bioRxiv preprint doi: https://doi.org/10.1101/2020.06.25.170423; this version posted June 25, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

1	Optimization of cerebrospinal fluid microbial metagenomic sequencing diagnostics
2	
3	Josefin Olausson ¹ , Sofia Brunet ¹ , Diana Vracar ^{1,2} , Yarong Tian ² , Sanna Abrahamsson ² , Sri
4	Harsha Meghadri ¹ , Per Sikora ³ , Maria Lind Karlberg ⁴ , Hedvig Engström Jakobsson ^{1*} , Ka-Wei
5	Tang ^{1,2}
6	
7	¹ Department of Clinical Microbiology, Sahlgrenska University Hospital, Region Västra
8	Götaland, Gothenburg, Sweden
9	² Wallenberg Centre for Molecular and Translational Medicine, Department of Infectious
10	Diseases, Institute of Biomedicine, University of Gothenburg, Gothenburg, Sweden
11	³ Clinical Genomics Gothenburg, Science for Life Laboratories, Gothenburg, Sweden
12 13	⁴ Department of Microbiology, Public Health Agency of Sweden, Solna, Sweden
14	JO: josefin.olausson@gu.se
15	SB: <u>sofia.brunet@gu.se</u>
16	YT: <u>yarong.tian@gu.se</u>
17	DV: <u>diana.vracar@gu.se</u>
18	SA: <u>sanna.abrahamsson@gu.se</u>
19	SHM: <u>harshameghadri@gu.se</u>
20	PS: <u>per.sikora@gu.se</u>
21	MLK: maria.lind.karlberg@folkhalsomyndigheten.se
22	KWT: <u>kawei.tang@gu.se</u>
23	HEJ: <u>hedvig.jakobsson@gu.se</u> *Corresponding author
24	Short title: Metagenomic sequencing of cerebrospinal fluid
25	
26	

27 Abstract

28 Background

Infection in the central nervous system is a severe condition associated with high morbidity and mortality. Despite ample testing, the majority of encephalitis and meningitis cases remain undiagnosed. Metagenomic sequencing of cerebrospinal fluid has emerged as an unbiased approach to identify rare microbes and novel pathogens. However, several major hurdles remains, including establishment of individual limits of detection, removal of false positives and implementation of universal controls.

35 **Results**

36 Twenty-one cerebrospinal fluid samples, in which a known pathogen had been positively 37 identified by available clinical techniques, were subjected to metagenomic DNA sequencing using massive parallel sequencing. Fourteen samples contained minute levels of Epstein-Barr 38 39 virus. Calculation of the detection threshold for each sample was made using total leukocyte 40 content in the sample and environmental contaminants found in bioinformatic classifiers. 41 Virus sequences were detected in all ten samples, in which more than one read was expected 42 according to calculations. Conversely, no viral reads were detected in seven out of eight 43 samples, in which less than one read was expected according to calculations. False positive pathogens of computational or environmental origin were readily identified, by using a 44 45 commonly available cell control. For bacteria additional filters including a comparison 46 between classifiers removed the remaining false positives and alleviated pathogen 47 identification.

48 Conclusions

Here we show a generalizable method for detection and identification of pathogen species using metagenomic sequencing. The sensitivity for each sample can be calculated using the leukocyte count and environmental contamination. The choice of bioinformatic method

53 Identification of pathogens require multiple filtering steps including read distribution,

sequence diversity and complementary verification of pathogen reads.

55

56 Keywords

57 Metagenomics, Cerebrospinal fluid, Pathogen classification, PaRCA, Epstein-Barr virus

58

59 Background

60 Infections in the central nervous system (CNS) are severe and despite extensive microbiological diagnostic analysis a causative pathogen cannot be identified in many of the 61 62 cases. A majority of CNS infections are caused by viruses, such as herpes simplex virus 1 63 (HSV1), varicella zoster virus (VZV or human herpesvirus 3) and enterovirus [1, 2]. Among 64 CNS infections, Streptococcus pneumoniae and Neisseria meningitidis are the most common 65 pathogens, while fungal or parasitic meningitis CNS infections are less common [3]. Epstein-66 Barr virus (EBV) has been implicated in recurrent meningitis and chronic encephalitis [4]. 67 However, due to the high prevalence of EBV and its ability to remain latent in B-lymphocytes 68 after primary infection and its role in tumorigenesis, assessing the clinical relevance of EBV 69 DNA detected in cerebrospinal fluid (CSF) is difficult and presence of EBV is often 70 considered to be an benign incidental finding [5, 6].

71

Current microbiological diagnostic methods include cultivation and nucleic acid detection of CSF, which are restricted to prior knowledge of the putative causing agent. Cultivation can detect a wide range of microorganisms, however, it is limited to viable and culturable pathogens. In contrast, nucleic acid detection is rapid and highly sensitive, but constrained to genetically conserved regions of known pathogens. Metagenomic sequencing using massive parallel sequencing, has the capability to discern multiple species and identify unknown bioRxiv preprint doi: https://doi.org/10.1101/2020.06.25.170423; this version posted June 25, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

species in samples. In metagenomics, the total nucleic acid present in the clinical sample is
sequenced, thus provides an unbiased tool to diagnose infections and unknown species in
samples [7-12].

81

82 Currently there is no standard for metagenomic sequencing in a clinical setting and the 83 technique is still faced with some major challenges [13]. Contrary to PCR, the sensitivity in 84 metagenomic sequencing is dependent on the fraction of pathogen sequences in the total 85 sequencing library. Furthermore, laboratory contaminations detected in sequencing have been 86 shown to differ greatly between laboratories and be dependent on the input biomass [14, 15]. 87 Nucleic acid derived from the host and environmental contaminants must therefore be taken 88 into account. Previous studies have calculated the sensitivity by using dilution of an 89 exogenous pathogen into a known quantity of host background. However, this does not take 90 into account the variability of clinical samples nor does it provide any guidance on how the 91 sensitivity of each sample should be calculated.

92

93 Bioinformatic pathogen identification is a second major obstacle. Several publically available 94 bioinformatic tools for classification are available, such as Centrifuge, Kraken and PathSeq 95 [16-18]. Two conceptual different methods are frequently used, alignment of single reads (e.g. 96 BLAST), or assemblies (k-mers), against pathogen databases. The list of pathogens generated 97 by these applications are often long and requires exhaustive examinations in order to discern 98 the true pathogen from bioinformatic misclassification and environmental contaminations. 99 Criteria for identifying the causative pathogen include sequences disseminated throughout the 100 microbial genome of the proposed pathogen, a threshold for number of pathogen reads in 101 relation to total number of reads, and confirmation using several alignment algorithms have been suggested to increase the specificity [19, 20]. Each laboratory does however apply theirown criteria.

104

105 We investigated the robustness of microbial metagenomics for clinical diagnostics of CNS-106 infections. To evaluate the diagnostic performance of the method, 21 CSF samples with 107 variable levels of known pathogens were sequenced with the aim to identify factors important 108 for calculating sensitivity. Also, four different taxonomic classifiers were assessed for their 109 efficiency to identify pathogens as well as the number of false positive pathogens identified. 110 Two commonly available cell lines were implemented as a positive and negative control to 111 support the removal of environmental contaminants and bioinformatic misclassifications. 112 Pathogen detection in DNA metagenomic sequencing in CSF is mainly limited by the 113 leukocyte count which affects the sensitivity and bioinformatic missclassifications which 114 affects the efficiency of pathogen identification.

115

116 **Results**

117 We implemented a metagenomic DNA sequencing methodology to unbiasedly detect 118 microbial species in CSF samples from patients with CNS symptoms in which a pathogen or 119 EBV had been detected (Additional Table 1). Samples positively identified with pathogen-120 specific quantitative PCR (qPCR), 16S rRNA gene sequencing or bacterial/mycotic culture in 121 CSF were included. Different pathogen types and variation of viral loads were chosen. CSF 122 samples containing low levels of EBV were chosen to establish the sensitivity of the method. 123 DNA from each sample was extracted and fragmented before library preparation and 124 sequenced using massive parallel sequencing. Datasets were processed using five 125 bioinformatic tools (Additional Figure 1).

126

127 **Bioinformatic classifiers**

Four bioinformatic classifiers were included, Kraken2, Centrifuge, our in-house developed PaRCA (Pathogen detection for Research and Clinical Applications) and CosmosID. CosmosID was tested mainly for its ability to generate concise pathogen lists, but the format of the platform prevented a detailed analysis of the raw data and was therefore not included in all comparisons in the manuscript. The four bioinformatic classifiers diverged with regards to fraction of processed reads (from 85%-100%, Additional Table 2-3). However, the ability to identify the primary pathogen was similar comparing the classifiers.

135

136 Sensitivity

137 Initially, three CSF samples (Sample 1-3) with high virus load of herpesvirus were analyzed. 138 HSV1 and VZV were detected by all bioinformatic classifiers (Table 1). In sample 1, HSV1 was positively identified at 1×10^4 genome equivalents per milliliter (Geq/ml) using qPCR. 139 140 The sequencing library consisting of more than 15 million reads contained 6.2-7.2 HSV1 141 reads per million sequences analyzed (parts per million; ppm). The following two samples 142 originated from patients with similar values of VZV DNA levels quantified by qPCR (1.9 and 3.9×10^5 Geq/ml). Despite equivalent levels a ten-fold difference in detected VZV reads was 143 144 observed between sample two (15-16 ppm) and sample three (135-147 ppm). Sample 2 contained 272×10^6 white blood cells (WBC) per liter compared with sample 3 which 145 contained 17×10^6 WBC per liter (Table 1). We hypothesized that the difference in sensitivity 146 147 was related to variations in leukocyte composition in the sample.

To further test the sensitivity, two CSF samples containing JC polyomavirus (JCV), a DNA virus with a relatively small genome, were processed. One sample contained high virus levels $(1.9 \times 10^5 \text{ Geq/ml})$ and the other low virus levels $(4.3 \times 10^3 \text{ Geq/ml})$ (Sample 4-5). JCV DNA was readily detected in both samples ranging from 1757-2096 ppm in sample
4 and 40-57 ppm in sample 5.

153	In order to verify that the methodology was applicable for bacterial agents, we
154	sequenced CSF from two patients with pneumococcal meningitis, diagnosed by cultivation
155	and/or 16S rRNA gene Sanger sequencing (Sample 6-7). DNA from Streptococcus
156	pneumoniae (S. pneumoniae) was classified with a range between 30,704-60,661 ppm
157	(Sample 6), and 679-804 ppm (Sample 7). In addition to the bacterial samples, we included
158	two CSF samples from patients with RNA viral enterovirus CNS infection (Sample 8-9). As
159	expected, no DNA reads were identified. Enterovirus was, however, found using
160	metagenomic RNA sequencing (Additional figure 2)
161	Samples with co-infections, where EBV was detected along with a primary
162	infectious agent (Enterovirus sample 9, VZV sample 10-11 and Cryptococcus sp. sample 12),
163	were analyzed. Neither the EBV nor the enterovirus was detected in sample 9. VZV and EBV
164	was detected in sample 10, and only VZV was detected in sample 11. Neither yeast nor EBV
165	DNA was detected in sample 12. The results where expected when the following equation
166	was applied for calculating the sensitivity for each agent.

167 The theoretically expected number of pathogen reads was calculated according 168 to pathogen genome size (G_P), the diploid human genome size of 6.5 billion basepairs (G_H), 169 pathogen copy according to PCR per milliliter (C_P), whole blood cell count per milliliter (C_H), 170 and adjusted according to the volume (V), sequencing library size (L) and mappability in 171 percent (M) to remove major contaminants.

172

$$Pathogen\,read = L / \left(\frac{C_H \times G_H \times V}{C_P \times G_P} \times M^{-1} \right)$$

Thus, the detection limit of a single read of a pathogen with a 1 million basepair genome in CSF with normal WBC count (5×10^3 per milliliter) using an input volume of 0.3 milliliter and >95% mappability require a sequencing library of approximately 10 million reads.

177 We included additional nine CSF samples with low levels of EBV DNA (50-178 2000 Geq/ml) (Sample 13-21). With the exception of sample 13 (patient diagnosed with CNS 179 Hodgkin's lymphoma type Post-Transplant Lymphoproliferative Disorder), and sample 16, 180 where EBV was considered the cause of the symptoms, the EBV findings were clinically 181 interpreted as benign incidental findings i.e. not the causative agent for the symptoms of 182 infection. The EBV DNA detected in the majority of samples is likely to originate from 183 latently infected B-lymphocytes recruited into the CSF. Despite the limitations for absolute 184 quantification using qPCR and the stochasticity of distribution of low level pathogen particles, 185 with one exception the calculated reads correlated with the detected reads in the sequencing 186 data (Table 1). In ten samples, more than 1 viral reads was expected and pathogen sequences 187 were found in all samples (Additional Figure 3). In seven samples where less than 1 read was 188 expected to be found, EBV reads were only detected in one dataset (sample 17). Sixteen 189 copies of EBV per milliliter was detected in sample 17 using qPCR and 11 reads were 190 detected using metagenomic sequencing even though 0.3 reads were expected. The 191 discrepancy between the calculation and and sequencing results is most likely due to the 192 stochastic distribution of the few viral particles in the sample. In sample 20, 0.99 reads were 193 expected to be detected in the dataset and a single EBV-read was identified in two of the four 194 classifiers (Kraken2 and Centrifuge). This read was further confirmed using BLAST. The 195 WBC count in sample 18 was below the reference interval of the leukocyte cell counter and 196 was therefore omitted.

197 All pathogen reads from PaRCA were mapped against the corresponding 198 genome sequences using CLC genomics workbench (Figure 1a-e, Additional Figure 4). A dispersed distribution of the reads to the corresponding genomes was observed for all samples, except sample 10, where 5 of the 7 VZV reads (1 overlapping read) originate from a repetitive region within the genome and is therefore expected to be detected at a higher rate, and the last 2 reads map to a downstream gene (no overlap) (Additional Figure 4d). Each sequencing library was subjected to BLAST using the respective reference pathogen genome. The variation of the absolute number pathogen reads comparing the different classifiers detected was lower than 25% (Table 1).

Qualitative and quantitative detection of a known pathogen can thus reproducibly be carried out using the different types of bioinformatic classifiers. Furthermore, an estimation of sensitivity for pathogens can be generated for each sample which can guide the clinician whether the sequencing depth is sufficient to find a certain type of pathogen (Additional Table 4). Notably however, each classifier produced diverse quantities of false positive hits.

Samp	leVerified Pathogen	Clinical Method	qPCR (Geq/ml)	PaRC A (reads)	Kraken2 (reads)	Centrifug e (reads)	Cosmosl D (reads)	BLAST (reads)	Calculate d reads	Range (ppm)	Leukocy (x10 ⁶ /l)
1	HSV1	qPCR	1.0x10 ⁴	97	105	107	107	108	90	6.2-7.2	41
2	VZV	qPCR	3.9x10⁵	213	219	223	211	213	365	14.9-16.0	₂₇₂ 214
3	VZV	qPCR	1.9x10⁵	2,196	2,234	2,251	2,170	2,197	3,072	134.8-147.1	17
4	JCV	qPCR	1.9x10⁵	23,766	24,018	24,190	22,318	23,847	N/A	1,757-2,096	_{N/A} 215
5	JCV	qPCR	4.3x10 ³	496	512	515	484	498	N/A	39.8-57.1	N/A
6	S.	Cultivation/16S rRN/	N/A	766,74	699,662	575,646	701,304	643,083	N/A	30,704-60,611	₅₅ 216
7	S. pneumoniae	16S rRNA qPCR	N/A 3.7x10 ²	12,988 -	11,762 -	12,511 -	12,277 -	12,274 -	N/A 0.1	679-804 Undet.	1064 217
8	Enterovirus	qPCR	6.6x10 ⁴	-	-	-	-	-	N/A	Undet.	95
9	Enterovirus EBV	qPCR qPCR	5.8x10 ⁴ 4.1x10 ²	-	-	-	-	-	N/A 0.1	Undet. Undet.	⁸¹⁴ 218
10	EBV VZV	qPCR qPCR	1.9x10 ³ 4.7x10 ³	10 7	9 7	9 7	8 7	9 7	2.5 4.5	0.8-1.1 0.7-0.8	¹⁸¹ 219
11	EBV VZV	qPCR qPCR	5.0 x10 ¹ 2.9x10 ³	- 15	- 15	- 15	- 12	- 15	0.1 5.5	Undet. 1.2-1.7	90 220
12	EBV Yeast sp.	qPCR Cultivation/Filmarray	9.1 x10 ² N/A	-	-	-	-	-	0.2 N/A	Undet. Undet.	164
13	EBV	qPCR	1.9x10 ³	81	85	82	79	82	20.5	6.7-7.5	₂₆ 221
14	EBV	qPCR	3.7x10 ²	-	-	-	-	-	0.6	Undet.	253
15	EBV	qPCR	3.2x10 ²	6	6	6	6	6	2.5	0.4-0.5	44 222
16	EBV	qPCR	2.7x10 ²	232	228	225	213	223	18.5	21.2-22.8	4
17	EBV	qPCR	1.6x10 ²	11	10	11	11	11	0.3	1.0-1.2	148223
18	EBV	qPCR	1.6x10 ²	-	-	-	-	-	N/A	Undet.	<4
19	EBV	qPCR	8.1 x10 ¹	-	-	-	-	-	0.6	Undet.	31 224
20	EBV	qPCR	5.0 x10 ¹	-	1	1	-	1	0.99	0-0.1	14
21	EBV	qPCR	5.0 x10 ¹	8	8	8	8	9	1.5	0.7-0.8	9 225

211 **Table 1.** Metagenomic sequencing pipeline results.

226 Reads from each classifier from verified pathogen. Calculated reads in accordance with the presented algorithm. N/A: leukocyte count

227 missing for sample 4 and 5, leukocyte count for sample 18 is below reference value, calculation is not applicable for bacteria, fungi and

228 RNA virus. 16S rRNA: 16S rRNA gene Sanger sequencing, HSV1: Herpes simplex virus 1, VZV: Varicella Zoster virus, JCV: JC polyomavirus,

229 EBV: Epstein-Barr virus

230 False positive pathogens

The diversity of viral species detected in metagenomic sequencing libraries were relatively low and recurrent. PaRCA, Kraken2, Centrifuge and CosmosID identified 2-31, 5-13, 17-96 and 0-4 viral species in each sample respectively (Figure 2a, Additional Table 5). Many of the most abundant viral species identified were found in multiple samples (Figure 3). Two samples (4 and 13) contained human virus which were not detected in multiple samples and not a previously confirmed pathogen (see below).

237 The non-pathogen/EBV viral reads were either of human origin, misclassified or 238 contaminations. Human endogenous retrovirus K was identified in all samples, except for the 239 water control, which was expected as the reads originates from the human genome (Figure 3 240 bottom, Additional Table 5). Another ubiquitously detected virus was the BeAN 58058 virus, 241 which was detected in all samples, except for the water control. An additional BLAST 242 examination identified these hits as human reads. Low levels of phage sequences known to 243 infect bacteria from the *Enterobacterales* order were detected in a few samples and in the 244 water control, most likely derived from bacteria purified enzymes used in the various steps of 245 library preparation. A conspicuous pseudomonas phage contaminant in sample 4, 5 and the 246 water control are likely derived from a bacterial contaminant at one of the sequencing sites. 247 Streptococcus phage species were detected in sample 6, from a patient with S. pneumoniae 248 meningitis. Importantly, the most prominent viral species identified in patient samples were 249 also present in the cell controls at similar levels and displayed a similar sequence identity and 250 could therefore be discarded as a pathogen.

Compared with the relatively few viral agents detected by the classifiers, bacterial species were abundant; 61-712 bacterial species were identified using PaRCA, 370-1408 in Kraken2, 845-2826 in Centrifuge and 0-14 in CosmosID (Figure 2b). Two samples originated from patients with a known *S. pneumoniae* meningitis (sample 6 and 7) and

255 bacteria were detected at 69,088 ppm and 803 ppm resepectively (PaRCA). With the omission 256 of the positive samples 6 and 7, trace levels (3.4-18.2 ppm) of S. pneumoniae was 257 ubiquitously detected in all samples. A known environmental contamination of Pseudomonas 258 was detected in the majority of the samples. In two samples (4 and 5) Pseudomonas 259 constituted 389,480 ppm (39%) and 590,195 ppm (59%) of the entire sequencing library 260 respectively, while the prevalence in other samples were lower 6.6-75,279 ppm (0.0007-261 7.5%). A large fraction of the detected bacteria are still left when using previously suggested 262 fixed cut-off at 100 ppm (0.01%) (Figure 2) and unlike the virus species the 263 contaminants/misclassifications cannot be entirely removed using the control samples. 264 However, when further applying an additional filter of comparison of the detected bacterial 265 species between the three classifiers (PaRCA, Kraken2 and Centrifuge) only the known 266 pathogen (S. pneumoniae) or environmental contaminants (Pseudomonas and Escherichia 267 coli) was left. Similarly no eukaryotic species were found in all three classifiers.

268 Considering the ubiquitous presence of viral misclassifications and 269 contaminants in samples as well as controls, a viral pathogen is easily identifiable, but require 270 additional analyses including read distribution and BLAST analysis, for verification in a 271 clinical setting (discussed below). In contrast, the large number of bacterial species identified 272 pose a bioinformatic challenge as the bacterial sequence can be derived from kit 273 contaminants, lab environment or bioinformatic misclassifications which obscure the 274 pathogen reads. As with the virus hits, removal of bacterial contaminants using cell controls 275 can efficiently remove the majority of species, but additional filters are required (Figure 4).

276

277 **Controls**

Two types of controls, water and cell control, were tested for their ability to mirror the bioinformatic missclassifications and contaminations observed in samples. In the water

control the dataset consisted of 99.6% bacterial sequences and 0.06% viral sequences (Additional Table 5). The cell controls originating from EBV-transformed cancer cells had a composition more similar to the samples with 99.2-99.4% human sequences. The number of viral and bacterial strains detected in the water control was 12 and 568 respectively. In contrast the cell controls contain sequences ranging from 3-4 viral and 61-177 bacterial strains.

The viral strains in the water control were mainly of phage origin. In contrast the viral strains detected in the cell controls were similar to the CSF samples, mainly Human endogenous retrovirus K and BeAN 58058 virus. Both cell lines originate from EBVtransformed cancer cells and harbours EBV DNA. The ppm-values of each cell line between sequencing runs was reproducible and no significant difference was found between the classifiers (Additional Figure 5, Additional Table 6).

In the water control, 98% of the sequencing library consisted of reads from *Pseudomonas* and the second most abundant bacterial strain found was *Escherichia coli* (0.1%), which is to be expected as most enzymes are produced in this bacterial system. In contrast, none of the bacterial strains in the cell controls constituted more than 0.1% of the sequencing library.

Thus, the water control efficiently amplified the environmental and kit contaminants, but in contrast to the cell control did not find human misclassifications. Also, since the water control consist entirely of contaminants, the absolute or proportionate content did not allow for a direct comparison with the patient samples. The cell control allowed for direct quantitive and qualitative subtraction of the majority of contaminants and putative pathogens were identified.

303 Unexpected virus findings

In sample 2 and 3 we identified 29-34 EBV reads in both samples in all classifiers (Additional Table 5). The reads were dispersed throughout the genome and displayed minor sequence variability with the reference genome in accordance with previous EBV findings (Additional Figure 6a-b). Due to the limited sample volume we were unable to verify and quantify this finding using qPCR.

309 In sample 4 we identified three viruses which were unexpected, mastadenovirus, 310 papillomavirus and torque teno virus (Additional Figure 6c-e). PaRCA identified 32 reads 311 matching human mastadenovirus C (HAC), Kraken2 32 reads, Centrifuge 30 reads and 312 CosmosID did not report any HAC sequences. The majority of reads, 25 out of 32 where 198 313 bp long, 5 reads where shorter and 2 were longer. BLAST-analysis showed that all reads 314 shared the same 3'-end. Four reads had mismatches in comparison with reference sequence. 315 Considering the size and distribution of the reads our findings are most likely a laboratory 316 amplicon contamination. Human papillomavirus (HPV) reads were detected in PaRCA (12 317 reads), Kraken2 (2 reads), but not by Centrifuge and CosmosID. Ten of the 12 reads were 105 318 bp long and the remaining two, 104 bp and 106 bp respectively. All reads aligned to the 3'-319 end of the virus genome in the L1 gene. Examination of BLAST results showed a high 320 similarity with HPV98 with a one or two base-pair mismatch. As above, considering the size 321 and distribution of the reads our findings were most likely a laboratory amplicon 322 contamination. CosmosID has an inbuilt function to filter out hits that are considered to be 323 amplicons, therefore the software did not report these reads. Different strains of 324 Anellovirus/Torque teno virus (TTV) were detected in the classifiers. PaRCA identified 75 325 reads, Kraken2 25 reads, Centrifuge 55 reads, while CosmosID did not detect any TTV reads. 326 Five distinct consensus reads/contigs were formed from the 75 reads identified in PaRCA. 327 Thirty-one reads formed a consensus reads of 196bp. BLAST analysis of this read displayed a 328 97% identity with TTV14, but only for 91bp of the fragment. The remaining parts of the 329 contig did not show any alignment with any viral species. The origin of this read is therefore 330 unknown. BLAST analysis of the remaining 4 reads/contigs showed alignment (>95% query 331 cover and identity) to an Anellovirus isolate previously identified in metagenomics. The 332 alignment showed an unusual coverage of the 5'-end of the genome and all the reads were 333 aligned to the first half of the genome. The reason for this unusual coverage is unknown, but 334 considering that TTV is widely detected in metagenomic sequencing and the multiple reads 335 aligning to a clinical isolate it is probable that these four contigs/reads originate from the 336 patient sample.

337 In sample 13, we detected 10 reads corresponding to hepatitis C virus (HCV) in 338 PaRCA. Kraken2, Centrifuge and CosmosID detected 5, 6 and 6 reads respectively. The 10 339 reads were concentrated to the 5'-end of the genome, but spread within the initial half of the 340 genome (Additional Figure 6f). An analysis of the BLAST results showed alignment with 341 HCV genotype 1. Synonymous mutations were found in multiple reads as well as gaps. Two 342 reads had a fusion between sequences from different regions of the HCV genome. The 343 sequence diversity indicates that the virus is from a patient, but the frameshift and fusion 344 reads indicates that they are of an artificial origin. Also, the patient had undergone HCV 345 serology analysis which was negative. Finally, considering that HCV is a RNA virus this 346 finding is most likely a laboratory amplicon contamination.

347

348 Discussion

In this study we subjected 21 CSF samples from patients with suspected or confirmed CNS infection to metagenomic DNA sequencing. Pathogen detection accuracy and efficiency was evaluated using five bioinformatic tools. Using 12 samples with minute levels of EBV we concluded that the sensitivity of detection was mainly affected by leukocyte content in the

353 samples and to lesser degree environmental contamination. Bioinformatic classifiers were 354 essentially equally efficient in terms of sensitivity, but produced vastly different number of 355 false positive hits, which inhibited efficient clinical pathogen identification. The removal of 356 these false positive hits originating from contaminants and bioinformatics classifications were 357 alleviated by using a EBV-containing cell control which served as a positive as well as a 358 negative control. A number of criteria have been suggested for how to identify a causative 359 agent in clinical samples e.g. by calculating the fraction of pathogen reads and/or an absolute 360 number of reads. However, using these methods the majority of samples used in this study 361 would be considered negative and/or contain a large number of agents which would be 362 considered falsely positive dependent on the choice of classifier. The lower detection limit 363 could be generalized and compared between studies/laboratories if the leukocyte count was 364 provided. In a similar manner, a general quantification of viral content using ppm is an 365 efficient reference point for comparison between studies [21, 22]. Furthermore, it is evident 366 that local contaminants greatly impact the sequencing library constitution. Therefore, it is 367 necessary that findings in negative controls from each study is presented in its entirety. Nine 368 CSF-samples were identified at the clinic to only contain EBV, and we did not identify any 369 additional pathogen, confirming the results from the clinic. Importantly, using our algorithm a 370 lower detection limit could be determined for pathogens. An alternative to metagenomic 371 sequencing is removal of the dominating host background using various methods including 372 centrifugation and nuclease treatment [23, 24]. However, this will deplete the majority of 373 nucleic acid and only minute amounts of nucleic acid will be left, which complicates the 374 library preparation. Sensitivity would also be reduced, especially for intracellular virus, and 375 bacteria which might precipitate if centrifugation is used. Likewise the specificity would be 376 impaired by the overwhelming number of environmental contaminants as seen in our water 377 control.

378 Our bioinformatic classifier PaRCA, which uses a combination of single reads 379 alignment and assemblies was able to detect more reads from HAC, HPV, TTV and HCV, but 380 failed to detect the single EBV read in sample 20. Bioinformatic classifiers for clinical 381 practice should not only quantify the pathogen reads, but also provide information of read 382 distribution, sequence diversity and subtraction of environmental contaminants and 383 bioinformatic misclassifications, facilitating pathogen detection as shown in this study. Novel 384 pathogens will also require classifiers to detect diverse sequences, as well as enable 385 investigation of sequences which might not classify completely to a genus. Our finding of a 386 novel TTV strain shows that there is a large difference between bioinformatics classifiers 387 ability to identify divergent sequences.

388 In this study we have used archived material, which impair a proper RNA 389 analysis due to degradation. Future studies using fresh CSF-samples where RNA integrity and 390 quantity is measured may provide similar guidelines for RNA pathogen detection. We only 391 included two verified bacterial CSF-samples in this study, one which was detected by 392 culturing and 16S rRNA gene sequencing, and the second one detected by 16S rRNA gene 393 sequencing. A limit of using metagenomic sequencing of CSF from bacterial meningitis 394 patients is the high levels of leukocytes, but this may be compensated by the higher amount of 395 bacterial nucleic acid compared with viral genomes. Here, we applied a fraction cutoff for 396 bacterial findings (>0.01%) in order to decrease the amount of false positive bacterial species 397 findings. This cutoff value should not be considered fixed and future studies with larger 398 bacterial cohort would provide additional guidelines for bacterial species identification.

399

400 Conclusions

We suggest that prior to clinical metagenomic DNA sequencing, an estimation of sequencingdepth is made by adjusting it to the leukocyte content in the sample. Also, a pathogen-

403 containing cell control sequenced at the same depth should be included in the same 404 sequencing run in order to generate the same type of reproducible background. Bioinformatic 405 processing should include a comparison between the pathogens detected in the cell control 406 and the sample as well as between multiple classifiers. Further candidate pathogens reads 407 should be confirmed by using BLAST and mapped against a reference genome to identify 408 read distribution and sequence diversity. A comprehensive evaluation including a theoretical 409 estimation on sensitivity of the metagenomics test as well as other clinical microbiological 410 assays e.g. serology and PCR should assist the clinician in interpreting the final results.

411

412 Methods

413 Sample collection

414 Included in this retrospective study were cerebrospinal fluid samples from patients with CNS 415 symptoms of infection, in which the Department of Clinical Microbiology at Sahlgrenska 416 University Hospital or the The Public Health Agency of Sweden previously had verified the 417 infectious agent during 2015-2017. The sample cohort was chosen to include a variety of 418 microorganisms (DNA/RNA virus, bacteria or fungi) with varying concentration of the 419 pathogens as determined by confirmatory testing using qPCR, cultures, 16S rRNA gene 420 Sanger sequencing or FilmArray (Additional Methods). The samples were stored at -20°C 421 after clinical testing. The cell lines P3HR1 (HTB-62, American Type Culture Collection, 422 ATCC, USA) and Namalwa (CRL-1432, American Type Culture Collection, ATCC, USA), 423 were used as combined negative controls as well as positive controls, due to its inherent EBV 424 genome. The controls were processed in parallel with the patient samples during all the 425 laboratory steps.

427 Sample processing

428 For samples processed at the Department of Clinical Microbiology at Sahlgrenska University 429 Hospital, total nucleic acid was extracted from 400 µl of cerebrospinal fluid using the MagNA 430 Pure Compact Nucleic Acid Isolation Kit I (Roche Diagnostics, Indianapolis, IN, USA) on the 431 MagNA Pure compact automated extractor. For samples processed at The Public Health 432 Agency of Sweden, total nucleic acid was extracted from 200 µl of cerebrospinal fluid sample 433 using the MagDEA® Dx SV (Precision System Science Co Ltd, Matsudo-city, Chiba, Japan) 434 on the magLEAD® 12gC automated extractor (Precision System Science Co Ltd). DNA 435 concentrations were determined using the Qubit Fluorometer (Thermo Fisher Scientific, 436 Waltham, MA, USA) using the dsDNA HS Assay Kit (Thermo Fisher).

437

438 Library preparation and sequencing

439 DNA libraries were prepared according to the modified protocol for metagenomic samples, 440 developed at the Public Health Agency of Sweden, using the Ion Xpress Plus Fragment 441 Library Kit (Thermo Fisher) on the AB Library Builder System (Thermo Fisher). Samples 442 were fragmented to 200 bp, followed by ligation of Ion P1 Adapter as well as Ion Xpress 443 Barcode adapters. The protocol was adjusted to suit low-input samples (<50 ng DNA) by 444 using a reduced volume of P1 adapter and barcodes (0.5 μ l). The libraries were amplified, 445 selecting the number of amplification cycles according to the sample input concentration, 446 varying between 14 to 20 cycles. Amplified libraries were size selected choosing an optimal 447 size range for each individual sample to ensure removal of small-sized PCR concatemers, 448 varying between 100 to 320 bp (including adapters). Size selection was performed using the 449 Pippin Prep platform (Sage Science, Beverly, MA, USA) with 2 % Dye free Agarose Gel 450 Cassette. Following visualization and an estimation of the concentration using the High 451 Sensitivity D1000 DNA Kit on the Agilent 2200 TapeStation system (Agilent Technologies,

452 CA, USA), the samples were pooled according to concentration. Subsequently, libraries were 453 purified using Agencourt AMPure XP (Beckman Coulter, Brea, CA, USA). Finally, libraries 454 were quantified using qPCR with the Ion Library TaqMan Quantitation Kit (Thermo Fisher) 455 and the size estimated using High Sensitivity D1000 DNA Kit on Agilent 2200 TapeStation 456 system (Agilent Technologies). For template preparation, libraries were pooled to a final 457 concentration of 50 pM, if obtainable. For libraries with lower concentration than 50 pM, 458 libraries were pooled to the available concentration. Thereafter, the Ion Chef Platform was 459 used to ligate the libraries onto spheres using the Ion 540 Kit-Chef (Thermo Fisher). 460 Following clonal amplification, libraries were loaded onto Ion 540 Chip and sequencing was 461 performed on the S5 System (XL, Prime; Thermo Fisher) according to the manufacturer's 462 protocol for 200 bp read length.

463

464 **Bioinformatic analysis**

465 **Quality Control**

BAM-files were converted into fastq files using the Torrent Suite Software provided for Ion S5 system. Reads were processed with FASTX toolkit [25] to fasta files. Fastqc was used to identify low-quality reads. Sequences were then subjected to the individual pipelines described below.

470

471 Pathogen detection for Research and Clinical Applications (PaRCA)

Databases were created using built-in tools in Kraken2 and Kaiju. Briefly, databases, corresponding to bacteria, viruses and eukaryotes were created at DNA, RNA and protein level resulting in nine total k-mer databases. The viral databases were comprised of all viral data in GenBank, the bacterial database consisted of the full Progenomes data [26] and eukaryotic databases were composed of the GenBank data for vertebrates, parasites and fungi. After download, the Progenomes database was continuously updated using scripts to reflect
changes within the NCBI taxonomy. Reads were initially trimmed at both directions using
BBDuk (BBMap 37.50) using an entropy mask of 0.9, trim quality of 16 and a minimum
length of 40. Reads were corrected using Fiona (0.2.9) with id=3 for substitution errors.

Reads were classified using Kraken2 and Kaiju by using individual databases. Kraken2 results were filtered using the kraken-filter with a threshold of 0.15 for eukaryotes and 0.05 for viruses and bacteria (a higher threshold indicates higher stringency). Thresholds for Kaiju: score and minimum matches were set to 85.20 for eukaryotes, 80.18 for bacteria and 75.15 for viruses.

486 After initial classification and filtering, Kraken2 results were individually compared and reads with hits in multiple databases were evaluated based on k-mer score with the highest scoring 487 488 match being retained for further downstream analysis. Kaiju scores were internally compared 489 and the hit with the longest protein alignment was preserved. Reads with both Kraken2 and 490 Kaiju hits were then compared and the lowest common ancestor of the two results was 491 selected using mergeOutputs with "-c lowest" from the Kaiju package. Reads where the 492 lowest common ancestor was a species designation were directly counted and saved while 493 reads with a higher lowest common ancestor were further processed in the pipeline. Reads 494 only classified by a single k-mer classifier were labelled as "singletons" and further 495 processed.

Reads were ordered by taxonomic ID, which then were regressed through the taxonomic tree until either a genus-level or kingdom-level was reached. Reads without genuslevel information or reads with a classification above genus level were stored separately for further analysis. After ordering into genus, all taxonomic IDs corresponding to a member of the genus were automatically downloaded from NCBI and corresponding accession identifiers were parsed from the NCBI accession dump file. Accession identifiers were then used to

create a slice of the BLAST nt-database for that specific genus. Reads classified as belonging
to the order "primates" was not processed further and received the taxonomic ID 9606 (Homo
Sapiens).

505 Reads were analyzed in BLAST within the genus using a threshold of an e-value of 10^{-3} and the ten best hits were then retained. The ten results per read were parsed and the 506 507 bit-score per taxon in the hits were aggregated. The taxon with the highest aggregate bit score 508 was then selected as the putative taxon ID for the read. After taxon identification, results were 509 merged and regressed in order to identify the species level classification of the putative taxon. 510 If the kingdom level was reached before a species identification was found, the original taxon 511 identifier was used in its place. Finally, any reads that were not successfully classified within 512 a genus in the BLAST database creation step were collected and subjected to BLAST against the full NT-database with an e-value of $>10^{-5}$ and a minimum query coverage of 20% as 513 514 threshold, again the ten best hits were preserved. The results from both BLAST analyses were 515 aggregated based on bit score and the resulting taxon ID regressed to species level if possible.

516 Classified reads were collected and presented using a krona-graph and tables in 517 an html format. Tables were reorderable on name, taxonomic id and read count. Tables were 518 also filterable, including wildcard functionality. FASTQ-files containing reads classified to an 519 individual species and aggregates corresponding to kingdoms and unclassified reads were 520 directly downloadable.

521

522 Kraken2

523 We used Kraken2 with a dustmasker included in the package.

524

525 Centrifuge

We subjected our samples to Centrifuge with the inbuilt quality control and repeatmasker based on dustmasking from NCBI tools. Briefly, the dustmasker converts the low-quality regions into N's so the aligner skips aligning these sequences [16]. In order to obtain reads from all pathogens included in this study, the total of both leaf and genus levels were incorporated from the Centrifuge reports, thus leading to higher amounts of total classified reads, however, since not all species were converted into the ETE3 toolkit, and some stops at genus level, this does not affect final results of classified pathogens.

533

534 CosmosID

535 Unassembled sequencing reads were directly analyzed using the commercially available 536 genomic platform CosmosID to achieve identification of microbes at species level [27]. Each 537 uploaded sample was searched and cleared from host sequences by the platform prior to 538 analysis. CosmosID automatically filters out phages and amplicon-originated sequences.

539

540 BLAST

BLAST analysis was performed with reference genomes for the pathogens. The cutoff was set to \geq 95% sequence identity and an e-value of \leq 10⁻³. Following standard steps for preprocessing reads, a BLAST search was performed with reads set as subjects and reference genomes set as queries. Reference genomes used were NC_001806 (HSV1), NC_001348 (VZV), NC_00196 (JCV), NC_003098 (*S. pneumoniae*), NC_007605 (EBV), NC_001405 (Human Mastadenovirus C; MAVC), FM_955837.2 (Human Papillomavirus 98; HPV98), MH_649255.1 (Anellovirus), and NC_004102.1 (HCV).

548

549 Calculations and statistical analysis

550	CLC genomics	workbench (Ve	er. 11,	Qiagen)	was used to	perform and	d plot	coverage a	analysis.
-----	--------------	---------------	---------	---------	-------------	-------------	--------	------------	-----------

- 551 Classified sequences from Kraken2 and Centrifuge were visualized using Pavian [28]. Ratio
- between sample ppm and control ppm were calculated, where an ratio ≤ 10 were considered a
- 553 contamination.
- 554 GrapPad Prism Ver. 7.0c was utilized to perform statistical analysis. Kruskal-Wallis test with
- 555 Dunn's multiple comparison tests was applied to compare reproducibility through pipelines.
- 556 A *p*-value ≤ 0.05 were considered significant.
- 557

559 **References**

- 560 1. Granerod J, Ambrose HE, Davies NW, et al. Causes of encephalitis and differences in their
- 561 clinical presentations in England: a multicentre, population-based prospective study. Lancet
- 562 Infect Dis. 2010;10(12):835-44.
- 563 2. Hong HL, Lee EM, Sung H, et al. Clinical features, outcomes, and cerebrospinal fluid
- findings in adult patients with central nervous system (CNS) infections caused by varicella-
- zoster virus: comparison with enterovirus CNS infections. J Med Virol. 2014;86(12):2049-54.
- 566 3. Okike IO, Ribeiro S, Ramsay ME, et al. Trends in bacterial, mycobacterial, and fungal
- 567 meningitis in England and Wales 2004-11: an observational study. Lancet Infect Dis.
 568 2014;14(4):301-7.
- 4. Maeda E, Akahane M, Kiryu S, et al. Spectrum of Epstein-Barr virus-related diseases: a
 pictorial review. Jpn J Radiol. 2009;27(1):4-19.
- 5. Martelius T, Lappalainen M, Palomaki M, Anttila VJ. Clinical characteristics of patients
 with Epstein Barr virus in cerebrospinal fluid. BMC Infect Dis. 2011;11:281.
- 573 6. Siddiqi OK, Ghebremichael M, Dang X, et al. Molecular diagnosis of central nervous
 574 system opportunistic infections in HIV-infected Zambian adults. Clin Infect Dis.
 575 2014;58(12):1771-7.
- 576 7. Salzberg SL, Breitwieser FP, Kumar A, et al. Next-generation sequencing in
 577 neuropathologic diagnosis of infections of the nervous system. Neurol Neuroimmunol
 578 Neuroinflamm. 2016;3(4):e251.
- 579 8. Chiu CY, Miller SA. Clinical metagenomics. Nat Rev Genet. 2019;20(6):341-55.
- 580 9. Palacios G, Druce J, Du L, et al. A new arenavirus in a cluster of fatal transplant-associated
- 581 diseases. N Engl J Med. 2008;358(10):991-8.
- 10. Wilson MR, Naccache SN, Samayoa E, et al. Actionable diagnosis of neuroleptospirosis
- by next-generation sequencing. N Engl J Med. 2014;370(25):2408-17.

- 584 11. Naccache SN, Peggs KS, Mattes FM, et al. Diagnosis of neuroinvasive astrovirus
- 585 infection in an immunocompromised adult with encephalitis by unbiased next-generation
- sequencing. Clin Infect Dis. 2015;60(6):919-23.
- 587 12. Wilson MR, Sample HA, Zorn KC, et al. Clinical Metagenomic Sequencing for Diagnosis
- of Meningitis and Encephalitis. N Engl J Med. 2019;380(24):2327-40.
- 589 13. Gu W, Miller S, Chiu CY. Clinical Metagenomic Next-Generation Sequencing for
 590 Pathogen Detection. Annu Rev Pathol. 2019;14:319-38.
- 591 14. Strong MJ, Xu G, Morici L, et al. Microbial contamination in next generation sequencing:
- 592 implications for sequence-based analysis of clinical samples. PLoS Pathog.593 2014;10(11):e1004437.
- 594 15. Zinter MS, Mayday MY, Ryckman KK, Jelliffe-Pawlowski LL, DeRisi JL. Towards
- 595 precision quantification of contamination in metagenomic sequencing experiments.596 Microbiome. 2019;7(1):62.
- 597 16. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive
 598 classification of metagenomic sequences. Genome Res. 2016;26(12):1721-9.
- 599 17. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using
- exact alignments. Genome Biol. 2014;15(3):R46.
- 18. Kostic AD, Ojesina AI, Pedamallu CS, et al. PathSeq: software to identify or discover
- microbes by deep sequencing of human tissue. Nat Biotechnol. 2011;29(5):393-6.
- 603 19. Breitwieser FP, Lu J, Salzberg SL. A review of methods and databases for metagenomic
- classification and assembly. Brief Bioinform. 2019;20(4):1125-36.
- 20. Nooij S, Schmitz D, Vennema H, Kroneman A, Koopmans MPG. Overview of Virus
- 606 Metagenomic Classification Methods and Their Biological Applications. Front Microbiol.
- 607 2018;9:749.

- 608 21. Tang KW, Alaei-Mahabadi B, Samuelsson T, Lindh M, Larsson E. The landscape of viral
- expression and host gene fusion and adaptation in human cancer. Nat Commun. 2013;4:2513.
- 610 22. Tang KW, Larsson E. Tumour virology in the era of high-throughput genomics. Philos
- 611 Trans R Soc Lond B Biol Sci. 2017;372(1732).
- 612 23. Song Y, Giske CG, Gille-Johnson P, et al. Nuclease-assisted suppression of human DNA
- background in sepsis. PLoS One. 2014;9(7):e103610.
- 614 24. Lewandowska DW, Zagordi O, Geissberger FD, et al. Optimization and validation of
- sample preparation for metagenomic sequencing of viruses in clinical samples. Microbiome.
- 616 2017;5(1):94.
- 617 25. FASTX-Toolkit. <u>http://hannonlab.cshl.edu/fastx_toolkit/</u> Accessed 18 May 2020.
- 618 26. Progenomes2. <u>http://progenomes.embl.de/</u>. Accessed 18 May 2020.
- 619 27. CosmosID. <u>http://www.cosmosid.com/</u> Accessed 18 May 2020.
- 620 28. Breitwieser FP, Salzberg SL. Pavian: Interactive analysis of metagenomics data for
- 621 microbiome studies and pathogen identification. Bioinformatics. 2019.
- 622

623 **Declarations**

- 624 Ethics approval and consent to participate
- 625 The study design and methods were approved by the Regional ethical review board in
- 626 Gothenburg (191-18).

627 Availability of data and materials

628 Will be available on the European Genome-phenome Archive upon publication.

629 **Competing interests**

630 Authors declare no competing interests.

631 Funding

- This project was supported by funding from the Sahlgrenska University Hospital Fund C4A,
- 633 FoU Laboratoriemedicin, the Konrad and Helfrid Johanssons Foundation, the Olle Engkvist
- 634 foundation, the Längmanska foundation and the Wilhelm and Martina Lundgren foundation.
- 635 Author contributions
- 636 This study was designed by HEJ, MLK, SB and KWT. Cells were cultured by YT. Samples
- 637 were selected by KWT, SB and DV. Metagenomic sequencing was performed by MLK, SB,
- and JO. PS, SA and SHM executed bioinformatic analysis. Calculations was done by JO, SA
- and KWT. DV and KWT provided clinical expertise. Manuscript was written by JO, SB, HEJ
- and KWT. Figures and tables was prepared by JO, YT and KWT. All authors read and
- 641 approved the final manuscript.
- 642

643 Acknowledgements

- We thank the Bioinformatics Core Facility at the Sahlgrenska Academy for bioinformaticsanalyses.
- 646

bioRxiv preprint doi: https://doi.org/10.1101/2020.06.25.170423; this version posted June 25, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

648 Additional information

- 649 Additional table 1. Clinical data (docx)
- 650 Additional table 2. Dataset species classification (docx)
- 651 Additional table 3. Pathogen detection by bioinformatic classifier (docx)
- **Additional table 4.** Patient report with estimation of sensitivity for pathogens (xlsx)
- 653 Additional table 5. Species identified in bioinformatic classifiers (xlsx)
- 654 Additional table 6. Cell control reproducibility (docx)
- 655 Additional method. Description of clinical methodology used for comparison (docx)

656

- 657 Additional figure 1. Overview of the sample and bioinformatic processing (pdf)
- 658 DNA from cerebrospinal fluid specimens was extracted and followed by library construction
- and sequencing. Datasets generated by the Ion S5 were processed by four different
- bioinformatics classifiers to profile the microbiome. BLAST was used for verification.

661

662 Additional figure 2. Enterovirus samples (pdf)

Results of viral species detected in RNA sequencing datasets of sample 8 and sample 9 inPaRCA.

- 665
- 666 Additional figure 3. Correlation between detected and calculated reads (pdf)
- 667 The reads detected by PaRCA correlated to the calculated reads using the algorithm,

668 Spearman Correlation coefficient, n=10

669

670 Additional figure 4. Coverage density plot of microbial species in CSF samples (pdf)

671	Reads from samples not shown in main figure mapped to reference genomes of (a) VZV
672	(NC_001348), (b) JCV (NC_00196), (c) S. pneumoniae (NC_003098), (d, f) VZV
673	(NC_001348), and (e, g-j) EBV (NC_007605) using CLC Genomics Workbench. Number of
674	reads (y-axis) at each nucleotide position of the genome (x-axis) depicted in blue. Dark blue
675	represents peak, bright blue average and light blue minimum coverage for respective section
676	of the genome.

677

678 Additional figure 5. Cell control coverage density plot and reproducibility (pdf)

679 Coverage analysis of EBV reads detected in cell controls Namalwa (a) and P3HR1 (b) 680 mapped to EBV reference genome (NC_007605) using CLC Genomics Workbench. Number 681 of reads (y-axis) at each nucleotide position of the genome (x-axis) depicted in blue. Dark 682 blue represents peak, bright blue average and light blue minimum coverage for respective 683 section of the genome. EBV reads shown as parts per million reads (ppm) in each of the cell 684 line controls for each of the bioinformatic classifier (c), n=4 (Namalwa) or n=5 (P3HR1); 685 Kruskal-Wallis test with Dunn's multiple comparisons show no significant difference 686 between the pipelines.

687

688 Additional figure 6. Coverage analysis for unexpected findings (pdf)

Reads from samples with ambigous findings mapped to reference genomes of EBV
NC_007605 (a-b), Human Mastadenovirus C (MAVC) NC_001405 (c), Human
Papillomavirus 98 (HPV98) FM_955837.2 (d), Anellovirus MH_649255.1 (e), and HCV
NC_004102.1 (f), using CLC Genomics Workbench.

693 Figure Captions

694 Figure 1. Pathogen genome alignment

695 Coverage density plot of sequencing reads from respective sample and control detected in

696 PaRCA aligned to reference genomes of HSV1 (a), VZV (b), JCV (c), S. pneumoniae (d) and

697 EBV (e-f). Number of reads (y-axis) at each nucleotide position of the genome (x-axis)

698 depicted in blue. Dark blue represents peak, bright blue average and light blue minimum

699 coverage for respective section of the genome.

700

701 Figure 2. Detected pathogens in bioinformatic classifiers

Number of viral (a) and bacterial species (b) classified in each of the samples and controls

vs using the different bioinformatic classifiers. Dark blue bars shows number of total number of

species classified, bright blue bars shows amount of bacterial species over the fraction cutoff

 $(\geq 0.01\%$ of the dataset), light blue bars shows number of species not removed using controls.

706

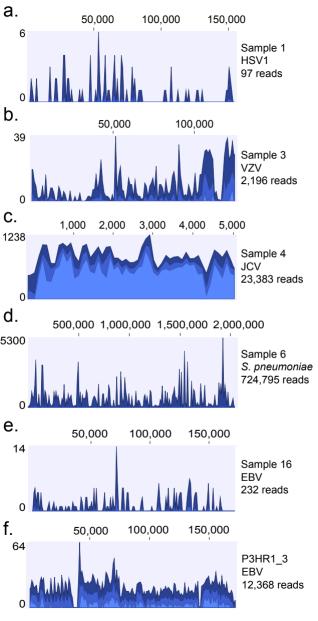
707 Figure 3. Viral species identified in datasets

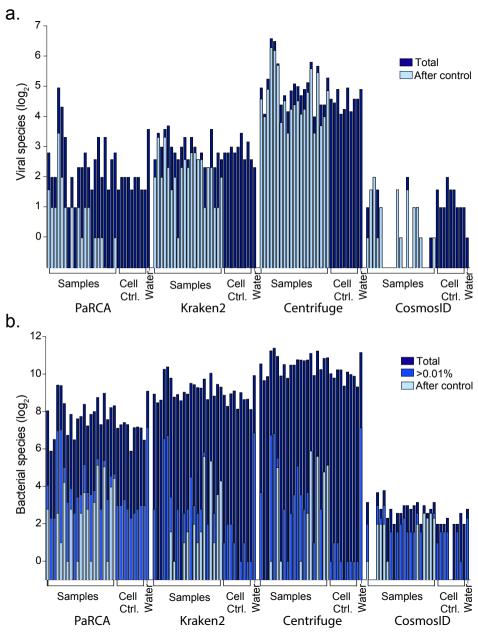
Heatmap showing the ten most abundant viral species in each sample detected using PaRCA.
AcMNPV: Autographa californica multiple nucleopolyhedrovirus. Controls: P; P3HR1, N;
Namalwa, W; water.

711

712 Figure 4. Discerning microbial pathogens from contaminations and misclassifications

Flowchart for identification of pathogens by removing false positive species. Virus contaminants can be removed by comparison of datasets with controls and manual examination of remaining viral reads. Phages can be disregarded as these virus do not infect human cells. Bacterial species require additional filters including a cutoff value and comparison between classifiers.





Virus		Samples														Controls															
		1	2	3	4	5	6	7	8 9	9 10) 11	1 12	2 13	14	15	16	17	18	19	20	21	P1	P2	P3	P4	P5	N1	N2	N3	N4	W
Pathogens	Human alphaherpesvirus 1																														
	Human alphaherpesvirus 3																														
	Human gammaherpesvirus 4																														
	JC polyomavirus 2																														
	Hepatitis C virus																														
	Erwinia phage 1																														
	Escherichia phage 1																														
	Escherichia phage 2																														
	Lactococcus phage 1																														
	Pseudomonas phage 1																														
	Pseudomonas phage 2																														
	Pseudomonas phage 3																														
Phages	Pseudomonas phage 4																														
i nages	Pseudomonas phage 5																														
	Pseudomonas phage 6																														
	Pseudomonas phage 7																														
	Pseudomonas phage 8																														
	Salmonella phage1																														
	Streptococcus phage 1																														
	Streptococcus phage 2																														
	Streptococcus phage 3																														
	AcNMPV																														
Other	BeAn 58058 virus																														
	Human endogenous retrovirus K																														

10⁻⁶ 0.05 % of total classified reads

