1   # Cont-ID:

2   # Detection of samples cross-contamination in viral metagenomic

3   # data

4

5   Johan Rollin*[1,2], Wei Rong*[1] and Sébastien Massart[1]

6

7   *These authors contributed equally to this work

8

9   1. *University of Liège, Gembloux Agro-Bio Tech, Plant Pathology Laboratory, 5030, Gembloux, Belgium*

10  2. *DNAVision, 6041, Gosselies, Belgium*

11  **Background:** High Throughput sequencing (HTS) technologies completed by the bioinformatic analysis
12  of the generated data are becoming an important detection technique for virus diagnostics. They have
13  the potential to replace or complement the current PCR-based methods thanks to their improved
14  inclusivity and analytical sensitivity, as well as their overall good repeatability and reproducibility. Cross-
15  contamination is a well-known phenomenon in molecular diagnostics and corresponds to the exchange
16  of genetic material between samples. Cross-contamination management was a key drawback during the
17  development of PCR-based detection and is now adequately monitored in routine diagnostics. HTS
18  technologies are facing similar difficulties due to their very high analytical sensitivity. As a single viral
19  read could be detected in millions of sequencing reads, it is mandatory to fix a detection threshold that
20  will be influenced by cross-contamination. Cross-contamination monitoring should therefore be a
21  priority when detecting viruses by HTS technologies.

22  **Results:** We present Cont-ID, a bioinformatic tool designed to check for cross-contamination by
23  analysing the relative abundance of virus sequencing reads identified in sequence metagenomic datasets
24  and their duplication between samples. It can be applied when the samples in a sequencing batch have
25  been processed in parallel in the laboratory and with at least one external alien control. Using 273 real
26  datasets, including 68 virus species from different hosts (fruit tree, plant, human) and several library
27  preparation protocols (Ribodepleted total RNA, small RNA and double stranded RNA), we demonstrated
28  that Cont-ID classifies with high accuracy (91%) viral species detection into (true) infection or (cross)
29  contamination. This classification raises confidence in the detection and facilitates the downstream
30  interpretation and confirmation of the results by prioritising the virus detections that should be
31  confirmed.

32

33  **Conclusions:** Cross-contamination between samples when detecting viruses using HTS can be monitored
34  and highlighted by Cont-ID (provided an alien control is present). Cont-ID is based on a flexible
35  methodology relying on the output of bioinformatics analyses of the sequencing reads and considering the
36  contamination pattern specific to each batch of samples. The Cont-ID method is adaptable so that each
37  laboratory can optimise it before its validation and routine use.

38  **Keywords:**(1)

39  Bioinformatic, virus, detection, sequencing, contamination, metagenomic

40

# **Background**

The advent of high-throughput sequencing (HTS) technologies coupled with the development of powerful bioinformatics approaches has improved our ability to detect viruses in a non-targeted way from any sample collected from diverse sources. Noteworthy, detecting viruses by HTS technologies relies on many steps in the laboratory: sampling, transport and storage, nucleic acid extraction, library preparation, and sequencing (1). Compared to other molecular tests like (RT)-PCR, these steps are much more numerous and complex (2).

The analytical sensitivity, e.g. the ability to detect viral species at very low concentration in a sample, has been demonstrated to be similar to or even better than RT-PCR for animals (3) or plant viruses(4,5)). In addition, the inclusivity of HTS technologies, e.g. the ability to detect all isolates from a species and all species whose nucleic acids are present in enough quantity in a nucleic acid extract, is particularly high compared to any other detection test (2,6). Consequently, the use of HTS technologies is currently expanding at a rapid pace in research and is also progressively used for the diagnostic of viruses threatening humans (7), including SARS-Cov-2 (8), livestock (9) or plant health (10).

The broader application of HTS technologies for virus detection, with the simultaneous analysis of tens to hundreds of samples, is raising a significant challenge that needs to be addressed: the management of cross-contamination between samples. Scientists and diagnosticians already faced such challenges decades ago during the development of PCR-based techniques for detecting plants (11,12) or animal viruses (13,14), and this phenomenon might worsen with the use of HTS for virus detection (2). The higher complexity of laboratory operations, the intrinsically very high inclusivity, and the very low limit of detection (few viral reads are enough to detect the virus) of HTS make cross-contamination a more pressing issue. This is a frequently observed but, until recently, rarely reported observation in many, if not all, laboratories that have tested these technologies for virus detection. In many cases, these problems are frequently limited to a low number of reads and are of little consequence. Still, the specifics of the diagnostics field, with the need to detect viruses that can be at very low titre in the sample, clearly give more impact to such potential contamination problems (2). The occurrence of contamination is, therefore, a key element to consider when interpreting the viruses detected in HTS datasets.

The consequences of erroneous detection due to cross-contamination between samples can be catastrophic, as described for tuberculosis prior to HTS (15) but also using HTS for human and plant viruses (16,17).

So, even if the cross-contamination issue of HTS is long known and discussed in the scientific community, proper methodologies and dedicated algorithms are still missing to address it. Until now, the burden of detection confirmation relied on the virologist's expertise and the use of laboratory tests to independently confirm the presence of the virus in the sample, which is a fastidious, costly, and time-consuming task. To minimise the confirmation burden, arbitrary thresholds (like 5 or 10 reads) (4,18) have been proposed to consider a detection valid. Still, these thresholds are subjectively fixed based on the sequencing/detection tools or the scientist's experience. In addition, it has been shown recently that the cross-sample contamination burden can be very variable between sequencing batches and that an adaptative threshold is required(5). Therefore, the need for formal bioinformatic pipelines for HTS-based data that consider the possible cross-contamination is growing (19).

To handle cross-contamination, several laboratory protocol improvements have been implemented over time: laboratory or reagent decontamination, alternate dual indexes, inter-run washing (20,21) or, more recently, the use of alien control. An alien control is defined as "a matrix infected by a target (called alien target) which belongs to the same group as the target organism to be tested in the samples, but that cannot be present in the samples of interest." (22). It is processed as external control alongside the sample

88  to be analysed. It is preferably the same type of matrix as the analysed samples: plant tissue, water …
89  Ideally, the alien target, in our case a virus, should be at a high concentration in the alien sample as it
90  allows a better analysis of cross-contamination between samples. Indeed, the probability of detecting any
91  virus at a low level due to cross-contamination rises if this virus is very abundant in at least one of the
92  processed samples. A high abundance of the alien virus will therefore allow better monitoring of
93  contaminations, including for other viruses highly abundant in at least one tested sample. The presence of
94  sequencing reads from the alien virus in any tested sample can be considered the consequence of
95  contamination from the alien control to this sample. Such information can be used to monitor the cross-
96  contamination level between samples within the sequencing batch.
97  Many generalist bioinformatic tools, such as Kraken (23) or BLAST (24) can detect the presence of viruses
98  in HTS datasets with very high analytical sensitivity, as the detection is possible from a single viral read or
99  contig. Some of them, like VirHunter (25), VirAnnot (26) or VirusDetect (27), have been specifically
100 developed for that purpose. Nevertheless, they have not been designed to detect cross-contamination in
101 the input datasets. Instead, they will detect a virus, whatever its origin: virus infection in the biological
102 sample or contamination from another sample. The risk of contamination is particularly acute for viruses
103 in very high abundance in one of the samples sequenced as a few contaminating reads can be detected by
104 the bioinformatic tools in other samples prepared in parallel. The situation's impact is growing, especially
105 in the diagnostic field (2), as false positive results due to contamination can lead to inaccurate data
106 interpretation, which can cause tremendous health and trade issues.

107 According to EDAM ontology (28), tools that address cross-contamination issues should be labelled as
108 "Sequence contamination filtering". We were looking for tools using EDAM terms and the usual ones (virus
109 reads contamination, cross-contamination, …). Some tools address similar issues like contamination on
110 bacterial isolates (ConFindr- (29) or bacterial metagenome (GUNC - (30). They both use methods relying
111 on operons organisation of genes that are not applicable for viruses. Croco (31) uses an approach mainly
112 based on bacterial quantitative data. Finally, DecontaMiner (32) can be applied to metagenome data,
113 including viruses but is based on a combination of detection methods (mainly mapping and Blast) that try
114 to assign the dark matter (reads from unknown origin) more than formally detecting the cross-
115 contamination material. To our knowledge, there is no tool specifically addressing cross-contamination
116 during virus detection in metagenome datasets. It means that some risks of false positive results remain
117 unmonitored for virologists, and the burden of confirmation of detection in case of false positive is still not
118 addressed.
119 To solve this issue, we present Cont-ID, a method designed to check sample cross-contamination for
120 viruses previously identified in metagenomic datasets. It relies on a simple requirement: every sample in
121 a sequencing batch should have been processed at the same time and followed the same steps in the
122 laboratory with at least one alien control as external control. Cont-ID uses a voting system to classify every
123 species prediction on each sample of the sequencing batch into (true) infection or (cross)
124 contamination. This tool will help the virologist to distinguish virus presence and virus cross-contamination
125 in HTS data improving the reliability of viral detection and the efficiency of downstream confirmation and
126 characterisation analyses. It can also help to improve feedback on upstream steps that might be linked to
127 cross-contamination events. Cont-ID is an open-source python (v3) based script method freely available
128 here: https://github.com/johrollin/viral_contamination.
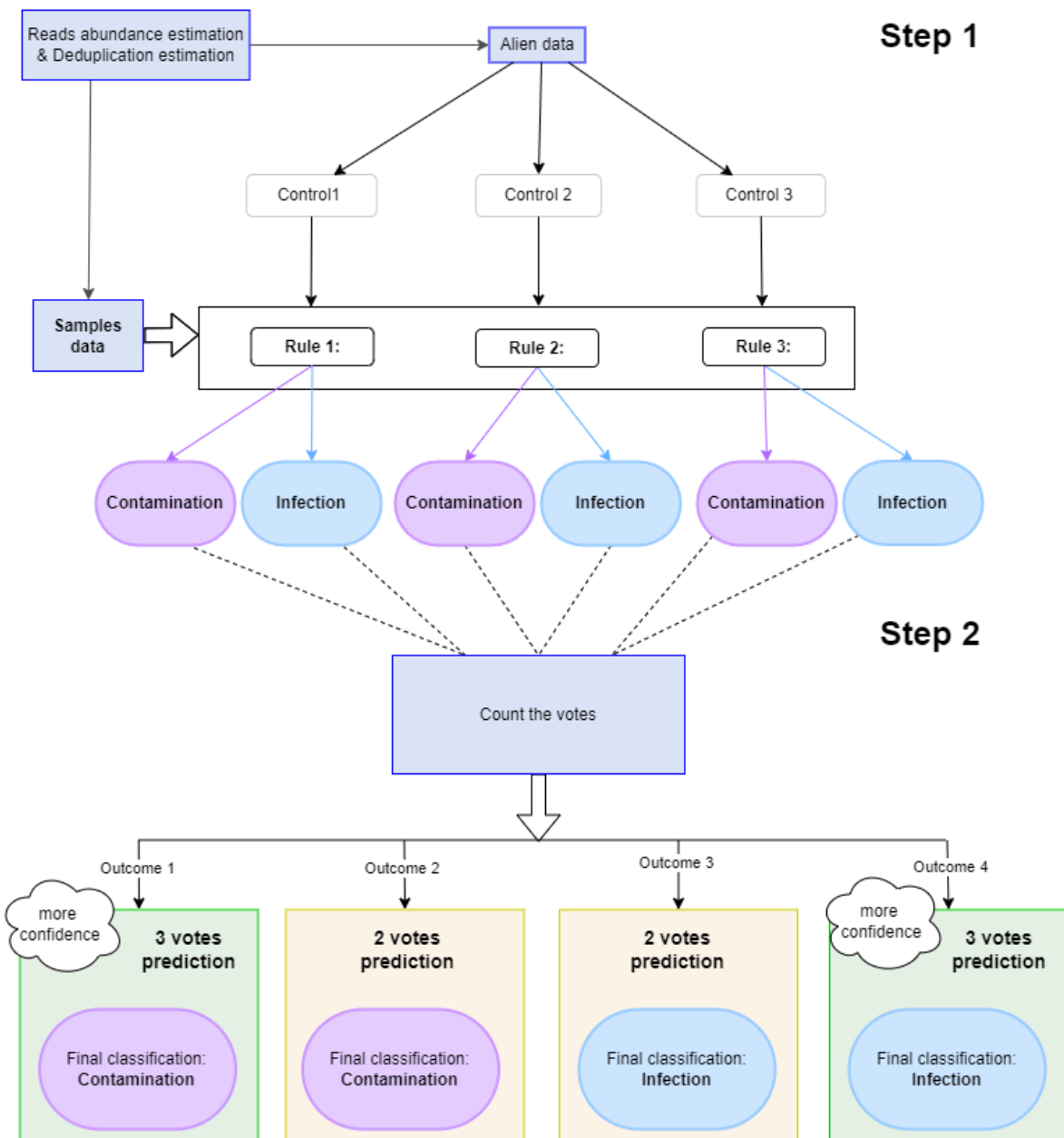129
130

# **<u>Methods:</u>**

132

## Implementation

In viral metagenomics, detecting multiple viral species in the same sample is frequent, and a virus species can be seen with different confidence levels in several samples of the same sequencing batch. Therefore, Cont-ID aims to determine whether a given detected virus in a sample is likely to be a contaminant or not by comparing it to the results from the other samples of the same sequencing batch, e.g. samples processed in parallel and following the same laboratory steps.

Cont-ID does not require any development or maintenance of database as it only relies on data generated by usual bioinformatics tools for HTS dataset analyses and, most specifically, two elements: (i) the normalised abundance estimation (number of reads assigned to each detected virus species on each sample) and (ii) the number of identical reads among pairs of samples (deduplication ratio). These input metrics are easy to obtain as the abundance estimation can be calculated by using any mapping tool like BWA (33) or a read classifier like Kraken/Bracken (23,34), and the number of identical reads from a virus between two samples can be obtained by running any deduplication tool like BBduck (35). A tabulated file containing these numbers associated with the detected virus name and the unique ID for each batch sample is used as input for Cont-ID, as shown in **Figure 2**. Each virus predicted on each batch sample is considered a distinct element and corresponds to a line in the generated table. A separate table is generated for the alien virus.

Computing the two elements mentioned above into three different metrics for the alien virus and each detected virus, Cont-ID can predict through three rules if a given viral species detection is likely a cross-contaminant or not in the sequencing batch, as described in **Figure 1**.
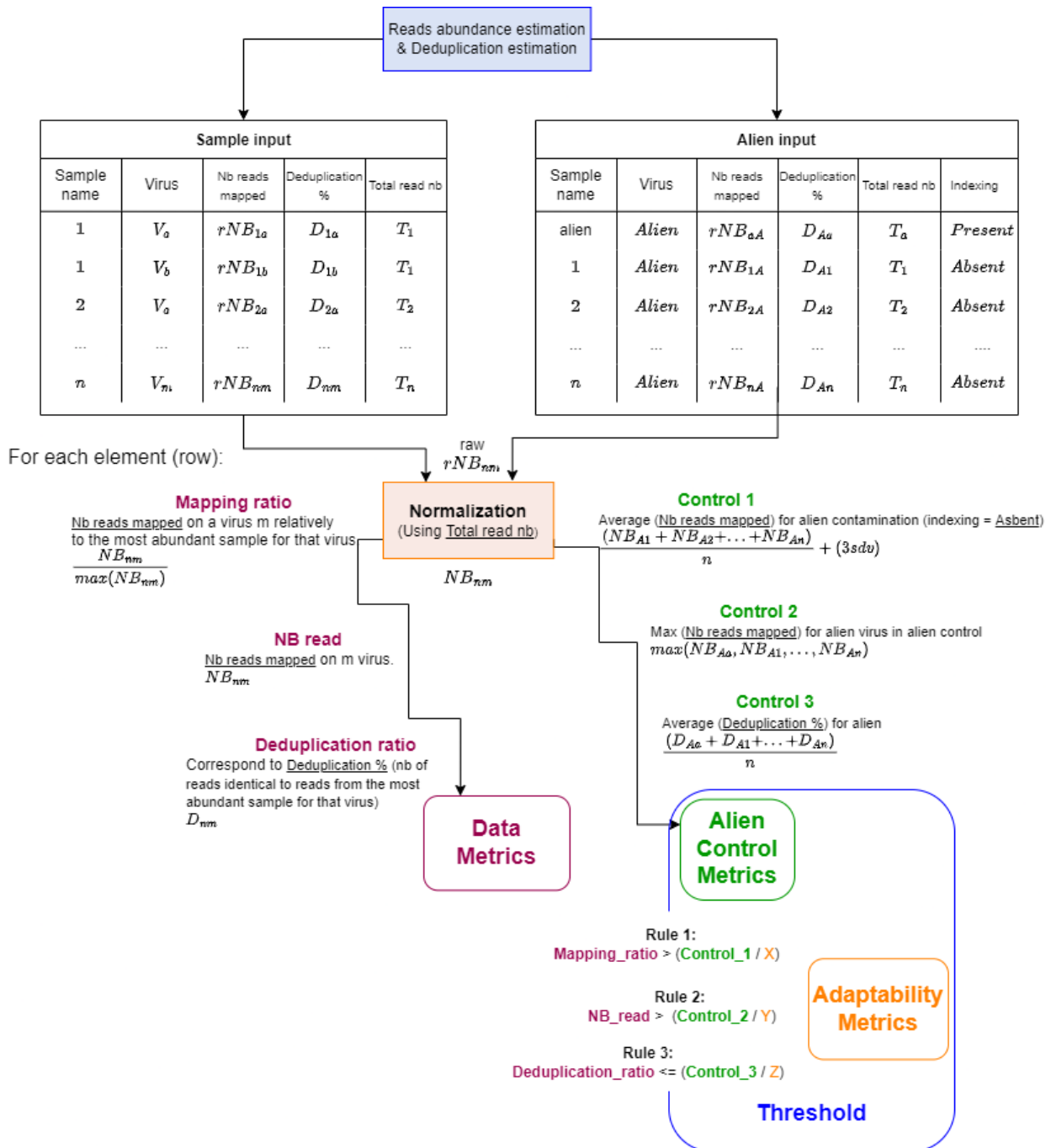
155

**Figure 1:** C**ross-contamination prediction with Cont-ID.**
*There are two input files, one for the alien data and one for the samples data. Alien file is used to calculate control thresholds which are then used along with the sample data to apply rules to a voting system (step 1). The votes are then counted to decide for each virus on each sample (element) either if it is a (cross)contamination or an infection (step 2).*

The three rules classify as contamination or infection each element according to the pattern of reads number observed among the samples and the alien control for the alien virus and the considered viral

163    species. Rules one and two both use the (normalised) reads abundance estimation, while rule three uses

164    the assessment of unique (identical) reads. Rules are calculated after normalising the number of reads

165    per sample and are described more precisely in **Figure 2**.

166



167
168
169    **Figure 2: Cont-ID rules explanation.**

170    *There are two input files, one for the alien data and one for the samples data. Alien file is used to calculate*

171    *each alien control metric after normalisation. The sample file is used to calculate each data metric after*

172    *normalisation. Each alien control metric is associated with a user (manually) designed adaptability metric*

173 *(X, Y or Z) to compose each rule's threshold. Finally, each Data Metric is compared to the corresponding*
174 *threshold in order to obtain the three rules used in Cont-ID.*
175
176 The first rule uses the mapping ratio of each virus in each sample (corresponding to an element): the
177 number of reads of each element is divided by the maximum number of reads of the corresponding virus
178 in one of the samples. This first rule compares this mapping ratio for the element with the Control 1 metrics
179 calculated for the alien virus and corresponds to the average number of reads mapped on the alien virus
180 in the samples for which the alien is a cross-contaminant, with three times the standard deviation of this
181 average.
182 The second rule relies on the number of reads of the element in the sample. The rule compares it with
183 Control 2, corresponding to the number of alien reads identified in the alien control. The third rule is based
184 on each element's deduplication ratio, which is compared with Control 3. The average deduplication ratio
185 of the alien virus reads between each tested sample and the alien control.
186
187 We aimed to find the most reliable formula for threshold calculation on each rule while allowing a part of
188 adaptability according to the biological system used. As the system variability can come from the
189 laboratory using HTS, the host and type of sample (fruit tree, herbaceous plant, human, animal …), the
190 type of virus (integrated or non) or the extraction protocol used (dsRNA, total RNA, small RNA...), each
191 rule includes a third number (represented by X, Y or Z) that is called adaptability metrics (see **Figure 2**).
192 The X will impact the first rule that considers the relative proportion of reads of a virus in this sample
193 compared to the sample with the maximum read of this virus. This threshold is a refinement of the "alien
194 threshold" described earlier(5). The default value proposed is 2. The Y divides the number of reads from
195 the alien virus in the alien control for comparing it to the number of reads of each virus in each sample. In
196 this publication, a default value of 1,000 has been fixed for Y, and it was in the range of the expected
197 (cross) contamination ratio (number of reads in the truly infected sample versus the number of reads in
198 contamination one). The Z metric impacts Control 3 and the evaluation of the proportion of identical reads
199 between different samples. The proportion of identical reads can be influenced by different factors
200 (mutation rate, respective genome length …). The role of Z is to consider those different factors. A default
201 value of 1.5 is proposed.
202 Default values of the three adaptability metrics have been provided in this publication after their
203 optimisation on the banana datasets and their evaluation of other datasets. Nevertheless, users can
204 independently modify them during the evaluation or validation of Cont-ID applied to their datasets. A
205 careful evaluation of the adaptability metrics by the user is recommended to evaluate their impact on the
206 diagnostic performance of the test. In addition, several sets of adaptability metrics can be run in parallel
207 for further improvements in diagnostics performance. The value given to the adaptability metrics and
208 controls resulting is always recorded in an additional log file (see Supplementary File 1 [log_file]). This log
209 file help to ensure traceability allowing the user to check the pertinence of the chosen numbers and to
210 adapt them when needed. As each of the three rules has two possible decisions (contamination or
211 infection), a majority vote will be obtained with two or three votes. The decision of each vote is available
212 in the generated result to support the result interpretation and let the user decides on the confidence to
213 give to each individual rule according to the biological system tested.
214  In addition, the proper quantitative comparison of sequencing reads datasets relies on normalising the
215 number of reads per sample, for example, as always done for transcriptomic or microbiome studies. This
216 assertion is also true for Cont-ID and corresponds to an adaptative parameter. To limit some bias due to
217 the difference in sequencing depth between samples in the same batch, we also normalise by default to 5
218 000 000 reads in this publication. Still, it is manually changeable by the user.
219 Finally, Cont-ID also has another level of flexibility: the script is made to ease the change of rules in the
220 code that can complete or replace the existing ones.

221

## Conditions of application for Cont-ID

223

Three main conditions are essential to run Cont-ID. First, an alien control should be used, alien control should contain a high concentration of the alien virus, so reads from that viral species are more prone to be detected in other samples when cross-contamination occurs from the alien control to the other samples. Similarly, if another virus is found in the alien control sample, that is also an indication of potential contamination (although not used so far by Cont-ID). The alien control is bioinformatically processed exactly as the samples of interest to generate the alien metrics for each sample (in a separate tabulated file). In the absence of external alien control, it is still possible to analyse sequencing batches if they include samples from different host species and some detected viruses, preferably at high abundance, are known to infect only some of the host species. In such a case, the alien file should be filled with the selected virus as if it was an alien (with the status of alien present/absent in the file). Nevertheless, the threshold set-up and the results will be less accurate and include fewer samples (the samples corresponding to species that can be hosts for the virus could not be considered). In addition, a high degree of confidence is needed regarding the actual infection of the sample selected as alien control by the virus selected as an alien virus. Cont-ID always requires at least one (cross) contamination in the alien file to be reported; otherwise, the threshold calculation will fail; in that case, the tool will state it.

239

The second application condition is related to the processing of the samples and the alien control. The alien control and all the other samples in a given batch should have been processed together in parallel for all the laboratory steps (RNA/DNA extraction, library preparation, sequencing) and bioinformatics (Reads cleaning, host removing …). This is a good diagnostic practice, but it is even more important here as the goal is to observe cross-contamination levels. The assumption is that the level seen with the alien represents what could have happened in samples of interest. Therefore, this assumption depends on processing all samples and control in parallel.

247

The third condition is that, once the user has fixed the adaptability metrics, the analysis should be carried out batch per batch. The calculation of sample and alien metrics is dynamically done for each batch as cross-contamination patterns can strongly vary between batches, as recently shown for banana samples (5).

252

253

## Sequencing reads datasets

255

The first datasets (batches A to D) were generated in our laboratory by total RNA sequencing protocol with ribodepletion applied to RNA extracted from banana plants (belonging to the *Musa* genus) (5). These data were generated to compare the test performance criteria of high throughput sequencing with classical virus testing protocols that include ImmunoCapture (IC)-(RT)-PCR and electron microscopy (36). The alien control corresponded to wheat plants infected by two species of barley yellow dwarf virus (BYDV-PAS and BYDV-PAV)(5). In total, four sequencing datasets (called A, B, C, and D) composed respectively of 27, 20, 27 and 25 samples were generated independently. A fifth batch generated during this validation experiment using diluted samples for evaluating the limit of detection (analytical sensitivity) was not included in our analysis according to the recent guidelines proposed for statistical analysis of validation datasets for plant pest detection (22). A total of 10 different viral species were infecting these samples, including banana mild mosaic virus (BanMMV), banana bract mosaic virus (BBrMV), banana bunchy top

267    virus (BBTV), cucumber mosaic virus (CMV), and five species belonging to the banana streak virus (BSV)
268    species complex. In addition, two other sequencing protocols were applied to some banana plants, small
269    RNA sequencing (5) starting from the same RNA extract as total RNA sequencing (for 21 samples in a single
270    batch) and double-stranded RNA (dsRNA) enrichment and sequencing protocol (37) applied from plant
271    tissue of 13 samples in a single sequencing batch.

272

| Batch_ID | NB of samples | Host type | Extraction Method | Extraction kit | Library kit | Sequencing | Data link | Publication |
|---|---|---|---|---|---|---|---|---|
| A (1) | 27 | Plant (Musa) | Total RNA extraction | RNeasy Plus Mini Kit (Qiagen, Venlo, Netherlands) | Stranded Total RNA library Prep Human/Mouse/Rat illumine CA, United States) & Ribo-Zero™ Plant Leaf Kit (illumine CA, United States) | Illumina NextSeq 500 2X150 | BioProject: PRJNA777477 samples starting with (1-XXX) | Wei et al., accepted. |
| B (2) | 20 | Plant (Musa) | Total RNA extraction | RNeasy Plus Mini Kit (Qiagen, Venlo, Netherlands) | Stranded Total RNA library Prep Human/Mouse/Rat illumine CA, United States) & Ribo-Zero™ Plant Leaf Kit (illumine CA, United States) | Illumina NextSeq 500 2X150 | BioProject: PRJNA777477 samples starting with (2-XXX) | Wei et al., accepted. |
| C (3) | 27 | Plant (Musa) | Total RNA extraction | RNeasy Plus Mini Kit (Qiagen, Venlo, Netherlands) | Stranded Total RNA library Prep Human/Mouse/Rat illumine CA, United States) & Ribo-Zero™ Plant Leaf Kit (illumine CA, United States) & TruSeq® Stranded Total RNA Library Prep Plant (illumine CA, United States) | Illumina NextSeq 500 2X150 | BioProject: PRJNA777477 samples starting with (3-XXX) | Wei et al., accepted. |
| D (5) | 25 | Plant (Musa) | Total RNA extraction | RNeasy Plus Mini Kit (Qiagen, Venlo, Netherlands) | TruSeq® Stranded Total RNA Library Prep Plant (illumine CA, United States) | Illumina NovaSeq 6000 2X150 | BioProject: PRJNA777477 samples starting with (5-XXX) | Wei et al., accepted. |
| E (6 sRNA) | 31 | Plant (Musa) | Total RNA extraction | RNeasy Plus Mini Kit (Qiagen, Venlo, Netherlands) | SMARTer smRNA-Seq Kit (Clonetech) | Illumina NovaSeq 6000 2X150 | BioProject: PRJNA777477 samples starting with (1sR-XXXX) | Wei et al., accepted. |
| F (5 dsRNA) | 9 | Plant (Musa) | Double stranded RNA | see article (Armelle Marais) | NEBNext Ultra II DNA library prep kit (New England BioLabs, US) | Illumina NovaSeq 6000 2X150 | BioProject: PRJNA777477 samples starting with (5ds-XXX) | Method: Marais et al., 2018 https://doi.org/10.1007/978-1-4939-7683-6_4 |
| G (Queensland university of technology) | 5 | Plant (mix) | Total nucleic acid | Maxwell™ Rapid Sample Concentrator instrument using SimplyRNA Tissue kit (AS1340, Promega) | TruSeq Stranded Total RNA | Illumina NovaSeq 6000 2X150 | BioProject: PRJNA752836 | Gauthier, M.-E. A.,et all https://doi.org/10.3390/BIOLOGY11020263 |
| H | 49 | Human | Total nucleic acid + amplification | NucliSENS EasyMAG platform (bioMérieux, Marcy l'Etoile, France) | Nextera XT (Illumina, San Diego, CA, USA) | Illumina NextSeq 500 2X150 | bioproject: PRJNA494633 | Bal et al., 2018 https://doi.org/10.1186/S12879-018-3446-5 |
| I | 25 | Human | Total nucleic acid | MagNAPure 96 DNA and Viral NA Small Volume Kit (Roche Diagnostics, Almere, the Netherlands) | EBNext Ultra Directional RNA Library Prep Kit for Illumina (New England Biolabs, Ipswich, MA) | Illumina HiSeq 4000 and NextSeq 500 depth: 10 million 2X150 | bioproject: PRJNA560243 | Boheemen et al., 2020 https://doi.org/10.1016/J.JMOLDX.2019.10.007 |
| J | 55 | Human | Total nucleic acid | TRIzol LS reagent (Invitrogen, USA) | SMARTer® Stranded Total RNA-Seq Kit v2 - Pico Input Mammalian (Takara Bio, USA) and the Trio RNA-Seq kit (NuGEN Technologies, USA) | Illumina HiSeq 2X150 | bioproject: PRJNA540900 | Li, et al., 2020 https://doi.org/10.1038/s41598-020-60992-6 |

273
274    **Table 1: list of datasets used on Cont-ID**
275

276    BSV is a species complex (genus: *Badnavirus*, family: *Caulimoviridae*) among which five species were
277    included in our samples: banana streak CA virus (BSCAV), banana Goldfinger virus (BSGFV), banana streak
278    IM virus (BSIMV), banana streak Mysore virus (BSMYV), and banana streak OL virus (BSOLV). Notably, some
279    species of this complex have their genome fully or partially integrated into the plant genome as
280    endogenous viral elements (EVE), most specifically in B genomes originating from *M. balbisiana*. These
281    EVE can be transcribed in the plant, and for some BSV species, they can even trigger an infection with viral
282    particles of BSV in the plant (38). It is well documented that BSGFV, BSIMV and BSOLV are constitutive of
283    *Musa balbisiana* (B genome) but can be activated in some conditions (39). In addition, BSMyV is also
284    integrated into the Musa B genome, although the ability to produce infectious viral particles is not yet
285    demonstrated. This brings additional complexity as EVE can be transcribed without the presence of a viral

286      particle. It has been recently demonstrated that the detection of BSV transcripts by HTS tests must be
287      confirmed by an independent test such as immunocapture (IC)-PCR for confirming the presence of viral
288      particles (5).

289      The other datasets used in this work came from publicly available datasets listed in **Table 1** and were
290      already included in peer-review publications. They were selected because they fit two criteria: (i) having
291      all virus presence checked in all the samples and (ii) having a virus species that could act like an alien
292      control for the input file. First, another set generated to detect viruses from diverse plant samples by high
293      throughput sequencing of total RNA extraction was kindly provided by Queensland University of
294      Technology  (17), corresponding to a total of 19 plant viruses and viroid in 5 samples. In addition, the
295      datasets generated from human samples came from published data from 3 different sources, with a total
296      of 129 samples containing 39 viral species (40–42). These three human datasets allowed us to test Cont-
297      ID with a large diversity of viruses, with different extraction and sequencing methods listed in
298      Supplementary File 2.
299
300      In total, ten sequencing batches, including 273 samples and the presence of 68 viral species, were used to
301      test the potential impact of a different host, extraction, and sequencing method on Cont-ID performances.
302      All the data generated are available with the link and procedure applied to obtain them described in **Table**
303      **1**; the indexing status of each virus in each sample is also available in Supplementary File 2.

304

## Bioinformatic analyses

306

### Quality control and mapping of sequencing reads

308

309      For all datasets, read quality control (quality trimming, reads deduplication) was performed using a
310      standard procedure described elsewhere (5). The cleaned reads were then mapped to a custom-built
311      database (DB) containing all complete genome sequences from previously detected viruses in the datasets.
312      For banana samples, all the complete genome sequences of the viruses were downloaded from NCBI nt
313      database on (12/12/2020) to serve as mapping DB. While the BYDV reference (KU170668 – for the alien
314      control) was selected as it was the closest sequence from our isolate. More information on the
315      composition of each mapping DB is available elsewhere (5).

316      The reads were mapped on the custom DB using Geneious mapper (Prime 2020.0.5, Biomatters). First, the
317      profile parameters "Low sensitivity / Fastest" were selected (with 20% mismatch and a maximum of 3
318      nucleotides gap allowed). To improve the results by aligning reads to each other in addition to the
319      reference sequence, the fine-tuning for mapping was set to "Iterate 2 times". The "multiple best matches"
320      option was set to "Randomly" (no multiple best matches between two different viruses were observed in
321      any sample processed). In the coming result section, we will refer to these parameters as "relax". A second
322      mapping referred to later on as "strict" was carried out using the same parameters except for the
323      mismatch allowance that was lower than 10%. Only the second mapping was carried out for small RNA
324      (20% mismatch is too much for small RNA). Indeed, using mismatches up to 20% should allow better
325      inclusivity of the analysis by mapping reads from isolates that can be genetically distant from the reference
326      sequences, especially if few reference genomes are available in the literature. Mapping with a strict
327      parameter was done to use small RNA and confirm this hypothesis. The tolerance of mismatches of 20%
328      is also close to many ICTV demarcation criteria to distinguish two different species (although these criteria

329  are often considered for only one or a few genes and might vary between families). Another test with
330  more relaxed parameters would increase the risk of adding non-specific reads (e.i. not generated from the
331  viral genomes) and was not considered.

332

## Deduplication of identical reads between samples

334

335  To investigate cross-contamination between samples, additional deduplication of identical reads between
336  samples was performed using dedupe V38.37 (from BBMap) embedded in Geneious (Prime 2020.0.5,
337  Biomatters) and with the parameters kmer seed length, maximum edit, and maximum substitutions set as
338  "31", "0", and "0", respectively. For each virus and sample, the mapped reads from each tested sample
339  and the sample with the highest number of mapped reads in the batch were grouped into a single pool
340  (using "Group sequences into a list" in Geneious) and deduplicated. The deduplication percentage equalled
341  the number of reads removed as duplicates divided by the lower number of reads between the two tested
342  samples. The deduplication percentage was not calculated on samples if less than 5 reads were mapped
343  to target viruses. For those samples, the rule (number three) automatically votes contamination. While for
344  the samples with the highest number of reads for a given virus, the deduplication ratio was set as reference
345  (i.e. "RF"), and the vote for rule three is infection.

## Confusion matrix and performance criteria calculation

347

348  We used a confusion matrix for each batch's results to have standard metrics for comparing batches and
349  samples. We compared the tool prediction for each element to the indexing status of the dataset
350  assimilating infection as a positive result and contamination as a negative result, as explained in **Table 2**.

351  **A**

| Cont-ID confusion matrix | | Prediction | |
|---|---|---|---|
| | | Infection (Positive) | Contamination (Negative) |
| **Indexing status** | Infection | True Positive (TP) | False Negative (FN) |
| | Contamination | False Positive (FP) | True Negative (TN) |

352  **B**

| | |
|---|---|
| **Diagnostic sensitivity (DSE)** | TP/(TP+FN) |
| **Diagnostic specificity (DSP)** | TN/(TN+FP) |
| **False omission rate (FOR)** | FN/(FN+TN) |
| **False discovery rate (FDR)** | FP/(FP+TP) |
| **Accuracy** | (TP+TN)/(TP+TN+FP+FN) |

353  **Table 2: (A) Confusion matrix based on Cont-ID results. (B) The formula is used to calculate tool**
354  **performance criteria.**

355

356  Based on the confusion matrix, we have four possibilities after a prediction: False Positive (FP) when the
357  tool wrongly predicted an infection, True Positive (TP) when the tool correctly predicted an infection,
358  False Negative (FN) when the tool wrongly predicted contamination and True Negative (TN) when the

359  tool correctly predicted contamination. In addition, we calculated several performance criteria
360  commonly used in diagnostics to evaluate our tool. To calculate those performance criteria
361  automatically, we used an automated script available on the same GitHub
362  (https://github.com/johrollin/Cont_ID/tree/master/further_analysis).

363

# Results

365

366  We used Cont-ID on ten metagenomic datasets, including a total of 273 samples, as a proof of concept
367  (see details in method). These datasets covered a broad range of use for Cont-ID as they were generated
368  from plant or human samples according to three library preparation protocols (total RNA, small RNA and
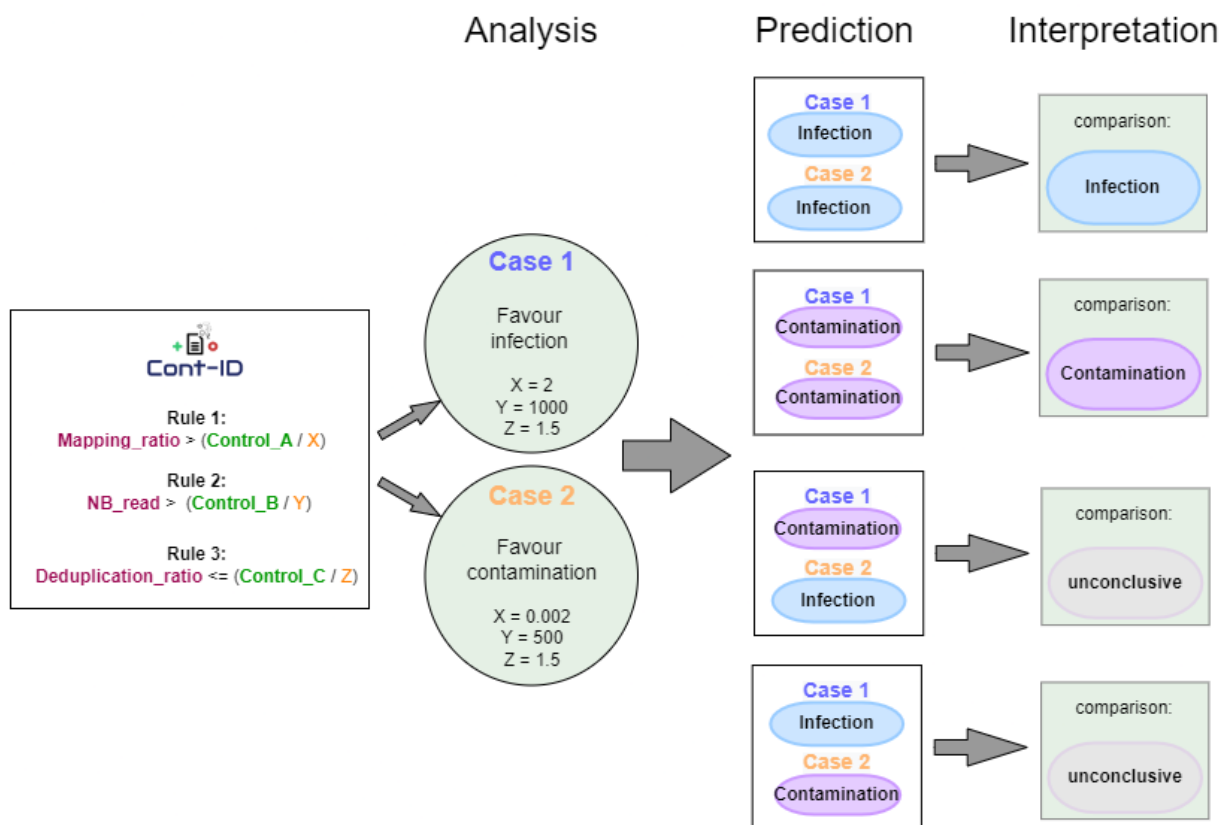369  double stranded RNA).

370

## Set up adaptability metrics datasets on the banana datasets

372

373  When applying for the first time Cont-ID on banana datasets generated from reference samples with
374  known viral status, the first objective was to determine the most appropriate values for the adaptability
375  metrics (X, Y and Z), allowing to minimise both FP (over-prediction of infection) and FN (over-prediction of
376  contamination). This was particularly complex as raising the value of an adaptability metric could lead to
377  an over-prediction of either contamination or infection by the rule, while lowering it had the opposite
378  effect.

379



380
381  **Figure 3: Cont-ID prediction when using the two default cases**

382
383 During the set-up of the method, we looked for the most adapted set of values to balance our rule
384 prediction on *Musa* datasets A, B and C. We tested several ranges of values aiming at limiting both wrong
385 predictions (FP and FN). The optimised single set of values maintaining FP and FN low in the three datasets
386 was not found. Indeed, variability was observed between batches, as any set limiting FP and FN in one or
387 several batches was not optimal for the other batch(es).
388 Indeed, the uneven proportion and pattern of cross-contaminations observed in different sequencing
389 batches made it very difficult to decide on a unique set of values. Instead, it seemed more efficient to
390 apply two different sets of values (called "case 1" and "case 2" further on) that favoured the prediction of
391 either true infection (TP - case 1) or true contamination (TN - case 2) from the same datasets. The
392 combination of the prediction from both cases would give additional information for interpretation. We
393 proposed values that gave the best performance criteria on our training datasets on bananas, and the
394 purpose of the diagnostic test was to minimise the risk of false negatives (priority 1) while keeping the
395 confirmation burden manageable (priority 2). Importantly, those sets of values can be manually adapted
396 by the user to improve one or several performance criteria of the test, to better fit the purpose of the HTS
397 tests carried out and its associated risks (risk of false positive or false negative) or to limit the "grey zone"
398 of inconclusive results (see under).
399 Therefore, we propose to run Cont-ID with two sets of adaptability metrics every time to compare the
400 results. Therefore, a high level of confidence is reached for the elements with identical predictions
401 between both cases. The combination can also highlight elements for which the prediction changed; they
402 correspond to the "grey zone" with metrics of abundance and/or duplication close to thresholds. In such
403 cases, the automated prediction might not be accurate. At this stage, it is mandatory to carry on additional
404 verification, such as checking the confidence (2 or 3 votes) for each prediction or comparing the threshold
405 numbers (also provided in the result) with the sample metrics. Cont-ID provides the list of votes for each
406 rule in each case to facilitate this additional verification. Then according to the additional information and
407 the test's purpose, the user can decide on the status (infection or contamination) or keep it inconclusive
408 but decide to test the virus presence independently by another test. For the presentation of the result,
409 the result is mentioned as "inconclusive" when both cases disagree.
410

## Evaluation of the method accuracy on banana samples

412
413 Based on the results obtained with the two sets of adaptability metrics, the tool predictions were
414 compared with the biological status of each reference banana sample (batches A, B, C and D
415 Supplementary File 2), allowing us to predict the cross-contamination on the four tested batches with an
416 average accuracy of 90%, excluding 23% of elements classified as "inconclusive" (see table 3A).
417

| | | | A - Relax mapping | | | | B - Strict mapping | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Batch | | A | B | C | D | A | B | C | D |
| | Sequencing method | | TotalRNA | | | | TotalRNA | | | |
| | Total element tested | | 105 | 93 | 143 | 128 | 76 | 65 | 93 | 68 |
| | Expected Infection/Contamination | | 28/77 | 14/79 | 34/109 | 19/109 | 28/48 | 14/51 | 34/59 | 17/51 |
| | Case 1 | 3 votes accuracy | 91% | 87% | 86% | 90% | 100% | 92% | 94% | 87% |
| | | 2 votes accuracy | 60% | 60% | 45% | 80% | 69% | 62% | 46% | 69% |
| | | overall accuracy | 73% | 67% | 66% | 85% | 78% | 68% | 62% | 75% |
| | Case 2 | 3 votes accuracy | 99% | 95% | 100% | 96% | 95% | 92% | 100% | 94% |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | **2 votes accuracy** | 58% | 69% | 87% | 67% | 70% | 59% | 64% | 44% |
| | | **overall accuracy** | 85% | 85% | 93% | 91% | 83% | 78% | 80% | 82% |
| | | **accuracy %** | 88% | 82% | 94% | 95% | 82% | 78% | 81% | 90% |
| | | **inconclusive %** | 23% | 20% | 33% | 16% | 5% | 17% | 32% | 28% |
| | **Cases combination** | **correct prediction** | 71 | 61 | 90 | 103 | 59 | 42 | 51 | 44 |
| | | **wrong prediction** | 10 | 13 | 6 | 5 | 13 | 12 | 12 | 5 |
| | | **inconclusive (occurrence)** | 24 | 19 | 47 | 20 | 4 | 11 | 30 | 19 |

**Table 3:** *Percentage of the accuracy of case 1 and case 2 analysed alone or in combination on banana samples sequenced by ribodepleted totalRNA sequencing. Each case is presented with the proportion of correct or wrong predictions according to the number of votes obtained (2 or 3). The percentage is given by three votes confidence count only the result with three votes while the overall accuracy aggregates the 2 and 3 vote results. When combining results from both cases, the percentage of inconclusive results and the number of correct or wrong predictions are stated. Two different mapping parameters were tested, allowing respectively 20% of mismatches (part A) or 10% of mismatches (part B)*

The predictions with the three votes using default mapping parameters ("relax mapping") are very trustworthy as the accuracy is higher than 86% and 95% for cases 1 and 2, respectively. These promising results are obtained on the fraction of the elements representing 25-50% and 48-84% for case 1 and case 2, respectively. The remaining elements are classified with two votes (more information is available in Supplementary File 3). The prediction accuracy with two votes is much lower, whatever the case. So, knowing the number of votes obtained by each element is crucial when the results need to be interpreted (and this number is always given in the report generated by Cont-ID). For case 1, most elements were predicted with two votes meaning that one of the (three) rules had the opposite prediction, which might explain why the accuracy was lower. While for case 2, the majority of the elements were predicted with three votes. The explanation is probably in the "expected Infection/contamination" row in **Table 3A**: for all batches, there is more contamination than infection (from 28 infections for 77 contaminations – batch A to only 19 infections for 109 contaminations - batch D). As stated above in the text and **Figure 3**, Case 2 is designed to favour contamination detection at the expense of infections occurring at a low concentration that tend to be considered contamination (FN). Nevertheless, as a direct consequence, true contamination (TN) detection is high (see confusion_matrix in Supplementary File 3).

Overall, case 2 presented a higher accuracy (85-91% relax mapping) than case 1 (66-85%), while the combination of the two cases reached a similar one (82-95 % relax mapping). Those good results from combination accuracy mean that very few predictions are wrong (5 – 13) in both cases, but 16-33% of elements are not counted in the accuracy percentage because they are inconclusive. The combination's importance relies on maintaining a high accuracy while highlighting the inconclusive prediction to prioritise them for manual expertise.

## The mapping parameters impacted the input files and the Cont-ID performance

In **Table 3**, we explored the impact on the prediction of two levels of mismatch tolerance (20% and 10%) when mapping the sequencing reads on the viral genome DB. The goal was to explore if changing a

454 parameter from the primary bioinformatics step delivering the input files of Cont-ID could have an impact
455 on the prediction. Strict mapping tends to lower the total number of elements tested due to a decrease in
456 the number of samples for which we have very few reads mapped to a candidate virus. Cont-ID has more
457 samples to process with a relaxed mapping, which should be better for threshold calculation. Logically, the
458 elements lost by the strict mapping parameter should be predicted as "contamination" and present a
459 relatively low number of reads. Indeed, those elements are most likely more distant reads (between 20
460 and 10 % mismatch with the reference genome) mapped on the virus. They could correspond to non-viral
461 reads wrongly mapped in datasets from samples tested negative by classical indexing. For example, for
462 batch B, on the 23 differential mapping results, the number of reads mapped ranged from 1 to 10. Of those
463 23 elements, 21 are classified as "contamination" the two remaining are labelled inconclusive. In batch C,
464 there are 50 differential elements between the two parameters, with 37 correct (13 from BSV), 11
465 inconclusive (10 from BSV) and 2 wrong (2 from BSV) classified elements. In total, 25 elements (on the 50
466 - batch C) are from the non-integrated virus, of which 24 are labelled contamination (one inconclusive).
467 The separation between integrated and non-integrated viruses is explained in **Table 4** and another
468 publication (5).
469 Using the relaxed mapping parameter seems beneficial for prediction as the accuracy is better (82-95%
470 relax, 78-90% strict). Moreover, thanks to the combination strategy, we can focus on the proportion of
471 inconclusive; it is uneven with an important increase, 5% (strict) to 23% (relax) for batch A, while in batch
472 D, it decreases from 28% to 16 %. However, when we look closely at the accuracy improvement, most
473 comes from differential elements (present only with relax) that are 'obvious' contamination with few
474 reads. So, most of the accuracy improvement did not come from very informative elements, except in
475 some rare occurrences where it helped classify well elements in relax parameters that were inconclusive
476 with strict or classified inconclusive elements in relax that were wrong with strict parameters. As an
477 example, in batch C, on the 24 elements for BanMMV, BBRMV, BBTV and CMV common in both conditions
478 (relax and strict), elements prediction is improved (from inconclusive [strict] to correct [relax]) for three of
479 them (sample 3B1, 3B2 and 3B14 with BanMMV).
480 There is, therefore, a slight improvement with relaxed mapping parameters, and we set these parameters
481 by default to generate the input files. Indeed, with the relaxed parameters, the number of reads for each
482 element (including alien) increases along the rise of the number of elements in the batch. This means that
483 we change the rule's threshold (see Figure 2), which is critical for the threshold calculation in a way that
484 seems more representative of reality than strict mapping. In these batches, some element metrics are very
485 close to the threshold used for the rules and slightly changing those metrics or the alien metrics (the alien
486 control metrics are obviously changed by the mapping parameters) can modify the prediction.
487 As we did not know the divergence of the virus genomes between different samples and the reference
488 genomes, it seemed more logical to use relaxed mapping parameters by default. According to the virus
489 system the user is working on and the ICTV demarcation criteria that go with it, these parameters should
490 or could be adapted.
491

## The virus biology can impact Cont-ID performance: the case of virus integration in the host genome

494

| | A- Non-integrated Virus | | | | B- Integrated Virus (BSV) | | | |
|---|---|---|---|---|---|---|---|---|
| **Batch** | **A** | **B** | **C** | **D** | **A** | **B** | **C** | **D** |
| **Sequencing method** | TotalRNA | | | | TotalRNA | | | |

| | | | 40 | 31 | 51 | 55 | 65 | 62 | 92 | 73 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Relax mapping** | **Total element tested** | | 40 | 31 | 51 | 55 | 65 | 62 | 92 | 73 |
| | **Expected Infection/Contamination** | | 19/21 | 9/22 | 22/29 | 10/45 | 9/56 | 5/57 | 12/80 | 9/64 |
| | **Case 1** | **3 votes accuracy** | 100% | 91% | 96% | 100% | 82% | 83% | 68% | 54% |
| | | **2 votes accuracy** | 94% | 70% | 75% | 100% | 47% | 56% | 43% | 78% |
| | | **overall accuracy** | 98% | 77% | 94% | 100% | 58% | 61% | 50% | 74% |
| | **Case 2** | **3 votes accuracy** | 100% | 85% | 100% | 100% | 97% | 100% | 100% | 93% |
| | | **2 votes accuracy** | 50% | 91% | 82% | 67% | 62% | 58% | 89% | 67% |
| | | **overall accuracy** | 88% | 87% | 92% | 96% | 83% | 84% | 93% | 88% |
| | **Cases combination** | **accuracy %** | 100% | 88% | 96% | 100% | 79% | 79% | 92% | 91% |
| | | **inconclusive %** | 15% | 16% | 6% | 4% | 28% | 23% | 48% | 25% |
| | | **correct prediction** | 34 | 23 | 46 | 53 | 37 | 38 | 44 | 50 |
| | | **wrong prediction** | 0 | 3 | 2 | 0 | 10 | 10 | 4 | 5 |
| | | **inconclusive (occurrence)** | 6 | 5 | 3 | 2 | 18 | 14 | 44 | 18 |

**Table 4**: *Percentage of the accuracy of case 1 and case 2 analysed alone or in combination on banana samples sequenced by ribodepleted totalRNA sequencing with relaxed mapping parameters. Each case is presented with the proportion of correct or wrong predictions according to the number of votes obtained (2 or 3). The percentage given by three votes confidence count only the result with three votes while the overall accuracy aggregates the 2 and 3 vote results. When combining results from both cases, the percentage of inconclusive results and the number of correct or wrong predictions are stated. Two types of viruses were tested, Non-integrated virus (part A) or integrated virus (part B).*

To highlight the potential impact of the virus biology on the results of Cont-ID, the analysis of banana batches was split between integrated and non-integrated viruses. Indeed, several species of BSV are integrated into its host genome, which complicates the reliable detection of BSV infection from sequencing datasets. Consequently, it has been recommended to confirm any detection of BSV reads by an independent PCR test combined with immunocapture of viral particles (5).

Table 4 shows better accuracy and a lower proportion of inconclusive results for non-integrated viruses compared to BSV. More elements with contamination status are obtained when looking for BSV than non-integrated viral species. This over-representation of contaminants might be caused by the transcription of integrated sequences of BSV even without viral particles, which will raise the number of detected reads. These are two points that reduced the efficiency of our method on BSV and, by extension, might also concern any other viral species integrated into the host genome and able to produce transcripts.

The global accuracy is lower for BSV species (79-92%) compared to the other viruses (88-100%), even if the maximum accuracy obtained with batch C (92%) was high. In addition, the proportion of inconclusive results should also be considered, and this proportion was much higher for BSV (23-45%) than for the other viruses (4-16%). So, the overall performance of Cont-ID is lower when applied on BSV and did not solve the issues of appropriate detection in sequencing data of viral infection from viruses integrated into the plant genome. Consequently, BSGFV, BSIMV, BSMYV, and BSOLV, which correspond to different but closely related species of Banana streak virus (BSV) integrated into the *Musa* genome, were excluded from the calculation of performance criteria for the banana datasets. BSCAV was also excluded (despite not being integrated) because of its similarity with other BSV species.

## Performance of Cont-ID on diverse datasets

The performance of Cont-ID using the two cases was further evaluated while diversifying the hosts (fruit trees, grasses, humans) and the sequencing protocols (total RNA, small RNA, dsRNA).

| Batch | | A | B | C | D | E | F | G | H | I | J | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Origin** | | Banana (own sequencing) | | | | | | Plant Mix (Gauthier et al., 2022) | Human (Bal et al 2018) | Human (Boheemen et al 2019) | Human (Li et al 2020) | |
| **Sequencing method** | | TotalRNA | | | | SmallRNA | dsRNA | TotalRNA | | | | **Average** |
| **Total element tested** | | 40 | 31 | 51 | 55 | 20 | 12 | 51 | 112 | 62 | 206 | 64 |
| **Expected Infection/Contamination** | | 19/21 | 9/22 | 22/29 | 10/45 | 19/1 | 5/7 | 18/33 | 37/75 | 25/37 | 50/156 | |
| **Case 1** | **3 votes accuracy** | 100% | 91% | 96% | 100% | 100% | 100% | 100% | 96% | 100% | 90% | 97% |
| | **2 votes accuracy** | 94% | 70% | 75% | 100% | 36% | 67% | 86% | 71% | 54% | 63% | 72% |
| | **overall accuracy** | 98% | 77% | 94% | 100% | 55% | 92% | 92% | 87% | 73% | 77% | 84% |
| **Case 2** | **3 votes accuracy** | 100% | 85% | 100% | 100% | 33% | 100% | 97% | 92% | 92% | 95% | 89% |
| | **2 votes accuracy** | 50% | 91% | 82% | 67% | 63% | 50% | 92% | 68% | 84% | 78% | 72% |
| | **overall accuracy** | 88% | 87% | 92% | 96% | 45% | 92% | 96% | 86% | 89% | 87% | 86% |
| **Cases combination** | **accuracy** | 100% | 88% | 96% | 100% | 50% | 100% | 98% | 93% | 93% | 92% | 91% |
| | **inconclusive %** | 15% | 16% | 6% | 4% | 20% | 17% | 8% | 15% | 29% | 24% | 15% |
| | **correct prediction** | 34 | 23 | 46 | 53 | 8 | 10 | 46 | 88 | 41 | 144 | 49,3 |
| | **wrong prediction** | 0 | 3 | 2 | 0 | 8 | 0 | 1 | 7 | 3 | 12 | 3,6 |
| | **inconclusive (occurrence)** | 6 | 5 | 3 | 2 | 4 | 2 | 4 | 17 | 18 | 50 | 14,4 |

**Table 5:**

*Percentage of the accuracy of case 1 and case 2 analysed alone or in combination from sequencing with relaxed mapping parameters (except for small RNA). Each case is presented with the proportion of correct or wrong predictions according to the number of votes obtained (2 or 3). The percentage given by three votes confidence count only the result with three votes while the overall accuracy aggregates the 2 and 3 vote results. When combining results from both cases, the percentage of inconclusive results and the number of correct or wrong predictions are stated. Several virus datasets were tested, banana samples (only non-BSV viruses are considered), a mix of plants, and human datasets.*

**Table 5** shows the method's accuracy on all datasets with relaxed mapping parameters (except for small RNA, see method). Overall, the accuracy of Cont-ID was 94%, with 15% of inconclusive results. The sRNA dataset provides a poor accuracy (50%) with 20% inconclusive; this can be explained by the (almost) absence of contamination (Expected Infection/Contamination 19/1) by the low level of reads found (see Supplementary Files 2 & 3 for more information). Apart from small RNA, the worst accuracy (88%) has been obtained from the batch B sequencing dataset of banana. Noteworthy, this protocol was independently evaluated for virus testing in banana, but its performance for virus detection was much lower than total RNA sequencing (5). The accuracy calculated from the single batch of dsRNA, with only 9 samples and 12 elements, was 100%. Even if not enough representative dataset was used for dsRNA, the method accuracy seems not too far from what we obtained in Total RNA, indicating that, Cont-ID is independent of the extraction method. On Total RNA, for banana samples, the accuracy ranged from 88% to 100%, with 4% to 15% of inconclusive results. The accuracy of the plant mix (G) was also very high (98%), with 8% of inconclusive results. On human datasets, the accuracy remained high (92-93%), but the

553 inconclusive results reached up to 15 - 29%. Overall, the application of Cont-ID on human datasets reached
554 similar performance in accuracy; the slightly worse inconclusive metrics can be explained by the fact that
555 the adaptability metrics might not be the best ones for the human dataset and underlined the importance
556 given at Cont-ID for the flexible adaptation of metrics and parameters.

557 For most of the datasets, case 1 performed worse than case 2, probably due to the design of the case
558 metrics (see **Figure 2**), where case 1 values were determined to favour infection. The expected
559 infection/contamination ratio showed that for all the datasets but E (small RNA), there was a lot more
560 contamination than infection; therefore, case 1 overpredicted infection, lowering its accuracy. In the E
561 dataset, case 1 (55%) performed better than case 2 (45%) as expected; it is also the case for the human
562 dataset H (97% case 1, 96% case 2), even if the ratio (37/75) leans toward contamination.

563

564 Those results indicated that Cont-ID performed well in classifying cross-contamination in very different
565 virus-host systems, even if some adjustments may be needed in some cases in the future. The different
566 levels of flexibility of Cont-ID made such adjustments possible. To provide an example of analysis, all the
567 information regarding batch C from the input file to the analysis file (including raw results) is available in
568 Supplementary File 3.

569

## 570 **Discussion**

571

572 Despite significant efforts to limit cross-contamination (dual indexes, inter-run washing …), this still
573 represent a concern and the appropriate distinction between low-level infection, and cross-sample
574 contamination is crucial for the large-scale development of HTS technologies as a diagnostic test.
575 Furthermore, it should be adequately managed because identifying and monitoring the cross-
576 contaminations improves the detection results' reliability. In other words, it can help to find the source of
577 contamination in the laboratory, take appropriate measures to minimise it, and raise confidence in the
578 detected viruses.

579 This publication improved a preliminary work on determining an adaptative contamination threshold for
580 the detection of plant viruses (5), which uses the maximal number of alien virus reads contaminating a
581 sample as the threshold of detection for each sequencing batch. So, instead of using a fixed number for
582 the contamination threshold as done in the literature, the threshold is adapted to the level of
583 contamination monitored in the batch thanks to the alien control. The former publication used a single
584 threshold corresponding to the maximum number of alien virus reads in a sample. Some limitations of this
585 previous threshold, for example, overestimating contamination when viral reads are in low number for a
586 virus, underlined the need for improvements. This was achieved with Cont-ID through the definition of
587 multiple formal rules, the automation of calculation and the ability to adapt the thresholds and rules by
588 the user. The tool's prediction relies on basic and usual information generated by bioinformatics analysis
589 of sequencing data (mapping and duplication numbers) and the use of external alien control. The criteria
590 based on reads (relative) abundance of each virus in each sample and the (approximation of) number of
591 identical reads for a virus between samples performed well while being relatively easy to generate. Our
592 objective with this tool was to show that exploring data generated by standard bioinformatic procedures
593 can facilitate the identification of cross-contamination between samples.

594 Cont-ID discriminated virus infection and cross-contamination between samples with a global accuracy of
595 91 % (median=95%) on the diverse range of datasets included in its evaluation. The diversity of situations
596 included viral species belonging to diverse viral families with cellular hosts belonging to plant or animal
597 kingdoms and three different library preparation protocols. Importantly, the default values of adaptability
598 metrics determined from banana dataset predicted cross-contamination with high accuracy (96%, on
599 banana excluding small RNA) and remain high even on human datasets (94%). To further help the user in

600  the analysis, we provide the detailed votes prediction in the result file (see Supplementary File 3). This is,
601  therefore, a solid basis for the diagnostician to check the level of confidence in the generated results.
602  Indeed, each prediction made by the method uses at least two rules to determine the classification of the
603  element for each case. A prediction with three votes is more confident than with two votes. But all
604  predictions with two votes do not provide the same confidence as it depends on which rules predicted
605  what. Of our three rules, two rely more or less directly on abundance estimation, which means that when
606  that metric is not obtainable in a reliable way, the tools' predictions will be impacted, and predictions with
607  those rules might be less confident. On the other hand, rule three (deduplication ratio) is less effective
608  when the read numbers are low. Depending on the scenario, the user should consider the relative
609  confidence of each rule when trying to confirm Cont-ID prediction. This underlines again the importance
610  of proper interpretation of the obtained results based on the virus biology

611

612  The prediction quality depends on the input data quality, meaning that the deduplication and mapping
613  parameters are essential and should be carefully considered while evaluating their impact on the results.
614  For example, some deduplication tools remove reads if a (small) read is contained in another (larger) read;
615  having that option active or not will significantly impact the deduplication ratio. As shown in the results,
616  mismatch parameters are very impactful for the mapping. Considerations like ICTV demarcation criteria
617  or what parameters the biologist would use to reconstruct the whole viral genome are helpful in deciding
618  the ones to use for Cont-ID input. In that regard, testing and expertise in bioinformatics analysis are heavily
619  beneficial. Here, the 20% mismatch parameter performed well; it might be different in other datasets (viral
620  composition) configurations or when working with databases containing many reference genomes.
621  Indeed, independently of the mismatch parameters used, using more genome references for each
622  expected species could also improve the ability to detect sequences from distant isolates by better
623  covering the genetic diversity of the virus.
624  The biology of the virus should also be considered, as shown by the results obtained with viruses with
625  functionally integrated genomes in the host, like BSV species. Our conclusion is that they should be
626  considered independently from the non-integrated viruses. It was challenging to extract a reliable metric
627  for BSV as the differentiation between reads from integrated genomes and reads from viral particles is
628  impossible. Indeed, the biology of viruses integrated into the host genome differs from non-integrated
629  viruses, as viral genome transcription can happen without viral particle production. We have not tested
630  our method on species with different biological behaviour like viroids or phages. But optimisation of the
631  adaptative metrics might likely be required in order the use Cont-ID with high accuracy. Viroid genomes
632  are generally smaller than viruses, while the phage genome tends to be much larger and has specific
633  biological features. For example, a different level of identical reads and abundance (calculation based on
634  reads number) could be obtained between the different scales of genome size.

635  For these reasons, Cont-ID allows the evaluation of other values for adaptability metrics (X, Y, Z) by each
636  user to adapt the tool and optimise its diagnostic performance depending on the biological matrix, the
637  protocol and the purpose of the test. Independently, the user can also adapt the metrics to reach the
638  appropriate balance between FN and FP by deciding if, for the purpose of the test and the available
639  resource for confirming detection, it is preferable to be overpredicting contamination to be confident that
640  all the virus detection remaining are true infection or the opposite (overpredicting infection to be sure not
641  to miss any).

642  In our tests, the analysis of the wrong predictions showed that none of the proposed rules (and
643  adaptability metrics values) allowed us to reach satisfactory accuracy with a proper balance between FN
644  and FP (see Supplementary File 1). We have observed that using two sets of adaptability metrics (one to
645  favour contamination and the other, infection prediction) gave a higher accuracy. In a real scenario (with

646  infection status not known for the samples), it is difficult to know if HTS virus detection (at low
647  concentration) is in the majority due to true infection or cross-contamination. The two-case strategy
648  allows the biologist to predict both scenarios with at least one case accurately. Indeed,  if the expected
649  ratio of infection/contamination is unknown, the relative performance of cases 1 or 2 will be unknown,
650  so it seems preferable to use the combination results instead of the individual.

651  If both cases agree, the assumption is that the prediction is correct. Nevertheless, combining the results
652  will provide a list of interesting inconclusive results. Each inconclusive result means that the two cases
653  delivered opposite predictions. Therefore, the scientist should address those results when analysing
654  Cont-ID prediction by checking the number of rules for each prediction, for example, knowing if 2 or 3
655  rules agreed and checking the results of each rule :  How close to the threshold was the read
656  abundance/ratio and/or the duplication rate? Spotting the few errors that may occur requires excellent
657  manual expertise as the usual manual verification methods may also indicate the wrong decision (if there
658  are many reads from cross-contamination, the mapping results can be wrongly positive while the and/or
659  the (RT)-PCR  can also be wrongly positive if the contamination occurred at an early stage and the (RT)-
660  PCR was carried out on the same nucleic acids extract). Other information about the virus-plant
661  interaction should be considered, like virus-species-cultivar compatibility or geographical virus
662  distribution (see investigation on unexpected viruses (5)).

663
664  Cont-ID also presents some limitations that need to be discussed. First, the number of identical reads
665  estimation comes from the deduplication procedure, which is an approximation, and that can be a problem
666  because it can consider the non-specific reads (reads that are not coming from cross-contamination but
667  that are identical to another sample from a common area of the genome) as identical to the probable
668  source of contamination by mistake. Indeed, this can be the case if, for example, two samples are infected
669  by the same virus isolate at very different concentrations. The presence of duplicated reads might suggest
670  contamination instead of a low-level infection. The risk of such an extreme situation is limited using two
671  other rules, although interpreting the data will require good expertise in virus genomic variability and
672  detailed information on the sample origins and virus prevalence and diversity.
673  In addition, the duplication metric assumes that contamination (if any) comes from the sample with the
674  highest number of reads. This theory seems logical since the more reads in a sample, the higher the
675  probability of detecting a few reads from it contaminating other samples (potential of contamination).
676  Nevertheless, it can create a bias when a virus is highly abundant in two (or more) samples and detected
677  with a low frequency in others. In that case, it is difficult to determine the true origin of cross-
678  contamination. Such a case could be a fundamental limit of our current method. If several samples with a
679  very high abundance of reads are present in a batch, as developed here, Cont-ID should be applied as many
680  times as the number of highly abundant samples. Ideally, Cont-ID should include the read duplication
681  comparison of each sample to all other samples for a virus, but this can raise additional issues (like
682  contamination from several origins at the same time), and, at this stage, it was not implemented.
683  We must also keep in mind that the relative quantity of genetic material between samples might change
684  because the biologist normalises the quantity of DNA/RNA at two steps of the process: before starting
685  library preparation and during the pooling of the prepared libraries. Meaning that the differential in
686  genomic material concentration (potential of contamination of a sample) is resettled. If cross-
687  contamination happens before that step, it can cause less accurate predictions from Cont-ID. This bias in
688  the estimation of abundance is another limitation of our method.
689  Using an (alien) control helps to know the expected level of contamination but is also impacted by the limit
690  of detection inherent to the standard bioinformatic procedures. Indeed, working with very few reads for
691  some viruses makes some analyses impossible when below their detection limit. For example, the

692 calculation of the duplication rate below a minimal number of reads (in this study, we chose 5) of a virus
693 did not make sense. The limit of calculation of the input metrics is another limitation of Cont-ID.
694
695 Cont-ID accuracy was high, but additional improvements can probably be explored. For example, by
696 exploiting the ability of other metrics generated during bio-informatic analyses (like RKPM, genome
697 coverage percentage, relative coverage depth repartition, …) to help detect contamination. In fact, some
698 of these metrics with several thresholds were tested for Cont-ID before selecting the three rules described
699 in **Figure 2** that provided the highest accuracy (in both contamination and infection determination).
700 Importantly, values leading to a perfect scenario were not identified, and a two-cases classification system
701 was set up (more information in Supplementary File 1).
702 Nevertheless, adding more metrics will also complexify the decision system. If more metrics are considered
703 for cross-contamination prediction, other implementations (decision tree, machine learning …) might be
704 envisioned to replace the current voting system. On the other hand, the detection in the alien control of
705 sequencing reads of other viruses detected in the tested samples is also the consequence of contamination
706 from one of the tested samples toward the alien control. This information is not used now but could also
707 be considered for future improvements as it requires less complex to implement. In addition, it might allow
708 refinement of Cont-ID, potentially introducing an adaptation of threshold per virus instead of a single
709 threshold for all samples from the sequencing batch. The idea is that two viruses present in the same batch
710 may have different relative abundance behaviour in the samples, so setting up a limit that can adapt for
711 each virus should improve the tool's ability to distinguish real infection from cross-contamination. Finally,
712 working with the combination of all/some viruses profile instead of each individually for contamination
713 check (similarly to what is used in metabarcoding of bacteria) can also be considered. Indeed, when a
714 sample contaminates another, it is expected that all the viruses (highly frequent) from the contaminating
715 sample can be found in the contaminated samples. Monitoring the virus detection profile of samples can
716 provide additional information for cross-contamination (and ease the quest for contamination origin).
717 Even if there is still improvement to be made, Cont-ID has already delivered an excellent ability to consider
718 the level of contamination genuinely present in the batch.
719
720 In conclusion, detection of cross-contamination is complex; in the age of sequencing, the contaminant
721 issue is increasingly important; therefore, Cont-ID will facilitate the interpretation of results by the
722 virologist/diagnostician and reduces the confirmation burden. We demonstrated that simple metrics like
723 relative abundance estimation and redundancies of genetic material (reads duplicates) could help monitor
724 contamination occurring in the laboratory. The method accurately distinguished cross-contamination from
725 infection in very diverse HTS viral datasets. Our standard parameters allowed very good accuracy (median
726 = 95%); in addition, Cont-ID has several levels of flexibility and can be adapted by each user to take into
727 account the specificities of the detection test (purpose of the test, type of samples, viruses to be detected,
728 laboratory work, available resources….). We believe this is the first significant step toward increasing the
729 monitoring and management of sample cross-contamination when using HTS technologies for virus
730 detection.
731
732

## Availability and requirements:

734
735 Project name: Cont-ID
736
737 Project home page: https://github.com/johrollin/Cont_ID
738

739    Operating system(s): Platform independent
740
741    Programming language: Python (v3.7)
742
743    Other requirements: pandas; NumPy
744
745    License: GNU GPL-3.0
746
747    Any restrictions to use by non-academics: none
748

# Declarations

## Ethics approval and consent to participate
751    Not applicable.
752
## Consent for publication
754    Not applicable.

755
## Availability of data and materials
757    The datasets generated and/or analysed during the current study are available in the NCBI Sequence Read
758    Archive (SRA) repository at https://www.ncbi.nlm.nih.gov/sra (See Table 1). In addition, cont-ID is freely
759    available   and   can   be   downloaded   with   the   following   command   (without   <>):
760    <https://github.com/johrollin/Cont_ID>. It can be used as a command line application on a personal
761    computer on any operating system (Linux, MacOSX or Windows) with python.
762
## Competing interests
764    The authors declare that they have no competing interests.

765
771
## Authors' contributions
773    WR and SM designed the data preparation and sequencing procedure. WR, JR and SM designed the
774    bioinformatics analysis. JR implemented the program. JR ran CroCo analyses. WR and JR validated Cont-ID
775    results with previous PCR indexing results and re-analysed outlier. JR and SM drafted the manuscript. All
776    authors read and approved the final manuscript.
777

783 Emilie Gauthier and Roberto Barrero for providing dataset G and discussing cross-contamination in viral
784 metagenomes with us.

785

## **<u>Reference</u>**

787

788 1. Lebas B, Adams I, Rwahnih M al, Baeyen S, Bilodeau GJ, Blouin AG, et al. Facilitating the adoption
789 of high-throughput sequencing technologies as a plant pest diagnostic test in laboratories: A step-
790 by-step description. EPPO Bulletin [Internet]. 2022 Aug 1;52(2):394–418. Available from:
791 https://onlinelibrary.wiley.com/doi/full/10.1111/epp.12863

792 2. Massart S, Olmos A, Jijakli H, Candresse T. Current impact and future directions of high
793 throughput sequencing in plant virus diagnostics. Virus Res. 2014 Aug 8;188:90–6.

794 3. Charlebois RL, Sathiamoorthy S, Logvinoff C, Gisonni-Lex L, Mallet L, Ng SHS. Sensitivity and
795 breadth of detection of high-throughput sequencing for adventitious virus detection. npj Vaccines
796 2020 5:1 [Internet]. 2020 Jul 17;5(1):1–8. Available from:
797 https://www.nature.com/articles/s41541-020-0207-4

798 4. Soltani N, Stevens KA, Klaassen V, Hwang MS, Golino DA, al Rwahnih M. Quality Assessment and
799 Validation of High-Throughput Sequencing for Grapevine Virus Diagnostics. Vol. 13, Viruses .
800 2021.

801 5. Rong W, Rollin J, Hanafi M, Roux N, Massart S. Validation of high throughput sequencing as virus
802 indexing test for Musa germplasm: performance criteria evaluation and contamination
803 monitoring using an alien control. PhytoFrontiers. 2022;

804 6. Maree HJ, Fox A, al Rwahnih M, Boonham N, Candresse T. Application of hts for routine plant
805 virus diagnostics: state of the art and challenges. Front Plant Sci [Internet]. 2018 Aug [cited 2019
806 Nov 15];9. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6119710/

807 7. Ng SH, Braxton C, Eloit M, Feng SF, Fragnoud R, Mallet L, et al. Current perspectives on high-
808 throughput sequencing (HTS) for adventitious virus detection: Upstream sample processing and
809 library preparation [Internet]. Vol. 10, Viruses. 2018. Available from:
810 https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6213814/

811 8. Kumar R, Nagpal S, Kaushik S, Mendiratta S. COVID-19 diagnostic approaches: different roads to
812 the same destination. VirusDisease 2020 31:2 [Internet]. 2020 Jun 13 [cited 2021 Oct
813 20];31(2):97–105. Available from: https://link.springer.com/article/10.1007/s13337-020-00599-7

814 9. Vereecke N, Carnet F, Pronost S, Vanschandevijl K, Theuns S, Nauwynck H. Genome Sequences of
815 Equine Herpesvirus 1 Strains from a European Outbreak of Neurological Disorders Linked to a
816 Horse Gathering in Valencia, Spain, in 2021 [Internet]. Vol. 10, Microbiology Resource
817 Announcements. American Society for Microbiology; 2021 [cited 2021 Oct 20]. Available from:
818 https://doi.org/

819 10. Olmos A, Boonham N, Candresse T, Gentit P, Giovani B, Kutnjak D, et al. High-throughput
820 sequencing technologies for plant pest diagnosis: challenges and opportunities. EPPO Bulletin.
821 2018 Aug;48(2):219–24.

822   11.   Lau HY, Botella JR. Advanced DNA-based point-of-care diagnostic methods for plant diseases
823          detection. Front Plant Sci. 2017 Dec 6;8:2016.

824   12.   Grosdidier M, Aguayo J, Marçais B, Ioos R. Detection of plant pathogens using real-time PCR: how
825          reliable are late Ct values? Plant Pathol [Internet]. 2017 Apr 1 [cited 2022 Jul 11];66(3):359–67.
826          Available from: https://onlinelibrary.wiley.com/doi/full/10.1111/ppa.12591

827   13.   Moonen P, Boonstra J, Hakze- Van Der Honing R, Boonstra- Leendertse C, Jacobs L, Dekker A.
828          Validation of a LightCycler-based reverse transcription polymerase chain reaction for the
829          detection of foot-and-mouth disease virus. J Virol Methods. 2003 Oct 1;113(1):35–41.

830   14.   Watzinger F, Ebner K, Lion T. Detection and monitoring of virus infections by real-time PCR. Mol
831          Aspects Med. 2006 Apr 1;27(2–3):254–98.

832   15.   Martínez M, de Viedma DG, Alonso M, Andrés S, Bouza E, Cabezas T, et al. Impact of Laboratory
833          Cross-Contamination on Molecular Epidemiology Studies of Tuberculosis. J Clin Microbiol
834          [Internet]. 2006 Aug [cited 2021 Oct 26];44(8):2967–9. Available from:
835          /pmc/articles/PMC1594630/

836   16.   Bukowska-Ośko I, Perlejewski K, Nakamura S, Motooka D, Stokowy T, Kosińska J, et al. sensitivity
837          of next-generation sequencing metagenomic analysis for detection of RNA and DNA viruses in
838          cerebrospinal fluid: The confounding effect of background contamination. Adv Exp Med Biol
839          [Internet]. 2017 Nov 1 [cited 2022 Jul 11];944:53–62. Available from:
840          https://link.springer.com/chapter/10.1007/5584_2016_42

841   17.   Gauthier MEA, Lelwala R v, Elliott CE, Windell C, Fiorito S, Dinsdale A, et al. Side-by-Side
842          Comparison of Post-Entry Quarantine and High Throughput Sequencing Methods for Virus and
843          Viroid Diagnosis. Biology 2022, Vol 11, Page 263 [Internet]. 2022 Feb 8 [cited 2022 Feb
844          14];11(2):263. Available from: https://www.mdpi.com/2079-7737/11/2/263

845   18.   Bloom JS, Sathe L, Munugala C, Jones EM, Gasperini M, Lubock NB, et al. Swab-Seq: A high-
846          throughput platform for massively scaled up SARS-CoV-2 testing. medRxiv. 2021
847          Mar;2020.08.04.20167874.

848   19.   Ballenghien M, Faivre N, Galtier N. Patterns of cross-contamination in a multispecies population
849          genomic project: detection, quantification, impact, and solutions. BMC Biol [Internet]. 2017 Mar
850          29 [cited 2021 Oct 26];15(1). Available from: /pmc/articles/PMC5370491/

851   20.   Costello M, Fleharty M, Abreu J, Farjoun Y, Ferriera S, Holmes L, et al. Characterisation and
852          remediation of sample index swaps by non-redundant dual indexing on massively parallel
853          sequencing platforms. BMC Genomics [Internet]. 2018 May 8;19(1):1–10. Available from:
854          https://bmcgenomics.biomedcentral.com/articles/10.1186/s12864-018-4703-0

855   21.   Champlot S, Berthelot C, Pruvost M, Andrew Bennett E, Grange T, Geigl EM. An Efficient
856          Multistrategy DNA Decontamination Procedure of PCR Reagents for Hypersensitive PCR
857          Applications. PLoS One [Internet]. 2010;5(9):e13042. Available from:
858          https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0013042

859    22.    Massart S, Lebas B, Chabirand A, Chappé AM, Dreo T, Faggioli F, et al. Guidelines for improving
860           statistical analyses of validation datasets for plant pest diagnostic tests. EPPO Bulletin [Internet].
861           2022 Aug 1;52(2):419–33. Available from:
862           https://onlinelibrary.wiley.com/doi/full/10.1111/epp.12862

863    23.    Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome Biol
864           [Internet]. 2019 Sep [cited 2020 Jan 14];20(1):762302. Available from:
865           https://www.biorxiv.org/content/10.1101/762302v1

866    24.    Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture
867           and applications. BMC Bioinformatics [Internet]. 2009 Dec 15 [cited 2022 Jul 12];10. Available
868           from: https://pubmed.ncbi.nlm.nih.gov/20003500/

869    25.    Sukhorukov G, Khalili M, Gascuel O, Candresse T, Marais-Colombel A, Nikolski M. VirHunter: A
870           Deep Learning-Based Method for Detection of Novel RNA Viruses in Plant Sequencing Data.
871           Frontiers in Bioinformatics. 2022 May 13;2.

872    26.    Lefebvre M, Theil S, Ma Y, Candresse T. The VirAnnot Pipeline: A Resource for Automated Viral
873           Diversity Estimation and Operational Taxonomy Units Assignation for Virome Sequencing Data.
874           https://doi.org/101094/PBIOMES-07-19-0037-A [Internet]. 2019 Oct 21 [cited 2021 Oct
875           25];3(4):256–9. Available from: https://apsjournals.apsnet.org/doi/abs/10.1094/PBIOMES-07-19-
876           0037-A

877    27.    Zheng Y, Gao S, Padmanabhan C, Li R, Galvez M, Gutierrez D, et al. VirusDetect: An automated
878           pipeline for efficient virus discovery using deep sequencing of small RNAs. Virology [Internet].
879           2017 Jan 1 [cited 2021 Oct 25];500:130–8. Available from:
880           https://linkinghub.elsevier.com/retrieve/pii/S0042682216303166

881    28.    Ison J, Kalaš M, Jonassen I, Bolser D, Uludag M, McWilliam H, et al. EDAM: an ontology of
882           bioinformatics operations, types of data and identifiers, topics and formats. Bioinformatics
883           [Internet]. 2013 May 15 [cited 2021 Oct 25];29(10):1325–32. Available from:
884           https://academic.oup.com/bioinformatics/article/29/10/1325/255660

885    29.    Low AJ, Koziol AG, Manninger PA, Blais B, Carrillo CD. ConFindr: rapid detection of intraspecies
886           and cross-species contamination in bacterial whole-genome sequence data. PeerJ [Internet]. 2019
887           May 31 [cited 2021 Feb 11];7(5):e6995. Available from: /pmc/articles/PMC6546082/

888    30.    Orakov A, Fullam A, Coelho LP, Khedkar S, Szklarczyk D, Mende DR, et al. GUNC: detection of
889           chimerism and contamination in prokaryotic genomes. Genome Biology 2021 22:1 [Internet].
890           2021 Jun 13 [cited 2021 Jul 24];22(1):1–19. Available from:
891           https://genomebiology.biomedcentral.com/articles/10.1186/s13059-021-02393-0

892    31.    Simion P, Belkhir K, François C, Veyssier J, Rink JC, Manuel M, et al. A software tool "CroCo"
893           detects pervasive cross-species contamination in next generation sequencing data. BMC Biol.
894           2018;16(1):1–9.

895    32.    Sangiovanni M, Granata I, Thind AS, Guarracino MR. From trash to treasure: detecting
896           unexpected contamination in unmapped NGS data. BMC Bioinformatics [Internet]. 2019 Apr 18
897           [cited 2021 Oct 25];20(Suppl 4). Available from: /pmc/articles/PMC6472186/

898  33.  Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
899       Bioinformatics [Internet]. 2009 Jul [cited 2022 Apr 28];25(14):1754–60. Available from:
900       https://pubmed.ncbi.nlm.nih.gov/19451168/

901  34.  Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: Estimating species abundance in
902       metagenomics data. PeerJ Comput Sci [Internet]. 2017 Jan 2 [cited 2022 Jul 12];2017(1):e104.
903       Available from: https://peerj.com/articles/cs-104

904  35.  Kechin A, Boyarskikh U, Kel A, Filipenko M. CutPrimers: A New Tool for Accurate Cutting of
905       Primers from Reads of Targeted Next Generation Sequencing. Journal of Computational Biology
906       [Internet]. 2017 Nov 1 [cited 2022 Mar 27];24(11):1138–43. Available from:
907       https://www.liebertpub.com/doi/full/10.1089/cmb.2017.0096

908  36.  de Clerck C, Crew K, van den houwe I, McMichael L, Berhal C, Lassois L, et al. Lessons learned
909       from the virus indexing of Musa germplasm: insights from a multiyear collaboration. Annals of
910       Applied Biology [Internet]. 2017 Jul 1 [cited 2022 Jun 19];171(1):15–27. Available from:
911       https://onlinelibrary.wiley.com/doi/full/10.1111/aab.12353

912  37.  Marais A, Faure C, Bergey B, Candresse T. Viral Double-Stranded RNAs (dsRNAs) from Plants:
913       Alternative Nucleic Acid Substrates for High-Throughput Sequencing. Methods in Molecular
914       Biology [Internet]. 2018 [cited 2021 Nov 19];1746:45–53. Available from:
915       https://link.springer.com/protocol/10.1007/978-1-4939-7683-6_4

916  38.  Chabannes M, Gabriel M, Aksa A, Galzi S, Dufayard JF, Iskra-Caruana ML, et al. Badnaviruses and
917       banana genomes: a long association sheds light on Musa phylogeny and origin. Mol Plant Pathol.
918       2021 Feb;22(2):216–30.

919  39.  Ricciuti E, Laboureau N, Noumbissié G, Chabannes M, Sukhikh N, Pooggin MM, et al.
920       Extrachromosomal viral DNA produced by transcriptionally active endogenous viral elements in
921       non-infected banana hybrids impedes quantitative PCR diagnostics of banana streak virus
922       infections in banana hybrids. Journal of General Virology [Internet]. 2021 Nov 2 [cited 2021 Nov
923       19];102(11):001670. Available from:
924       https://www.microbiologyresearch.org/content/journal/jgv/10.1099/jgv.0.001670

925  40.  Bal A, Pichon M, Picard C, Casalegno JS, Valette M, Schuffenecker I, et al. Quality control
926       implementation for universal characterisation of DNA and RNA viruses in clinical respiratory
927       samples using single metagenomic next-generation sequencing workflow. BMC Infectious
928       Diseases 2018 18:1 [Internet]. 2018 Oct 29 [cited 2021 Oct 25];18(1):1–10. Available from:
929       https://link.springer.com/articles/10.1186/s12879-018-3446-5

930  41.  Li CX, Li W, Zhou J, Zhang B, Feng Y, Xu CP, et al. High resolution metagenomic characterisation of
931       complex infectomes in paediatric acute respiratory infection. Scientific Reports 2020 10:1
932       [Internet]. 2020 Mar 3 [cited 2021 Oct 25];10(1):1–11. Available from:
933       https://www.nature.com/articles/s41598-020-60992-6

934  42.  Boheemen S van, Rijn AL van, Pappas N, Carbo EC, Vorderman RHP, Sidorov I. Since January 2020
935       Elsevier has created a COVID-19 resource centre with free information in English and Mandarin

936       on the novel coronavirus COVID- research that is available on the COVID-19 resource centre -

937       including this with acknowledgement of the origin. 2020;(January).

938