria and Fusobacteria.



**Figure 8. Differential metabolic gene expression.** Metabolic pathway reconstruction in metagenome and metatranscriptome were analyzed using the KEGG mapper. Functions identified in the metagenome (MG-HS100+MS151+MS301) and metatranscriptome (MT-HS100+MS151+MS301). Blue: predicted functions exclusive in metagenome; Purple: Common in metagenome and metatranscriptome; Red: Exclusive in metatranscriptome. Black arrow head represents the functions in MT. Function in individual data are shown in Supplementary Fig. 20.



**Figure 9. Correlation between the metagenome and metatranscriptome.** Linear regression analysis was applied to the MT and MG data examined from the perspective of species and function. Spearman's rank correlation between MG and MT (A) Bacterial species, (B) Functions at KEGG Level 4.

**Table 1. Metatranscriptome sequence statistics.**

Sample name	Read			Contig		
	Number of PE reads (M)	Average length (bp)	Total bases (Mb)	% of reads assembled in contig	Number of contig	Average length (bp)
MS151_library 1	4.8	145	690.5	98.4	7,253	207
MS151_library 2	4.8	143	690.2	98.2	7,517	209
MS151_library 3	5.4	143	765.2	98.4	8,291	202
MS151_library 4	4.5	143	649.9	98.4	7,364	209
MS151_library 5	4.8	144	686.8	98.0	7,072	212
MS151_library 6	4.3	145	625.0	98.0	6,183	213
MS151_library 7	4.9	143	696.9	98.4	7,635	204
MS151_library 8	3.6	147	525.9	97.3	5,889	222
MS151_library 9	4.9	146	707.4	97.9	7,875	208
MS151_library 10	4.5	145	653.9	98.4	6,779	215
MS151_library 11	4.8	145	698.7	98.1	7,117	211
MS151_library 12	5.3	145	768.9	98.0	8,249	212
HS100	50.4	100	5,039.2	98.6	11,713	203
MS151 (Lib-all)	56.5	144	8,159.3	99.3	56,491	208
MS301	32.7	178	5,837.7	98.3	108,905	314
HS100+MS151+MS301	139.6	136	19,036.1	99.1	216,712	268

HS100: HiSeq 2000 - 100 PE; MS151: MiSeq - 151 PE; MS151: MiSeq - 301 PE; M: Million; bp: basepair; Mb: Mega bases; PE: Paired-end sequencing.

## SUPPLEMENTARY FIGURE AND TABLE LEGENDS

**Supplementary Figure 1. Fecal metatranscriptome library preparation.** High quality total RNA from a fecal sample was isolated and analyzed by agarose gel electrophoresis and Bioanalyzer (A); Total RNA was enriched for mRNA by depleting the rRNA by subtractive hybridization method (B), the enriched mRNA was fragmented by Covaris (C); A library was prepared using Illumina compatible adaptor (D); In addition, 12 libraries from the same mRNA were prepared for multiplexing (E). The quality of RNA, mRNA and the libraries was analyzed on 2100 Bioanalyzer Instrument.

**Supplementary Figure 2. Phylum level analysis of multiplexed libraries using read and contig based analysis.** The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for bacterial taxonomic assignment at phylum level using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib-all.

**Supplementary Figure 3. Genus level analysis of multiplexed libraries using read and contig based analysis.** The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for bacterial taxonomic assignment at genus level using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib-all, and top 1% genus are shown (data sorted high to low abundance in Lib-all).

**Supplementary Figure 4. Species level analysis of multiplexed libraries using read and contig based analysis.** The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for bacterial taxonomic assignment at species level using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib-all, and top 1% bacterial species are shown (data sorted high to low abundance in Lib-all).

**Supplementary Figure 5. Taxonomic analysis of the metatranscriptome.** The read and contig based analysis of HS100, MS151, MS301, and HS100-MS151-MS301 (A). (B) The MT for each sequencing strategy (HS100, MS151 and MS301) was sampled for 30M reads. The reads were assembled into contigs and analyzed for taxonomic annotations based on read and contig. Data is sorted high to low on MS301\_read dataset.

**Supplementary Figure 6. Functional analysis of metatranscriptome at level 1 of multiplexed libraries using read and contig based analysis.** The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for functional assignment at Level 1 using MGRAST KEGG module using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib-all, and all the six Level 1 functions are shown.

**Supplementary Figure 7. Functional analysis of metatranscriptome at level 2 of multiplexed libraries using read and contig based analysis.** The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for functional assignment at Level 2 using



MGRAST KEGG module using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib-all, and top 10 Level 2 functions are shown. The data is sorted high to low on Lib1.

**Supplementary Figure 8. Functional analysis of metatranscriptome at level 3 of multiplexed libraries using read and contig based analysis.**

The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for functional assignment at Level 3 using MGRAST KEGG module using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib1-12, and top 10 Level 3 functions are shown. The data is sorted high to low on Lib1.

**Supplementary Figure 9. Functional analysis of metatranscriptome at functional level 4 of multiplexed libraries using read and contig based analysis.**

The twelve metatranscriptome libraries were sequenced on Illumina MiSeq (151 PE) and analyzed for functional assignment at Level 4 using MGRAST KEGG module using sequence read (A) and assembled contigs (B). Also, all the twelve libraries were combined in-silico and called as Lib1-12, and top 1% Level 4 functions are shown. The data is sorted high to low on Lib1.

**Supplementary Figure 10. Functional analysis of metatranscriptome based on read and contig.**

(A) Level 1, (B) Level 2, (C) Level 3, and (4) Functional. The MT for each sequencing strategy (HS100, MS151 and MS301) was sampled for 30M reads. The reads were assembled into contigs and analyzed for taxonomic annotations based on read and contig. Data is sorted high to low on MS301\_read dataset. For Level 1 all functional categories are shown, for Levels 2-4, only top 10 functions are shown.

**Supplementary Figure 11. Abundance of bacterial species in different phyla in MG and MT.**

Abundance of bacterial species in different phyla - Firmicutes, Bacteroidetes, Actinobacteria, Proteobacteria, Fusobacteria and Verrucomicrobia. Note the higher abundance percentage in metatranscriptome compared to metagenome data, indicating that some species are more metabolically active. Only top 10 species are shown for MT-HS100-MS151-MS301 and MG-HS100-MS151-MS301 (data sorted on MT-HS100-MS151-MS301).

**Supplementary Figure 12. Functional analysis at level 1.**

Percent abundance of the predicted (based on metagenome) and expressed (metatranscriptome) function. Data is sorted high to low on MT-HS100+MS151+MS301.

**Supplementary Figure 13. Functional analysis at Level 2:**

Percent abundance of the predicted (based on metagenome) and expressed (metatranscriptome) function. Functions are sorted high to low on MT-HS100+MS151+MS301 and above 1% are reported.

**Supplementary Figure 14. Functional analysis at Level 3.**

Percent abundance of the predicted (based on metagenome) and expressed (metatranscriptome) function. Data is sorted high to low on MT\_HS100+MS151+MS301, and top 10 functions and above 1% are reported.

**Supplementary Figure 15. Functional analysis at functional level.**

Percent abundance of the predicted (based on metagenome) and expressed (metatranscriptome) function. Functions is sorted high to low on MT\_HS100+MS151+MS301 and top 10 functions are reported.

**Supplementary Figure 16. Functional analysis at level 1 in individual phylum.** The functions in individual phylum were analyzed in the metatranscriptome (MT-HS100+MS151+MS301) data.

**Supplementary Figure 17. Functional analysis at level 2 in individual phylum.** The functions in individual phylum were analyzed in the metatranscriptome (MT-HS100+MS151+MS301) data.

**Supplementary Figure 18. Functional analysis at level 3 in individual phylum.** The functions in individual phylum were analyzed in the metatranscriptome (MT-HS100+MS151+MS301) data and sorted high to low on Firmicutes and above 1% functions are shown.

**Supplementary Figure 19. Diversity indices for bacterial species (A) and functions (B), in MG and MT.** The Shannon diversity and evenness are calculated for MG using the contig assembly of data MG-HS100, MG-MS151, MG-MS301 and MG-HS100+MS151+MS301, and MT using the contig assembly of data MT-HS100, MT-MS151, MT-MS301 and MT-HS100+MS151+MS301.

**Supplementary Figure 20. KEGG metabolic pathway in metagenome and metatranscriptome.** Functions identified in the metagenome (MG-HS100+MS151+MS301) and metatranscriptome (MT-HS100+MS151+MS301). Blue: predicted functions exclusive in metagenome; Red: Exclusive in metatranscriptome.

**Supplementary Table 1.** List of bacterial species identified based on read based analysis. Only above 1% are mentioned and sorted high to low on Lib1.

**Supplementary Table 2.** List of bacterial species identified based on contig based analysis. Only above 1% are mentioned and sorted high to low on Lib1.

**Supplementary Table 3.** Random sampling of the metatranscriptome sequence read and de-novo assembly of contigs.

**Supplementary Table 4:** Abundance of bacterial species in metatranscriptome data based on read and contig analysis.

**Supplementary Table 5.** List of accession numbers.

**Supplementary Table 6.** List of bacterial species/sequences identified in the metatranscriptomics data.