

# 1 Clinical and biological insights from viral genome sequencing

2  
3 Authors: Charlotte J. Houldcroft<sup>1§</sup>, Mathew A. Beale<sup>2</sup> & Judith Breuer<sup>2,3\*</sup>

4  
5 \*corresponding author

## 6 7 Affiliations:

- 8 1. Infection, Immunity and Inflammation, Great Ormond Street Institute of Child Health,  
9 University College London, UK
- 10 2. Division of Infection and Immunity, University College London, UK
- 11 3. Great Ormond Street Hospital for Children NHS Foundation Trust, London, UK.

12 § Present address: Division of Biological Anthropology, University of Cambridge, UK

## 13 14 Introduction

15 Since the publication of the first shotgun sequenced genome (cauliflower mosaic virus<sup>1</sup>), the  
16 draft human genome<sup>2</sup> and the first bacterial genomes (*Haemophilus influenzae*<sup>3</sup> and  
17 *Mycoplasma genitalium*<sup>3</sup>), combined with the rapidly falling cost of high-throughput  
18 sequencing<sup>4</sup>, genomics has become a major contributor to our understanding of human and  
19 pathogen biology. Multiple large scale systematic pathogen genome projects have been  
20 recently completed or are on-going (e.g. sequencing thousands of microbiomes and fungal  
21 genomes<sup>5,6</sup>); these projects are shaping our knowledge of the genetic variation present in human  
22 and pathogen populations, the nature of genetic changes that underlie disease, and the sheer  
23 diversity of microorganisms with which we share our environments.

24 The methods and data from whole genome sequencing are increasingly being applied to clinical  
25 medicine, both from a human<sup>7</sup> and pathogen perspective. For example whole-pathogen genome  
26 sequencing has been used to identify new routes of *Mycobacterium abscessus*<sup>8</sup> nosocomial  
27 transmission and to understand *Neisseria meningitidis* epidemics in Africa<sup>9</sup>, while partial  
28 genome sequencing has been used to detect drug resistance in RNA viruses such as influenza<sup>10</sup>  
29 and DNA viruses such as human cytomegalovirus (HCMV)<sup>11</sup>. Viral genome sequencing has

30 gained considerable traction, often focused on research or epidemiology. Whole pathogen  
31 genome sequencing has the advantage of detecting all known drug resistance mutations in a  
32 single test while deep sequencing can identify low level drug resistance mutations early enough  
33 for clinical intervention<sup>12,13</sup>. Whole genomes also provide good data with which to identify  
34 linked infections for public health and infection control purposes<sup>14,15</sup>. Notwithstanding,  
35 progress in whole-genome sequencing (WGS) of viruses for clinical practice has been slow. In  
36 contrast whole-genome sequencing of bacteria is now well accepted particularly for outbreak  
37 tracking and for the management of nosocomial transmission of antimicrobial resistant  
38 bacteria<sup>16,17</sup>.

39 This review will address the challenges and opportunities for making WGS, using modern next  
40 generation sequencing (NGS) methods, a standard part of clinical virology practice. We will  
41 discuss the strengths, weaknesses and technical challenges inherent to different viral WGS  
42 laboratory methods (Table 1). The importance of deeply sequencing certain viral pathogens  
43 will be addressed. We will also explore two areas in which viral WGS has recently proven its  
44 clinical utility: metagenomic sequencing to identify viruses causing encephalitis (box 1); and  
45 role of WGS in molecular epidemiology and public health management of the pan-American  
46 Zika virus outbreak (box 2). Finally, we will briefly consider the ethical and data analysis  
47 challenges which clinical viral WGS presents.

48

#### 49 **Why sequence viruses in clinical practice?**

50 For small viruses such as HIV, influenza, HBV and HCV, partial genome sequencing has been  
51 widely used for research purposes, but also has important clinical applications. For example,  
52 the management of highly active anti-retroviral therapy (HAART) for HIV relies heavily on  
53 viral sequencing for detection of mutations conferring drug resistance. HAART has  
54 dramatically improved survival of patients with HIV, but successful therapy requires long-term

55 suppression of viral replication with anti-retroviral drugs, which may be prevented by impaired  
56 host immunity, sub-optimal drug penetration to host tissue compartments and incomplete  
57 patient adherence to therapy<sup>18</sup>. Where viral replication continues to occur, the high mutation  
58 rate of HIV enables resistance variants to emerge. It has become standard practice in many  
59 parts of the world to sequence the HIV *pol* gene, which encodes the main viral enzymes, for  
60 mutations conferring resistance to inhibitors of reverse transcriptase, integrase and protease<sup>19</sup>,  
61 particularly when patients are first diagnosed and when viral loads indicate treatment failure.  
62 Sequencing resistance mutations has allowed more targeted alterations in treatment with  
63 significantly greater reductions in virus loads compared with standard care (undetectable HIV  
64 load in 32% vs 14% of patients after six months)<sup>20,21</sup>. Thus sequencing resistance mutations to  
65 guide HIV treatment improves disease outcomes. Similar approaches have been taken for  
66 identifying HCV<sup>22</sup>, HBV<sup>23</sup>, and influenza<sup>24</sup> resistance mutations.

67

### 68 **Why sequence whole genomes?**

69 Limited sequencing of the small number of genes that are targeted by anti-viral agents, such as  
70 the HIV polymerase gene, has hitherto been the norm in clinical practice. For detecting a  
71 limited number of antiviral resistance mutations, WGS has been too costly and labour-intensive  
72 to justify. However, the increase in numbers of antivirals targeting genes that are located across  
73 the genome, coupled with falling costs of sequencing and the use of sequence data for  
74 transmission studies, are driving a reappraisal of the need for WGS. For example, antiviral  
75 treatment for HCV now targets four gene products (NS3, NS4A, NS5A, NS5B) encoded by  
76 more than 50% of the viral genome<sup>25</sup>. Separate targeted sequencing for each of these can be as  
77 expensive and time consuming as WGS<sup>26</sup>. Partial genome sequencing is particularly  
78 problematic for larger viral genomes most notably those of the herpesviruses HCMV<sup>11</sup>, VZV<sup>27</sup>,  
79 HSV-1<sup>28</sup> and -2<sup>29</sup>. These have traditionally been treated with drugs targeting the

80 protein/thymidine kinase and DNA polymerase genes. However the growing numbers of drugs  
81 in development that interact with different proteins encoded by viral genes scattered across the  
82 genome, means that the targeted sequencing of multiple genes required for resistance testing is  
83 costly and less tractable<sup>30</sup>. Sequencing the whole genome captures all resistance mutations  
84 simultaneously and obviates the need to design and optimise new PCR assays for detecting  
85 resistance to new drugs. A good example of this is HCMV, where WGS can simultaneously  
86 capture the genes with products targeted by the licensed therapies such as UL27 (unknown  
87 function), UL54 (DNA polymerase), and UL97 (protein kinase), as well as newer drugs such  
88 as letermovir which targets UL56 (terminase complex), enabling comprehensive anti-viral  
89 resistance testing in a single test<sup>11</sup>. At the same time WGS has the potential to provide  
90 information on epitopes, evolution of sequences within a patient over time<sup>11</sup>, and evidence of  
91 recombination between HCMV strains<sup>31</sup>. WGS can highlight putative novel drug resistance  
92 mutations, or predicted changes to epitopes, although phenotypic testing of any findings in a  
93 model system is required to confirm clinical resistance (e.g. <sup>32</sup>) or to map epitope changes (e.g.  
94 <sup>33</sup>).

95 As pre-existing resistance to anti-viral drugs (for example, protease inhibitor-resistant HCV<sup>34</sup>  
96 and nucleoside analog reverse transcriptase inhibitor-resistant HBV<sup>35</sup>) increases, whole  
97 genome sequences will provide the comprehensive resistance data required for selecting  
98 appropriate treatment to achieve good patient outcomes. A complete knowledge of all  
99 resistance mutations can also support more radical management decisions. In a recent case  
100 report, identification of extensive genome-wide HCMV drug resistance within a patient  
101 supported the clinical decision to change to immunotherapeutic treatments, specifically  
102 autologous cytomegalovirus-specific T cells<sup>36</sup>.

103 Whole genomes may also better identify transmission events and outbreaks, which is not  
104 always possible with sub-genomic fragments. For example, sequencing respiratory syncytial

105 virus (RSV) genomes demonstrated that variation was present outside the gene traditionally  
106 used for genotyping, and could be used to help track outbreaks within households, where there  
107 had been insufficient time for single genes to accumulate enough genetic variability to be used  
108 for transmission studies<sup>37</sup>. The increased number of phylogenetically informative variant sites  
109 obtained from generating full or near full length genomes has been shown to obviate the need  
110 for high quality sequences, allowing robust linking of Ebola cases and public health  
111 interventions in real time during the recent epidemic<sup>38</sup>. This also applies to Zika virus, and Box  
112 2 explores the role of WGS in public health efforts to control the outbreak in South America.  
113 The increased use of whole pathogen sequencing routinely for diagnostic purposes<sup>39</sup> is likely  
114 to have wider clinical and research benefits. For example HIV genome sequencing to identify  
115 resistance mutations, can also be used to explore questions related to viral evolution<sup>40</sup>, public  
116 health<sup>41</sup> and viral genetic association with disease. This includes well-powered genotype-  
117 phenotype association studies or genome-to-genome association studies, which look for  
118 associations between viral genetic variants, host genetic variants, and outcomes of infection,  
119 such as viral load set point in HIV infection<sup>42</sup>.

120

### 121 **Why do we need deep sequencing?**

122 Modern methods which make use of massively parallel sequencing provide better opportunities  
123 to examine pathogen diversity through analysis of viral populations within or between hosts  
124 that contain nucleotide variants or haplotypes at low (sub-consensus, less than 50%)  
125 frequencies. Minority variant analysis is particularly powerful for RNA or retro-transcribing  
126 viruses, because they typically have high within-host nucleotide diversity. HIV is the classic  
127 example; the viral replication cycle utilises an error-prone reverse transcriptase enzyme that  
128 introduces mutations at an extremely high rate ( $4.1 \pm 1.7 \times 10^{-3}$  per base per cell)<sup>43</sup>. This results  
129 in a given patient containing not one, but many closely related viruses each bearing subtly

130 different variants, sometimes described as a quasispecies or cloud of intra-host viral diversity.  
131 The presence of a mixed population of viruses introduces problems for determining the true  
132 consensus ‘majority’ sequence, but these minority (non-consensus) variants may also alter the  
133 clinical phenotype of the virus, or predict changes in genotype, tropism or drug resistance. For  
134 example, a minor variant conferring drug resistance in HIV present at only 2.1% of sequencing  
135 reads in a baseline patient sample can rapidly rise to become a majority (consensus) variant  
136 under the selective pressure of drug treatment<sup>44</sup>. Investigators have observed similar changes  
137 in frequency of resistance-associated alleles during treatment of viruses such as HBV<sup>45</sup>, HCV<sup>46</sup>,  
138 HCMV<sup>11</sup> and influenza<sup>47</sup>.

139 Sensitive deep-sequencing of viruses is not only required to detect drug resistance: for HIV, it  
140 is also key in genotypic prediction of receptor tropism, which has clinical implications in  
141 treatment of HIV. HIV can be grouped genotypically by its cellular co-receptor usage as R5  
142 (CCR5-using), X4 (CXCR4-using) or R5X4 (dual tropism). Maraviroc is a CCR5 receptor  
143 antagonist, blocking infection by R5-tropic HIV genotypes, but contraindicated in HIV+  
144 individuals who have X4 or R5X4 HIV genotypes. Just a 2% frequency of X4 or R5X4  
145 genotypes is predictive of maraviroc treatment failure<sup>48</sup>. Sub-consensus frequencies of X4 or  
146 R5X4 HIV are also important to the success<sup>49</sup> or failure<sup>50</sup> of bone marrow transplants from  
147 CCR5-deleted (CCR5- $\Delta$ 32) donors. This may influence decisions to continue or stop anti-viral  
148 therapy in these patients<sup>49</sup>.

149 Detection of minority variants and haplotype identification may also detect mixed infections.  
150 In HCMV mixed-genotype infections or super-infections<sup>51</sup> are associated with poor clinical  
151 outcomes - detection of these by WGS might support a decision to treat disease in these patients  
152 more aggressively.

153 Establishing the clinical associations of minority variants is clearly important, as with  
154 maraviroc treatment failure; Sanger sequencing of a virus population can detect minority

155 variants down to frequencies of between 10 and 40% (e.g. <sup>52</sup>), whilst NGS can sequence those  
156 same PCR amplicons to a much greater depth<sup>53</sup>, and consequently capture more of the  
157 variability present. Thresholds of sensitivity and specificity established need to be specific to  
158 the virus in question, and reflect the potential biases of the sequencing methods used. Many  
159 studies of HIV drug resistance utilising deep-sequencing of PCR amplicons require minority  
160 variants to be present at >1%, to reduce the possibility of false positives<sup>54,55</sup>. This may lead to  
161 a failure to detect true drug resistance mutations at frequencies of 0.1%-1%, which may  
162 ultimately be associated with poor treatment outcome on some drug regimes<sup>55</sup>. While a 1-2%  
163 frequency threshold (or lower) may be clinically relevant to drug resistance in HIV, it is less  
164 clear whether the same degree of sensitivity would be required for monitoring vaccine escape  
165 in HBV or drug resistance in herpesviruses (discussed below). Large cohorts of patients will  
166 need to be followed with samples collected before, during and after treatment<sup>44,48</sup>, to establish  
167 clinical significance thresholds for minority drug resistance<sup>11</sup> and vaccine escape variants for  
168 each virus.

169 Direct deep sequencing of clinical material, either by shotgun or RNAseq methods (so called  
170 metagenomic methods) also provides the opportunity for unbiased detection of pathogen  
171 sequences and thus primary diagnosis of viral and other infections, thereby providing an  
172 alternative to culture, electron microscopy and qPCR. This is discussed further below.

173

#### 174 **Practical considerations for sequencing virus genomes**

175 As previously alluded to, sequencing viral nucleic acid whether cultured or directly from  
176 clinical specimens, is complicated by the presence of contaminating host DNA<sup>56</sup>. This makes  
177 it different from bacterial sequencing which is easily carried out using clinical isolates and thus  
178 sample preparation is relatively straightforward (Table 2). Currently, genome sequencing of  
179 viruses can be achieved by ultradeep sequencing or by enriching for viral nucleic acid prior to

180 sequencing either directly or through prior concentration of viral particles. All approaches  
181 have their own costs and complexities.

182

183 The three primary methods currently used for viral genome sequencing are summarised in  
184 Figure 1.

185

186 (i) **Metagenomics - ultra deep sequencing**

187

188 Metagenomic approaches have been extensively used for pathogen discovery and for  
189 characterising microbial and general pathogen diversity in environmental and clinical  
190 samples<sup>57,58</sup>. Total DNA and/or RNA from a sample, including from host, bacteria, viruses,  
191 fungi and other pathogens present are extracted, put through library preparation and sequenced  
192 by ‘shotgun’ or RNA-seq methods (see Box 1). These approaches have proven to be very  
193 powerful for detecting viral<sup>59,60,61</sup> and other causes<sup>62</sup> of encephalitis where other conventional  
194 methods such as PCR have failed. Box 1 explores the growing diagnostic applications for  
195 metagenomics and RNAseq, for example in encephalitis of unknown aetiology (e.g.<sup>63-65</sup>). In  
196 addition, a number of whole viral genomes have been sequenced in this manner, including  
197 Epstein-Barr virus (EBV)<sup>66</sup> and HCV<sup>26</sup>. However, these methods may be insensitive, because  
198 of the presence of contaminating host and commensal pathogen nucleic acid<sup>56</sup> (Table 2) in  
199 clinical specimens. For example on-target read yields (the proportion of reads matching the  
200 target genome) from metagenomic WGS of 0.008% (EBV genome from the blood of a healthy  
201 adult<sup>67</sup>), 0.0003% (lassa virus genomes from clinical samples<sup>68</sup>) and 0.3% (a filtration and  
202 centrifugation enriched Zika virus sample<sup>69</sup>) have been reported. The read depths obtained are  
203 often inadequate for robust resistance calling<sup>26</sup> and the cost is high. Thus the method has  
204 typically only been performed on a small number of samples for research purposes (e.g.<sup>69,70</sup>).  
205 To improve read depths, concentration of viral particles prior to sequencing (as for example in



206 the Zika case<sup>69</sup>), depletion of host material or ultra-deep sequencing have been employed, all  
207 of which add to the cost. Concentrating viral particles from clinical specimens by antibody-  
208 mediated pulldown (e.g. VIDISCA), filtration, or ultracentrifugation, to isolate a fragment size  
209 profile, and depletion of free nucleic acid<sup>71-74</sup> have all been tried. These host nucleic acid  
210 depletion methods may result in there being insufficient viral nucleic acid for sequencing  
211 library preparations. To overcome this, non-specific amplification methods (e.g. multiple  
212 displacement amplification; MDA) which make use of random primers and phi 29 polymerases  
213 may be effective in increasing DNA load. However, these approaches are time consuming,  
214 costly, and may increase the risk of biases, error and contamination without necessarily  
215 improving the sensitivity of sequencing<sup>75,76</sup>. Moreover, there are often still a high proportion  
216 of host reads present in treated samples<sup>77</sup>.

217 Where metagenomic methods are used for pathogen discovery or diagnosis, appropriate  
218 bioinformatic tools and databases capable of evaluating whether detected pathogens sequences  
219 are truly likely to be the cause of infection, innocent bystanders or contaminants are critical.  
220 Bioinformatic analyses of large metagenomic datasets places an increased burden on high  
221 performance computational resources.

222 The fact that metagenomics requires no prior knowledge of the viral genome, can be considered  
223 a strength<sup>26</sup> in that it allows novel viruses to be sequenced without the need for primer or probe  
224 design and synthesis. This is particularly apposite for rapid responses to emerging threats such  
225 as Zika<sup>78</sup>. Metagenomic viral genome sequencing may also ‘piggy back’ on projects to  
226 sequence virus-associated cancer genomes, which informs clinical care of the cancer or  
227 provides further information on cancer evolution, while generating high coverage of integrated  
228 virus genomes as part of the process<sup>66</sup>. However, the presence of incidental findings (human  
229 genome sequences with potential disease associations, pathogens which were not part of the  
230 question that prompted the initial sequencing) may also present ethical (and even diagnostic)

231 dilemmas for some applications of clinical metagenomics (discussed below and reviewed in  
232 <sup>79</sup>). A recent case in point was a cluster of acute flaccid myelitis cases associated with  
233 enterovirus D68<sup>80</sup>. The analysis of the metagenomic datasets derived from patients was the  
234 subject of discussion through formal<sup>81</sup> and informal scientific channels  
235 (<http://omicsomics.blogspot.co.uk/2015/07/leaky-clinical-metagenomics-pipelines.html>), with  
236 different groups disagreeing over the interpretation of the same data, especially as some of the  
237 alternative pathogens detected can cause treatable bacterial disease. Regulation and reporting  
238 frameworks will be important to resolve future issues of this kind.

239

240 (ii) **PCR amplicon enrichment,**

241 An alternative to metagenomic approaches is to enrich for the specific viral genome prior to  
242 sequencing. PCR amplification of hundreds to thousands of base pairs of viral genetic material  
243 using primers that are complementary to a known nucleotide sequence has been the most  
244 common approach to enriching for small viral genomes such as HIV and influenza, prior to  
245 NGS sequencing for diagnostic and public health purposes. Recent examples of this approach  
246 being applied for public health include sequencing measles virus by PCR-WGS to provide  
247 maximum phylogenetic resolution of an outbreak at the 2010 Winter Olympics<sup>82</sup>, sequencing  
248 of Ebola virus genome to study epidemic dynamics<sup>38</sup>, and Zika virus genome sequencing  
249 (explored in Box 2). PCR whole-genome sequencing of norovirus (7.5kb genome) has been  
250 used to understand norovirus transmission in community<sup>83</sup> and hospital<sup>84</sup> settings. For example,  
251 this research showed that some cases within a hospital with plausible epidemiological linkage  
252 were in fact independent introductions of the pathogen; but that other cases were the result of  
253 transmission, despite infection control practices being in place<sup>84</sup>. Other PCR-based deep  
254 sequencing studies have generated multiple whole genomes for influenza<sup>85</sup> (~13.5kb), dengue<sup>86</sup>  
255 (~11kb), and HCV<sup>87</sup> (9.6kb). This method is feasible (as with PCR and Sanger sequencing)

256 because these viruses all have relatively small genomes, requiring only a small number of PCR  
257 amplicons to assemble whole genome sequences. RNA virus heterogeneity may however  
258 necessitate the use of multiple overlapping primer sets to ensure comprehensive amplification  
259 of all genotypes, for example HCV<sup>26</sup>, norovirus<sup>83</sup>, rabies<sup>88</sup> and RSV<sup>37</sup>. PCR amplicon  
260 sequencing is also more successful for WGS of samples with low virus loads than metagenomic  
261 methods<sup>26</sup>, although other methods such as target enrichment of viral sequences may work  
262 equally well in low copy number samples (e.g. low copy norovirus samples<sup>89</sup>).

263 Overlapping PCRs combined with NGS have been used to sequence the whole genomes of  
264 larger viruses such as HCMV<sup>90</sup>, but this overlapping amplicon method has limited scalability,  
265 since many primers are needed<sup>90</sup> and a greater amount of starting DNA to allow for each  
266 additional PCR, which may not be available from clinical samples. This limits the number of  
267 suitable samples available and also the genomes which can be studied with this method. A  
268 molecular epidemiology study of the relatively small Ebola genome required between 8 and 19  
269 PCR products to amplify the genome for MinION nanopore sequencing<sup>38</sup>, whilst 14<sup>83</sup> and 22  
270 pairs<sup>84</sup> of primers were needed to amplify and Illumina sequence norovirus genomes. This  
271 becomes less practical in a clinical rather than research setting because of the high laboratory  
272 workload associated with large numbers of discrete PCR reactions, the necessity for  
273 individually normalising concentrations of different PCR amplicons prior to pooling, the  
274 increasing probability of reaction failure due to primer mismatch, particularly in very variable  
275 genomes and the increasing labour and consumables cost associated with multiple PCR  
276 reactions<sup>91</sup>. Therefore, although PCR-based sequencing of viruses as large as 250 Kb is  
277 technically possible, the proportional relationship between genome size and technical  
278 complexity make PCR sequencing of sequencing viral genomes beyond 20 - 50 Kb impractical  
279 with current technologies, particularly with regards to large multi-sample studies or routine  
280 diagnostics. Another consideration is that increasing PCR reactions require a corresponding

281 increase in available sample, and this is not always possible where clinical specimens are  
282 limited. Improvements in microfluidic technologies may help to overcome some of these  
283 barriers to PCR-based methods, for example Fluidigm, RainDance and other ‘droplet’  
284 sequencing technologies. Microfluidics-based PCR and pooling of multiple amplicons have  
285 been used successfully to sequence multiple anti-microbial resistance loci, for example<sup>92</sup>, and  
286 can also applied to viral genomes, potentially down to the single-cell sequencing level.  
287 PCR may encounter problems in amplifying highly variable pathogens such as HCV<sup>93</sup> and  
288 norovirus where there are many different genotypes, with some genotypes encountering primer  
289 amplification issues<sup>26,89</sup>, or where there is insufficient characterisation of intra-genotypic  
290 diversity, leading to primer mis-matches<sup>83</sup>. Careful design of degenerate primers may help to  
291 mitigate these problems, but novel variants still present a risk to detection and amplification.

292

293 **(iii) Target Enrichment methods**

294 Target enrichment (TE) methodologies (also known as pulldown, capture or specific  
295 enrichment methods) represent one solution to problems of PCR or metagenomic sequencing  
296 of virus genomes. We and a number of other groups have been developing methods that can be  
297 used to sequence whole viral genomes directly from clinical samples without the need for prior  
298 culture or PCR<sup>94-96</sup>. These methods typically involve small RNA/DNA probes designed to be  
299 complementary to the pathogen reference sequence (or panel of references). Unlike in specific  
300 PCR amplicon based methods, the entire genome can be covered by a single tube of  
301 overlapping probes which are used in a hybridisation reaction to capture or ‘pull down’  
302 complementary DNA sequences bound to a solid phase (e.g. streptavidin-labelled magnetic  
303 beads) from the total nucleic acids present in a sample, followed by sequencer-specific (e.g.  
304 Illumina) adaptor ligation and a small number of PCR cycles to enrich for successfully ligated  
305 fragments. This has been used successfully to characterise large and small clinically relevant  
306 viruses such as HCV<sup>26</sup>, HSV1<sup>97</sup>, VZV<sup>96</sup>, EBV<sup>98</sup>, CMV<sup>66</sup>, HHV6<sup>99</sup> and HHV7<sup>100</sup>. The reaction

307 is performed in a single well and, like microfluidics-based PCR reactions, is amenable to high  
308 throughput automation<sup>98</sup>. The lack of a culture step means that the sequences obtained are more  
309 representative of original virus than cultured viral isolates, with fewer mutations than observed  
310 in PCR amplified templates<sup>66,96</sup>. The success of this method is in part based on the number of  
311 available reference sequences for the virus of interest: specificity increases when baits are  
312 designed against a larger panel of reference sequences, leading to better capture of the breadth  
313 of within and between sample diversity. TE probe design allows for limited mismatching  
314 between template and probe, but whilst PCR requires only knowledge of flanking regions of a  
315 target region, TE requires knowledge of the internal sequence in order to design baits. This is  
316 balanced by the fact that TE is less vulnerable to a single amplicon failure due to mismatch as  
317 internal and overlapping regions may still be captured even if one probe fails<sup>66,96</sup>. As such TE  
318 is not suitable for characterisation of novel viruses with low homology to known viruses, where  
319 metagenomics (or in some cases, PCR using degenerate primers), may be more appropriate.  
320 As with all methods, the technique is also subject to constraints with regard to starting viral  
321 load. We have shown that although capable of sequencing virus from viral loads as low as 2000  
322 IU/ml (HCV) or 2500 IU/ml (HCMV), targets could only be enriched so much, leading to  
323 reduced depth of coverage in sequencing data at lower viral concentration<sup>26,66</sup>. With  
324 metagenomics, the proportion of sequencing data mapping to the pathogen genome (the on-  
325 target read percentage) that can be expected from unenriched sequencing of clinical samples is  
326 small. Depending upon the starting pathogen load in a sample, TE can enrich percentage on-  
327 target viral reads from 0.01% up to 80% or more<sup>66</sup>. This allows a higher degree of multiplexing  
328 than unenriched metagenomics, and brings an accompanying decrease in the price of  
329 sequencing, albeit with a relative increase in the cost of library preparation. There are  
330 alternative approaches to enriching viral reads which include pulse-field gel electrophoresis<sup>101</sup>,

331 which separates large viral genomes from smaller host DNA fragments, allowing for  
332 sequencing libraries composed of a smaller proportion of contaminating host DNA.

333 Enrichment techniques which make use of degenerate RNA or DNA probes to hundreds of  
334 viral species to pull viral nucleic acid out of samples and sequence them, e.g.the VirCapSeq  
335 method, have also been developed<sup>102</sup>. This method is designed for detection of both known and  
336 novel viruses, although its performance remains to be evaluated.

337

### 338 **Comparison of all three methods**

339 To date, there has been very little direct comparison between the three methods for viral  
340 genome sequencing in clinical practice, with only one paper evaluating relative performance  
341 for HCV sequencing<sup>26</sup>. Results from this study, in which three different enrichment protocols,  
342 two metagenomic methods and one overlapping PCR method were evaluated, showed that  
343 metagenomic methods were the least sensitive, yielding the lowest genome coverage for  
344 comparable sequencing effort and were more prone to yield incomplete genome assemblies.  
345 The PCR method was the least tractable and most labour intensive, requiring repeated  
346 amplification and was the most likely to miss mixed infections, but where reactions were  
347 successful, yielded the most consistent read depth, whereas metagenomics and TE yielded read  
348 depths in proportion to virus copy number. Some HCV genotypes (particularly genotype 2)  
349 were more prone to generate incomplete sequences when PCR was used instead of  
350 metagenomics or TE. Targeted enrichment was the most consistent method, achieving full  
351 genomes and identical consensus sequences. The ease of library preparation for metagenomic  
352 and TE sequencing of HCV was considered a major advantage for clinical sequencing, but PCR  
353 may still be appropriate for very low virus load samples.

354 Similar results were achieved in a study comparing norovirus sequencing from PCR amplicons  
355 and target enrichment<sup>89</sup>. TE generated 100% genome coverage in 164/164 samples, while PCR-

356 based capsid sequencing was only possible in 158/164 samples, with PCR failures attributable  
357 to low virus titres and PCR primer mismatches, suggesting TE is more sensitive than PCR for  
358 norovirus sequencing and better accommodates between-strain sequence heterogeneity<sup>89</sup>. TE  
359 has also been used as a fall-back method for samples with lower virus loads which do not give  
360 WGS after metagenomic sequencing<sup>103</sup>. Both metagenomic and TE methods have the  
361 advantage that they are applicable to all size pathogen genomes, whereas PCR based methods  
362 are less tractable for sequencing larger viral genomes or for non-viral (e.g. bacterial, fungal,  
363 parasite) pathogen genomes.  
364 These direct comparisons of different methods<sup>26,89</sup> will be important in demonstrating the  
365 situations in which each method should be used, based on their sensitivity and specificity, as  
366 well as factors which are relevant to clinical diagnostic labs such as cost, scalability and turn-  
367 around time (summarised in Table 1).

368

### 369 **Challenges of analysis and interpretation**

370 Beyond the technical challenges of method choice for viral WGS, there are a number of other  
371 roadblocks which may slow the advance of WGS in the clinic. They may be considered in three  
372 groups: ethical issues, including incidental host and microbiological findings; regulatory  
373 issues, such as the establishment of standards, good laboratory practice and sensitivity and  
374 specificity thresholds for sequencing; and analytical issues, regarding data interpretation and  
375 the proliferation of analysis options.

376

### 377 **Ethical issues and incidental findings**

378 In many clinical tests (e.g. MRI scans, host genome sequencing), there is a risk of detecting a  
379 disease association that was not part of the original investigation yet may have clinical  
380 significance for the individual or their family. These so called ‘incidental findings’ remain a

381 topic of intense medical ethical debate<sup>104</sup>. The risk of incidental findings in pathogen  
382 sequencing (e.g. discovery of HIV infection during metagenomic sequencing for other  
383 pathogens) is not unique and has been resolved in clinical virology laboratories, where  
384 multiplex PCRs are used and only one of the tests has been requested. In these cases it is the  
385 practice of the laboratory to suppress the result that has not been requested (personal  
386 communication, J Breuer). In UK laboratories, the clinical virologist who interprets the test  
387 results is part of the team managing the patient and as such may decide to discuss an unexpected  
388 result with the physician-in-charge. Incidental *host* genetic findings (e.g. detection of variants  
389 that predispose to cancer risk) from a pathogen metagenomics study are not reported to the  
390 individual in the UK, because this reporting is only permissible with patient consent. In regard  
391 to both host and virus incidental findings, targeted enrichment and PCR have an advantage as  
392 they target only the pathogen of interest. The ethical and privacy concerns associated with the  
393 presence of host genetic data in publically available metagenomic datasets have been well  
394 reviewed by Hall and colleagues<sup>79</sup> and represent a separate challenge.

395

### 396 **Regulatory challenges**

397 Regulation, as well as helping to address some of the concerns addressed above, will also be  
398 important in standardising WGS of viruses. The framework required to make viral WGS  
399 sufficiently robust and reproducible in clinical practice will come from a number of areas.

400 The framework of laboratory accreditation and benchmark testing already available (for  
401 example CLIA in the USA, or accreditation against medical laboratory quality and competence  
402 standardisation criteria for ISO 15189) will support the development of viral WGS standards  
403 if there is sufficient pressure from hospitals, journals and funding agencies.

404 Lessons learned from the use of PCR in diagnostics may be useful here, beginning with  
405 ensuring good clinical laboratory and molecular practices<sup>105,106</sup>. This will mean including



406 negative samples in every sequencing run, to assess contamination thresholds, spiking samples  
407 with a known virus to provide a sensitivity threshold and including positive controls and  
408 controls for batch-to-batch variation<sup>1077</sup>, all of which will increase sequencing costs and are  
409 likely to deter adoption of pathogen genome sequencing by laboratories sequencing small  
410 batches of samples. The result may be to drive centralisation of virus WGS to ensure adequate  
411 standards are kept, ensure large batches of samples and keep costs down.

412 The issues of sensitivity and contamination are especially important in WGS because of the  
413 risk of both false-negative and false-positive detection of pathogens. Highly sensitive  
414 sequencing (whether metagenomic, PCR or TE based) may detect low-level contaminating  
415 viral nucleic acid (reviewed in<sup>108,109</sup>). For example murine leukaemia virus<sup>110,111</sup> and  
416 parvovirus-like sequences<sup>112,113</sup> are just two of many contaminants that have been recognised  
417 to come from common laboratory reagents such as nucleic acid extraction columns<sup>114</sup>. As with  
418 other highly sensitive technologies, robust laboratory practices and protocols are needed to  
419 minimize contamination. It is also important to remember that detection of viral nucleic acid  
420 does not necessarily identify the cause of illness, and it is good practice when using NGS  
421 methods for diagnosis of viral infections to confirm the findings with alternative, independent  
422 methods which do not rely on nucleic acid testing. For example in cases of encephalitis of  
423 unknown origin, positive NGS findings can be confirmed by immunohistochemical analysis of  
424 the affected tissue<sup>59,115</sup>, or identification of the virus by electron microscopy or tissue culture<sup>79</sup>.

425 The standardisation of methods, including bioinformatics approaches will be key to the  
426 successful use of NGS and WGS in clinical virology. Software packages that use a graphical  
427 user interface (GUI) rather than requiring command-line expertise, with strict version control  
428 of software and analysis pipelines to make results reproducible, best practices easily shareable,  
429 and to allow accreditation of analysis software will be necessary, whilst retaining an  
430 appreciation that best-practice analysis methods are continually evolving and prematurely

431 standardising in an overly rigid manner may inhibit innovation. Commercialisation and  
432 regulation may help, providing financial and regulatory incentives to ensure that analysis tools  
433 and technologies meet the needs of clinical sequencing for virology. Finally for drug resistance,  
434 the development of well curated databases of which mutations are truly indicative of drug  
435 resistance will be critical for accurate clinical interpretation. Such databases have already been  
436 created for HIV<sup>116</sup>, HBV<sup>117,118</sup> and HCV<sup>119</sup>, but without recognition of their value by funding  
437 agencies, and corresponding centralised funding to ensure their continued maintenance and  
438 upkeep, tools may become swiftly outdated or unusable.

439

#### 440 **Financial barriers to the use of viral WGS in a clinical setting**

441 While there are good reasons for sequencing whole genomes, and the general use of NGS, if  
442 diagnostic or hospital-based laboratories are to be persuaded to make the transition away from  
443 sequencing sub-genomic fragments, they need to see not only that the additional information  
444 gained from WGS is really of benefit to patient care; but that WGS is (or will become) as  
445 scalable and automatable as sub-genomic fragment sequencing, that the regulatory framework  
446 is suitable and that the price of sequencing whole genomes is competitive with sequencing  
447 fragments.

448 Currently the costs of sequencing viral genomes, notwithstanding their small size, remain  
449 generally higher than sequencing of sub-genomic target resistance genes. Equally, whole  
450 genome information may provide important additional knowledge, as discussed above. The  
451 cost difference between sequencing a target region and the whole virus genome is largely  
452 governed by the size of the genome versus the size and number of target loci.

453

#### 454 **What does the future hold? Long-read sequencing and host depletion**

455 Current generation NGS technologies based around Illumina, 454, Ion Torrent or Sanger  
456 methodologies as described above have the ubiquitous problem of generating short-read data  
457 which presents challenges for haplotype phasing of intra-host minor variants, which aims to  
458 identify whether a set of genetic variants occur on the same genetic background (clonal  
459 population) or on related, highly-similar but different genetic backgrounds within the same  
460 population (sometimes called a viral swarm or cloud); as well as sequencing across repetitive,  
461 recombinatorial or mobile genetic regions which are more difficult to resolve using short reads  
462 due to problems such as mapping ambiguities. The clinical implications of understanding  
463 whether, for example, multi-drug resistance occurs on a clonal genetic background or in a  
464 mixed population of viruses with different drug resistance profiles is currently unclear.

465 While there are computational tools (e.g.<sup>120</sup>) to help resolve these issues, especially of interest  
466 to researchers, there are also new technologies available. Newer single molecule sequencers  
467 such as PacBio (Pacific Biosciences) and MinION (Oxford Nanopore) are capable of extremely  
468 long read sequencing, and in some cases whole viral sequences (for example viruses with  
469 genomes under 20kb, such as Ebola virus, norovirus and influenza A) could theoretically be  
470 obtained from single reads. The MinION also has the advantage of being very fast, taking in  
471 some cases as little as four hours to go from sample receipt to reporting of analysed data<sup>121</sup>.

472 Data on viral read lengths achieved from MinION sequencing have been relatively modest (e.g.  
473 mean read lengths of: 751bp (Modified Vaccinia Ankara), 758bp (cowpox virus)<sup>122</sup>, 455bp  
474 [range 126–1477] (chikungunya virus), 358bp [220–672] (Ebola virus), 1576bp<sup>123</sup> or 6895bp  
475 (HCMV)[personal communication, M Beale] and 572bp [range 318–792] (HCV)<sup>124</sup>). Results  
476 from the better-established PacBio technology are more promising, including a recent report  
477 of a pseudorabies virus genome sequenced with a mean read length of 12,777bp (against a  
478 double-stranded DNA genome ~142kb in length)<sup>125</sup>, and 9.2kb reads have been achieved in

479 PacBio HCV genome sequencing (where only 9.2kb of the 9.6kb had been pre-amplified by  
480 PCR)<sup>126</sup>.

481 A drawback of both NGS and single-molecule sequencing however is the need for high  
482 coverage to minimize the impact of sequencing errors, particularly in the context of drug  
483 resistance studies, as drug resistance most frequently results from single nucleotide mutations  
484 or small deletions (1-3 bases), especially in lower-fidelity RNA viruses<sup>127</sup>. This can be a  
485 challenge where the amount of viral genome is dwarfed by the presence of host DNA, and  
486 when the error profile of a technology makes point mutations particularly hard to detect  
487 accurately<sup>121</sup>. At the time of writing, MinION sequencing (R9 pore chemistry) has raw 2D read  
488 error rates of ~5% [personal communication, Josh Quick], which compares unfavourably with  
489 Illumina (<0.1%), Ion Torrent (~1%) and PacBio (13% single pass, <1% with circular  
490 consensus read) error rates<sup>128</sup>.

491 However, demonstration of the potential for using these long read technologies with target  
492 enrichment provides a potential way forward<sup>123,129</sup>, as ambiguities can be resolved if sufficient  
493 depth of sequence is achieved for the target pathogen, and errors rates for all methodologies  
494 may be reduced with technological and analytical improvements. Products or methods which  
495 can deplete the host genetic background but not the viral nucleic acids within a sample would  
496 be an alternative solution, meaning a higher proportion of virus reads would be recovered from  
497 each sequencing run. While there are already solutions in place for bacterial sequencing (e.g.  
498 human ribosome RNA or mitochondrial depletion, selective depletion of DNA with a certain  
499 methylation pattern), there are no dedicated products for viral sequencing.

500

## 501 **Conclusion**

502 Whole virus genome sequencing is of growing importance in a clinical context, for diagnosis,  
503 disease management and molecular epidemiology (including infection control). There are a

504 number of methods available to achieve WGS of viruses from clinical samples. Currently the  
505 choice of methodology (amplicon sequencing, target enrichment or metagenomics) is specific  
506 to both the virus and the clinical question. Metagenomic sequencing is most appropriate for  
507 diagnostic sequencing of unknown or poorly characterised viruses, PCR works well where viral  
508 genomes are short and diversity in primer binding sites is low, while target enrichment works  
509 for all pathogen sizes, but is particularly advantageous for large viruses and for viruses with  
510 diverse but well characterised genomes. Two obvious areas of innovation currently exist: firstly  
511 for methods that can effectively deplete host DNA whilst preserving viral DNA, and secondly  
512 for further development in the long-read technology market in order to achieve the range of  
513 flexibility and competitive pricing that exists in the short-read market. New technologies are  
514 needed to unite the strengths of these different methods and allow healthcare providers to invest  
515 in a single technology which is suitable for all viral WGS applications.

516

## 517 **Acknowledgements**

518 The authors would like to thank Julianne Brown and Kimberly Gilmour (GOSH) and Ronan  
519 Doyle (UCL) for their helpful discussions, and Josh Quick (University of Birmingham) for  
520 sharing unpublished MinION statistics.

521

## 522 **Funding**

523 CJH was funded by Action Medical Research grant GN2424. MAB was funded through the  
524 European Union's Seventh Programme for research, technological development and  
525 demonstration under grant agreement No 304875 held by JB. This work was supported by the  
526 National Institute for Health Research Biomedical Research Centre at Great Ormond Street  
527 Hospital for Children NHS Foundation Trust and University College London. JB receives  
528 funding from the UCLH/UCL National Institute for Health Research Biomedical Research

529 Centre. The authors acknowledge infrastructure support for the UCL Pathogen Genomics Unit  
530 from the UCL MRC Centre for Molecular Medical Virology and the UCLH/UCL National  
531 Institute for Health Research Biomedical Research Centre. The funders had no role in study  
532 design, data collection and interpretation, or the decision to submit the work for publication.

533

534 **Box 1. RNA-seq and metagenomics diagnostics.**

BOX:RNASEQ AND METAGENOMICS DIAGNOSTICS

In cases of encephalitis of unknown origin, metagenomic techniques are becoming increasingly promising diagnostic tools. There are a variety of protocols in use, but the clearest distinction is between RNA-seq and metagenomics. RNA-seq is the sequencing of either the total RNA or a subset of RNA extracted from a sample (cerebrospinal fluid or brain biopsy, for example), converted to cDNA and sequenced. Metagenomics is generally used to describe the same procedure for DNA, but may also include simultaneous sequencing of DNA and RNA by incorporating a cDNA synthesis step. RNA-seq methodologies may improve detection of pathogenic viruses, as many viruses have RNA genomes; the expression of viral genes in the CSF or brain is indicative of both the presence of the virus, and which viral genes are being transcribed. However, DNA viruses which experience low-level transcription may be poorly detected using RNA-seq and read numbers for DNA viruses may be higher in metagenomic datasets<sup>64</sup>.

Both methods have successfully identified new or known viral pathogens implicated in encephalitis of unknown origin. Metagenomics has been used to aid in diagnosis and characterisation of enterovirus D68 in cases of acute flaccid paralysis<sup>80</sup>. Metagenomics identified herpesviruses in the CSF of four patients with suspected viral meningoencephalitis<sup>130</sup>. RNA-seq also successfully identified HSV1 in an encephalitis case, although the use of a DNase I digestion (intended to lower the amount of host nucleic acid

in the subsequent sequencing library) lowered the number of HSV1 reads<sup>64</sup>. Mumps vaccine virus has also been detected a chronic encephalitis case using RNAseq [Morfopoulou, S. Deep sequencing reveals persistence of cell-associated mumps vaccine virus in chronic encephalitis. *Acta Neuropathologica* (In Press.)].

RNA-seq has been very successful in identifying encephalitis caused by astroviruses<sup>131,132</sup> and coronaviruses<sup>59</sup>. The deaths of three squirrel breeders from encephalitis was linked to a novel squirrel bornavirus through the use of a metagenomic protocol in which DNA and RNA were separately extracted and sequenced as discrete libraries, providing complementary data<sup>63</sup>. Ultimately, metagenomics provides more information about the virus genome present in a sample than PCR alone, which may be important for molecular epidemiology, while RNA-seq has the power to selectively capture information on which sequences are present, as well as informing researchers about viral gene expression of relevance to pathology.

535

536 Box 2: The role of whole-genome sequencing in Zika virus epidemiology and infection control

Box 2: The role of whole-genome sequencing in Zika virus epidemiology and infection control

Zika whole-genome sequencing is being used to understand the epidemiology of the outbreak (where did the virus come from? When did it enter Brazil?); to understand the connection between the virus and microcephaly; and to inform control measures, by stopping importation or interrupting transmission from a reservoir, and informing blood safety measures in hospitals, for example by demonstrating transfusion transmissibility of the virus. Whole-genome (or near whole-genome) sequence is required from flavivirus genomes to give molecular epidemiology studies sufficient power<sup>41</sup>. WGS, phylogenetic analysis and molecular clock dating, combined with other epidemiological data, were useful in excluding

hypotheses about the introduction of Zika virus to South America<sup>41</sup>. For example, the most recent common ancestor of strains circulating in Brazil predates the 2014 football World Cup, making it highly unlikely that this event was responsible for introducing Asian-lineage Zika virus to South America<sup>41</sup>.

WGS is also central to understanding Zika virus pathogenesis, and could be used to interrogate the whole genome of Zika virus for changes associated with microcephaly, as not enough of the virus's biology is currently understood to allow studies to limit themselves to smaller regions of the genome. It's likely that a wide sample of Zika whole genome sequences, from around the world and from microcephaly and asymptomatic cases, will be needed to give confidence to any studies linking particular mutations to the birth defects seen in the recent Zika virus outbreak. No changes in the Zika virus genome have yet been unambiguously associated with microcephaly<sup>41,69,78</sup>.

Whole-genome and fragment sequencing were used to identify a case of probably transfusion transmission of Zika virus through a platelet donation. This has significant public health and infection control relevance as it suggests asymptomatic donors are capable of transmitting the virus to immunocompromised individuals, although PCR-based testing had already established the presence of Zika virus in the blood supply in a previous outbreak, in this case without molecular epidemiology to demonstrate cases of Zika virus in blood product recipients<sup>133</sup>. Blood products may need to be routinely screened for Zika virus<sup>134</sup>.

Finally, whole-genome sequencing of Zika isolates has found sequence polymorphisms within primer-binding sites<sup>135</sup>, which may make PCR-based diagnosis and virus load quantification more difficult. This highlights the need to characterise population-level diversity, especially in epidemics, where the locally circulating virus sequence may have diverged significantly from related sequences from other locations or time periods. A number of projects are underway to achieve these goals, including the ZIBRA mobile laboratory



project<sup>136</sup>, employing portable metagenomic sequencing of Zika virus (<http://zibraproject.github.io/>) and real-time reporting of results<sup>103</sup>.

537

538 **Figure 1: Major methods for sequencing viral genomes from clinical specimens.**

539 All specimens originally comprise a mix of host (in blue) and pathogen (in red) sequences.  
 540 Direct metagenomic sequencing provides an accurate representation of the sequences within  
 541 the sample at the cost of high sequencing and data analysis/storage costs. PCR amplicon  
 542 sequencing uses many discrete PCR reactions to enrich the viral genome, significantly  
 543 increasing the workload for large genomes, but reducing the sequencing costs. Target  
 544 enrichment sequencing uses virus-specific nucleotide probes bound to a solid phase to enrich  
 545 the viral genome in a single reaction, reducing workload, but increasing library cost (relative  
 546 to PCR).

547

548 **Table 1. Advantages and disadvantages of different viral sequencing sample**  
 549 **preparation approaches**

	Advantages	Disadvantages
<b>Metagenomics</b>	<ul style="list-style-type: none"> <li>• Simple, cost-effective sample preparation</li> <li>• Can sequence novel/poorly characterised genomes</li> <li>• Effective in ‘pathogen fishing’ approaches to identify potential underlying pathogen</li> <li>• Low number of PCR cycles limits introduction of amplification mutations</li> <li>• Preservation of minor variant frequencies</li> </ul>	<ul style="list-style-type: none"> <li>• High sequencing cost to obtain sufficient pathogen sequence</li> <li>• Relatively low sensitivity to target pathogen, and coverage proportional to viral load</li> <li>• High proportion of non-pathogen reads increases computational challenges</li> <li>• Incidental sequencing of human and off-target pathogens raises ethical/diagnostic issues</li> </ul>

	<p>reflects <i>in vivo</i> variation</p> <ul style="list-style-type: none"> <li>• No primer/probe design enables rapid response to novel pathogens or sequence variants.</li> </ul>	
<b>PCR</b>	<ul style="list-style-type: none"> <li>• Tried and trusted – large technical resource of well-established methods and trained staff</li> <li>• Highly specific – most sequencing reads will be pathogen, reducing sequencing costs</li> <li>• Highly sensitive, with good coverage achievable even at low pathogen load</li> <li>• Relatively straightforward to introduce new primer designs for novel sequences</li> </ul>	<ul style="list-style-type: none"> <li>• Labour intensive and difficult to scale for large genomes</li> <li>• Iterating standard PCRs across large genomes requires high sample volume</li> <li>• PCR reactions subject to primer mismatch, particularly in poorly characterised or highly diverse pathogens, or those with novel variants</li> <li>• Limited ability to sequence novel pathogens</li> <li>• High number of PCR cycles may introduce amplification mutations</li> <li>• Uneven amplification of different PCR amplicons may influence minor variant and haplotype reconstruction</li> </ul>
<b>Target Enrichment</b>	<ul style="list-style-type: none"> <li>• Single tube sample preparation suited to high throughput automation and sequencing of large genomes</li> <li>• Increased specificity over metagenomics reduces sequencing costs</li> <li>• Overlapping tiling of probes increases tolerance for individual primer mismatches</li> </ul>	<ul style="list-style-type: none"> <li>• High cost and technical expertise for sample preparation</li> <li>• Unable to sequence novel pathogens and requires well characterised reference genomes for probe design</li> <li>• Sensitivity is comparable to PCR but coverage is proportional to pathogen load – low pathogen load yields</li> </ul>

	<ul style="list-style-type: none"> <li>• Reduced number of PCR cycles (relative to PCR) limits introduction of amplification mutations</li> <li>• Preservation of minor variant frequencies reflects <i>in vivo</i> variation</li> </ul>	<p>low/incomplete coverage</p> <ul style="list-style-type: none"> <li>• Cost and time to generate new probe sets limits rapid response to emerging/novel sequences</li> </ul>
--	--	---

550

551 **Table 2. Limitations of viral sequencing**

	<b>Bacteria</b>	<b>Viruses.</b>	<b>Challenges</b>
<b>Genome</b>	dsDNA	dsDNA, ssDNA, partially dsDNA, ssRNA, dsRNA	Different extraction protocols for different viruses, use of cDNA synthesis in RNA viruses or second strand synthesis for ssDNA viruses,
<b>Gene Conservation</b>	Bacteria have highly conserved genes essential for life (e.g. 16s) allowing broad microbiome studies and surveys of taxa	No homologous genes between viruses of different phyla	Lack of conserved homology between viral phyla prevents universal primer based surveys of virome.
<b>Culture</b>	Often straightforward to culture and obtain pure, highly enriched bacterial DNA/RNA	Challenging to culture, and requires a host cell for replication	Cultured virus is heavily contaminated with host cell genome/transcriptome, reducing equivalent viral sequencing output
<b>Clinical specimens</b>	Hardy bacterial cells with cell walls can often be separated from human cells in clinical specimens using differential lysis methods	Viruses are intracellular pathogens, and cannot easily be separated from clinical samples prior to extraction	Clinical specimens are heavily contaminated with host genome/transcriptome, reducing equivalent viral sequencing output

	prior to extraction		
<b>Bacterial methylation patterns</b>	Bacteria are prokaryotes and use different methylation patterns from eukaryotes. Host DNA can be depleted post-extraction using restriction endonucleases directed against CpG methylation	DNA viruses are often methylated by host intracellular machinery, and may possess similar methylation patterns	DNA digestion according to methylation patterns is less effective as a means of host-depletion for viral sequencing post-extraction

552

553

554

555

556 **References**

557

558 1 Gardner, R. C. *et al.* The complete nucleotide sequence of an infectious clone of  
559 cauliflower mosaic virus by M13mp7 shotgun sequencing. *Nucleic Acids Res* **9**, 2871-  
560 2888 (1981).

561 2 Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**,  
562 860-921, doi:10.1038/35057062 (2001).

563 3 Fraser, C. M. *et al.* The minimal gene complement of *Mycoplasma genitalium*.  
564 *Science* **270**, 397-403 (1995).

565 4 Hayden, E. C. Technology: The \$1,000 genome. *Nature* **507**, 294-295,  
566 doi:10.1038/507294a (2014).

567 5 Turnbaugh, P. J. *et al.* The human microbiome project. *Nature* **449**, 804-810,  
568 doi:10.1038/nature06244 (2007).

569 6 Grigoriev, I. V. *et al.* MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic*  
570 *Acids Res* **42**, D699-704, doi:10.1093/nar/gkt1183 (2014).

571 7 Worthey, E. A. *et al.* Making a definitive diagnosis: successful clinical application of  
572 whole exome sequencing in a child with intractable inflammatory bowel disease.  
573 *Genet Med* **13**, 255-262, doi:10.1097/GIM.0b013e3182088158 (2011).

574 8 Bryant, J. M. *et al.* Whole-genome sequencing to identify transmission of  
575 *Mycobacterium abscessus* between patients with cystic fibrosis: a retrospective cohort  
576 study. *Lancet* **381**, 1551-1560, doi:10.1016/S0140-6736(13)60632-7 (2013).

- 577 9 Lamelas, A. *et al.* Emergence of a new epidemic *Neisseria meningitidis* serogroup A  
578 Clone in the African meningitis belt: high-resolution picture of genomic changes that  
579 mediate immune evasion. *MBio* **5**, e01974-01914, doi:10.1128/mBio.01974-14  
580 (2014).
- 581 10 Zaraket, H. *et al.* Genetic makeup of amantadine-resistant and oseltamivir-resistant  
582 human influenza A/H1N1 viruses. *J Clin Microbiol* **48**, 1085-1092,  
583 doi:10.1128/JCM.01532-09 (2010).
- 584 11 Houldcroft, C. J. *et al.* Detection of Low Frequency Multi-Drug Resistance and Novel  
585 Putative Maribavir Resistance in Immunocompromised Pediatric Patients with  
586 Cytomegalovirus. *Frontiers in Microbiology* **7**, doi:10.3389/fmicb.2016.01317  
587 (2016).
- 588 12 Witney, A. A. *et al.* Clinical application of whole-genome sequencing to inform  
589 treatment for multidrug-resistant tuberculosis cases. *J Clin Microbiol* **53**, 1473-1483,  
590 doi:10.1128/JCM.02993-14 (2015).
- 591 13 Simen, B. B. *et al.* Low-abundance drug-resistant viral variants in chronically HIV-  
592 infected, antiretroviral treatment-naive patients significantly impact treatment  
593 outcomes. *J Infect Dis* **199**, 693-701, doi:10.1086/596736 (2009).
- 594 14 Smith, G. J. *et al.* Origins and evolutionary genomics of the 2009 swine-origin H1N1  
595 influenza A epidemic. *Nature* **459**, 1122-1125, doi:10.1038/nature08182 (2009).
- 596 15 Gire, S. K. *et al.* Genomic surveillance elucidates Ebola virus origin and transmission  
597 during the 2014 outbreak. *Science* **345**, 1369-1372, doi:10.1126/science.1259657  
598 (2014).
- 599 16 Koser, C. U. *et al.* Routine use of microbial whole genome sequencing in diagnostic  
600 and public health microbiology. *PLoS Pathog* **8**, e1002824,  
601 doi:10.1371/journal.ppat.1002824 (2012).
- 602 17 Cartwright, E. J., Koser, C. U. & Peacock, S. J. Microbial sequences benefit health  
603 now. *Nature* **471**, 578, doi:10.1038/471578d (2011).
- 604 18 Paredes, R. & Clotet, B. Clinical management of HIV-1 resistance. *Antivir Res* **85**,  
605 245-265, doi:10.1016/j.antiviral.2009.09.015 (2010).
- 606 19 Van Laethem, K., Theys, K. & Vandamme, A. M. HIV-1 genotypic drug resistance  
607 testing: digging deep, reaching wide? *Curr Opin Virol* **14**, 16-23,  
608 doi:10.1016/j.coviro.2015.06.001 (2015).
- 609 20 Durant, J. *et al.* Drug-resistance genotyping in HIV-1 therapy: the VIRADAPT  
610 randomised controlled trial. *Lancet* **353**, 2195-2199 (1999).
- 611 21 Clevenbergh, P. *et al.* Persisting long-term benefit of genotype-guided treatment for  
612 HIV-infected patients failing HAART. The Viradapt Study: week 48 follow-up.  
613 *Antivir Ther* **5**, 65-70 (2000).
- 614 22 Khudyakov, Y. Molecular surveillance of hepatitis C. *Antivir Ther* **17**, 1465-1470,  
615 doi:10.3851/IMP2476 (2012).
- 616 23 Kim, J. H., Park, Y. K., Park, E. S. & Kim, K. H. Molecular diagnosis and treatment  
617 of drug-resistant hepatitis B virus. *World journal of gastroenterology* **20**, 5708-5720,  
618 doi:10.3748/wjg.v20.i19.5708 (2014).
- 619 24 McGinnis, J., Laplante, J., Shudt, M. & George, K. S. Next generation sequencing for  
620 whole genome analysis and surveillance of influenza A viruses. *J Clin Virol* **79**, 44-  
621 50, doi:10.1016/j.jcv.2016.03.005 (2016).
- 622 25 Pawlotsky, J. M. Hepatitis C Virus Resistance to Direct-Acting Antiviral Drugs in  
623 Interferon-Free Regimens. *Gastroenterology* **151**, 70-86,  
624 doi:10.1053/j.gastro.2016.04.003 (2016).

- 625 26 Thomson, E. *et al.* Comparison of next generation sequencing technologies for the  
626 comprehensive assessment of full-length hepatitis C viral genomes. *J Clin Microbiol*,  
627 doi:10.1128/JCM.00330-16 (2016).
- 628 27 Brunnemann, A. K. *et al.* Drug resistance of clinical varicella-zoster virus strains  
629 confirmed by recombinant thymidine kinase expression and by targeted resistance  
630 mutagenesis of a cloned wild-type isolate. *Antimicrob Agents Chemother* **59**, 2726-  
631 2734, doi:10.1128/AAC.05115-14 (2015).
- 632 28 Karamitros, T. *et al.* De Novo Assembly of Human Herpes Virus Type 1 (HHV-1)  
633 Genome, Mining of Non-Canonical Structures and Detection of Novel Drug-  
634 Resistance Mutations Using Short- and Long-Read Next Generation Sequencing  
635 Technologies. *PLoS One* **11**, e0157600, doi:10.1371/journal.pone.0157600 (2016).
- 636 29 Piret, J. & Boivin, G. Antiviral drug resistance in herpesviruses other than  
637 cytomegalovirus. *Rev Med Virol* **24**, 186-218, doi:10.1002/rmv.1787 (2014).
- 638 30 Melendez, D. P. & Razonable, R. R. Letemovir and inhibitors of the terminase  
639 complex: a promising new class of investigational antiviral drugs against human  
640 cytomegalovirus. *Infect Drug Resist* **8**, 269-277, doi:10.2147/IDR.S79131 (2015).
- 641 31 Lassalle, F. *et al.* Islands of linkage in an ocean of pervasive recombination reveals  
642 two-speed evolution of human cytomegalovirus genomes. *Virus Evolution* **2**,  
643 doi:10.1093/ve/vew017 (2016).
- 644 32 Lanier, E. R. *et al.* Analysis of Mutations in the Gene Encoding Cytomegalovirus  
645 DNA Polymerase in a Phase 2 Clinical Trial of Brincidofovir Prophylaxis. *J Infect*  
646 *Dis* **214**, 32-35, doi:10.1093/infdis/jiw073 (2016).
- 647 33 Kaverin, N. V. *et al.* Epitope mapping of the hemagglutinin molecule of a highly  
648 pathogenic H5N1 influenza virus by using monoclonal antibodies. *J Virol* **81**, 12911-  
649 12917, doi:10.1128/JVI.01522-07 (2007).
- 650 34 Franco, S. *et al.* Detection of a sexually transmitted hepatitis C virus protease  
651 inhibitor-resistance variant in a human immunodeficiency virus-infected homosexual  
652 man. *Gastroenterology* **147**, 599-601 e591, doi:10.1053/j.gastro.2014.05.010 (2014).
- 653 35 Fujisaki, S. *et al.* Outbreak of infections by hepatitis B virus genotype A and  
654 transmission of genetic drug resistance in patients coinfecting with HIV-1 in Japan. *J*  
655 *Clin Microbiol* **49**, 1017-1024, doi:10.1128/JCM.02149-10 (2011).
- 656 36 Pierucci, P. *et al.* Novel autologous T-cell therapy for drug-resistant cytomegalovirus  
657 disease after lung transplantation. *J Heart Lung Transplant*,  
658 doi:10.1016/j.healun.2015.12.031 (2016).
- 659 37 Agoti, C. N. *et al.* Local evolutionary patterns of human respiratory syncytial virus  
660 derived from whole-genome sequencing. *J Virol* **89**, 3444-3454,  
661 doi:10.1128/JVI.03391-14 (2015).
- 662 38 Quick, J. *et al.* Real-time, portable genome sequencing for Ebola surveillance. *Nature*  
663 **530**, 228-232, doi:10.1038/nature16996 (2016).
- 664 39 Aanensen, D. M. *et al.* Whole-Genome Sequencing for Routine Pathogen  
665 Surveillance in Public Health: a Population Snapshot of Invasive *Staphylococcus*  
666 *aureus* in Europe. *MBio* **7**, doi:10.1128/mBio.00444-16 (2016).
- 667 40 Mbisa, J. L. *et al.* Evidence of Self-Sustaining Drug Resistant HIV-1 Lineages  
668 Among Untreated Patients in the United Kingdom. *Clin Infect Dis* **61**, 829-836,  
669 doi:10.1093/cid/civ393 (2015).
- 670 41 Faria, N. R. *et al.* Zika virus in the Americas: Early epidemiological and genetic  
671 findings. *Science* **352**, 345-349, doi:10.1126/science.aaf5036 (2016).
- 672 42 Bartha, I. *et al.* A genome-to-genome analysis of associations between human genetic  
673 variation, HIV-1 sequence diversity, and viral control. *eLife* **2**, e01123,  
674 doi:10.7554/eLife.01123 (2013).

- 675 43 Cuevas, J. M., Geller, R., Garijo, R., Lopez-Aldeguer, J. & Sanjuan, R. Extremely  
676 High Mutation Rate of HIV-1 In Vivo. *PLoS Biol* **13**, e1002251,  
677 doi:10.1371/journal.pbio.1002251 (2015).
- 678 44 Vandenhende, M. A. *et al.* Prevalence and evolution of low frequency HIV drug  
679 resistance mutations detected by ultra deep sequencing in patients experiencing first  
680 line antiretroviral therapy failure. *PLoS One* **9**, e86771,  
681 doi:10.1371/journal.pone.0086771
- 682 PONE-D-13-40833 [pii] (2014).
- 683 45 Zhou, B. *et al.* Composition and Interactions of Hepatitis B Virus Quasispecies  
684 Defined the Virological Response During Telbivudine Therapy. *Sci Rep* **5**, 17123,  
685 doi:10.1038/srep17123 (2015).
- 686 46 Itakura, J. *et al.* Resistance-Associated NS5A Variants of Hepatitis C Virus Are  
687 Susceptible to Interferon-Based Therapy. *PLoS One* **10**, e0138060,  
688 doi:10.1371/journal.pone.0138060 (2015).
- 689 47 Rogers, M. B. *et al.* Intrahost dynamics of antiviral resistance in influenza A virus  
690 reflect complex patterns of segment linkage, reassortment, and natural selection.  
691 *MBio* **6**, doi:10.1128/mBio.02464-14 (2015).
- 692 48 Swenson, L. C., Daumer, M. & Paredes, R. Next-generation sequencing to assess HIV  
693 tropism. *Curr Opin HIV AIDS* **7**, 478-485, doi:10.1097/COH.0b013e328356e9da  
694 (2012).
- 695 49 Hutter, G. *et al.* Long-term control of HIV by CCR5 Delta32/Delta32 stem-cell  
696 transplantation. *N Engl J Med* **360**, 692-698, doi:10.1056/NEJMoa0802905 (2009).
- 697 50 Kordelas, L. *et al.* Shift of HIV tropism in stem-cell transplantation with CCR5  
698 Delta32 mutation. *N Engl J Med* **371**, 880-882, doi:10.1056/NEJMc1405805 (2014).
- 699 51 Coquette, A. *et al.* Mixed cytomegalovirus glycoprotein B genotypes in  
700 immunocompromised patients. *Clin Infect Dis* **39**, 155-161, doi:10.1086/421496  
701 (2004).
- 702 52 Solmone, M. *et al.* Use of massively parallel ultradeep pyrosequencing to characterize  
703 the genetic diversity of hepatitis B virus in drug-resistant and drug-naive patients and  
704 to detect minor variants in reverse transcriptase and hepatitis B S antigen. *J Virol* **83**,  
705 1718-1726, doi:10.1128/JVI.02011-08 (2009).
- 706 53 Chou, S. *et al.* Improved detection of emerging drug-resistant mutant cytomegalovirus  
707 subpopulations by deep sequencing. *Antimicrob Agents Chemother* **58**, 4697-4702,  
708 doi:10.1128/AAC.03214-14 (2014).
- 709 54 Fonager, J. *et al.* Identification of minority resistance mutations in the HIV-1  
710 integrase coding region using next generation sequencing. *J Clin Virol* **73**, 95-100,  
711 doi:10.1016/j.jcv.2015.11.009 (2015).
- 712 55 Kyeyune, F. *et al.* Low-Frequency Drug Resistance in HIV-Infected Ugandans on  
713 Antiretroviral Treatment Is Associated with Regimen Failure. *Antimicrob Agents*  
714 *Chemother* **60**, 3380-3397, doi:10.1128/AAC.00038-16 (2016).
- 715 56 Liu, P. *et al.* Direct sequencing and characterization of a clinical isolate of Epstein-  
716 Barr virus from nasopharyngeal carcinoma tissue by using next-generation  
717 sequencing technology. *J Virol* **85**, 11291-11299, doi:10.1128/JVI.00823-11 (2011).
- 718 57 Venter, J. C. *et al.* Environmental genome shotgun sequencing of the Sargasso Sea.  
719 *Science* **304**, 66-74, doi:10.1126/science.1093857 (2004).
- 720 58 Mulcahy-O'Grady, H. & Workentine, M. L. The Challenge and Potential of  
721 Metagenomics in the Clinic. *Frontiers in immunology* **7**, 29,  
722 doi:10.3389/fimmu.2016.00029 (2016).
- 723 59 Morfopoulou, S. *et al.* Human Coronavirus OC43 Associated with Fatal Encephalitis.  
724 *N Engl J Med* **375**, 497-498, doi:10.1056/NEJMc1509458 (2016).

725 60 Naccache, S. N. *et al.* Diagnosis of neuroinvasive astrovirus infection in an  
726 immunocompromised adult with encephalitis by unbiased next-generation  
727 sequencing. *Clin Infect Dis* **60**, 919-923, doi:10.1093/cid/ciu912 (2015).

728 61 Huang, W. *et al.* Whole-Genome Sequence Analysis Reveals the Enterovirus D68  
729 Isolates during the United States 2014 Outbreak Mainly Belong to a Novel Clade. *Sci*  
730 *Rep* **5**, 15223, doi:10.1038/srep15223 (2015).

731 62 Wilson, M. R. *et al.* Actionable diagnosis of neuroleptospirosis by next-generation  
732 sequencing. *N Engl J Med* **370**, 2408-2417, doi:10.1056/NEJMoa1401268 (2014).

733 63 Hoffmann, B. *et al.* A Variegated Squirrel Bornavirus Associated with Fatal Human  
734 Encephalitis. *N Engl J Med* **373**, 154-162, doi:10.1056/NEJMoa1415627 (2015).

735 64 Perlejewski, K. *et al.* Next-generation sequencing (NGS) in the identification of  
736 encephalitis-causing viruses: Unexpected detection of human herpesvirus 1 while  
737 searching for RNA pathogens. *J Virol Methods* **226**, 1-6,  
738 doi:10.1016/j.jviromet.2015.09.010 (2015).

739 65 Duncan, C. J. *et al.* Human IFNAR2 deficiency: Lessons for antiviral immunity. *Sci*  
740 *Transl Med* **7**, 307ra154, doi:10.1126/scitranslmed.aac4227 (2015).

741 66 Depledge, D. P. *et al.* Specific capture and whole-genome sequencing of viruses from  
742 clinical samples. *PLoS One* **6**, e27805, doi:10.1371/journal.pone.0027805 (2011).

743 67 Allen, U. D. *et al.* The genetic diversity of Epstein-Barr virus in the setting of  
744 transplantation relative to non-transplant settings: A feasibility study. *Pediatr*  
745 *Transplant*, doi:10.1111/ptr.12610 (2015).

746 68 Matranga, C. B. *et al.* Enhanced methods for unbiased deep sequencing of Lassa and  
747 Ebola RNA viruses from clinical and biological samples. *Genome Biol* **15**, 519,  
748 doi:10.1186/PREACCEPT-1698056557139770 (2014).

749 69 Calvet, G. *et al.* Detection and sequencing of Zika virus from amniotic fluid of fetuses  
750 with microcephaly in Brazil: a case study. *Lancet Infect Dis*, doi:10.1016/S1473-  
751 3099(16)00095-5 (2016).

752 70 Lei, H. *et al.* Epstein-Barr virus from Burkitt Lymphoma biopsies from Africa and  
753 South America share novel LMP-1 promoter and gene variations. *Sci Rep* **5**, 16706,  
754 doi:10.1038/srep16706 (2015).

755 71 Kohl, C. *et al.* Protocol for metagenomic virus detection in clinical specimens. *Emerg*  
756 *Infect Dis* **21**, 48-57, doi:10.3201/eid2101.140766 (2015).

757 72 Sauvage, V. & Eloit, M. Viral metagenomics and blood safety. *Transfusion clinique*  
758 *et biologique : journal de la Societe francaise de transfusion sanguine* **23**, 28-38,  
759 doi:10.1016/j.tracli.2015.12.002 (2016).

760 73 Lecuit, M. & Eloit, M. The diagnosis of infectious diseases by whole genome next  
761 generation sequencing: a new era is opening. *Frontiers in cellular and infection*  
762 *microbiology* **4**, 25, doi:10.3389/fcimb.2014.00025 (2014).

763 74 Oude Munnink, B. B. *et al.* Autologous antibody capture to enrich immunogenic  
764 viruses for viral discovery. *PLoS One* **8**, e78454, doi:10.1371/journal.pone.0078454  
765 (2013).

766 75 Sabina, J. & Leamon, J. H. Bias in Whole Genome Amplification: Causes and  
767 Considerations. *Methods in molecular biology* **1347**, 15-41, doi:10.1007/978-1-4939-  
768 2990-0\_2 (2015).

769 76 Jensen, R. H. *et al.* Target-dependent enrichment of virions determines the reduction  
770 of high-throughput sequencing in virus discovery. *PLoS One* **10**, e0122636,  
771 doi:10.1371/journal.pone.0122636 (2015).

772 77 Denesvre, C., Dumarest, M., Remy, S., Gourichon, D. & Eloit, M. Chicken skin  
773 virome analyzed by high-throughput sequencing shows a composition highly different  
774 from human skin. *Virus Genes* **51**, 209-216, doi:10.1007/s11262-015-1231-8 (2015).



775 78 Mlakar, J. *et al.* Zika Virus Associated with Microcephaly. *N Engl J Med* **374**, 951-  
776 958, doi:10.1056/NEJMoa1600651 (2016).

777 79 Hall, R. J., Draper, J. L., Nielsen, F. G. & Dutilh, B. E. Beyond research: a primer for  
778 considerations on using viral metagenomics in the field and clinic. *Front Microbiol* **6**,  
779 224, doi:10.3389/fmicb.2015.00224 (2015).

780 80 Greninger, A. L. *et al.* A novel outbreak enterovirus D68 strain associated with acute  
781 flaccid myelitis cases in the USA (2012-14): a retrospective cohort study. *Lancet*  
782 *Infect Dis* **15**, 671-682, doi:10.1016/S1473-3099(15)70093-9 (2015).

783 81 Breitwieser, F. P., Pardo, C. A. & Salzberg, S. L. Re-analysis of metagenomic  
784 sequences from acute flaccid myelitis patients reveals alternatives to enterovirus D68  
785 infection. *F1000Res* **4**, 180, doi:10.12688/f1000research.6743.2 (2015).

786 82 Gardy, J. L. *et al.* Whole-Genome Sequencing of Measles Virus Genotypes H1 and  
787 D8 During Outbreaks of Infection Following the 2010 Olympic Winter Games  
788 Reveals Viral Transmission Routes. *J Infect Dis* **212**, 1574-1578,  
789 doi:10.1093/infdis/jiv271 (2015).

790 83 Cotten, M. *et al.* Deep sequencing of norovirus genomes defines evolutionary patterns  
791 in an urban tropical setting. *J Virol* **88**, 11056-11069, doi:10.1128/JVI.01333-14  
792 (2014).

793 84 Kundu, S. *et al.* Next-generation whole genome sequencing identifies the direction of  
794 norovirus transmission in linked patients. *Clin Infect Dis* **57**, 407-414,  
795 doi:10.1093/cid/cit287 (2013).

796 85 Watson, S. J. *et al.* Molecular Epidemiology and Evolution of Influenza Viruses  
797 Circulating within European Swine between 2009 and 2013. *J Virol* **89**, 9920-9931,  
798 doi:10.1128/JVI.00840-15 (2015).

799 86 Parameswaran, P. *et al.* Genome-wide patterns of intrahuman dengue virus diversity  
800 reveal associations with viral phylogenetic clade and interhost diversity. *J Virol* **86**,  
801 8546-8558, doi:10.1128/JVI.00736-12 (2012).

802 87 Newman, R. M. *et al.* Whole genome pyrosequencing of rare hepatitis C virus  
803 genotypes enhances subtype classification and identification of naturally occurring  
804 drug resistance variants. *J Infect Dis* **208**, 17-31, doi:10.1093/infdis/jis679 (2013).

805 88 Jakava-Viljanen, M. *et al.* Evolutionary trends of European bat lyssavirus type 2  
806 including genetic characterization of Finnish strains of human and bat origin 24 years  
807 apart. *Arch Virol* **160**, 1489-1498, doi:10.1007/s00705-015-2424-0 (2015).

808 89 Brown, J. R. *et al.* Norovirus whole genome sequencing by SureSelect target  
809 enrichment: a robust and sensitive method. *J Clin Microbiol*,  
810 doi:10.1128/JCM.01052-16 (2016).

811 90 Renzette, N., Bhattacharjee, B., Jensen, J. D., Gibson, L. & Kowalik, T. F. Extensive  
812 genome-wide variability of human cytomegalovirus in congenitally infected infants.  
813 *PLoS Pathog* **7**, e1001344, doi:10.1371/journal.ppat.1001344 (2011).

814 91 Bialasiewicz, S. *et al.* Detection of a divergent Parainfluenza 4 virus in an adult  
815 patient with influenza like illness using next-generation sequencing. *BMC Infect Dis*  
816 **14**, 275, doi:10.1186/1471-2334-14-275 (2014).

817 92 Johnson, T. A. *et al.* Clusters of Antibiotic Resistance Genes Enriched Together Stay  
818 Together in Swine Agriculture. *MBio* **7**, e02214-02215, doi:10.1128/mBio.02214-15  
819 (2016).

820 93 Bonsall, D. *et al.* ve-SEQ: Robust, unbiased enrichment for streamlined detection and  
821 whole-genome sequencing of HCV and other highly diverse pathogens. *F1000Res* **4**,  
822 1062, doi:10.12688/f1000research.7111.1 (2015).

823 94 Wylie, T. N., Wylie, K. M., Herter, B. N. & Storch, G. A. Enhanced virome  
824 sequencing using targeted sequence capture. *Genome Res*, doi:10.1101/gr.191049.115  
825 (2015).

826 95 Tsangaras, K. *et al.* Hybridization capture using short PCR products enriches small  
827 genomes by capturing flanking sequences (CapFlank). *PLoS One* **9**, e109101,  
828 doi:10.1371/journal.pone.0109101 (2014).

829 96 Depledge, D. P. *et al.* Deep sequencing of viral genomes provides insight into the  
830 evolution and pathogenesis of varicella zoster virus and its vaccine in humans. *Mol*  
831 *Biol Evol* **31**, 397-409, doi:10.1093/molbev/mst210 (2014).

832 97 Ebert, K., Depledge, D. P., Breuer, J., Harman, L. & Elliott, G. Mode of virus rescue  
833 determines the acquisition of VHS mutations in VP22-negative herpes simplex virus  
834 1. *J Virol* **87**, 10389-10393, doi:10.1128/JVI.01654-13 (2013).

835 98 Palser, A. L. *et al.* Genome diversity of Epstein-Barr virus from multiple tumor types  
836 and normal infection. *J Virol* **89**, 5222-5237, doi:10.1128/JVI.03614-14 (2015).

837 99 Tweedy, J. *et al.* Complete Genome Sequence of the Human Herpesvirus 6A Strain  
838 AJ from Africa Resembles Strain GS from North America. *Genome Announc* **3**,  
839 doi:10.1128/genomeA.01498-14 (2015).

840 100 Donaldson, C. D., Clark, D. A., Kidd, I. M., Breuer, J. & Depledge, D. D. Genome  
841 Sequence of Human Herpesvirus 7 Strain UCL-1. *Genome Announc* **1**,  
842 doi:10.1128/genomeA.00830-13 (2013).

843 101 Kamperschroer, C., Gosink, M. M., Kumpf, S. W., O'Donnell, L. M. & Tartaro, K. R.  
844 The genomic sequence of lymphocryptovirus from cynomolgus macaque. *Virology*  
845 **488**, 28-36, doi:10.1016/j.virol.2015.10.025 (2016).

846 102 Briese, T. *et al.* Virome Capture Sequencing Enables Sensitive Viral Diagnosis and  
847 Comprehensive Virome Analysis. *MBio* **6**, e01491-01415, doi:10.1128/mBio.01491-  
848 15 (2015).

849 103 Naccache, S. N. *et al.* Distinct Zika Virus Lineage in Salvador, Bahia, Brazil. *Emerg*  
850 *Infect Dis* **22**, 1788-1792, doi:10.3201/eid2210.160663 (2016).

851 104 Hofmann, B. Incidental findings of uncertain significance: To know or not to know--  
852 that is not the question. *BMC Med Ethics* **17**, 13, doi:10.1186/s12910-016-0096-2  
853 (2016).

854 105 England, P. H. Good Laboratory Practice when Performing Molecular Amplification  
855 Assays. Q 4 Issue 4.4. *UK Standards for Microbiology Investigations*. **Q 4** (2013).

856 106 Viana, R. V. & Wallis, C. L. *Good Clinical Laboratory Practice (GCLP) for*  
857 *Molecular Based Tests Used in Diagnostic Laboratories*. (INTECH Open Access  
858 Publisher, 2011).

859 107 Blomquist, T., Crawford, E. L., Yeo, J., Zhang, X. & Willey, J. C. Control for  
860 stochastic sampling variation and qualitative sequencing error in next generation  
861 sequencing. *Biomol Detect Quantif* **5**, 30-37, doi:10.1016/j.bdq.2015.08.003 (2015).

862 108 Houldcroft, C. J. & Breuer, J. Tales from the crypt and coral reef: the successes and  
863 challenges of identifying new herpesviruses using metagenomics. *Front Microbiol* **6**,  
864 188, doi:10.3389/fmicb.2015.00188 (2015).

865 109 Munro, A. C. & Houldcroft, C. Human cancers and mammalian retroviruses: should  
866 we worry about bovine leukemia virus? *Future Virology* **11**, 163-166,  
867 doi:10.2217/fvl.16.5 (2016).

868 110 Hue, S. *et al.* Disease-associated XMRV sequences are consistent with laboratory  
869 contamination. *Retrovirology* **7**, 111, doi:10.1186/1742-4690-7-111 (2010).

870 111 Erlwein, O. *et al.* DNA extraction columns contaminated with murine sequences.  
871 *PLoS One* **6**, e23484, doi:10.1371/journal.pone.0023484 (2011).

872 112 Rosseel, T., Pardon, B., De Clercq, K., Ozhelvaci, O. & Van Borm, S. False-positive  
873 results in metagenomic virus discovery: a strong case for follow-up diagnosis.  
874 *Transbound Emerg Dis* **61**, 293-299, doi:10.1111/tbed.12251 (2014).

875 113 Naccache, S. N. *et al.* The perils of pathogen discovery: origin of a novel parvovirus-  
876 like hybrid genome traced to nucleic acid extraction spin columns. *J Virol* **87**, 11966-  
877 11977, doi:10.1128/JVI.02323-13 (2013).

878 114 Salter, S. J. *et al.* Reagent and laboratory contamination can critically impact  
879 sequence-based microbiome analyses. *BMC Biol* **12**, 87, doi:10.1186/s12915-014-  
880 0087-z (2014).

881 115 Lipkin, W. I. A Vision for Investigating the Microbiology of Health and Disease. *J*  
882 *Infect Dis* **212 Suppl 1**, S26-30, doi:10.1093/infdis/jiu649 (2015).

883 116 Shafer, R. W. Rationale and uses of a public HIV drug-resistance database. *J Infect*  
884 *Dis* **194 Suppl 1**, S51-58, doi:10.1086/505356 (2006).

885 117 Gnaneshan, S., Ijaz, S., Moran, J., Ramsay, M. & Green, J. HepSEQ: International  
886 Public Health Repository for Hepatitis B. *Nucleic Acids Res* **35**, D367-370,  
887 doi:10.1093/nar/gkl874 (2007).

888 118 Rhee, S. Y. *et al.* Hepatitis B virus reverse transcriptase sequence variant database for  
889 sequence analysis and mutation discovery. *Antiviral Res* **88**, 269-275,  
890 doi:10.1016/j.antiviral.2010.09.012 (2010).

891 119 Kuiken, C., Yusim, K., Boykin, L. & Richardson, R. The Los Alamos hepatitis C  
892 sequence database. *Bioinformatics* **21**, 379-384, doi:10.1093/bioinformatics/bth485  
893 (2005).

894 120 Hong, L. Z. *et al.* BAsE-Seq: a method for obtaining long viral haplotypes from short  
895 sequence reads. *Genome Biol* **15**, 517, doi:10.1186/PREACCEPT-6768001251451949  
896 (2014).

897 121 Schmidt, K. *et al.* Identification of bacterial pathogens and antimicrobial resistance  
898 directly from clinical urines by nanopore-based metagenomic sequencing. *J*  
899 *Antimicrob Chemother*, doi:10.1093/jac/dkw397 (2016).

900 122 Kilianski, A. *et al.* Bacterial and viral identification and differentiation by amplicon  
901 sequencing on the MinION nanopore sequencer. *Gigascience* **4**, 12,  
902 doi:10.1186/s13742-015-0051-z (2015).

903 123 Eckert, S. E., Chan, J. Z.-M., Houniet, D., Breuer, J. & Speight, G. Enrichment of  
904 long DNA fragments from mixed samples for Nanopore sequencing. *bioRxiv*,  
905 doi:10.1101/048850 (2016).

906 124 Greninger, A. L. *et al.* Rapid metagenomic identification of viral pathogens in clinical  
907 samples by real-time nanopore sequencing analysis. *Genome Med* **7**, 99,  
908 doi:10.1186/s13073-015-0220-9 (2015).

909 125 Mathijs, E., Vandenbussche, F., Verpoest, S., De Regge, N. & Van Borm, S.  
910 Complete Genome Sequence of Pseudorabies Virus Reference Strain NIA3 Using  
911 Single-Molecule Real-Time Sequencing. *Genome Announc* **4**,  
912 doi:10.1128/genomeA.00440-16 (2016).

913 126 Bull, R. A. *et al.* A method for near full-length amplification and sequencing for six  
914 hepatitis C virus genotypes. *BMC Genomics* **17**, 247, doi:10.1186/s12864-016-2575-8  
915 (2016).

916 127 Kimberlin, D. W. & Whitley, R. J. Antiviral resistance: mechanisms, clinical  
917 significance, and future implications. *J Antimicrob Chemother* **37**, 403-421 (1996).

918 128 Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of  
919 next-generation sequencing technologies. *Nat Rev Genet* **17**, 333-351,  
920 doi:10.1038/nrg.2016.49 (2016).

921 129 Karamitros, T. & Magiorkinis, G. A novel method for the multiplexed target  
922 enrichment of MinION next generation sequencing libraries using PCR-generated  
923 baits. *Nucleic Acids Res* **43**, e152, doi:10.1093/nar/gkv773 (2015).

924 130 Guan, H. *et al.* Detection of virus in CSF from the cases with meningoencephalitis by  
925 next-generation sequencing. *J Neurovirol* **22**, 240-245, doi:10.1007/s13365-015-  
926 0390-7 (2016).

927 131 Brown, J. R. *et al.* Astrovirus VA1/HMO-C: an increasingly recognized neurotropic  
928 pathogen in immunocompromised patients. *Clin Infect Dis* **60**, 881-888,  
929 doi:10.1093/cid/ciu940 (2015).

930 132 Fremont, M. L. *et al.* Next-Generation Sequencing for Diagnosis and Tailored  
931 Therapy: A Case Report of Astrovirus-Associated Progressive Encephalitis. *J*  
932 *Pediatric Infect Dis Soc* **4**, e53-57, doi:10.1093/jpids/piv040 (2015).

933 133 Musso, D. *et al.* Potential for Zika virus transmission through blood transfusion  
934 demonstrated during an outbreak in French Polynesia, November 2013 to February  
935 2014. *Euro Surveill* **19** (2014).

936 134 Barjas-Castro, M. L. *et al.* Probable transfusion-transmitted Zika virus in Brazil.  
937 *Transfusion*, doi:10.1111/trf.13681 (2016).

938 135 Ellison, D. W. *et al.* Complete Genome Sequences of Zika Virus Strains Isolated from  
939 the Blood of Patients in Thailand in 2014 and the Philippines in 2012. *Genome*  
940 *Announc* **4**, doi:10.1128/genomeA.00359-16 (2016).

941 136 Faria, N. R. *et al.* Mobile real-time surveillance of Zika virus in Brazil. *Genome Med*  
942 **8**, 97, doi:10.1186/s13073-016-0356-2 (2016).

943