

1 **Internal control for process monitoring of clinical metagenomic next-generation sequencing of urine samples**

2

3 Victoria A. Janes ^{a#}, Jennifer S. van der Laan ^a, Sébastien Matamoros ^a, Daniel R. Mende ^a, Menno D. de Jong ^a and Constance Schultsz ^{a,b}

4

5 ^a Amsterdam UMC, University of Amsterdam, Medical Microbiology, Amsterdam, The Netherlands

6 ^b Amsterdam UMC, University of Amsterdam, Global Health - Amsterdam Institute for Global Health and Development (AIGHD), Amsterdam,
7 The Netherlands

8

9 Running head: Internal control for clinical mNGS of urine sample

10

11 #Corresponding author: Victoria A. Janes, v.a.janes@amsterdamumc.nl

12

13

14

15

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

ABSTRACT

Background

Process control for clinical metagenomic next-generation sequencing (mNGS) is not yet widely applied, while technical sources of bias are plentiful. We present an easy-to-use internal control (IC) method focussing on technical process control applied to metagenomics in clinical diagnostics.

Methods

DNA of nine urine samples was sequenced in the absence and presence of *Thermus thermophilus* DNA as IC in incremental concentrations (0.5-2-5%). Between aliquots of each sample, we compared the IC relative abundance (RA), and after *in silico* subtraction of IC reads, the microbiota and the RA of pathogens. The optimal IC spike-in concentration was defined as the lowest concentration still detectable in all samples.

Results

The RA of IC correlated linearly with the spiked IC concentration ($r^2=0.99$). IC added in a concentration of 0.5% of total DNA concentration was detectable in all samples, regardless of human/bacterial composition and after *in silico* removal gave the smallest difference in RA of pathogens compared to the unspiked aliquot of the sample. The microbiota of sample aliquots sequenced in the presence and absence of IC was highly similar after *in silico* removal of IC reads (median BC-dissimilarity per sample of 0.059), provided samples had sufficient bacterial read counts.

Conclusion

T. thermophilus DNA at a percentage of 0.5% of the total DNA concentration can be applied for the process control of mNGS of urine samples. We demonstrated negligible alterations in sample microbial composition after *in silico* subtraction of IC sequence reads. This approach contributes toward implementation of mNGS in the clinical microbiology laboratory.

39 INTRODUCTION

40 Metagenomic next-generation sequencing (mNGS) holds potential as a rapid pathogen detection tool for clinical microbiology diagnostics(1).

41 In mNGS, all DNA in a sample has an equal chance of being sequenced, including DNA from host and contaminant cells. This means the
42 interpretation of mNGS results can be challenging. For example, common reagent contaminants were previously incorrectly identified as
43 causative of infection(2, 3). Differences in sample type, starting DNA concentration of a sample, library preparation efficiency for GC-rich
44 organisms, and sequencing depth can all potentially bias the mNGS readout(4–7). If mNGS is to be safely used for clinical diagnostics,
45 potential sources of error or variation in library preparation and sequencing need to be considered and controlled for.

46 While blank extraction controls and mock community and/or positive control samples can be used to assess the level of contamination and
47 library preparation efficiency at batch level, external controls are not suitable for process monitoring of individual samples that vary greatly in
48 microbial and host composition, DNA concentration and possible inhibitors. To this end, an internal control (IC) consisting of low
49 concentration exogenous DNA that is completely unrelated to potential pathogens and the host microbiota, can be spiked into extracted DNA
50 from each sample and subsequently detected nested within the diagnostic test procedure. Detection of a constant relative abundance (RA) of IC
51 across all samples would ensure the library preparation and sequencing to be technically successful, ruling out technical error in samples where
52 no pathogen reads are detected.

53 The application of ICs as process control is standard practise in molecular diagnostic microbiology such as for pathogen-specific PCR(8), but is
54 still lacking in reported diagnostic applications of mNGS(2, 9–14). Studies that do describe ICs in metagenomic analyses typically applied fixed
55 amounts of short synthetic DNA aimed at quantification of microorganisms or quantification of competition between target DNA and host
56 DNA in library preparation protocols that include a DNA amplification step(2, 6, 15–18). It is unclear if the type and fragment length of IC
57 used for such protocols are directly applicable to PCR-free sequencing protocols, such as the protocol used for diagnostic mNGS of clinical
58 urine samples(19). Spike-in quantification standards are often patented and expensive while simpler and cheaper spike-ins may suffice for
59 technical process control. Most importantly, because the DNA concentration of clinical urine samples varies greatly, spiking samples with a
60 fixed amount of IC will result in overrepresentation of IC in low concentration samples and underrepresentation of IC reads in high
61 concentration samples making such a strategy unsuitable for process control.

62 We describe the design and validation of an IC procedure for the process control of PCR-free library preparation and mNGS of clinical urine
63 samples using full-length *Thermus thermophilus* DNA in a concentration that was titrated according to the DNA concentration of the urine
64 sample.

65

66 METHODS

67 *Urine sample collection*

68 As part of a larger study of diagnostic mNGS of urinary tract infections (UTIs), we collected consecutive urine samples that were sent to the
69 clinical microbiology laboratory of the Academic Medical Center, Amsterdam UMC for routine diagnostic semi-quantitative culture(19).
70 Personal data were anonymised and handled in compliance with the General Data Protection Regulation and medical-ethical guidelines of the
71 Academic Medical Center Amsterdam for anonymised use of patient materials. DNA extraction was performed after in-house enzymatic lysis,
72 using the automated NucliSENS easyMag platform (Biomérieux) according to manufacturer's instructions(19). These samples were sequenced
73 without IC and the DNA was stored at -80°C until further use. Nine representative samples of this collection were selected based on variations
74 in bacterial versus human read counts, culture results, and DNA concentration (table 1) and were sequenced again in the presence of
75 incremental concentrations of IC. The selected samples represented the extremes, i.e. both culture positive and negative high human/low
76 bacterial read samples (A, B, C), high bacterial/low human read samples (F, G), equal human/bacterial read samples (E), and culture negative
77 samples with a low DNA concentration (H, I).

78

79 *IC procedure*

80 We designed an IC with comparable GC-content and genome length to the common Gram-negative uropathogen *Escherichia coli* (~5.1Mb;
81 GC% 50.6; NCBI:taxid 562). Further, the IC should not be a human pathogen or part of the human microbiota or have genomic similarity to
82 these. The full-length bacterial genome of the extremophile *Thermus thermophilus* met these criteria and was previously successfully used as IC
83 in analyses of environmental river microbiota(20).

84 To allow for detection of IC DNA at a low constant RA across all samples, we titrated the IC concentration relative to the sample DNA
85 concentration. Therefore, DNA was first extracted from urine and quantified before addition of IC DNA at an appropriate concentration relative
86 to the total DNA concentration.

87 DNA concentrations of stored DNA aliquots (-80 C) of each sample were measured using the Qubit HS dsDNA quantification kit (Thermo
88 Fisher Scientific). The IC DNA concentrations were titrated to a final concentration of 0.5%, 2%, and 5% of the total DNA concentration of
89 the sample. Not all concentrations were tested in each sample if limited amounts of DNA were available. IC DNA and sample DNA were
90 mixed to a final DNA concentration of 1 ng/ μ l in a volume of 100 μ l (supplementary table 1). To assess variation between sequencing runs, two
91 aliquots of the same sample were sequenced without IC in both runs (B and B1 + 0% IC sample, Table 1). In addition, B1+0% IC sample
92 enabled us to see which proportion of reads were derived through contamination as this was the only sample sequenced without IC on the 2nd
93 run. A sample only containing IC (*T. thermophilus* DNA suspended in TE buffer 1x, 1ng/ μ l) was sequenced to confirm that *T. thermophilus*
94 reads were not misidentified as other taxa.

95

96 *Library preparation and sequencing*

97 Sample aliquots were sheared into 200bp fragments using the S220 Focused-ultrasonicator (Covaris). The Ion Xpress™ Plus Fragment Library
98 Kit (Thermo Fisher Science) was used for PCR-free, library preparation according to manufacturer's specifications. Sequencing was performed
99 on the Ion Torrent Proton and PGM (Thermo Fisher Scientific). This platform was used because of its sequencing speed which is required for
100 clinical diagnostics. The platform was set to produce 2 million reads per sample for run 1 (10 samples sequenced as part of 55 samples in total)
101 and 3 million reads per sample for run 2, both 200bp single-end (SE).

102

103 *Sequence analysis*

104 We used FastQC to check the read quality(21), Trimmomatic V0.38 to remove low quality reads(22), and BBMap to quantify and remove
105 human reads by aligning to the GRCh38 Human Genome obtained from Genbank(23, 24). The IC and bacterial reads were identified using
106 Kraken2 against the MiniKraken2_v1_8GB database, downloaded on the 2nd July 2019(25). The RA of IC was defined as the *Thermus* read
107 count at genus level divided by the total number of reads per sample (including human and unclassified reads). We assessed the RA of the IC

108 per sample and per spiked IC DNA concentration to identify the lowest IC DNA concentration that was still detectable across all samples. We
109 assessed the RA of IC in the B1+0% sample to assess the proportion of IC reads derived through cross-contamination. Next, all reads mapping
110 to the genus *Thermus* were removed *in silico* before we used Bracken to calculate the RA of all classified bacterial species by Bayesian re-
111 estimation of sequence reads(26). Bacterial species were annotated as urinary pathogens based on clinical microbiology common practice as
112 well as the results of a scoping review of the literature (Supplementary table 2).

113 For each sample, the RA of urinary pathogens was compared between aliquots sequenced in the presence and absence of IC. In case of a
114 polymicrobial composition, the cumulative RA of pathogens was used. The cumulative RA of pathogens was compared between sample
115 aliquots sequenced in the presence and absence of IC after *in silico* subtraction of IC reads. To assess whether sequencing samples in the
116 presence of IC and subsequent *in silico* subtraction of IC reads affected the sample microbiota, the RA of each detected species was compared
117 between aliquots by calculating the Bray-Curtis dissimilarity (BC-dissimilarity). The optimal spike-in concentration of IC DNA was defined as
118 the lowest still detectable IC concentration with minimum impact on the microbiota and cumulative RA of pathogens.

119

120

121 RESULTS

122

123 *Sequencing results*

124 We sequenced 28 aliquots of DNA obtained from 9 urine samples, spiked with incremental IC concentrations, in addition to one sample
125 containing 100% IC DNA. The mean read count per sample aliquot was 1,486,600 (SD=252,004) for the unspiked samples of the first run, and
126 4,093,176 reads (SD = 615,999) for the spiked samples sequenced in the second run (table 1). The read distributions per sample are depicted in
127 figure 1. In the samples with high human read count (samples A-D, H and I), the mean human read count varied between 1,2-1,4 million reads
128 and the mean bacterial read count varied between 2,442-244,708 per sample. For the samples with high bacterial read count (samples E and F)
129 the mean human read count per sample was 44,517 and 138,446 versus 3 and 3,4 million bacterial reads respectively. Resequencing the sample
130 with equal bacterial/human read count (sample G) again produced an equal distribution of human and bacterial reads (figure 1). Two aliquots of

131 sample B were sequenced without IC on both runs. The RA of bacterial species of the two aliquots correlated strongly (BC-dissimilarity 0.061),
132 indicating results from the two sequencing runs were comparable.

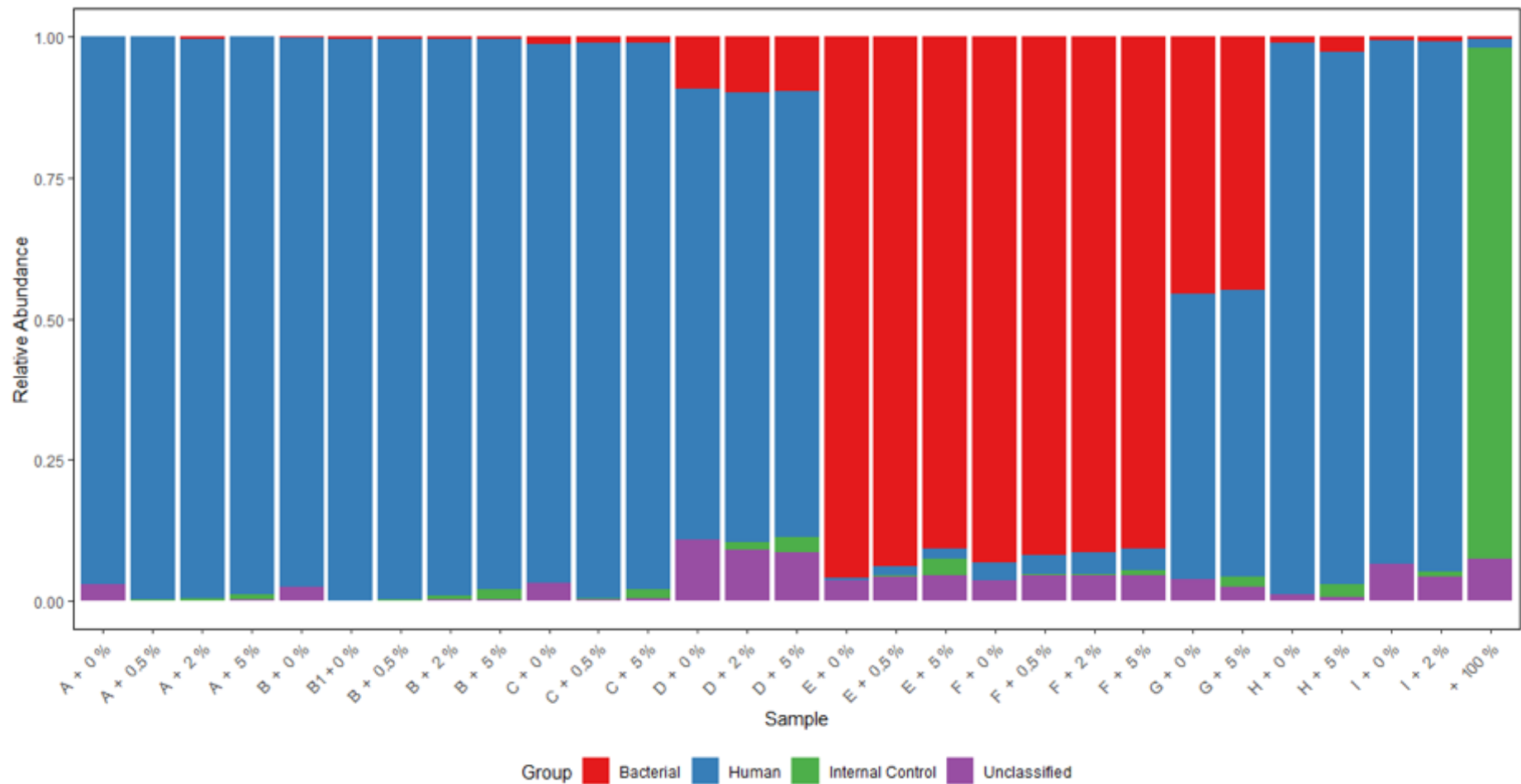
| Sample | Culture result | DNA concentration (ng/μl) | % spiked IC DNA | Total no. of reads | No. of IC reads | No. of unclassified reads | No. of bacterial reads | No. of human reads |
|--------|--|---------------------------|-----------------|--------------------|-----------------|---------------------------|------------------------|--------------------|
| A | <i>S. aureus</i> >10 ⁴ CFU/ml | 12,8 | 0 | 1800197 | 0 | 54004 | 686 | 1745507 |
| | | | 0.5 | 2834871 | 2598 | 3687 | 2761 | 2825825 |
| | | | 2 | 2906634 | 13076 | 4006 | 11769 | 2877783 |
| | | | 5 | 3476412 | 36171 | 7003 | 2124 | 3431114 |
| B | <i>P. aeruginosa</i> >10 ⁴ CFU/ml, <i>E. cloacae</i> >10 ⁴ CFU/ml | 17,68 | 0 | 1217911 | 0 | 31239 | 3697 | 1182975 |
| | | | 0 | 4230673 | 42 | 5609 | 16060 | 4208962 |
| | | | 0.5 | 3669184 | 6844 | 5283 | 13648 | 3643409 |
| | | | 2 | 3146673 | 22956 | 5870 | 12104 | 3105743 |
| C | Commensal flora 10 ⁴ CFU/ml | 10,68 | 5 | 3142635 | 58897 | 9333 | 11404 | 3063001 |
| | | | 0 | 1222137 | 0 | 39799 | 16904 | 1165434 |
| | | | 0.5 | 4060689 | 8404 | 10986 | 45521 | 3995778 |
| | | | 5 | 3525448 | 62633 | 14092 | 38950 | 3409773 |
| D | Commensal flora >10 ⁴ CFU/ml | 12,26 | 0 | 1018979 | 0 | 111370 | 94437 | 813172 |
| | | | 2 | 2469179 | 32993 | 225690 | 244708 | 1965788 |
| | | | 5 | 3107498 | 83231 | 271172 | 298285 | 2454810 |
| E | <i>E. coli</i> >10 ⁴ CFU/ml | 11,68 | 0 | 1276427 | 2 | 45820 | 1222817 | 7788 |
| | | | 0.5 | 3283575 | 8834 | 143238 | 3081986 | 49517 |
| | | | 5 | 3891596 | 117764 | 176733 | 3534129 | 62970 |
| F | <i>E. coli</i> >10 ⁵ CFU/ml, <i>K. pneumoniae</i> >10 ⁵ CFU/ml | 22,00 | 0 | 1529193 | 0 | 57574 | 1423548 | 48071 |
| | | | 0.5 | 3935399 | 3266 | 181040 | 3612043 | 139050 |
| | | | 2 | 4386781 | 12872 | 201483 | 4013053 | 159373 |
| | | | 5 | 3572208 | 29667 | 162892 | 3241806 | 137843 |
| G | <i>P. mirabilis</i> >10 ⁵ CFU/ml, <i>M. morgani</i> >10 ⁵ CFU/ml, <i>P. aeruginosa</i> <10 ⁴ CFU/ml | 9,78 | 0 | 1427110 | 0 | 56789 | 649441 | 720880 |
| | | | 5 | 2860345 | 54974 | 71919 | 1285320 | 1448132 |
| H | Commensal flora 10 ⁴ CFU/ml | 8,2 | 0 | 1395035 | 0 | 17319 | 15325 | 1362391 |
| | | | 5 | 3628947 | 76747 | 29465 | 96711 | 3426024 |
| I | Commensal flora 10 ⁴ CFU/ml | 7,94 | 0 | 1624302 | 0 | 108709 | 12326 | 1503267 |
| | | | 2 | 2780217 | 29602 | 117812 | 23174 | 2609629 |
| IC | <i>T. thermophilus</i> | 1,00 | 100 | 2680871 | 2425350 | 199119 | 10558 | 45844 |

133

134 **Table 1. Characteristics of the sequenced urine samples.** The total number of reads, IC, unclassified, bacterial and human reads per sample

135 aliquot are stated. Semi-quantitative culture results are given in colony forming units per ml of urine (CFU/ml).

136



137

138

139 **Figure 1. Read distributions.** Shown is the total read distribution of each sequenced aliquot, indicating the RA of bacterial (red), human
140 (blue), IC (green) and unclassified (purple) reads. Following the sample name (A – IC) is the spiked IC concentration (%).

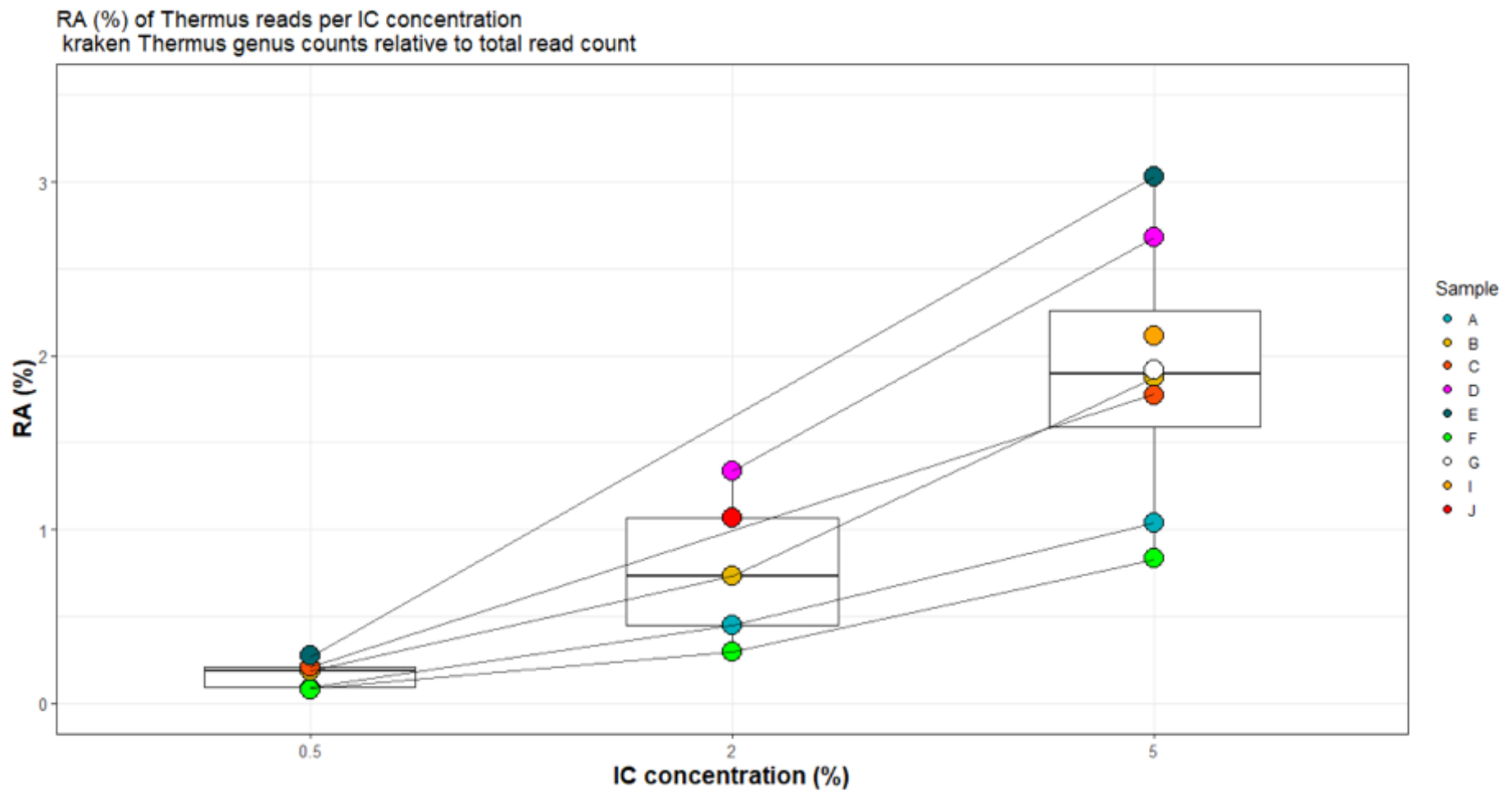
141

142 *Detection of IC*

143 The resequencing of an aliquot of sample B without IC on the 2nd run (amidst samples sequenced in the presence of IC) not only allowed
144 assessment of between-run variability but was also used to assess the proportion of IC reads derived through contamination. Out of 4,230,673
145 reads in this aliquot (aliquot name B1+0%), 42 reads (0.001%) identified as *Thermus* at genus level of which 39 identified as *T. thermophilus*.

146 To assess whether *T. thermophilus* reads mapped against other bacterial species, a sample containing 100% IC DNA was sequenced. Of all
147 2,635,027 reads, 2,425,350 reads (92%) mapped to the *Thermus* genus and of these 1,885,710 (71,6% of total) were classified as *T.*

148 *thermophilus*. 2,544 (0.096%) reads mapped to other bacterial species, predominantly to *E. coli* and *Klebsiella pneumoniae*, 8,014 reads (0.3%)
149 mapped to the domain Bacteria, 45,844 reads identified as human (1.7%), and 199,119 (7.5%) reads remained unclassified.
150



151

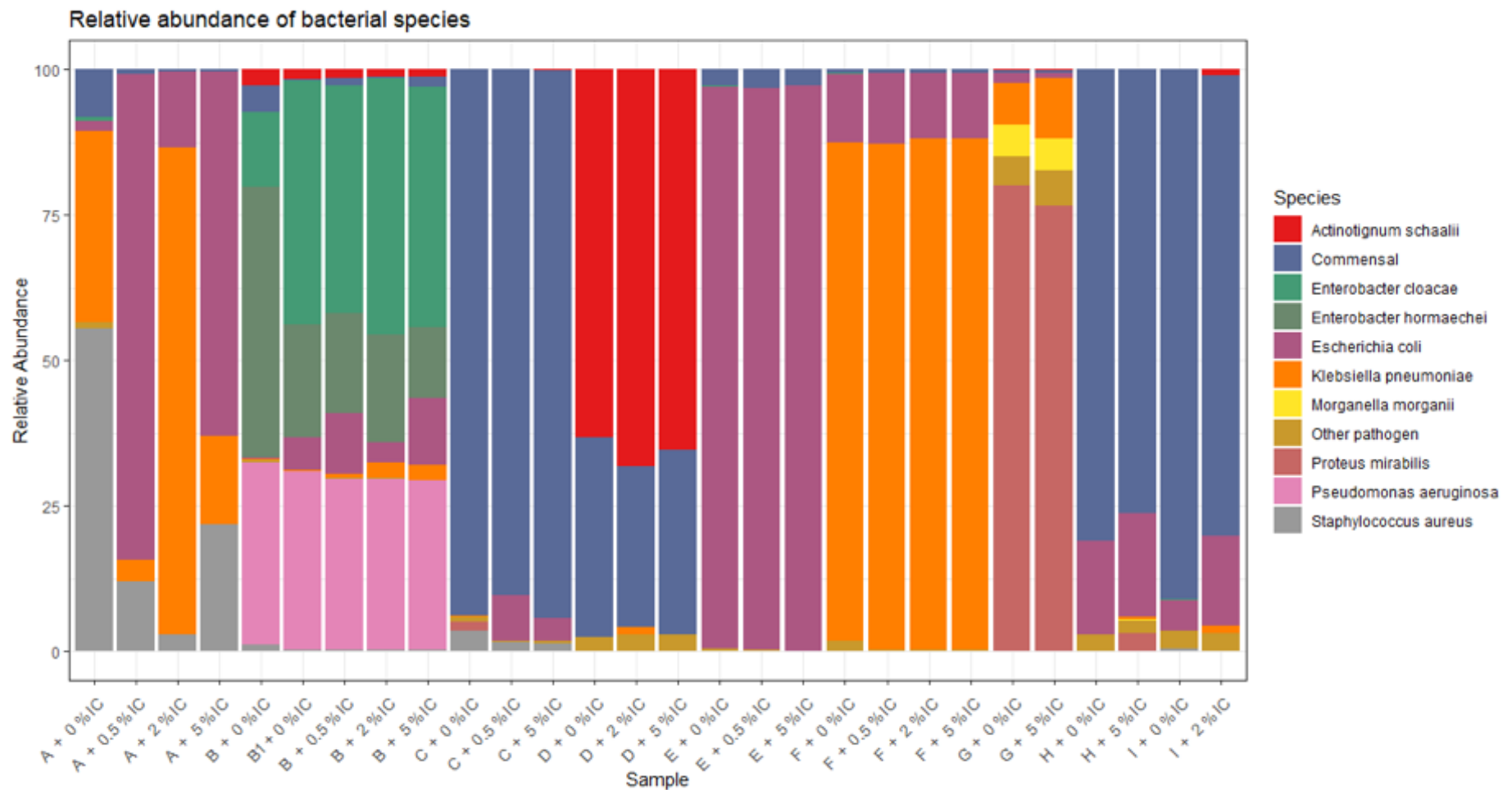
152 **Figure 2 RA of IC per sample and per IC concentration.**

153

154 Next, we assessed the RA of IC reads per sample at IC DNA concentrations of 0.5, 2 and 5% of the total DNA concentration. The RA of IC
155 reads rose linearly with increasing concentrations of spiked IC DNA (Pearson's $r^2=0.99$) and was 10.24 times higher in the samples spiked with
156 IC DNA at 5% of the total DNA concentration compared to the samples spiked with IC DNA at 0.5% of the total DNA concentration (figure 2).
157 There was no significant difference in RA of IC between samples with a high or equal bacterial read count (E, F and G) compared to samples
158 with a low bacterial and high human read count (A-D, H and I) (Welch two sample t-test, $p = 0.8$).

159 IC DNA spiked into samples at concentration of 0.5% of the total DNA concentration was detectable in all samples (median RA 0.19%; IQR
160 0.12%; range 0.09-0.27%). The median RA of IC reads at this concentration was 187 times higher than the RA of the IC reads in the aliquot of
161 sample B that was sequenced in the absence of IC on run 2, and thus clearly exceeded the RA of IC expected to be derived through cross-
162 contamination.

163 After *in silico* subtraction of IC reads, the RA of the bacterial species reads was highly similar between sample aliquots sequenced in the
164 presence and absence of IC, illustrated by the similar species distribution between aliquots of the same sample seen in figure 3 and by the
165 median BC-dissimilarity per sample of 0.059 (IQR 0.04, range 0.008-0.61). The exception was sample A with a BC-dissimilarity of 0.61. The
166 aliquots of this sample had a median bacterial read count of only 2442 (range 686-11,769), over 61 times lower than the median bacterial read
167 count of all other samples (median 150,363; range 12,043 (sample B) – 3,426,924 (sample F)).



168

169 **Figure 3. Relative abundance of bacterial pathogenic species and commensals reads per sample aliquot.** Shown is the RA of the 9 most
170 abundant bacterial pathogens classified after *in silico* removal of IC. Commensal bacterial species were grouped (ochre yellow).

171

172 Since the main outcome measure for diagnostic mNGS of urine samples was the RA of bacterial pathogens, we calculated the difference in RA
173 of bacterial pathogens between spiked and unspiked sample aliquots after *in silico* subtraction of IC reads (table 2). An IC DNA concentration
174 of 0.5% gave the smallest difference in cumulative RA of bacterial pathogens between spiked and unspiked sample aliquots, and was still
175 detected in all samples. The median difference in RA was 1.15% (table 2).

176

| Relative Abundance of IC (%) | | | | |
|------------------------------|------|------|--------|------|
| IC concentration (%) | Min | Max | Median | IQR |
| 0.5 | 0.17 | 6.72 | 1.15 | 0.82 |
| 2.0 | 0.28 | 7.94 | 4.14 | 4.53 |
| 5.0 | 0.15 | 4.16 | 1.22 | 2.54 |

177

178 **Table 2. The median difference in cumulative RA of pathogens between spiked and unspiked aliquots of a sample after *in silico***
179 **removal of IC reads.** Results are shown per IC concentration.

180

181

182 DISCUSSION

183

184 In this study we successfully applied *T. thermophilus* DNA as IC for the process control of mNGS of clinical urine samples. DNA aliquots
185 extracted from urine samples were sequenced in the absence and presence of incremental concentrations of IC DNA. IC DNA added at a
186 concentration of 0.5% of the total sample DNA concentration was detectable in all samples, did not alter the microbial composition of the
187 sample DNA or the RA of detected bacterial pathogens substantially after *in silico* removal of IC reads. The results presented here are directly
188 relevant to clinical microbiology practise as they are derived from clinical urine samples, representing the full array of variation seen in
189 microbial and host DNA composition and concentration.

190 Other studies described the use of fixed amounts of short fragment synthetic DNA as IC aimed at quantification or as a measure for pathogen
191 detection sensitivity(2, 6, 15–17). For such uses of IC, the RA of IC in the readout is integral to the analysis, making this a different approach to
192 using IC as process control, where it is desirable to detect a constant low RA of IC in each sample to ensure the sequencing process was
193 technically successful. The DNA concentration of urine samples is highly variable, meaning spiking a fixed amount of IC, as done in these

194 other studies, can easily lead to over- or underrepresentation of IC depending on the sample DNA concentration. Moreover, short fragments of
195 IC could be sheared to even shorter fragments during library preparation and could consequently be lost during size selection when targeting
196 species that require longer fragmentation time than the IC(20). The spiking of samples with IC in a concentration relative to the total DNA
197 concentration of that sample, ensured detection of IC regardless of the concentration of the sample and thus allowed for consistent quality
198 control of each sample.

199

200 Many bacteria that cause UTI, such as *E. coli*, are commensal to the genito-urinary tract and are only considered causative of disease when
201 present in concentrations above a certain threshold(27). This study was not designed to address quantification, but quantification of urine
202 mNGS results is needed to align with current routine culture-based microbiology diagnostics whereby a diagnosis is established based on semi-
203 quantitative cultures.

204

205 Two reads were initially identified as *T. thermophilus* in sample D+0% sequenced on run 1 in the absence of IC. No samples were spiked with
206 IC on this run. We hypothesised that this finding was caused by misclassification of Taq-polymerase DNA, derived from *Thermus aquaticus*,
207 which may have regions of genomic similarity to *T. thermophilus*. However, BLASTn identified these 2 reads as *E. coli* plasmids p94EC-5.

208

209 When assessing the clinical relevance of species present at very low relative abundance, the likelihood that these reads are present as a result of
210 cross-contamination should be taken into consideration. Relative abundances of 0.05-2.78% have been described for droplet cross-
211 contamination during sample preparation(28). Barcode hopping, the incorrect assignment of library molecules from the expected barcode to a
212 different barcode in a multiplexed pool, was observed in biological mock community samples at a rate of 0.033% on the IonTorrent PGM
213 platform which has similar technology to the IonTorrent Proton(29). We demonstrate that the inclusion of IC DNA allowed for the assessment
214 of cross-contamination by sequencing replicates in the presence and absence of IC. The impact of cross-contamination and barcode hopping
215 appeared minimal in our setting. We propose that such a sample sequenced in the absence of IC should be included in each library preparation
216 and sequencing run, to assess cross-contamination and to determine the minimum threshold of the RA of IC to consider a sample successfully
217 sequenced.

218

219 By resequencing biological replicates in the presence and absence of IC divided over 2 sequencing runs, we demonstrated negligible alterations
220 in sample microbial and pathogen composition after *in silico* subtraction of IC reads for 8 out of 9 samples. The incongruent sample A
221 consisted of >95% human reads, leaving only a median of 2,442 (range 686-11,769) classified bacterial reads per sequenced aliquot. The low
222 bacterial read count can explain the poor correlation between sequenced aliquots and this data may help set guidelines for a minimum
223 sequencing depth for clinical samples and stresses the need for removal of host cells prior to sequencing(28, 29). Despite the variation in RA of
224 bacterial species in sample A, the IC was detected in all aliquots, indicating sequencing and library preparation processes had been technically
225 successful.

226 Our study has some limitations. Not every sample could be sequenced with each IC concentration, because limited volumes of DNA were
227 present for some samples. However the overall results were consistent and RA of IC correlated linearly with spiked IC concentrations, implying
228 our results can be extrapolated to other concentrations. Even lower spike-in concentrations could be suitable but were not tested in sufficient
229 numbers. One aliquot of sample B was spiked with 0.05% IC DNA, yielding a RA of 0.025% (747 reads; data not shown).

230 Our approach does not control for DNA extraction efficiency. To achieve this, samples could be spiked with bacterial cells prior to DNA
231 extraction. However, this comes with additional pitfalls. Spiking samples with a quantified bacterial cell suspension in a concentration of 0.5%
232 relative to the sample will not necessarily result in 0.5% IC DNA. The human genome is several orders of magnitude longer than bacterial
233 genomes, meaning each human cell produces approximately 1000 times more DNA than each bacterial cell(6). Given the number of
234 sequencing reads produced per sequencing run is finite, host DNA will easily outcompete bacterial DNA, leading to difficulty in interpreting
235 the mNGS read out: in cases where no pathogen reads are detected, distinguishing between a true negative (at the given sequencing depth) and a
236 technical fail would be impossible. The spiking of samples with IC DNA in a concentration relative to the total DNA concentration of that
237 sample mitigated this problem and ensured detection of IC regardless of the concentration or composition of the sample and thus allowed for
238 consistent process control of each sample.

239

240 In conclusion we developed an IC for the process control of mNGS of urine samples that at a percentage of 0.5% of the total DNA
241 concentration was detectable in all samples, regardless of the sample composition or DNA concentration. By resequencing sample aliquots in
242 the presence and absence of IC, we demonstrated negligible alterations in sample microbial and pathogen composition after *in silico* subtraction

243 IC. We showed how a lower detection threshold for detection of IC can be established taking contamination into account. This approach could
244 be used as process control for diagnostic mNGS and contributes toward implementation of mNGS in the clinical microbiology laboratory.

245

246 Acknowledgements

247 This study was supported by the COMPARE Consortium, which received funding from the European Union's Horizon 2020 research and
248 innovation programme under grant agreement No. 643476. The authors declare no conflict of interest.

249

250

251 REFERENCES

252

- 253 1. Chiu CY, Miller SA. 2019. Clinical metagenomics. *Nat Rev Genet* 20:341–355.
- 254 2. Wilson MR, O’Donovan BD, Gelfand JM, Sample HA, Chow FC, Betjemann JP, Shah MP, Richie MB, Gorman MP, Hajj-Ali RA,
255 Calabrese LH, Zorn KC, Chow ED, Greenlee JE, Blum JH, Green G, Khan LM, Banerji D, Langelier C, Bryson-Cahn C, Harrington W,
256 Lingappa JR, Shanbhag NM, Green AJ, Brew BJ, Soldatos A, Strnad L, Doernberg SB, Jay CA, Douglas V, Josephson SA, DeRisi JL.
257 2018. Chronic Meningitis Investigated via Metagenomic Next-Generation Sequencing. *JAMA Neurol* 75:947–955.
- 258 3. Hornung BVH, Zwittink RD, Kuijper EJ. 2019. Issues and current standards of controls in microbiome research. *FEMS Microbiol Ecol*
259 95.
- 260 4. Videnska P, Smerkova K, Zwinsova B, Popovici V, Micenkova L, Sedlar K, Budinska E. 2019. Stool sampling and DNA isolation kits
261 affect DNA quality and bacterial composition following 16S rRNA gene sequencing using MiSeq Illumina platform. *Sci Rep* 9:13837.
- 262 5. Bowers RM, Clum A, Tice H, Lim J, Singh K, Ciobanu D, Ngan CY, Cheng J-F, Tringe SG, Woyke T. 2015. Impact of library
263 preparation protocols and template quantity on the metagenomic reconstruction of a mock microbial community. *BMC Genomics*
264 16:856.
- 265 6. Schlaberg R, Chiu CY, Miller S, Procop GW, Weinstock G. 2017. Validation of Metagenomic Next-Generation Sequencing Tests for
266 Universal Pathogen Detection. *Arch Pathol Lab Med* 141:776–786.
- 267 7. Ye SH, Siddle KJ, Park DJ, Sabeti PC. 2019. Benchmarking Metagenomics Tools for Taxonomic Classification. *Cell* 178:779–794.
- 268 8. Deer DM, Lampel KA, González-Escalona N. 2010. A versatile internal control for use as DNA in real-time PCR and as RNA in real-
269 time reverse transcription PCR assays. *Lett Appl Microbiol* 50:366–372.
- 270 9. Mongkolrattanothai K, Naccache SN, Bender JM, Samayoa E, Pham E, Yu G, Dien Bard J, Miller S, Aldrovandi G, Chiu CY. 2017.
271 Neurobrucellosis: Unexpected Answer From Metagenomic Next-Generation Sequencing. *J Pediatric Infect Dis Soc* 2017/01/06. 6:393–
272 398.
- 273 10. Hu Z, Weng X, Xu C, Lin Y, Cheng C, Wei H, Chen W. 2018. Metagenomic next-generation sequencing as a diagnostic tool for

- 274 toxoplasmic encephalitis. *Ann Clin Microbiol Antimicrob* 17:45.
- 275 11. Ivy MI, Thoendel MJ, Jeraldo PR, Greenwood-Quaintance KE, Hanssen AD, Abdel MP, Chia N, Yao JZ, Tande AJ, Mandrekar JN,
276 Patel R. 2018. Direct Detection and Identification of Prosthetic Joint Infection Pathogens in Synovial Fluid by Metagenomic Shotgun
277 Sequencing. *J Clin Microbiol* 56.
- 278 12. Jun C, Huan H, Wei F, Duozhi S, Chen L, Yang S, GuoFeng G, Hao W, Qian Z, LiQing W, HongLong W, Long H, Luyao C, Jin Z,
279 Shela L, FeiYan W, Zhou Z. 2019. Detection of pathogens from resected heart valves of patients with infective endocarditis by next-
280 generation sequencing. *Int J Infect Dis*.
- 281 13. Wilson MR, Naccache SN, Samayoa E, Biagtan M, Bashir H, Yu G, Salamat SM, Somasekar S, Federman S, Miller S, Sokolic R,
282 Garabedian E, Candotti F, Buckley RH, Reed KD, Meyer TL, Seroogy CM, Galloway R, Henderson SL, Gern JE, DeRisi JL, Chiu CY.
283 2014. Actionable diagnosis of neuroleptospirosis by next-generation sequencing. *N Engl J Med* 2014/06/04. 370:2408–2417.
- 284 14. Hasman H, Saputra D, Sicheritz-Ponten T, Lund O, Svendsen CA, Frimodt-Møller N, Aarestrup FM. 2014. Rapid Whole-Genome
285 Sequencing for Detection and Characterization of Microorganisms Directly from Clinical Samples. *J Clin Microbiol* 52:139 LP – 146.
- 286 15. Stämmler F, Gläsner J, Hiergeist A, Holler E, Weber D, Oefner PJ, Gessner A, Spang R. 2016. Adjusting microbiome profiles for
287 differences in microbial load by spike-in bacteria. *Microbiome* 4:28.
- 288 16. Turlousse DM, Yoshiike S, Ohashi A, Matsukura S, Noda N, Sekiguchi Y. 2017. Synthetic spike-in standards for high-throughput 16S
289 rRNA gene amplicon sequencing. *Nucleic Acids Res* 45:e23–e23.
- 290 17. Bal A, Pichon M, Picard C, Casalegno JS, Valette M, Schuffenecker I, Billard L, Vallet S, Vilchez G, Cheynet V, Oriol G, Trouillet-
291 Assant S, Gillet Y, Lina B, Brengel-Pesce K, Morfin F, Josset L. 2018. Quality control implementation for universal characterization of
292 DNA and RNA viruses in clinical respiratory samples using single metagenomic next-generation sequencing workflow. *BMC Infect Dis*
293 18:537.
- 294 18. Hardwick SA, Chen WY, Wong T, Kanakamedala BS, Deveson IW, Ongley SE, Santini NS, Marcellin E, Smith MA, Nielsen LK,
295 Lovelock CE, Neilan BA, Mercer TR. 2018. Synthetic microbe communities provide internal reference standards for metagenome
296 sequencing and analysis. *Nat Commun* 9:3096.

- 297 19. Janes VA, Matamoros S, Willemse N, Visser CE, de Wever B, Jakobs ME, Subramanian P, Hasan NA, Colwell RR, de Jong MD,
298 Schultsz C. 2017. Metagenomic sequencing to replace semi-quantitative urine culture for detection of urinary tract infections: a proof of
299 concept. bioRxiv.
- 300 20. Satinsky BM, Gifford SM, Crump BC, Moran MA. 2013. Use of Internal Standards for Quantitative Metatranscriptome and Metagenome
301 Analysis. *Methods Enzymol* 531:237–250.
- 302 21. Andrews S. 2010. FastQC: a quality control tool for high throughput sequence data.
- 303 22. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* 30:2114–2120.
- 304 23. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Wheeler DL. 2006. GenBank. *Nucleic Acids Res* 34:D16–D20.
- 305 24. Bushnell B. 2014. BBMap : A Fast , Accurate , Splice-Aware Aligner. *LBNL Dep Energy Jt Genome Inst* 3–5.
- 306 25. Wood DE, Salzberg SLS, Venter C, Remington K, Heidelberg J, Halpern A, Rusch D, Eisen J, Wu D, Paulsen I, Nelson K, Nelson W,
307 Fouts D, Levy S, Knap A, Lomas M, Nealon K, White O, Peterson J, Hoffman J, Parsons R, Baden-Tillson H, Pfannkoch C, Rogers Y-
308 H, Smith H, Tyson G, Chapman J, Hugenholtz P, Allen E, Ram R, Richardson P, Solovyev V, Rubin E, Rokhsar D, Banfield J,
309 Huttenhower C, Gevers D, Knight R, Abubucker S, Badger J, Chinwalla A, Creasy H, Earl A, FitzGerald M, Fulton R, Giglio M,
310 Hallsworth-Pepin K, Lobos E, Madupu R, Magrini V, Martin J, Mitreva M, Muzny D, Sodergren E, Versalovic J, Wollam A, Worley K,
311 Wortman J, Young S, Zeng Q, Aagaard K, Abolude O, Allen-Vercoe E, Alm E, Alvarado L, Altschul S, Gish W, Miller W, Myers E,
312 Lipman D, Brady A, Salzberg SLS, Huson D, Auch A, Qi J, Schuster S, Brady A, Salzberg SLS, Rosen G, Garbarine E, Caseiro D,
313 Polikar R, Sokhansanj B, Liu B, Gibbons T, Ghodsi M, Treangen T, Pop M, Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson
314 O, Huttenhower C, Treangen T, Koren S, Sommer D, Liu B, Astrovsкая I, Ondov B, Darling A, Phillippy A, Pop M, Ames S, Hysom
315 D, Gardner S, Lloyd G, Gokhale M, Allen J, Kindblom C, Davies J, Herzberg M, Svensäter G, Wickström C, Foweraker J, Cooke N,
316 Hawkey P, Könönen E, Saarela M, Karjalainen J, Jousimies-Somer H, Alaluusua S, Asikainen S, Camacho C, Coulouris G, Avagyan V,
317 Ma N, Papadopoulos J, Bealer K, Madden T, Zimin A, Marçais G, Puiu D, Roberts M, Salzberg SLS, Yorke J, Pruitt K, Tatusova T,
318 Brown G, Maglott D, Marçais G, Kingsford C, Roberts M, Hayes W, Hunt B, Mount S, Yorke J, Magoc T, Pabinger S, Canzar S, Liu X,
319 Su Q, Puiu D, Tallon L, Salzberg SLS, Mavromatis K, Ivanova N, Barry K, Shapiro H, Goltsman E, McHardy A, Rigoutsos I, Salamov
320 A, Korzeniewski F, Land M, Lapidus A, Grigoriev I, Richardson P, Hugenholtz P, Kyrpides N, Ondov B, Bergman N, Phillippy A. 2014.

- 321 Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol* 15:R46.
- 322 26. Lu J, Breitwieser FP, Thielen P, Salzberg SL. 2017. Bracken: estimating species abundance in metagenomics data. *PeerJ Comput Sci*
323 3:e104.
- 324 27. KASS EH. 1957. Bacteriuria and the Diagnosis of Infections of the Urinary Tract: With Observations on the Use of Methionine as a
325 Urinary Antiseptic. *AMA Arch Intern Med* 100:709–714.
- 326 28. Pereira-Marques J, Hout A, Ferreira RM, Weber M, Pinto-Ribeiro I, van Doorn L-J, Knetsch CW, Figueiredo C. 2019. Impact of Host
327 DNA and Sequencing Depth on the Taxonomic Resolution of Whole Metagenome Sequencing for Microbiome Analysis. *Front*
328 *Microbiol* 10:1277.
- 329 29. Charalampous T, Kay GL, Richardson H, Aydin A, Baldan R, Jeanes C, Rae D, Grundy S, Turner DJ, Wain J, Leggett RM, Livermore
330 DM, O’Grady J. 2019. Nanopore metagenomics enables rapid clinical diagnosis of bacterial lower respiratory infection. *Nat Biotechnol*
331 37:783–792.

332