

1 **Quality control implementation for universal characterization of**
2 **DNA and RNA viruses in clinical respiratory samples using single**
3 **metagenomic next-generation sequencing workflow**

4 A. Bal^{1, 2, 3, 4} antonin.bal@chu-lyon.fr
5 M. Pichon^{1, 2, 3} maxime.pichon01@chu-lyon.fr
6 C. Picard⁵ caroline.picard@pasteur.fr
7 JS. Casalegno^{1, 2, 3} jean-sebastien.casalegno@chu-lyon.fr
8 M. Valette^{1, 2, 3} martine.valette@chu-lyon.fr
9 I. Schuffenecker¹ isabelle.schuffenecker@chu-lyon.fr
10 L. Billard⁶ billardl29@gmail.com
11 S. Vallet^{6, 7} sophie.vallet@chu-brest.fr
12 G. Vilchez⁴ gaelle.vilchez@chu-lyon.fr
13 V. Cheynet⁴ valerie.cheynet@biomerieux.com
14 G. Oriol⁴ guy.oriol@biomerieux.com
15 S. Assant⁴ sophie.assant@chu-lyon.fr
16 Y. Gillet⁸ yves.gillet@chu-lyon.fr
17 B. Lina^{1, 2, 3} bruno.lina@chu-lyon.fr
18 K. Brengel-Pesce⁴ karen.brengel-pesce@biomerieux.com
19 F. Morfin^{1, 2, 3} florence.morfin-sherpa@chu-lyon.fr
20 L. Josset^{1, 2, 3} laurence.josset@chu-lyon.fr

21 ¹Laboratoire de Virologie, Institut des Agents Infectieux, Groupement Hospitalier Nord, Hospices Civils de
22 Lyon, Lyon, France

23 ²Univ Lyon, Université Lyon 1, Faculté de Médecine Lyon Est, CIRI, Inserm U1111 CNRS UMR5308, Virpath,
24 Lyon, France

25 ³Hospices Civils de Lyon, Centre National de Référence des virus respiratoires France Sud, Lyon, France

26 ⁴Laboratoire Commun de Recherche HCL-bioMerieux, Centre Hospitalier Lyon Sud, Pierre-Bénite, France

27 ⁵Unité de Biologie des Infections Virales Emergentes, Institut Pasteur, Lyon, France; CIRI Inserm U1111,
28 CNRS 5308, ENS, UCBL, Faculté de Médecine Lyon Est, Université de Lyon, Lyon, France.

29 ⁶INSERM UMR1078 "Génomique, Génétique Fonctionnelle et Biotechnologies", Axe Microbiota, Univ Brest,
30 Brest, France

31 ⁷Département de Bactériologie-Virologie, Hygiène et Parasitologie-Mycologie, Pôle de Biologie-Pathologie,
32 Centre Hospitalier Régional et Universitaire de Brest, Hôpital de la Cavale Blanche, Brest, France

33 ⁸Hospices Civils de Lyon, Urgences pédiatriques, Hôpital Femme Mère Enfant, Bron, France

34 **Corresponding author:** Laurence Josset, Pharm.D, Ph.D

35 Associate Professor
36 Hospices Civils de Lyon
37 National reference center for respiratory viruses
38 103 Grande-Rue de la Croix Rousse
39 69317, Lyon
40 France
41 Telephone: +33 (0)4 72 07 10 22
42 Email: laurence.josset@chu-lyon.fr

43 **Keywords:** clinical virology, quality control, next-generation sequencing; viral metagenomics; respiratory
44 viruses

45 **Abstract**

46 **Background**

47 In recent years, metagenomic Next-Generation Sequencing (mNGS) has increasingly been
48 used for an accurate assumption-free virological diagnosis. However, the systematic
49 workflow evaluation on clinical respiratory samples and implementation of quality controls
50 (QCs) is still lacking.

51 **Methods**

52 A total of 3 QCs were implemented and processed through the whole mNGS workflow: a no-
53 template-control to evaluate contamination issues during the process; an internal and an
54 external QC to check the integrity of the reagents, equipment, the presence of inhibitors, and
55 to allow the validation of results for each sample. The workflow was then evaluated on 37
56 clinical respiratory samples from patients with acute respiratory infections previously tested
57 for a broad panel of viruses using semi-quantitative real-time PCR assays (28 positive
58 samples including 6 multiple viral infections; 9 negative samples). Selected specimens
59 included nasopharyngeal swabs (n = 20), aspirates (n = 10), or sputums (n = 7).

60 **Results**

61 The optimal spiking level of the internal QC was first determined in order to be sufficiently
62 detected without overconsumption of sequencing reads. According to QC validation criteria,
63 mNGS results were validated for 34/37 selected samples. For valid samples, viral genotypes
64 were accurately determined for 36/36 viruses detected with PCR (viral genome coverage
65 ranged from 0.6% to 100%, median = 67.7%). This mNGS workflow allowed the detection of
66 DNA and RNA viruses up to a semi-quantitative PCR Ct value of 36. The six multiple viral
67 infections involving 2 to 4 viruses were also fully characterized. A strong correlation between
68 results of mNGS and real-time PCR was obtained for each type of viral genome (R^2 ranged
69 from 0.72 for linear single-stranded (ss) RNA viruses to 0.98 for linear ssDNA viruses).

70 **Conclusions**

71 Although the potential of mNGS technology is very promising, further evaluation studies are
72 urgently needed for its routine clinical use within a reasonable timeframe. The approach
73 described herein is crucial to bring standardization and to ensure the quality of the generated
74 sequences in clinical setting. We provide an easy-to-use single protocol successfully
75 evaluated for the characterization of a broad and representative panel of DNA and RNA
76 respiratory viruses in various types of clinical samples.

77 **Background**

78 Since the development of Next Generation-Sequencing (NGS) technologies in 2005,
79 the use of metagenomic approaches has grown considerably. It is now considered as an
80 efficient unbiased tool in clinical virology[1,2], in particular for the characterization of viral
81 acute respiratory infections (ARIs). Several advantages of metagenomic NGS (mNGS)
82 compared to conventional real-time Polymerase Chain Reaction (PCR) assays have been
83 highlighted. Firstly, the full viral genetic information is immediately available allowing the
84 investigation of respiratory outbreaks, viral epidemiological surveillance, or identification of
85 specific mutations leading to antiviral resistance or higher virulence [3–5]. Secondly, a
86 significant improvement in viral ARIs diagnosis has been reported [4,6–9]; as the process is
87 sequence independent, mNGS is able to identify highly divergent viral genomes, rare
88 respiratory pathogens, and to discover respiratory viruses missed by targeted PCR [1,4,7].

89 However, the diversity in viral nucleic acid types has impaired the development of a
90 unique viral metagenomic workflow allowing the comprehensive characterization of viruses
91 present in a clinical sample. Most of the published viral metagenomic protocols have been
92 optimized for the detection either of DNA viruses or RNA viruses [4,5,10–13]. In addition,
93 despite the growing number of studies using a metagenomic process in clinical virology,
94 evaluation of workflows has not systematically included both clinical samples and quality
95 control (QC) implementation. A metagenomic protocol involves a large number of steps and
96 all of these have to be controlled to ensure the quality of the generated sequences [6,14–16].
97 Furthermore, specimen to specimen, environmental, and reagent contaminations are also a
98 major concern in metagenomic setting and must be accurately evaluated [6,17–19].

99 The objective of this study was to implement QCs in a single metagenomic protocol and to
100 evaluate it for the detection of a broad panel of DNA and RNA viruses in clinical respiratory
101 samples.

102 **Methods**

103 **Clinical samples**

104 A total of 37 respiratory samples collected from patients hospitalized in the university
105 hospital of Lyon (Hospices Civils de Lyon, HCL) were retrospectively selected to evaluate
106 our metagenomic approach. Selected specimens included various types of clinical samples;
107 nasopharyngeal swabs (n=20), aspirates (n=10), or sputums (n=7). These samples were
108 initially sent to our laboratory for routine viral diagnosis of ARI using semi-quantitative real-
109 time PCR assays targeting a comprehensive panel of DNA and RNA viruses (r-gene,
110 bioMérieux, Marcy l'étoile, France). This panel included: influenza virus type A and B,
111 adenovirus, cytomegalovirus, Epstein-Barr virus, human herpes virus 6, human bocavirus
112 (HBoV), human rhinovirus, respiratory syncytial virus, human parainfluenza virus, human
113 coronavirus (HCoV), human metapneumovirus, and measles virus. Twenty-two samples were
114 positive for only one targeted virus, 6 were characterized by a multiple viral infection and 9
115 were negative for all the targeted viruses. These 9 samples were also found to be negative
116 using the FilmArray Respiratory Panel (FA RP, bioMérieux). After PCR testing, the rest of
117 samples were stored at -20°C until mNGS analysis.

118 **Metagenomic workflow**

119 For sample viral enrichment, a 3-step method was applied to 200µl of thawed and vortexed
120 sample [20]: low-speed centrifugation (6000g, 10 min, 4 °C), followed by filtration of the
121 supernatant using 0.80 µm filter (Sartorius, Göttingen, Germany) to remove eukaryotic and
122 bacterial cells, without loss of large viruses [21] and then Turbo DNase treatment (0.1U/µL,
123 37 °C, 90 min; Life Technologies, Carlsbad, CA, USA). Total nucleic acid was extracted
124 using the NucliSENS EasyMAG platform (bioMérieux, Marcy l'Etoile, France) followed by
125 an ethanol precipitation (2 hours at -80°C). As previously described, modified whole
126 transcriptome amplification was performed to amplify both DNA and RNA viral nucleic acids

127 (WTA2, Sigma-Aldrich, Darmstadt, Germany) [21]. Amplified DNA and cDNA were then
128 purified using a QiaQuick column (Qiagen, Hilden, Germany) and quantified using the Qubit
129 fluorometer HS dsDNA Kit (Life Technologies, Carlsbad, CA, USA). Nextera XT
130 DNA Library preparation and Nextera XT Index Kit were used to prepare paired-end
131 libraries, according to the manufacturer's recommendations (Illumina, San Diego, CA, USA).
132 After normalization, a pool of libraries (V/V) was made and quantified using universal KAPA
133 library quantification kit (Kapa Biosystems, Wilmington, MA, USA); 1% PhiX genome was
134 added to the quantified library before sequencing with Illumina NextSeq 500™ platform
135 (Fig. 1). In addition, it should be noticed that our wet-lab process was designed to prevent
136 contaminations as much as possible: reagents were stored and prepared in a DNA-free room;
137 patient samples were opened in a laminar flow hood in a pre-PCR room; after the
138 amplification step, tubes were handled and stored in a post-PCR room.

139 **Bioinformatic analysis**

140 A stepwise bioinformatic filtering pipeline was used to quality filter reads using cutadapt and
141 sickle; and to remove human, archaeal, bacterial, and fungal sequences by aligning reads with
142 bwa mem. The databases used were GRCh38.p2, RefSeq archaea, RefSeq bacteria, and
143 RefSeq fungi. Remaining reads were aligned on ezVIR viral database v0.1 [22] and
144 bacteriophage genomes from the RefSeq database (downloaded on 17 February 2017) using
145 bwa mem. Normalization for comparing viral genome coverage values was performed using
146 reads per kilobase of virus reference sequence per million mapped reads (RPKM) ratio [4,23].
147 RPKM ratio corrects differences in both sample sequencing depth and viral gene length. Viral
148 reads (expressed in RPKM) from the No-Template Control (NTC) were subtracted from viral
149 reads (in RPKM) of each sample within the batch prior to further analysis. A sample was
150 considered to be positive for a particular virus when the RPKM of this virus was positive. No
151 threshold regarding genome coverage pattern was applied nor requirement to cover a

152 particular region of the genome. This latter requirement could be important to correctly
153 identify RNA virus subtypes with high recombination frequencies within a species, but has to
154 be implemented specifically for each viral family.

155 **Quality control implementation**

156 All respiratory specimens were spiked with internal quality control (IQC) before sample
157 preparation. MS2 bacteriophage from a commercial kit (MS2, IC1 RNA internal control; r-
158 gene, bioMérieux) was selected as the IQC. As positive external quality control (EQC), we
159 used viral transport medium spiked with MS2 at the same concentration used for the IQC. A
160 No-Template Control (NTC) was implemented to evaluate contamination during the process.
161 NTC was constituted of viral transport medium (Sigma-virocult, MWE, Corsham, UK) that
162 was processed through all mNGS steps. Two QC testing (QCT) were performed: QCT1 which
163 was the semi-quantitative detection of MS2 using a commercial real-time PCR assay (IC1
164 RNA internal control, r-gene, bioMérieux,) after amplification step (Fig. 1). QCT1 validation
165 criteria were: MS2 semi-quantitative PCR Cycle threshold (Ct) below 37 Ct for IQC and
166 EQC, and no MS2 detection for NTC. QCT2 evaluated the sequencing performance by
167 quantifying the number of reads aligned on the MS2 genome (in RPKM) and MS2 genome
168 coverage (MS2 genome accession number: NC_001417.2; Fig. 1). QCT2 validation criteria
169 were MS2 genome coverage >95% for positive EQC, and an MS2 RPKM > 0 for IQC.

170 **Statistical analysis**

171 Statistical analyses were performed using GraphPrism version 5.02 applying the appropriate
172 statistical test (associations between mNGS and viral real-time PCR assay were determined
173 by applying the Pearson's correlation coefficient and differences between median and
174 distributions were evaluated by the Mann–Whitney U test). A p-value less than 0.05 was
175 considered to be statistically significant.

176 **Results**

177 **Determination of optimal internal quality control spiking**

178 MS2 bacteriophage (MS2), a single-stranded RNA virus (ssRNA), was used as the IQC to
179 validate the whole metagenomic process for each sample. In order to optimize IQC spiking
180 level, the sensitivity of the metagenomic analysis workflow for MS2 detection was first
181 evaluated with a ten-fold serial dilutions of MS2 (from 10^{-2} to 10^{-5}) in a nasopharyngeal swab
182 tested negative using FA RP (bioMérieux). MS2 was detected in internal QCT1 (IQCT1) for
183 all levels of MS2 spiking (Ct ranged from 17.5 at the 10^{-2} dilution to 26.4 Ct at the 10^{-5}
184 dilution). Full to partial MS2 genome coverage was obtained for all MS2 spiking levels in
185 internal QCT2 (IQCT2; coverage ranged from 98% at the 10^{-2} dilution to 69% at the 10^{-5}
186 dilution). For the highest spiking level, 66.0% of the total number of viral reads was mapped
187 to MS2; for the lowest spiking level, 0.9% were so (Fig. 2). To limit the number of NGS reads
188 consumed for IQC detection, the optimal spiking condition was determined to be the 10^{-5}
189 dilution and was used for the rest of the study.

190 **Validation of mNGS results**

191 A total of 37 clinical respiratory samples from patients with ARIs caused by a broad panel of
192 DNA and RNA viruses or of unknown etiology were analyzed in a single mNGS workflow.
193 Libraries were sequenced to a mean of 5,139,248 million reads passing quality filters (range:
194 270,975 to 13,586,456 reads). Human sequences represented the main part of NGS reads for
195 both positive samples (mean = 61.3%) and negative samples (mean = 67.1%), but not of NTC
196 which was mainly composed of bacterial reads (67.8%). The proportion of viral reads ranged
197 from 0.006% to 85.2% (mean = 9.6% for positive samples and 0.6 % for negative samples,
198 Additional file 1). Viral metagenomic results were then validated according to the criteria
199 described in the Methods section. QCT1 (MS2 molecular detection performed before library
200 preparation) was negative for NTC. After sequencing, viral contamination represented 0.13%

201 (4,245/3,215,616) of the total reads generated from NTC including 2 MS2 reads (MS2 RPKM
202 = 173). For targeted viruses, 21 reads (RPKM = 480) and 185 reads (RPKM = 1.1E+04)
203 mapping to influenza A(H3N2) and HBoV were detected, respectively. The positive EQC was
204 successfully detected at QCT1 (MS2 PCR positive at 25 Ct) and after the sequencing step
205 (QCT2; MS2 genome coverage = 99.7%, MS2 RPKM = 5.5E+05). Regarding IQC results,
206 37/37 samples passed QCT1 (MS2 PCR Ct values <37) and were therefore further processed.
207 A total of 33/37 samples passed QCT2 (MS2 RPKM > 0; Fig. 3). For these 33 samples, MS2
208 genome coverage ranged from 15% to 100% (Additional file 2).

209 The 4 samples that did not pass IQCT2 included one sputum that was previously tested
210 negative using real-time PCR (sample # 37), one HCoV positive sputum (sample # 11, Ct =
211 32), one HBoV positive nasopharyngeal swab (sample # 19, Ct = 30), and one
212 nasopharyngeal aspirate tested positive for HBoV and CMV (sample # 23, Ct = 15 and 31,
213 respectively). For sample # 37 and sample # 19, none of the real-time PCR targeted viruses
214 were detected after bioinformatic analysis. For sample # 19, we sequenced a replicate which
215 similarly failed both IQC and HBoV detection. We could not test any replicate for sample #
216 37 owing to insufficient quantity. Viral metagenomics results for sample # 23 were validated
217 as viral reads represented 85.2% (9,489,578/11,144,324) of the total reads generated (Fig. 3).
218 For sample # 11, the number of reads mapping to HCoV was 9/5,125,947 with a HCoV
219 genome coverage of 0.2%. Results were therefore not validated for this sample. Overall,
220 mNGS results were validated for 34/37 samples including 26/28 positive samples and 8/9
221 negative samples.

222 **Metagenomic workflow evaluation according to viral genome type**

223 The evaluation of the metagenomic workflow was performed using the 26 previously
224 validated respiratory samples tested positive with viral real-time PCR targeting a
225 representative panel of DNA and RNA viruses. For all 26 samples tested, viral metagenomic

226 sequencing allowed the identification of the 36/36 viral genotypes matching targeted PCR
227 results and on-target viral genome coverage ranged from 0.6 to 100% (median = 67.7%). For
228 these 36 targeted viruses, the real-time PCR Ct values ranged from 15 to 37 Ct (median = 28
229 Ct). The six multiple viral infections involving from 2 to 4 different viruses were also fully
230 characterized (Table 1). For sample # 25 (sample tested positive for 2 DNA viruses and 2
231 RNA viruses using real-time PCR), mNGS results were cross-checked on a duplicate which
232 reported RPKM deviations lower than 0.5 log for each targeted virus (mNGS results for the 2
233 replicates are summarized in Additional file 3). Regarding mNGS results obtained from the 8
234 negative samples validated with IQC, no clinically relevant virus was detected. A strong
235 correlation between mNGS and real-time PCR results was obtained for each viral genome
236 type (R^2 ranged from 0.72 for linear ssRNA viruses to 0.98 for linear ssDNA viruses, Fig. 4a).
237 Normalized read counts were significantly lower for linear dsDNA viruses than for other viral
238 genome types (Fig. 4b).

239 **Discussion**

240 Over the last few years, a growing number of viral metagenomic protocols have been
241 published but systematic evaluation on clinical respiratory samples and validation by QC is
242 still lacking. In the present study, we describe a process allowing the sensitive detection of
243 both DNA and RNA viruses in a single assay and implemented several QCs to validate the
244 whole metagenomic workflow.

245 First, IQC was implemented to control the integrity of the reagents, equipment, the presence
246 of inhibitors, and to allow the validation of mNGS results for each sample. The MS2
247 bacteriophage was selected as IQC for three main reasons; firstly MS2 is widely used as IQC
248 during viral real-time PCR assays to control both extraction and inhibition [24], secondly, an
249 RNA virus was required to control the random reverse transcription and second strand
250 synthesis steps, and thirdly MS2 is a ssRNA virus with a small genome (3,569-bp) that is
251 perfectly characterized and therefore can be easily detected after bioinformatic analysis
252 without the need for extensive NGS reads. The use of MS2 as an IQC has been previously
253 reported for metagenomic analysis of cerebrospinal fluid specimens [25]. In another
254 metagenomic study, RNA of MS2 was included after extraction as an IQC but the use of
255 purified RNA does not validate the viral enrichment step [26]. In the protocol described
256 herein, whole MS2 virions were added to each clinical sample from the beginning of the
257 workflow. QCT1 was implemented to control the first steps of the process and to avoid
258 unnecessary library preparation when these steps fail. At the end of the workflow, QCT2 was
259 able to invalidate 2 samples as neither MS2 nor viruses causing ARIs were significantly
260 detected after metagenomic analysis while routine PCR screenings detected a HBoV and a
261 HCoV. The re-testing of these 2 samples found the same findings suggesting an inhibition or
262 a competition issue during the process. Without the use of IQC, these samples would have
263 been mistakenly classified as false negatives by mNGS. However, the expected competition

264 between viruses and MS2 during the process could lead to a non-detection of IQC reads in
265 case of high viral load. Thus, the interpretation of IQC results should consider the proportion
266 of viral reads of each sample. Although not observed, IQC reads may also be reduced in
267 samples with a greater numbers of patient cells which may affect the sensitivity of the assay
268 [25].

269 In addition to IQC, we implemented negative external control because contamination issues
270 are frequently reported in metagenomic studies and may lead to misinterpretation in clinical
271 practice [17]. mNGS reads in this negative control were mainly composed of bacterial reads.
272 However, viral reads (mainly derived from prokaryote viruses) were also detected which
273 could be present in reagents (“kitome”) or may represent laboratory contaminants or bleed-
274 over contaminations from highly positive samples within the batch. Such contamination was
275 observed in the present study from the highly positive HBoV sample (sample # 23, Ct=15)
276 which contaminated the NTC (HBoV: 185 reads, RPKM = 1.1E+04 RPKM). In the clinical
277 setting, subtracting NTC viral reads prior to interpretation of each sample result is therefore
278 required.

279 To evaluate the workflow, clinical respiratory samples tested for a representative panel of
280 DNA and RNA viruses using real-time PCR were selected. This workflow is based on a
281 previous publication where a single protocol had been specifically developed for stool
282 specimens and evaluated on mock communities containing high concentrations of spiked
283 viruses [21]. Interestingly, 6 multiple viral infections involving both DNA and RNA viruses
284 were fully characterized highlighting the power of our mNGS approach as a universal method
285 for virus characterization despite the lack of common viral sequence. In addition to viruses
286 targeted by PCR, viral reads deriving from the commensal virome, including viruses from the
287 *Anelloviridae* family, were generated both in PCR negative and positive samples but not in
288 the NTC.

289 Regarding the sensitivity of the mNGS approach, a wide range of semi-quantitative real-time
290 PCR Ct values was covered. Thorburn *et al.*, compared mNGS to conventional real-time
291 PCR for the detection of RNA viruses on nasopharyngeal swabs and reported a detection cut-
292 off of 32 Ct for the mNGS approach [27]. Our workflow allowed the characterization of both
293 DNA and RNA viruses up to a semi-quantitative real-time PCR Ct value of 36 which is
294 considered to be a low viral load. A major critical point in viral metagenomics is to reduce
295 host and bacterial components. In comparison with similar studies, viral reads herein were
296 highly represented (mean = 7.4%); for example, a study on 16 nasopharyngeal aspirates tested
297 positive with viral PCR assays found a mean of 0.05% of viral reads [12]. In addition, a
298 strong correlation between results of mNGS and conventional real-time PCR was obtained by
299 regrouping viruses according to their genome types. Similar findings were reported
300 elsewhere, suggesting that mNGS results could be used for semi-quantitative measurement of
301 the viral load in clinical samples [3–5,28]. A lower RPKM values for dsDNA viruses
302 compared to the other viral genome types were noticed. As previously described for EBV and
303 CMV, the necessary use of DNase to reduce host contamination may affect these fragile large
304 dsDNA viruses [9,10]. As the detection limit of mNGS analysis is mainly dependent on viral
305 load and total number of reads per sample, this effect could be overcome by increasing
306 sequencing depth; however, we chose to limit the costs of the workflow.

307 The reagent cost of this mNGS approach is relatively low and was estimated to ~€150 thanks
308 to our viral enrichment process and the amplification method using a commercial kit which is
309 diluted 5-fold [21]. The use of a universal workflow for both DNA and RNA viruses also
310 reduces the reagent cost compared with metagenomic protocols targeting DNA and RNA
311 viruses separately. In contrast, targeted NGS of specific viruses following their specific
312 amplification by PCR can be up to 2 times cheaper based on our experience (e.g. influenza
313 virus sequencing [29]). Due to several limitations, including its cost and a long turnaround

314 time, viral metagenomics is currently considered to be a second-line approach and is not used
315 as a primary routine diagnostic tool. However, with the improvement of sequencing
316 technologies allowing real-time sequencing such as MinION sequencers (Oxford nanopore,
317 Oxford, United Kingdom), it could be envisioned that mNGS will gradually be used for
318 primary diagnosis in the mid-term. In case of high viral load and sufficient DNA input after
319 amplification our workflow might be used with a MinION sequencer.

320 The approach described in this preliminary work is crucial to bring standardization for the
321 routine clinical use of mNGS process within a reasonable timeframe. Further evaluation
322 studies with a greater number of samples are urgently needed to establish IQC cut-off
323 according to the number of viral, human and bacterial reads, and to define the performance of
324 the workflow, including repeatability, reproducibility, as well as the detection limit for each
325 virus. In addition, improvement of the bioinformatics pipeline are being explored, including
326 implementation of threshold regarding genome coverage pattern [25], but their impact on
327 performance of the workflow has to be established.

328 **Conclusion**

329 The potential of mNGS is very promising but several factors such as inhibition, competition,
330 and contamination can lead to a dramatic misinterpretation in the clinical setting. Herein, we
331 provide an efficient and easy to use mNGS workflow including quality controls successfully
332 evaluated for the comprehensive characterization of a broad and representative panel of DNA
333 and RNA viruses in various types of clinical respiratory samples.

334 **Abbreviations**

335 NGS: Next-Generation Sequencing, mNGS: metagenomic Next-Generation Sequencing,
336 ARIs: Acute Respiratory Infections, PCR: Polymerase Chain Reaction, QC: quality controls,
337 HCL: Hospices Civils de Lyon, IQC: Internal Quality Control, MS2: MS2 bacteriophage,
338 EQC: External Quality Control, NTC: No-Template Control, QCT Quality Control Testing,
339 Ct: Cycle threshold, RPKM: reads per kilobase of virus reference sequence per million
340 mapped reads

341 **Ethics approval and consent to participate**

342 This single center retrospective study received approval from HCL board of the French data
343 protection authority (*Commission Nationale de l'Informatique et des Libertés*) and is
344 registered with the national data protection agency (number 17-024). Respiratory samples
345 were collected for regular disease management during hospital stay and no additional samples
346 were taken for this study. In accordance with French legislation relating to this type a study, a
347 written informed consent from participants was not required for the use of de-identified
348 collected clinical samples (Bioethics law number 2004-800 of August 6, 2004). During their
349 hospitalization in the HCL, patients are made aware that their de-identified data including
350 clinical samples may be used for research purposes, and they can opt out if they object to the
351 use of their data.

352 **Consent for publication**

353 Not applicable.

354 **Availability of data and materials**

355 Data generated during this study are included in supplementary information files. Sequencing
356 datasets used and/or analysed during the current study are available from the corresponding
357 author on reasonable request.

358 **Competing interests**

359 AB has served as consultant to bioMérieux. KB, VC, GO are employees of bioMérieux.

360 **Funding**

361 This study was funded by a metagenomic grant received in 2014 from the French foundation
362 of innovation in infectious diseases (FINOVI, *fondation innovation en infectiologie*).

363 **Authors' contributions**

364 AB, LJ, FM, KB SA conceived the study. AB, MP, LB, VC performed the sample
365 preparations and sequencing. LJ, GO, GV performed bioinformatic analysis. LJ is the
366 guarantor for the NGS data. YG, MV, IS, BL, SV, FM are the guarantor for clinical data and
367 sample collection. AB was the main writer of the manuscript. All authors reviewed and
368 approved the final version of the manuscript.

369 **Acknowledgments**

370 We thank Audrey Guichard, Gwendolyne Burfin, Delphine Falcon and Cecile Darley for their
371 technical assistance as well as Philip Robinson (DRCI, Hospices Civils de Lyon) for his
372 excellent help in manuscript preparation. Part of these data has been presented at the
373 International Conference of Clinical Metagenomic held in Geneva in October 2017.

374 **References**

- 375
- 376 1. Mokili JL, Rohwer F, Dutilh BE. Metagenomics and future perspectives in virus
377 discovery. *Curr Opin Virol.* 2012;2:63–77.
 - 378 2. Capobianchi MR, Giombini E, Rozera G. Next-generation sequencing technology in
379 clinical virology. *Clin Microbiol Infect.* 2013;19:15–22.
 - 380 3. Prachayangprecha S, Schapendonk CME, Koopmans MP, Osterhaus ADME, Schürch AC,
381 Pas SD, et al. Exploring the potential of next-generation sequencing in detection of
382 respiratory viruses. *J Clin Microbiol.* 2014;52:3722–30.
 - 383 4. Graf EH, Simmon KE, Tardif KD, Hymas W, Flygare S, Eilbeck K, et al. Unbiased
384 Detection of Respiratory Viruses by Use of RNA Sequencing-Based Metagenomics: a
385 Systematic Comparison to a Commercial PCR Panel. *J Clin Microbiol.* 2016;54:1000–7.
 - 386 5. Fischer N, Indenbirken D, Meyer T, Lütgehetmann M, Lellek H, Spohn M, et al.
387 Evaluation of Unbiased Next-Generation Sequencing of RNA (RNA-seq) as a Diagnostic
388 Method in Influenza Virus-Positive Respiratory Samples. *J Clin Microbiol.*
389 2015;53:2238–50.
 - 390 6. Schlager R, Queen K, Simmon K, Tardif K, Stockmann C, Flygare S, et al. Viral
391 Pathogen Detection by Metagenomics and Pan-Viral Group Polymerase Chain Reaction
392 in Children With Pneumonia Lacking Identifiable Etiology. *J Infect Dis.* 2017;215:1407–
393 15.
 - 394 7. Xu L, Zhu Y, Ren L, Xu B, Liu C, Xie Z, et al. Characterization of the Nasopharyngeal
395 Viral Microbiome from Children with Community-Acquired Pneumonia but Negative for
396 Luminex xTAG Respiratory Viral Panel Assay Detection. *J Med Virol.* 2017
397 Dec;89(12):2098-2107.
 - 398 8. Lewandowska DW, Schreiber PW, Schuurmans MM, Ruehe B, Zagordi O, Bayard C, et
399 al. Metagenomic sequencing complements routine diagnostics in identifying viral
400 pathogens in lung transplant recipients with unknown etiology of respiratory infection.
401 *PloS One.* 2017;12:e0177340.
 - 402 9. Parize P, Muth E, Richaud C, Gratigny M, Pilmis B, Lamamy A, et al. Untargeted next-
403 generation sequencing-based first-line diagnosis of infection in immunocompromised
404 adults: a multicentre, blinded, prospective study. *Clin Microbiol Infect.* 2017;23:574.e1-
405 574.e6.
 - 406 10. Lewandowska DW, Zagordi O, Geissberger F-D, Kufner V, Schmutz S, Böni J, et al.
407 Optimization and validation of sample preparation for metagenomic sequencing of viruses
408 in clinical samples. *Microbiome.* 2017;5:94.
 - 409 11. Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, et al. Viruses in the
410 faecal microbiota of monozygotic twins and their mothers. *Nature.* 2010;466:334–8.
 - 411 12. Yang J, Yang F, Ren L, Xiong Z, Wu Z, Dong J, et al. Unbiased parallel detection of viral
412 pathogens in clinical samples by use of a metagenomic approach. *J Clin Microbiol.*
413 2011;49:3463–9.
 - 414 13. Kim K-H, Bae J-W. Amplification methods bias metagenomic libraries of uncultured
415 single-stranded and double-stranded DNA viruses. *Appl Environ Microbiol.*
416 2011;77:7663–8.
 - 417 14. Kozyreva VK, Truong C-L, Greninger AL, Crandall J, Mukhopadhyay R, Chaturvedi V.
418 Validation and Implementation of Clinical Laboratory Improvements Act-Compliant
419 Whole-Genome Sequencing in the Public Health Microbiology Laboratory. *J Clin*
420 *Microbiol.* 2017;55:2502–20.
 - 421 15. Simner PJ, Miller S, Carroll KC. Understanding the Promises and Hurdles of
422 Metagenomic Next-Generation Sequencing as a Diagnostic Tool for Infectious Diseases.
423 *Clin Infect Dis.* 2018 Feb 10;66(5):778-788.

- 424 16. Ruppé E, Schrenzel J. Messages from the second International Conference on Clinical
425 Metagenomics (ICCMg2). *Microbes Infect.* 2018 Apr;20(4):222-227.
- 426 17. Miller RR, Uyaguari-Diaz M, McCabe MN, Montoya V, Gardy JL, Parker S, et al.
427 Metagenomic Investigation of Plasma in Individuals with ME/CFS Highlights the
428 Importance of Technical Controls to Elucidate Contamination and Batch Effects. *PLoS*
429 *One.* 2016;11:e0165691.
- 430 18. Thoendel M, Jeraldo P, Greenwood-Quaintance KE, Yao J, Chia N, Hanssen AD, et al.
431 Impact of Contaminating DNA in Whole-Genome Amplification Kits Used for
432 Metagenomic Shotgun Sequencing for Infection Diagnosis. *J Clin Microbiol.*
433 2017;55:1789–801.
- 434 19. Gargis AS, Kalman L, Lubin IM. Assuring the Quality of Next-Generation Sequencing in
435 Clinical Microbiology and Public Health Laboratories. *J Clin Microbiol.* 2016;54:2857–
436 65.
- 437 20. Hall RJ, Wang J, Todd AK, Bissielo AB, Yen S, Strydom H, et al. Evaluation of rapid and
438 simple techniques for the enrichment of viruses prior to metagenomic virus discovery. *J*
439 *Virol Methods.* 2014;195:194–204.
- 440 21. Conceição-Neto N, Zeller M, Lefrère H, De Bruyn P, Beller L, Deboutte W, et al.
441 Modular approach to customise sample preparation procedures for viral metagenomics: a
442 reproducible protocol for virome analysis. *Sci Rep.* 2015;5:16532.
- 443 22. Petty TJ, Cordey S, Padioleau I, Docquier M, Turin L, Preynat-Seauve O, et al.
444 Comprehensive Human Virus Screening Using High-Throughput Sequencing with a
445 User-Friendly Representation of Bioinformatics Analysis: a Pilot Study. *J Clin Microbiol.*
446 2014;52:3351–61.
- 447 23. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying
448 mammalian transcriptomes by RNA-Seq. *Nat Methods.* 2008;5:621–8.
- 449 24. Dreier J, Störmer M, Kleesiek K. Use of Bacteriophage MS2 as an Internal Control in
450 Viral Reverse Transcription-PCR Assays. *J Clin Microbiol.* 2005;43:4551–7.
- 451 25. Schlager R, Chiu CY, Miller S, Procop GW, Weinstock G, Professional Practice
452 Committee and Committee on Laboratory Practices of the American Society for
453 Microbiology, et al. Validation of Metagenomic Next-Generation Sequencing Tests for
454 Universal Pathogen Detection. *Arch Pathol Lab Med.* 2017;141:776–86.
- 455 26. Zhou Y, Fernandez S, Yoon I-K, Simasathien S, Watanaveeradej V, Yang Y, et al.
456 Metagenomics Study of Viral Pathogens in Undiagnosed Respiratory Specimens and
457 Identification of Human Enteroviruses at a Thailand Hospital. *Am J Trop Med Hyg.*
458 2016;95:663–9.
- 459 27. Thorburn F, Bennett S, Modha S, Murdoch D, Gunson R, Murcia PR. The use of next
460 generation sequencing in the diagnosis and typing of respiratory infections. *J Clin Virol*
461 *Off Publ Pan Am Soc Clin Virol.* 2015;69:96–100.
- 462 28. Yang J, Yang F, Ren L, Xiong Z, Wu Z, Dong J, et al. Unbiased Parallel Detection of
463 Viral Pathogens in Clinical Samples by Use of a Metagenomic Approach[▽]. *J Clin*
464 *Microbiol.* 2011;49:3463–9.
- 465 29. Pichon M, Gaymard A, Josset L, Valette M, Millat G, Lina B, et al. Characterization of
466 oseltamivir-resistant influenza virus populations in immunosuppressed patients using
467 digital-droplet PCR: Comparison with qPCR and next generation sequencing analysis.
468 *Antiviral Res.* 2017;145:160–7.

469 **Table 1. Metagenomic NGS results for the validated respiratory samples tested positive with viral real-**
 470 **time PCR.**

Sample No.	Real-time PCR Ct values		Viral genome type	mNGS results for targeted viruses ^a			
				Identification	No. of reads	RPKM	Coverage(%)
1	HRV/EV	25	linear ssRNA	HRV-A19	13,061	5.5E+06	97.6
2		24		HRV-A19	29,743	8.2E+06	98.2
3		29		HRV-A63	2,672	1.4E+06	58.1
4		34		HRV-A56	453	1.4E+04	75.2
5	RSV	27		RSV-B	14,218	1.9E+06	91.2
6		36		RSV-A	187	1.5E+03	22.0
7	MPV	33		HMPV-A	44,556	9.1E+05	100.0
8	HCoV	20		HCoV NL63	73,878	2.4E+06	94.2
9		24		HCoV 229E	19,615	1.1E+06	99.8
10		28		HCoV 229E	20,666	2.4E+05	100.0
12		36		HCoV NL63	1,815	1.3E+04	9.6
13	MV	23		Measles Virus	289,019	9.1E+06	98.1
14	IBV	23		Influenza B	42,212	1.1E+06	97.2
15	IAV	27	fragmented ssRNA	Influenza A(H3N2)	24,234	1.9E+05	78.6
16		34		Influenza A(H3N2)	1,559	1.9E+04	21.2
17		35		Influenza A(H3N2)	258	1.8E+03	26.5
18	HBoV	24	linear ssDNA	HBoV-1	79,504	2.7E+06	100.0
20	AdV	17	linear dsDNA	HAdVC-1	245,2476	1.6E+07	99.8
21		36		HAdVD-51	18	8.0E+01	0.6
22 ^b		30		HAdVC-6	284	1.0E+03	6.2
	HHV-6	28		HHV-6B	18,411	1.4E+04	54.8
23 ^b	HBoV	15	linear ssDNA	HBoV-1	9,470,426	1.6E+08	100.0
	CMV	31	linear dsDNA	CMV	653	2.5E+02	5.3
24 ^b	HBoV	17	linear ssDNA	HBoV-1	7,966,089	1.1E+08	100
	MPV	29	linear ssRNA	HMPV-A	10,629	5.9E+04	95.7
25 ^{b, c}	AdV	26	linear dsDNA	HAdVC-2	2,165	6.8E+03	12.4
	HPIV	26	linear ssRNA	HPIV-3	17,576	1.3E+05	66.7
	HRV/EV	34		HRV-C	446	7.0E+03	9.2
	CMV	27	linear dsDNA	CMV	34,577	1.7E+04	24.8
26 ^b	HRV/EV	26	linear ssRNA	HRV-A78	114,684	1.4E+07	99.9
	AdV	30	linear dsDNA	HAdVC-2	65	1.6E+03	9.6
	RSV	30	linear ssRNA	RSV-A	586	3.5E+04	68.7
27 ^b	AdV	32	linear dsDNA	HAdVC-2	24	1.3E+02	3.2
	HPIV	37	linear ssRNA	HPIV-2	50	6.3E+02	2.3
28 ^b	HRV/EV	31		HRV-A71	1,309	3.5E+04	61.3
	EBV	23	linear dsDNA	EBV	2,556	3.0E+03	39.3

471 HRV: human rhinovirus, EV: enterovirus, RSV: respiratory syncytial virus, HCoV: human coronavirus, HMPV: human
 472 metapneumovirus, HPIV: human parainfluenza virus, MV: measles virus, HBoV: human bocavirus, AdV: adenovirus, HHV:
 473 human herpes virus, CMV: cytomegalovirus, EBV: Epstein-Baar virus, Ct: Cycle threshold, RPKM: reads per kilobase of
 474 virus reference sequence per million mapped reads (normalization of the number of reads mapping to a targeted viral
 475 genome).

476 ^aTargeted viruses: viruses detected with real-time PCR.

477 ^bMultiple viral infections.

478 ^cCross-checked on duplicate sample (deviation <0.5 log).

479 **Figure titles and legends**

480 **Fig. 1. Schematic representation of the metagenomic workflow and quality control steps.**

481 The whole process is summarized in the middle. On the left side, internal control (MS2 bacteriophage) is represented
482 in blue, and external controls are represented in red, including positive control (MS2 bacteriophage spiked in viral
483 transport medium) and No-Template Control (NTC: viral transport medium). Quality control testing 1 corresponds to
484 MS2 bacteriophage molecular detection with commercial PCR assay. Quality control testing 2 corresponds to control
485 by sequencing metrics (number of MS2 reads normalized with RPKM ratio and MS2 genome coverage). On the right,
486 each technique used by phases is indicated black. In addition, on the far right the duration of each step is indicated.

487 **Fig. 2. Determination of optimal spiking level for internal quality control.**

488 The sensitivity of the metagenomic analysis workflow for MS2 bacteriophage (Internal Quality Control, IQC)
489 detection was evaluated with a MS2 ten-fold serial dilutions in a nasopharyngeal swab tested negative with multiplex
490 viral PCR. Relative abundance of MS2 bacteriophage and viral families are represented depending on the MS2 spiked-
491 in concentration. IQCT1 corresponds to MS2 molecular detection with commercial real-time PCR assay after
492 amplification step. IQCT2 corresponds to control by sequencing metrics (number of MS2 reads normalized with
493 RPKM ratio and MS2 genome coverage).

494 **Fig. 3. Internal quality control detection after metagenomic analysis of the respiratory samples selected.**

495 Distribution of normalized read counts (RPKM) for MS2 bacteriophage (internal quality control, IQC) depending on
496 the proportion of viral reads for the 37 respiratory samples selected. MS2 RPKM was determined after subtracting of
497 NTC MS2 RPKM. IQC was not detected for 4/37 samples (represented in red); among them 3 samples were tested
498 positive with viral real-time PCR.

499 **Fig. 4. Evaluation of the metagenomic NGS workflow according to the viral genome type.**

500 a) Correlation between the results of metagenomic NGS and viral real-time PCR for validated respiratory samples
501 tested positive with viral PCR. Normalized number of reads (RPKM) obtained for targeted virus are displayed against
502 the real-time PCR Ct values for fragmented ssRNA virus (influenza virus) linear dsDNA virus (adenovirus, Epstein-
503 Baar virus, cytomegalovirus, human herpes virus-6) linear ssDNA (human bocavirus) and linear ssRNA (human
504 rhinovirus, respiratory syncytial virus, parainfluenza virus, human coronavirus, human metapneumovirus and measles
505 virus). The correlation coefficients are shown for each viral genome type. b) RPKM normalized by Ct for each viral
506 genome type of validated respiratory samples tested positive with viral PCR. Bars show median and interquartile
507 ranges, p-values calculated with the Mann-Whitney U test are shown.

508 **Titles and legends for additional files**

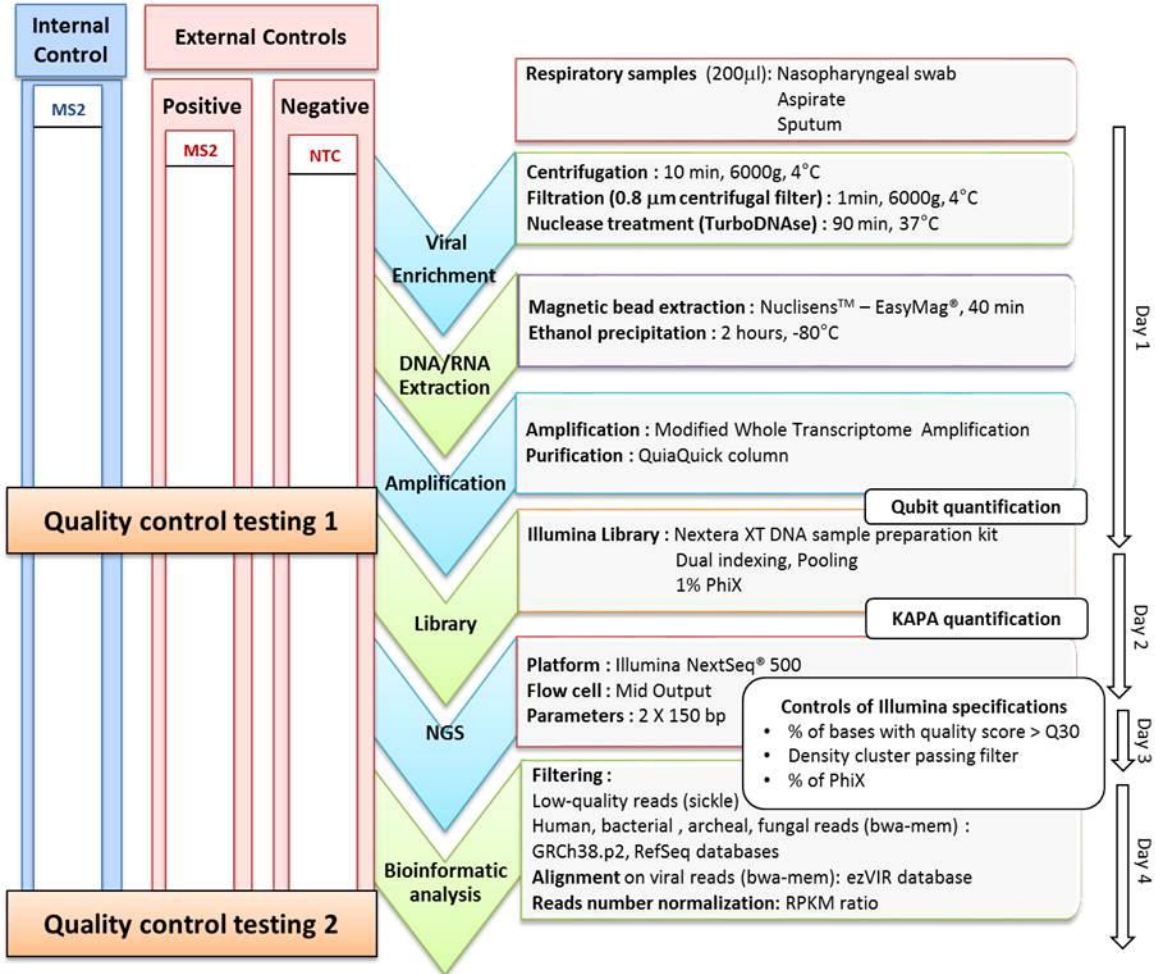
509 **Additional file 1. (Table.xls) Summary of clinical samples and metagenomic NGS information.**

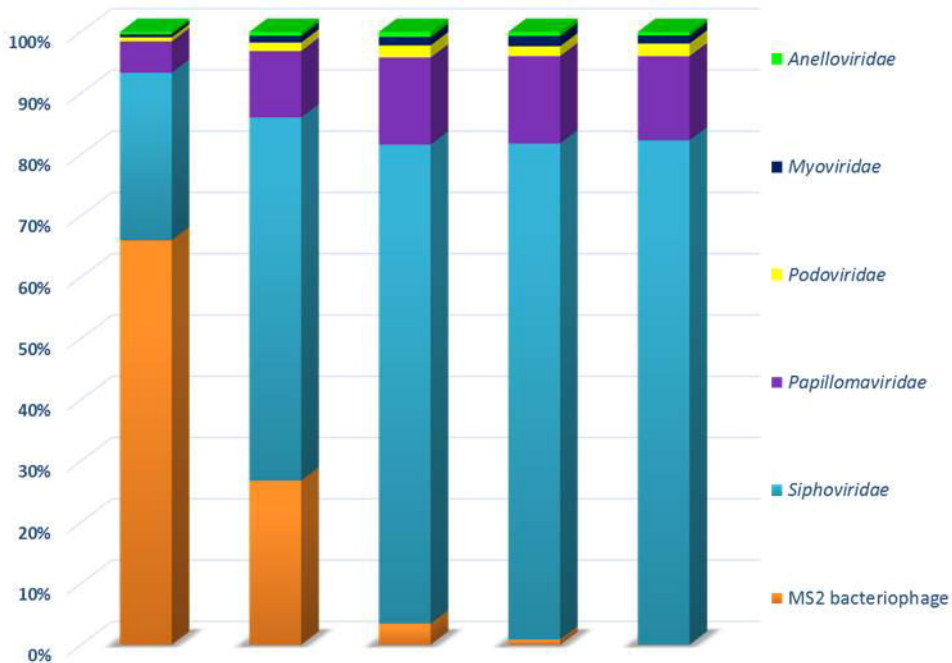
510 **Additional file 2. (Table.xls) Quality control testing results.**

511 QCT1 corresponds to MS2 bacteriophage molecular detection with commercial real-time PCR assay. QCT2
512 corresponds to control by sequencing metrics (number of MS2 reads normalized with RPKM ratio and MS2 genome
513 coverage). MS2 RPKM for the 37 selected clinical samples was determined after subtracting of NTC MS2 RPKM.

514 **Additional file 3. (Figure.ppt) Metagenomic NGS results for duplicates of sample # 25.**

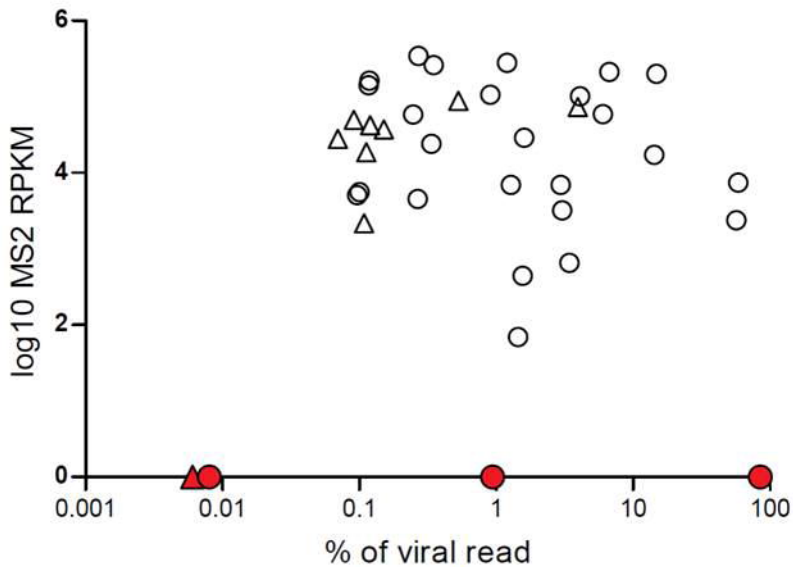
515 Sample # 25 corresponds to a clinical respiratory sample tested positive for 2 DNA viruses (adenovirus,
516 cytomegalovirus) and 2 RNA viruses (human parainfluenza virus, human rhinovirus) using real-time PCR. This
517 sample was analyzed twice using our single metagenomic workflow (replicate 1 and replicate 2). a) Pie charts show
518 classification of reads into human, bacteria, viruses, fungi, archea and unknown categories (unassigned reads). b)
519 Normalized read counts (RPKM) for each targeted virus (viruses detected with real-time PCR) and for internal quality
520 control (MS2 bacteriophage). c) Coverage plot of targeted viral genomes and internal quality control (MS2
521 bacteriophage). Sequencing reads were mapped on ezVIR viral database that identified human adenovirus C-2
522 (accession number: KF268130.1), cytomegalovirus (accession number: GQ396662.1), human parainfluenza virus 3
523 (accession number: KF687321.1), human rhinovirus C (accession number: JF317014.1) and MS2 bacteriophage
524 (accession number: NC_001417.2).



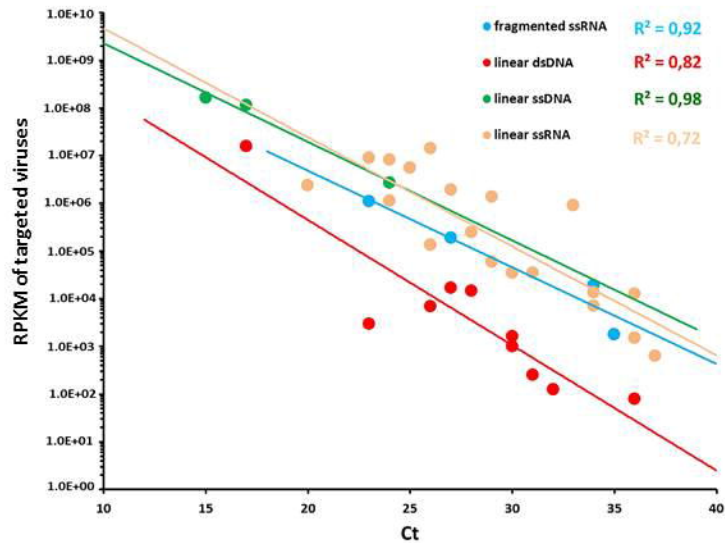
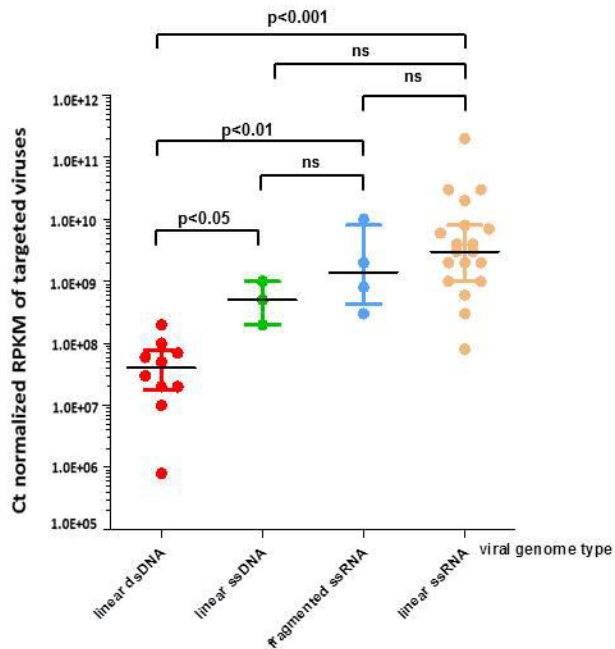


MS2 spiking level		10 ⁻²	10 ⁻³	10 ⁻⁴	10 ⁻⁵	Control
	%MS2	66.0	26.8	3.4	0.9	0
	IQCT1(Ct)	18	19	22	26	No ct
IQCT2	MS2 RPKM	8.9E+06	2.4E+06	9.7E+05	4.1E+05	0
	MS2 Coverage(%)	98.5	95.7	95.2	69.6	0

%MS2, ratio of number of reads mapping to the MS2 bacteriophage genome to the total number of viral reads; IQCT, internal quality control testing; MS2 RPKM: reads per kilobase of MS2 bacteriophage sequence per million mapped reads; Ct, Cycle threshold. Control, no spike in.



MS2 RPKM, reads per kilobase of MS2 bacteriophage sequence per million mapped reads; triangles indicate samples tested negative with viral real-time PCR, circles indicate samples tested positive; % of viral reads, ratio of number of reads mapping to viral genome to the total number of reads.

a**b**

RPKM, reads per kilobase of targeted viral sequence per million mapped reads; Targeted viruses: viruses detected with viral real-time PCR; Ct, Cycle threshold; ns, no significant.