1

# Application of a bioinformatic pipeline to RNA-seq data identifies novel virus-like sequence in human blood

Marko Melnick[*], Patrick Gonzales[*], Thomas J. LaRocca[†], Robin D. Dowell[‡],
Yuping Song[**], Joanne Wuu[‡‡], Michael Benatar[‡‡], Björn Oskarsson[§§], Leonard Petrucelli[**,††],
Christopher D. Link[*,§] Mercedes Prudencio[**,††]

[*] Integrative Physiology, University of Colorado, Boulder, Colorado, 80303, USA

[†] Department of Health and Exercise Science, Center for Healthy Aging, Colorado State University, Fort Collins, Colorado, 80523, USA

[‡] BioFrontiers Institute and Department of Molecular, Cellular and Developmental Biology, University of Colorado, Boulder, Colorado, 80303, USA

[§] Institute for Behavioral Genetics, University of Colorado, Boulder, Colorado, 80303, USA

[**] Department of Neuroscience, Mayo Clinic, 4500 San Pablo Road, Jacksonville, Florida, 32224, USA

[††] Neuroscience Graduate Program, Mayo Clinic Graduate School of Biomedical Sciences, Jacksonville, Florida, 32224, USA

[‡‡] Department of Neurology, University of Miami, Miami, Florida, 33136, USA

[§§] Department of Neurology, Mayo Clinic, 4500 San Pablo Road, Jacksonville FL, 32224, USA

44

# Application of a bioinformatic pipeline to RNA-seq data identifies novel virus-like sequence in human blood

47

48 Keywords: ALS, Transcriptomics, RNA-seq, Microbiome, Virome

49

50 Corresponding author
51 Phone: 303-735-5112
52 Department of Integrative Physiology
53 354 UCB
54 Boulder Colorado, 80303
55 Email: Marko.Melnick@colorado.edu

56

57 Mercedes Prudencio
58 4500 San Pablo Road,
59 Griffin Building Rm 221
60 Jacksonville, FL 32224
61 Phone: 904-953-6638; Fax: 904-953-7370
62 Email: Prudencio.Mercedes@Mayo.edu

63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87

# ABSTRACT

Numerous reports have suggested that infectious agents could play a role in neurodegenerative diseases, but specific etiological agents have not been convincingly demonstrated. To search for candidate agents in an unbiased fashion, we have developed a bioinformatic pipeline that identifies microbial sequences in mammalian RNA-seq data, including sequences with no significant nucleotide similarity hits in GenBank. Effectiveness of the pipeline was tested using publicly available RNA-seq data. We then applied this pipeline to a novel RNA-seq dataset generated from a cohort of 120 samples from amyotrophic lateral sclerosis (ALS) patients and controls, and identified sequences corresponding to known bacteria and viruses, as well as novel virus-like sequences. The presence of these novel virus-like sequences, which were identified in subsets of both patients and controls, were confirmed by quantitative RT-PCR. We believe this pipeline will be a useful tool for the identification of potential etiological agents in the many RNA-seq data sets currently being generated.

# INTRODUCTION

**Background of organisms in neurodegeneration**

Infection has been proposed to play a role in multiple neurodegenerative diseases[1], including amyotrophic lateral sclerosis (ALS)[2]. ALS is the most common motor neuron disease in adults, with the majority of individuals dying within 3-5 years of symptom onset. The disease is defined by the degeneration and death of motor neurons in the brain and spinal cord, resulting in progressive weakness and eventually death, typically from respiratory muscle weakness[3]. Around 5-10% of ALS cases are inherited, termed familial ALS (fALS), with the remaining cases considered sporadic ALS (sALS). After decades of study, the etiology of sALS remains a mystery, although suspected risk factors for ALS include exposure to heavy metals, pesticides, chemical solvents, cigarette smoke, and unidentified factors related to US military service[4–7].

Along with these environmental risk factors, there has been a long history, with variable success, in the search for pathogens that might contribute to ALS[8–12] and other neurodegenerative diseases such as Alzheimer's disease (AD)[13–15], Parkinson's disease (PD)[16–18], and multiple sclerosis (MS)[19].

Diverse pathogens have been reported in the blood, cerebrospinal fluid (CSF) and central nervous system (CNS) from ALS patients. For example, bacteria that have been detected include *Cutibacterium acnes, Corynebacterium sp, Fusobacterium nucleatum, Lawsonella clevelandesis,* and *Streptococcus thermophilus* in CSF[20], and mycoplasma in blood[21]. Fungi, including *Candida famata*, *Candida albicans*, *Candida parapsilosis*, *Candida glabrata*, and *Penicillium notatum,* have been detected in CSF, while *Malassezia globosa*, *Cryptococcus neoformans*[11], and *Candida albicans* have been found in various regions of the CNS[11,22,23]. The search for viruses that contribute to ALS pathology is much more extensive and includes studies on herpes virus[9,24], enterovirus[9,25–28], human immunodeficiency virus (HIV)[29,30], and human endogenous retrovirus (HERV-K)[31–33]. Importantly, multiple studies using immunohistochemistry have shown an

134 increased load of various pathogens in ALS samples compared to controls in multiple tissues
135 suggesting these pathogens are present and cannot be simply attributed to contamination[9,11,20,22,23].
136 Ultimately, the presence of ALS dysbiosis is unresolved and remains an active area of investigation
137 with evidence for[34–38] and against[39] it.
138 　　　The biological role that these alternative microbiotas play in ALS is also unclear. ALS
139 patients may have a compromised blood brain barrier (BBB) or blood spinal cord barrier (BSCB)
140 function[40,41]. It has been reported that ALS patients also have elevated Gram negative
141 endotoxin/lipopolysaccharide (LPS) in the blood[42]. Patients with ALS also display activation of
142 the innate immune system along with changes in blood[43,44], spinal cord and motor neurons[45], but
143 if and how bacteria might influence activation is an active area of research. A "dual hit" hypothesis
144 by Correia et al. suggests inflammation via LPS may contribute to mis-localization and
145 aggregation of ALS-implicated protein TAR DNA-binding protein 43 (TDP-43)[46].
146 　　　Numerous studies have looked for biomarkers of ALS[47] using metabolomics[48,49],
147 neuroinflammation[50,51], DNA methylation[52,53], gene expression[54], microRNA expression[55,56] and
148 our previous study which analyzed protein levels of poly(GP) in *C9ORF72*-associated ALS
149 (c9ALS)[57]. The search for pathogens using sequencing data from blood samples in ALS patients
150 has been conducted before[58–61], but previous efforts have not looked for novel pathogens.
151 　　　Next-generation sequencing (NGS) technologies have shown broad detection of pathogens
152 in a target-independent unbiased fashion[62–65]. However, designing a microbial detection
153 experiment is non-trivial considering the variety of methods[66] and algorithms[67] that can be applied.
154 Our primary goal when designing a new pipeline was to be conservative and unbiased with regards
155 to discovery of novel pathogens. Furthermore, we wanted our pipeline to allow for the
156 quantification of both novel and known pathogens. While other pipelines have used reads that do
157 not map to the host genome (unmapped reads) for microbial identification and quantification, these
158 pipelines cannot be used for discovery as they rely on existing databases of microbial genomes[68–
159 71]. Thus, we opted for de-novo assembly of unmapped reads into contigs, followed by alignment
160 of unmapped reads back to these contigs for quantification. A similar pipeline known as IMSA[72]
161 uses this strategy, but disregards contigs that might be identified by translated amino acid sequence
162 similarity using BLASTX (a set we call the "dark biome") as well as subsequent contigs with no
163 BLASTN or BLASTX hit (a set we call the "double dark biome").
164 　　　We validated our pipeline by using datasets with known bacterial or viral infections. We
165 also examined the differences in microbial identification between polyA and total RNA recovery
166 in multiple tissues, and investigated the effects of globin depletion of blood samples. We then used
167 our pipeline on a novel blood dataset (termed "Our Study") as well as on five other published ALS
168 datasets from blood or spinal cord samples, analyzed each dataset individually, and analyzed
169 across datasets for changes in microbiota. While we did not identify any microbes enriched in the
170 blood of ALS patients, we did identify and validate a novel virus-like sequence, demonstrating the
171 potential of the bioinformatic pipeline we have established.
172
173
174
175 **MATERIALS AND METHODS**
176
177 **Blood Collection and RNA Extraction**
178 　　　A total of 120 RNA whole blood samples constitute Our Study, which included 30 healthy
179 controls (from general population that do not have blood relatives suffering from ALS, CTL), 30

180  pre-symptomatic *C9ORF72* mutant carriers (C9A), 30 symptomatic *C9ORF72* ALS cases (C9S),
181  and 30 symptomatic *C9ORF72*-negative ALS cases (SYM). PAXgene blood RNA tubes were
182  collected at Mayo Clinic Jacksonville and at University of Miami.  All 120 RNA samples selected
183  for RNA-seq were received and processed at Mayo Clinic Jacksonville using PAXgene blood RNA
184  kit following manufacturer's recommendations (Qiagen). Blood RNA was of high quality,
185  assessed in an Agilent Bioanalyzer (Agilent), with RNA integrity values ranging from 7.4 to 9.8,
186  with a median value of 8.7. RNA samples were then sent to The Jackson Laboratory for globin
187  depletion, library preparation and sequencing of total blood RNA.
188
189  **Globin Depletion**
190      Due to the abundance of large haemoglobin RNA transcripts present in the blood, a globin
191  depletion step, using the Ambion GLOBINclear kit (AM1980), was performed before sequencing
192  of the blood RNA samples in order maximize coverage on non-globin genes. In brief, one
193  microgram of total RNA was used as starting material, and specific biotinylated oligos were used
194  to capture globin mRNA transcripts. The capture oligos were hybridized with total RNA samples
195  at 50°C for 30 min. Streptavidin magnetic beads were then used to bind to the biotinylated capture
196  oligos hybridized to globin mRNA by incubating at 50°C for 30 min. The magnetic streptavidin
197  beads-biotin complex were then captured to the side of the tubes by a magnet, and the resulting
198  supernatant is free of globin mRNA. The globin depleted RNA was further purified by RNA
199  binding beads and finally eluted in elution buffer. The resulting RNA free of >95% globin mRNA
200  transcripts was then processed for next generation sequencing. Of note, to assess the efficiency of
201  the globin RNA depletion, 10% of the samples processed were selected randomly and amplified
202  using a Target-Amp Nano labeling kit (Epicentre). Samples were normalized to 100 ng input and
203  reverse transcribed. First strand cDNA was generated by incubating at 50°C for 30 min with first
204  strand premix and Superscript III. This was followed by second strand cDNA synthesis through
205  DNA polymerase by incubating at 65°C for 10 min and at 80°C for 3 min. In-vitro transcription
206  was then performed at 42°C for 4 hours followed by purification using RNeasy mini kit (Qiagen).
207      Due to the large number of samples, the globin depletion step was performed in two batches.
208  We provided guidelines on how samples would be divided among the batches and also for how
209  samples would be grouped in the sequencing runs in order to minimize technical variability. The
210  Jackson Laboratory personnel were blinded to the identity of the samples.
211      RNA-seq of total blood RNA (globin and ribosomal RNA depleted) was performed in an
212  Illumina HiSeq4000 with >70 million read pairs per sample. Raw reads were then sent back to us
213  for bioinformatics analyses.
214
215  **Quantitative RT-PCR for blood RNA samples**

216      A total of 500 ng of total blood RNA was used for reverse transcription polymerase chain
217  reaction (RT-PCR), using the High Capacity cDNA Transcription Kit with random primers
218  (Applied Biosystems). Quantitative real-time PCR (qRT-PCR) was performed using SYBR
219  GreenER qPCR SuperMix (Invitrogen). Samples were run in triplicate, and qRT-PCRs were run
220  on a QuantStudio 7 Flex Real-Time system (Applied Biosystems).
221
222  List of primers and their sequences:
223  *RDRP* forward 5'-GCTGTCAAATCGGTTTCCAAC-3';
224  *RDRP* reverse 5'-CTGCCTTCGTCATCTTGGAG-3';
225  *GAPDH* forward 5'-GTTCGACAGTCAGCCGCATC-3';

226    *GAPDH* reverse 5'-GGAATTTGCCATGGGTGGA-3'.

227

228    **Transcriptomics**
229    See pipeline description in results for an overview of the pipeline; see bioinformatics
230    supplement File S1 for a more detailed description of the analysis pipeline, versions, and statistical
231    quantification. All data in this study was processed identically using the pipeline.

232

233

234    **Statistical Analysis**
235    To assess statistical differences between conditions, a two tailed Student's *t*-test is
236    calculated using normalized coverage for contigs or binned normalized coverage for
237    species/genus, etc. The number of contigs or genus/species is used to obtain an adjusted p-value
238    using scipy in Python. Cutoff for statistical significance is less than an adjusted p-value of 0.05
239    unless otherwise stated.

240

241

242    **Data availability**
243    Raw sequencing data for Our Study dataset is available in the NCBI Sequence Read
244    Archive under the accession number (PRJN). All other datasets are publicly available and all of
245    the code used in this manuscript is available at https://github.com/Senorelegans/MysteryMiner.
246    Supplemental material available at figshare: https:// doi.org/(INSERT).

247

248

249

250

251    **RESULTS**
252    **Pipeline description**
253    Mystery Miner is written as a Nextflow pipeline. Below is a short overview of the Mystery
254    Miner pipeline (Fig1). A more in-depth explanation, list of software and versions used, and typical
255    parameters of each step is described in the bioinformatics supplement, and all of the code used in
256    this manuscript can be found at https://github.com/Senorelegans/MysteryMiner.
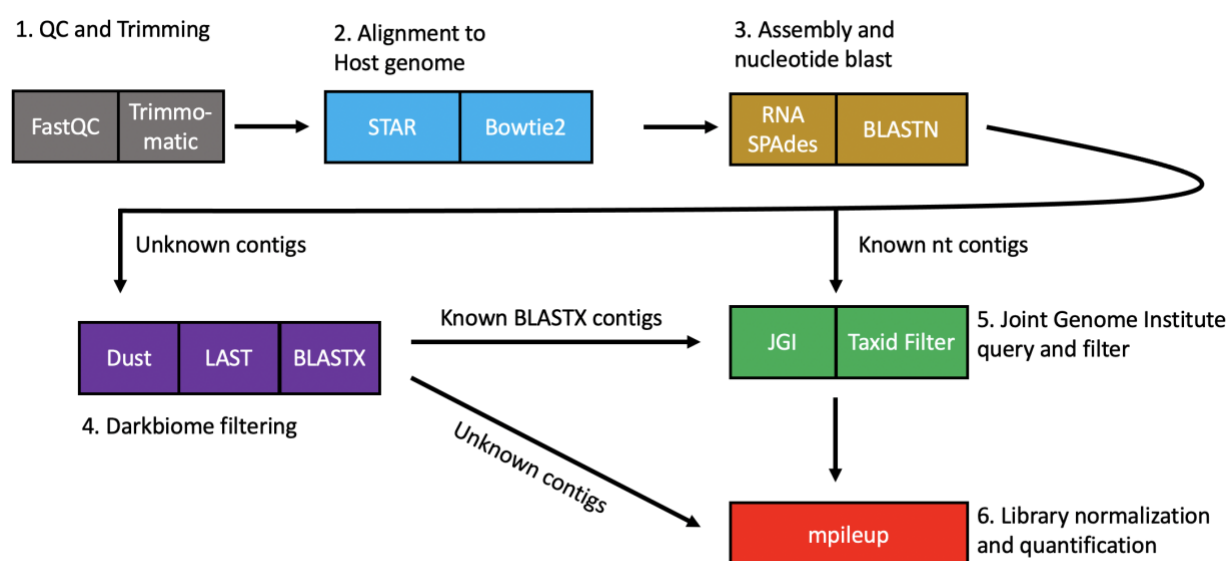257    Raw reads were first checked for quality using FastQC then trimmed to remove both
258    adaptor contamination and low quality basecalls using Trimmomatic. Trimmed reads were then
259    mapped to the host genome using multiple alignment algorithms in series (STAR, Bowtie2) and
260    unmapped reads were retained for contig assembly. Filtering out host reads made downstream
261    assembly faster and required less memory. We assembled contigs from unmapped reads with the
262    SPAdes assembler (with "-rna" setting). This assembler was chosen for its recent use in studies of
263    microbial diversity[73] and proven robustness to biological and technical variation[74]. The species
264    each contig belongs to was identified with BLASTN using default settings, and the top hit for each
265    contig was retained (a set we call "regular biome"). Contigs with no BLASTN hits were then
266    filtered to remove highly repetitive regions (DUST) and retained if they had a greater than 60%
267    pairwise alignment (LAST) between contigs assembled from a single sample, group/condition, or
268    all samples.
269    We then identified contigs that lacked detectable nucleotide similarity to any GenBank
270    entry but showed similarity at the amino acid level using BLASTX ("dark biome"). Furthermore,
271    contigs with no BLASTN or BLASTX hits were labelled as "double dark biome" (also filtered by

272 DUST and LAST). Every contig in the "regular biome" and "dark biome" were then queried
273 against the Joint Genome Institute Server for additional taxonomic information. As Mystery Miner
274 is an agnostic tool, it cannot distinguish between true tissue or cell-associated microbes and
275 experimentally introduced contamination.
276      For quantification, we mapped the non-host reads using Bowtie2 to the contigs obtained
277 from SPAdes. Next, we mapped reads to contigs using samtools mpileup (default mapq score) to
278 calculate the amount of reads over each base pair in a contig. We then calculated coverage on the
279 contigs by summing all of the counts for each base pair in a contig and dividing by the length of
280 the contig. We then calculated normalized coverage by library size using the number of mapped
281 reads to the host genome. This gave us normalized coverage (NC) for a contig or binned
282 normalized coverage (BNC) for multiple contigs within a species/genus, etc. To assess statistical
283 differences between conditions, a Student's $t$-test was calculated through NC or BNC, using the
284 number of contigs or genus/species to obtain an adjusted p-value using scipy in Python.
285
286
287



288
289
290 **Figure 1. Diagram of Mystery Miner Pipeline**
291 Reads were first checked with FastQC and trimmed using Trimmomatic (1. grey). Reads were then
292 aligned to the host genome using various aligners (2. blue). Non-host (unmapped) reads were
293 assembled into contigs with RNA SPAdes and regular biome contigs were identified with
294 BLASTN (3. yellow). Unidentified contigs were filtered for repetitive sequences with Dust, filter
295 by single, group or all with LAST, and dark biome contigs were identified with BLASTX. Double
296 dark biome unidentified BLASTX contigs were sent directly to quantification (4. purple). Dark
297 biome and regular biome contigs were assigned complete taxonomy using the JGI server and
298 filtered one last time to remove mammalian/host genome contigs (5. Green). Non-host reads were
299 then mapped to all contigs and normalized coverage was calculated for subsequent statistical
300 analysis.
301
302

**Validating Mystery Miner on datasets with known bacterial or viral infection**

To confirm that Mystery Miner is able to recover and quantify known bacterial infections from sequencing data, we utilized an *in vitro* model of *Chlamydia trachomatis* infection (Humphrys 2016)[75]. In this study, epithelial cell monolayers were infected with *Chlamydia trachomatis*; and polyA RNA (6 samples) and total RNA (6 samples) were sequenced 1 hour and 24 hours post infection (hpi). Using the Mystery Miner pipeline, out of $5.32 \times 10^6$ reads from all of the samples, $6.04 \times 10^5$ reads remained unmapped (~11.34%) after trimming and mapping to the host genome (File S2). From these non-host reads, 3,257 contigs were assembled and 1,199 of these contigs were identified as regular biome (File S3). An additional 27 contigs had no BLASTN hit. Of these, we identified 2 dark biome (BLASTX identified) and no double dark biome (no BLASTX or BLASTN hit) contigs (File S4 and File S5).

Using Mystery Miner we successfully identified, and found significantly elevated levels, of *Chlamydia trachomatis* (BNC by species) in 24 hours post infection (hpi) samples compared to 1 hpi samples in both polyA (Padj = 0.004) and total RNA (Padj = 0.0005). In addition to *Chlamydia trachomatis,* we identified 6 additional bacterial species and one viral species (Alphapapillomavirus 7) in the samples (Figure 2A), including significantly elevated levels of *Mycoplasma hyorhinis* contigs in total RNA samples. No significant differences were observed in the dark or double dark contigs (File S6).

To confirm that the pipeline can detect known viral infections, we ran Mystery Miner on a total RNA dataset from an *in vitro* model of severe acute respiratory syndrome coronavirus (SARS-CoV) 1 or 2 infection (Emanuel2020[76]). In this study human epithelial Calu3 cells were infected with SARS-CoV-1 or SARS-CoV-2 (4, 12, or 24 hours), mock (4 or 24 hours), or untreated (4 hours).
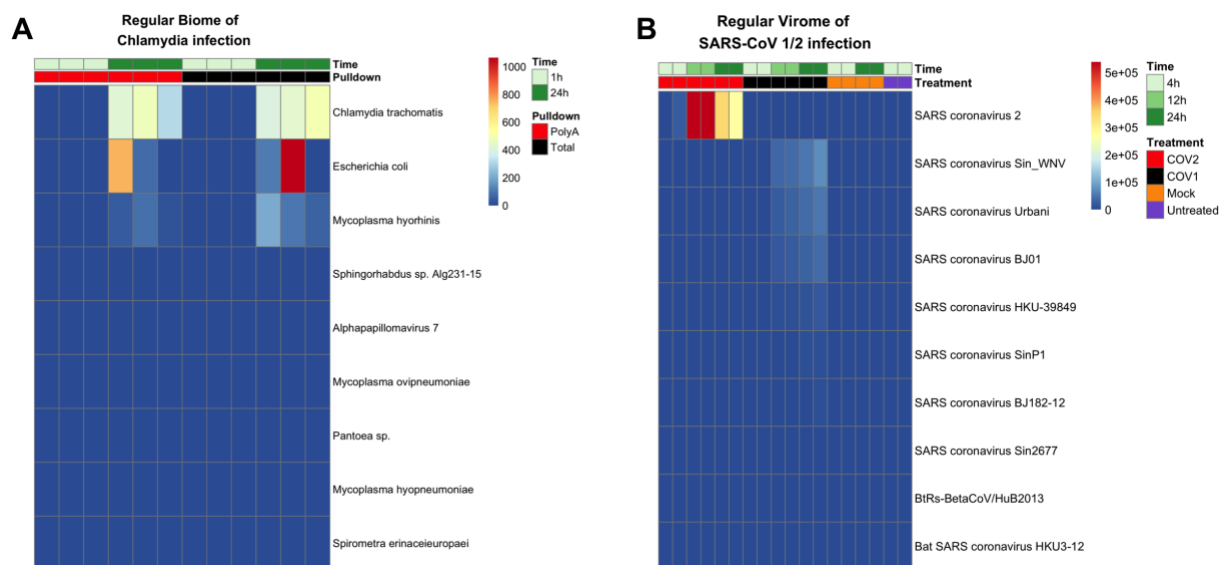
Out of the $2.81 \times 10^8$ reads obtained from all of the samples, $8.23 \times 10^7$ reads remained unmapped (~29%) after trimming and mapping to the host genome (File S2). From these non-host reads, 42,816 contigs were assembled, of which 1,346 regular biome, 27 dark biome, and 7 double dark biome contigs passed the filtering steps (File S2, File S3, File S4, File S5)

Mystery Miner successfully identified both SARS-CoV-2 and SARS-CoV-1 isolates and found significantly elevated levels of each virus compared to controls (Figure 2B). Hereafter we refer to SARS-CoV-1 or SARS-CoV-2 infected cells as COV1 or COV2 to avoid confusion with recovered names of contigs. Consistent with the viruses being similar, we identified both SARS-CoV-2 and SARS-CoV-1 in both the COV1-24hr and COV2-24hr samples when compared to mock-24hr. However, when we compared COV2-24hr to COV1-24hr, we were able to distinguish SARS-CoV-1 isolates from SARS-CoV-2 in the appropriate samples (i.e., SARS-CoV-2 was significantly elevated in COV2). Similar results were seen in the 12 hour samples but the 4 hour samples did not have sufficient viral reads to detect either SARS-CoV virus (File S7). To simulate a novel virus, we ran Mystery Miner on versions of the BLASTN and BLASTX databases obtained before SARS-CoV-2 was discovered and were able to properly identify SARS-CoV-2 as a bat related coronavirus[77] (Figure S1) (File S7).

Collectively, these data show that Mystery Miner is able to identify potential bacterial and viral infections, properly identify infected samples using quantification, and detect significant differences between infected samples and controls for bacteria, viruses, and isolates of a virus.

349



350
351
352
353 **Figure 2. Heatmap of binned normalized coverage for bacterial or viral infected datasets. A.**
354 Regular biome contigs binned by species from Humphrys et al., 2016. Time refers to 1or 24 hours
355 post infection (hpi) of epithelial cell monolayers with *Chlamydia trachomatis* (green). Pulldown
356 refers to library enrichment for polyA RNA (red) or total RNA (black). **B**. Regular virome of
357 contigs binned by name from Emanuel et al., 2020 for SARS-CoV-2 infected cells (COV2) (red),
358 or SARS-CoV-1infected cells (COV1) (black), mock virus (orange), or untreated sample (purple).
359 Time refers 4,12, or 24 hpi of Calu3 cells with indicated virus (green). Top 10 hits per experiment
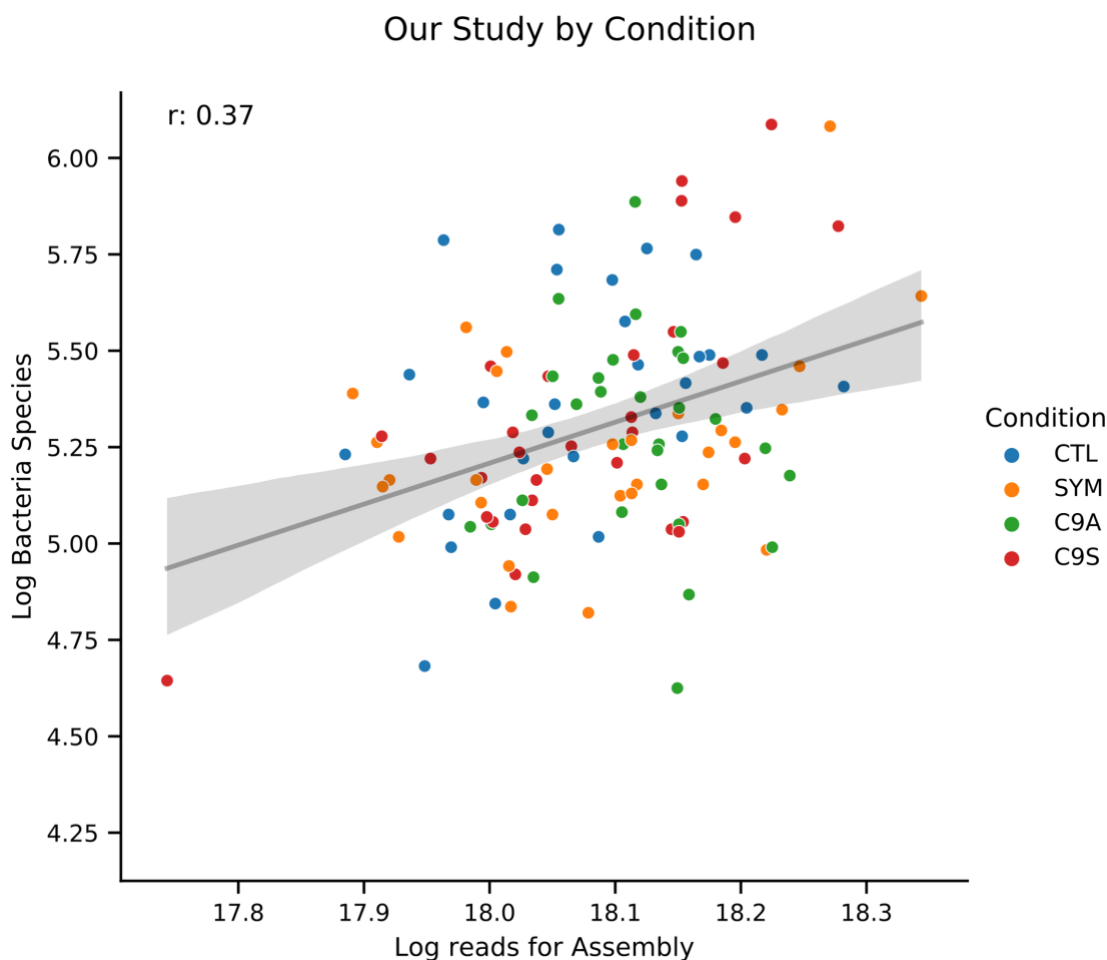360 shown for brevity.

361
362
363
364
365
366
367 **Effects of library pulldown or globin depletion in RNA-seq datasets**
368
369      In order to compare effects of library enrichment or depletion, we compared recovered
370 pathogens in a dataset that has polyA enrichment or rRNA depleted total RNA from blood or
371 colonic tissue (VonSchack2018)[78]. When we compared polyA RNA vs total RNA and looked at
372 BNC by superkingdom of bacteria we found significantly more reads map to bacteria for total
373 RNA than polyA RNA (Padj = 0.0349), in blood but not in colon (Padj=0.11709) (Figures S2 and
374 File S8). We found similar amounts of significant BNC by species for polyA RNA vs total RNA
375 in blood (29) and in colon (26). We then looked at significant BNC by genus and found double the
376 amount in blood (14) compared to colon (7), with only one significant genus (*Actinomyces*) found
377 in both comparisons. We did not find any significant differences in coverage when we looked at
378 the species, genus or superkingdom level for viruses (File S8). We conclude that library

379 enrichment with total RNA compared to polyA RNA increases bacterial recovery and diversity in
380 blood.
381     We next looked at a RNA-seq dataset from whole blood with globin depleted (GD) vs non-
382 globin depleted (NGD) total RNA (Shin2014[79]). With BNC by superkingdom, we found
383 significantly increased levels in globin depleted vs. not-depleted samples for both bacteria (Padj =
384 0.004) (Figure S3) and viruses (Padj = 0.030) (Figure S4). We also found significant differences
385 in BNC by species (Figure S5) or genus (Figure S6) primarily from *E. coli* with elevated levels in
386 globin-depleted blood RNA. We did not find any significant differences when we looked for
387 viruses at the species or genus level (File S9).
388



389
390 **Figure 3. Log number of bacterial species vs Log reads for Assembly in Our Study.** Scatterplot
391 where each dot is a sample from a dataset with log number of bacterial contigs assembled on the
392 Y-axis and Log reads used in SPAdes on the X-axis. Samples show a modest correlation (Pearson's
393 r=0.37) between library size and bacterial species recovered. Data fit with a regression (black line)
394 and 95% confidence interval (gray area).

395
396
397
398
399

**Analysis of Our Study**

We used Mystery Miner on our novel RNA-seq dataset of globin depleted and rRNA depleted total blood RNA from 120 individuals. These samples were from four subject groups including healthy control participants (CTL), ALS symptomatic *C9ORF72* negative patients (SYM), *C9ORF72* positive ALS symptomatic patients (C9S) and *C9ORF72* positive asymptomatic individuals (C9A).
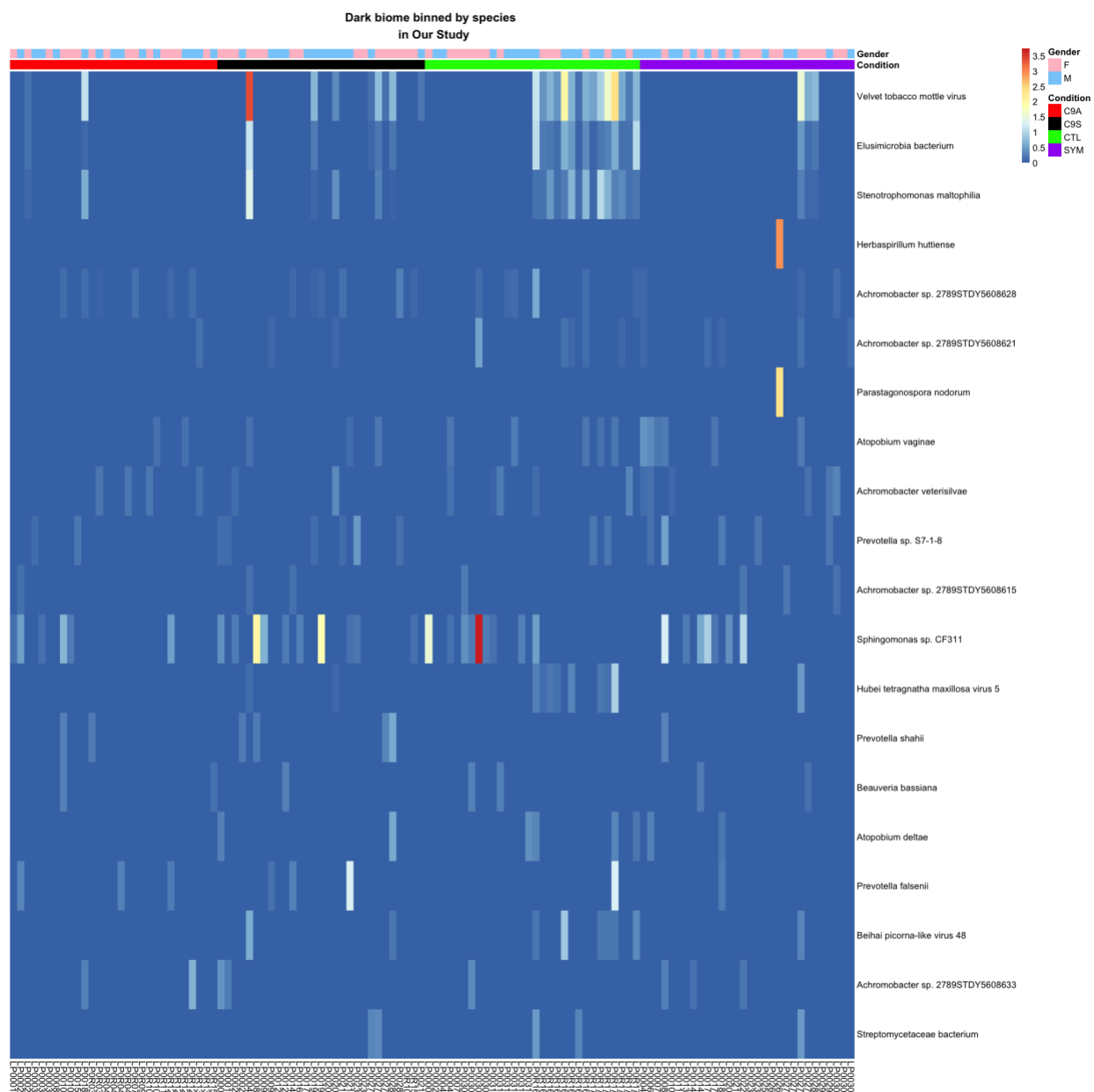
The entire dataset contains a combined $8.64 \times 10^9$ reads. Approximately 2.7% ($2.34 \times 10^8$) of the reads did not map to the human genome. From these non-host reads 2,976,988 contigs were assembled and 17,047 BLASTN contigs (regular biome) were identified. A total of 25,815 contigs had no BLASTN hit and after filtering we identified 2,980 dark biome (BLASTX identified) and 859 double dark biome (no BLASTX or BLASTN hit) contigs (File S2, File S3, File S4, File S5).

In general, we found a modest positive correlation between library size and number of bacterial contigs assembled, species detected (Figure 3), and genera detected for each sample as well as a homogenous mixture of samples with respect to disease status. No specific taxonomy or contig sequence correlated with participant class within the dataset. Yet, by pooling bacterial read counts across all of the samples, we found *alpha proteo-bacteria*, *Actinobacteria, Firmicutes,* and *Bacteroidetes* as the most highly represented taxonomies, consistent with other blood biome studies[80] (Figure S6). Most of the bacterial genera (~65%) assigned to the dark biome contigs were also found in the regular biome, however this was not the case in the viral sets, as only 5% (4/78) of dark viral contigs were observed in the regular biome (File S10). This observation suggested that our pipeline might be identifying novel viral sequences.

Within the dark biome contigs, we noted numerous contigs with a region of protein sequence similarity to RNA-dependent RNA polymerase (RdRP) from multiple RNA viruses, showing highest similarity to the velvet tobacco mottle virus (first row in heatmap of Figure 4, complete metadata shown in Figure S7). Our attention was drawn to the largest (~5 kb) dark biome contig (one of the contigs showing similarity to the velvet tobacco mottle virus) hereafter labeled as "RDRP contig". To confirm the presence of the RDRP contig in the original samples, we designed primers to the RDRP contig and performed reverse transcriptase polymerase chain reaction (RT-PCR) on seven samples, four of which had high coverage (predicted present) and three with zero coverage (predicted absent). We found typical levels for detection of a virus[81] in the samples with high coverage and detected nothing in samples with zero coverage (Table 1). We conclude that Mystery Miner is biologically validated and can recover unknown pathogens from human subjects.
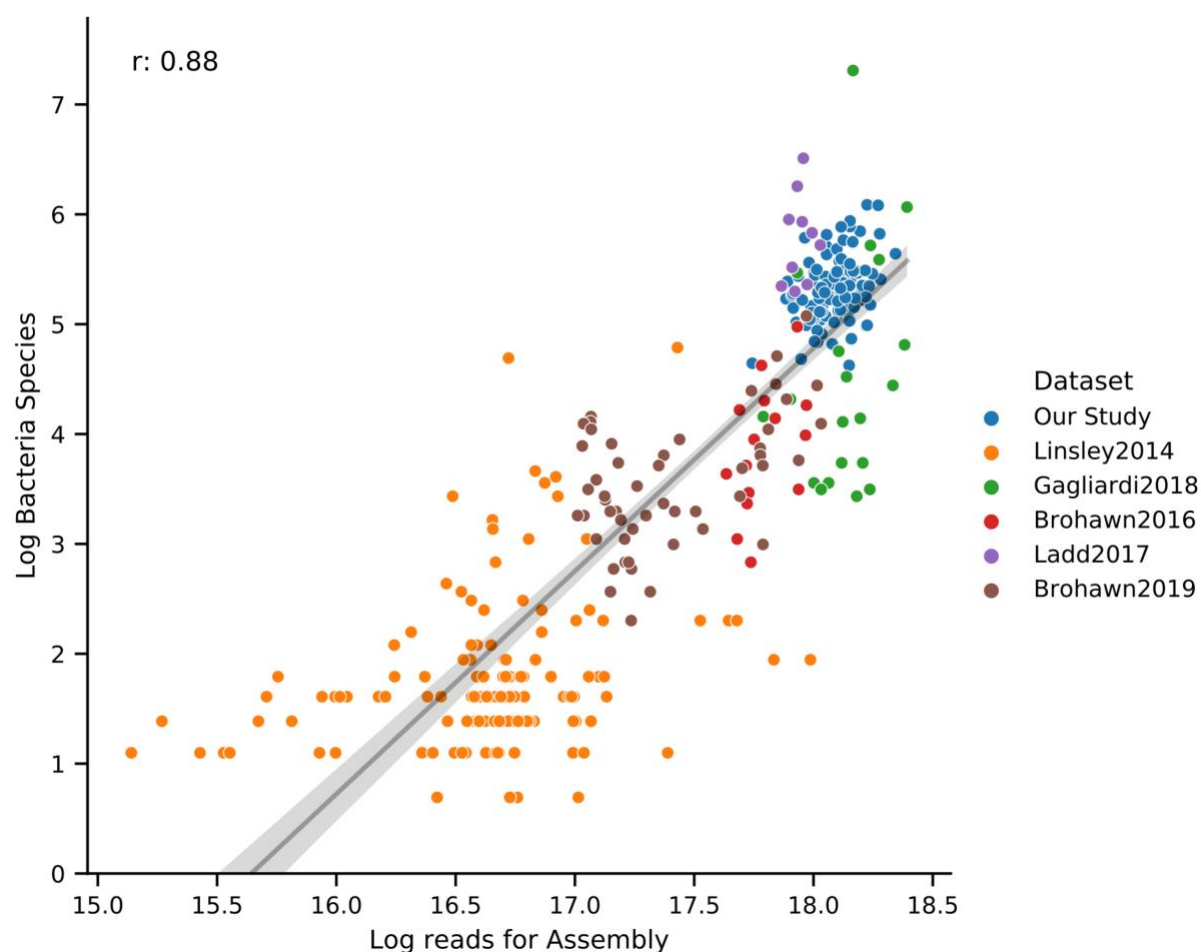
446
447
448



449
**Figure 4. Heatmap of dark biome contigs binned by species in Our Study**.
Heatmap of normalized coverage of dark biome contigs binned by species. The highest coverage belongs to contigs that show high similarity to velvet tobacco mottle virus. Zero coverage is blue and goes to red with increasing values. These samples were from four subject groups including healthy controls [(CTL) green], *C9ORF72* negative ALS symptomatic [(SYM) purple], *C9ORF72* positive ALS symptomatic [(C9S) blue] and *C9ORF72* positive asymptomatic [(C9A) red] patients. Sex indicated as light blue (male) and pink (female). Top 20 species sorted by binned normalized coverage was shown for brevity.

458
459

460
461
462
463
464
465
466
467



468

**Figure 5. Log number of bacterial species vs Log reads for Assembly for ALS Datasets.** Scatterplot where each dot is a sample from a dataset with log number of bacterial contigs assembled on the Y-axis and Log reads used in SPAdes on the X-axis. ALS datasets show a high correlation (Pearson's r=0.88) between library size and bacterial species recovered. Data fit with a regression (black line) and 95% confidence interval (gray area).
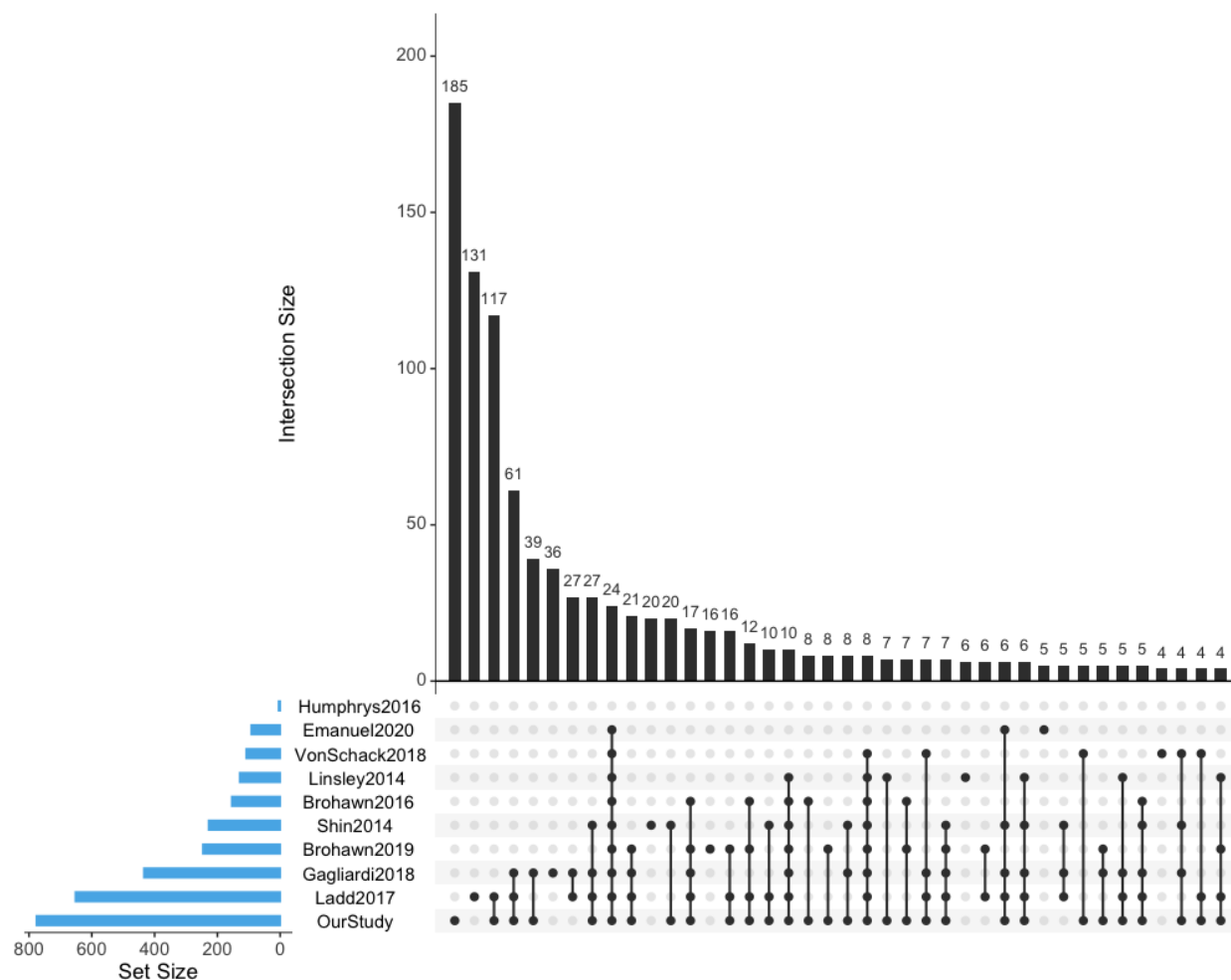
474
475

**Analysis of published ALS datasets**

We next sought to explore whether similar results would be obtained from other ALS datasets. To this end, we examined five other publicly available ALS datasets, consisting of two that used total RNA from blood (Linsley2014[82], Gagliardi 2018[58]), and three datasets from spinal cord (Brohawn2016[83], Ladd2017[84], Brohawn 2019[85]). We have provided a summary table of

481  datasets for all studies used in this paper (Table 2). As we observed in Our Study, we first noted
482  that increased library size correlated with an increased number of bacterial contigs assembled,
483  species detected, and genera detected (Figure 5, and Figure S8-10 show all datasets used in this
484  study).
485       We then looked at the total overlap of genus or species to determine if there are similarities
486  in recovered microbial or viral sequences between datasets. For genus in the regular bacteriome,
487  our dataset had the highest number of unique genus (185), followed by Ladd2017 (117), and
488  Gagliardi2018 (38). The highest number of overlapping bacterial genus was between our dataset
489  and Ladd2017 (133) followed by the intersection between our dataset, Ladd2017 and
490  Gagliardi2018 (61) and there was a modest overlap (24) for 9/10 datasets (Figure 6). We observed
491  roughly the same trend in the regular bacterial biome at the species level and in the dark bacterial
492  biome (S Figure 11, File S11). In contrast, the regular virome of each dataset was relatively unique
493  with very low amounts of overlap (<= 3) between datasets (species and genus shows a similar
494  pattern). Interestingly, the highest overlap for species in the dark virome was between our dataset
495  and Ladd2017 (13), one of which is similar to RDRP viruses, although the contigs in Ladd's data
496  were not similar to the velvet tobacco mottle virus in our dataset (Figure S12, File S12).
497       In addition to comparing datasets using taxonomy, we also compared contigs between
498  datasets for nucleotide similarity (> 70%) using LAST (File S1 for methods). We found that in
499  general, datasets in the regular biome with the largest amount of contigs have the most overlap.
500  Unsurprisingly, in the dark biome we observed less overlap by nucleotide similarity and that our
501  RDRP contig does not share nucleotide similarity with contigs from any dataset. In addition, we
502  also compared the nucleotide similarity of double dark biome contigs and found there is not a large
503  percentage of similar contigs between datasets (File S13).

**Figure 6. Upset plots of overlapping genus in the regular bacteriome between datasets.**
Upset plots are Venn diagram-like plots. A set refers to a dataset used in this study and each set is on a row with total amounts in a set as a blue bar plot on the left (ordered by set size). The black histogram on top shows the counts that are in the intersection of sets (a single dot for one dataset or connected dots for overlap of multiple datasets). Intersections less than 4 are removed for visualization purposes.
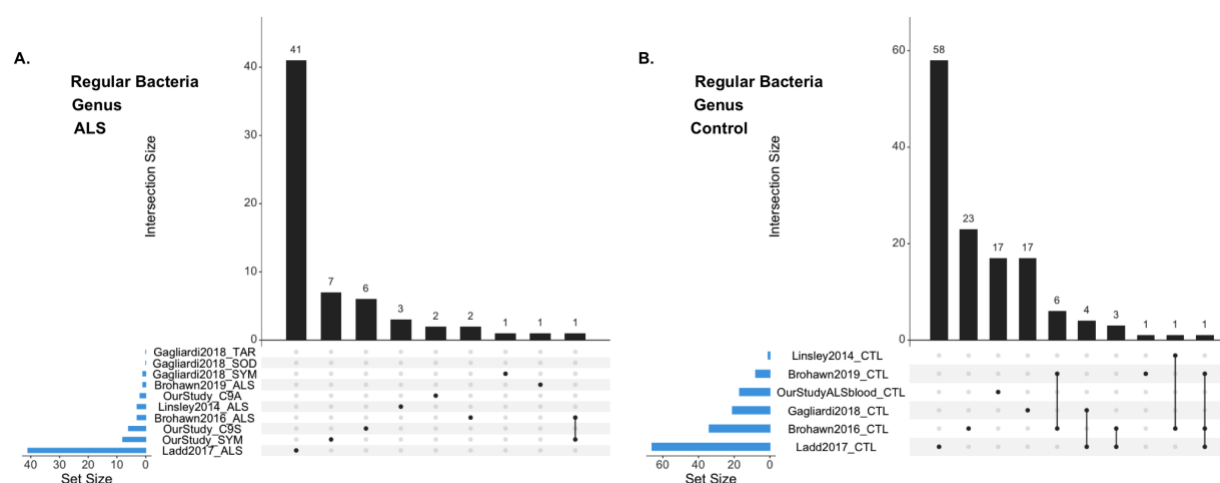
**Comparison of taxonomy by condition within ALS datasets**

Finally, we looked for differences in ALS vs control samples for each dataset. In the Gagliardi2014 dataset, when we compared ALS patients with the *FUS* mutation to controls, we found 3 significant differences in BNC by species in the regular bacteriome (*Neisseria sp.*, *Pseudomonas sp.*, *Sphingomonas sp.*) and one significant difference in BNC by genus in the dark bacteriome (*photobacterium)*. In ALS patients with mutations in *SOD1* compared to controls, we found two species significantly different in the regular bacteriome (*Hydrogenophaga crassostreae, Sphingomonas hengshuiensis) (Gagliardi FUS and SOD1 supplement)*. We did not find anything significant in sporadic ALS, or in ALS patients with *TARDBP* mutations with regards to

524 genus/species (regular or dark biome or viruses) for Gagliardi2014. We found no significant
525 statistical differences between ALS and control samples for genus/species of viruses/bacteria in
526 the regular/dark biome for any of the remaining ALS datasets.
527
528 **Meta analysis between datasets**
529
530 Since our dataset and many others had no significant comparisons for ALS vs control
531 groups within each dataset, a meta-analysis between datasets using this criteria would be difficult.
532 As a second pass analysis we created a less stringent filtering method in order to compare the
533 presence of microbes for each group between datasets (ALS vs. ALS; or controls vs. controls)
534 (Figure 7). We assigned a contig to a condition if $\geq 2$ samples from that condition contain at least
535 90% of the summed normalized coverage (from all samples) to the contig. This filtering greatly
536 reduced the number of comparable genus/species for each dataset and, for example, reduced the
537 genus of the regular bacteriome in our dataset from 305 for all samples to 33 (SYM:8, C9S:6,
538 C9A:2, CTL:17) (File S14).
539 When we looked at ALS or control contigs in the regular bacteriome, the highest number
540 of unique genus or species was from Ladd2017, and in general there was a small amount of overlap
541 between datasets ($\leq 1$ for ALS or $\leq 8$ for controls) (Figure 7). When we looked at genus in the dark
542 bacteriome we saw no overlap for ALS contigs and low overlap ($\leq 1$) between control conditions
543 (species was similar) (File S14). In the regular virome there was no overlap between datasets and
544 only our study (one contig from ALS) and Ladd2017 (three from ALS, five from controls) had
545 contigs that passed the filter (similar values for species). When we looked in the dark virome by
546 genus there was no overlap between datasets, and our dataset had only one genus (*Sobemovirus*
547 from controls*)* with the rest coming from Ladd2017 (18 from controls, 5 from ALS) (File S15). In
548 conclusion, despite our conservative and loose approaches, we did not find any convincing
549 evidence in ALS samples that the presence (or lack of presence) of an organism (or multiple
550 organisms) was different compared to control samples.
551



**Fig 7. Upset plots of overlapping genus between datasets in the regular biome for ALS or controls.**
Upset plots are Venn diagram-like plots. A set refers to a contig that was assigned to a condition from a dataset. Each set is on a row with total amounts in a set as a blue bar plot on the left (ordered by set size). The black histogram on top shows the counts that are in the intersection of sets (a

558 single dot for one dataset or connected dots for overlap of multiple datasets). A. ALS contigs in
559 the regular bacteriome. B. Control contigs from the regular bacteriome.
560
561
562 **Discussion**
563
564      We have created Mystery Miner to search for and quantify known and unknown microbes
565 in RNA-seq datasets as a tool for researchers to study dysbiosis and identify infectious agents. We
566 validated the pipeline by recovering and quantifying *Chlamydia* and SARS-CoV reads from RNA-
567 seq datasets from intentionally infected cells. Interestingly, we also identified *Mycoplasma* reads
568 in the *Chlamydia* dataset, suggesting that Mystery Miner may also be able to identify unsuspected
569 bacterial infections. We also use published data to investigate the difference of polyA vs total RNA
570 recovery of bacterial species in multiple tissues. Perhaps surprisingly, we did not see a consistent
571 difference in the recovery of bacterial reads between the two types of RNA-seq libraries,
572 considering that bacterial transcripts are not expected to be polyadenylated. However, it is well-
573 recognized that polyA selection is imperfect, and libraries constructed from polyA-selected RNA
574 routinely contain transcripts thought not to be polyadenylated (e.g., rRNA). We also found
575 increased recovery of bacterial species with globin RNA depletion in blood. This could be the
576 result of an effective increase in read depth for bacteria when not sequencing globin, or an increase
577 in contamination from the globin depletion step. We stress that our bioinformatic approach alone
578 cannot distinguish between contamination and the true existence of microbial sequences in human
579 tissue.
580      We then used Mystery Miner on Our Study dataset consisting of 8.64 X $10^9$ reads. This
581 dataset was generated from whole blood total RNA that was depleted for both ribosomal and globin
582 transcripts. It encompasses samples from controls, participants with a *C9ORF72* hexanucleotide
583 expansion (symptomatic and pre-symptomatic), and *C9ORF72* negative ALS patients. We found
584 no statistical difference in microbial sequence read coverage between controls and any class of
585 ALS patients, examining either individual contigs or using various taxonomical binning
586 approaches. We also did not detect any batch effects or obvious age- or sex- biases in the recovery
587 of microbial reads (Figure S7). Overall, we found no evidence of blood microbial sequences
588 contributing to either *C9ORF72* negative ALS or symptomatic patients harboring the *C9ORF72*
589 hexanucleotide expansion. However, ALS is a CNS disease, so our findings in these blood samples
590 do not necessarily preclude a role for microbes in this disease.
591
592      A unique aspect of Mystery Miner is that it tracks non-human reads that do not have
593 significant BLASTN hits in GenBank. We were intrigued by the identification of a large contig
594 (>5kb) in the dark biome of our ALS dataset that showed protein sequence similarity to RNA-
595 dependent RNA polymerases, the essential replicase of RNA viruses. To validate that this virus-
596 like sequence was not an artifact of contig assembly or a contaminant introduced during library
597 construction or sequencing, we used RT-PCR of the original patient samples to demonstrate that
598 this sequence was present in positive samples identified through the RNA-seq analysis and not
599 detectable in negative samples. We cannot extrapolate from this specific example to determine
600 what fraction of the "dark" and "double dark" sequences represent true novel microbial sequences
601 present in human blood, but we note that analysis of human cell free blood DNA has reported the
602 identification of thousands of novel bacterial sequences[86]. We suggest that Mystery Miner is a
603 generally useful tool for the identification of novel microbial sequences in RNA-seq data.

604

605       To extend our analysis we processed publicly available blood and spinal cord ALS datasets
606  through our pipeline. As observed in our dataset, library size generally correlated with number of
607  bacterial contigs assembled and number of bacterial genera/species recovered. When the microbial
608  sequences we found in our dataset were compared to the other datasets we found similar
609  genera/species and, not surprisingly, larger datasets generally had greater overlap. One dataset
610  (Ladd2017) yielded comparable recovery of bacteria and viruses for the regular biome but a far
611  greater recovery bacteria and viruses in the dark biome compared to all the other datasets. This
612  study used laser capture microdissection (LCM) to isolate cervical spinal cord motor neurons and
613  had comparable read amounts per sample to other studies and was conducted in the same
614  laboratory as two other studies (Brohawn2016, Brohawn2019). We are unsure why this dataset
615  yielded a much larger dark biome compared to the other datasets. Potentially these differences are
616  a byproduct of using LCM to acquire samples.

617

618       We then analyzed several publicly available ALS datasets for statistically significant
619  differences between recovered microbial sequences in ALS and control samples. Only one dataset
620  (Gagliardi2018) had any significant microbial sequence differences between control and ALS
621  samples, specifically ALS patients with *FUS* or *SOD1* mutations. However, the sample number
622  in this sub-study was small (N = 2-3), and in the case of the *SOD1* patients the excess microbial
623  reads were in the control samples. In the absence of additional information (e.g., batch assignments
624  for the samples) it is difficult to conclude that these sequence/sample correlations are meaningful.
625  Finally, we compared identified microbial sequences in the control and ALS samples across the
626  datasets and did not identify any genera/species that were preferentially recovered in either sample
627  type.

628

629       Using our bioinformatic analysis pipeline Mystery Miner, we have not identified an
630  association between ALS pathology and sequences corresponding to known or unknown microbial
631  species. However, there are intrinsic limitations in using "repurposed" RNA-seq data to assay
632  tissue-associated microbial sequences, including the relatively small number of non-human reads
633  (<1% of total) upon which the analysis is based. This limited sequence signal could preclude
634  identification of rarer microbes. Perhaps more problematic is the possibility that contaminating
635  sequences obscure true tissue-associated microbial sequences. Any candidate microbes identified
636  using Mystery Miner that correlate with human phenotypes will necessarily require independent
637  validation. Despite these limitations, we believe Mystery Miner will be a useful tool for future
638  researchers investigating known and unknown microbes that could contribute to disease, as our
639  analyses have shown it to be sensitive to bacterial/viral agents in sequencing data.

640

641

642

643

653
654

| Condition | Sample | GAPDH RT-PCR Ct Value | RDRP RT-PCR Ct Value | RDRP RNA-seq Normalized Coverage |
|---|---|---|---|---|
| SYM | LP00274 | 20.562019 | 36.401 | 1.56 |
| C9S | LP00041 | 20.783213 | 36.346 | 3.39 |
| C9S | LP00192 | 20.899612 | 35.636 | 0.67 |
| C9A | LP000180 | 19.982108 | 34.832 | 1.11 |
| C9S | LP000183 | 20.176418 | undetermined | 0 |
| C9S | LP000197 | 20.125161 | undetermined | 0 |
| C9A | LP000157 | 20.062433 | undetermined | 0 |

655
656
657   **TABLE 1. RT-PCR AND NORMALIZED COVERAGE RESULTS FOR RDRP CONTIG**
658   Quantitative RT-PCR and normalized coverage results from the 5180 bp RDRP contig. For the RDRP contig positive
659   samples (top 4) with high normalized coverage and detectable Ct values and negative samples (bottom 3) with no
660   normalized coverage and undetectable Ct values. GAPDH was used as a positive control for qRT-PCR and shows
661   comparable levels for all samples. These samples were from three conditions *C9ORF72* negative ALS symptomatic
662   patients (SYM), *C9ORF72* positive ALS symptomatic patients (C9S) and *C9ORF72* positive asymptomatic
663   individuals (C9A).

664
665
666
667
668
669
670
671
672

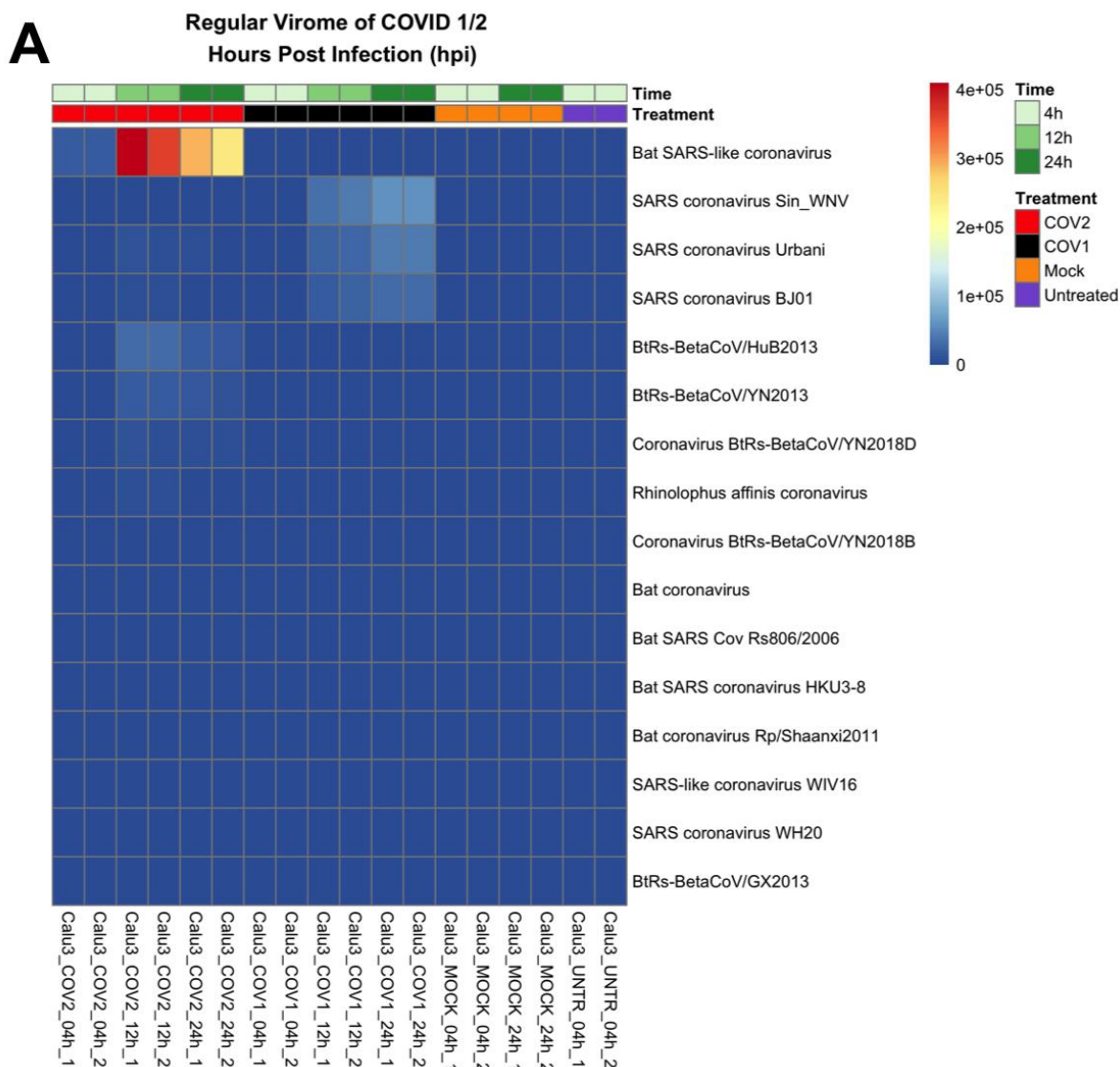| Name | Groups | # Samples | Tissue | Pulldown |
|---|---|---|---|---|
| Humphrys2016 | 1- or 24-hours post infection with *Chlamydia trachomatis* | 12 | Cultured epithelial cell monolayers | PolyA Total RNA |
| VonSchack2018 | PolyA or Total RNA from blood or colon | 16 | Whole Blood Colon | PolyA RNA Total RNA |
| Shin2014 | Globin depleted Not globin depleted | 24 | Whole Blood | Total RNA |
| Emanuel2020 | Severe acute respiratory syndrome coronavirus 1 or 2 infection Controls | 18 | Calu3 cells | Total RNA |
| Our Study | *C9ORF72* negative ALS, *C9ORF72* positive and symptomatic ALS, *C9ORF72* positive asymptomatic participants Controls | 120 | Whole Blood | Total RNA hemoglobin and rRNA depleted |
| Linsley2014 | ALS type 1 diabetes, sepsis, multiple sclerosis patients before and 24 hours after the first treatment with IFN-beta Controls | 134 | Whole blood | Total RNA |
| Gagliardi2018 | Sporadic ALS, ALS with mutations in *FUS*, *SOD1*, *TARDBP* Controls | 20 | Peripheral blood mononuclear cells | Total RNA |
| Brohawn2016 | ALS Controls | 15 | Cervical spinal cord | Total RNA rRNA depleted |

| Ladd2017 | ALS Controls | 10 | Laser capture microdissection (LCM) to isolate cervical spinal cord motor neurons | Total RNA |
|---|---|---|---|---|
| Brohawn2019 | ALS, Alzheimer's disease (AD), Parkinson's disease (PD) Controls | 53 | Cervical spinal cord | Total RNA |

673
674
675  **TABLE 2. STUDY DESIGN FOR DATASETS USED IN THIS PAPER**
676  Overview of the datasets used in this paper. The first three studies are only used to validate our pipeline. The six
677  subsequent studies are ALS related from both blood and spinal cord.
678
679
680
681  **Supplemental Figures**
682
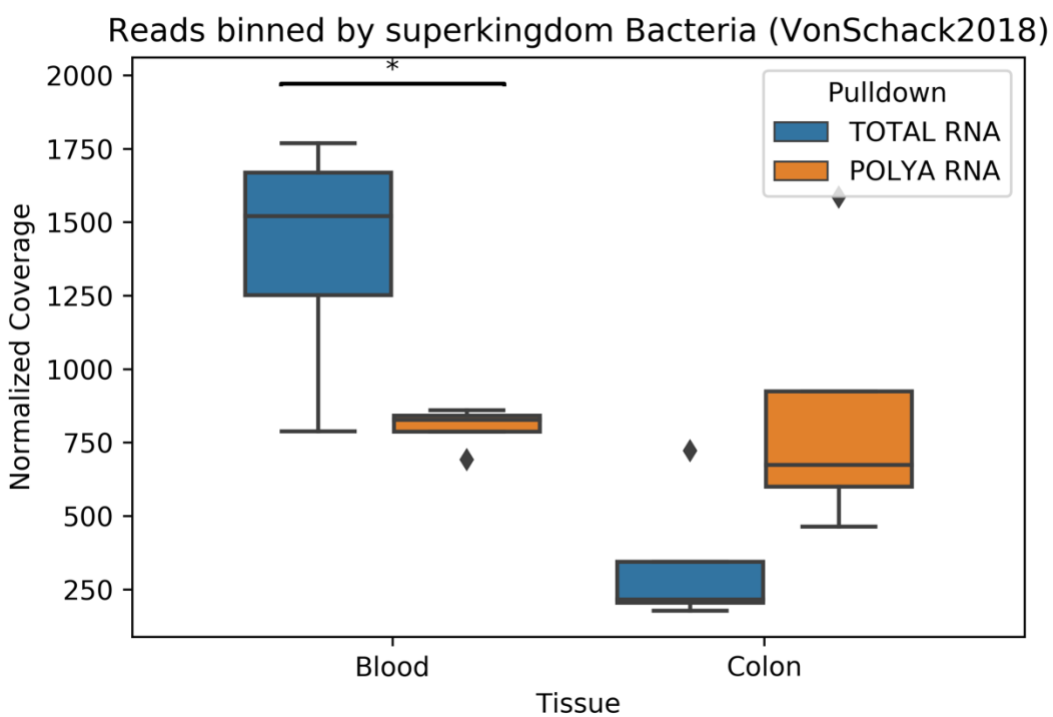
683
684
**Figure S1. Heatmap of normalized coverage of regular Virome from Emanuel2020 with**
**BLAST to nt database from 05/10/2019**
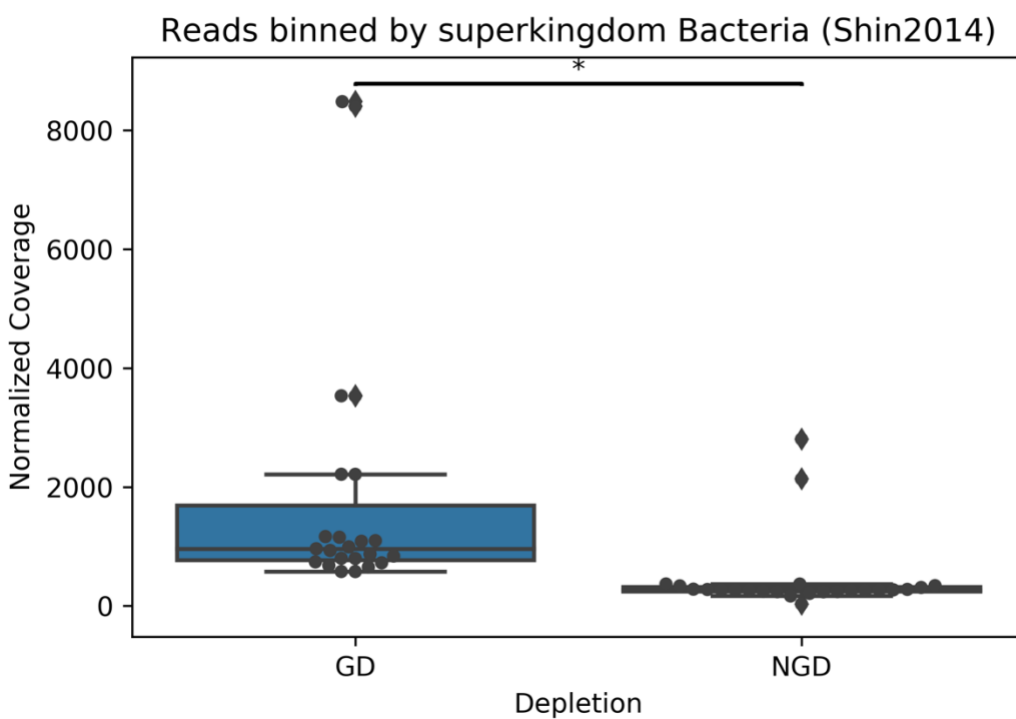Heatmap of normalized coverage of dark biome contigs binned by species (top 30 species). The
nucleotide database was from 5/10/2019 before the discovery of SARS-CoV-2. The top row
shows the same row from the main text but identified as a bat SARS like coronavirus.
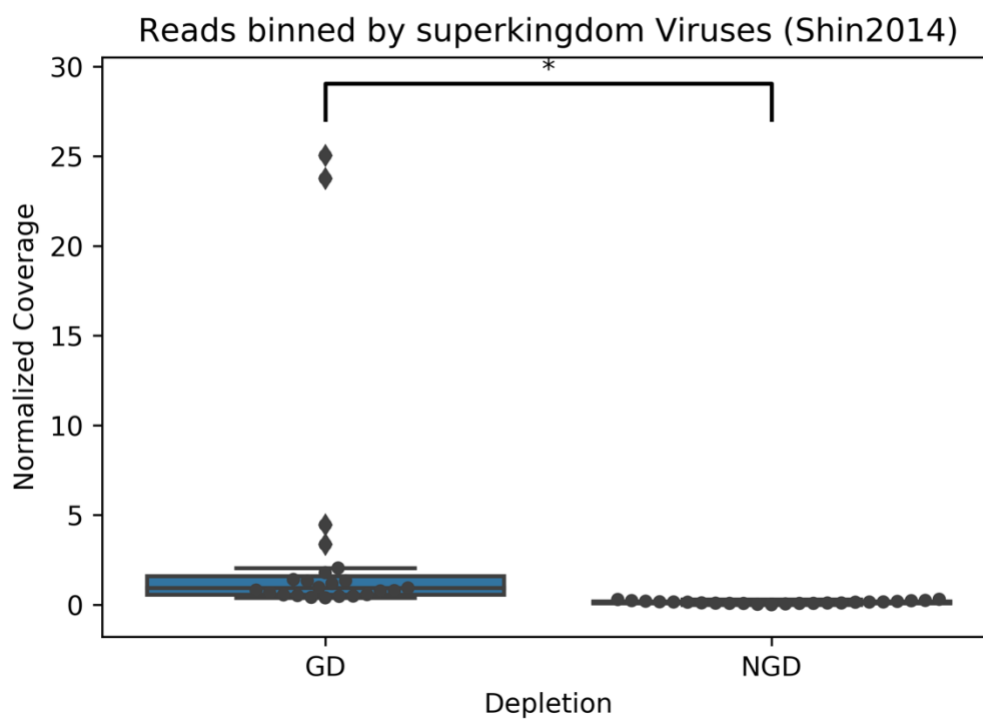
690
691
692
693
694
695
696

**Figure S2. Boxplot of normalized coverage for superkingdom Bacteria in VonSchack2018**
Boxplot of normalized coverage of regular biome contigs binned by superkingdom Bacteria.
Blood shows significantly more reads in total RNA vs polyA RNA compared to Colon tissue.

705
**Figure S3. Boxplot of normalized coverage for superkingdom Bacteria in Shin2014**
Boxplot of normalized coverage of regular biome contigs binned by superkingdom Bacteria.
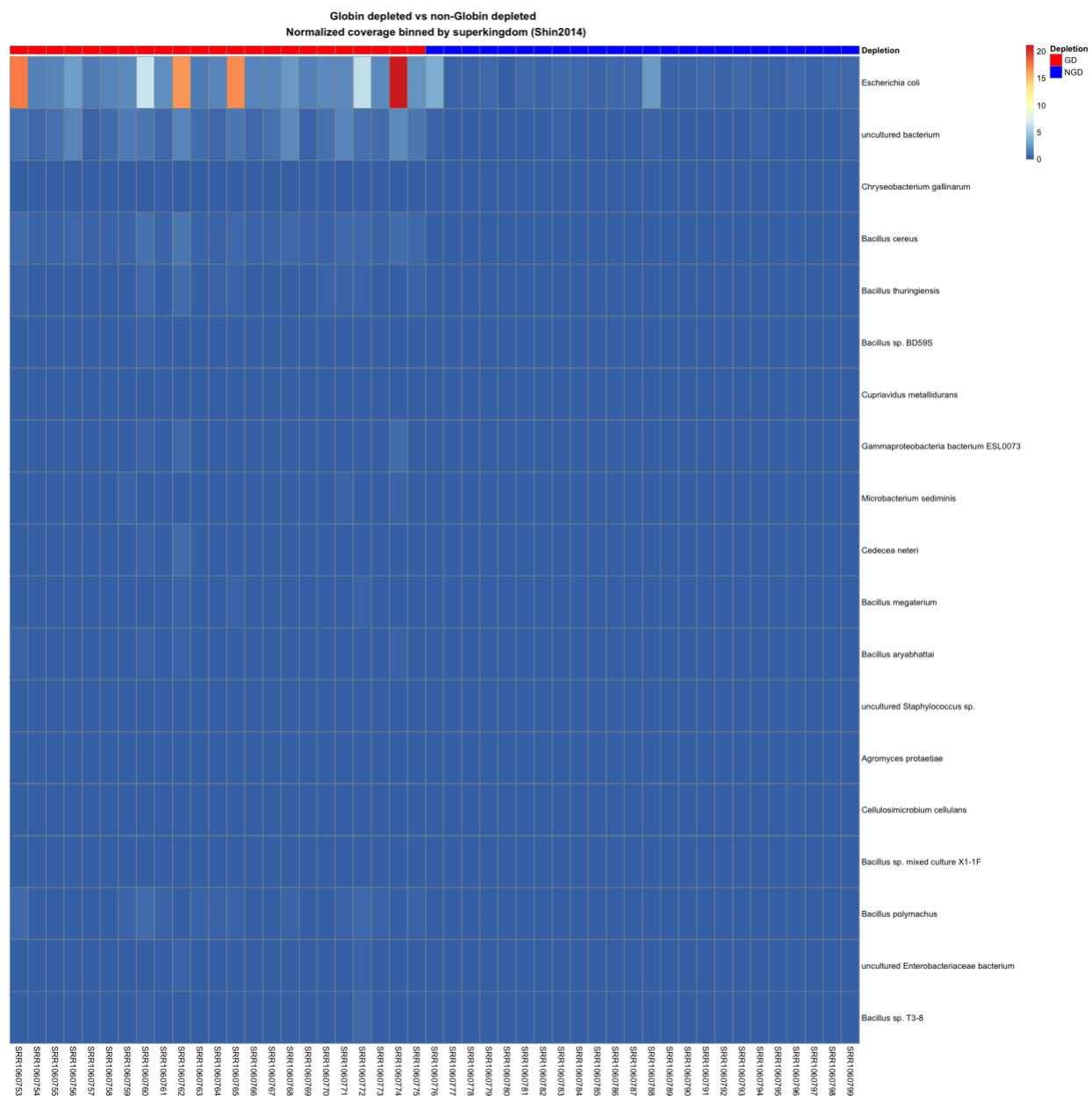Globin depletion (GD) has significantly more coverage than non-globin depleted (NGD) blood.

709

710

711

**Figure S4. Boxplot of normalized coverage for superkingdom Viruses in Shin2014**
Boxplot of normalized coverage of regular biome contigs binned by superkingdom Viruses.
Globin depletion (GD) has significantly more coverage than non-globin depleted (NGD) blood.

**Figure S5. Heatmap of normalized coverage of regular Bacteriome binned by species from Shin2014**

Heatmap of normalized coverage of regular biome contigs binned by bacteria species (top 20 species shown for brevity). Globin depletion (GD) is red and non-globin depletion is blue (NGD).

**Figure S5. Heatmap of normalized coverage of regular Bacteriome binned by genus from Shin2014**

Heatmap of normalized coverage of regular biome contigs binned by bacteria genus. Globin depletion (GD) is red and non-globin depletion is blue (NGD).

732
**Figure S6. Log coverage binned by phylum from our ALS dataset**
Coverage is summed for all of the samples and *alpha proteo-bacteria*, *Actinobacteria,*
*Firmicutes,* and *Bacteroidetes* are the most highly represented.
736
737
738
739
740

741
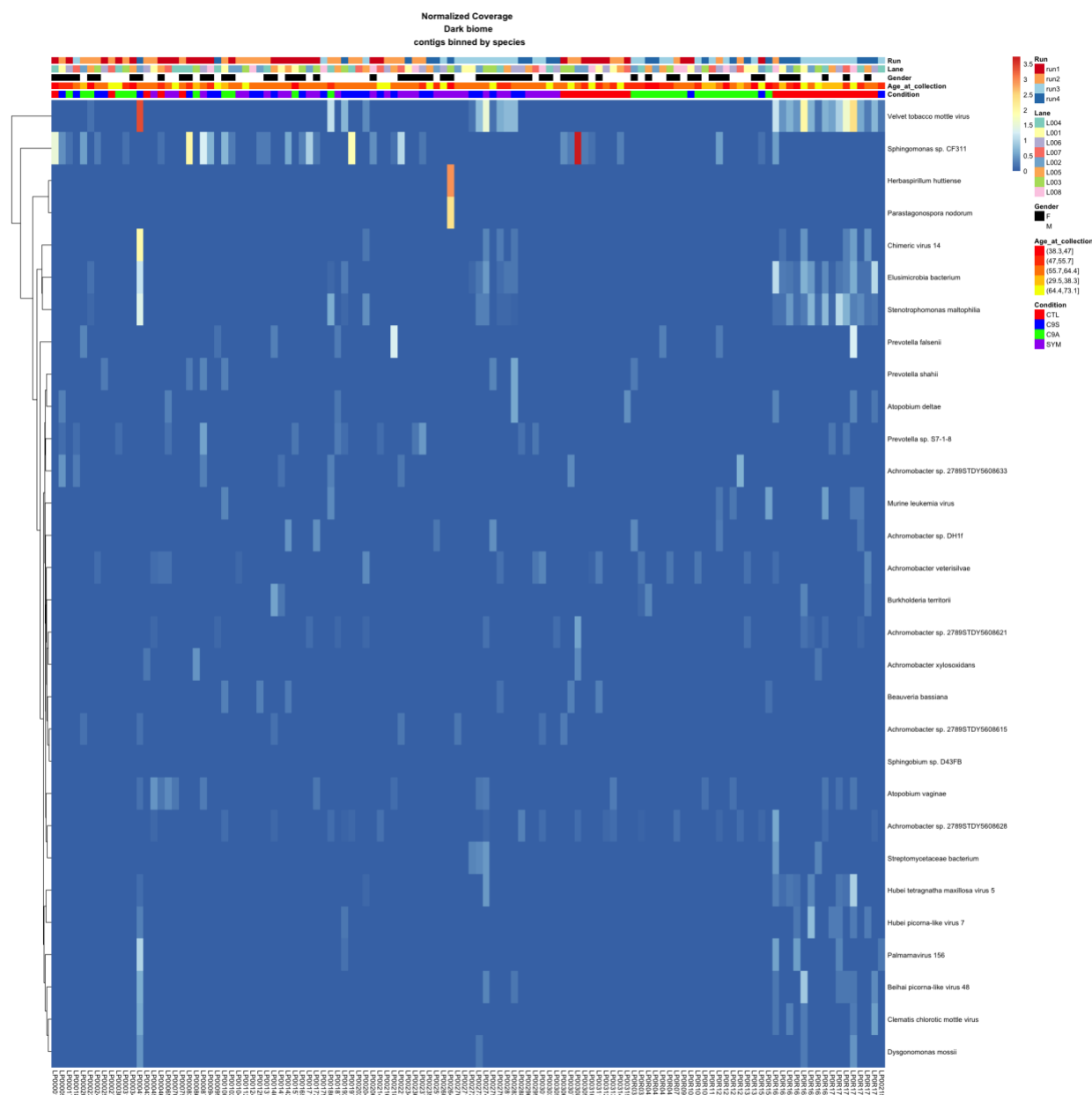**Figure S7. Heatmap of normalized coverage of dark biome contigs binned by species with metadata**
Heatmap of normalized coverage of dark biome contigs binned by species (top 30 species shown for brevity). The highest coverage belongs to contigs that show high similarity to velvet tobacco mettle virus. Zero coverage is blue and goes to red with increasing values. These samples were from four conditions including control patients [(CTL) green], ALS symptomatic patients [(SYM) purple], C9-ORF positive ALS symptomatic patients [(C9S) blue] and C9-ORF positive asymptomatic patients [(C9A) red]. Other metadata include gender, lane, run, and age at collection.

**Figure S8. Log Bacterial contigs vs log reads for Assembly.** Scatterplot where each dot is a sample from a dataset with log number of Bacterial contigs assembled on the Y-axis and Log reads used in SPAdes on the X-axis. Aside from the Shin, Humphrys, and Emanuel datasets there is a general trend of increased number of bacterial contigs with amount of reads.

**Figure S9. Log number of bacterial species vs log reads for Assembly.** Scatterplot where each dot is a sample from a dataset with log number of number of bacterial species detected on the Y-axis and Log reads used in SPAdes on the X-axis.

**Figure S10. Log number of bacterial genus vs log reads for Assembly.** Scatterplot where each dot is a sample from a dataset with log number of number of bacterial genus detected on the Y-axis and Log reads used in SPAdes on the X-axis.

**Figure S11. Upset plots of Bacteria for genus/species of regular/dark genome**

Upset plots are venn diagram-like plots. Each set is on a row with total amounts in a set as a blue bar plot on the left. The black histogram on top shows the counts that are in the intersection of sets (a single dot for one set or connected dots for multiple sets). The highest number of overlapping bacterial genus is between our dataset and Ladd2017 (133) followed by the intersection between our dataset, Ladd2017 and Gagliardi2018 (61) and there is a modest overlap (24) for 9/10 datasets. This is roughly similar in the Bacterial species figure and in general the larger datasets have more unique and the highest number of overlap.

**Figure S12. Upset plots of Viruses for genus/species of regular/dark genome**
Upset plots are venn diagram-like plots. Each set is on a row with total amounts in a set as a blue
bar plot on the left. The black histogram on top shows the counts that are in the intersection of
sets (a single dot for one set or connected dots for multiple sets). The regular virome of each
dataset is relatively unique with very low amounts of overlap ($<= 3$) between datasets (species
and genus shows a similar pattern). Interestingly, the highest overlap for species in the dark
virome is between our dataset and Ladd2017 (13).

**Figure S13. Upset plots of Bacteria in the regular biome for genus/species in ALS and Control contigs**

Upset plots are venn diagram-like plots. Each set is on a row with total amounts in a set as a blue bar plot on the left. The black histogram on top shows the counts that are in the intersection of sets (a single dot for one set or connected dots for multiple sets). We assigned a contig to a condition if >= 2 samples from that condition contain at least 90% of the summed normalized coverage (from all samples) to the contig. In the genus and species from ALS samples there is a low amount of overlap between datasets ( <= 1). When we look at control samples there is a much higher overlap for both genus and species.

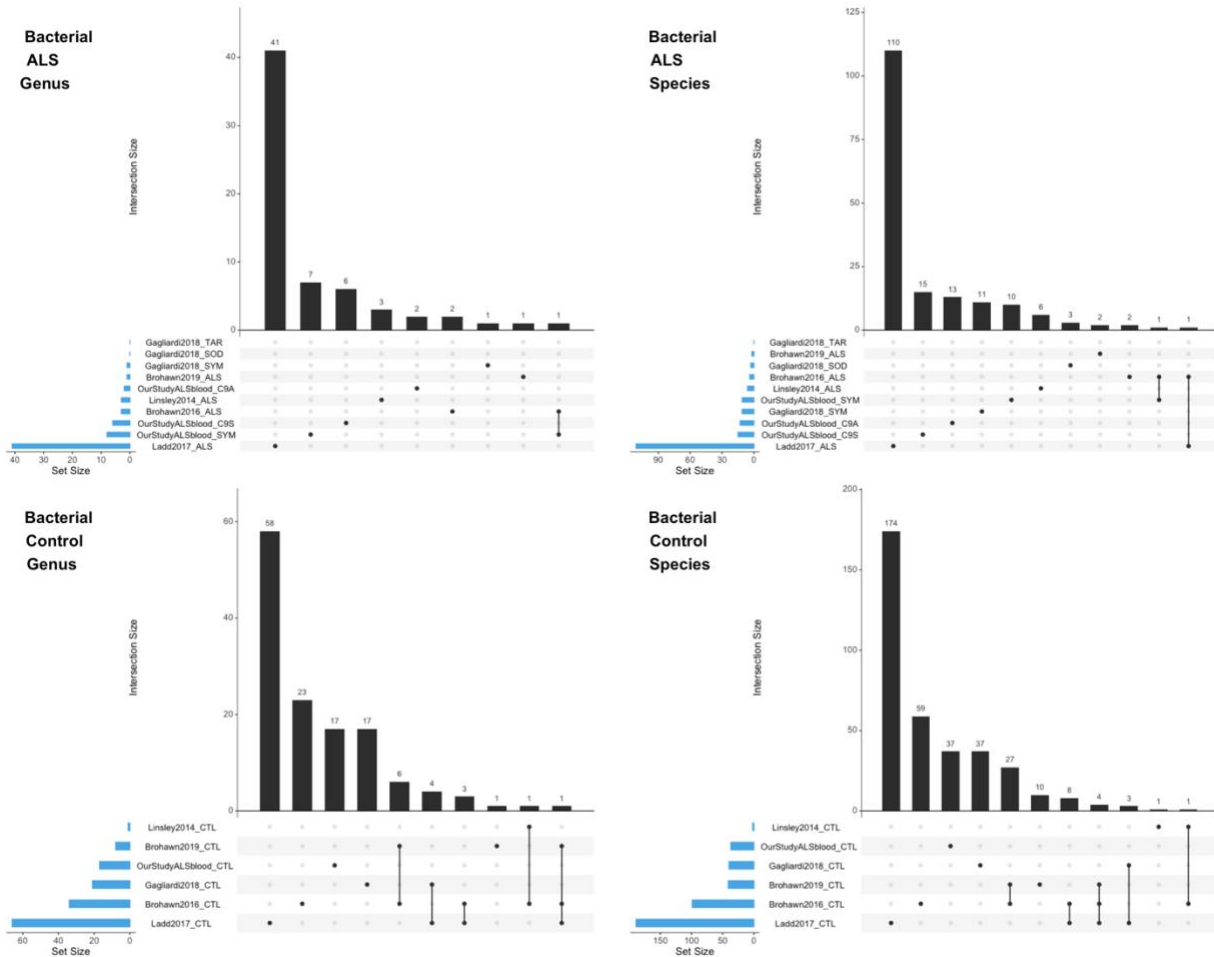**Figure S13. Upset plots of Bacteria in the dark biome for genus/species in ALS and Control contigs**

Upset plots are venn diagram-like plots. Each set is on a row with total amounts in a set as a blue bar plot on the left. The black histogram on top shows the counts that are in the intersection of sets (a single dot for one set or connected dots for multiple sets). We assigned a contig to a condition if >= 2 samples from that condition contain at least 90% of the summed normalized coverage (from all samples) to the contig. Conditions with no recovered viruses have been omitted for clarity. Similarly to the regular bacteriome, there is no overlap in ALS samples and a small amount of overlap in the conditions.

847
## References
849

850  1.   Patrick KL, Bell SL, Weindel CG, Watson RO. Exploring the "multiple-hit hypothesis" of
851       neurodegenerative disease: Bacterial infection comes up to bat. *Front Cell Infect*
852       *Microbiol*. 2019. doi:10.3389/fcimb.2019.00138

853  2.   Castanedo-Vazquez D, Bosque-Varela P, Sainz-Pelayo A, Riancho J. Infectious agents and
854       amyotrophic lateral sclerosis: another piece of the puzzle of motor neuron degeneration.
855       *J Neurol*. 2019. doi:10.1007/s00415-018-8919-3

856  3.   Mehta P, Kaye W, Raymond J, et al. Prevalence of amyotrophic lateral sclerosis — United
857       States, 2015. *Morb Mortal Wkly Rep*. 2018;67(46):1285-1289.
858       doi:10.15585/mmwr.mm6746a1

859  4.   Talbott EO, Malek AM, Lacomis D. The epidemiology of amyotrophic lateral sclerosis. In:
860       *Handbook of Clinical Neurology*. Vol 138. Elsevier B.V.; 2016:225-238. doi:10.1016/B978-
861       0-12-802973-2.00013-6

862  5.   Ingre C, Roos PM, Piehl F, Kamel F, Fang F. Risk factors for amyotrophic lateral sclerosis.
863       *Clin Epidemiol*. 2015. doi:10.2147/CLEP.S37505

864  6.   Zhan Y, Fang F. Smoking and amyotrophic lateral sclerosis: A mendelian randomization
865       study. *Ann Neurol*. 2019. doi:10.1002/ana.25443

866  7.   Opie-Martin S, Wootton RE, Budu-Aggrey A, et al. Relationship between smoking and
867       ALS: Mendelian randomisation interrogation of causality. *J Neurol Neurosurg Psychiatry*.
868       2020. doi:10.1136/jnnp-2020-323316

869  8.   Kohne DE, Gibbs CJ, White L, Tracy SM, Meinke W, Smith RA. Virus detection by nucleic
870       acid hybridization: Examination of normal and ALS tissues for the presence of poliovirus.
871       *J Gen Virol*. 1981;56(2):223-233. doi:10.1099/0022-1317-56-2-223

872  9.   Pertschuk LP, Broome JD, Brigati DJ, et al. JEJUNAL IMMUNOPATHOLOGY IN
873       AMYOTROPHIC LATERAL SCLEROSIS AND MULTIPLE SCLEROSIS IDENTIFICATION OF VIRAL
874       ANTIGENS BY IMMUNOFLUORESCENCE. *Lancet*. 1977;309(8022):1119-1123.
875       doi:10.1016/S0140-6736(77)92382-0

876  10.  Xue YC, Feuer R, Cashman N, Luo H. Enteroviral Infection: The Forgotten Link to
877       Amyotrophic Lateral Sclerosis? *Front Mol Neurosci*. 2018;11:63.
878       doi:10.3389/fnmol.2018.00063

879  11.  Alonso R, Pisa D, Fernández-Fernández AM, Rábano A, Carrasco L. Fungal infection in
880       neural tissue of patients with amyotrophic lateral sclerosis. *Neurobiol Dis*. 2017;108:249-
881       260. doi:10.1016/j.nbd.2017.09.001

882  12.  Andrade FC, Vergetti V, Cozza G, Falcao MC, Azevedo G. Amyotrophic Lateral Sclerosis-
883       like Syndrome after Chikungunya. *Cureus*. October 2019. doi:10.7759/cureus.5876

884  13.  Deutsch SI, Mohs RC, Davis KL. A rationale for studying the transmissibility of Alzheimer's
885       disease. *Neurobiol Aging*. 1982;3(2):145-147. doi:10.1016/0197-4580(82)90011-2

886  14.  Taylor GR, Crow TJ, Markakis DA, Lofthouse R, Neeley S, Carter GI. Herpes simplex virus
887       and Alzheimer's disease: A search for virus DNA by spot hybridisation. *J Neurol Neurosurg*
888       *Psychiatry*. 1984;47(10):1061-1065. doi:10.1136/jnnp.47.10.1061

889  15.  Sochocka M, Zwolińska K, Leszek J. The Infectious Etiology of Alzheimer's Disease. *Curr*
890       *Neuropharmacol*. 2017;15(7). doi:10.2174/1570159x15666170313122937

891   16.   Irkeç C. [Virologic and immunologic considerations in Parkinson's disease]. *Mikrobiyol*
892          *Bul*. 1982;16(4):293-296. http://www.ncbi.nlm.nih.gov/pubmed/6304477. Accessed
893          December 9, 2019.

894   17.   Abushouk AI, El-Husseny MWA, Magdy M, et al. Evidence for association between
895          hepatitis C virus and Parkinson's disease. *Neurol Sci*. 2017;38(11):1913-1920.
896          doi:10.1007/s10072-017-3077-4

897   18.   Parashar A, Udayabanu M. Gut microbiota: Implications in Parkinson's disease. *Park*
898          *Relat Disord*. 2017;38:1-7. doi:10.1016/j.parkreldis.2017.02.002

899   19.   Libbey JE, Cusick MF, Fujinami RS. Role of pathogens in multiple sclerosis. *Int Rev*
900          *Immunol*. 2014;33(4):266-283. doi:10.3109/08830185.2013.823422

901   20.   Alonso R, Pisa D, Carrasco L. Searching for Bacteria in Neural Tissue From Amyotrophic
902          Lateral Sclerosis. *Front Neurosci*. 2019;13:171. doi:10.3389/fnins.2019.00171

903   21.   Gil C, González AAS, León IL, et al. Detection of Mycoplasmas in Patients with
904          Amyotrophic Lateral Sclerosis. *Adv Microbiol*. 2014;04(11):712-719.
905          doi:10.4236/aim.2014.411077

906   22.   Alonso R, Pisa D, Marina AI, et al. Evidence for fungal infection in cerebrospinal fluid and
907          brain tissue from patients with amyotrophic lateral sclerosis. *Int J Biol Sci*.
908          2015;11(5):546-558. doi:10.7150/ijbs.11084

909   23.   Pisa D, Alonso R, Rábano A, Carrasco L. Corpora Amylacea of Brain Tissue from
910          Neurodegenerative Diseases Are Stained with Specific Antifungal Antibodies. *Front*
911          *Neurosci*. 2016;10:86. doi:10.3389/fnins.2016.00086

912   24.   Cermelli C, Vinceti M, Beretti F, et al. Risk of sporadic amyotrophic lateral sclerosis
913          associated with seropositivity for herpesviruses and echovirus-7. *Eur J Epidemiol*.
914          2003;18(2):123-127. doi:10.1023/a:1023067728557

915   25.   Berger MM, Kopp N, Vital C, Redl B, Aymard M, Lina B. Detection and cellular localization
916          of enterovirus RNA sequences in spinal cord of patients with ALS. *Neurology*.
917          2000;54(1):20-25. doi:10.1212/wnl.54.1.20

918   26.   Vandenberghe N, Leveque N, Corcia P, et al. Cerebrospinal fluid detection of enterovirus
919          genome in ALS: A study of 242 patients and 354 controls. *Amyotroph Lateral Scler*.
920          2010;11(3):277-282. doi:10.3109/17482960903262083

921   27.   Xue YC, Feuer R, Cashman N, Luo H. Enteroviral infection: The forgotten link to
922          amyotrophic lateral sclerosis? *Front Mol Neurosci*. 2018;11.
923          doi:10.3389/fnmol.2018.00063

924   28.   Giraud P, Beaulieux F, Ono S, Shimizu N, Chazot G, Lina B. Detection of enteroviral
925          sequences from frozen spinal cord samples of Japanese ALS patients. *Neurology*.
926          2001;56(12):1777-1778. doi:10.1212/wnl.56.12.1777

927   29.   Verma A, Berger JR. ALS syndrome in patients with HIV-1 infection. *J Neurol Sci*.
928          2006;240(1-2):59-64. doi:10.1016/j.jns.2005.09.005

929   30.   Moodley K, Bill PLA, Bhigjee AI, Patel VB. A comparative study of motor neuron disease in
930          HIV-infected and HIV-uninfected patients. *J Neurol Sci*. 2019;397:96-102.
931          doi:10.1016/J.JNS.2018.12.030

932   31.   Douville R, Liu J, Rothstein J, Nath A. Identification of active loci of a human endogenous
933          retrovirus in neurons of patients with amyotrophic lateral sclerosis. *Ann Neurol*.
934          2011;69(1):141-151. doi:10.1002/ana.22149

935    32.    Li W, Lee M-H, Henderson L, et al. Human endogenous retrovirus-K contributes to motor
936           neuron disease. *Sci Transl Med*. 2015;7(307):307ra153-307ra153.
937           doi:10.1126/scitranslmed.aac8201

938    33.    Arru G, Mameli G, Deiana GA, et al. Humoral immunity response to human endogenous
939           retroviruses K/W differentiates between amyotrophic lateral sclerosis and other
940           neurological diseases. *Eur J Neurol*. 2018;25(8):1076-e84. doi:10.1111/ene.13648

941    34.    Blacher E, Bashiardes S, Shapiro H, et al. Potential roles of gut microbiome and
942           metabolites in modulating ALS in mice. *Nature*. 2019;572(7770):474-480.
943           doi:10.1038/s41586-019-1443-5

944    35.    Fang X, Wang X, Yang S, et al. Evaluation of the microbial diversity in amyotrophic lateral
945           sclerosis using high-throughput sequencing. *Front Microbiol*. 2016;7(SEP).
946           doi:10.3389/fmicb.2016.01479

947    36.    Sun J, Zhan Y, Mariosa D, et al. Antibiotics use and risk of amyotrophic lateral sclerosis in
948           Sweden. *Eur J Neurol*. 2019;26(11):1355-1361. doi:10.1111/ene.13986

949    37.    Zhang Y guo, Wu S, Yi J, et al. Target Intestinal Microbiota to Alleviate Disease
950           Progression in Amyotrophic Lateral Sclerosis. *Clin Ther*. 2017;39(2):322-336.
951           doi:10.1016/j.clinthera.2016.12.014

952    38.    Obrenovich M, Jaworski H, Tadimalla T, et al. The role of the microbiota–gut–brain axis
953           and antibiotics in ALS and neurodegenerative diseases. *Microorganisms*. 2020;8(5).
954           doi:10.3390/microorganisms8050784

955    39.    Brenner D, Hiergeist A, Adis C, et al. The fecal microbiome of ALS patients. *Neurobiol
956           Aging*. 2018;61:132-137. doi:10.1016/j.neurobiolaging.2017.09.023

957    40.    Henkel JS, Beers DR, Wen S, Bowser R, Appel SH. Decreased mRNA expression of tight
958           junction proteins in lumbar spinal cords of patients with ALS. *Neurology*.
959           2009;72(18):1614-1616. doi:10.1212/WNL.0b013e3181a41228

960    41.    Garbuzova-Davis S, Sanberg PR. Blood-CNS barrier impairment in ALS patients versus an
961           animal model. *Front Cell Neurosci*. 2014;8(FEB):21. doi:10.3389/fncel.2014.00021

962    42.    Zhang R, Miller RG, Gascon R, et al. Circulating endotoxin and systemic immune
963           activation in sporadic amyotrophic lateral sclerosis (sALS). *J Neuroimmunol*. 2009;206(1-
964           2):121-124. doi:10.1016/j.jneuroim.2008.09.017

965    43.    Mantovani S, Garbelli S, Pasini A, et al. Immune system alterations in sporadic
966           amyotrophic lateral sclerosis patients suggest an ongoing neuroinflammatory process. *J
967           Neuroimmunol*. 2009;210(1-2):73-79. doi:10.1016/j.jneuroim.2009.02.012

968    44.    Murdock BJ, Zhou T, Kashlan SR, Little RJ, Goutman SA, Feldman EL. Correlation of
969           peripheral immunity with rapid Amyotrophic lateral sclerosis progression. *JAMA Neurol*.
970           2017;74(12):1446-1454. doi:10.1001/jamaneurol.2017.2255

971    45.    Sta M, Sylva-Steenland RMR, Casula M, et al. Innate and adaptive immunity in
972           amyotrophic lateral sclerosis: Evidence of complement activation. *Neurobiol Dis*.
973           2011;42(3):211-220. doi:10.1016/j.nbd.2011.01.002

974    46.    Correia AS, Patel P, Dutta K, Julien JP. Inflammation induces TDP-43 mislocalization and
975           aggregation. *PLoS One*. 2015;10(10). doi:10.1371/journal.pone.0140248

976    47.    Verber NS, Shepheard SR, Sassani M, et al. Biomarkers in motor neuron disease: A state
977           of the art review. *Front Neurol*. 2019;10(APR):291. doi:10.3389/fneur.2019.00291

978    48.    Blasco H, Corcia P, Moreau C, et al. 1H-NMR-Based metabolomic profiling of CSF in early

979        amyotrophic lateral sclerosis. *PLoS One*. 2010;5(10). doi:10.1371/journal.pone.0013223

980    49.    Blasco H, Veyrat-Durebex C, Bocca C, et al. Lipidomics Reveals Cerebrospinal-Fluid
981        Signatures of ALS. *Sci Rep*. 2017;7(1). doi:10.1038/s41598-017-17389-9

982    50.    Mitchell RM, Freeman WM, Randazzo WT, et al. A CSF biomarker panel for identification
983        of patients with amyotrophic lateral sclerosis. *Neurology*. 2009;72(1):14-19.
984        doi:10.1212/01.wnl.0000333251.36681.a5

985    51.    Guo J, Yang X, Gao L, Zang D. Evaluating the levels of CSF and serum factors in ALS. *Brain
986        Behav*. 2017;7(3). doi:10.1002/brb3.637

987    52.    Young PE, Jew SK, Buckland ME, Pamphlett R, Suter CM. Epigenetic differences between
988        monozygotic twins discordant for amyotrophic lateral sclerosis (ALS) provide clues to
989        disease pathogenesis. *PLoS One*. 2017;12(8). doi:10.1371/journal.pone.0182638

990    53.    Coppedè F, Stoccoro A, Mosca L, et al. Increase in DNA methylation in patients with
991        amyotrophic lateral sclerosis carriers of not fully penetrant SOD1 mutations. *Amyotroph
992        Lateral Scler Front Degener*. 2018;19(1-2):93-101. doi:10.1080/21678421.2017.1367401

993    54.    Swindell WR, Kruse CPS, List EO, Berryman DE, Kopchick JJ. ALS blood expression profiling
994        identifies new biomarkers, patient subgroups, and evidence for neutrophilia and hypoxia.
995        *J Transl Med*. 2019;17(1):170. doi:10.1186/s12967-019-1909-0

996    55.    Waller R, Wyles M, Heath PR, et al. Small RNA sequencing of sporadic amyotrophic
997        lateral sclerosis cerebrospinal fluid reveals differentially expressed miRNAs related to
998        neural and glial activity. *Front Neurosci*. 2018;11(JAN). doi:10.3389/fnins.2017.00731

999    56.    Waller R, Goodall EF, Milo M, et al. Serum miRNAs miR-206, 143-3p and 374b-5p as
1000        potential biomarkers for amyotrophic lateral sclerosis (ALS). *Neurobiol Aging*.
1001        2017;55:123-131. doi:10.1016/j.neurobiolaging.2017.03.027

1002    57.    Gendron TF, Chew J, Stankowski JN, et al. Poly(GP) proteins are a useful
1003        pharmacodynamic marker for C9ORF72-associated amyotrophic lateral sclerosis. *Sci
1004        Transl Med*. 2017;9(383):7866. doi:10.1126/scitranslmed.aai7866

1005    58.    Gagliardi S, Zucca S, Pandini C, et al. Long non-coding and coding RNAs characterization
1006        in Peripheral Blood Mononuclear Cells and Spinal Cord from Amyotrophic Lateral
1007        Sclerosis patients. *Sci Rep*. 2018;8(1):2378. doi:10.1038/s41598-018-20679-5

1008    59.    Zucca S, Gagliardi S, Pandini C, et al. RNA-Seq profiling in peripheral blood mononuclear
1009        cells of amyotrophic lateral sclerosis patients and controls. *Sci Data*. 2019;6(1):190006.
1010        doi:10.1038/sdata.2019.6

1011    60.    Rahman MR, Islam T, Huq F, Quinn JMW, Moni MA. Identification of molecular signatures
1012        and pathways common to blood cells and brain tissue of amyotrophic lateral sclerosis
1013        patients. *Informatics Med Unlocked*. 2019;16:100193. doi:10.1016/J.IMU.2019.100193

1014    61.    van Rheenen W, Diekstra FP, Harschnitz O, et al. Whole blood transcriptome analysis in
1015        amyotrophic lateral sclerosis: A biomarker study. *PLoS One*. 2018;13(6):e0198874.
1016        doi:10.1371/journal.pone.0198874

1017    62.    Parker J, Chen J. Application of next generation sequencing for the detection of human
1018        viral pathogens in clinical specimens. *J Clin Virol*. 2017;86:20-26.
1019        doi:10.1016/j.jcv.2016.11.010

1020    63.    Bouquet J, Gardy JL, Brown S, et al. RNA-Seq Analysis of Gene Expression, Viral Pathogen,
1021        and B-Cell/T-Cell Receptor Signatures in Complex Chronic Disease. *Clin Infect Dis*.
1022        2017;64(4):476-481. doi:10.1093/cid/ciw767

1023  64.  Westermann AJ, Barquist L, Vogel J. Resolving host-pathogen interactions by dual RNA-
1024       seq. *PLoS Pathog*. 2017;13(2):e1006033. doi:10.1371/journal.ppat.1006033
1025  65.  Moore RA, Warren RL, Freeman JD, et al. The Sensitivity of Massively Parallel Sequencing
1026       for Detecting Candidate Infectious Agents Associated with Human Tissue. Jordan IK, ed.
1027       *PLoS One*. 2011;6(5):e19838. doi:10.1371/journal.pone.0019838
1028  66.  Poussin C, Sierro N, Boué S, et al. Interrogating the microbiome: experimental and
1029       computational considerations in support of study reproducibility. *Drug Discov Today*.
1030       2018;23(9):1644-1657. doi:10.1016/j.drudis.2018.06.005
1031  67.  Roumpeka DD, Wallace RJ, Escalettes F, Fotheringham I, Watson M. A review of
1032       bioinformatics tools for bio-prospecting from metagenomic sequence data. *Front Genet*.
1033       2017;8(MAR). doi:10.3389/fgene.2017.00023
1034  68.  Mangul S, Yang HT, Strauli N, et al. ROP: dumpster diving in RNA-sequencing to find the
1035       source of 1 trillion reads across diverse adult human tissues. *Genome Biol*. 2018;19(1):36.
1036       doi:10.1186/s13059-018-1403-7
1037  69.  Cavadas B, Ferreira J, Camacho R, Fonseca NA, Pereira L. QmihR: Pipeline for
1038       Quantification of Microbiome in Human RNA-seq. In: Springer, Cham; 2017:173-179.
1039       doi:10.1007/978-3-319-60816-7_21
1040  70.  Simon LM, Karg S, Westermann AJ, et al. MetaMap: an atlas of metatranscriptomic reads
1041       in human disease-related RNA-seq data. *Gigascience*. 2018;7(6).
1042       doi:10.1093/gigascience/giy070
1043  71.  Gihawi A, Rallapalli G, Hurst R, Cooper CS, Leggett RM, Brewer DS. SEPATH:
1044       benchmarking the search for pathogens in human tissue whole genome sequence data
1045       leads to template pipelines. *Genome Biol*. 2019;20(1):208. doi:10.1186/s13059-019-
1046       1819-8
1047  72.  Cox JW, Ballweg RA, Taft DH, Velayutham P, Haslam DB, Porollo A. A fast and robust
1048       protocol for metataxonomic analysis using RNAseq data. *Microbiome*. 2017;5(1):7.
1049       doi:10.1186/s40168-016-0219-5
1050  73.  Almeida A, Mitchell AL, Boland M, et al. A new genomic blueprint of the human gut
1051       microbiota. *Nature*. 2019;568(7753):499-504. doi:10.1038/s41586-019-0965-1
1052  74.  Papudeshi B, Haggerty JM, Doane M, et al. Optimizing and evaluating the reconstruction
1053       of Metagenome-assembled microbial genomes. *BMC Genomics*. 2017;18(1):915.
1054       doi:10.1186/s12864-017-4294-1
1055  75.  Humphrys MS, Creasy T, Sun Y, et al. Simultaneous transcriptional profiling of bacteria
1056       and their host cells. Ramsey K, ed. *PLoS One*. 2013;8(12):e80597.
1057       doi:10.1371/journal.pone.0080597
1058  76.  Emanuel W, Kirstin M, Vedran F, et al. Bulk and single-cell gene expression profiling of
1059       SARS-CoV-2 infected human cell lines identifies molecular targets for therapeutic
1060       intervention. *bioRxiv*. May 2020:2020.05.05.079194. doi:10.1101/2020.05.05.079194
1061  77.  Boni MF, Lemey P, Jiang X, et al. Evolutionary origins of the SARS-CoV-2 sarbecovirus
1062       lineage responsible for the COVID-19 pandemic. *Nat Microbiol*. July 2020:1-10.
1063       doi:10.1038/s41564-020-0771-4
1064  78.  Zhao S, Zhang Y, Gamini R, Zhang B, von Schack D. Evaluation of two main RNA-seq
1065       approaches for gene quantification in clinical RNA sequencing: polyA+ selection versus
1066       rRNA depletion. *Sci Rep*. 2018;8(1):4781. doi:10.1038/s41598-018-23226-4

1067   79. Shin H, Shannon CP, Fishbane N, et al. Variation in RNA-Seq Transcriptome Profiles of
1068       Peripheral Whole Blood from Healthy Individuals with and without Globin Depletion.
1069       Wang K, ed. *PLoS One*. 2014;9(3):e91041. doi:10.1371/journal.pone.0091041

1070   80. Castillo DJ, Rifkin RF, Cowan DA, Potgieter M. The healthy human blood microbiome: Fact
1071       or fiction? *Front Cell Infect Microbiol*. 2019;9(MAY):148. doi:10.3389/fcimb.2019.00148

1072   81. Mackay IM, Arden KE, Nitsche A. Real-time PCR in virology. *Nucleic Acids Res*.
1073       2002;30(6):1292-1305. doi:10.1093/nar/30.6.1292

1074   82. Linsley PS, Speake C, Whalen E, Chaussabel D. Copy number loss of the interferon gene
1075       cluster in melanomas is linked to reduced T cell infiltrate and poor patient prognosis.
1076       *PLoS One*. 2014;9(10). doi:10.1371/journal.pone.0109760

1077   83. Brohawn DG, O'Brien LC, Bennett JP. RNAseq analyses identify tumor necrosis factor-
1078       mediated inflammation as a major abnormality in ALS spinal cord. *PLoS One*.
1079       2016;11(8):e0160520. doi:10.1371/journal.pone.0160520

1080   84. C Ladd A, G Brohawn D, P Bennett J. Laser-captured spinal cord motorneurons from ALS
1081       subjects show increased gene expression in vacuolar ATPase networks. *J Syst Integr*
1082       *Neurosci*. 2017;3(6). doi:10.15761/jsin.1000182

1083   85. Bennett JP, Keeney PM, Brohawn DG. RNA Sequencing Reveals Small and Variable
1084       Contributions of Infectious Agents to Transcriptomes of Postmortem Nervous Tissues
1085       From Amyotrophic Lateral Sclerosis, Alzheimer's Disease and Parkinson's Disease
1086       Subjects, and Increased Expression of Genes From Disease-Activated Microglia. *Front*
1087       *Neurosci*. 2019;13. doi:10.3389/fnins.2019.00235

1088   86. Kowarsky M, Camunas-Soler J, Kertesz M, et al. Numerous uncharacterized and highly
1089       divergent microbes which colonize humans are revealed by circulating cell-free DNA.
1090       *Proc Natl Acad Sci U S A*. 2017;114(36):9623-9628. doi:10.1073/pnas.1707009114

1091