

1 Supplemental methods

2 A. Reference databases

3 SEAR has the option of removing contamination by subtracting reads using the BWA
4 index for the human genome (HG19 build, [1]), this is not supplied in the SEAR
5 package to reduce file size but can be downloaded (<https://genome.ucsc.edu/>).
6 SEAR requires a reference database for the read-clustering step, the default supplied
7 is the ARGannot database [2] but other options are available (such as CARD [3]) and
8 the user can supply any multiFASTA file as a database. The supplied ARGannot
9 database was customized as follows: ARGannot ARGs were clustered at 97%
10 identity using USEARCH [4] and the representative sequence for each cluster was
11 added to the pipeline's ARG database. Each cluster and representative sequence is
12 annotated with gene type and the class of antimicrobial to which the gene confers
13 resistance. The *Shigella* reference database for the benchmarking test was made by
14 downloading the FASTA files for each ARG tested (n=19) in the *Holt et al.* study,
15 removing duplicate entries and creating a multiFASTA file (n=16).

16 B. SEAR pipeline

17 **Dependencies:** SEAR is available as both a command-line and web-based
18 application and a full list of scripts, reference databases, external software and
19 modules are listed in Supplemental Table S1. SEAR is distributed with a SEAR_bin
20 directory that contains all the required dependencies that must be installed on the
21 system. The USEARCH package [4] is used for FASTQ quality checking, conversion
22 to FASTA format and clustering. Due to the memory restrictions with the free
23 distribution of USEARCH, the input files are split into smaller, temporary files and
24 piped into USEARCH fastq_filter and subsequently into USEARCH usearch_global.

25 In the local alignment step, SEAR uses the supplied commandline BLAST package
26 [5] and requires an internet connection. Currently, SEAR uses remote blast to NCBI
27 databases and uses cURL [6] to create up to date, local databases for other online
28 resources during every SEAR run, but local databases could be installed and used.
29 **Running:** SEAR can be run via the website interface or in a single command line
30 script. Adjustable parameters and the default settings are listed in Supplemental
31 Table S2. When using the website, files are checked and uploaded to the server
32 using the SEAR CGI script, which subsequently launches the command line script
33 using the supplied files and parameters.

34 **C. Relative abundance calculation**

35 The proportion of mapped bases relative to the length of each reference sequence is
36 calculated and the total number of successfully mapped reads is retrieved to allow for
37 ARG annotation and calculation of relative abundance. When calculating relative
38 abundance, the total number (n) of ARGs that have been annotated are used to
39 calculate a relative abundance (RA) percentage for each ARG. Firstly, an abundance

40 value (A) is calculated for each gene according to: $A = \frac{X}{Y} \cdot \frac{1}{L}$, where X denotes the
41 number of reads that successfully mapped, Y denotes the total number of reads in
42 the input file/s and L denotes the length of the reference gene (in bases). Relative

43 abundance is then calculated using: $RA = \left(\frac{A}{\sum_{i=1}^n} \right) * 100$

44 In this way, the relative abundance measure describes the proportion of sequence
45 reads that have built the consensus sequence of each annotated ARG from a single
46 pipeline run.

47 **D. Metagenome sample collection**

48 Samples were collected from two faecal sources within the River Cam Catchment,
49 Cambridge, UK on the 21st June 2012. The waste effluent of the University of
50 Cambridge dairy farm (latitude: 52.22259, longitude: 0.02603) was sampled prior to it
51 being applied to the surrounding fields as fertiliser. The effluent of the municipal
52 wastewater treatment works (WWTW) (latitude: 52.234469, longitude: 0.154614) was
53 collected from the effluent discharge pipe that enters the River Cam. Samples were
54 collected in 10L sterile polypropylene containers. Sample volumes were based on the
55 microbial abundances, as previously determined for these sites using a DNA
56 extraction series (data not shown). Samples were transported at 4°C to the
57 laboratory and processed within 2 hours.

58 **E. Sample filtration, metagenomic DNA extraction and** 59 **sequencing**

60 Samples were vacuum pre-filtered through 3.0 µm membranes (Millipore) to remove
61 debris and eukaryotic cells before being filtered at 2 Bar through 0.22 µm
62 membranes (Millipore) to capture the prokaryotic cells. Metagenomic DNA was
63 extracted by vortexing membranes in phosphate buffered saline with Tween20 (2%)
64 before enzymatic lysis (Meta-G-Nome DNA isolation kit; Epicentre). Assessment of
65 DNA quality and concentration was made by TBE agarose (2%) gel electrophoresis
66 and spectrophotometry (Nanodrop ND-1000; ThermoScientific). For each sample, 2
67 µg of DNA was sequenced by the Eastern Sequence and Informatics Hub,
68 Cambridge, UK. Seventy-five base pair paired-end libraries were prepared from the
69 samples and were sequenced using an Illumina HiSeq2000.

70 **F. NGS datasets**

71 **Environmental:** The FASTQ files for the WWTW and Farm effluent metagenomes
72 are available via the European Nucleotide Archive (ENA) (study: ERP003955).
73 Sample accession numbers are as follows: farm effluent (ERS786322), WWTW
74 effluent (ERS781558).

75 **Human Microbiome Project (HMP):** The FASTQ files for 32 microbiomes from
76 Spanish patients were downloaded from the HMP via the ENA website (study:
77 PRJEB1220) (accessed: 02.03.2015) [7].

78 ***Shigella sonnei*:** A global whole-genome-sequencing dataset of 126 clinical isolates
79 of *Shigella sonnei* (an enteric pathogen) [8] was used to test the utility of the pipeline
80 for identifying ARGs in clinical isolates. The FASTQ files for the 126 isolates were
81 downloaded from the Sanger FTP site (study: PRJEB2128) (accessed: 02.03.2015).

82

83 **G. Supplemental References**

- 84 1. Consortium GR. hg19. In: Consortium GR, editor. UCSC2009.
- 85 2. Gupta SK, Padmanabhan BR, Diene SM, Lopez-Rojas R, Kempf M, Landraud
86 L, et al. ARG-ANNOT, a New Bioinformatic Tool To Discover Antibiotic Resistance
87 Genes in Bacterial Genomes. *Antimicrob Agents Chemother.* 2014;58(1):212-20. doi:
88 10.1128/aac.01310-13.
- 89 3. McArthur AG, Waglechner N, Nizam F, Yan A, Azad MA, Baylay AJ, et al. The
90 comprehensive antibiotic resistance database. *Antimicrob Agents Chemother.*
91 2013;57(7):3348-57. Epub 2013/05/08. doi: 10.1128/AAC.00419-13. PubMed PMID:
92 23650175; PubMed Central PMCID: PMC3697360.

- 93 4. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
94 Bioinformatics. 2010;26(19):2460-1. Epub 2010/08/17. doi:
95 10.1093/bioinformatics/btq461. PubMed PMID: 20709691.
- 96 5. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment
97 search tool. Journal of Molecular Biology. 1990;215(3):403-10. doi:
98 [http://dx.doi.org/10.1016/S0022-2836\(05\)80360-2](http://dx.doi.org/10.1016/S0022-2836(05)80360-2).
- 99 6. Stenberg D. cURL. 1996.
- 100 7. HMP. Structure, function and diversity of the healthy human microbiome.
101 Nature. 2012;486(7402):207-14. Epub 2012/06/16. doi: 10.1038/nature11234.
102 PubMed PMID: 22699609; PubMed Central PMCID: PMC3564958.
- 103 8. Holt KE, Baker S, Weill FX, Holmes EC, Kitchen A, Yu J, et al. Shigella sonnei
104 genome sequencing and phylogenetic analysis indicate recent global dissemination
105 from Europe. Nat Genet. 2012;44(9):1056-9. Epub 2012/08/07. doi:
106 10.1038/ng.2369. PubMed PMID: 22863732; PubMed Central PMCID:
107 PMC3442231.
- 108