

1 An improved experimental pipeline for preparing circular ssDNA
2 viruses for next-generation sequencing

3 Catherine D. Aimone¹, J. Steen Hoyer², Anna E. Dye¹, David O. Deppong¹, Siobain Duffy²,
4 Ignazio Carbone³, Linda Hanley-Bowdoin^{1*}

5

6

7 ¹Department of Plant and Microbial Biology, North Carolina State University, Raleigh NC
8 27695 USA

²Department of Ecology, Evolution, and Natural Resources, Rutgers University, New
Brunswick, NJ 08901USA

9 ³Department of Entomology and Plant Pathology, North Carolina State University, Raleigh NC
10 27695 USA

11

12 *Address correspondence to Catherine D. Aimone, cddoyle@ncsu.edu

13

14

15 **Abstract**

16 We present an optimized protocol for enhanced amplification and enrichment of viral DNA for
17 Next Generation Sequencing of begomovirus genomes. The rapid ability of these viruses to
18 evolve threatens many crops and underscores the importance of using next generation
19 sequencing efficiently to detect and understand the diversity of these viruses. We combined
20 enhanced rolling circle amplification (RCA) with EquiPhi29 polymerase and size selection to
21 generate a cost-effective, short-read sequencing method. This optimized protocol produced short-
22 read sequencing with at least 50% of the reads mapping to the viral reference genome. We
23 provide other insights into common misconceptions about RCA and lessons we have learned
24 from sequencing single-stranded DNA viruses. Our protocol can be used to examine viral DNA
25 as it moves through the entire pathosystem from host to vector, providing valuable information
26 for viral DNA population studies, and would likely work well with other CRESS DNA viruses.

Keywords: phi29, EquiPhi29, MiSeq, whiteflies, viral DNA sequencing

Highlights

- Protocol for short-read, high throughput sequencing of single-stranded DNA viruses using random primers
- Comparison of the sequencing of total DNA versus size-selected DNA
- Comparison of phi29 and EquiPhi29 DNA polymerases for rolling circle amplification of viral single-stranded DNA genomes

27

28

29 **1. Introduction**

30 *Begomoviruses*, one of the nine genera in the *Geminiviridae*, are single-stranded DNA
31 (ssDNA) viruses that infect a wide variety of plant species, including many important crops.
32 They are also classified as CRESS DNA viruses, a large group of circular ssDNA viruses that
33 encode replication-associated proteins (Rep) originating from a common ancestor (Zhao et al.,
34 2019). Eukaryotic CRESS DNA viruses impact a wide range of plant and animal hosts and
35 evolve rapidly (Zhao et al., 2019).

36 Cassava and tomato are among the important crops whose yields are severely impacted
37 by begomovirus diseases. Cassava is an important root crop in Africa, Asia, and Latin America,
38 with African farmers producing over half of the total cassava worldwide (FAOSTAT, 2016). In
39 Africa and more recently in Asia, cassava yields have been reduced by Cassava mosaic disease
40 (CMD), which is caused by a complex of 11 begomoviruses collectively referred to as cassava
41 mosaic begomoviruses (CMBs). Annual cassava losses in Africa have been estimated to be 15-
42 24% or 12-23 million tons (US \$1.2-2.3 billion) (Thresh J. M. , 1997; Uzokwe et al., 2016). In
43 some regions of Africa, cassava farmers have experienced losses of up to 95%. Tomato, an
44 important vegetable crop that is grown around the world, is a host for over 100 begomovirus
45 species, with the most devastating being tomato yellow leaf curl virus (Moriones and Navas-
46 Castillo, 2000). In the eastern United States, tomato production is also negatively impacted by a
47 second begomovirus, tomato mottle virus (ToMoV), which causes widespread disease with yield
48 losses of up to 50% (Abouzid, Polston, and Hiebert, 1992; Polston and Anderson, 1997). Given
49 their significant impact on agriculture, it is important to understand how begomoviruses change
50 over time and adapt to new hosts and environments. Next generation sequencing (Jeske, 2018) is
51 an important approach for gaining insight into begomovirus populations.

52 Begomoviruses fall into two classes – the Old World viruses and the New World viruses
53 (Lefeuvre et al., 2011). Their genomes consist of either one or two circular DNAs. CMBs are
54 Old World viruses, while ToMoV is a New World virus. The genomes of the CMBs and ToMoV
55 consist of two components designated as DNA-A and DNA-B that together total 5-6 Kb in size.
56 Both components are required for systemic infection (Stanley and Gay, 1983). The genome
57 components contain divergent transcription units separated by a 5' intergenic sequence that
58 contains the origin of replication and promoters for gene transcription (Hanley-Bowdoin et al.,
59 2013). DNA-A encodes 5-6 proteins necessary for replication, transcription, encapsidation, and
60 combatting host defenses (Hanley-Bowdoin et al., 2013). DNA-B encodes two proteins essential
61 for movement (Hanley-Bowdoin et al., 2013).

62 Begomoviruses are encapsidated into double icosahedral virions and transmitted by
63 whiteflies (*Bemisia tabaci*) (Hanley-Bowdoin et al., 2013). When a whitefly feeds on the phloem
64 of an infected plant, it acquires virions that can be transmitted to a healthy plant during the next
65 feeding cycle. Structural studies have shown that a begomovirus virion only contains one ssDNA
66 molecule, such that the DNA-A and DNA-B components of bipartite viruses are packaged
67 separately into virions (Bottcher et al., 2004). As a consequence, successful transmission of a
68 bipartite begomovirus requires acquisition and transmission of at least two virions – one
69 containing DNA-A and another containing DNA-B. Once the virions enter a phloem-associated
70 cell, viral ssDNA is released, converted to double-stranded DNA (dsDNA), and replicated via a
71 rolling circle mechanism (Hanley-Bowdoin et al., 2013). As infection proceeds, nascent viral
72 ssDNA can undergo multiple rounds of replication or be packaged into virions for future
73 transmission by whiteflies. ToMoV is only transmitted by whiteflies, while CMBs can be

74 transmitted via vegetative propagation of infected stem cuttings as well as by whiteflies (Legg et
75 al., 2014).

76 Begomoviruses have been shown to evolve rapidly (Duffy and Holmes, 2009; Lima et al.,
77 2017; Rocha et al., 2013), making them good models for studying the evolution of ssDNA
78 viruses. Many factors can contribute to begomovirus evolution, including agricultural practices,
79 whitefly transmission, and abiotic stress. This underscores the importance of understanding how
80 begomoviruses evolve through an entire pathosystem from host to vector. Generally, ssDNA
81 viruses exist as genetically diverse populations, with variation similar to that of RNA virus
82 populations (Elena and Sanjuán, 2007; Safari and Roossinck, 2014). The high genetic diversity
83 of virus populations is linked to their rapid ability to evolve, to emerge in a new host, and to
84 break disease resistance (Duffy, 2008). Viral diversity is driven by high mutation and
85 recombination rates (Lefeuvre and Moriones, 2015; Sanjuán et al., 2010). The amount and type
86 of genetic variation within a viral population is a direct measure of evolvability and pathogenesis
87 (de la Iglesia and Elena, 2007; Elena, Fraile, and García-Arenal, 2014; Elena and Sanjuán, 2007).
88 However, our ability to study viral evolution in real-time has been limited by being able to
89 accurately describe the genetic structure of viral populations over time (Acevedo, Brodsky, and
90 Andino, 2014).

91 Deep sequencing technologies have advanced our understanding of the genetic variation
92 of evolving virus populations, beyond virus identification and characterization of viral species
93 (Acevedo et al., 2014). Next generation sequencing can provide insight into how viral
94 populations, as quasi-species, are highly variable with a range of beneficial or neutral mutations
95 that occur at low frequency (Dean et al.; Dickins and Nekrutenko, 2009). With NGS, we can
96 actively track naturally occurring viral variants through infection, adaptation to new hosts, and

97 host range expansion (Ruark-Seward et al., 2020). Yet, the ability of NGS to track viral variants
98 is restricted by several technical challenges, including biased amplification, errors introduced
99 during amplification and sequencing, and low viral read depth.

100 Current methods of viral amplification rely on the polymerase chain reaction (PCR) using
101 virus-specific primers and Taq polymerase or rolling circle amplification (RCA) using random
102 hexamers and phi29 DNA polymerase (Dean et al., 2001). Unlike RCA, sequence-specific PCR
103 can introduce sequence bias that masks viral diversity in a population (Sipos et al., 2010). RCA
104 amplifies circular episomes like begomovirus genomes more efficiently than linear DNA,
105 thereby enriching for begomovirus sequences (Idris et al., 2014). RCA is less susceptible to
106 sequence bias, less error-prone than traditional PCR, and does not fix errors in the sample to be
107 sequenced (Lou et al., 2013; Wang et al., 2014). However, RCA produces hyper-branched,
108 concatenated products (Lasken and Stockwell, 2007), that must be linearized by restriction
109 enzyme digestion or mechanical shearing before NGS sequencing (Inoue-Nagata et al., 2004).
110 Short-read sequencing in combination with RCA and size selection has improved viral read
111 depth for RNA viruses (Acevedo and Andino, 2014; Acevedo et al., 2014). For DNA viruses,
112 long-read sequencing of size-enriched viral DNA has been successful in reducing sequencing
113 error (Mehta et al., 2019).

114 Short-read sequencing has been used in combination with RCA to enrich for CMB viral
115 sequences in cassava (Kathurima, 2016) and *Nicotiana benthamiana* (Chen, Khatabi, and
116 Fondong, 2019). Other begomovirus species, including tomato leaf curl New Delhi virus
117 (Juárez et al., 2019) and euphorbia yellow mosaic virus (Richter et al., 2016) have also been
118 amplified by RCA for short-read sequencing. Likewise, RCA has been used to improve read
119 depth of mastreviruses, which constitute another genus in the *Geminiviridae* (Claverie et al.,

120 2019). Depending on the research question, the methods cited above can return sufficient viral
121 read depth for identifying new viruses and determining the prominent viruses in an infected
122 plant. However, to reliably detect subconsensus viral variants in a population, a higher level of
123 coverage is required (Juárez et al., 2019).

124 Here, we describe an experimental pipeline for analyzing begomovirus DNA population
125 dynamics across a complete pathosystem constituted by the plant host and the insect vector. The
126 pipeline combines size selection and linear amplification by an improved phi29 DNA
127 polymerase to increase viral read coverage for diversity studies (Fig. 1).

128 **2. Methods**

129 *2. 1. Virus-infected plants and viruliferous whiteflies*

130 Cassava plants (*Manihot esculenta* cv. Kibandameno or Kibaha) were propagated from
131 stem cuttings and grown at 28°C under a 12-h light/dark cycle. Plants with ca. 8-10 nodes and
132 stems 1.5 cm in diameter (ca. 2 months after propagation) were inoculated at the apical meristem
133 using a hand-held micro sprayer (40 psi) to deliver gold particles coated with plasmid DNA (100
134 ng/plasmid/plant) (Ariyo et al., 2006; Cabrera-Ponce et al., 1997). The plasmids, which
135 contained partial tandem dimers of DNA-A and or DNA-B of *African cassava mosaic virus*
136 (ACMV; GenBank accessions MT858793.1 and MT858794.1) and *East African cassava mosaic*
137 *Cameroon virus* (EACMCV; AF112354.1 and FJ826890.1)(Chowda Reddy et al., 2012;
138 Fondong and Chen, 2011; Fondong et al., 2000; Hoyer et al., 2020). Three plants were co-
139 inoculated with both ACMV and EACMCV. Leaf punches from symptomatic cassava plants
140 were sampled at 28 days post-infection (dpi), flash-frozen in liquid nitrogen, and stored for
141 analysis.

142 Tomato seedlings (*Solanum lycopersicum* cv. Florida Lanai) were grown from seed at
143 25°C under a 12-h light/dark cycle. Plants with five true leaves (ca. 4 weeks old) were
144 agroinoculated with ToMoV DNA-A and DNA-B (Abouzid et al., 1992; Reyes et al., 2013) as
145 described by Rajabu et al. (2018). An infected plant was sampled at the third leaf below the
146 apical meristem at 21 dpi, immediately prior to the whitefly access period for the acquisition of
147 ToMoV. The leaf tissue (1 mg) was separated into two parts, one part for total DNA extraction
148 and the other part for virion extraction. *Bemisia tabaci* MEAM1 adult whiteflies between 2 and
149 10 days post-eclosion were allowed to acquire the virus by feeding on a symptomatic plant
150 infected with ToMoV for an Inoculation Access Period (IAP) of 72 h (Ng et al., 2011; Rajabu et
151 al., 2018). Whiteflies were collected via aspiration and stored in 70% ethanol for analysis.

152 2.2 DNA extraction and size selection

153 Frozen leaf tissue was ground using a homogenizer (Model# MM 301, RETSCH-
154 Laboratory Mills, Clifton, NJ), and total DNA was extracted from leaf samples using the
155 MagMax™ Plant DNA Isolation Kit according to manufacturer's instructions (Thermo Fisher
156 Scientific, Waltham, MA). Total DNA was extracted from groups of five whiteflies using the
157 Qiagen DNeasy Blood and Tissue kit according to the manufacturer's instructions (Qiagen,
158 Hilden, Germany). Total DNA from cassava, tomato, and whiteflies (250 ng) was size selected
159 for 1-6 Kb DNA on a 0.75% agarose gel at 25V DC for 3-8 h using the Blue Pippin Prep system
160 (Model # BDQ3010, Sage Science, Beverly MA). The amount of size-selected output DNA was
161 typically less than 1 ng.

162 Virion DNA was generated by homogenizing five whiteflies or resuspending 1 mg of
163 frozen cassava or tomato leaf tissue in 50 mM Tris, 10 mM MgSO₄, 0.1 M NaCl, pH 7.5,
164 followed by low-speed centrifugation. The supernatant was subjected to 0.22 µM filtration

165 followed by DNase I digestion (2.5 U for 3 h at 37°C). Virion DNA was isolated using the
166 QIAamp MinElute Virus Spin Kit (Qiagen, Hilden, Germany)(Ndunguru et al., 2016; Ng et al.,
167 2011; Rosario et al., 2015).

168 2.3. *Viral levels*

169 The concentration of ACMV DNA-A (primer pair - P3P-AA2F and P3P-AA2R+4R;
170 Table 1), DNA-B (primer pair - ACMVBdiv4 and ACMVBfor1; Table 1), and EACMCV DNA-
171 A (primer pair – EACMVQ1 and EACMVQ; Table 1) and DNA-B (primer pair –
172 EACMVBREV4 and EACMVBfor1.2; Table 1) were measured by quantitative PCR (qPCR) in
173 total DNA samples (0.01 µg) extracted from cassava leaf tissue and analyzed in 96-well plates
174 on a Max3000P System (Stratagene, San Diego CA). Primers were tested in conventional PCR to
175 optimize annealing temperature and amplification efficiency for qPCR. For ACMV DNA-A,
176 qPCR was performed using Power SYBR Green PCR Master Mix (Applied Biosystems, Foster
177 City CA), starting with a 2 min denaturing step at 94°C, followed by 30 cycles consisting of 15
178 sec at 94°C, 1 min at 60°C, 30 sec at 72°C. The PCR conditions for EACMCV DNA-A were 10
179 min at 95°C, followed by 30 cycles of 30 sec at 95°C, 30 sec at 60°C, and 30 sec at 72°C.
180 Reactions were performed in three technical replicates. ACMV DNA-B and EACMCV DNA-B
181 were run following the above conditions respectively with an annealing temperature of 58°C.
182 Viral DNA was quantified using a qPCR standard curve generated by amplification of a 10-fold
183 dilution series (10^{-10} to 10^{-16} g/µL) of plasmid DNA with a single copy of ACMV DNA-A or
184 EACMCV DNA-A following the protocol described by Rajabu et al. (2018). The concentration
185 of the template DNA in the reaction mix was converted from ng/µL to copy number/µL using the
186 following formula; $(C \times 10^{-9}/MW) \times NA$ where C = template concentration ng/µL,
187 MW = template molecular weight in Daltons, and NA = Avogadro's constant 6.022×10^{23} . MW

188 was obtained by multiplying the number of base pairs of a plasmid by the average molecular
189 mass of one base pair (660 g/mol). A base 10 logarithmic graph of copy number versus the
190 threshold cycle (Li et al., 2009) for the dilution factor was plotted and used as a standard curve to
191 determine the amount of viral DNA (copy number/ μ L) of total DNA in a reaction mix (Rajabu et
192 al., 2018).

193 The concentration of ToMoV genomic components was quantified by qPCR in total
194 DNA samples (0.01 μ g) from tomato leaf tissue and in whiteflies (2 ng) using the DNA-A primer
195 pair, ToMoVA6-F, and ToMoVA6-R, and the DNA-B primer pair, ToMoVB4-F, and
196 ToMoVB4-R (Table 1). The qPCR protocol was the same as described above for ACMV with an
197 annealing temperature of 57°C for 1 min. Viral DNA was quantified using a qPCR standard
198 curve generated by amplification of plasmid DNA with a single copy of each ToMoV segment
199 (pNSB1691 and pNSB1692) following the method described above.

200 *2.4. Rolling circle amplification*

201 Total DNA (100 ng) from symptomatic Kibaha leaf tissue was amplified using the
202 TempliPhi Amplification Kit (GE Healthcare, Chicago IL), which contains phi29 DNA
203 polymerase (Dean et al., 2002), according to the manufacturer's instruction at 30°C for 18 h. The
204 reaction buffer solution of the TempliPhi Amplification Kit included random hexamer primers.
205 Separately, 2 μ L of total DNA was denatured at 95°C for 3 min, then cooled on ice for 3 min for
206 amplification with the EquiPhi29 kit (Thermo Fisher Scientific, Waltham MA). The cooled
207 reaction was mixed with 0.5 μ L of 10X EquiPhi29 Reaction Buffer, 1.0 μ L of Exo-resistant
208 random primers, and 1.5 μ L of nuclease-free water. The denatured DNA product (5 μ L) was
209 amplified using 1 μ L (10 U) of EquiPhi29 DNA polymerase, 1.5 μ L of 10X EquiPhi29 Reaction
210 Buffer, 0.2 μ L of 100 mM DTT, 2 μ L of 10 mM dNTP mix, 1.0 μ L (0.1 U) of pyrophosphatase

211 and 9.3 μL of nuclease-free water (Povilaitis et al., 2016) according to the manufacturer's
212 instructions, except that the reactions were performed at 40°C for 2 h. EquiPhi29 conditions
213 were optimized to retain the highest amount of dsDNA based on Povilaitis et al., 2016, the
214 manufacturer's report, and an optimization experiment (Povilaitis et al., 2016) (Supp. 1A).

215 Total DNA (100 ng) and RCA products (100 ng) generated using the TempliPhi
216 Amplification Kit were treated with 1 μL (1 U) of Mung Bean nuclease (New England Biolabs,
217 Ipswich, MA), 3 μL of CutSmart Buffer (New England Biolabs) in a 30 μL reaction volume for
218 30 min at 30°C. The solution was inactivated with 3 μL of SDS (0.01%), and DNA was
219 recovered by ethanol precipitation, according to the manufacturer's specifications. RCA products
220 (11 μL) generated using the TempliPhi Amplification Kit were also treated with 1.0 μL (1 U) of
221 Klenow (large subunit) (New England Biolabs) and 1.0 μL (1 U) of T4 DNA polymerase (New
222 England Biolabs) in a mixture of 5.0 μL 10X NEB Buffer #2 (New England Biolabs), 0.5 μL 10
223 mM dNTP, and 31.5 μL nuclease-free water at 25°C for 1 h. The reaction was inactivated by
224 incubation at 75°C for 1 h. After the repair reaction, residual salt and enzyme were removed by
225 suspending 60 μL of SPRIselect beads (Beckman Coulter, Pasadena CA) in the 50 μL
226 RCA/repair reaction. The mixture was incubated at room temperature for 5 min, placed on a
227 DynaMag[™]-2 magnetic rack (Thermo Fisher Scientific, Waltham MA) for 5 min at room
228 temperature or until the liquid was clear. The liquid was removed and the bead pellet was
229 washed in 200 μL of 80% ethanol on the magnetic rack for 30 sec. The ethanol was removed and
230 air-dried for 5 min on the magnetic rack. The pellet was resuspended in 17 μL of nuclease-free
231 water.

232 The concentrations of the RCA products after the various treatments described above
233 were measured using a Qubit 3.0 fluorometer with the dsDNA HS assay kit or ssDNA assay kit,

234 according to manufacturer's instructions (Thermo Fisher Scientific, Waltham MA). The dsDNA
235 HS assay kit only detects dsDNA, while the ssDNA assay kit detects both ssDNA and dsDNA
236 (i.e. total DNA). For all dsDNA concentration measurements (e.g., Figure 2), Qubit fluorometer
237 readings were used to calculate the amount of DNA. The amount of total dsDNA (in nanograms)
238 was calculated directly using the high sensitivity dsDNA buffer to measure the concentration
239 (ng/ μ L) and multiplying by the volume of the RCA reaction. Total ssDNA mass was estimated
240 by multiplying the concentration (ng/ μ L) determined using ssDNA buffer by the volume of the
241 RCA reaction, and subtracting the total mass of dsDNA determined using high sensitivity
242 dsDNA HS assay kit.

243 *2.5. Library Preparation*

244 Total DNA, size-selected DNA, or virion DNA (2 μ L) from Kibandameno leaves, tomato
245 leaves, or whiteflies was amplified using EquiPhi29 DNA polymerase as described above. Two
246 separate reactions were set up for each sample. Each RCA reaction was treated with Klenow and
247 T4 DNA polymerases followed by a purification step using SPRIselect beads, as described
248 above. After the clean-up step, each RCA reaction was diluted to 0.2 ng/ μ L (1 ng total in 5 μ L)
249 for library construction. Libraries were prepared using the Nextera XT DNA Library Prep Kit
250 (Illumina, San Diego CA) following manufacturer's instructions using Unique Dual Index
251 adaptors (Integrated DNA Technologies, San Jose CA) and 12 rounds of PCR. The Nextera XT
252 DNA Library Prep Kit was chosen because of its rapid library preparation, low input of DNA,
253 and optimization for small genomes. The libraries were cleaned with 40 μ L of SPRIselect beads,
254 as described above. The libraries were analyzed for size distribution (400-800 bp), yield, and
255 quality using a 2100 Bioanalyzer Instrument (Agilent Technologies, Santa Clara, CA). The
256 concentration of each library was determined using a Qubit 3.0 fluorometer using the dsDNA HS

257 assay kit as described above. Libraries were diluted to 15 nM, and equal molar amounts were
258 pooled for sequencing on an Illumina MiSeq platform.

259 The molarity of each library was determined using the following formula:

260 $(\text{ng}/\mu\text{L})/(\text{nmol}/\mu\text{L}) \times (1 \times 10^6)$. The molecular weight of each library was determined by taking
261 the average base-pair length from the Bioanalyzer profile multiplied by the average mass of the
262 four nucleotide bases plus the weight of 5'-PO₄ ((average base-pair length x 607.8) + 157.9).

263 The average mass of nucleotide bases and the MW of 5'-PO₄ was taken from Thermo Fisher
264 Scientific DNA and RNA Molecular Weight and Conversion guide (Scientific). Libraries were
265 pooled for sequencing on an Illumina MiSeq instrument (Fig. 1). A full printable version of our
266 protocol is available at cassavavirusevolution.vcl.ncsu.edu.

267 2.6. Data Analysis

268 Raw sequencing data were processed using Cutadapt (v.1.16.3) to remove the universal
269 3' adapters from the paired-end reads and to trim the 5' ends to give fastq quality scores > 30
270 (Martin, 2011). The quality-controlled reads were aligned to the DNA-A and DNA-B
271 components of ACMV and EACMCV or ToMoV using BWA mem (v. 0.8) with default
272 parameters (Li, 2013; Li et al., 2009). Quality-controlled reads were also aligned to the cassava
273 reference genome v.7 (Bredeson et al., 2016) for cassava plant samples, to the tomato reference
274 SL4.0 assembly (Hosmani et al., 2019) for tomato samples, and to MEAM1 assembly (GenBank
275 ASM185493v1) (Chen et al., 2016) for whitefly samples. Duplicate reads were discarded using
276 Picard MarkDuplicates (v. 2.18.2.1, <http://broadinstitute.github.io/picard/>). Samtools idxstats (v.
277 2.0.3) was used to generate mapping statistics (Li et al., 2009). Sufficient read coverage was
278 designated as 1000X fold coverage (~20,000 reads/genome) based on (Juárez et al., 2019).
279 Coverage was calculated using the following formula from Illumina (coverage=read length X

280 number of reads/genome length). For workflow and full parameters see Galaxy workflow,
281 ViralSeq (cassavavirusevolution.vcl.ncsu.edu) (Giardine et al., 2005). Raw Illumina data are
282 available at the NCBI Sequence Read Archive (PRJNA658475).

283 **3. Results**

284 *3.1. Analysis of RCA variability*

285 Begomoviruses have circular ssDNA genomes that are converted to dsDNA during viral
286 replication in plants. Viral ssDNA accumulates to high levels during infection, while viral
287 dsDNA occurs at much lower levels. Many library protocols for short-read sequencing involve
288 the ligation of adapters to dsDNA ends or transposase-mediated fragmentation and tagging of
289 dsDNA. Thus, the *in vitro* conversion of viral ssDNA to dsDNA is the first step during library
290 construction for ssDNA viruses. RCA is used frequently to amplify circular viral DNA genomes
291 and to convert viral ssDNA to dsDNA (Inoue-Nagata et al., 2004; Dean et al., 2001). We
292 examined the efficiency of RCA to convert ssDNA to dsDNA and several parameters that might
293 influence the amount of virus-specific dsDNA available for library construction.

294 Total DNA isolated at 28 dpi from four symptomatic Kibaha plants infected with ACMV
295 and EACMCV was incubated in RCA reactions containing phi29 and random hexamers
296 (TempliPhi Amplification Kit). The RCA reactions and an equal amount of unamplified total
297 DNA were digested with Mung Bean nuclease (MB) to remove ssDNA from the samples. RCA
298 increased the amount of total DNA 10-fold relative to input (Fig. 2A). The amounts of the input
299 DNA and the DNA after RCA were both greatly reduced by MB treatment, indicating that most
300 of the DNA before and after RCA is single-stranded and therefore cannot be ligated to library
301 adapters (Fig. 2A).

302 For successful library construction, we sought to increase the amount of viral dsDNA
303 after RCA. During the RCA de-branching step, ssDNA overhangs may be present on dsDNA
304 products, preventing adaptor ligation and leading to exclusion from the final library. To decrease
305 ssDNA overhangs after RCA, we used an end repair reaction to convert ssDNA overhangs to
306 dsDNA. Using the same total DNA sample in Fig. 2A, RCA was performed followed by an end-
307 repair by T4 polymerase and DNA polymerase I (Klenow fragment). The end-repair reaction
308 increased the amount of dsDNA 2-fold (Fig. 2B), indicating that repairing the debranched ends
309 increased the amount of dsDNA. We also tested random primers versus virus-specific primers in
310 the RCA reaction and found that random primers resulted in equal or higher levels of dsDNA
311 depending on the concentration of virus-specific primers, supporting findings by Dean *et al.*,
312 2001 (Supp. 1B).

313 To further increase the amount of dsDNA after RCA, we tested a modified form of phi29
314 DNA polymerase marketed as EquiPhi29, which has been reported to increase dsDNA output up
315 to 7-fold (Povilaitis et al., 2016). Total DNA was amplified in RCA reactions containing the
316 phi29 or the EquiPhi29 DNA polymerase. The amounts of dsDNA and ssDNA were measured
317 using the Qubit dsDNA HS assay kit and ssDNA assay kit as described in the methods. The
318 EquiPhi29 DNA polymerase yielded ~175-fold more ssDNA and 5-fold more dsDNA than the
319 phi29 DNA polymerase (Fig. 2C). Even though most of the increase in RCA products was
320 ssDNA, the 5-fold increase in dsDNA when combined with the 2-fold increase after the DNA
321 end-repair reaction resulted in more dsDNA available for ligation to library adapters.

322 During the process of testing different parameters, we noticed that RCA was highly
323 variable in DNA output. Four total DNA samples (1-4) isolated from symptomatic Kibandameno
324 leaves were amplified using RCA with EquiPhi29 and phi29 (Fig. 2D). This process was

325 repeated twice generating two technical replicates for each sample. The dsDNA concentration of
326 each sample and its technical replicate were measured using a Qubit dsDNA HS assay kit (Fig.
327 2D). Two of the four samples had technical replicates that differed in concentration by more than
328 10-fold (samples 1 and 4, Fig. 2D). Variability was also observed with phi29 (data not shown).
329 This problem was overcome by standardizing the amount the RCA product (2 μ L at 5 ng/ μ L, i.e.
330 10 ng DNA) used for library construction. If the yield of the RCA product from a given reaction
331 was insufficient for dilution to 5 ng/ μ L, that reaction was repeated.

332 *3.2. Total DNA versus size-selected DNA from leaf tissue*

333 Our goal was to develop a protocol where we could achieve sufficient read coverage to
334 detect viral variants across a complete transmission cycle. Based on the coverage reported by
335 Juárez et al., (2019) to detect viral variants and our calculations, we set 1000X coverage (ca.
336 20,000 150-bp reads/genome component) as our minimum coverage for detecting low viral
337 variants occurring at 3% and 1% frequency in the population (for 30 and 10 variant-supporting
338 reads, respectively). To achieve this goal, we examined additional methods to increase the
339 number of reads mapping to the viral genomes.

340 The DNA-A and DNA-B components of the ACMV and EACMCV genomes are ca. 2.8
341 Kb in size and, as such, are much smaller than plant genomic DNA. Hence, we asked if size
342 selection would increase the number of reads mapping to the viral genomes. We used the
343 BluePippin system to select for DNA < 6 Kb from total DNA samples isolated from
344 symptomatic Kibandameno leaves from two plants. The starting total DNA and the size-selected
345 DNA were amplified using the optimized RCA protocol. Libraries were generated using our
346 experimental pipeline (Fig. 1) and sequenced on the Illumina MiSeq platform. We sequenced
347 two technical replicates for each sample. After processing and mapping the resulting reads to the

348 ACMV and EACMCV reference genomes and the cassava reference, we found that nearly all of
349 the reads mapped to the viral reference genomes when the DNA was size-selected (blue)
350 compared to half of the reads for total DNA (light grey) (Fig. 3A). Generally, the use of size-
351 selected DNA (blue) for library construction increased the number of reads mapping to each of
352 the viral genome components compared to total DNA (grey), improving mapping by ~ 2-fold
353 (Fig. 3B). Nearly all of the reads that did not map to the viral genome components but mapped to
354 the cassava reference genome (Fig. 3A), indicating that size selection is an effective method for
355 separating viral DNA from host DNA. Size-selection was effective in increasing the viral read
356 counts in samples with both high and low levels of CMB genome components (Supp. 2A and B),
357 indicating that size selection can produce reliable results over a 10-fold range. However, the
358 virus titer before RCA does not reflect the resulting coverage and read count after RCA and NGS
359 sequencing (compare Supp. 2A and 3). This result also underscores the variability of RCA (Fig.
360 3C).

361 Size-selection also improved read coverage compared to coverage from total DNA for
362 both cassava (Supp. 3A-D) and whitefly (Supp. 3E-H) samples. The coverage was relatively
363 even across the genomes. Dips in coverage were seen in some profiles at the 3' ends of the
364 convergent transcription units (at the ends of the converging arrows) and in the 5' intergenic
365 regions (at the ends of the linear maps), but the coverage was still above 1000X.

366 After observing that size selection increased viral read count, we evaluated whether
367 virion DNA containing only packaged viral ssDNA would also increase viral read count (Fig.
368 3C). The resulting average read count (yellow) from three bioreplicates was highly variable and
369 500-1000 fold lower than the average read counts of total or size-selected DNA (Fig. 3B).

370 *3.3. Sequencing viral DNA from viruliferous whiteflies*

371 We used the ToMoV-tomato-whitefly pathosystem to assess if our optimized protocols
372 could be applied to another begomovirus. As with the cassava pathosystem, sequencing libraries
373 generated from total DNA (grey) isolated from a ToMoV-infected tomato plant (source plant)
374 resulted in fewer viral reads than libraries constructed from size-selected DNA (blue; Fig. 4A).
375 The average reads from virions for ToMoV DNA-A were similar to total DNA and less than
376 size-selected DNA, while the average read counts from virions for ToMoV DNA-B were lower
377 than total and size-selected DNA (Fig. 4A). This underscores the variability of sequencing from
378 virions. Examining the average percent reads mapping to ToMoV versus host DNA (source
379 plant), we found that over 99% of the reads mapped to ToMoV when using size-selected DNA
380 (Fig. 4B). Using total DNA and virion DNA resulted in ~30% and ~25% of the reads mapping to
381 ToMoV, respectively. This confirms the advantage of size-selection seen in the cassava
382 pathosystem (Fig. 3).

383 We also sequenced DNA from groups of 5 whiteflies that had acquired ToMoV virions
384 from the sampled source plant. We sequenced groups of whiteflies with low, medium, and high
385 viral loads as determined by qPCR. The low, medium and high load groups were split into total
386 DNA and size-selected DNA treatments and sequenced. Read counts for the total DNA samples
387 increased with viral load (Fig. 4C) and were lower than read counts for total DNA from infected
388 tomato (Fig. 4A). In contrast, read counts for size-selected DNA samples (blue) were similar
389 across the different viral load levels and were 20-fold higher than the total whitefly DNA
390 samples (grey; Fig. 4B). Overall, the percent of reads mapping to the viral genome was greater
391 than 90% for size-selected DNA (Fig. 4D). Size-selection also resulted in a 6-fold increase in
392 read coverage compared to total DNA (Supp 3E-H).

393 We also attempted to sequence virion DNA from ToMoV-infected tomato (Fig. 4A) and
394 three sets of five viruliferous whiteflies (Fig. 4E). The resulting reads for virion DNA (yellow)
395 from the ToMoV-infected tomato were variable between the DNA-A and DNA-B components,
396 and the average read count was 100-fold lower than total DNA (grey) and size-selected DNA
397 (blue). The resulting reads from virion DNA (yellow) from viruliferous whiteflies were less than
398 1000X coverage and varied between bioreplicates, indicating that sequencing from whitefly
399 virions was not reproducible. It is also important to note that sequencing total and size-selected
400 DNA provides information both viral ssDNA and dsDNA while sequencing virion DNA only
401 yields sequence data about ssDNA.

402 **4. Discussion**

403 In this study, we examined the impact of RCA reaction conditions and size selection on
404 short-read sequencing of begomovirus DNA. Our optimized protocol effectively enriched for
405 viral DNA and produced short-read sequence data from enhanced RCA reaction products across
406 a range of viral loads. Using size-selected DNA, over 90% of the reads mapped to the viral
407 reference genomes from cassava (Fig. 3A). Without size selection, ca. 50% of the reads mapped
408 to the viral genome and the remaining 50% mapped to the host genome (Fig. 3A). Given that
409 both approaches can yield high numbers of viral reads, using total DNA for library construction
410 may be the preferred approach for plant samples when access to size fractionation
411 instrumentation, cost, and/or time are limiting. A recent study similarly concluded that size
412 selection can be beneficial for long-read sequencing of begomoviruses, which has advantages for
413 certain applications (Mehta et al. 2019).

414 Previous studies using RCA and short-read sequencing to characterize begomovirus
415 sequences reported less than 50% of the reads mapping to the viral genome and even lower

416 proportions when the RCA step was omitted (Kathurima, 2016). Sequencing of CMBs yielded
417 mapping ranges of 0.87-6.9% from cassava (Kathurima, 2016) and 0.9-29.6% from *Nicotiana*
418 *benthamiana* (Chen et al., 2019). Low read counts following RCA were also observed for tomato
419 leaf curl New Delhi virus (0.47-1.05 %; (Juárez et al., 2019) and euphorbia yellow mosaic virus
420 (1.24-1.35%; (Richter et al., 2016). Our improved sequence method resulted in at least 7X higher
421 viral mapping reads for CMBs and 66X higher than other begomoviruses reported in the
422 literature.

423 Our methods can be used to characterize viral DNA sequences in whiteflies (Fig. 4B).
424 Size selection had a much larger impact on viral read counts from viruliferous whitefly samples
425 compared to infected plant samples, resulting in at least 1000X coverage across a range of viral
426 loads. Sequencing total DNA from whiteflies results in 1000X read coverage when viral loads
427 are high but not when they are low. In contrast, a wide range of viral loads was successfully
428 sequenced using both size selection and total DNA approaches from infected plants. Sequencing
429 viral sequences from whitefly virions produced average viral read counts that were highly
430 variable and 100-fold lower than that average read count from total and size-selected DNA (Fig.
431 4C and E). Currently, most NGS sequencing from whiteflies is from enriched virions for vector-
432 enabled metagenomic surveys (Ng et al., 2011; Rosario et al., 2015). We found that sequencing
433 from virions did not provide enough coverage and read depth for studying virus population
434 diversion in pools of 5 whiteflies.

435 We found that RCA was not efficient at converting ssDNA to dsDNA (Fig. 2A). Even
436 with an improved DNA polymerase (EquiPhi29), most of the RCA product is ssDNA. RCA has
437 a preference to amplify ssDNA as linear, concatenated copies with low conversion of ssDNA to
438 dsDNA, which can be improved by increasing the reaction incubation time to over 25 hours

439 (Ducani, Bernardinelli, and Högberg, 2014; Zhang and Tanner, 2017). However, the Equiph29
440 DNA polymerase resulted in ~5-fold more dsDNA after RCA with only a 2-hour incubation
441 time, increasing the amount of viral DNA template available for library construction in much
442 shorter reaction time (Fig. 2C). We also found that the amount of DNA produced by RCA varied
443 even when reactions contained the same input DNA (Fig. 2D). Thus, it is important to
444 standardize the amount of RCA product used for library preparation.

445 **6. Conclusion**

446 We established a short-read sequencing protocol for ssDNA viruses that provides high
447 numbers of reads that map to viral reference genomes. The method can use total DNA or size-
448 selected DNA from leaf and whitefly samples that are amplified using random primers and a
449 modified phi29 DNA polymerase prior to library construction. We also found that RCA is
450 variable and is strongly biased towards the amplification of ssDNA products. We cannot fully
451 explain the poor conversion rate to dsDNA, and further studies are needed to understand how
452 ssDNA is converted to dsDNA during RCA. In summary, we have developed an improved tool
453 for begomovirus DNA diagnostics and studying begomovirus population dynamics. Our
454 approach should be applicable to other CRESS DNA viruses.

455 **Author contributions**

456 CDA - Experimental design, execution, and analysis; manuscript preparation

457 JSH - Experimental design and data analysis; manuscript preparation

458 AED - Experimental design, execution, and data analysis; manuscript preparation

459 DOD - Experimental execution

460 IC - Data analysis

461 SD- Experimental design

462 LHB - Experimental design and manuscript preparation

463 **Funding**

464 This work was supported by the National Science Foundation grant OISE-1545553 to LHB and
465 SD.

466

467 **Declaration of Competing Interest**

468 The authors declare that they have no conflict of interest.

469 **Acknowledgments**

470 We thank Mary Beth Dallas for her help growing cassava and tomato plants. The Office of
471 Advanced Research Computing (OARC) at Rutgers, The State University of New Jersey,
472 provided access to and maintenance of the Amarel cluster. We also thank the NC State
473 University Genome Science Laboratory for their support.

474 **References**

475 Abouzeid, A.M., Polston, J.E. and Hiebert, E., 1992. The nucleotide sequence of tomato mottle
476 virus, a new geminivirus isolated from tomatoes in Florida. *J Gen Virol* 73, 3225-3229.

477 Acevedo, A. and Andino, R., 2014. Library preparation for highly accurate population
478 sequencing of RNA viruses. *Nat Protoc* 9, 1760-1769.

479 Acevedo, A., Brodsky, L. and Andino, R., 2014. Mutational and fitness landscapes of an RNA
480 virus revealed through population sequencing. *Nature* 505, 686-90.

481 Ariyo, O.A., Atiri, G.I., Dixon, A.G.O. and Winter, S., 2006. The use of biolistic inoculation of
482 cassava mosaic begomoviruses in screening cassava for resistance to cassava mosaic
483 disease. *J of Virol Methods* 137, 43-50.

484 Bottcher, B., Unseld, S., Ceulemans, H., Russell, R.B. and Jeske, H., 2004. Geminata structures
485 of African cassava mosaic virus. *J Virol* 78, 6758-65.

- 486 Bredeson, J.V., Lyons, J.B., Prochnik, S.E., Wu, G.A., Ha, C.M., Edsinger-Gonzales, E.,
487 Grimwood, J., Schmutz, J., Rabbi, I.Y., Egesi, C., Nauluvula, P., Lebot, V., Ndunguru, J.,
488 Mkamilo, G., Bart, R.S., Setter, T.L., Gleadow, R.M., Kulakow, P., Ferguson, M.E.,
489 Rounsley, S. and Rokhsar, D.S., 2016. Sequencing wild and cultivated cassava and
490 related species reveals extensive interspecific hybridization and genetic diversity. *Nature*
491 *Biotechnol* 34, 562-570.
- 492 Cabrera-Ponce, J., López, L., Assad-Garcia, N., Medina-Arevalo, C., Bailey, A. and Herrera-
493 Estrella, L., 1997. An efficient particle bombardment system for the genetic
494 transformation of asparagus (*Asparagus officinalis* L.). *Plant Cell Rep* 16, 255-260.
- 495 Chen, K., Khatabi, B. and Fondong, V.N., 2019. The AC4 Protein of a Cassava Geminivirus Is
496 Required for Virus Infection. *Mol Plant-Microbe Interactions* 32, 865-875.
- 497 Chen, W., Hasegawa, D.K., Kaur, N., Kliot, A., Pinheiro, P.V., Luan, J., Stensmyr, M.C., Zheng,
498 Y., Liu, W., Sun, H., Xu, Y., Luo, Y., Kruse, A., Yang, X., Kontsedalov, S., Lebedev, G.,
499 Fisher, T.W., Nelson, D.R., Hunter, W.B., Brown, J.K., Jander, G., Cilia, M., Douglas,
500 A.E., Ghanim, M., Simmons, A.M., Wintermantel, W.M., Ling, K.-S. and Fei, Z., 2016.
501 The draft genome of whitefly *Bemisia tabaci* MEAM1, a global crop pest, provides novel
502 insights into virus transmission, host adaptation, and insecticide resistance. *BMC Biology*
503 14, 110.
- 504 Chowda Reddy, R.V., Dong, W., Njock, T., Rey, M.E.C. and Fondong, V.N., 2012. Molecular
505 interaction between two cassava geminiviruses exhibiting cross-protection. *Virus Res*
506 163, 169-177.
- 507 Claverie, S., Ouattara, A., Hoareau, M., Filloux, D., Varsani, A., Roumagnac, P., Martin, D.P.,
508 Lett, J.-M. and Lefeuvre, P., 2019. Exploring the diversity of Poaceae-infecting
509 mastreviruses on Reunion Island using a viral metagenomics-based approach. *Sci Rep* 9,
510 12716.
- 511 de la Iglesia, F. and Elena, S.F., 2007. Fitness Declines in Tobacco Etch
512 Virus upon Serial Bottleneck Transfers. *J of Vir* 81, 4941.

- 513 Dean, F.B., Hosono, S., Fang, L., Wu, X., Faruqi, A.F., Bray-Ward, P., Sun, Z., Zong, Q., Du,
514 Y., Du, J., Driscoll, M., Song, W., Kingsmore, S.F., Egholm, M. and Lasken, R.S., 2002.
515 Comprehensive human genome amplification using multiple displacement amplification.
516 Proc Natl Acad Sci U S A 99, 5261-6.
- 517 Dean, F.B., Nelson, J.R., Giesler, T.L. and Lasken, R.S., 2001. Rapid amplification of plasmid
518 and phage DNA using Phi 29 DNA polymerase and multiply-primed rolling circle
519 amplification. Genome Res 11, 1095-1099.
- 520 Dickins, B. and Nekrutenko, A., 2009. High-resolution mapping of evolutionary trajectories in a
521 phage. Genome Biol Evol 1, 294-307.
- 522 Ducani, C., Bernardinelli, G. and Högberg, B., 2014. Rolling circle replication requires single-
523 stranded DNA binding protein to avoid termination and production of double-stranded
524 DNA. Nucleic Acids Res 42, 10596-604.
- 525 Duffy, S. and Holmes, E.C., 2009. Validation of high rates of nucleotide substitution in
526 geminiviruses: phylogenetic evidence from East African cassava mosaic viruses. J Gen
527 Virol 90, 1539-47.
- 528 Duffy, S., Shackelton, L. A. & Holmes, E. C., 2008. Rates of evolutionary change in viruses:
529 patterns and determinants. Nat. Rev. Genet 9, 267-276.
- 530 Elena, S.F., Fraile, A. and García-Arenal, F. 2014. Chapter Three - Evolution and Emergence of
531 Plant Viruses. In: Maramorosch, K. and Murphy, F.A. (Eds), Adv Virus Res, Academic
532 Press, pp. 161-191.
- 533 Elena, S.F. and Sanjuán, R., 2007. Virus Evolution: Insights from an Experimental Approach.
534 Annu Rev Ecol, Evol, and S 38, 27-52.
- 535 FAOSTAT, 2016. FAO Statistical Databases. <http://faostat.fao.org/> Food and Agriculture
536 Organization (FAO) of the United Nations, Rome, Italy
- 537 Fondong, V.N. and Chen, K., 2011. Genetic variability of East African cassava mosaic
538 Cameroon virus under field and controlled environment conditions. Virol 413, 275-282.

- 539 Fondong, V.N., Pita, J.S., Rey, M.E., de Kochko, A., Beachy, R.N. and Fauquet, C.M., 2000.
540 Evidence of synergism between African cassava mosaic virus and a new double-
541 recombinant geminivirus infecting cassava in Cameroon. *J Gen Virol* 81, 287-97.
- 542 Giardine, B., Riemer, C., Hardison, R.C., Burhans, R., Elnitski, L., Shah, P., Zhang, Y.,
543 Blankenberg, D., Albert, I., Taylor, J., Miller, W., Kent, W.J. and Nekrutenko, A., 2005.
544 Galaxy: a platform for interactive large-scale genome analysis. *Genome Res* 15, 1451-5.
- 545 Hanley-Bowdoin, L., Bejarano, E.R., Robertson, D. and Mansoor, S., 2013. Geminiviruses:
546 masters at redirecting and reprogramming plant processes. *Nat Rev Microbiol* 11, 777-
547 788.
- 548 Hosmani, P.S., Flores-Gonzalez, M., van de Geest, H., Maumus, F., Bakker, L.V., Schijlen, E.,
549 van Haarst, J., Cordewener, J., Sanchez-Perez, G., Peters, S., Fei, Z., Giovannoni, J.J.,
550 Mueller, L.A. and Saha, S., 2019. An improved de novo assembly and annotation of the
551 tomato reference genome using single-molecule sequencing, Hi-C proximity ligation and
552 optical maps. *bioRxiv*, 767764.
- 553 Hoyer, J.S., Fondong, V.N., Dallas, M.M., Aimone, C.D., Deppong, D.O., Duffy, S. and Hanley-
554 Bowdoin, L., 2020. Deeply sequenced infectious clones of key cassava begomovirus
555 isolates from Cameroon. *bioRxiv*, 2020.08.10.244335.
- 556 Idris, A., Al-Saleh, M., Piatek, M.J., Al-Shahwan, I., Ali, S. and Brown, J.K., 2014. Viral
557 metagenomics: analysis of begomoviruses by illumina high-throughput sequencing.
558 *Viruses* 6, 1219-36.
- 559 Inoue-Nagata, A.K., Albuquerque, L.C., Rocha, W.B. and Nagata, T., 2004. A simple method for
560 cloning the complete begomovirus genome using the bacteriophage ϕ 29 DNA
561 polymerase. *J Virol Methods* 116, 209-211.
- 562 Jeske, H., 2018. Barcoding of Plant Viruses with Circular Single-Stranded DNA Based on
563 Rolling Circle Amplification. *Viruses* 10.

- 564 Juárez, M., Rabadán, M.P., Martínez, L.D., Tayahi, M., Grande-Pérez, A. and Gómez, P., 2019.
565 Natural Hosts and Genetic Diversity of the Emerging Tomato Leaf Curl New Delhi Virus
566 in Spain. *Front in Microbiol* 10.
- 567 Kathurima, T.M., Ateka, E. M., Nyende, A. B., & Holton, T. A., 2016. The rolling circle
568 amplification and next generation sequencing approaches reveal genome wide diversity
569 of Kenyan cassava mosaic geminivirus. *Afr J Biotechnol* 15, 2045-2052.
- 570 Lasken, R.S. and Stockwell, T.B., 2007. Mechanism of chimera formation during the Multiple
571 Displacement Amplification reaction. *BMC Biotechnol* 7, 19-19.
- 572 Lefeuvre, P., Harkins, G.W., Lett, J.M., Briddon, R.W., Chase, M.W., Moury, B. and Martin,
573 D.P., 2011. Evolutionary time-scale of the begomoviruses: evidence from integrated
574 sequences in the Nicotiana genome. *PLoS One* 6, e19193.
- 575 Lefeuvre, P. and Moriones, E., 2015. Recombination as a motor of host switches and virus
576 emergence: geminiviruses as case studies. *Curr Opin Virology* 10, 14-19.
- 577 Legg, J.P., Shirima, R., Tajebe, L.S., Guastella, D., Boniface, S., Jeremiah, S., Nsami, E.,
578 Chikoti, P. and Rapisarda, C., 2014. Biology and management of Bemisia whitefly
579 vectors of cassava virus pandemics in Africa. *Pest manag sci*.
- 580 Li, H., 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
581 arXiv, 1303.3997.
- 582 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,
583 Durbin, R. and Genome Project Data Processing, S., 2009. The Sequence Alignment/Map
584 format and SAMtools. *Bioinfo* 25, 2078-2079.
- 585 Lima, A.T.M., Silva, J.C.F., Silva, F.N., Castillo-Urquiza, G.P., Silva, F.F., Seah, Y.M.,
586 Mizubuti, E.S.G., Duffy, S. and Zerbini, F.M., 2017. The diversification of begomovirus
587 populations is predominantly driven by mutational dynamics. *Virus evolution* 3, vex005.

- 588 Lou, D.I., Hussmann, J.A., McBee, R.M., Acevedo, A., Andino, R., Press, W.H. and Sawyer,
589 S.L., 2013. High-throughput DNA sequencing errors are reduced by orders of magnitude
590 using circle sequencing. *Proc Natl Acad Sci U S A* 110, 19872-19877.
- 591 Martin, M., 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
592 2011 17, 3.
- 593 Mehta, D., Hirsch-Hoffmann, M., Were, M., Patrignani, A., Zaidi, S.S.-E.A., Were, H.,
594 Gruissem, W. and Vanderschuren, H., 2019. A new full-length circular DNA sequencing
595 method for viral-sized genomes reveals that RNAi transgenic plants provoke a shift in
596 geminivirus populations in the field. *Nucleic Acids Res* 47, e9-e9.
- 597 Moriones, E. and Navas-Castillo, J., 2000. Tomato yellow leaf curl virus, an emerging virus
598 complex causing epidemics worldwide. *Virus Res* 71, 123-34.
- 599 Ndunguru, J., De León, L., Doyle, C.D., Sseruwagi, P., Plata, G., Legg, J.P., Thompson, G.,
600 Tohme, J., Aveling, T., Ascencio-Ibáñez, J.T. and Hanley-Bowdoin, L., 2016. Two Novel
601 DNAs That Enhance Symptoms and Overcome CMD2 Resistance to Cassava Mosaic
602 Disease. *J Virol* 90, 4160-4173.
- 603 Ng, T.F.F., Duffy, S., Polston, J.E., Bixby, E., Vallad, G.E. and Breitbart, M., 2011. Exploring
604 the Diversity of Plant DNA Viruses and Their Satellites Using Vector-Enabled
605 Metagenomics on Whiteflies. *PLOS ONE* 6, e19050.
- 606 Polston, J.E. and Anderson, P.K., 1997. The Emergence Of Whitefly-Transmitted Geminiviruses
607 in Tomato in the Western Hemisphere. *Plant Disease* 81, 1358-1369.
- 608 Povilaitis, T., Alzbutas, G., Sukackaite, R., Siurkus, J. and Skirgaila, R., 2016. In vitro evolution
609 of phi29 DNA polymerase using isothermal compartmentalized self replication
610 technique. *Protein engineering, design & selection : PEDS* 29, 617-628.
- 611 Rajabu, C.A., Kennedy, G.G., Ndunguru, J., Ateka, E.M., Tairo, F., Hanley-Bowdoin, L. and
612 Ascencio-Ibáñez, J.T., 2018. Lanai: A small, fast growing tomato variety is an excellent
613 model system for studying geminiviruses. *J Virol Methods* 256, 89-99.

- 614 Reyes, M.I., Nash, T.E., Dallas, M.M., Ascencio-Ibáñez, J.T. and Hanley-Bowdoin, L., 2013.
615 Peptide Aptamers That Bind to Geminivirus Replication Proteins Confer a Resistance
616 Phenotype to *Tomato Yellow Leaf Curl Virus* and *Tomato Mottle*
617 *Virus* Infection in Tomato. *J Virol* 87, 9691-9706.
- 618 Richter, K.S., Götz, M., Winter, S. and Jeske, H., 2016. The contribution of translesion synthesis
619 polymerases on geminiviral replication. *Virology* 488, 137-148.
- 620 Rocha, C.S., Castillo-Urquiza, G.P., Lima, A.T.M., Silva, F.N., Xavier, C.A.D., Hora-Júnior,
621 B.T., Beserra-Júnior, J.E.A., Malta, A.W.O., Martin, D.P., Varsani, A., Alfenas-Zerbini,
622 P., Mizubuti, E.S.G. and Zerbini, F.M., 2013. Brazilian Begomovirus Populations Are
623 Highly Recombinant, Rapidly Evolving, and Segregated Based on Geographical
624 Location. *J Virol* 87, 5784-5799.
- 625 Rosario, K., Seah, Y.M., Marr, C., Varsani, A., Kraberger, S., Stainton, D., Moriones, E.,
626 Polston, J.E., Duffy, S. and Breitbart, M., 2015. Vector-Enabled Metagenomic (VEM)
627 Surveys Using Whiteflies (Aleyrodidae) Reveal Novel Begomovirus Species in the New
628 and Old Worlds. *Viruses* 7, 5553-5570.
- 629 Ruark-Seward, C.L., Bonville, B., Kennedy, G. and Rasmussen, D.A., 2020. Evolutionary
630 dynamics of *Tomato spotted wilt virus*; within and between
631 alternate plant hosts and thrips. *bioRxiv*, 2020.01.13.904250.
- 632 Safari, M. and Roossinck, M.J., 2014. How does the genome structure and lifestyle of a virus
633 affect its population variation? *Curr Opin Virol* 9, 39-44.
- 634 Sanjuán, R., Nebot, M.R., Chirico, N., Mansky, L.M. and Belshaw, R., 2010. Viral Mutation
635 Rates. *J Virol* 84, 9733.
- 636 Scientific, T.F. DNA and RNA Molecular Weights and Conversions.
- 637 Sipos, R., Székely, A., Révész, S. and Márialigeti, K. 2010. Addressing PCR Biases in
638 Environmental Microbiology Studies. In: Cummings, S.P. (Ed), *Bioremediation:*
639 *Methods and Protocols*, Humana Press, Totowa, NJ, pp. 37-58.

- 640 Stanley, J. and Gay, M.R., 1983. Nucleotide sequence of of cassava latent virus DNA. Nature
641 301, 2660-2662.
- 642 Thresh J. M. , D.F.a.G.W.O.-N., 1997. African cassava mosaic disease: the magnitude of the
643 problem. Afr J Root Tuber Crops 2, 13-19.
- 644 Uzokwe, V.N.E., Mlay, D.P., Masunga, H.R., Kanju, E., Odeh, I.O.A. and Onyeka, J., 2016.
645 Combating viral mosaic disease of cassava in the Lake Zone of Tanzania by
646 intercropping with legumes. Crop Prot 84, 69-80.
- 647 Wang, X., Li, Y., Ni, T., Xie, X., Zhu, J. and Zheng, Z.M., 2014. Genome sequencing accuracy
648 by RCA-seq versus long PCR template cloning and sequencing in identification of human
649 papillomavirus type 58. Cell biosci 4, 5.
- 650 Zhang, Y. and Tanner, N.A., 2017. Isothermal Amplification of Long, Discrete DNA Fragments
651 Facilitated by Single-Stranded Binding Protein. Sci Rep 7, 8497.
- 652 Zhao, L., Rosario, K., Breitbart, M. and Duffy, S. 2019. Chapter Three - Eukaryotic Circular
653 Rep-Encoding Single-Stranded DNA (CRESS DNA) Viruses: Ubiquitous Viruses With
654 Small Genomes and a Diverse Host Range. In: Kielian, M., Mettenleiter, T.C. and
655 Roossinck, M.J. (Eds), Adv Virus Res, AP, pp. 71-133.

656

657 **Table 1.** qPCR Primer Sequences

Primers	Sequence	qPCR target
P3P-AA2F	TCTGCAATCCAGGACCTACC	ACMV DNA-A
P3P-AA2R+4R	GGCTCGCTTCTTGAATTGTC	ACMV DNA-A
ACMVBdiv4	ATTGAGCACCAGGCGATAT	ACMV DNA-B
ACMVBfor1	CACATAGAGGCAGTAGCCATAAA	ACMV DNA-B
EACMVQ1	GTACCATGCGTCGTTTGAATA	EACMCV DNA-A
EACMVQ2	GCAAGTCCCAGAGGAAATAGA	EACMCV DNA-A
EACMVBREV4	GCATCGACTGTGATCGCATAC	EACMCV DNA-B
EACMVBfor1.2	CCAAGGATACACAAAAGATTGC	EACMCV DNA-B
ToMoVA6F	TCAGGTTGTGGTTGAACCGT	ToMoV DNA-A
ToMoVA6R	TTAGACTGTGCGGGACATGG	ToMoV DNA-A
ToMoVB4F	CGACGAGCTATTTGGTGCA	ToMoV DNA-B
ToMoVB4R	TCTCAACTGAGAGCACTCGC	ToMoV DNA-B

658

660 **Figure Legends**

661 *Fig 1. Workflow of viral DNA sequencing for short-read sequencing platforms*

662 *Fig 2. Improvements to RCA reactions*

663 (A) The amount of total DNA (-RCA) and total DNA amplified with RCA (+RCA) treated
664 (+MB) or not treated (-MB) with Mung Bean nuclease (MB). (B) The amount of DNA (ng) in an
665 RCA reaction before -end repair and after +end repair. (C) Amount of ssDNA and dsDNA after
666 RCA amplification of total DNA with either Phi29 or EquiPhi29 DNA polymerases. (D) The
667 concentration of the total DNA of technical replicates from 4 cassava leaf DNA samples
668 amplified using RCA.

669 *Fig 3. Size selection increases viral DNA read counts for cassava samples.*

670 (A) The average percent reads mapping to viral DNA (ACMV and EACMCV) and host DNA,
671 and unmapped reads for total DNA (grey) and size-selected DNA (blue) samples. (B) The
672 average number of reads corresponding to ACMV DNA-A (Hosmani et al.), ACMV DNA-B
673 (AB), EACMCV DNA-A (EA), and EACMCV DNA-B for total (grey) and size-selected (blue)
674 DNA samples. (C) The average number of reads corresponding to ACMV DNA-A (Hosmani et
675 al.), ACMV DNA-B (AB) for virion DNA (yellow). The bars in A, B correspond to 2 standard
676 errors from two bio-samples and two technical reps each. The bars in C correspond to 2 standard
677 errors from three bio-samples with two technical replicates each. The asterisks indicate a
678 significant difference (P-value < 0.05) between the viral read counts in size-selected and total
679 DNA in Student's t-tests.

680 *Fig 4. Size selection increases viral DNA read counts for tomato and whitefly samples.*

681 A) The average number of reads corresponding to ToMoV DNA-A (TA) and ToMoV DNA-B
682 (TB) for total (grey), size selected (blue), virion (yellow) DNA samples from the tomato source

683 plant. (B) The average percent reads mapping to viral DNA (ToMoV) and host DNA for total
684 DNA (grey), size-selected DNA (blue), and virion DNA (yellow) samples. (C) The average
685 number of reads corresponding to ToMoV DNA-A (TA) and ToMoV DNA-B (TB) for total
686 (grey) and size-selected (blue) DNA samples from pools of 5 whiteflies with low (30,000-40,000
687 DNA-A copies/ng total DNA), medium (45,000-60,000 DNA-A copies/ng total DNA), and high
688 viral loads (150,000-200,000 DNA-A copies/ng total DNA) virus loads. (D) The average percent
689 reads mapping to viral DNA (ToMoV), host DNA, and whitefly DNA for total DNA (grey), size-
690 selected DNA (blue), and virion DNA (yellow) samples. (E) The average number of reads
691 corresponding to ToMoV DNA-A (TA) and ToMoV DNA-B (TB) for virion DNA (yellow). The
692 bars in A and C correspond to 2 standard errors for three bio-samples with two technical
693 replicates. The asterisks indicate a significant difference (P-value < 0.05) between the viral read
694 counts in size-selected and total DNA in Student's t-tests.

695 *Supp 1. Optimization of RCA conditions*

696 (A) Amount of ssDNA and dsDNA after RCA amplification of total DNA with EquiPhi29
697 polymerase for 2 h and 3 h at 40°C. (B) Amount of dsDNA after RCA amplified with random
698 hexamers (hx) alone or combined with different amounts of virus-specific primers (vs).

699 *Supp 2. Virus levels and virus-mapping read count for cassava biological replicates.*

700 (A) Log viral copy number of ACMV-A (AA0), ACMV-B (AB), EACMCV-A (EA), and
701 EACMCV-B (Bernardo et al.) for biological replicate 1 (green) and biological replicate 2 (blue).
702 (B) Average read count of two technical replicates for biological replicate 1 (green) and 2 (blue)
703 for ACMV-A (Hosmani et al.), ACMV-B (AB), EACMCV-A (EA), and EACMCV-B (Bernardo
704 et al.) without size-selection. (C) Average read count of two technical replicates for biological

705 replicate 1 (green) and 2 (blue) for ACMV-A (Hosmani et al.), ACMV-B (AB), EACMCV-A
706 (EA), and EACMCV-B (Bernardo et al.) with size-selection.
707 *Supp 3. Plots of Illumina read depth across virus segments.*
708 Each row corresponds to one library, from leaf tissue from a cassava plant (A to D) or a single
709 pool of whiteflies (E to H). Pairs of rows correspond to technical duplicate libraries made from
710 total DNA (A and B, E and F) or to size-selected DNA (C and D, G and H). Canonical virus
711 genes are drawn as gray arrows below each set of graphs, left to right for virus sense (AV1, AV2
712 [for ACMV and EACMCV], BV1), and right to left for complementary sense (AC1 to AC4,
713 BC1).
714

Figure 1 Aimone et al

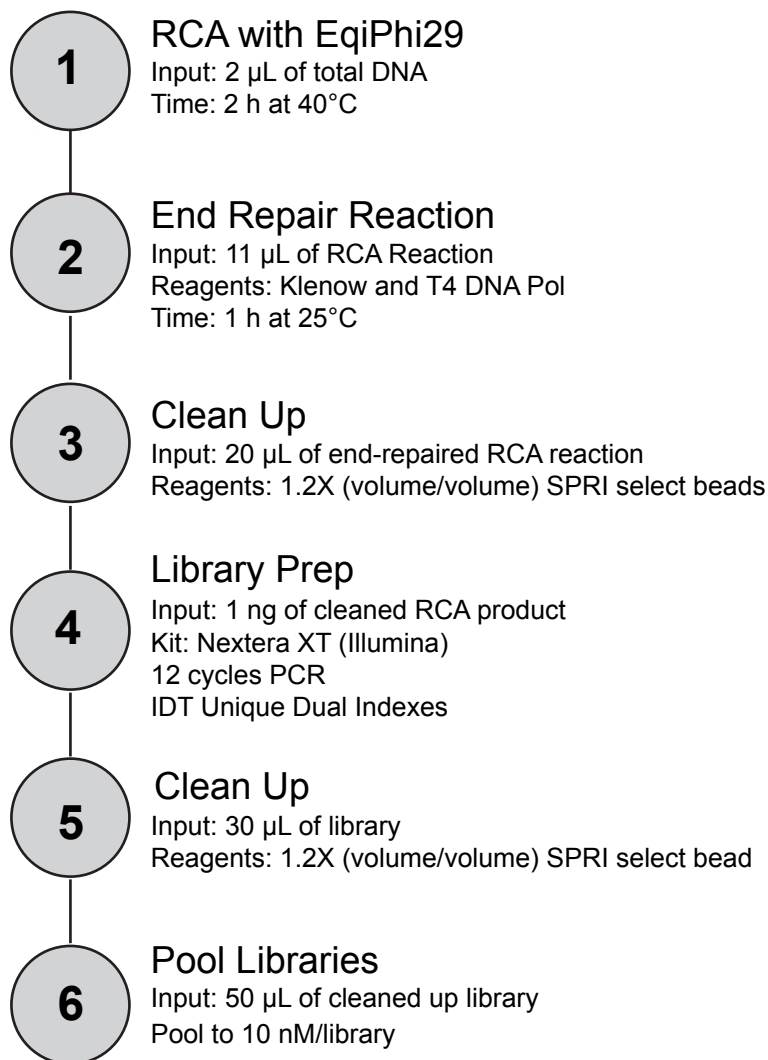


Figure 2 Aimone et al

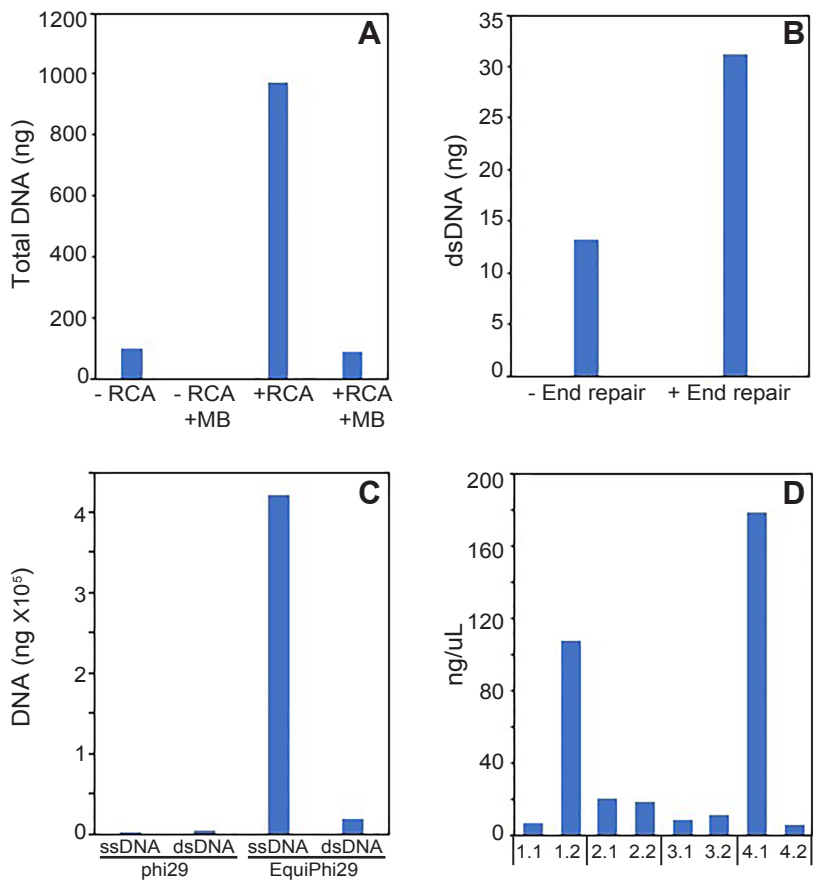


Figure 3 Aimone et al

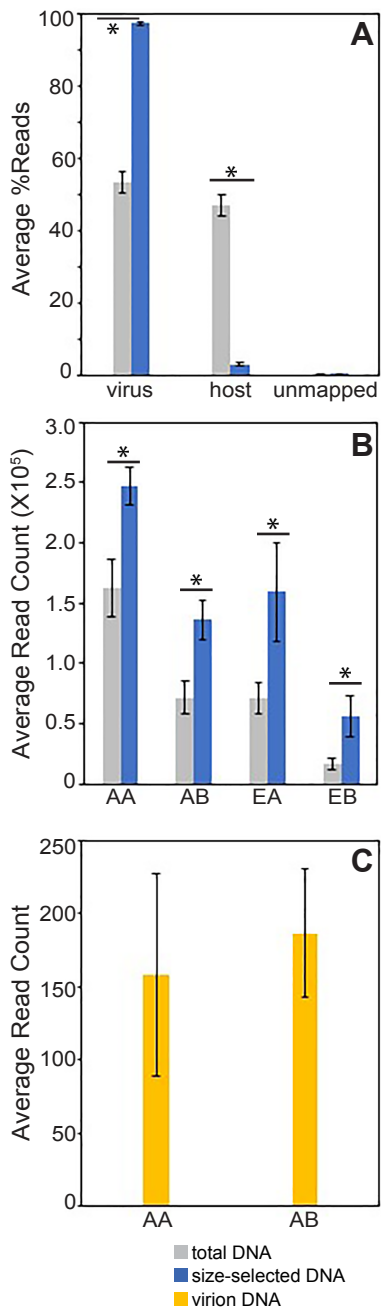


Figure 4 Aimone et al

