

1 Exploring Viral Diversity In A Unique South African Soil Habitat

2

3 *Jane Segobola^{1, 2}, Evelien Adriaenssens², Tsepo Tsekod¹, Konanani Rashamuse¹ and *Don*
4 *Cowan²*

5 *¹Biosciences Unit, Council for Scientific and Industrial Research (CSIR), Pretoria, South*
6 *Africa*

7 *²Centre for Microbial Ecology and Genomics, University of Pretoria, Pretoria, South Africa*

8 *Corresponding Author: Don Cowan, Centre for Microbial Ecology and Genomics,*
9 *Department of Genetics, University of Pretoria, Hatfield 0028, Pretoria, South Africa.*

10 *Tel: +27 (12) 420 5873, don.cowan@up.ac.za*

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25 ABSTRACT

26

27 The Kogelberg Biosphere Reserve in the Cape Floral Kingdom in South Africa is known for
28 its unique plant biodiversity. The potential presence of unique microbial and viral
29 biodiversity associated with this unique plant biodiversity led us to explore the fynbos soil
30 using metaviromic techniques. In this study, metaviromes of a soil community from the
31 Kogelberg Biosphere Reserve has been characterised in detail for the first time. Metaviromic
32 DNA was recovered from soil and sequenced by Next Generation Sequencing. The MetaVir,
33 MG-RAST and VIROME bioinformatics pipelines were used to analyse taxonomic
34 composition, phylogenetic and functional assessments of the sequences. Taxonomic
35 composition revealed members of the order Caudovirales, in particular the family
36 *Siphoviridae*, as prevalent in the soil samples and other compared viromes. Functional
37 analysis and other metaviromes showed a relatively high frequency of phage-related and
38 structural proteins. Phylogenetic analysis of *PolB*, *PolB2*, *terL* and *T7gp17* genes indicated
39 that many viral sequences are closely related to the order Caudovirales, while the remainder
40 were distinct from known isolates. The use of single virome which only includes double
41 stranded DNA viruses limits this study. Novel phage sequences were detected, presenting an
42 opportunity for future studies aimed at targeting novel genetic resources for applied
43 biotechnology.

44

45

46

47

48

49

50

51

52

53

54 **1. INTRODUCTION**

55

56 The Cape Floristic Region situated in the Western Cape province of South Africa is one of
57 five Mediterranean-type ecosystems in the world ¹ and is recognized as one of the world's
58 biodiversity hotspots ². Fynbos (fine bush) is the main vegetation type of this region with the
59 *Proteaceae*, *Ericaceae* and *Restionaceae* families dominating Kogelberg Biosphere Reserve
60 Fynbos vegetation. Within this region, the fynbos comprises approximately 9000 plant
61 species of which 70% are endemic to the region ^{1,3}. Fynbos vegetation types survive on
62 highly heterogeneous, acidic, sandy, well-leached and infertile soils. The fynbos plants also
63 survive invasions by foreign plants ⁴ and seasonal drought conditions ⁵.

64 Microorganisms make up a great proportion of the living population in the biosphere. They
65 provide important ecosystem services in edaphic habitats ⁶ and form complex symbiotic
66 relationships with plants ⁷. Plant-associated microorganism studies have shown high
67 microbial diversity in fynbos soils ², where they play a role in sustaining plant communities ⁸.
68 A study focusing on the linkage between fynbos soil microbial diversity and plant diversity
69 showed the presence of novel taxa and of bacteria specifically associated with the
70 rhizospheric zone ⁹. Studies on ammonium-oxidizing bacteria demonstrated that plant-species
71 specific and monophyletic ammonium oxidizing bacterial clades were present in fynbos soils
72 ¹⁰, where abundance might be driven by the acidic and oligotrophic nature of these soils ¹¹.
73 There is evidence that above-ground floral communities are implicated in shaping microbial
74 communities ^{12,13}, and that some microbial clades show a high level of plant–host specificity
75 ¹⁰. This is consistent with the general concept of the mutualistic relationships between the
76 plants and the microbial communities in fynbos soils ¹⁴.

77 Soil-borne viruses, including phages, are of great importance in edaphic habitats due to their
78 ability to transfer genes from host to host and as a potential cause of microbial mortality
79 (leading to changes in turnover and concentration of nutrients and gases), processes that can
80 profoundly influence the ecology of soil biological communities ¹⁵. Virus diversity associated
81 with fynbos plants from Kogelberg Biosphere Reserve fynbos soil has never been thoroughly
82 investigated ¹⁶. The difficulty of culturing viruses, which are absolutely dependent on a cell
83 host to provide the apparatus for replication and production of progeny virions, presents a
84 barrier to fully accessing viral biodiversity. This is a particular issue in poorly studied

85 habitats, such as fynbos soil, where the true microbial (host) diversity is largely unknown and
86 most microbial phylotypes have never been cultured¹⁷. The biodiversity and ecology of
87 viruses in many soils therefore remain poorly investigated and poorly understood¹⁸.

88 Metaviromic surveys of terrestrial environments such as hot desert soil¹⁸, rice paddy soil
89^{19,20}, Antarctic cold desert soil^{21,22} and hot desert hypolithic niche communities²³ have been
90 reported in recent years and have significantly advanced the field of soil viral ecology^{20,24}.
91 These studies have also facilitated the discovery of novel virus genomes^{20,22,23} and novel
92 viral enzymes²⁵.

93 However, surveys of viral diversity using NGS sequencing techniques in conjunction with
94 metaviromic databases have focused principally on aquatic environments^{26–28}. Studies on
95 taxonomic composition using public metaviromic databases for viral diversity estimations
96 have shown that a majority of environmental virus sequences are unknown¹⁹: ~70% of
97 sequences have no homologs in public databases and are therefore typically labelled “viral
98 dark matter”^{29,30}. Bacteriophages constitute the largest known group of viruses found in both
99 aquatic^{24,31} and soil environments^{32,33}.

100 Here we report the first investigation of virus diversity in a unique soil type (fynbos soil)
101 using metaviromic approaches. The metavirome of Kogelberg Biosphere Reserve fynbos soil
102 was characterised in terms of diversity and functional composition and adds a new level of
103 understanding to the exceptional biodiversity of this habitat.

104

105 **2. RESULTS AND DISCUSSION**

106

107 **2.1. VIRAL MORPHOLOGY**

108

109 Analysis of the morphology of viruses identified in Kogelberg Biosphere Reserve fynbos soil
110 was carried out by transmission electron microscopy (TEM). TEM analysis of the virus
111 preparations showed that the majority of the isolated virus particles were morphologically
112 similar to known virus taxonomic groups³⁴. The isolated virus particles from the fynbos soil
113 were tailed, spherical or filamentous (Supplementary Fig S1 online). Various particles with
114 head-tail morphology, typically belonging to the families *Myoviridae*, *Siphoviridae* or
115 *Podoviridae*, were observed.

116 These results are in a good agreement with previously published findings showing the high
117 dominance of tailed phages in soils from various geographic areas ^{24,33,35}. The undetermined
118 spherical or filamentous morphologies in TEM micrographs could be *bona fide* but
119 uncharacterised viral structures. Spherical particles resembling capsid structures could be
120 members of the *Leviviridae*, *Partitiviridae*, *Chrysoviridae*, *Totiviridae* or *Tectiviridae*
121 families, or small plant viruses ³⁴. Filamentous particles may possibly correspond to the virus
122 structures of the *Inovirus* genus, the members of which contain circular ssDNA within
123 flexible filamentous virions. The presence of spherical types and filamentous type of virus
124 particles was also reported for Delaware soils ³². The aggressive extraction procedure used in
125 the current study may have resulted in a high incidence of phage tail breakage and the
126 generation of tailless phages ³⁶.

127 2.2. METAVIROME ASSEMBLY

128 Assembly of the DNA sequence reads yielded 13,595 contigs larger than 500 bp, with an
129 average length of 2,098 bp, accounting for a total of 28,526,478 bp (Table 1). Two different
130 metagenomics pipelines; MetaVir ³⁷ and VIROME ³⁸, were used for analysis of the contigs,
131 while MG-RAST ³⁹ was used for the analysis of the uploaded reads (Table 2). The MetaVir
132 pipeline predicted 51,274 genes, with 5,338 affiliated contigs (i.e., contigs with at least one
133 BLAST hit) and 7880 unaffiliated contigs (Table 2). MetaVir compares reads/contigs to
134 complete viral genomes from the Refseq database and is specifically designed for the analysis
135 of environmental viral communities ³⁷. The VIROME pipeline ³⁸ predicted 51,242 protein
136 coding regions. Of these, 9555 were assigned as functional proteins, and 31,109 were
137 unassigned (Table 2). Comparisons of functional and taxonomic analysis between Virome
138 and MetaVir indicate that many of the predicted genes were overlapping between the two
139 pipelines with MetaVir on average having a higher predictive potential (Supplementary Table
140 S1 online). The MG-RAST pipeline predicted 2,555,524 protein coding regions. Of these
141 predicted protein features, 119,220 were assigned a functional annotation using protein
142 databases (M5NR) ⁴⁰ and 2,362,076 had no significant similarities to sequences in the protein
143 databases (ORFans). MG-RAST core analysis and annotation depends heavily on the SEED
144 database which is largely comprised of bacterial and archaeal genomes ⁴¹. The majority of the
145 annotated sequences in MG-RAST were mapped to bacterial genomes. This high percentage
146 of bacterial sequences in metaviromes may be due to the presence of unknown prophages in
147 bacterial genomes, phages carrying host genes, relatively large size of bacterial genomes
148 compared to viral genomes and larger size of the microbial genome database which is

149 statistically increasing the chance of matching bacterial sequences. The MG-RAST pipeline
 150 was used to analyse the reads, not the contigs and shows, therefore a higher number of
 151 predicted features, including more partial CDSs ⁴² No rDNA sequences were found with the
 152 MG-RAST and VIROME pipelines, confirming the viral origins of the DNA. The fact that
 153 more than 80% of the hits in this study, consistent with previous viral metagenomics studies
 154 ^{31,43,44}, were assigned as hypothetical proteins derived from unknown viruses suggests the
 155 presence of a substantial pool of novel viruses.

Features	CLC
#Pre-QC Sequence reads	7,019,527
#Pre-QC sequence in base pairs	1,488,462,918
#post-QC average read length	212.05
#contigs	13,595
#contigs/reads in bp	28,526,478 bp

156

157 **Table 1: Next Generation sequencing data analysis.** Representation of the assembly,
 158 annotation, and diversity statistics produced by CLC Genomics

159

Features	MetaVir	MG-RAST	VIROME
#predicted CDS	51,274	2,555,524	51,242
#affiliated CDS*	5,868	119,220	9,555
#ORFans*	45,406	2,362,076	31,109
#rRNAs	NA	0	0
Database used for CDS annotation	RefSeq pfam	virus, GenBank, IMG, KEGG, PATRIC, RefSeq, SEED, SwissProt, TrEMBL, eggNOG, COG,	KEGG, SEED, COG, GO, UniRef100, PHGSEED, MgOI, ACLAME

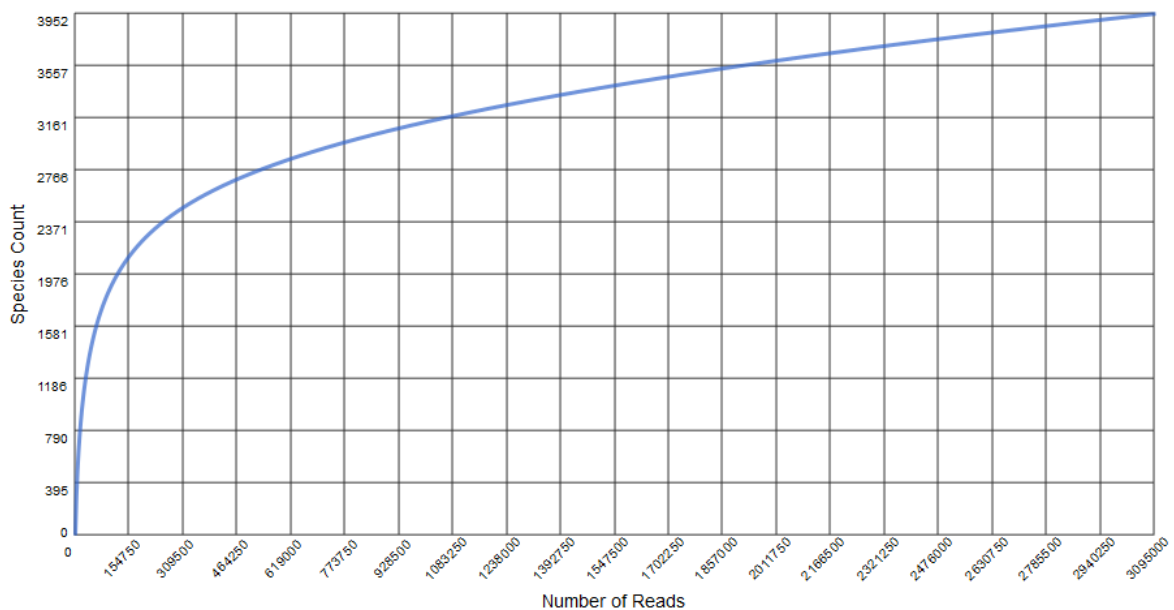
160

161 **Table 2: Comparison of the automated pipelines;** such as MetaVir (contigs), VIROME
 162 (contigs) and MG-RAST (reads), used to characterize the Kogelberg Biosphere Reserve.*
 163 Affiliated CDS are CDS with homologues in at least one of the databases used, while
 164 ORFans are predicted ORFs which have no database homologue.

165

166 2.3. VIRAL DIVERSITY ESTIMATION AND TAXONOMIC COMPOSITION

167 The rarefaction curve computed by MG-RAST showed 3952 species clusters at 90%
 168 sequence identity for the 3,095,000 reads. The curve did not reach an asymptote (Fig 1),
 169 although extrapolation suggested that approximately 78% of the viral diversity was covered
 170 by the metavirome sequence dataset.



171

172 **Figure 1: Rarefaction curve of the Kogelberg Biosphere Reserve fynbos soil**
 173 **metavirome.** Clustering was set at 90% similarity.

174 MetaVir was used for viral taxonomic composition analysis of the contigs. The taxonomic
 175 composition was computed from a BLASTp comparison of the predicted proteins in the
 176 contigs with the Viral Refseq protein database (release of 2016-01-19). The results revealed
 177 that 37.6% of the contigs represented a significant hit (threshold of 50 on the BLAST bit
 178 score). MetaVir identified 18 virus families, in which prokaryotic viruses were the most

179 abundant and dominated by the order *Caudovirales*, consistent with the TEM observations.
 180 The relative abundance ranking of the different families was as follows: tailed bacteriophage
 181 families *Siphoviridae* > *Myoviridae* > *Podoviridae*, followed by the algae-infecting family
 182 *Phycodnaviridae*, the archaeal virus family *Ampullaviridae* and the amoeba-infecting family
 183 *Mimiviridae* (Table 3). Surprisingly, large viruses belonging to the families *Phycodnaviridae*
 184 and *Mimiviridae* were detected, which should have been removed during the filtration
 185 process due to the use of a 0.22- μm filtration step to remove bacterial cells. The identification
 186 of *Mimiviridae* suggests that this filtration process allowed partial mimivirus particles or free-
 187 floating DNA to pass through the membrane. Mimiviruses appear to infect only species
 188 of *Acanthamoeba*, which are ubiquitous in nature and have been isolated from diverse
 189 environments including freshwater lakes, river waters, salt water lakes, sea waters, soils and
 190 the atmosphere^{45,35,46,47}. This suggests the existence of Mimivirus relatives in the KBR soil.

191 Other viral families and unclassified viruses (dsDNA and ssDNA) were found in low
 192 numbers. Putative contamination of *Enterobacteria* phage phiX174 was also detected in our
 193 metavirome sequences. This phage is used for quality control in sample preparation for high-
 194 throughput sequencing. Seven sequences from this dataset are similar to the phiX174 genome
 195 and were thus disregarded in the taxonomic composition as an artefact of sample processing.
 196 Plant viruses were not identified in the dataset, most probably because the majority of plant
 197 viruses are RNA viruses which were not sampled in this study.

Virus Order and family	Hosts	Relative abundance of taxa
<i>Caudovirales</i>		
<i>Myoviridae</i>	Bacteria, Archaea	29
<i>Podoviridae</i>	Bacteria	23
<i>Siphoviridae</i>	Bacteria, Archaea	45
<i>Herpesvirales</i>		
<i>Herpeviridae</i>	Vertebrates	0.04
Virus Family and groups not assigned in to Order		

<i>Phycodnaviridae</i>	Algae	2
<i>Ampullaviridae</i>	Archaea	0.9
<i>Mimiviridae</i>	Amoebae	0.8
<i>Salterprovirus</i>	Archaea	0.7
<i>Tectiviridae</i>	Bacteria, Archaea	0.5
<i>Iridoviridae</i>	Vertebrates (Amphibians, Fishes), Invertebrates	0.1
<i>Marseilleviridae</i>	Amoeba	0.04
<i>Nudiviridae</i>	Arthropods	0.04
<i>Poxviridae</i>	Human, Arthropods, Vertebrates	0.02
<i>Baculoviridae</i>	Invertebrates	0.02
<i>Bicaudaviridae</i>	Archaea	0.02
<i>Turriviridae</i>	Archaea	0.02
<i>Asfarviridae</i>	Swine	0.02
<i>Retroviridae</i>	Vertebrates	0.02
Virus not assigned into Family		
Unclassified dsDNA phages	Bacteria	2
Unclassified dsDNA virus	NA	4
Unclassified ssDNA Viruses	NA	0.07
Unclassified phages	Bacteria	2

198

199 **Table 3: Taxonomic abundance.** Representation of taxonomic abundance of identified viral
200 ORFs BLASTp with threshold of E value 10^{-5} identified by MetaVir.

201 The viral composition of Kogelberg Biosphere Reserve fynbos soil was compared to 12
202 previously published metaviromes from both similar and dissimilar environments, including
203 fresh water⁴⁸, soil and hypolithic niche communities^{22,23}, pond water²⁷ and sea water⁴⁹ (Fig

204 2). A comparative metaviromics approach was used to investigate the assumption that certain
205 environments will select for specific viruses ^{50,51}.

206

SAMPLE TYPES	SOIL		HYPOLITH		DEEP SEA		FRESH WATER						POND
	KBR (s)	AOS (s)	AH	NH	ALOHA (ds)	B47 (ds)	LB (fw)	LP (fw)	57th-1 (fw)	M1 (fw)	M2 (fw)	Far (fw)	SP (p)
Viruses	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
Retro-transcribing viruses	0.000	0.000	0.020	0.000	0.300	0.010	0.000	0.000	0.010	0.030	0.080	0.000	0.000
Satellites	0.000	0.000	0.000	0.000	0.000	0.000	0.020	0.040	0.000	0.000	0.000	0.000	0.000
dsDNA viruses, no RNA stage	97.00	96.00	97.00	85.00	96.00	98.00	95.00	94.00	98.00	98.00	98.00	100.00	100.00
Caudovirales	89.00	80.00	89.00	81.00	52.00	81.00	77.00	79.00	61.00	66.00	65.00	72.00	24.00
Siphoviridae	39.00	37.00	59.00	47.00	16.00	28.00	27.00	28.00	19.00	20.00	20.00	9.00	19.00
Myoviridae	26.00	30.00	18.00	16.00	52.00	27.00	26.00	28.00	33.00	34.00	27.00	56.00	5.00
Podoviridae	20.00	11.00	10.00	5.00	7.00	24.00	21.00	20.00	7.00	11.00	16.00	6.00	0.000
Ampullavirus	1.00	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Mimiviridae	0.800	0.000	1.00	0.000	0.000	1.00	0.080	1.00	0.000	0.000	0.000	0.000	0.000
Salteprovirus	0.700	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.000	0.000
Tectivirus	0.500	0.070	0.030	0.080	0.000	0.010	0.040	0.000	0.020	0.004	0.040	0.000	0.000
Iridoviridae	0.100	0.400	0.090	0.080	0.600	0.200	0.200	0.300	0.500	0.600	0.600	0.000	0.000
Marseilleviridae	0.040	0.300	0.050	0.040	0.090	0.060	0.050	0.040	0.300	0.400	0.300	0.000	0.000
Herpesvirales	0.040	0.100	0.100	0.000	0.500	0.040	0.050	0.040	0.500	0.400	0.300	0.000	0.000
Nudiviridae	0.020	0.070	0.000	0.000	0.000	0.000	0.000	0.000	0.300	0.200	0.100	0.000	0.000
Poxviridae	0.020	0.400	0.060	0.040	2.00	0.200	0.070	0.100	0.000	1.00	1.00	0.000	0.000
Baculoviridae	0.020	0.100	0.080	0.000	0.200	0.070	0.020	0.000	0.200	0.200	0.300	0.000	0.000
Asfivirus	0.020	0.070	0.020	0.000	0.000	0.010	0.000	0.000	0.040	0.020	0.300	0.000	0.000
Polydnaviridae	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.008	0.000	0.000	0.000
Adenoviridae	0.000	0.000	0.020	0.000	0.000	0.000	0.000	0.000	0.000	0.004	0.005	0.000	0.000
White spot syndrome virus	0.000	0.070	0.000	0.000	0.000	0.000	0.020	0.000	0.020	0.020	0.010	0.000	0.000
Fuselloviridae	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.000	0.000
Ascovirus	0.000	0.300	0.030	0.000	0.000	0.100	0.000	0.000	0.800	0.020	0.500	0.000	0.000
Polyomaviridae	0.000	0.000	0.000	0.000	0.000	0.010	0.000	0.000	0.000	0.000	0.009	0.000	0.000
Bicaudaviridae	0.000	0.000	0.050	0.000	0.300	0.000	0.000	0.000	0.070	0.040	0.070	0.000	0.000
Corticovirus	0.000	0.000	0.000	0.000	0.000	0.010	0.020	0.000	0.000	0.000	0.005	0.000	0.000
Ligamenvirales	0.000	0.000	0.030	0.000	0.000	0.040	0.020	0.040	0.030	0.020	0.020	0.000	0.000
Plasmavirus	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.010	0.000	0.000	0.000
dsRNA viruses	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.040	0.005	0.008	0.009	0.000	0.000
Environmental samples	0.000	0.000	0.000	0.000	0.000	0.000	0.030	0.100	0.000	0.000	0.000	0.000	0.000
ssDNA viruses	0.080	0.700	0.800	0.500	0.090	0.100	2.00	3.00	0.400	0.300	0.200	0.000	0.000
ssRNA viruses	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Unassigned viruses	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Unclassified archaeal viruses	0.000	0.000	0.000	0.000	0.090	0.000	0.000	0.000	0.000	0.000	0.009	0.000	0.000
Unclassified phages	3.00	3.00	3.00	0.000	3.00	2.00	3.00	2.00	2.00	2.00	2.00	0.000	0.000
Unclassified virophages	0.000	0.400	0.020	0.000	0.000	0.500	0.300	0.800	0.020	0.030	0.040	0.000	0.000
Unclassified viruses	0.000	0.000	0.000	0.000	0.000	0.020	0.000	0.000	0.000	0.000	0.000	0.000	0.000

208 **Figure 2: Comparison of the Kogelberg Biosphere Reserve metavirome taxonomic**
209 **composition with selected publically available metaviromes.** Abundances normalized
210 according to predicted genome size with the GAAS tool. Blue colour represents 0.000 taxon,
211 yellow represents 0.01 – 19.00, mustard represents 20.00 – 29.00, light red represents 30.00 –
212 49.00, and red represents 50.00 – 100.00 taxon. More details on the description of
213 metaviromes are described in Supplementary Table S3 online.

214 The *Caudovirales* taxon dominated all metaviromes. In particular, members of the family
215 *Siphoviridae* were dominant in most metaviromes except for some of the freshwater samples,
216 in which myoviruses were dominant. Within the dsDNA viruses, members of rare taxonomic
217 groupings such as the genera *Tectivirus*, *Asfivirus* and *Salterprovirus*, the families
218 *Mimiviridae*, *Iridoviridae*, *Marselleviridae*, *Nudiviridae*, *Poxviridae* and *Baculoviridae* and
219 the order *Herpesvirales* were detected in soil samples as well as in hypolith, deep sea, and
220 freshwater metaviromes. Archaeal virus signatures belonging to the family *Ampullaviridae*
221 have been observed only in the Kogelberg Biosphere Reserve fynbos soil. This family
222 contains viruses with pleomorphic morphologies and a dsDNA genome, and the type species
223 infects the thermoacidophile *Acidianus convivator*, isolated from Italian hot springs⁵². Fresh
224 Water Lake, Antarctic soil and coral metaviromes showed a high abundance of ssDNA
225 viruses, results possibly biased by the use of phi29 polymerase amplification (MDA) of the
226 metaviromic DNA during library construction. The amplification of metaviromic DNA using
227 phi29 polymerase amplification (Multiple Displacement Amplification) has been reported to
228 be biased towards ssDNA templates¹⁹. It is notable, however, that a high abundance of
229 ssDNA viruses has been observed in beach freshwater samples⁵³, where amplification was
230 not used in the preparation of metagenomic DNA. However, in general, other metaviromes
231 which were not amplified using MDA showed a very low number of ssDNA viruses. In
232 general, soils or soil-associated habitats seem to harbour relatively fewer ssDNA viruses and
233 more tailed phages than aquatic ecosystems.

234 Consistent with other data^{22,24,43}, it was found that bacteriophage sequences in Kogelberg
235 Biosphere Reserve fynbos soil made up the majority of the virus fraction. Bacteriophages are
236 common in the environment and are the dominant viral type recovered from metaviromics
237 analyses in soil environments^{18,20,23,30}. This finding was not surprising, given the
238 observations from previous studies^{35,54} which showed high prokaryotic abundances in the

239 Kogelberg soil environment. Nevertheless, signature sequences from large dsDNA eukaryotic
240 virus families such as *Mimiviridae*⁵⁵ were represented in the Kogelberg
241 Biosphere Reserve library despite the use of small pore size filters in sample preparation.
242 Mimivirus signatures have been reported previously in other soil habitats²². Sequences that
243 were found to be most similar to mimivirus ORFs were also obtained from Sargasso sea
244 water samples, suggesting that these viruses, and their hosts, have a rather cosmopolitan
245 distribution⁴⁶.

246

247 2.4. PHYLOGENY OF THE KOGELBERG BIOSPHERE RESERVE FYNBOS SOIL 248 METAVIROME

249 Specific markers targeting virus families or species were used to analyse the taxonomic
250 affiliations of the annotated ORFs and analyse the diversity within the group (reviewed in⁵⁶ .
251 Phylogenetic trees were drawn from metavirome sequences on the basis of homology to
252 marker gene reference sequences from the PFAM database. Sequences homologous to the
253 marker genes (*polB*, *polB2*, *T7gp17* and *terL* (Supplementary Fig S2, S3, S4 and S5 online)
254 and reference sequences were used to draw phylogenetic trees.

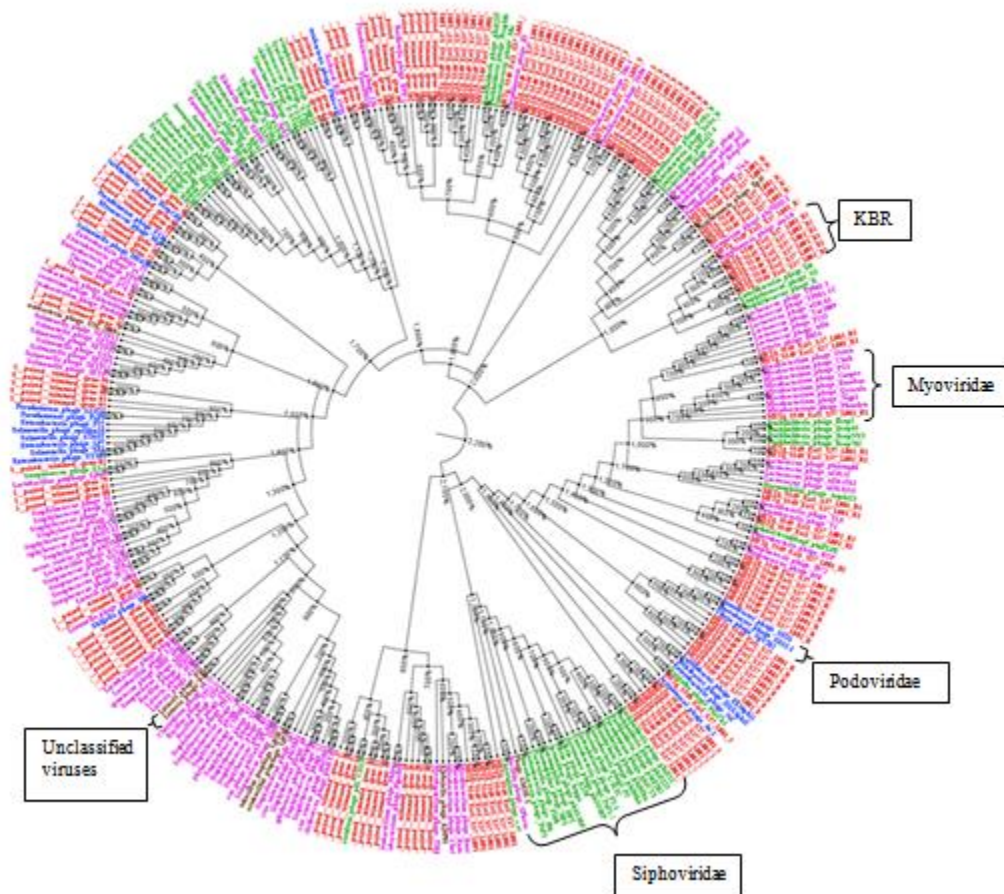
255 Using the DNA polymerase family B (*polB*) marker gene, conserved in all dsDNA viruses,
256 Kogelberg Biosphere Reserve sequences appeared to be distantly related to *Rhodothermus*
257 phage RM378 (order *Caudovirales*, family *Myoviridae*). This phage is the only sequenced
258 representative of the “Far T4” group of myoviruses (i.e., distantly related to *Escherichia virus*
259 *T4*) found in a previous diversity analysis of sequences from French lakes²⁸ . The Kogelberg
260 *polB* sequences from this study as well as the *gp23* and *gp20* marker gene sequences from the
261 French lake study contribute to the expansion of the “Far T4”-like phages dataset.

262 A DNA polymerase family B (*polB2*) marker gene, which is conserved in members of
263 *Adenoviridae*, *Salterprovirus*, and *Ampullaviridae* and *Podoviridae* family viral groups, was
264 analysed. The analysis showed a separate clade of sequences from the Kogelberg Biospheres
265 reserve soil samples. Other *polB2* sequences from our dataset were found to be distantly
266 related to members of the *Adenoviridae* family (isolated from a wide range of animal
267 sources), the *Podoviridae* family (such as *Mycoplasma* phage *P1*, *Clostridium* phage *phi24R*,
268 *Bacillus* phages *B103*, *phi39*, *Ga1*), the *Ampullaviridae* family (such as *Acidianus*-bottle-

269 shaped virus) and the *Tectiviridae* family (such as *Bacillus* phages *GIL16C*, *Bam35C* and
270 *AP50*).

271 Analysis of the metavirome sequence database using the marker gene *T7gp17* showed the
272 presence of members of the *Podoviridae* family, subfamily *Autographivirinae* and genus
273 *Phikmvvirus* and *T7virus*. Members of the genus *phikmvvirus* such as *Pseudomonas* phage
274 LKA1, and unclassified phiKMV phages such as *Ralstonia* phage RSB1, were found to be
275 closely related to the Kogelberg Biosphere Reserve sequences. Currently unclassified
276 members of the genus *T7virus*, such as *Klebsiella* phage K11 and *Yersinia* phage ϕ YeO3-12,
277 were also found to be closely related to sequences in the Kogelberg Biosphere Reserve
278 metavirome. The phages in the subfamily *Autographivirinae* are known to infect a wide range
279 of environmentally important bacteria⁵⁷.

280 Tailed phages of the order *Caudovirales* were the most commonly observed DNA viruses in
281 the Kogelberg Biosphere Reserve sequences, consistent with other environmental samples
282^{23,58,59}. A phylogenetic tree built from a *Caudovirales*-specific terminase large subunit marker
283 gene (*terL*) was used to visualise the diversity of the Kogelberg Biosphere Reserve fynbos
284 soil *Caudovirales* (Fig 3). The Kogelberg Biosphere Reserve sequences clustered with all
285 three families of tailed phages, indicating high phage richness in our sample set. These results
286 were consistent with the taxonomic affiliations of contigs in the virus families shown in Table
287 3.



288
289

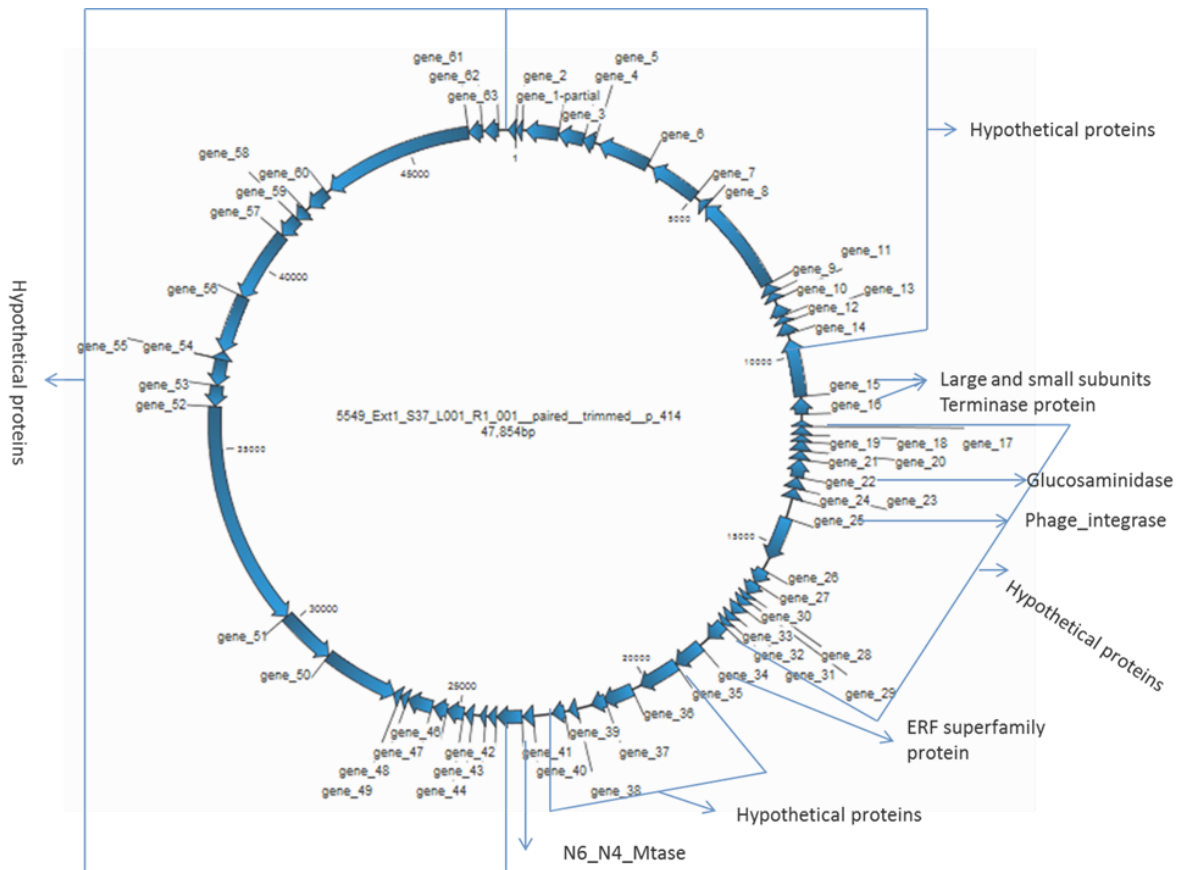
290 **Figure 3: terL phylogenetic tree.** Viral sequence origin of Caudovirales indicated with
291 different colours on the contigs names. Kogelberg Biosphere Reserve fynbos soil - Red,
292 Siphoviridae – green, Myoviridae – purple, Podoviridae - blue, unclassified viruses – grey.

293

294 2.5 ANALYSIS OF A NEAR-COMPLETE PHAGE GENOME

295

296 MetaVir assemblies predicted 352 genes from the 6 contigs larger than 40kb, as well as 758
297 genes predicted from 19 contigs of between 20kb and 40kb. The 6 largest contigs were
298 predicted to be linear, double stranded genomes. The sizes of the genomes were predicted to
299 be 47kb long with 63 genes for the largest contig (Fig 4), followed by 44kb with 58 genes,
300 42kb with 61 genes, 42kb with 53genes, 40kb with 68 genes and 40kb with 49 genes. The
301 genes in these contigs were predicted to show similarity to members of the order
302 *Caudovirales*.



303

304 **Figure 4: Gene annotation of contig 414.** Arrowed blocks are open reading frames (ORFs),
 305 showing their orientation. Numbers within the contiguous genome are nucleotide positions,
 306 starting within gene number 1 and onwards in a clockwise orientation

307

308 The largest contig represents a near-complete phage genome in the family *Podoviridae*.
 309 Members of this family typically contain double stranded and linear genomes of around 40 -
 310 45kb in length with approximately 55 genes⁶⁰. Four of the genes in this assembled genome
 311 (genes 15, 16, 34 and 41) showed similarity to members of both *Podoviridae* and
 312 *Siphoviridae* families. The translated products of two of these genes (15 and 16) were
 313 identified as putative terminase large subunit (gene 15) and terminase small subunit (gene 16)
 314 genes, with 88% and 89% amino acid identity to *Puniceispirillum* phage HMO-2011 and
 315 *Pseudomonas* phage vB_PaeP_Tr60_Ab31, respectively. Both *Puniceispirillum* phage HMO-
 316 2011 and *Pseudomonas* phage vB_PaeP_Tr60_Ab31 belong to the family *Podoviridae*. The
 317 *terL* phylogenetic tree (Supplementary Fig S4 online) showed a distant relatedness to
 318 members of the *Podoviridae* clade. Both terminase large and small subunits, together termed
 319 the terminase complex, are involved in the cleavage and packaging of concatemeric phage
 320 dsDNA⁶¹. The large terminase subunit is involved in DNA cleavage and translocation into

321 the procapsid while the small terminase subunit is involved in packaging initiation and
322 stimulation of the ATPase activity of the large terminase. These DNA packaging mechanisms
323 are used by most members of the *Caudovirales*.

324

325 The translated product of gene 34 was identified as a putative ERF superfamily protein and
326 showed 55% amino acid identity to a homologue encoded by the unclassified *Clostridium*
327 phage phiCP340 (order *Caudovirales*, family *Siphoviridae*). The ERF superfamily proteins
328 are involved in the recombination of phage genomes⁶². The translated product of gene 41
329 was identified as a putative gp77 and showed 95% amino acid similarity to a homologue
330 encoded by *Mycobacterium* phage Che9d (order *Caudovirales*, family *Siphoviridae*, genus
331 *Che8likevirus*). gp77 proteins are known to function as shut-off genes during early stages of
332 phage replication⁶³.

333

334 Fifty nine of the translated products of genes in the assembled phage genome showed identity
335 to hypothetical proteins. Of these hypothetical proteins, 56 showed no sequence similarity to
336 known virus families in BLASTp comparison to the RefseqVirus protein database. Three of
337 the genes were predicted to encode glucosaminidase (a hydrolytic enzyme), Phage integrase
338 (a site-specific recombinase that mediates controlled DNA integration and excision) and
339 PDDEXK_1 (nuclease superfamily). Members of this PDDEXK_1 family belong to the PD-
340 (D/E) XK nuclease superfamily. The PD-(D/E)XK nuclease superfamily contains type II
341 restriction endonucleases and many other enzymes involved in DNA recombination and
342 repair⁶⁴.

343

344 The protein sequences identified in this analysis indicated the presence of a putative ERF
345 superfamily protein, Phage integrase and PDDEXK_1 family; all proteins implicated in DNA
346 recombination. The ERF superfamily protein encoded by gene 34, whose sequences are
347 expressed during recombination of temperate phages, catalyses annealing of single-stranded
348 DNA chains and pairing of ssDNA with homologous dsDNA, which may function in RecA-
349 dependent and RecA-independent DNA recombination pathways⁶⁵.

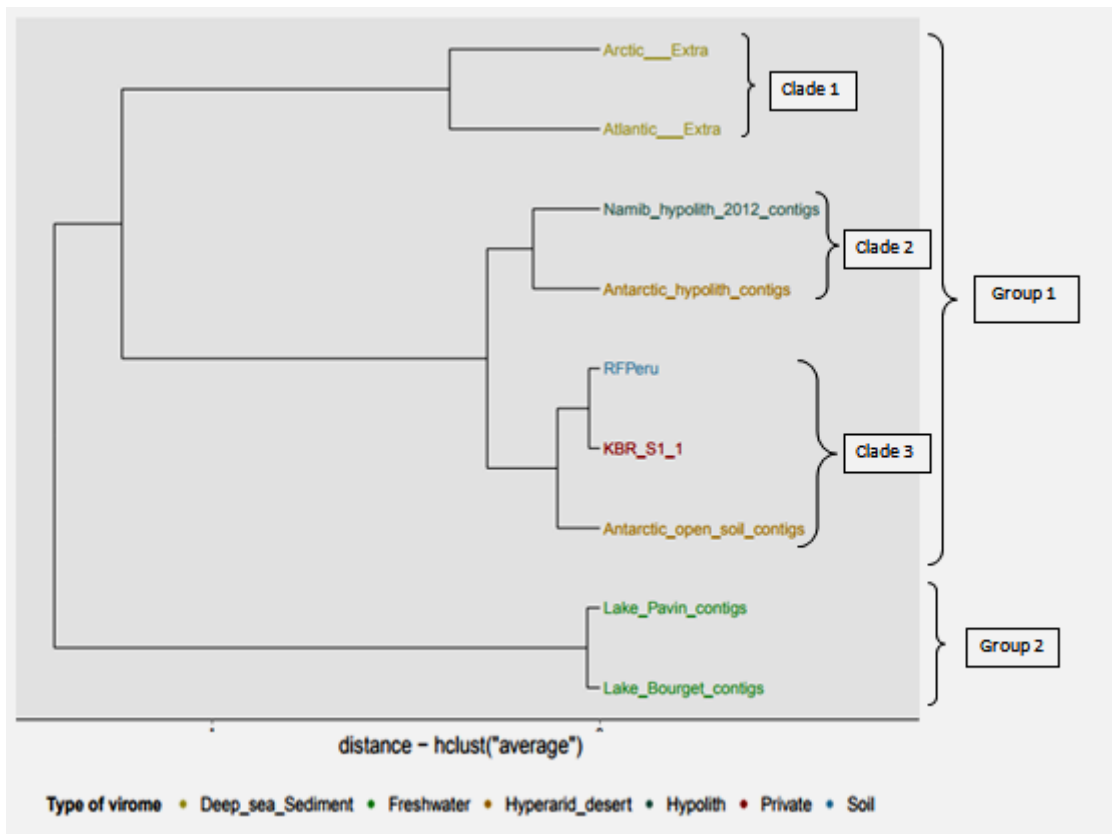
350 A few large contigs contained some predicted ORFs with similarities to phage sequences and
351 coding for specific conserved phage proteins, including terminases, structural proteins
352 (mainly related to *Caudovirales* tail structures) and phage DNA polymerases (Supplementary
353 Table S2 online).

354

355 2.6 CLUSTER ANALYSIS

356 Contig datasets from nine metaviromes from various aquatic and soil habitats were selected
357 for dinucleotide frequency comparisons ⁶⁶.

358 A comparison of the dinucleotide frequencies of the 9 metaviromes shows a clear bimodal
359 clustering (Fig 5). Group 1, composed of soil-associated habitat and deep sea sediment
360 metaviromes, is further subdivided into soil, hypolith and sediments clades. Group 2 was
361 restricted to freshwater habitats. The Arctic and Atlantic deep sea sediment and freshwater
362 lake ²⁸ metaviromes clustered in single independent nodes. Such clustering reflects
363 significant genetic similarity between these metaviromes, despite the geographical distances
364 between sample locations.



365

366 **Figure 5: Hierarchical clustering of nine metaviromes (assembled into contigs) based on**
367 **dinucleotide frequencies.** The types of biome are differentiated by colour with Kogelberg
368 Biosphere Reserve – red, freshwater – dark green, hyperarid desert – light blue, hypersaline –
369 yellow, hypolith – dark blue, seawater – light green and unknown biomes – gold. The x-axis
370 denotes eigenvalues distances. The tree was constructed using MetaVir server pipeline

371 according to the method in ⁶⁶. More details on sample names are described in supplementary
372 Table S3 online.

373

374 Both hypolithic metaviromes (i.e., cold Antarctic and hot Namib Desert hypolithic biomass
375 samples) clustered as a single node, despite their widely differing habitat-associated
376 environmental characteristics (dominated by an est. 50°C mean annual temperature
377 difference) and substantial spatial separation (approx. 55 degrees of latitude), suggesting that
378 aridity and not temperature may be the dominant driver of host and viral diversity ^{67,22}.
379 Interestingly, soil related metaviromes (from Kogelberg Biosphere Reserve fynbos soil,
380 Peruvian rainforest soil and Antarctic Dry Valley desert soil) clustered together and were
381 clearly distinct from soils which were geographically much closer.

382

383 The Kogelberg Biosphere Reserve soil metavirome clustered at a single sub-node with the
384 Peruvian rainforest soil metavirome. Both of these habitats experience high annual rainfall
385 and warm temperatures and are characterised by heavily leached and low nutrient status soils,
386 suggesting that soil composition and/or nutrient status may be the strong driver of the host
387 and viral diversity ^{68,69}. These observations suggest a niche-dependent pattern, where
388 spatially distinct niche environments cluster together and separate from their geographically
389 closer soil counterparts ⁶⁷.

390

391 Previous study reported that cluster analysis of hypolith and open soil metaviromes from
392 Antarctic and Namib Desert soil has shown that both hypolith metaviromes clustered at a
393 single node and also that both open soil metaviromes displayed an identical pattern ⁶⁷.
394 Similarly to our study, related habitat types harboured more closely related viral
395 communities, despite the great geographic distances or differing environmental conditions.
396 The common factor in these hyperarid environments is water scarcity, which may be a key
397 driver of community speciation and recruitment in these environments. We conclude that
398 these adaptations and the nature of soil habitat compared to the 'refuge' habitat of quartz
399 stones for hypolithic communities, may be the driving force between both communities not to
400 cluster together.

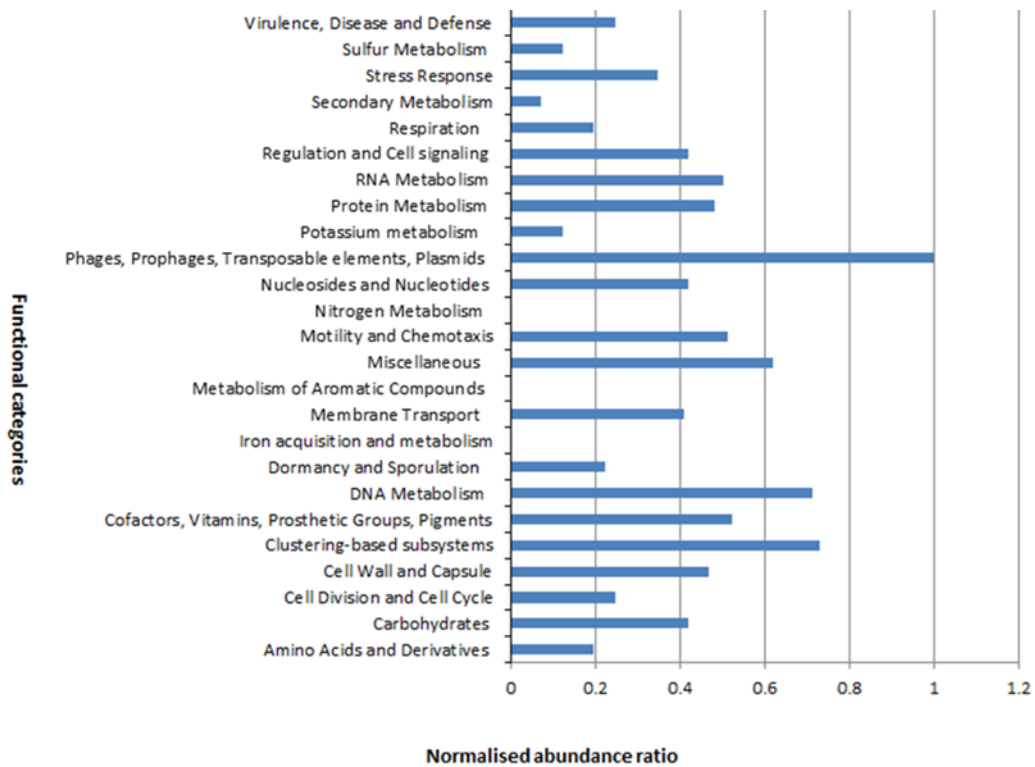
401

402 2.7 FUNCTIONAL PROPERTIES OF THE KOGELBERG BIOSPHERE RESERVE
403 FYNBOS SOIL METAVIROME

404

405 The functional implication of the reads was explored using MG-RAST. The Kogelberg
406 Biosphere Reserve metavirome sequences exhibited a high proportion of uncharacterized
407 ORFs, with 2,362,076 sequences showing no significant similarities to proteins in the
408 databases (ORFans). Twelve functional categories were annotated by MG-RAST, each
409 subdivided into distinct subsystems (Fig 6). The database searches against SEED in the MG-
410 RAST subsystem resulted in 9360 hits. The highest percentage hits (20.3%) in the functional
411 annotation belonged to the “Phage, prophages, transposable elements and plasmids”
412 subsystem category, with r1t-like streptococcal phages, phage packaging machinery and
413 phage replication annotations most commonly identified.

414



415

416

417 **Figure 6: Functional assignment of predicted ORFs.** Functional annotation was performed
418 at 60% similarity cut-off as predicted by MG-RAST.

419

420 The other functional subsystem categories showed “Clustering-based subsystems (e.g.,
421 biosynthesis of galactoglycans and related lipopolysaccharides; catabolism of an unclassified
422 compound etc., and other clusters identified as unclassified). The “Protein metabolism” and

423 “DNA metabolism” functional categories were also dominant annotations. Many proteins in
424 these functional categories, such as terminases, HNH homing endonucleases, DNA helicases,
425 DNA polymerases and DNA primases, could potentially be of phage origin. These functional
426 groups have also been found to be highly represented in previous metaviromic datasets
427 ^{23,70,71}.

428

429 Analysis of the metavirome reads using the KEGG Orthology (KO) database showed
430 metabolism protein families (carbohydrate metabolism, amino acid metabolism and
431 nucleotide metabolism) to be the most commonly identified. Members of the genetic
432 information procession protein family, including replication and repair, transcription and
433 translation proteins, were also commonly identified. Deeper analysis of a subset of annotated
434 contigs identified genes encoding numerous virus structures (e.g., phage capsid, terminase,
435 tail fibre protein etc.) and DNA manipulating enzymes (e.g., endonuclease, DNA methylase,
436 primase-polymerase, DNA primase/helicase, DNA polymerase I, integrase, ssDNA annealing
437 protein, exonuclease, transferase, site-specific DNA methylase, ligase, recombinase etc.).

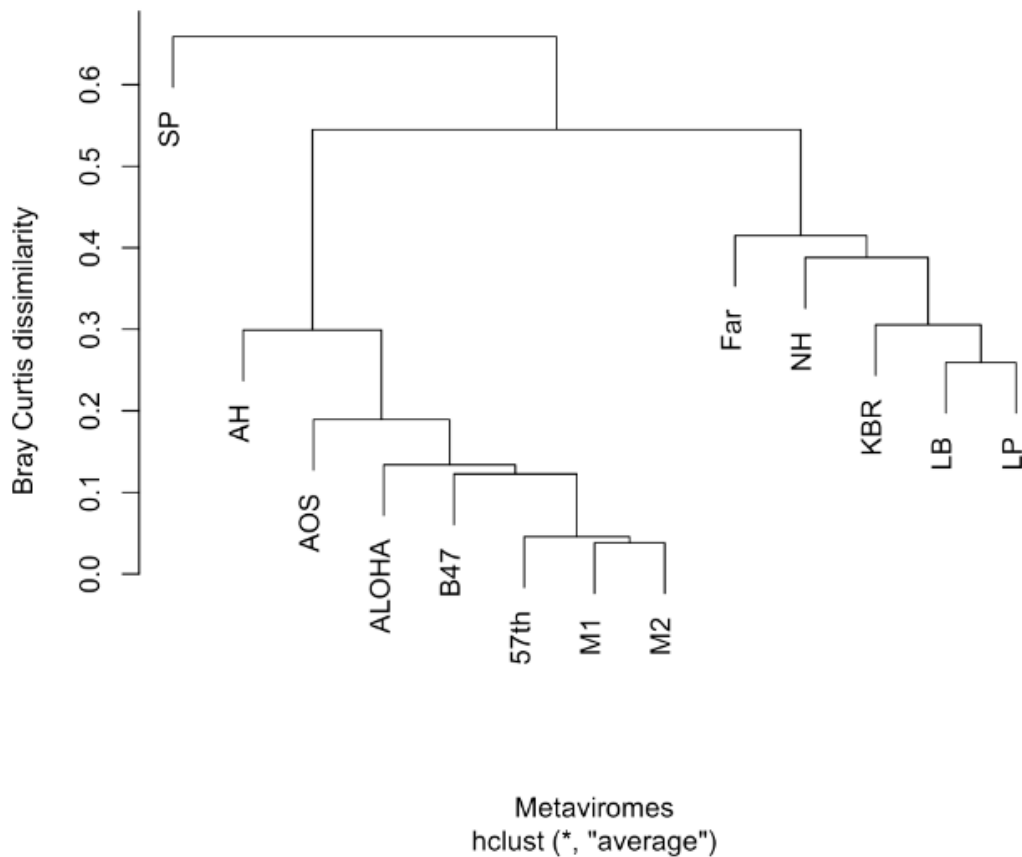
438

439 From this analysis, we demonstrate that phage-related genes and metabolic genes are highly
440 represented. The virome displayed a strong enrichment in phage-like genes (e.g. phages,
441 prophages, transposable elements, plasmids) and lacked typical cellular categories rarely
442 observed in sequenced phages (e.g. ‘cofactors, vitamins, prosthetic groups, pigments’).
443 Cellular categories commonly identified in known phages were retrieved (e.g. ‘nucleosides
444 and nucleotides’, ‘DNA metabolism’). The highly abundance of virome-associated metabolic
445 genes shows that the phages may have the potential to interfere with the metabolism of their
446 hosts. Our virome analysis, consistent with other virome studies, demonstrate the unexpected
447 picture of global ‘viral’ metabolism, suggesting that viruses might actively dictate the
448 metabolism of infected cells on a global scale ⁷¹.

449

450 The functional assignments from the SEED database of Kogelberg Biosphere Reserve fynbos
451 soil was clustered with SEED database functional assignments of the 12 previously published
452 metaviromes from both similar and dissimilar environments (fresh water ⁴⁸, soil and
453 hypolithic niche communities ^{22,23}, pond water ²⁷ and sea water ⁴⁹ mentioned in Fig 2. A
454 cluster analysis of the SEED database subsystem classification revealed different functional
455 patterns between the metaviromes and no clear soil clustering (Fig 7). The sequences from

456 Kogelberg Biosphere Reserve clustered amongst the sequences from three of the fresh water
 457 lakes and the Namib hypolith metaviromes. Antarctic samples (Antarctic open soil and
 458 Antarctic hypolith) were more distinct and formed a heterogeneous clade with the other fresh
 459 water samples. This can be potentially be explained by the larger number of cellular
 460 contamination in some of these metaviromes. This finding suggests that different biomes can
 461 share similar functional patterns and, conversely, that taxonomically similar viromes can
 462 encode different functional genes. It may also indicate that certain phage groups are more
 463 prevalent in certain biogeographic regions.



464
 465 **Figure 7: Cluster analysis of functional assignment of predicted ORFs.** Viromes were
 466 clustered with the hclust algorithm in R according to the abundance of SEED database
 467 functional categories present. SEED categories were assigned using Megan6 after blastp-
 468 based comparison with the non-redundant protein database of NCBI. More details on the
 469 description of metaviromes are described in Supplementary Table 2 online.

470 This study is not without limitations. The major limitation to this study is the use of only a
 471 single virome that includes only double stranded DNA viruses.

472

473 **3. CONCLUSION**

474 We have successfully used the metaviromics approach to explore the diversity and functional
475 composition of a previously unexplored Kogelberg Biosphere Reserve fynbos soil virome.
476 Our quantitative comparison of taxonomic and functional composition of the Kogelberg soil
477 metavirome with other published viromes is a valuable and novel contribution that will
478 enhance the repertoire of publicly available datasets and advance our understanding of viral
479 ecology. Furthermore, contigs corresponding to novel virus genomes were assembled in the
480 current work; this presents an opportunity for future studies aimed at targeting these novel
481 genetic resources for applied biotechnology.

482

483 **4. EXPERIMENTAL DESIGN**

484 **4.1 SAMPLE SITE LOCATION**

485 Samples were collected from the Kogelberg Biosphere Reserve, situated to the east of Cape
486 Town, South Africa in the Boland Mountains (GPS coordinates: 34°19'48".0 S, 18°57'21.0"
487 E). Open soil samples were collected aseptically during the winter of 2014. Approximately
488 20kg of soil was collected at depth of 0 - 4cm. Soil samples were stored in sterile containers
489 at -80°C.

490 **4.2 SAMPLE PROCESSING, DNA EXTRACTION**

491

492 Samples were collected in the open soil. Only 3 samples were collected. The DNA of these
493 samples were pooled together for NGS sequencing. Soil samples were processed as
494 previously described⁷² with some modifications. 8 kg of soil and 1X SM buffer (8L) (0.1 M
495 NaCl, 8 mM MgSO₄, 50mM Tris-HCl, pH 7.5) were mixed and shaken vigorously in a sterile
496 container until soil was well suspended and left overnight at 4°C to settle. The supernatant
497 was centrifuged at 10000g for 15min to pellet any remaining soil particles and other debris
498 and passed through a 0.22µm filter (Millipore, streicup 500ml). The filtrate was treated with
499 DNase. Viral particles were precipitated with 10% (w/v) polyethylene glycol (PEG) 8000
500 overnight at 4°C and centrifuged for 15min at 11000g. After removing the supernatant the
501 viral pellet was resuspended in TE buffer, pH 7.6⁷².

502 The absence of bacterial and eukaryotic DNA was confirmed by PCR with primers pairs E9F
503 (5'-GAG TTT GAT CCT GGC TCA G-3') and U1510R (5'-GGT TAC CTT GTT ACG ACT

504 T-3') and ITS1 (5'- TCCGTAGGT GAACCTGCGG-3') and ITS4 (5'-
505 TCCTCCGCTTATTGATATGC-3')⁷³.

506 4.3 TRANSMISSION ELECTRON MICROSCOPY

507 Aliquots of viral suspensions isolated from soil were fixed with 2 % glutaraldehyde for three
508 hours at 4°C and 10 µl of the phage suspension was overlaid on a carbon coated grid of 200
509 Mesh⁷⁴. The suspension was allowed to dry on the grid, which was then negatively stained
510 with 2% uranyl acetate. Excess stain was removed using filter paper and allowed to air-dry
511 prior to examination using a Philips (FEI) CM100 TEM.

512 4.4 DNA EXTRACTION SEQUENCING

513 DNA was extracted from virus particle preparations using a ZR soil microbe DNA
514 MidiPrepTM kit according to manufacturer's instructions (Zymo Research). Extracted
515 metaviromic DNA (unamplified) was sequenced using an Illumina MiSeq platform (Inqaba
516 Biotechnical Industries). Briefly, following DNA quantification using NanoDrop
517 Fluorospectrometer 3300, 1 ng of isolated metavirome DNA was used to prepare 4
518 individually indexed NexteraXT libraries. They were then sequenced using the MiSeq v3
519 (600 cycles) sequencing kit, generating 2 x 300 bp reads. The raw reads were trimmed and
520 demultiplexed, resulting in four fastq files.

521 4.5 SEQUENCE DATA ANALYSIS

522 The quality of the raw read files was checked with CLC Genomics Workbench version 6.0.1
523 (CLC, Denmark). The reads were then filtered and trimmed, with the removal of low quality
524 (sequence limit of 0.05), ambiguous reads (maximal of 2 and minimum length of 15). This
525 yielded 1,488,462,918 reads with an average length of 212.05 bp. The post-QC reads were
526 assembled using CLC Genomics Workbench as paired files (3 X 2 read files per site). The
527 assembly resulted in 28,511,204 contigs with a minimum length of 1,002 bases at an N50 of
528 2,047 and a maximum of 47,854 bases.

529 The processed reads were assembled *de novo* using CLC Genomics Workbench version 6.0.1
530 using the default settings. Reads and contigs were uploaded to the MetaVir³⁷ ([http://metavir-
531 meb.univ-bpclermont.fr](http://metavir-
531 meb.univ-bpclermont.fr)), VIROME (<http://virome.dbi.udel.edu/>)³⁸ and MG-RAST
532 (<http://metagenomics.anl.gov/>)³⁹ servers for virus diversity estimations. The viromes were
533 uploaded in 2015 and analysed in 2017. The taxonomic composition was computed from a

534 BLAST comparison with the Refseq complete viral genomes protein sequence database from
535 NCBI (release of 2016-01) using BLASTp with a threshold of 50 on the BLAST bitscore.
536 The assembled sequences were searched for open reading frames (ORFs) and compared to
537 the RefSeq complete viral database using MetaVir and MG-RAST. Functional and organism
538 assignments were based on annotation and other information obtained from the following
539 databases: GenBank, Integrated Microbial Genomes (IMG), Kyoto Encyclopaedia of Genes
540 and Genomes (KEGG), Pathosystems Resource Integration Center (PATRIC), RefSeq,
541 SEED, Swiss-Prot, tremble, and egnog; and for the assignment of functional hierarchy,
542 COG (clusters of orthologous groups), KEGG Orthology (KO), and NOG databases were
543 used. The Genome relative Abundance and Average Size (GAAS) ⁷⁵ tools were used for
544 normalization of the total composition, estimation of the mean genome length and for the
545 estimation of relative abundance and size for each taxon. The phylogenetic tree were
546 generated by an open-source JavaScript library called jsPhyloSVG ⁷⁶. The phylogenetic trees
547 were based on the reference sequences and the Kogelberg Biosphere Reserve virome
548 sequences, and computed with 100 bootstraps. Further analysis of the sequences was
549 performed using METAGENassist (a web server that provides a broad range of statistical
550 tools for comparative metagenomics) ⁷⁷. Functional assignments produced by VIROME using
551 120 identified functional subsystems were used for the statistical analysis with
552 METAGENassist.

553 Clustering analysis comparison was plotted as a clustering tree and computed with pvclust
554 computed by MetaVir (an R package for assessing the uncertainty in hierarchical clustering)
555 ⁷⁸ (Fig 6). Hierarchical clustering using dinucleotide comparisons was used to quantify the
556 grouping behaviour of nine published metaviromes and the comparison were plotted and
557 demonstrated as a clustering dendrograms. Only metaviromes containing more than 50,000
558 sequences and with an average sequence length of over 100 bp were used, as this comparison
559 is based on a normalised virome sub-sample. Metaviromes that did not match these criteria
560 were not listed for nucleotide composition bias comparison. Hence, only 9 metaviromes were
561 suitable for comparison using dinucleotide frequencies in the MetaVir sever. The largest
562 contigs were analysed by MetaVir. The SEED classification clustering of the 12 metaviromes
563 was assessed using BLASTp against the nr database of NCBI (release 2017-05) ⁷⁹.
564 Differences between the virome SEED functional components were transformed into a Bray
565 Curtis dissimilarity matrix using the vegan package in RStudio, clustered using the hclust
566 algorithm (method = average), and represented as a dendrogram ^{80,81}.

567 DATA AVAILABILITY

568

569 Viral sequences from Kogelberg Biosphere Reserve fynbos soil sample are available on
570 MetaVir under the project KBR under the names “KBR 1 and KBR 2”.

571

572 REFERENCE

- 573 1. Bergh, E. W. & Compton, J. S. South African Journal of Botany A one-year post- fire
574 record of macronutrient cycling in a mountain sandstone fynbos ecosystem , South
575 Africa. *South African J. Bot.* **97**, 48–58 (2015).
- 576 2. Slabbert, E., Kongor, R. Y., Esler, K. J. & Jacobs, K. Microbial diversity and
577 community structure in Fynbos soil. *Mol. Ecol.* **19**, 1031–1041 (2010).
- 578 3. Van Wyk, A. E. & Smith, G. Regions of Floristic Endemism in Southern Africa.
579 *Umdaus* 5–17 (2001).
- 580 4. Sprent, J. I. & Parsons, R. Nitrogen fixation in legume and non-legume trees. *F. Crop.*
581 *Res.* **65**, 183–196 (2000).
- 582 5. Mucina, L. & Wardell-Johnson, G. W. Landscape age and soil fertility, climatic
583 stability, and fire regime predictability: beyond the OCBIL framework. *Plant Soil* **341**,
584 (2011).
- 585 6. Jeanbille, M. *et al.* Soil Parameters Drive the Structure, Diversity and Metabolic
586 Potentials of the Bacterial Communities Across Temperate Beech Forest Soil
587 Sequences. *Microb. Ecol.* **71**, 482–493 (2016).
- 588 7. Rappé, M. S. & Giovannoni, S. J. The uncultured microbial majority. *Annu. Rev.*
589 *Microbiol.* **57**, 369–394 (2003).
- 590 8. van der Heijden, M. G. A., Bardgett, R. D. & van Straalen, N. M. The unseen majority:
591 soil microbes as drivers of plant diversity and productivity in terrestrial ecosystems.
592 *Ecol. Lett.* **11**, 296–310 (2008).
- 593 9. Stafford, W. H. L., Baker, G. C., Brown, S. A., Burton, S. G. & Cowan, D. A.
594 Bacterial diversity in the rhizosphere of Proteaceae species. *Environ. Microbiol.* **7**,
595 1755–1768 (2005).
- 596 10. Ramond, J. B., Lako, J. D. W., Stafford, W. H. L., Tuffin, M. I. & Cowan, D. A.
597 Evidence of novel plant-species specific ammonia oxidizing bacterial clades in acidic

- 598 South African fynbos soils. *J. Basic Microbiol.* **55**, 1040–1047 (2015).
- 599 11. Prosser, J. I. & Nicol, G. W. Relative contributions of archaea and bacteria to aerobic
600 ammonia oxidation in the environment. *Environ. Microbiol.* **10**, 2931–41 (2008).
- 601 12. Hamilton, E. W. & Frank, D. A. Can plants stimulate soil microbes and their own
602 nutrient supply? Evidence from a grazing tolerant grass. *Ecology* **82**, 2397–2402
603 (2001).
- 604 13. Nüsslein, K. & Tiedje, J. M. Soil bacterial community shift correlated with change
605 from forest to pasture vegetation in a tropical soil. *Appl. Environ. Microbiol.* **65**, 3622–
606 6 (1999).
- 607 14. Keluskar, R., Nerurkar, A. & Desai, A. Mutualism between autotrophic ammonia-
608 oxidizing bacteria (AOB) and heterotrophs present in an ammonia-oxidizing colony.
609 *Arch. Microbiol.* **195**, 737–747 (2013).
- 610 15. Kimura, M., Jia, Z. J., Nakayama, N. & Asakawa, S. Ecology of viruses in soils: Past,
611 present and future perspectives. *Soil Sci. Plant Nutr.* **54**, 1–32 (2008).
- 612 16. Cowan, D. A., Rybicki, E. P., Tuffin, M. I., Valverde, A. & Wingfield, M. J.
613 Biodiversity: So much more than legs and leaves. *S. Afr. J. Sci.* **109**, 1–9 (2013).
- 614 17. Schoenfeld, T. *et al.* Functional viral metagenomics and the next generation of
615 molecular tools. *Trends Microbiol.* **18**, 20–9 (2010).
- 616 18. Zablocki, O., Adriaenssens, E. M. & Cowan, D. Diversity and Ecology of Viruses in
617 Hyperarid Desert Soils. *Appl. Environ. Microbiol.* **82**, 770–777 (2016).
- 618 19. Kim, K. H. *et al.* Amplification of uncultured single-stranded DNA viruses from rice
619 paddy soil. *Appl. Environ. Microbiol.* **74**, 5975–5985 (2008).
- 620 20. Kim, M.-S., Whon, T. W. & Bae, J.-W. Comparative viral metagenomics of
621 environmental samples from Korea. *Genomics Inform.* **11**, 121–8 (2013).
- 622 21. Srinivasiah, S. *et al.* Direct assessment of viral diversity in soils by random PCR
623 amplification of polymorphic DNA. *Appl. Environ. Microbiol.* **79**, 5450–7 (2013).
- 624 22. Zablocki, O. *et al.* High-Level Diversity of Tailed Phages, Eukaryote-Associated
625 Viruses, and Virophage-Like Elements in the Metaviromes of Antarctic Soils. *Appl.*
626 *Environ. Microbiol.* **80**, 6888–6897 (2014).
- 627 23. Adriaenssens, E. M. *et al.* Metagenomic analysis of the viral community in Namib
628 Desert hypoliths. *Environ. Microbiol.* **17**, 480–495 (2015).
- 629 24. Fancello, L. *et al.* Viruses in the desert: a metagenomic survey of viral communities in
630 four perennial ponds of the Mauritanian Sahara. *ISME J.* **7**, 359–69 (2013).

- 631 25. Gudbergsdóttir, S. R., Menzel, P., Krogh, A., Young, M. & Peng, X. Novel viral
632 genomes identified from six metagenomes reveal wide distribution of archaeal viruses
633 and high viral diversity in terrestrial hot springs. *Environ. Microbiol.* n/a-n/a (2015).
634 doi:10.1111/1462-2920.13079
- 635 26. Breitbart, M. Marine viruses: truth or dare. *Ann. Rev. Mar. Sci.* **4**, 425–48 (2012).
- 636 27. Rodriguez-Brito, B. *et al.* Viral and microbial community dynamics in four aquatic
637 environments. *Isme J* **4**, 739–751 (2010).
- 638 28. Roux, S. *et al.* Assessing the Diversity and Specificity of Two Freshwater Viral
639 Communities through Metagenomics. *PLoS One* **7**, e33641 (2012).
- 640 29. Roux, S., Hallam, S. J., Woyke, T. & Sullivan, M. B. Viral dark matter and virus-host
641 interactions resolved from publicly available microbial genomes. *Elife* **4**, (2015).
- 642 30. Hatfull, G. F. Dark matter of the biosphere: The amazing world of bacteriophage
643 diversity. *J. Virol.* **89**, 8107–8110 (2015).
- 644 31. Alhamlan, F. S., Ederer, M. M., Brown, C. J., Coats, E. R. & Crawford, R. L.
645 Metagenomics-based analysis of viral communities in dairy lagoon wastewater. *J.*
646 *Microbiol. Methods* **92**, 183–8 (2013).
- 647 32. Williamson, K. E., Radosevich, M. & Wommack, K. E. Abundance and diversity of
648 viruses in six Delaware soils. *Appl. Environ. Microbiol.* **71**, 3119–3125 (2005).
- 649 33. Reavy, B. *et al.* Distinct circular single-stranded DNA viruses exist in different soil
650 types. *Appl. Environ. Microbiol.* **81**, 3934–45 (2015).
- 651 34. Ackermann, H. W. 5500 Phages examined in the electron microscope. *Arch. Virol.*
652 **152**, 227–243 (2007).
- 653 35. Zablocki, O., Adriaenssens, E. & Cowan, D. Diversity and ecology of viruses in
654 hyperarid desert soils. *Appl. Environ. Microbiol.* **82**, AEM.02651-15- (2015).
- 655 36. Ackermann, H.-W. Classification of bacteriophages. *The bacteriophages* **2**, 8–16
656 (2006).
- 657 37. Roux, S., Tournayre, J., Mahul, A., Debroas, D. & Enault, F. Metavir 2: new tools for
658 viral metagenome comparison and assembled virome analysis. *BMC Bioinformatics*
659 **15**, 76 (2014).
- 660 38. Wommack, K. E. *et al.* VIROME: a standard operating procedure for analysis of viral
661 metagenome sequences. *Stand. Genomic Sci.* **6**, 427 (2012).
- 662 39. Meyer, F. *et al.* The metagenomics RAST server - a public resource for the automatic
663 phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* **9**, 386

- 664 (2008).
- 665 40. Wilke, A. *et al.* The M5nr: a novel non-redundant database containing protein
666 sequences and annotations from multiple sources and associated tools. *BMC*
667 *Bioinformatics* **13**, 141 (2012).
- 668 41. Overbeek, R., Disz, T. & Stevens, R. The SEED: A peer-to-peer environment for
669 genome annotation. *Commun. ACM* **47**, 47–51 (2004).
- 670 42. Mohiuddin, M. & Schellhorn, H. E. Spatial and temporal dynamics of virus occurrence
671 in two freshwater lakes captured through metagenomic analysis. *Front. Microbiol.* **6**,
672 960 (2015).
- 673 43. Breitbart, M. *et al.* Genomic analysis of uncultured marine viral communities. *Proc.*
674 *Natl. Acad. Sci. U. S. A.* **99**, 14250–5 (2002).
- 675 44. Cann, A. J., Fandrich, S. E. & Heaphy, S. Analysis of the virus population present in
676 equine faeces indicates the presence of hundreds of uncharacterized virus genomes.
677 *Virus Genes* **30**, 151–156 (2005).
- 678 45. Gascuel, O. *et al.* BIONJ: an improved version of the NJ algorithm based on a simple
679 model of sequence data. *Mol. Biol. Evol.* **14**, 685–695 (1997).
- 680 46. Ghedin, E. & Claverie, J.-M. M. Mimivirus relatives in the Sargasso sea. *Viol. J.* **2**, 62
681 (2005).
- 682 47. Short, S. M. & Short, C. M. Diversity of algal viruses in various North American
683 freshwater environments. *Aquat. Microb. Ecol.* **51**, 13–21(2008).
- 684 49. Angly, F. E. *et al.* The marine viromes of four oceanic regions. *PLoS Biol.* **4**, e368
685 (2006).
- 686 50. Dinsdale, E. a *et al.* Functional metagenomic profiling of nine biomes. *Nature* **452**,
687 629–632 (2008).
- 688 51. de Wit, R. & Bouvier, T. ‘Everything is everywhere, but, the environment selects’;
689 what did Baas Becking and Beijerinck really say? *Environ. Microbiol.* **8**, 755–758
690 (2006).
- 691 52. Häring, M., Rachel, R., Peng, X., Garrett, R. A. & Prangishvili, D. Viral Diversity in
692 Hot Springs of Pozzuoli, Italy, and Characterization of a Unique Archaeal Virus,
693 Acidianus Bottle-Shaped Virus, from a New Family, the Ampullaviridae. *J Virol* **79**,
694 9904–9911 (2005).
- 695 53. Watkins, S. C. *et al.* Assessment of a metaviromic dataset generated from nearshore
696 Lake Michigan. *Mar. Freshw. Res.* **67**, 1700 (2016).

- 697 54. Ashelford, K. E., Day, M. J. & Fry, J. C. Elevated abundance of bacteriophage
698 infecting bacteria in soil. *Appl. Environ. Microbiol.* **69**, 285–289 (2003).
- 699 55. Raoult, D., La Scola, B. & Birtles, R. The discovery and characterization of
700 Mimivirus, the largest known virus and putative pneumonia agent. *Clin. Infect. Dis.*
701 **45**, 95–102 (2007).
- 702 56. Adriaenssens, E. M. & Cowan, D. A. Using signature genes as tools to assess
703 environmental viral ecology and diversity. *Appl. Environ. Microbiol.* **80**, 4470–4480
704 (2014).
- 705 57. Adriaenssens, E. M. *et al.* Bacteriophages LIMelight and LIMEzero of *Pantoea*
706 agglomerans, belonging to the ‘phiKMV-Like Viruses’. *Appl. Environ. Microbiol.* **77**,
707 3443–3450 (2011).
- 708 58. Roux, S. *et al.* Analysis of metagenomic data reveals common features of halophilic
709 viral communities across continents. *Environ. Microbiol.* **18**, 889–903 (2015).
- 710 60. Hulo, C. *et al.* ViralZone: A knowledge resource to understand virus diversity. *Nucleic*
711 *Acids Res.* **39**, 576–582 (2011).
- 712 61. Penadés, J. R., Chen, J., Quiles-Puchalt, N., Carpena, N. & Novick, R. P.
713 Bacteriophage-mediated spread of bacterial virulence genes. *Curr. Opin. Microbiol.*
714 **23**, 171–178 (2015).
- 715 62. Wittmann, J., Gartemann, K.-H., Eichenlaub, R. & Dreiseikelmann, B. Genomic and
716 molecular analysis of phage CMP1 from *Clavibacter michiganensis* subspecies
717 *michiganensis*. **1**, 6–14 (2011).
- 718 63. Rybniker, J., Plum, G., Robinson, N., Small, P. L. & Hartmann, P. Identification of
719 three cytotoxic early proteins of mycobacteriophage L5 leading to growth inhibition in
720 *Mycobacterium smegmatis*. *Microbiology* **154**, 2304–2314 (2008).
- 721 64. Letunic, I. *et al.* SMART 4.0: towards genomic data integration. *Nucleic Acids Res.* **32**,
722 142D–144 (2004).
- 723 65. Dziewit, L., Oscik, K., Bartosik, D. & Radlinska, M. Molecular characterization of a
724 novel temperate sinorhizobium bacteriophage, Φ LM21, encoding DNA
725 methyltransferase with CcrM-like specificity. *J. Virol.* **88**, 13111–24 (2014).
- 726 66. Willner, D., Thurber, R. V. & Rohwer, F. Metagenomic signatures of 86 microbial and
727 viral metagenomes. *Environ. Microbiol.* **11**, 1752–1766 (2009).
- 728 67. Zablocki O., van Zyl L., Adriaenssens EM., Rubagotti E., Tuffin M., Cary SC., C. D.
729 Niche-dependent genetic diversity in Antarctic metaviromes. *Bacteriophage* **4**,

- 730 e980125 (2014).
- 731 68. Leigh, E. G. Structure and climate in tropical rain forest. *Source Annu. Rev. Ecol. Syst.*
732 **6**, 67–86 (1975).
- 733 69. Fry, M. A detailed characterization of soils under different Fynbos-climate-geology
734 combinations in the south-western Cape. (1987).
- 735 70. Cai, L., Zhang, R., He, Y., Feng, X. & Jiao, N. Metagenomic Analysis of
736 Virioplankton of the Subtropical Jiulong River Estuary, China. *Viruses* **8**, 35 (2016).
- 737 71. Roux, S., Krupovic, M., Debroas, D. & Forterre, P. Assessment of viral community
738 functional potential from viral metagenomes may be hampered by contamination with
739 cellular sequences. *Open Bio* **3**, (2013).
- 740 72. Casas, V. & Rohwer, F. Phage metagenomics. *Methods Enzymol.* **421**, 259–268
741 (2007).
- 742 73. Merseguel, K. B. *et al.* Genetic diversity of medically important and emerging *Candida*
743 species causing invasive infection. *BMC Infect. Dis.* **15**, 57 (2015).
- 744 74. Ackermann, H.-W. in *Bacteriophages* 113–126 (Springer, 2009).
- 745 75. Angly, F. E. *et al.* The GAAS metagenomic tool and its estimations of viral and
746 microbial average genome size in four major biomes. *PLoS Comput. Biol.* **5**, e1000593
747 (2009).
- 748 76. Smits, S. A. & Ouverney, C. C. jsPhyloSVG: A Javascript Library for Visualizing
749 Interactive and Vector-Based Phylogenetic Trees on the Web. *PLoS One* **5**, e12267
750 (2010).
- 751 77. Arndt, D. *et al.* METAGENassist: A comprehensive web server for comparative
752 metagenomics. *Nucleic Acids Res.* **40**, 1–8 (2012).
- 753 78. Suzuki, R. & Shimodaira, H. Pvclust: An R package for assessing the uncertainty in
754 hierarchical clustering. *Bioinformatics* **22**, 1540–1542 (2006).
- 755 79. Aziz, R. K. *et al.* SEED Servers: High-Performance Access to the SEED Genomes,
756 Annotations, and Metabolic Models. *PLoS One* **7**, e48053 (2012).
- 757 80. Oksanen J, Blanchet FG, Kindt R, L. P. and others. ‘vegan’ 2.3.4.4— community
758 ecology package. <http://CRAN.R-project.org/package=vegan>. (2016).
- 759 81. Racine, J. S. RStudio: A Platform-Independent IDE for R and Sweave. *J. Appl.*
760 *Econom.* **27**, 167–172 (2012).

761

762

763 **ACKNOWLEDGMENTS**

764 We gratefully acknowledge financial support from National Research Foundation (NRF), the
765 Department of Science and Technology, the University of Pretoria Genomics Research
766 Institute, the Claude Leon Foundation (to EMA), Cape Nature (Kogelberg Biosphere reserve)
767 and the Council for Scientific and Industrial Research (CSIR). We declare no conflicts of
768 interest.

769 **AUTHOR CONTRIBUTIONS**

770 D.C., E.A., T.T., and K.R. conceived and supervised the study. J.S., K.R., and T.T designed
771 the experiments. J.S. and K.R. performed the experiments. J.S., E.A., and K.R analysed data.
772 J.S., D.C., E.A., T.T., and K.R., wrote the manuscript.

773

774 **ADDITIONAL INFORMATION**

775 The sequences of metaviromes identified and verified in this project have been submitted to
776 MetaVir server with the project ID shown in Supplementary1 Table S1.

777 Supplementary information accompanies this paper

778 Competing Interests: The authors declare no competing financial interests.

779

780

781