**Title**

MetaMap: An atlas of metatranscriptomic reads in human disease-related RNA-seq data

**Authors**

Simon LM1, Karg S1, Westermann AJ2,3, Engel M1,4, Elbehery AHA5, Hense B1, Heinig

M1, Deng L5, Theis FJ1,6

**Affiliations**

1 Helmholtz Zentrum München, German Research Center for Environmental Health, Institute

of Computational Biology, Neuherberg, Germany

2 Institute for Molecular Infection Biology, University Würzburg, Würzburg, Germany

3 Helmholtz Institute for RNA-Based Infection Research (HIRI), Würzburg, Germany

4 Helmholtz Zentrum München, German Research Center for Environmental Health,

Scientific Computing Research Unit, Neuherberg, Germany

5 Helmholtz Zentrum München, German Research Center for Environmental Health, Institute

of Virology, Neuherberg, Germany

6 Department of Mathematics, Technische Universität München, Munich, Germany

**Corresponding authors**

Simon LM; lukas.simon@helmholtz-muenchen.de

Theis FJ; fabian.theis@helmholtz-muenchen.de

20 **Abstract**

21 Background: With the advent of the age of big data in bioinformatics, large volumes of data

22 and high performance computing power enable researchers to perform re-analyses of

23 publicly available datasets at an unprecedented scale. Ever more studies imply the

24 microbiome in both normal human physiology and a wide range of diseases. RNA

25 sequencing technology (RNA-seq) is commonly used to infer global eukaryotic gene

26 expression patterns under defined conditions, including human disease-related contexts, but

27 its generic nature also enables the detection of microbial and viral transcripts.

28 Findings: We developed a bioinformatic pipeline to screen existing human RNA-seq datasets

29 for the presence of microbial and viral reads by re-inspecting the non-human-mapping read

30 fraction. We validated this approach by recapitulating outcomes from 6 independent

31 controlled infection experiments of cell line models and comparison with an alternative

32 metatranscriptomic mapping strategy. We then applied the pipeline to close to 150 terabytes

33 of publicly available raw RNA-seq data from >17,000 samples from >400 studies relevant to

34 human disease using state-of-the-art high performance computing systems. The resulting

35 data of this large-scale re-analysis are made available in the presented MetaMap resource.

36 Conclusions: Our results demonstrate that common human RNA-seq data, including those

37 archived in public repositories, might contain valuable information to correlate microbial and

38 viral detection patterns with diverse diseases. The presented MetaMap database thus

39 provides a rich resource for hypothesis generation towards the role of the microbiome in

40 human disease.

41

42

43

44 **Keywords**

47

## Data Description

### Context

Recent studies have demonstrated the paramount importance of the microbiome for human health and disease [1]. For example, imbalance of the human gut microbiome was linked to non-communicable diseases such as obesity [2,3], diabetes [4], cardiovascular disease [5], chronic obstructive pulmonary disease [6], or colorectal carcinoma [7,8], to name just a few.

The advent of high-throughput sequencing technologies has revolutionized the life sciences. RNA-seq technology produces one of the most frequent next generation sequencing data types and has been applied to study a large number of biological samples relevant to human disease. The majority of the underlying raw data is freely accessible from data repositories such as the Gene Expression Omnibus (GEO) (>1,700 human RNA-seq data sets as of january 2018) or the Sequence Read Archive (SRA) [9].

However, these data are typically exclusively used for single species (i.e. human) transcriptomics such as differential gene expression or alternative splicing analysis [9,10]. Reads that do not map onto the human genome are considered noise or contamination and therefore generally ignored [11,12] (collectively about 9% of total reads, Fig. 1). Five years ago, it was postulated that interspecies interactions might be studied by simultaneous detection and quantification of RNA transcripts from a given host and a microbe via 'dual' RNA-seq [13]. Meanwhile this approach has been successfully applied to the interaction of mammalian cells with diverse bacterial [14] and viral pathogens [15–19].

Inspired by dual RNA-seq, in this study we hypothesize that reads in archived RNA-seq datasets derived from human primary cells or tissue samples that fail to map against the human reference genome may contain valuable information about the presence of certain microbes in the respective body niches and/or under defined disease conditions. To enable metatranscriptomic study of these data, we combined existing read alignment and metagenomic classification software into a two-step 'omni' RNA-seq pipeline to

75   comprehensively quantify archaeal, bacterial and viral reads in human RNA-seq data (Fig.

76   1).

77         In the first step of this so called 'Metamap' pipeline, all reads are aligned against the

78   human genome using the ultra-fast RNA-seq aligner STAR [20] and subsequently only the

79   fraction of unmapped reads is subjected to metatranscriptomic classification using CLARK-S

80   [21] (see Methods for details). The combination between scalability and accuracy was the

81   main motivation behind choosing these two software packages over competing methods

82   [22,23]. It is important to note that CLARK-S uses a set of uniquely discriminative short

83   sequences at the species level to classify reads. Therefore, reads containing non-

84   discriminative sequences that fail to be uniquely assigned to a single species, e.g. reads

85   originating from the bacterial ribosomal 16S rRNA gene, will be considered 'unclassified'

86   (altogether 8.6% in Fig. 1).

87         The output of CLARK-S is an operational taxonomic units (OTU) count matrix, where

88   rows correspond to viral, bacterial and archeal species and columns to (human) samples.

89   Each entry corresponds to the number of non-human reads classified to the respective

90   species. For convenience, in the following we refer to the set of microbial and viral species

91   profiled using our approach as 'metafeatures'.

92         By screening the study abstracts of the SRA for search terms prioritizing human

93   clinical datasets derived from polyA-independent sequencing protocols (see Methods) we

94   identified over 400 studies relevant to human disease comprising more than 17,000 cDNA

95   libraries (close to 150 terabytes of raw sequencing data). Raw sequencing reads from these

96   studies were downloaded and analyzed using the high performance computing system of the

97   Leibniz Supercomputing Centre (LRZ) of the Bavarian Academy of Sciences and Humanities

98   which facilitated ultra-fast processing with median speeds of 25 and 21 million reads per hour

99   per core per run for the STAR and CLARK-S steps, respectively. Overall, of the total over

100  500 billion RNA-seq reads processed, around 91% could be mapped to the human genome.

101  A fraction of 8.6% of all reads remained non-discriminative at the species level and defined

102  as "unclassified". 0.03%, 0.20% and 0.39% of all reads were assigned to archaeal, bacterial

103  or viral metafeatures, respectively. Despite these relatively low percentages, the absolute

104  numbers of reads classified were in the hundred millions to billions, enabling statistical

105  analyses.

106  **Methods**

107  _High performance computing environment._ Project computations including download,

108  alignment of reads onto the human genome and metafeature quantification were made on

109  the high performance Linux Cluster at the LRZ (www.lrz.de/services/compute/linux-cluster).

110  _RNA-seq data retrieval._ Raw next generation sequencing data were downloaded from the

111  SRA. The R package _SRAdb_ was downloaded on 23 May 2017 and used to query of the

112  SRA database. To identify SRA projects that contain transcriptomic analyses of human RNA-

113  seq data, the SRA attributes 'taxon_id', 'library_source', 'library_strategy', 'platform' were

114  searched for the terms '9606', 'TRANSCRIPT', 'RNA-seq', 'ILLUMINA', respectively. To

115  remove potential bias derived from different sequencing technologies we also restricted the

116  query to SRA runs annotated with 'ILLUMINA' in SRA attribute 'platform'. To exclude studies

117  with insufficient sample size for statistical analysis the query was restricted to SRA projects

118  containing more than five runs. To avoid concentrating the analysis on a small number of large

119  projects the query was restricted to SRA projects with less than 500 runs. To identify studies

120  focusing on phenotypes relevant to human disease, we restricted the query to runs

121  containing at least one or more of the terms 'disease', 'patient', 'primary' and 'clinical' in the

122  SRA attribute 'study_abstract'. To exclude _in vitro_ (cell-culture) experiments, but focus on

123  primary (clinical) samples, SRA runs containing the terms "mutant" or "cell-line" were

124  removed from our selection. Furthermore, SRA runs containing the terms "single cell" and

125  "GTEx" were removed. Finally, samples with less than 1 million total reads or read lengths

126  <50 base pairs were excluded. The described query resulted in 484 Short Read Projects

127  (SRPs) containing a total of 21,659 RNA-seq runs. Due to technical problems (i.e. missing

128  URLs, restricted access) we were unable to download a fraction of 4,078 samples.

129  _Human alignment._ Alignment of reads against the human reference genome (hg38) and

130  simultaneous human gene expression quantification was conducted with STAR (version

131    2.5.2). To increase mapping speed of a large number of samples, we used the --

132    *genomeLoad LoadAndKeep* function to load the STAR index once and keep it in memory for

133    subsequent alignments. The parameter --*quantmode GeneCounts* was used to generate the

134    human gene expression count tables. Unmapped reads were saved with the --

135    *outReadsUnmapped Fastx* parameter. To further increase mapping speed, multiple threads

136    were used as implemented with the parameter --*runThreadN 28*. Runs with less than 30

137    percent reads mapping to the human genome were excluded from downstream analysis. All

138    human alignments were conducted on the LRZ "CoolMUC2" Linux-Cluster. This cluster

139    contains 384 nodes with 64 GB RAM memory and 28 cores each.

140    *Metafeature quantification.* Metafeature quantification was conducted with CLARK-S (version

141    1.2.3). CLARK-S is a software method for fast and accurate sequence classification of

142    metagenomic next-generation sequencing data, including RNA-seq data. One major issue

143    during the classification of metagenomic data is the rising number of targets to align against.

144    CLARK-S solves this issue by building a large index file consisting of discriminative *k*-mers.

145    The metagenomic reference database was generated following the description of the CLARK

146    website using the following two commands: 1) *set_targets.sh bacteria virus --species* and 2)

147    *buildSpacedDB.sh*. This database contained a total of 16,551 genome sequences

148    corresponding to 6,979 unique species (additional file 1). To allow uniform processing,

149    paired-end sequencing experiments were analyzed independently. Each single unmapped

150    reads file was used as input for CLARK-S with the following parameters:

151    *classify_metagenome.sh --spaced –O* list of FASTQ files. To increase classification speed,

152    the CLARK-S express mode was selected and multiple threads were used with parameters --

153    *m 2* and --*n 32*, respectively. The output files of this step contain all input read identifiers with

154    the corresponding metafeature classification. In the subsequent step, total counts are

155    summarized for each feature with the *estimate_abundance.sh* command. To enable

156    comparison across single-end and paired-end experiments, metafeature counts from paired-

157    end experiments were averaged and subsequently rounded to conserve count distribution.

158    To account for varying sequencing depths, metafeature abundance was estimated as the

159     number of reads per million (RPM) total reads sequenced. Metafeature quantification was

160     conducted on the LRZ "Teramem" Linux-Cluster. This cluster contains one node with 6,144

161     GB RAM memory and 96 cores.

162     *BLAST based metafeature classification.* To validate results generated by the MetaMap

163     pipeline, the Basic Local Alignment Search Tool [24] was used as follows. A BLAST

164     database was created from the same genome sequences used in the CLARK-S approach.

165     Then, reads were aligned to this database using BLASTN with a threshold E-value of 1e-10.

166     Produced counts from paired-end experiments were averaged. For each file, BLAST was

167     done by running approximately 10 kilobase chunks (record separator ">") in parallel using

168     GNU parallel (28 jobs), each with 8 threads using one node on the LRZ "CoolMUC3" Linux

169     Cluster. This cluster contains 148 nodes with 96 GB RAM memory and 64 cores each.

170     Output was parsed to exclusively keep reads that could be assigned at the species level.

171     *Differential metafeature abundance.* Differential metafeature abundance analysis was

172     performed using the R package DESeq2 [25]. For each of the four published bona fide dual

173     RNA-seq studies we classified samples into two groups based on the provided annotations:

174     1) Samples expected to contain the known pathogen, such as human papillomavirus positive

175     head and neck tumors in the Zhang et al study, and 2) pathogen-free controls, such as

176     mock-treated cells in the Westermann et al study. Using this binary outcome we performed

177     differential expression analysis across all detected metafeatures. To account for sequencing

178     depth, library size factors were estimated from the total number of sequenced reads. The

179     dispersion for the negative binomial distribution was estimated using a local linear regression

180     as implemented in the *DESeq()* function via the *fitType* parameter 'local'.

181     **Data Validation and quality control**

182     We validated our approach by recovering the ground truth in bona fide dual RNA-seq

183     experiments performed with human cell lines and samples from patients with well-known

184     infection status. Of the four selected studies, one analyzed an infection model based on a

185     bacterial (*Salmonella enterica* serovar Typhimurium) and three based on distinct viral

186     pathogens (Human papillomavirus, Herpes simplex virus, Rhinovirus). As expected,

187    MetaMap detected the known pathogen at higher levels in the respective study compared to

188    the other studies and pathogens (Table 1). Moreover, using the annotation provided in the

189    respective study, we performed differential metafeature abundance analysis to identify those

190    metafeatures that show the largest difference in abundance levels between the infected and

191    control samples. The correct infection agent showed the most significant difference across all

192    metafeatures between infected and control samples for each study (Fig. 2). For example,

193    Westermann et al [26] generated dual RNA-seq data from HeLa cells infected with the

194    enteric bacterial pathogen *Salmonella enterica* serovar Typhimurium and compared them to

195    mock-treated control samples. Accordingly, we here observed *Salmonella enterica* as the

196    most differentially abundant metafeature between the infected and the control samples

197    (P<1e-75, Fig. 2A). Likewise we recovered *Alphapapillomavirus 9*, *Human alphaherpesvirus*

198    *1* (also known as herpes simplex virus 1) and *Rhinovirus A* as the most differentially

199    abundant metafeatures in the data from Zhang et al [27], Rutkowski et al [28] and Bai et al

200    [29], respectively. In the Westermann et al [26] and Rutkowski et al [28] studies, several

201    additional metafeatures showed a strong differential abundance effect (Fig. 2A & C). These

202    metafeatures were closely related to the true infection agent, i.e *Salmonella bongori* (P<1e-

203    67) and *Panine alphaherpesvirus 3* (P<1e-9) for the Westermann et al [26] or Rutkowski et al

204    [28] study, respectively. These findings confirm that our MetaMap pipeline recapitulates

205    results from dedicated dual RNA-seq studies, i.e. studies based on known infectious agents.

206    Therefore, MetaMap may be equally suited to detect previously unknown microbial and viral

207    species in human primary samples.

208

| Study | Infection agent | Total reads | *Salmonella enterica* | Alphapapillomavirus 9 | H. alphaherpesvirus 1 | Rhinovirus A |
|---|---|---|---|---|---|---|
| Westermann et al | *Salmonella enterica* serovar Typhimurium | 1.0e+07 | **6.3e+03** | 1.2e-01 | 1.5e-01 | 1.2e-01 |
| Zhang et al | Human papillomavirus | 4.6e+07 | 3.0e-02 | **5.1e+01** | 2.2e-02 | 2.2e-02 |
| Rutkowski et al | Herpes simplex virus | 3.5e+07 | 1.1e+00 | 3.1e-02 | **3.1e+04** | 3.0e-02 |

| Bai et al | Rhinovirus | 6.6e+06 | 2.0e-01 | 1.5e-01 | 1.5e-01 | **4.4e+01** |
|-----------|-----------|---------|---------|---------|---------|---------|

209    Table 1. Overview of four dual RNA-seq studies used to validate the MetaMap pipeline. Total
210    reads column depicts the average read depth per sample for each study. Average metafeature
211    abundance for *Alphapapillomavirus 9*, *Salmonella enterica*, *Human alphaherpesvirus 1* and *Rhinovirus*
212    *A* are shown in RPM. The correct infection agent for the respective study is highlighted in bold font.
213
214    As an additional control, we re-analysed two projects contained in our data collection

215    that are derived from the B lymphoblast cell line, under non-infectious conditions. However,

216    since Epstein-Barr virus is used for transfection and transformation of lymphocytes to

217    lymphoblasts, we expected to detect reads from this virus in these projects [30], but no

218    further viral or microbial reads [31]. Indeed the most abundant metafeatures in each project

219    were dominated by reads classified to *Gammaherpesvirus 4* (also known as Epstein-Barr

220    virus, EBV) and *Enterobacteria phage phiX174 sensu lato* (phiX), commonly used as spike-in

221    in Illumina sequencing runs [32] (Fig. 3A-B). On average 95% and 97% of all metafeature

222    reads were classified as phiX or EBV for projects SRP041338 and SRP091453, respectively

223    (Fig. 3C). Conversely, the abundance of reads mapping to bacterial species for these two

224    projects corresponds to the bottom percentile as compared to all other projects in the

225    MetaMap database, supporting sterility of this cell line (Fig. 3D). This demonstrates that

226    MetaMap not only is capable of re-discovering known pathogenic species (true positives) in

227    controlled infection experiments (Fig. 2), but it also minimizes the detection of false positives

228    or at least, provides measures such as abundance and significance allowing the user to

229    identify and counterselect those species.

230    As a technical validation, we compared our approach to an alternative

231    metatranscriptomic classification strategy for the Westermann et al [33] study. All non-human

232    reads were aligned using BLASTN to a BLAST database consisting of the same genomic

233    sequences used by CLARK-S (see Methods for details). The average metafeature

234    abundances across all 42 samples derived from the BLAST based approach and CLARK-S

235    correlated significantly (Spearman correlation, Rho: 0.16, P: 3.1e-10) (Fig. 4A). BLAST

236    showed higher sensitivity and detected more metafeatures compared to CLARK-S (indicated

237    by the accumulation of dots at value 0 on the X-axis in Fig. 4A). This is mostly observed for

238    low abundance metafeatures which could represent low counts derived from sequencing

239    and/or mapping errors. However, most importantly the true pathogen metafeature

240    '*Salmonella enterica*' showed very high correlation across samples between the BLAST and

241    CLARK-based abundance estimates (Fig. 4B). Noteworthy, the MetaMap pipeline processed

242    reads more than three orders of magnitude faster than BLAST, demonstrating a significant

243    speed advantage while generating comparable results (Fig. 4C).
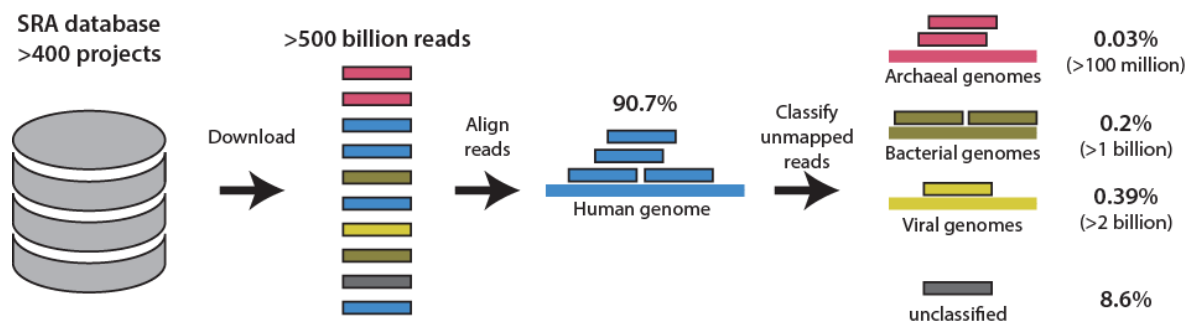
244    **Re-use potential**

245         Microbial and viral contamination in next-generation sequencing data was observed

246    before. It can be caused by incorrect mapping due to sequence similarity between different

247    species [34,35]. To minimize such effects, we encourage focussing on studies including

248    intra-project comparisons, such as exemplified in the differential metafeature abundance

249    analysis. Contaminating agents should affect all runs within a project to the same extent and

250    therefore not show a condition-specific effect. Alternatively, these "contaminations" might

251    actually reflect true biological factors. For example, in the Westermann et al study [33] we

252    detected substantial levels of phiX in both conditions (infected samples and mock-treated

253    controls), but only the '*Salmonella*' metafeature showed a condition-specific effect.

254         All the raw data described in the present study were publicly available before, yet

255    have been very cumbersome to extract individually. The presented MetaMap database now

256    makes these data easily accessible for a very broad community, thereby allowing for global

257    comparisons over hundreds of individual studies and thousands of sampled conditions. While

258    we attempted to minimize the risk of detecting false positives (Fig. 3), it should be noted that

259    not all metafeatures classified by MetaMap will necessarily refer to true biological factors.

260    Rather our pipeline provides the user with a scientific starting ground to validate the

261    presence/absence of defined microbial and viral species under defined conditions and

262    explore the underlying biology and significance in greater detail. As a potential use case of

263    these data, users can test for associations of microbial or viral metafeatures with a plethora

264    of human diseases, or between themselves. In addition, users with interest in a specific

265    bacterial or viral species can easily identify studies, and consequently disease contexts, in

266    which reads from this organism were detected. This could give an important first hint to

267    assess whether the respective species might be implicated in a given human disease

268    etiology. Furthermore, this resource provides the opportunity to validate findings derived from

269    standard microbiome profiling technologies, such as 16S rRNA gene based or shotgun

270    metagenomics [36]. Finally, metafeature detection in human clinical RNA-seq samples may

271    provide a critical advantage when studying microbes or viruses which are challenging to

272    isolate.

273    All generated metafeature OTU count tables from 17,278 cDNA libraries from 436 SRA

274    projects including annotation are provided for download. The MetaMap pipeline can be

275    accessed via the protocols.io website with digital object identifier

276    dx.doi.org/10.17504/protocols.io.msec6be.

277 **Figures**



278
279 Figure 1. Schematic illustrates the MetaMap pipeline. Over 400 projects from studies relevant to

280 human disease were identified in the SRA database. Over 500 billion RNA-seq reads were

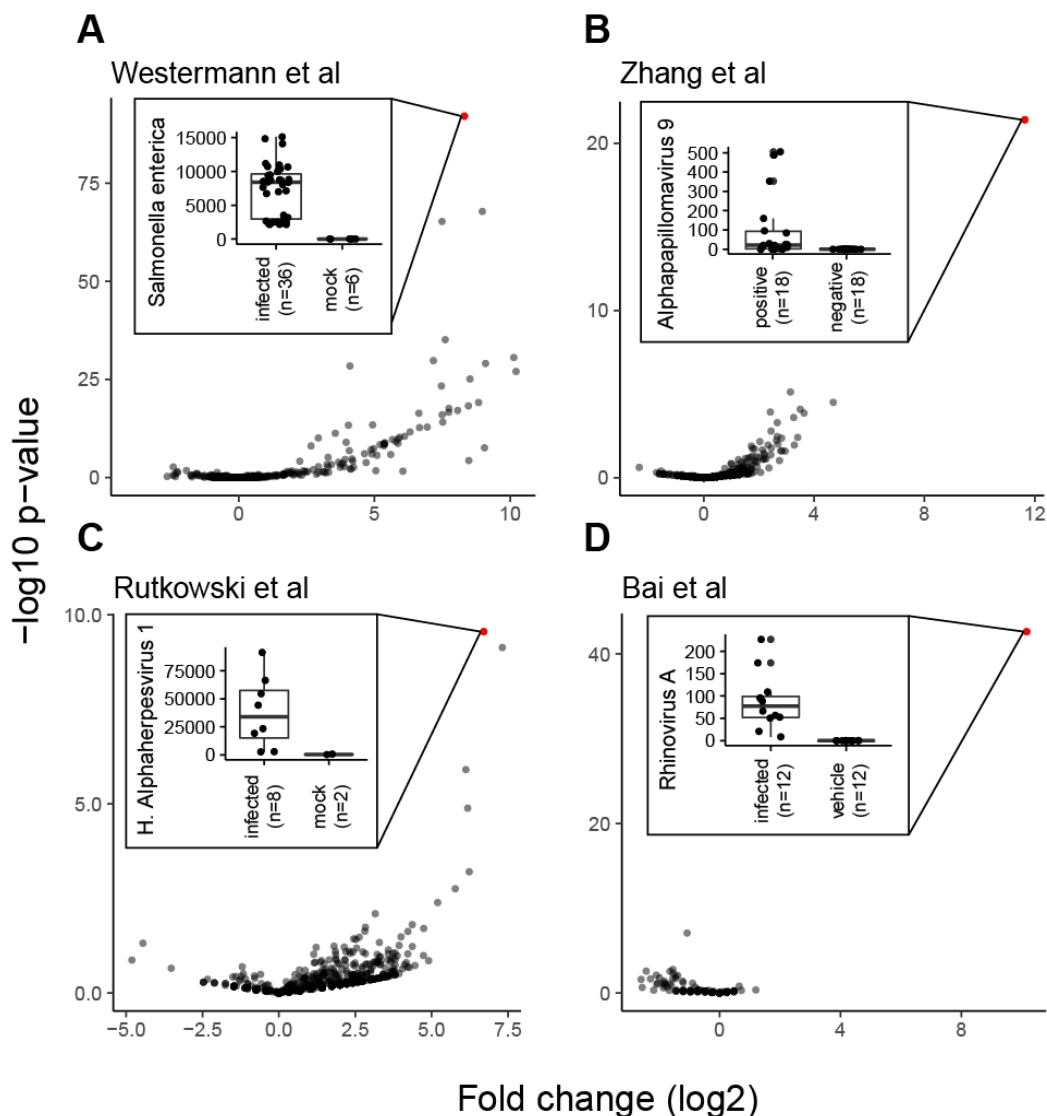281 downloaded and first filtered by mapping them onto the human genome and subsequently the

282 remaining reads underwent metafeature classification. 90.7% of all reads mapped to the human

283 genome. 0.03%, 0.20% and 0.39% of all reads were assigned to archaeal, bacterial or viral
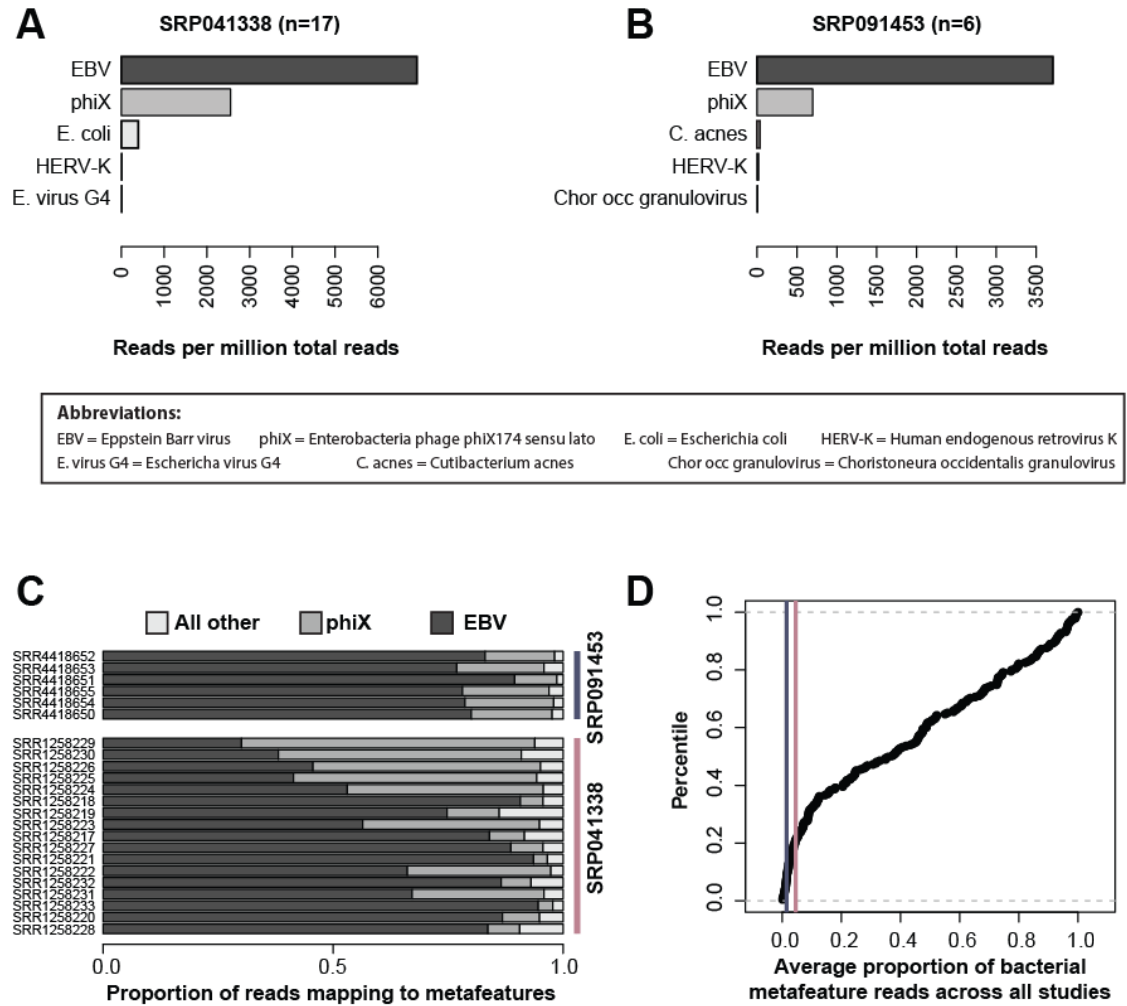
284 metafeatures, respectively. 8.6% of all reads remain non-discriminative at the species level

285 ('unclassified').

286

Figure 2. Differential metafeature abundance analysis of controlled infection experiments recovers ground truth. Panels A-D depict "volcano" plots showing fold change and inverted p-value on the X and Y axes, respectively. Each dot represents a metafeature. The most significant metafeature is colored in red. Insets display boxplots of the abundance levels in RPM of the top hit metafeature across conditions for each study. For all boxplots, the box represents the interquartile range, the horizontal line in the box is the median, and the whiskers represent 1.5 times the interquartile range.

295

Figure 3. Analysis of lymphoblast cell line experiments further supports the MetaMap pipeline. Panels

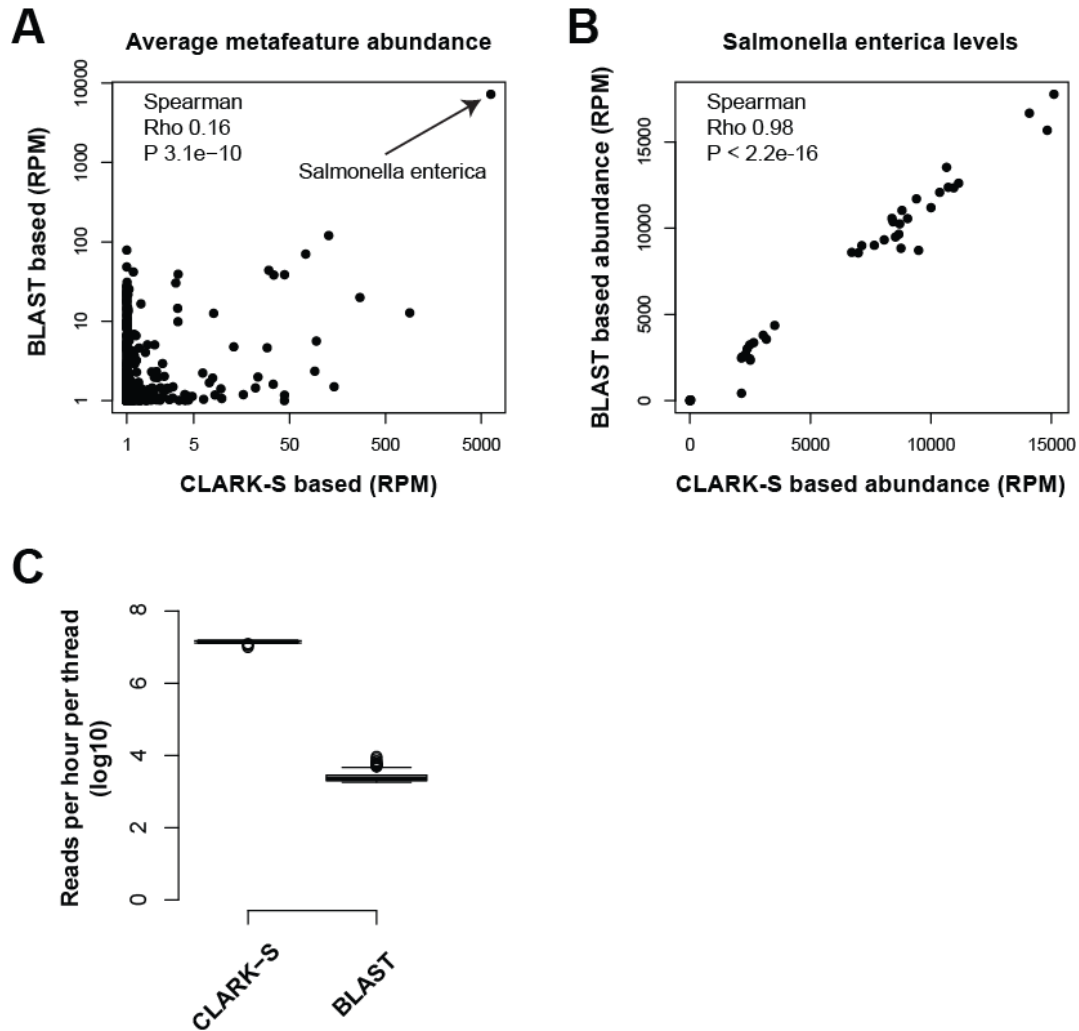A and B depict mean abundance levels across all samples of the top five metafeatures for projects

SRP041338 and SRP091453, respectively. Panel C shows relative proportion of reads mapping to

EBV, phiX and all other metafeatures across RNA-seq samples. Panel D depicts the cumulative

distribution plot of the average proportion of bacterial metafeature reads across all projects. Purple

and pink vertical lines highlight projects SRP041338 and SRP091453, respectively.

302

303

304  Figure 4. Alternative BLAST-based classification method validates metafeature abundance estimates

305  by MetaMap. Panel A depicts average metafeature RPM levels derived using the CLARK-S software,

306  as implemented in the MetaMap pipeline, and a BLAST-based alternative approach on the X- and Y-

307  axes, respectively. Panel B shows the correlation in *Salmonella enterica* abundance levels between

308  the two classification approaches. Panel C shows the difference in classification speed between the

309  BLAST and CLARK-S metatranscriptomic classification. Y axis shows the number of reads processed

310  per hour per thread in log10 space.

311

312  **Acknowledgments**

315

## References

1. Young VB. The role of the microbiome in human health and disease: an introduction for clinicians. BMJ. 2017;356: j831.

2. Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JI. An obesity-associated gut microbiome with increased capacity for energy harvest. Nature. 2006;444: 1027–1031.

3. Henao-Mejia J, Elinav E, Jin C, Hao L, Mehal WZ, Strowig T, et al. Inflammasome-mediated dysbiosis regulates progression of NAFLD and obesity. Nature. 2012;482: 179–185.

4. Cani PD, Bibiloni R, Knauf C, Waget A, Neyrinck AM, Delzenne NM, et al. Changes in gut microbiota control metabolic endotoxemia-induced inflammation in high-fat diet-induced obesity and diabetes in mice. Diabetes. 2008;57: 1470–1481.

5. Wang Z, Klipfell E, Bennett BJ, Koeth R, Levison BS, Dugar B, et al. Gut flora metabolism of phosphatidylcholine promotes cardiovascular disease. Nature. 2011;472: 57–63.

6. Engel M, Endesfelder D, Schloter-Hai B, Kublik S, Granitsiotis MS, Boschetto P, et al. Influence of lung CT changes in chronic obstructive pulmonary disease (COPD) on the human lung microbiome. PLoS One. 2017;12: e0180859.

7. Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, et al. Genomic analysis identifies association of Fusobacterium with colorectal carcinoma. Genome Res. 2012;22: 292–298.

8. Castellarin M, Warren RL, Freeman JD, Dreolini L, Krzywinski M, Strauss J, et al. Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma. Genome Res. 2012;22: 299–306.

9. Kodama Y, Shumway M, Leinonen R, International Nucleotide Sequence Database Collaboration. The Sequence Read Archive: explosive growth of sequencing data. Nucleic Acids Res. 2012;40: D54–6.

10. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, et al. A survey of best practices for RNA-seq data analysis. Genome Biol. 2016;17: 13.

11. Gouin A, Legeai F, Nouhaud P, Whibley A, Simon J-C, Lemaitre C. Whole-genome re-sequencing of non-model organisms: lessons from unmapped reads. Heredity . 2015;114: 494–501.

12. Peng X, Wang J, Zhang Z, Xiao Q, Li M, Pan Y. Re-alignment of the unmapped reads with base quality score. BMC Bioinformatics. 2015;16 Suppl 5: S8.

13. Westermann AJ, Gorski SA, Vogel J. Dual RNA-seq of pathogen and host. Nat Rev Microbiol. 2012;10: 618–630.

14. Westermann AJ, Barquist L, Vogel J. Resolving host-pathogen interactions by dual RNA-seq. PLoS Pathog. 2017;13: e1006033.

15. Juranic Lisnic V, Babic Cac M, Lisnic B, Trsan T, Mefferd A, Das Mukhopadhyay C, et al. Dual analysis of the murine cytomegalovirus and host cell transcriptomes reveal new aspects of the virus-host cell interface. PLoS Pathog. 2013;9: e1003611.

16. Xu G, Strong MJ, Lacey MR, Baribault C, Flemington EK, Taylor CM. RNA CoMPASS: a dual approach for pathogen and host transcriptome analysis of RNA-seq datasets. PLoS One. 2014;9: e89445.

17. Park S-J, Kumar M, Kwon H-I, Seong R-K, Han K, Song J-M, et al. Dynamic changes in host gene expression associated with H5N8 avian influenza virus infection in mice. Sci Rep. 2015;5: 16512.

18. Saxena K, Simon LM, Zeng X-L, Blutt SE, Crawford SE, Sastri NP, et al. A paradox of transcriptional and functional innate interferon responses of human intestinal enteroids to enteric virus infection. Proceedings of the National Academy of Sciences. 2017;114: E570–E579.

19. Wesolowska-Andersen A, Everman JL, Davidson R, Rios C, Herrin R, Eng C, et al. Dual RNA-seq reveals viral infections in asthmatic children without respiratory illness which are associated with changes in the airway transcriptome. Genome Biol. 2017;18: 12.

20. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics. 2012;29: 15–21.

21. Ounit R, Lonardi S. Higher classification sensitivity of short metagenomic reads with CLARK-S. Bioinformatics. 2016;32: 3823–3825.

22. Lindgreen S, Adair KL, Gardner PP. An evaluation of the accuracy and speed of metagenome analysis tools. Sci Rep. 2016;6: 19233.

23. Engström PG, Steijger T, Sipos B, Grant GR, Kahles A, Rätsch G, et al. Systematic evaluation of spliced alignment programs for RNA-seq data. Nat Methods. 2013;10: 1185–1191.

24. Altschul S. Basic Local Alignment Search Tool. J Mol Biol. 1990;215: 403–410.

25. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2 [Internet]. 2014. doi:10.1101/002832

26. Westermann AJ, Förstner KU, Amman F, Barquist L, Chao Y, Schulte LN, et al. Dual RNA-seq unveils noncoding RNA functions in host–pathogen interactions. Nature. 2016;529: 496–501.

27. Zhang Y, Koneva LA, Virani S, Arthur AE, Virani A, Hall PB, et al. Subtypes of HPV-Positive Head and Neck Cancers Are Associated with HPV Characteristics, Copy Number Alterations, PIK3CA Mutation, and Pathway Signatures. Clin Cancer Res. 2016;22: 4735–4745.

28. Rutkowski AJ, Erhard F, L'Hernault A, Bonfert T, Schilhabel M, Crump C, et al. Widespread disruption of host transcription termination in HSV-1 infection. Nat Commun. 2015;6: 7126.

bioRxiv preprint doi: https://doi.org/10.1101/269092; this version posted February 22, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

29. Bai J, Smock SL, Jackson GR Jr, MacIsaac KD, Huang Y, Mankus C, et al. Phenotypic responses of differentiated asthmatic human airway epithelial cultures to rhinovirus. PLoS One. 2015;10: e0118286.

30. Santpere G, Darre F, Blanco S, Alcami A, Villoslada P, Mar Albà M, et al. Genome-Wide Analysis of Wild-Type Epstein–Barr Virus Genomes Derived from Healthy Individuals of the 1000 Genomes Project. Genome Biol Evol. 2014;6: 846–860.

31. Mangul S, Olde Loohuis LM, Ori A, Jospin G, Koslicki D, Yang HT, et al. Total RNA Sequencing reveals microbial communities in human blood and disease specific effects [Internet]. 2016. doi:10.1101/057570

32. Mukherjee S, Huntemann M, Ivanova N, Kyrpides NC, Pati A. Large-scale contamination of microbial isolate genomes by Illumina PhiX control. Stand Genomic Sci. 2015;10: 18.

33. Westermann AJ, Förstner KU, Amman F, Barquist L, Chao Y, Schulte LN, et al. Dual RNA-seq unveils noncoding RNA functions in host-pathogen interactions. Nature. 2016;529: 496–501.

34. Strong MJ, Xu G, Morici L, Splinter Bon-Durant S, Baddoo M, Lin Z, et al. Microbial contamination in next generation sequencing: implications for sequence-based analysis of clinical samples. PLoS Pathog. 2014;10: e1004437.

35. Bonfert T, Csaba G, Zimmer R, Friedel CC. Mining RNA–Seq Data for Infections and Contaminations. PLoS One. 2013;8: e73071.

36. Cox MJ, W O C, Moffatt MF. Sequencing the human microbiome in health and disease. Hum Mol Genet. 2013;22: R88–R94.