**Hecatomb: An End-to-End Research Platform for Viral Metagenomics**

Michael J. Roach[1], Sarah J. Beecroft[2], Kathie A. Mihindukulasuriya[3,4], Leran Wang[3,4],

Anne Paredes[3], Kara Henry-Cocks[1], Lais Farias Oliveira Lima[5], Elizabeth A. Dinsdale[1],

Robert A. Edwards[1], Scott A. Handley[3,4*]


1) Flinders Accelerator for Microbiome Exploration, Flinders University, Adelaide, SA,

Australia

2) Harry Perkins Institute of Medical Research, Perth, WA, Australia

3) Department of Pathology & Immunology, Washington University School of Medicine,

St. Louis, MO, USA

4) The Edison Family Center for Genome Sciences & Systems Biology, Washington

University School of Medicine, St. Louis, MO, USA

5) Biology Department, San Diego State University, San Diego, CA, USA

**\*Corresponding Author:** *shandley@wustl.edu*

## Keywords:

## Abstract

**Background:** Analysis of viral diversity using modern sequencing technologies offers extraordinary opportunities for discovery. However, these analyses present a number of bioinformatic challenges due to viral genetic diversity and virome complexity. Due to the lack of conserved marker sequences, metagenomic detection of viral sequences requires a non-targeted, random (shotgun) approach. Annotation and enumeration of viral sequences relies on rigorous quality control and effective search strategies against appropriate reference databases. Virome analysis also benefits from the analysis of both individual metagenomic sequences as well as assembled contigs. Combined, virome analysis results in large amounts of data requiring sophisticated visualization and statistical tools.

**Results:** Here we introduce Hecatomb, a bioinformatics platform enabling both read and contig based analysis. Hecatomb integrates query information from both amino acid and nucleotide reference sequence databases. Hecatomb integrates data collected throughout the workflow enabling analyst driven virome analysis and discovery. Hecatomb is available on GitHub at https://github.com/shandley/hecatomb.

**Conclusions:** Hecatomb provides a single, modular software solution to the complex tasks required of many virome analysis. We demonstrate the value of the approach by applying Hecatomb to both a host-associated (enteric) and an environmental (marine) virome data set. Hecatomb provided data to determine true- or false-positive viral sequences in both data sets and revealed complex virome structure at distinct marine reef sites.

## Background

24 Viruses parasitize host cell molecular processes and as a result alter host (prokaryotic

25 and eukaryotic) cell physiology. Virus-host interactions can influence organismal

26 physiology and environmental ecosystems. Viruses are also the most dominant entity on

27 the planet with current global estimates as high as $10^{31}$ viral particles [1,2], and they are

28 omnipresence in all cellular life forms [3]. As such they exert significant influence on their

29 surroundings.

30 The effect of viruses on human life and society are dramatically demonstrated through

31 phenomena such as global pandemics. However, the true burden of viruses on human

32 health is incredibly varied in terms of breadth and severity. There are many well-known

33 acute viral diseases such as the "common cold" (rhinoviruses, adenoviruses and

34 enteroviruses) [4] which cause tremendous amounts of morbidity, but limited mortality. In

35 contrast, chronic Epstein-Barr virus (EBV) infection has recently been associated with the

36 onset of multiple sclerosis [5]. Consequential virus-host interactions are not limited to

37 humans. For example, Geminivirus infection of plants has resulted in nearly US$2 billion

38 loss in African cassava production [6]. Similar foot-and-mouth disease virus (FMDV), a

39 highly contagious disease of cloven-hoofed animals, is widespread in Africa with an

40 annual US$2.3 billion negative impact on livestock [7]. Virus 'spillover' infection from

41 animal to human ("zoonotic" viruses) is unfortunately an all too regular event [8]. Viral

42 zoonosis from viruses such as SARS-CoV-2, monkeypox, Ebola and Zika viruses have

43 tremendous negative impacts on human health and society and new zoonotic viruses are

44 constantly emerging presenting a persistent threat to human health [9].

45    Viral assemblages, often referred to as *viromes*, are also associated with human health

46    and disease [10]. Stool samples from patients with inflammatory bowel disease (IBD)

47    suffer dysbiosis of microbial populations, having expanded numbers of bacteriophage

48    (hereafter *phage*) from the order Caudovirales [11–15]. Enteric vertebrate virus expansion

49    occurs in both rhesus macaques and humans with acquired immunodeficiency syndrome

50    (AIDS) [16,17]. Thus, health is not only influenced by infection with single viruses, but

51    also viromes. A comprehensive virus analysis workflow enables the analysis of both.

52    Viruses also influence global ecosystems. For example, the release of intracellular iron

53    and sulphur from bacteria following lytic phage infection releases nutrients used by

54    phytoplankton into marine environments via a mechanism called a *viral shunt* [18]. These

55    phytoplankton are in-turn eaten by higher trophic levels altering the entire marine food

56    web. Many other aquatic nutrient cycling and biogeochemical processes are attributed,

57    both directly and indirectly, to viral modification of prokaryotic and protistan assemblages

58    [19–23]. Terrestrial environment carbon and nutrient cycling are also influenced by

59    bacteriophage [24–26]. Viral modification of both aquatic and terrestrial ecosystems

60    underlies the importance of environmental virome studies to comprehensively understand

61    climate, ecology and production. Virome analysis tools should enable detailed

62    interrogation of viruses from any ecosystem broadening our understanding of the global

63    virome.

64    Metagenomic sequencing offers a powerful tool to study viral diversity [27]. However,

65    there are currently many challenges associated with viral metagenomics. While viruses

66    are the most abundant and diverse biological entity on the planet, they represent a

67    minority of reference genomes in GenBank, largely due to difficulties associated with

68   studying them [28]. Recent efforts to populate new viral genomes into reference

69   databases are slowly closing this gap, and have yielded 10s to 100s of thousands of novel

70   metagenome-assembled viral genomes [29–35]. Other efforts have yielded many new

71   high-quality viral genomes by combining the laborious and time-consuming experimental

72   work with student-learning outcomes [36]. Despite these efforts, there is still a vast

73   amount of sequence information that remains taxonomically or functionally ill-defined.

74   These sequences are regularly referred to as "viral dark matter" and poses a significant

75   barrier to the annotation of viral sequences from metagenomic data (reviewed in [37]).

76   The success of viral annotation is directly impacted by the size and diversity of the

77   reference database. Sensitive search algorithms are better able to identify viral

78   sequences that are only distantly related to reference database sequences. More diverse

79   databases improve viral sequence annotation, but larger databases are less conducive

80   to these high sensitivity searches due to increased computational requirements.

81   Database limitations are further amplified when deciding to query sequences against

82   amino acid or nucleotide reference databases. Translated searches to amino acid

83   databases offer superior sensitivity, however, limiting searches solely to amino acid

84   databases risks missing sequences only available in nucleotide databases.

85   Another challenge to interpretation of reference based sequence annotation is that viral

86   metagenomes are often plagued with false positive classifications [38–40]. Viruses share

87   sequence homology with other domains of life, including 'stolen' genes incorporated from

88   their hosts' genomes, and repetitive or low-complexity regions that are also found in other

89   organisms, such as insertion elements or transposons [38–40]. These sequences are

90   present in reference databases and can result in false-classifications due to shared

91 sequence similarity across taxonomies. The presence of false-positive classifications

92 may influence data interpretation. For instance, mis-classification of viral sequences in

93 clinical samples could lead to incorrect hypotheses about virus-disease associations or

94 patient diagnosis. Similarly, an increased false-positive rate in any environment could

95 lead to over-estimates of species richness and diversity. False-positive taxonomic

96 assignments are largely unavoidable without highly-curated databases which require

97 tremendous resources and time at the risk of missing newly discovered viruses which

98 have yet to make there way through the curation process. Thus, it is important for virome

99 analysis bioinformatic tools to provide a system to classify the quality of taxonomic

100 assignments empowering researchers to make informed decisions.

101 Here we present Hecatomb, a bioinformatics platform designed to address many of these

102 issues. Hecatomb performs rigorous quality control followed by tiered taxonomic

103 assignment using MMseqs2 querying sequences against virus-specific and trans-

104 kingdom amino acid and nucleotide reference databases [41]. Hecatomb also performs

105 metagenomic assembly and contig taxonomic classification providing simultaneous

106 analysis of both read and contig based viral annotations. While hecatomb provides pre-

107 compiled databases and recommended settings, it is easily customizable and extensible.

108 The primary output of Hecatomb is a comprehensive annotation table containing data

109 generated throughout the workflow that is designed to be easily merged with sample data

110 for visualisation and statistical analysis. Hecatomb has been successfully applied to

111 several viral metagenomics projects and has accelerated the discovery of novel viruses

112 and characterisation of viral populations [42–46].

113   Hecatomb is open-source with the project files hosted on GitHub at

114   github.com/shandley/hecatomb [47], with full support available using GitHub issues.

115   Documentation and training vignettes are available at hecatomb.readthedocs.io.

116   Documentation covers installing and optional configuration of the software; detailed

117   information including databases, and output files; advanced usage cases; an FAQ; and a

118   tutorial covering some example analyses of the results. Hecatomb is available for

119   installation from the Bioconda [48] and is distributed under a permissive MIT licence.

120   Bioconda package information for Hecatomb is available at

121   anaconda.org/bioconda/hecatomb.

## Implementation

123   An overview of the Hecatomb pipeline is shown in Figure 1. Hecatomb processes reads

124   through four key modules (Figure 1A). First (module 1), reads are preprocessed to

125   remove low-quality or contaminating sequences (low-quality sequence, primers,

126   adapters, host, common laboratory contaminants and duplicates). Second, preprocessed

127   reads are passed through both a read-based analysis and an assembly module (modules

128   2 and 3). For taxonomic assignment, Hecatomb uses preprocessed databases

129   (Supplementary Methods). The final module (module 4) combines information obtained

130   from both the read-based and assembly modules. Results are stored throughout each

131   module, primarily as tab-separated value (tsv) files for universal compatibility and easy

132   data analysis with any framework (e.g. Python, R, Bash, Excel). Emphasis is placed on

133   data preservation at each stage to provide analysts with as much detail as possible to

134   inform interpretation of results.

135 Hecatomb is installed via Conda and it makes liberal use of Conda environments to

136 ensure portability and ease of installation (Fig 1B). All required and optional software

137 dependencies are summarised in Table S1 [41,49–59]. Users need only install Conda

138 which Hecatomb uses to automatically install all dependencies. Conda environments for

139 jobs are created automatically by Snakemake [60] [49]. The use of isolated Conda

140 environments for Hecatomb and the individual pipeline jobs minimises package version

141 conflicts, minimises overhead when rebuilding environments for updated dependencies,

142 and allows maintenance and customisation of different versions of Hecatomb and its

143 dependencies without interacting with installed programs and system modules.

144 A custom built launcher script is included to make running the pipeline as simple as

145 possible. The launcher populates the required file paths, the default configuration, and

146 offers a convenient way to modify parameters and customise options. The Snakemake

147 command generated and runtime configuration is printed to the terminal window for the

148 user's reference. Accessory scripts are also available from this launcher for installing

149 reference databases, as well as adding custom host genomes, and combining results

150 from multiple analyses.

151 Hecatomb is able to be deployed on an high-performance computing (HPC) cluster and

152 has makes use of Snakemake profiles for cluster job schedulers (e.g. Slurm, SGE, etc.).

153 Snakemake uses profiles to submit pipeline jobs to the job scheduler and monitor their

154 progress. Although optional, using the scheduler is highly recommended as it allows for

155 more efficient use of HPC resources compared to submitting the whole Hecatomb

156 pipeline as a local job. Profiles can be created manually, but Hecatomb has been

157 designed for compatibility with the official Cookiecutter

158    (https://github.com/cookiecutter/cookiecutter)                profiles            for            Snakemake

159    (https://github.com/Snakemake-Profiles/doc).


160    **Sequence data preprocessing.** Hecatomb can process both single and paired-end

161    Illumina or MGI sequencing reads as well as long-read technology from PacBio and

162    Oxford Nanopore platforms. Hecatomb can also process sequences obtained from other

163    library types with minor modifications to the Hecatomb configuration file and by supplying

164    library specific adapters or primer sequences. A preprocessing module is also available

165    for sequencing utilising the round A/B library protocol for viral metagenomics [61]. The

166    round A/B library protocol enables sequencing of all types of viral genomes (single and

167    double stranded RNA and DNA viruses) and requires the use of a combination of phased

168    PCR primers. The preprocessing module in Hecatomb removes these non-biological

169    sequence contaminants, along with additional common laboratory sequence

170    contaminants in the UniVec database [62].


171    For host-associated samples (e.g. stool, saliva, skin swabs from humans or mucus from

172    corals) Hecatomb implements a host-sequence removal strategy using Minimap2 and a

173    host reference genome specifically optimised to avoid removing potential viral sequences

174    [56]. To remove all potentially viral sequences in reference genomes all viral genomes

175    from the National Center for Biotechnology Information (NCBI) viral assembly database

176    (ncbi.nlm.nih.gov/assembly/?term=viruses) were downloaded and computationally split

177    into short fragments with an average length of 85 bases sharing a 30 base overlap using

178    shred.sh from the BBTools suite [52]. Shredded viral sequences were mapped and

179    masked from host-reference genomes using bbmap.sh requiring a minimum identity of

180    90% and at most, 2 insertions and deletions. In addition, low-entropy sequences were

181  masked from host genomes (entropy = 0.5) using bbmask.sh. This process results in a

182  set of host-associated reference genomes masked of 'virus-like' and low-entropy

183  sequences, limiting the likelihood that a real viral sequence will be removed. Pre-

184  computed masked reference genomes for the following host genomes: human, mouse,

185  rat, camel, *Caenorhabditis elegans*, dog, cow, macaque, mosquito, pig, rat and tick are

186  available within Hecatomb using the --host flag. Scripts are provided to generate new

187  masked host genomes.

188  For the final stage of the preprocessing module, Hecatomb removes sequence

189  redundancy by clustering each sample using linclust [63]. Clustering sequences reduces

190  the number of sequences requiring taxonomic classification to a single, representative

191  sequence from a cluster of similar sequences. Sequences are clustered requiring a

192  minimum sequence identity of 97% and 80% alignment coverage of target sequence to

193  the representative sequence (--min-seq-id 0.97 -c 0.8 --cov-mode 1). Hecatomb

194  maintains the size of each cluster in the annotation table as well as the counts normalised

195  to the total number of high-quality reads per individual sample (normalized as percent of

196  the non-host reads). Clustering settings are also easily adjustable in the Hecatomb

197  configuration file.

198  At the end of this process, non-redundant sequences have been removed and the

199  remaining sequences are free from non-biological (reagent) contaminants and likely host-

200  sequences. These high-quality sequences are then used for *de novo* metavirome

201  assembly and read-based annotation.

202 **Metavirome assembly.** A unique feature of Hecatomb is that it completes both individual

203 read and assembly-based analysis. The first step of the metavirome assembly module is

204 individual sample assemblies using MEGAHIT [53] (Figure 2). Long-reads are not

205 amenable to using short-read assemblers and are therefore assembled using Canu [54].

206 Contigs from individual sample assemblies are subsequently assembled into a population

207 assembly using Flye [64]. Per sample contig abundance are calculated by mapping

208 individual sample reads to the population assembly. Read counts are reported normalised

209 to library size and contig length using a variety of measures (reads per kilobase million

210 (RPKM), fragments per kilobase million (FPKM) and sequences per million (SPM)). SPM

211 is the same calculation as used for transcripts per kilobase million (TPM) except that the

212 sequences are not assumed to be transcripts, thus the nomenclature adjustment.

213 Calculations for RPKM, FPKM, and SPM are summarised in Supplementary Methods and

214 an explanation is available in [65]. Taxonomic assignment of contigs in the population

215 assembly is accomplished using MMseqs2 [41], queried against the secondary nucleotide

216 database. Contig properties (e.g. length, GC-content) are combined with taxonomic

217 assignments and sample abundance estimates into a final table. This contig table is

218 merged with data obtained through the read based analysis to supplement contig

219 mapping data with read-based taxonomic assignments and individual read properties.

220 **Read-based annotation.** Taxonomy (and functional information when available) is

221 assigned using an iterative query strategy against both amino acid and nucleotide

222 reference databases (Figure 3A). This strategy is designed to minimise false-positive viral

223 annotations while maintaining sensitivity and runtime performance. All queries are carried

224 out using MMseqs2 [41]. The strategy starts with a translated query of all sequences

225    against a database of all viral (taxonomy ID: 10239) amino acid sequences in UniProtKB

226    [66] clustered at 99% identity to reduce redundancy and target database size. Any

227    sequence that matches a known viral protein is subsequently cross-checked against the

228    complete multi-kingdom UniClust50 amino acid database [67]. The use of the well-

229    annotated UniClust50 database enables functional as well as taxonomic annotation. This

230    two-step query strategy captures all potential viral sequences in the first step, reducing

231    the number of queries required in the secondary confirmatory step against the larger

232    multi-kingdom database. The MMseqs2 searches can be time-consuming. Options are

233    provided to use the default slower high-sensitivity parameters (--start-sens 1 --sens-steps

234    3 -s 7), or fast search parameters (-s 4.0) that yield greatly improved runtime performance.

235    Sequences not identified as viral-like using translated queries to the amino acid database

236    are subject to a similar iterative search using untranslated queries against a viral

237    nucleotide sequence database consisting of all viral sequences in GenBank clustered at

238    100% identity to remove redundancy. This primary search is followed by a secondary

239    confirmatory query against a polymicrobial nucleotide database containing representative

240    RefSeq genomes from bacteria (n = 14,933), archaea (n = 511), fungi (n = 423), protozoa

241    (n = 90) and plant (n = 145) genomes [68]. This iterative strategy enables sequence

242    queries to target databases to be run on commodity hardware while still having

243    representation of a broad diversity of non-viral kingdoms to minimise false-positive

244    annotations.

245    Following secondary translated and untranslated searches Hecatomb augments

246    sequence annotations using the lowest common ancestor (LCA) 2b-LCA algorithm

247    described in [69]. This approach provides conservative taxonomic assignments at lower-

248   nodes of the tree when similarity is found across a heterogeneous collection of

249   taxonomies. However, the LCA algorithm fails when crossing higher taxonomic levels.

250   For example, sequences with similarity to both bacterial and viral taxa have a LCA of

251   "root" in the NCBI tree, while viruses from distinct viral domains are assigned to "virus

252   root". Hecatomb detects these instances and instead of classifying them to the root

253   lineages refactors to the top-hit annotation. While this sometimes results in the

254   reclassification of sequences to a non-virus lineage (e.g. if the tophit was bacterial) this

255   novel approach provides additional information about sequences with ambiguous

256   taxonomic assignments. This can be useful for instance in the identification of prophage

257   regions which remains a challenging area of research [70,71].

258   **Outputs.** Hecatomb output files are described in Supplementary Methods. Output tables

259   are all tab-separated value (.tsv) files to ensure ease of use with data analysis. This

260   tabular format is universally compatible with commonly used research software and

261   programming languages such as Python, R, Excel or Bash and is easily merged with data

262   from external sources, such as viral Baltimore classifications, International Committee on

263   Taxonomy of Viruses (ICTV) taxonomy, or other external sources. The read annotation

264   file is designed to acquire, preserve and organise data obtained throughout the pipeline

265   with both study specific sample information as well as external data sources (Figure 3B).

266   The process of investigating and removing false-positive annotations in viral

267   metagenomes can be complex, but the abundance of alignment metrics in this file is

268   designed to empower researchers to perform this step quickly and easily.

269 **Results**

270 **Re-evaluation of a mammalian host-associated enteric virome.** Hecatomb's data

271 structure (Figure 3B) integrates a large amount of information about individual sequences

272 including taxonomic lineages, alignment statistics (e.g. E-values, percent identity,

273 alignment length) and data from external virus information resources (e.g. Baltimore

274 classification). To assess how this data structure can be used to evaluate the content of

275 a complex virome we reanalysed a previously published data set (95 samples) obtained

276 from stool samples collected from SIV-infected rhesus macaques (*Macaca mulatta*)

277 (NCBI BioProject accession: PRJEB9503) [16]. Sequence data were generated using the

278 Illumina MiSeq 2×250 bp paired-end protocol on libraries of total nucleic acid (DNA and

279 cDNA to enable detection of both RNA and DNA viruses) extracted from stool samples.

280 This data set was selected as it contains sequences from viruses from multiple Baltimore

281 classifications (RNA and DNA genomes) that infect a variety of cell types (e.g. animal and

282 plant). In addition, the original study identified differences in enteric virus abundances

283 associated with SIV infection, enabling a comparative quantitative benchmark to evaluate

284 Hecatomb with previously published results.

285 For the reevaluation study, Hecatomb was run using default parameters. Hecatomb's

286 taxonomic assignments classified sequences into phylogenetically diverse groups (Figure

287 4A). Bacteriophage from the family Microviridae and the order Caudovirales,

288 (Siphoviridae, Myoviridae and Podoviridae), were the most abundantly classified viral

289 sequence in the study. Hecatomb also identified a large number of sequences belonging

290 to the Picornaviridae and Adenoviridiae, viral families regularly associated with

291 gastrointestinal disease. Picronaviruses and adenoviruses were also identified in the

292    original study with several adenoviruses having their full genomes sequenced as well as

293    plaque purified [72]. Hecatomb also classified sequences belonging to a diverse set of

294    viruses typically associated with infection of plants and protists (Figure S1).

295    Hecatomb assigns NCBI taxonomy [73] using MMseqs2 [41] to query metagenomic

296    sequences to relevant reference sequence databases. Taxonomic assignments relying

297    on sequence similarity are dependent on the thresholds chosen. A permissive threshold

298    risks increasing the rate of false-positives, while a stringent threshold may result in an

299    increased rate of false-negatives. A perfectly accurate threshold is unlikely to exist,

300    particularly given the high-variability in evolutionary histories across all viral types. In this

301    case, plots and additional statistical analysis can prove useful in evaluating true- and

302    false-positive viral annotations. Hecatomb collects alignment statistics (e.g. e-values,

303    percent identity, alignment length, etc.) in the taxonomic assignment module and

304    organises these data to assist in the identification of both true and false-positive

305    taxonomic classifications.

306    As an example of how the alignment statistics can be used to evaluate true- or false-

307    positive taxonomic assignments we examined percent identity and alignment lengths of

308    the four viral families identified in the original study (Circoviridae, Picornaviridae,

309    Adenoviridae and Parvoviridae). Hecatomb also annotated sequences to these same 4

310    viral families using both translated queries to amino acid (aa) databases and untranslated

311    queries to nucleotide (nt) databases (Figure 4B).

312    While the statistics underlying sequence similarity searches are well understood, the

313    application of thresholds to those statistics to infer taxonomy and function are more

314   nebulous. Therefore, Hecatomb provides some additional guidelines to aid with the

315   determination of true positives compared to false positives. For example, a quadrant

316   system can be used to evaluate individual per family (or other taxonomic level)

317   assignments (Figure 4B). Sequences in the upper two quadrants are highly similar to

318   sequences in the reference databases over short (upper left, quartile 1 (Q1)) or long

319   (upper right, Q2) alignment lengths, while sequences in the lower two quadrants have low

320   similarity over short (lower left, Q3) or long (lower right, Q4) alignment lengths. For this

321   analysis we arbitrarily selected 70% identity to represent the cut-off between low and

322   high-identity for translated (aa database) and 90% identity for untranslated (nt database)

323   alignments. Translated alignment length is reported in nucleotide base pairs rather than

324   amino acid length. Therefore, a cutoff of 150 base pairs for both translated and

325   untranslated alignment lengths was chosen (Figure 4B). Using this framework it is clear

326   that there are many query sequences with high-identity (both short and long alignments)

327   to sequences in both the aa and nt reference databases for the 4 families of previously

328   identified animal viruses (Figure S2).

329   There were also a large number of query sequences classified as having statistically

330   significant sequence similarity to reference sequences from viruses of protists (Figure

331   4B). Mimiviridae, that infect Acanthamoeba, and Phycodnaviridae, that infect algae, are

332   both dsDNA viruses with large genomes [74]. While it is conceivable that these viruses

333   may exist in the stool samples of rhesus macaques via water or food, using the quadrant

334   framework there is little or no evidence of high-identity alignments to any sequence in

335   either the aa or nt databases (Figure 4B, Figure S2). Hecatomb does not automatically

336   remove sequences from these families as they would be common in environmental

337   datasets. There is evidence for short and long low identity alignments (quadrant 4) to both

338   Phycodnaviridae and Mimiviridae reference sequences. Thus, these sequences should

339   be analysed using additional metrics (i.e. E-values, abundance across samples, etc.) to

340   determine if these represent potentially novel viral sequences. This would not have been

341   possible using stringent E-value filtering prior to data analysis.

342   Hecatomb also quantifies the normalised number of sequences (percent of host-removed

343   reads) at each taxonomic depth. The normalised percent abundances per sample can be

344   evaluated as the number of sequences assigned to a taxonomy per sample enabling

345   statistical comparisons. The original study found evidence for four families of animal

346   viruses (Circoviridae, Picornaviridae, Adenoviridae and Parvoviridae) in stool samples

347   obtained from macaques infected with SIV or uninfected controls. The abundance of

348   sequences from each viral family were similar between SIV-infected and uninfected

349   macaques early in the study, but the abundance increased significantly as SIV-infection

350   progressed while remaining the same in uninfected control animals. Evaluation of the

351   normalised abundances for each of these four viral families using Hecatomb confirmed

352   the findings of the original analysis (Figure 4C).

353   There were several viral families represented only using untranslated alignment to

354   Hecatombs nucleotide database, including the Herpesviridae (Figure 5). All of the

355   sequences assigned to the Herpesviridae aligned to only three target GenBank entries

356   (Figure 5B). One entry (AF191073) dominated the similarities. All three were assigned a

357   taxonomy with very low E-values suggesting statistically significant alignments (Figure

358   5C). However, all three of these entries belong to a single type of herpesvirus, Stealth

359   virus 1 clone 3B43 [75]. The Stealth virus 1 genome was originally described as

360    containing sections of both bacterial and viral genes. The three Stealth virus sequences

361    identified by Hecatomb are identical to the bacterial segments when queried against the

362    NCBI nt database (Figure 5D), suggesting that they are bacterial in origin. Very few

363    sequences were found with alignments to the viral portion of the Stealth Virus 1 genome,

364    which would be expected due to the random, shotgun sequencing process. This suggests

365    that these sequences were called viral by hecatomb due to their similarity to a bacterial

366    region of a viral genome, but that they are more likely bacterial false-positive

367    contamination. Indeed, the original study identified Herpesviridae and many other false-

368    positive sequences that were only removed following computationally-expensive blastn

369    and blastx searches of the Non-Redundant nucleotide and protein databases [76].

370    **Evaluation of an environmental dataset.** We assessed Hecatomb's ability to analyse

371    non-human associated viromes by processing a previously studied coral reef dataset

372    (NCBI  BioProject  accession:  PRJNA595374)  [77,78].  The  dataset  consists  of

373    metagenomic sequencing (Illumina MiSeq, paired 2x250) of both seawater and coral

374    mucus from inner and outer sections of a Bermuda reef system. The original studies

375    identified statistically significant differences in bacterial compositions between the coral

376    mucus and seawater microbiomes and the coral mucus microbiomes from the inner and

377    outer reefs. However, the viruses were not described in the original study. To further

378    interrogate the viruses in these samples, study sequences were downloaded from SRA

379    and run through Hecatomb using the fast parameters (--fast). Of the top 20 most abundant

380    viral families, 10 are bacteriophages (Figure 6A). The relative abundance of viral families

381    are mostly higher in inner reef samples compared to outer reef samples with exceptions

382    such as Herelleviridae, Adintoviridae, Inoviridae, and unclassified Cressdnaviricota.

383    Principal coordinate analysis (PCoA) of Bray-Curtis dissimilarity and a subsequent

384    analysis of variance (ANOVA) confirmed a non-homogenous distribution across samples

385    and groups (p = 0.053) (Figure 6B). Inner reef samples cluster closely together whereas

386    outer reef samples appear to be far more varied. To examine compositional differences

387    at each site, we performed permutational analysis of variance (PERMANOVA) of Bray-

388    Curtis dissimilarity. We find that samples differ based on position (inner vs. outer reef: p

389    = 0.001) and on the combined position and source (reef and mucus: p = 0.001), but not

390    on source alone (inner vs. outer mucus: p = 0.185).

391    We calculated similarity percentage (SIMPER) between inner and outer samples, and

392    between the outer reef samples only to identify viruses distinct to each group. SIMPER

393    analysis identified many viral species that were significantly more abundant in inner reef

394    samples, but none that were more abundant in the outer reef samples (Figure S3). In

395    particular, Hecatomb classified reads to over 20 species of Synechococcus phage as

396    being associated with outer reef samples. Viruses that contributed the largest fold

397    differences included a phage that infects Verrucomicrobia (a mucin-degrading bacteria),

398    and Namao virus (a Mimiviridae protozoan virus) which might infect Symbiodinium–

399    coral's endosymbiotic dinoflagellate.

400    When comparing outer reef coral mucus with outer reef water samples, we identified eight

401    viruses that were more abundant in the reef water samples, with a phage that infects

402    Halomonas bacteria as the largest fold difference (Figure S4). The largest fold differences

403    observed in the coral samples included a Pyramimonas algae virus, a Vibrio phage, a

404    Rhizobium phage, and a Pseudomonas phage.

## Discussion

Virome sequencing is the premier approach to evaluate the viral content of both host-derived and environmental samples. In the broadest terms, virome sequencing is used to answer two questions: i) What individual viruses are present in a sample or set of samples? ii) How does virome composition compare between groups of samples? The answers to these questions can be used to evaluate different biological questions. For example, knowing what individual viruses are present in a sample can be useful for identifying etiological agents of infectious disease. In contrast, analysis of the total virome or collection of viruses within a sample can be used to characterise ecological niches between groups. Both types of studies are dependent on effective computational tools not only to identify and classify viral reads within a metagenome, but also to assist in interpretation of complex virome data in association with study data.

Virome analysis is almost entirely dependent on sequence similarity queries against reference sequence databases. Historically, there have been two approaches to accomplishing this. The first is 'brute force' wherein all unclassified sequences are queried against a comprehensive, multi-kingdom reference sequence database (e.g. NCBI nt or nr). This approach relies on the search algorithm (e.g. BLAST, DIAMOND [79]) to pick the best or lowest-common ancestor of a group of hits to provide a final taxonomic assignment to an unknown query sequence. Hecatomb takes a different approach by first capturing all 'potentially viral' sequences by first querying sequences against a viral sequence database. These 'potentially viral' sequences typically represent only a small fraction of the full metagenomic data making subsequent computation more tractable. To confirm viral taxonomic assignment, all potentially viral sequences are cross-checked

428    against a curated small transkingdom reference database containing genomic

429    representatives from all kingdoms of life. Hecatomb completes this iterative search

430    approach using translated searches against amino acid databases as well as

431    untranslated searches against nucleotide databases, combining the results of each to

432    ensure detection of viral sequences is database independent. This iterative search

433    strategy uses databases orders of magnitude smaller than comprehensive, multi-kingdom

434    databases (such as nt and nr) increasing computational efficiency without limiting viral

435    detection.

436    Hecatombs' design philosophy recognizes that there are no 'perfect' databases or search

437    algorithms. Both the brute force and iterative search approaches against comprehensive

438    or curated databases will result in different rates of true/false positives/negatives. Instead,

439    Hecatomb relies on providing a compiled and rich set of data for search result evaluation.

440    We used this strategy to reassess the virome composition of SIV-infected and uninfected

441    rhesus macaques [16]. The original study used an iterative approach, but relied on

442    comprehensive, transkingdom databases (NCBI nt and nr) and identified associations

443    between four families of animal viruses (Circoviridae, Picornaviridae, Adenoviridae and

444    Parvoviridae) and SIV-infection. The new Hecatomb trans-kingdom database is 6 orders

445    of magnitude smaller than GenBank nt ($5.0 \times 10^6$ versus $1.3 \times 10^{12}$) which results in a

446    significant reduction in computational time and resources. Hecatomb identified the same

447    four viral families and their relationship to SIV mediated disease (Figure 4C). Similar to

448    our analysis of these samples using Hecatomb, the original study also classified a number

449    of sequences to the Mimiviridae and Phycodnaviridae. Statistical comparison of these

450    sequences between groups (e.g. SIV-infected vs. uninfected) did not reveal any

451    significant associations thus they were not discussed further. However, new evaluation

452    of results from Hecatomb indicates that there were likely false-positive classifications

453    reported in the original analysis (Figure 4B). This highlights how coordinated data such

454    as alignment statistics and taxonomy can be powerful tools for virome evaluation.

455    We were also able to evaluate the viromes of environmental (non-host associated)

456    viromes. This analysis was primarily designed to identify compositional changes in

457    viromes between reef types (inner or outer) and within coral mucosa and the surrounding

458    water from a previously published metagenomic data set [77,78]. The original study

459    identified elevated levels of Pelagibacter, Synechococcus, and unclassified Rickettsiales

460    in inner reef samples compared to outer reef samples. Indeed, we found many elevated

461    Synechococcus phages and other cyanophages in inner reef samples. However, we

462    found only a few Mimiviridae viruses that were elevated which might be associated with

463    Pelagibacter and unclassified Rickettsiales, despite Pelagibacter being identified as the

464    most abundant genus in the original study. It's possible that Synechococcus and other

465    cyanobacteria growth rates are high, and that this is offset by greater viral activity (a viral

466    shunt) that results in nutrient cycling to other microbes in the reef system. Heterotrophic

467    bacteria and archaea are significant sources of fixed-nitrogen in coral reefs (reviewed in

468    [80]), so viral activity of cyanobacteria would therefore be beneficial to the entire reef

469    ecosystem by supplying both organic nitrogen, and by feeding these nitrogen-fixing

470    bacteria.

471    The inner reef coral mucus and reef water viromes clustered tightly suggesting that there

472    was little difference in these viromes. The consistency in viral compositions between coral

473    mucus and reef water samples of the inner reef systems is interesting and suggests an

474     equilibrated flux of viral particles between coral mucus microbiomes and the surrounding

475     reef water. Conversely, differences were observed in viral abundances of outer reef

476     samples, and most were found to be species that were more elevated in coral mucus

477     compared to reef water samples. The greater differences in viral compositions between

478     the outer reef coral mucus and water samples could indicate that the greater exchange

479     of water between the reef system and open ocean may be depleting viruses from this

480     ecosystem. Furthermore, the greater thermal stability and reduced particulate load (from

481     terrestrial runoff) results in a reduced turnover of coral mucus in the outer reef samples

482     (described in [77,78]), which may also contribute to the higher relative abundances of

483     viruses in inner reef systems in general.

484     **Conclusions**

485     Virome analysis is complex and requires efficient computational tools to generate analyst

486     friendly results. Hecatomb provides a comprehensive and computationally efficient

487     solution for both read- and assembly-based viral annotation and virome analysis. The

488     pipeline is delivered with a convenient and easy-to-use front end and is compatible with

489     different sequencing technologies. Hecatomb's comprehensive collection of data

490     throughout the running of the pipeline, in particular the collection of alignment statistics,

491     empowers identification and interrogation of viral taxonomic assignments. We

492     demonstrate Hecatomb's utility for rapid processing and analysis of viral metagenomes

493     with a well-studied validation gut viral metagenome dataset. We also demonstrate its

494     utility for mining regular metagenome samples for virome analysis by analysing an

495     existing environmental dataset.

496

497 **Declarations**

498 **Ethics approval and consent to participate**: Not applicable.

499 **Consent for publication:** All authors have confirmed consent for publication.

500 **Availability of data and materials:**

501 **Project name**: Hecatomb

502 **Project home page**: github.com/shandley/hecatomb

503 **Project documentation:** hecatomb.readthedocs.io

504 **Operating system**: Linux

505 **Programming language**: Python

506 **Other requirements**: Conda

507 **Licence**: MIT

508 **Restrictions to use by non-academics**: None

509 The reanalysis with Hecatomb utilised pre-existing datasets which are available under

510 the NCBI BioProject accessions PRJEB9503 for the macaque SIV dataset [16] and

511 PRJNA595374 for the coral reef dataset [77,78]. The Hecatomb annotations are

512 available at doi.org/10.5281/zenodo.6388251, and all commands used for analysing the

513 results are available at

514 gist.github.com/beardymcjohnface/3d3245b2bf6d9544c524f412037d5065.

515

519

520 **Author's contributions:** MJR, RAE, and SAH conceived the pipeline and data

521 structures. KAM, LW, and AP provided suggestions about the pipeline and data

522 visualisations. MJR, SJB, RAE, and SAH coded the pipeline. KH-C contributed to

523 documentation and analysis. MJR and SAH performed the analysis and interpretation.

524 LFOL, RAE, and EAD helped with interpretation of results. MJR, RAE, and SAH drafted

525 the original manuscript. All authors reviewed and edited the manuscript.

526

532

533 **List of abbreviations:**

534 AIDS: acquired immunodeficiency syndrome

535 SIV: simian immunodeficiency virus

536 HPC: high-performance computing

537 NCBI: National Center for Biotechnology Information

538    RPKM: reads per kilobase million

539    FPKM: fragments per kilobase million

540    SPM: sequences per million

541    LCA: lowest common ancestor

542    ICTV: International Committee on Taxonomy of Viruses

543    PERMANOVA: permutational analysis of variance

544    PCoA: principal coordinate analysis

545    ANOVA: analysis of variance

546    SIMPER: similarity percentag

# References

1. Hendrix RW, Smith MC, Burns RN, Ford ME, Hatfull GF. Evolutionary relationships among diverse bacteriophages and prophages: all the world's a phage. Proc Natl Acad Sci U S A. Elsevier; 1999;96:2192–7.

2. Mushegian AR. Are There 1031 Virus Particles on Earth, or More, or Fewer? J Bacteriol [Internet]. 2020;202. Available from: http://dx.doi.org/10.1128/JB.00052-20

3. Koonin EV, Dolja VV, Krupovic M, Varsani A, Wolf YI, Yutin N, et al. Global organization and proposed megataxonomy of the virus world. Microbiol Mol Biol Rev [Internet]. American Society for Microbiology; 2020;84. Available from: https://journals.asm.org/doi/10.1128/MMBR.00061-19

4. Heikkinen T, Järvinen A. The common cold. Lancet. Elsevier; 2003;361:51–9.

5. Bjornevik K, Cortese M, Healy BC, Kuhle J, Mina MJ, Leng Y, et al. Longitudinal analysis reveals high prevalence of Epstein-Barr virus associated with multiple sclerosis. Science. science.org; 2022;375:296–301.

6. Patil BL, Fauquet CM. Cassava mosaic geminiviruses: actual knowledge and perspectives. Mol Plant Pathol. 2009;10:685–701.

7. Casey-Bryars M, Reeve R, Bastola U, Knowles NJ, Auty H, Bachanek-Bankowska K, et al. Waves of endemic foot-and-mouth disease in eastern Africa suggest feasibility of proactive vaccination approaches. Nat Ecol Evol. 2018;2:1449–57.

8. Prempeh H. Foot and mouth disease: the human consequences [Internet]. BMJ. 2001. p. 565–6. Available from: http://dx.doi.org/10.1136/bmj.322.7286.565

9. Grange ZL, Goldstein T, Johnson CK, Anthony S, Gilardi K, Daszak P, et al. Ranking the risk of animal-to-human spillover for newly discovered viruses. Proc Natl Acad Sci U S A [Internet]. 2021;118. Available from: http://dx.doi.org/10.1073/pnas.2002324118

10. Liang G, Bushman FD. The human virome: assembly, composition and host interactions. Nat Rev Microbiol. 2021;19:514–27.

11. Norman JM, Handley SA, Baldridge MT, Droit L, Liu CY, Keller BC, et al. Disease-specific alterations in the enteric virome in inflammatory bowel disease. Cell. 2015;160:447–60.

12. Pérez-Brocal V, García-López R, Nos P, Beltrán B, Moret I, Moya A. Metagenomic Analysis of Crohn's Disease Patients Identifies Changes in the Virome and Microbiome Related to Disease Status and Therapy, and Detects Potential Interactions and Biomarkers. Inflamm Bowel Dis. 2015;21:2515–32.

13. Fernandes MA, Verstraete SG, Phan TG, Deng X, Stekol E, LaMere B, et al. Enteric

Virome and Bacterial Microbiota in Children With Ulcerative Colitis and Crohn Disease. J Pediatr Gastroenterol Nutr. ncbi.nlm.nih.gov; 2019;68:30–6.

14. Clooney AG, Sutton TDS, Shkoporov AN, Holohan RK, Daly KM, O'Regan O, et al. Whole-Virome Analysis Sheds Light on Viral Dark Matter in Inflammatory Bowel Disease. Cell Host Microbe. 2019;26:764–78.e5.

15. Liang G, Conrad MA, Kelsen JR, Kessler LR, Breton J, Albenberg LG, et al. Dynamics of the Stool Virome in Very Early-Onset Inflammatory Bowel Disease. J Crohns Colitis. 2020;14:1600–10.

16. Handley SA, Desai C, Zhao G, Droit L, Monaco CL, Schroeder AC, et al. SIV Infection-Mediated Changes in Gastrointestinal Bacterial Microbiome and Virome Are Associated with Immunodeficiency and Prevented by Vaccination. Cell Host Microbe. Elsevier; 2016;19:323–35.

17. Monaco CL, Gootenberg DB, Zhao G, Handley SA, Ghebremichael MS, Lim ES, et al. Altered Virome and Bacterial Microbiome in Human Immunodeficiency Virus-Associated Acquired Immunodeficiency Syndrome. Cell Host Microbe. 2016;19:311–22.

18. Poorvin L, Rinta-Kanto JM, Hutchins DA, Wilhelm SW. Viral release of iron and its bioavailability to marine plankton. Limnol Oceanogr. Wiley; 2004;49:1734–41.

19. Wilhelm SW, Suttle CA. Viruses and Nutrient Cycles in the Sea [Internet]. BioScience. 1999. p. 781–8. Available from: http://dx.doi.org/10.2307/1313569

20. Fuhrman JA. Marine viruses and their biogeochemical and ecological effects. Nature. 1999;399:541–8.

21. Wommack KE, Colwell RR. Virioplankton: viruses in aquatic ecosystems. Microbiol Mol Biol Rev [Internet]. Am Soc Microbiol; 2000; Available from: https://journals.asm.org/doi/abs/10.1128/mmbr.64.1.69-114.2000

22. Weinbauer MG. Ecology of prokaryotic viruses. FEMS Microbiol Rev. 2004;28:127–81.

23. Suttle CA. Viruses in the sea. Nature. 2005;437:356–61.

24. Emerson JB, Roux S, Brum JR, Bolduc B, Woodcroft BJ, Jang HB, et al. Host-linked soil viral ecology along a permafrost thaw gradient. Nat Microbiol. 2018;3:870–80.

25. Trubl G, Jang HB, Roux S, Emerson JB, Solonenko N, Vik DR, et al. Soil Viruses Are Underexplored Players in Ecosystem Carbon Processing. mSystems [Internet]. 2018;3. Available from: http://dx.doi.org/10.1128/mSystems.00076-18

26. Starr EP, Nuccio EE, Pett-Ridge J, Banfield JF, Firestone MK. Metatranscriptomic reconstruction reveals RNA viruses with the potential to shape carbon cycling in soil. Proc Natl Acad Sci U S A. 2019;116:25900–8.

616   27. Dutilh BE, Cassman N, McNair K, Sanchez SE, Silva GGZ, Boling L, et al. A highly
617   abundant bacteriophage discovered in the unknown sequences of human faecal
618   metagenomes. Nat Commun. 2014;5:4498.

619   28. Krishnamurthy SR, Wang D. Origins and challenges of viral dark matter. Virus Res.
620   Elsevier; 2017;239:136–42.

621   29. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD, Lawley TD.
622   Massive expansion of human gut bacteriophage diversity. Cell. Elsevier;
623   2021;184:1098–109.e9.

624   30. Nayfach S, Páez-Espino D, Call L, Low SJ, Sberro H, Ivanova NN, et al.
625   Metagenomic compendium of 189,680 DNA viruses from the human gut microbiome.
626   Nat Microbiol. nature.com; 2021;6:960–70.

627   31. Nayfach S, Roux S, Seshadri R, Udwary D, Varghese N, Schulz F, et al. A genomic
628   catalog of Earth's microbiomes. Nat Biotechnol. 2021;39:499–509.

629   32. Soto-Perez P, Bisanz JE, Berry JD, Lam KN, Bondy-Denomy J, Turnbaugh PJ.
630   CRISPR-Cas System of a Prevalent Human Gut Bacterium Reveals Hyper-targeting
631   against Phages in a Human Virome Catalog. Cell Host Microbe. 2019;26:325–35.e5.

632   33. Gregory AC, Zablocki O, Zayed AA, Howell A, Bolduc B, Sullivan MB. The Gut
633   Virome Database Reveals Age-Dependent Patterns of Virome Diversity in the Human
634   Gut. Cell Host Microbe. 2020;28:724–40.e8.

635   34. Paez-Espino D, Roux S, Chen I-MA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR
636   v.2.0: an integrated data management and analysis system for cultivated and
637   environmental viral genomes [Internet]. Nucleic Acids Research. 2019. p. D678–86.
638   Available from: http://dx.doi.org/10.1093/nar/gky1127

639   35. Tisza MJ, Buck CB. A catalog of tens of thousands of viruses from human
640   metagenomes reveals hidden associations with chronic diseases. Proc Natl Acad Sci U
641   S A [Internet]. 2021;118. Available from: http://dx.doi.org/10.1073/pnas.2023202118

642   36. Hanauer DI, Graham MJ, SEA-PHAGES, Betancur L, Bobrownicki A, Cresawn SG,
643   et al. An inclusive Research Education Community (iREC): Impact of the SEA-PHAGES
644   program on research outcomes and student learning. Proc Natl Acad Sci U S A.
645   National Acad Sciences; 2017;114:13531–6.

646   37. Pargin E, Roach M, Skye A, Edwards R, Giles S. The human gut virome:
647   Composition, colonisation, interactions, and impacts on human health [Internet]. OSF
648   Preprints. 2022. Available from: https://doi.org/10.31219/osf.io/s9px2

649   38. Rosseel T, Pardon B, De Clercq K, Ozhelvaci O, Van Borm S. False-positive results
650   in metagenomic virus discovery: a strong case for follow-up diagnosis. Transbound
651   Emerg Dis. Wiley; 2014;61:293–9.

652  39. Skewes-Cox P, Sharpton TJ, Pollard KS, DeRisi JL. Profile hidden Markov models
653  for the detection of viruses within metagenomic sequence data. PLoS One.
654  journals.plos.org; 2014;9:e105067.

655  40. Ponsero AJ, Hurwitz BL. The Promises and Pitfalls of Machine Learning for
656  Detecting Viruses in Aquatic Metagenomes. Front Microbiol. frontiersin.org;
657  2019;10:806.

658  41. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for
659  the analysis of massive data sets. Nat Biotechnol. nature.com; 2017;35:1026–8.

660  42. Roach M, Cantu A, Vieri MK, Cotten M, Kellam P, Phan M, et al. No Evidence
661  Known Viruses Play a Role in the Pathogenesis of Onchocerciasis-Associated Epilepsy.
662  An Explorative Metagenomic Case-Control Study. Pathogens [Internet]. 2021;10.
663  Available from: http://dx.doi.org/10.3390/pathogens10070787

664  43. Hesse RD, Roach M, Kerr EN, Papudeshi B, Lima LFO, Goodman AZ, et al. Phage
665  Diving: An Exploration of the Carcharhinid Shark Epidermal Virome. Viruses.
666  Multidisciplinary Digital Publishing Institute; 2022;14:1969.

667  44. Adiliaghdam F, Amatullah H, Digumarthi S, Saunders TL, Rahman R-U, Wong LP,
668  et al. Human enteric viruses autonomously shape inflammatory bowel disease
669  phenotype through divergent innate immunomodulation. Sci Immunol.
670  2022;7:eabn6660.

671  45. Kim AH, Armah G, Dennis F, Wang L, Rodgers R, Droit L, et al. Enteric virome
672  negatively affects seroconversion following oral rotavirus vaccination in a longitudinally
673  sampled cohort of Ghanaian infants. Cell Host Microbe. 2022;30:110–23.e5.

674  46. Mihindukulasuriya KA, Mars RAT, Johnson AJ, Ward T, Priya S, Lekatz HR, et al.
675  Multi-Omics Analyses Show Disease, Diet, and Transcriptome Interactions With the
676  Virome. Gastroenterology. 2021;161:1194–207.e8.

677  47. Roach M, Handley S, Edwards R, SarahBeecroft, Roach M, henr, et al.
678  shandley/hecatomb: v1.1.0 [Internet]. 2022. Available from:
679  https://zenodo.org/record/7042227

680  48. Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, et al.
681  Bioconda: sustainable and comprehensive software distribution for the life sciences. Nat
682  Methods. nature.com; 2018;15:475–6.

683  49. Roach MJ, Pierce-Ward NT, Suchecki R, Mallawaarachchi V, Papudeshi B, Handley
684  SA, et al. Ten simple rules and a template for creating workflows-as-applications
685  [Internet]. OSF Preprints. 2022. Available from: http://dx.doi.org/10.31219/osf.io/8w5j3

686  50. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine.
687  Bioinformatics [Internet]. academic.oup.com; 2012; Available from:
688  https://academic.oup.com/bioinformatics/article-abstract/28/19/2520/290322

689    51. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
690    Bioinformatics. academic.oup.com; 2018;34:i884–90.

691    52. Bushnell B. BBMap: A fast, accurate, splice-aware aligner [Internet]. Lawrence
692    Berkeley National Lab. (LBNL), Berkeley, CA (United States); 2014 Mar. Report No.:
693    LBNL-7065E. Available from: https://www.osti.gov/biblio/1241166

694    53. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node
695    solution for large and complex metagenomics assembly via succinct de Bruijn graph.
696    Bioinformatics. academic.oup.com; 2015;31:1674–6.

697    54. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable
698    and accurate long-read assembly via adaptive k-mer weighting and repeat separation.
699    Genome Res. genome.cshlp.org; 2017;27:722–36.

700    55. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads
701    using repeat graphs. Nat Biotechnol. nature.com; 2019;37:540–6.

702    56. Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics.
703    academic.oup.com; 2018;34:3094–100.

704    57. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
705    Alignment/Map format and SAMtools. Bioinformatics. academic.oup.com;
706    2009;25:2078–9.

707    58. Shen W, Ren H. TaxonKit: A practical and efficient NCBI taxonomy toolkit. J Genet
708    Genomics. Elsevier; 2021;48:844–50.

709    59. Shen W, Le S, Li Y, Hu F. SeqKit: A Cross-Platform and Ultrafast Toolkit for
710    FASTA/Q File Manipulation. PLoS One. journals.plos.org; 2016;11:e0163962.

711    60. Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, et al.
712    Sustainable data analysis with Snakemake. F1000Res. 2021;10:33.

713    61. Finkbeiner SR, Holtz LR, Jiang Y, Rajendran P, Franz CJ, Zhao G, et al. Human
714    stool contains a previously unrecognized diversity of novel astroviruses. Virol J.
715    Springer; 2009;6:161.

716    62. Sayers EW, Bolton EE, Brister JR, Canese K, Chan J, Comeau DC, et al. Database
717    resources of the national center for biotechnology information. Nucleic Acids Res.
718    2022;50:D20–6.

719    63. Steinegger M, Söding J. Clustering huge protein sequence sets in linear time. Nat
720    Commun. 2018;9:2542.

721    64. Kolmogorov M, Bickhart DM, Behsaz B, Gurevich A, Rayko M, Shin SB, et al.
722    metaFlye: scalable long-read metagenome assembly using repeat graphs. Nat
723    Methods. 2020;17:1103–10.

724 65. Zhao Y, Li M-C, Konaté MM, Chen L, Das B, Karlovich C, et al. TPM, FPKM, or
725 Normalized Counts? A Comparative Study of Quantification Measures for the Analysis
726 of RNA-seq Data from the NCI Patient-Derived Models Repository. J Transl Med.
727 2021;19:269.

728 66. UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. Nucleic
729 Acids Res. 2021;49:D480–9.

730 67. Mirdita M, von den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M. Uniclust
731 databases of clustered and deeply annotated protein sequences and alignments.
732 Nucleic Acids Res. 2017;45:D170–6.

733 68. Kitts PA, Church DM, Thibaud-Nissen F, Choi J, Hem V, Sapojnikov V, et al.
734 Assembly: a resource for assembled genomes at NCBI. Nucleic Acids Res.
735 2016;44:D73–80.

736 69. Hingamp P, Grimsley N, Acinas SG, Clerissi C, Subirana L, Poulain J, et al.
737 Exploring nucleo-cytoplasmic large DNA viruses in Tara Oceans microbial
738 metagenomes. ISME J. 2013;7:1678–95.

739 70. Roach MJ, McNair K, Giles SK, Inglis L, Pargin E, Decewicz P, et al. Philympics
740 2021: Prophage Predictions Perplex Programs [Internet]. bioRxiv. 2021 [cited 2022 May
741 12]. p. 2021.06.03.446868. Available from:
742 https://www.biorxiv.org/content/biorxiv/early/2021/06/03/2021.06.03.446868

743 71. Inglis LK, Edwards RA. How Metagenomics Has Transformed Our Understanding of
744 Bacteriophages in Microbiome Research. Microorganisms. Multidisciplinary Digital
745 Publishing Institute; 2022;10:1671.

746 72. Abbink P, Maxfield LF, Ng'ang'a D, Borducchi EN, Iampietro MJ, Bricault CA, et al.
747 Construction and evaluation of novel rhesus monkey adenovirus vaccine vectors. J
748 Virol. 2015;89:1512–22.

749 73. Schoch CL, Ciufo S, Domrachev M, Hotton CL, Kannan S, Khovanskaya R, et al.
750 NCBI Taxonomy: a comprehensive update on curation, resources and tools. Database
751 [Internet]. 2020;2020. Available from: http://dx.doi.org/10.1093/database/baaa062

752 74. Sun T-W, Yang C-L, Kao T-T, Wang T-H, Lai M-W, Ku C. Host Range and Coding
753 Potential of Eukaryotic Giant Viruses. Viruses [Internet]. 2020;12. Available from:
754 http://dx.doi.org/10.3390/v12111337

755 75. Martin WJ. Bacteria-related sequences in a simian cytomegalovirus-derived stealth
756 virus culture. Exp Mol Pathol. 1999;66:8–14.

757 76. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.
758 BLAST+: architecture and applications. BMC Bioinformatics. 2009;10:421.

759 77. Lima LFO, Alker A, Papudeshi B, Morris M, Edwards R, de Putron S, et al. Coral

760    and Seawater Metagenomes Reveal Key Microbial Functions to Coral Health and
761    Ecosystem Functioning Shaped at Reef Scale. 2021;

762    78. Lima LFO, Weissman M, Reed M, Papudeshi B, Alker AT, Morris MM, et al.
763    Modeling of the Coral Microbiome: the Influence of Temperature and Microbial Network.
764    MBio [Internet]. 2020;11. Available from: http://dx.doi.org/10.1128/mBio.02691-19

765    79. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using
766    DIAMOND. Nat Methods. 2015;12:59–60.

767    80. Rädecker N, Pogoreutz C, Voolstra CR, Wiedenmann J, Wild C. Nitrogen cycling in
768    corals: the key to understanding holobiont functioning? Trends Microbiol. 2015;23:490–
769    7.

## Figure Legends

**Figure 1: Hecatomb pipeline and implementation**

(**A**) The Hecatomb pipeline is divided into four modules. Sequencing reads for each sample undergo preprocessing and clustering (*orange*); quality trimmed reads for each sample undergo assembly and assemblies for each sample are coalesced into a single assembly (*green*); clustered reads undergo annotation using viral and multi-kingdom protein databases and clustered reads not annotated by the protein search are annotated using viral and multi-kingdom nucleotide databases (*blue*); read-based annotations are combined with the assembly to provide contig annotations (*pink*). The assembly stages– *green* and *pink*–can optionally be skipped. (**B**) Hecatomb takes in command line arguments, data, configuration parameters and outputs both results for analysis and run information. Hecatomb interacts with the job scheduler in high-performance computing (HPC) environments. Hecatomb distributes individual tasks to the job queue. Command- line arguments, *grey*; files, *yellow*; Conda environments, *blue*; scripts/programs, *green*; workload manager, *pink*.

**Figure 2: Metavirome Assembly**

(**1**) High-quality kmer-normalised sequences from individual samples are assembled using either MEGAHIT or Canu. (**2**) The sequences for each sample are mapped to their respective assemblies. (**3**) The unmapped reads from all samples are pooled together. (**4**) The pooled unmapped reads are assembled using either MEGAHIT or Canu. (**5**) The contigs from all sample assemblies and the unmapped reads assembly are combined

791    together. **(6)** Overlapping contigs are joined together using Flye using the subassemblies

792    algorithm.

793    **Figure 3: Read-based annotation**

794    **(a)** Iterative taxonomic annotation strategy. All alignments are completed using

795    MMSeqs2. **(1)** High-quality representative sequences are queried against a viral amino

796    acid (aa, *green*) sequence database. **(2)** Potentially viral sequences are subjected to a

797    secondary, confirmatory query against a multi-kingdom amino acid sequence database.

798    **(3)** Representative sequences that do not match a known viral amino acid are subjected

799    to an untranslated query to a viral nucleic acid sequence database (nt, *purple*) **(4)** followed

800    by a secondary, confirmatory query against a multi-kingdom nucleotide database **(5)**.

801    Sequences that have been classified as either viral (*blue*) or nonviral (*pink*) in either the

802    translated (aa database) or untranslated (nt database) queries are combined into a final

803    taxonomy table. **(b)** Read annotation data structure. **(1)** Read Annotations are generated

804    using the clustered sequences (seqtable.fasta). **(2)** The clustered sequence IDs are

805    unpacked to yield the sample ID, the number of reads that sequence represents, and the

806    percent of host-removed reads that sequence represents. **(3)** The alignment metrics from

807    the annotation module are joined into the read annotations using the sequence ID as the

808    primary key. **(4)** Taxonomic annotations are calculated and joined into the read

809    annotations again using the sequence ID. **(5)** ICTV viral classifications are joined into the

810    read annotations by the Taxonomic Family annotation. **(6)** Sample metadata can be

811    joined into the read annotation table using the sample ID as the primary key. **(7)** The read

812    annotation table with sample metadata can be quickly and easily analysed.

813    **Figure 4: Reanalysis of rhesus macaque stool viromes**

814    **(A)** Abundance of reads classified by viral Phylum (colour) and Type (shape). Phyla

815    represented by fewer than 1,000 reads were excluded. **(B)** Percent identity and alignment

816    lengths of all sequences classified for the 4 animal viruses identified in the previous study

817    and two viruses of protists. Horizontal (70% identity) and vertical (150 base alignment

818    length) dashed lines indicate a user-defined quadrant space. Each point represents an

819    individual sequence colored by classification method (aa = classified via a translated

820    search to an amino acid database, nt = classified via an untranslated search to a

821    nucleotide database). Panels A and B represent data obtained from all 95 samples in the

822    study. **(C)** Comparison of the number of sequences in SIV-infected and uninfected

823    samples. Significance determined by the Wilcoxon signed-rank test. * = $P \leq 0.05$, ** $P \leq$

824    0.01, *** $P \leq 0.001$, **** $P \leq 0.0001$.

825    **Figure 5: Ambiguous classification of bacterial sequences as Herpesviridae**

826    **(A)** Percent identity and alignment length of all sequences assigned to the Herpesviridae.

827    Note, there are no reads that were assigned using a translated search to an amino acid

828    (aa) database. **(B)** Representation of GenBank accessions assigned to the

829    Herpesviridae. **(C)** Summary of e-values for the 3 Herpesviridae accessions. **(D)**

830    Summary counts of the taxonomic hits using blastn to the NCBI nucleotide (nt) database

831    for each accession.

832    **Figure 6: Reanalysis of Coral Reef Metagenomes**

833    **(a)** The 20 most abundant viral families across coral reef samples. The sum of percent

834    reads for each sample type are shown for each viral family. Viral families have been

835    ordered and coloured by their Phyla. Points have been colored by sample type with inner

836    and outer reef water samples coloured light- and dark- blue respectively, and inner and

837    outer coral mucus samples colored light- and dark-green respectively. **(b)** Principle

838    coordinate analysis (PCoA) of viral species abundance. Inner and outer reef water

839    samples are coloured light blue and dark blue respectively. Inner and outer coral mucus

840    samples are coloured light and dark green respectively. The Vegan package was used to

841    calculate a bray-curtis distance matrix from the viral species counts, followed by

842    multivariate dispersions with betadisp, and an Analysis of Variance (ANOVA) identified a

843    non-homogenous distribution (P = 0.053). Ellipses for sample groups are drawn at 95%

844    confidence levels for multivariate t-distribution.

845    **Figure S1: Taxonomic subsets of virus types**

846    Viral families present in the 95-sample SIV reanalysis study **(A)** Plant viruses, and **(C)**

847    Protist viruses

848    **Figure S2: Sequence per Quadrant Evaluation**

849    Percentage of reads per quadrant in Figure 5. **(A)** translated (aa reference database) and

850    **(B)** untranslated (nt reference database)

851    **Figure S3: Viral abundance for inner and outer reef samples**

852    Viruses more abundant by Similarity Percentage (SIMPER) analysis in inner reef samples

853    are colored red. Viral species constituting 95% of variance that are significantly different

854    (p<0.05, log2 fold difference > 2) are shown. Infinite values are capped at an absolute

855    log2 fold difference of 5.

856    **Figure S4: Viral abundance for outer reef coral mucus and outer reef water**

857    **samples**

858    Viruses more abundant by Similarity Percentage (SIMPER) analysis in outer reef coral

859    mucus samples and outer reef water samples are coloured red and blue respectively.

860    Viral species constituting 95% of variance that are significantly different (p<0.05, log2 fold

861    difference > 1) are shown. Infinite values are capped at an absolute log2 fold difference

862    of 5.

# Figure 1

# Figure 2

# Figure 3

# Figure 4

# Figure 5

# Figure 6