

1 **Title: Tracking infection dynamics at single-cell level reveals highly**
2 **resolved expression programs of a large virus infecting algal blooms**

3
4 **Authors:** Chuan Ku^{1,6‡}, Uri Sheyn^{1,‡}, Arnau Sebé-Pedrós^{2,3}, Shifra Ben-Dor⁴, Daniella
5 Schatz¹, Amos Tanay^{2,3}, Shilo Rosenwasser^{5,*}, Assaf Vardi^{1,*}

6
7 **Affiliations:**

8 ¹ Department of Plant and Environmental Sciences, Weizmann Institute of Science, Rehovot,
9 Israel

10 ² Department of Computer Science and Applied Mathematics, Weizmann Institute of Science,
11 Rehovot, Israel

12 ³ Department of Biological Regulation, Weizmann Institute of Science, Rehovot, Israel

13 ⁴ Department of Life Sciences Core Facilities, Weizmann Institute of Science, Rehovot, Israel

14 ⁵ Institute of Plant Sciences and Genetics in Agriculture, The Hebrew University of
15 Jerusalem, Rehovot, Israel

16 ⁶ Present address: Institute of Plant and Microbial Biology, Academia Sinica, Taipei, Taiwan

17 ‡ These authors contributed equally to this work: Chuan Ku, Uri Sheyn

18 * Corresponding authors. Email: shilo.rosenwasser@mail.huji.ac.il (S.R.);

19 assaf.vardi@weizmann.ac.il (A.V.)

20

21

22

23

24

25

26

27 **Abstract:**

28 Nucleocytoplasmic large DNA viruses have the largest genomes among all viruses and
29 infect diverse eukaryotes across various ecosystems, but their expression regulation and
30 infection strategies are not well understood. We profiled single-cell transcriptomes of the
31 worldwide-distributed microalga *Emiliana huxleyi* and its specific coccolithovirus
32 responsible for massive bloom demise. Heterogeneity in viral transcript levels detected among
33 single cells was used to reconstruct the viral transcriptional trajectory and to map cells along a
34 continuum of infection states. This enabled identification of novel viral genetic programs,
35 which are composed of five kinetic classes with distinct promoter elements. The infection
36 substantially changed the host transcriptome, causing rapid shutdown of protein-encoding
37 nuclear transcripts at the onset of infection, while the plastid and mitochondrial
38 transcriptomes persisted to mid- and late stages, respectively. Single-cell transcriptomics
39 thereby opens the way for tracking host-pathogen infection dynamics at high resolution within
40 microbial communities in the marine environment.

41

42

43 **Main text**

44 Nucleocytoplasmic large DNA viruses (NCLDV) are the largest viruses known today
45 in both genome and virion size. They have been found in most major lineages of eukaryotes
46 across diverse habitats (1–4), especially in the marine environment (5, 6). Among the
47 NCLDVs of special ecological importance are members of the family *Phycodnaviridae* that
48 infect a wide range of key algal hosts (1). These include the cosmopolitan calcifying
49 eukaryotic alga *Emiliana huxleyi* (Haptophyta), which forms massive annual blooms in the
50 oceans that have a profound impact on the carbon and sulfur biogeochemical cycles (7). *E.*
51 *huxleyi* blooms are frequently terminated by a large dsDNA virus — EhV (8), which

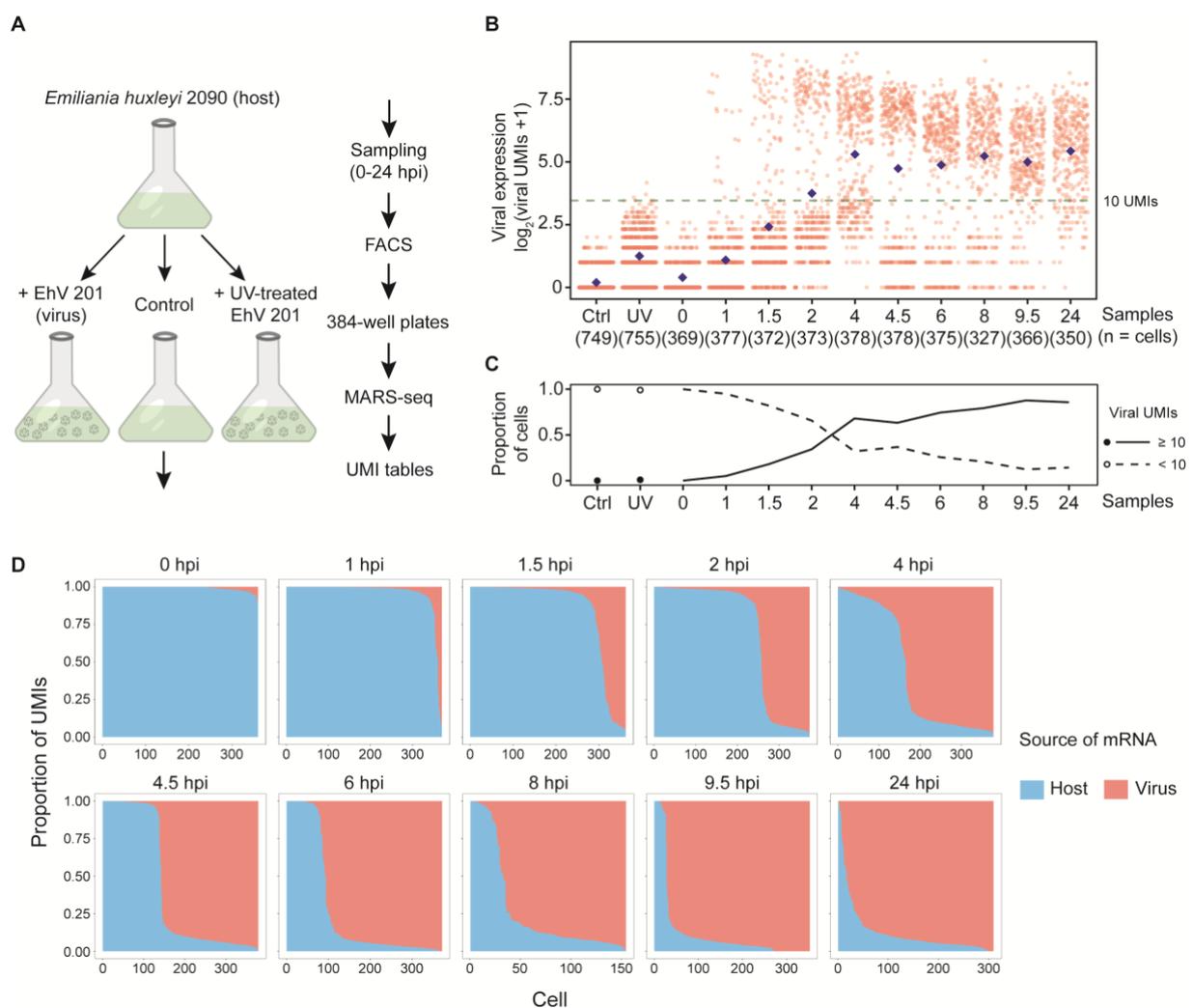
52 enhances nutrient cycling and carbon export to the deep ocean (9, 10). This host-virus model
53 provides a trackable system for understanding viral life cycle strategies and host responses.

54 High-throughput bulk RNA sequencing (RNA-seq) has been used for whole-genome
55 expression profiling of NCLDV during infection, shedding light on gene prediction,
56 transcript structure, and changes in metabolic pathways (11–13). However, bulk RNA-seq
57 profiles average gene expression levels across many cells, whereas infection states can be
58 variable among single cells. To overcome this limitation, single-cell RNA-seq (scRNA-seq)
59 approaches have been developed to probe the transcriptomes of individual cells in a highly
60 parallel manner. These methods have revolutionized our understanding of various
61 developmental and immunological processes (14, 15), including host-virus interactions in
62 mammalian systems (16, 17).

63 Here we employed scRNA-seq to study EhV infection of *E. huxleyi* at the single-cell
64 level, in order to characterize the temporal dynamics and regulation of viral and host
65 transcriptomes. *E. huxleyi* CCMP2090 cultures were infected with the lytic virus EhV201 (13,
66 18) (Fig. 1A, Table S1, and fig. S1). Individual cells were isolated into 384-well plates by
67 fluorescence-activated cell sorting (FACS) and processed using the MARS-seq protocol (19),
68 which uses cellular barcodes to tag RNA molecules from individual cells and enables
69 transcript quantification by using unique molecular identifiers (UMIs). By mapping reads to
70 the reference *E. huxleyi* transcriptome and to the EhV201 genome, we were able to profile
71 both viral and host transcripts in each individual cell.

72 We compared the total expression levels of all viral genes in 5,179 cells from infected
73 cultures, control cultures, and cultures mock-infected by UV-inactivated viruses (Fig. 1B).
74 While the average viral gene expression increased over time of infection, high variability in
75 viral expression levels between single cells was observed despite a high virion-to-cell ratio
76 (multiplicity of infection 5:1) that was used to reduce the possible variation due to encounter
77 rates. Throughout the first 10 hours of infection, two coexisting groups of cells were observed

78 — one with clear detection of active viral expression (≥ 10 viral UMIs) and the other with
 79 transcript abundance at the level of noise (cf. control or mock infection). The observed
 80 bimodal distribution of viral transcripts implies the existence of an all-or-none switch that
 81 rapidly turns on viral expression (Fig. 1B). The proportion of cells with at least 10 viral UMIs
 82 increases from 0% to nearly 90% within 10 hours (Fig. 1C). Interestingly, the mRNA pool in
 83 each individual cell was mostly dominated by either host or viral transcriptome, with very few
 84 cells in intermediate states, forming a sharp decline in host-virus transcript ratio (Fig. 1D and
 85 fig. S2), indicating EhV massively transforms the cellular transcriptome by taking over almost
 86 the entire mRNA pool (fig. S3).
 87

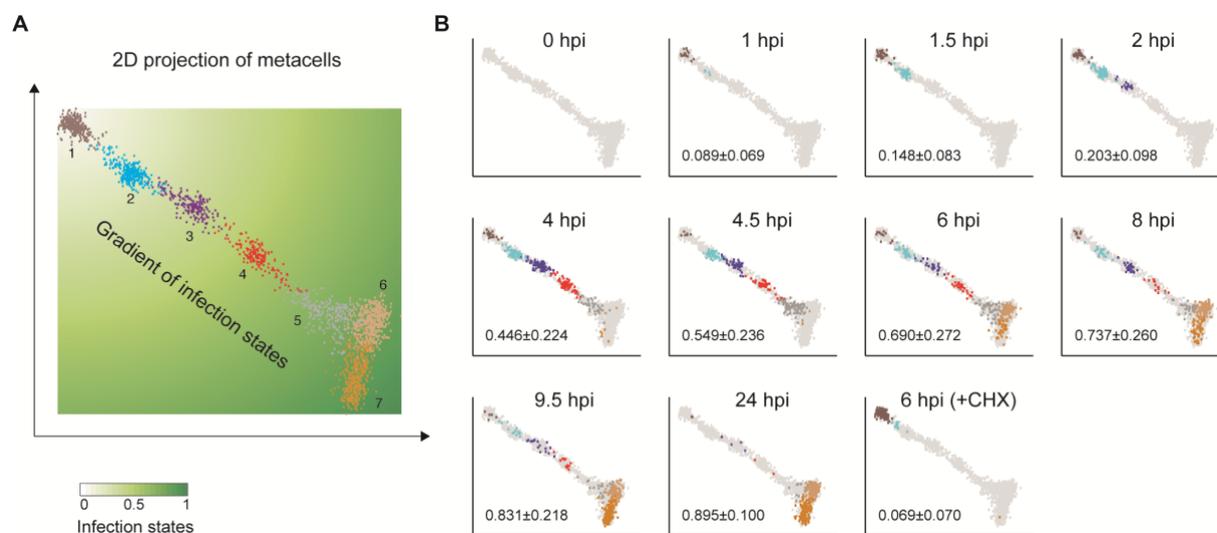


89 **Fig. 1 Transcriptome profiling of EhV infection in *E. huxleyi* at single-cell resolution.**

90 (A) The experimental setup for using single-cell RNA-seq to study EhV infection of *E.*
91 *huxleyi*. Individual cells were isolated by FACS during a time course of infection. The single-
92 cell transcriptomes were sequenced using the MARS-seq protocol, which allows absolute
93 quantification of viral and host transcripts by tagging each individual cell and each transcript
94 by cellular barcodes and unique molecular identifiers (UMIs), respectively. (B) Highly
95 heterogeneous total viral UMI counts among single cells across samples. Each red dot
96 represents an individual cell and the blue diamonds the average expression of all cells at each
97 time point. A threshold of 10 UMIs (fig. S3) highlights the bimodality of overall viral
98 expression and excludes cells only with low- or noise-level viral UMIs, as in mock infection
99 by UV-deactivated EhV (UV; 4 and 6 hpi) and control (Ctrl) samples. (C) The proportion of
100 cells with at least 10 viral UMIs (active viral expression) increased rapidly during early hours
101 of infection and approached 90% at 9.5 hpi. (D) Cliff diagrams represent the relative mRNA
102 levels of host and virus across the infection time course. The cells (columns) are sorted by the
103 relative proportions of host and viral mRNA, showing the sharp decrease in host mRNA
104 relative to viral mRNA (“cliff”) and the within-population heterogeneity through time. hpi:
105 hours post infection.

106
107 To further explore the variation among infected single cells, we applied the MetaCell
108 analysis (20) based on viral gene expression, which divided the cells into seven groups
109 (metacells) that defined the different phases of infection (Fig. 2A). Infected cells were spread
110 along an infection continuum, reaching maximum heterogeneity at 6 hpi (Fig. 2B). In
111 contrast, cells pre-treated with cycloheximide, an inhibitor of eukaryotic translation (21), were
112 mainly localized in the initial phase (metacell 1) of infection at 6 hpi (Fig. 2B). This suggests
113 that newly synthesized proteins are required for later viral transcription programs, but not for
114 the initial phase. The continuum of cells allows us to map infection progression on a

115 pseudotemporal scale between of 0 to 1 (Fig. 2A), providing a new dimension to determine
116 infection states of individual cells within heterogeneous populations.
117



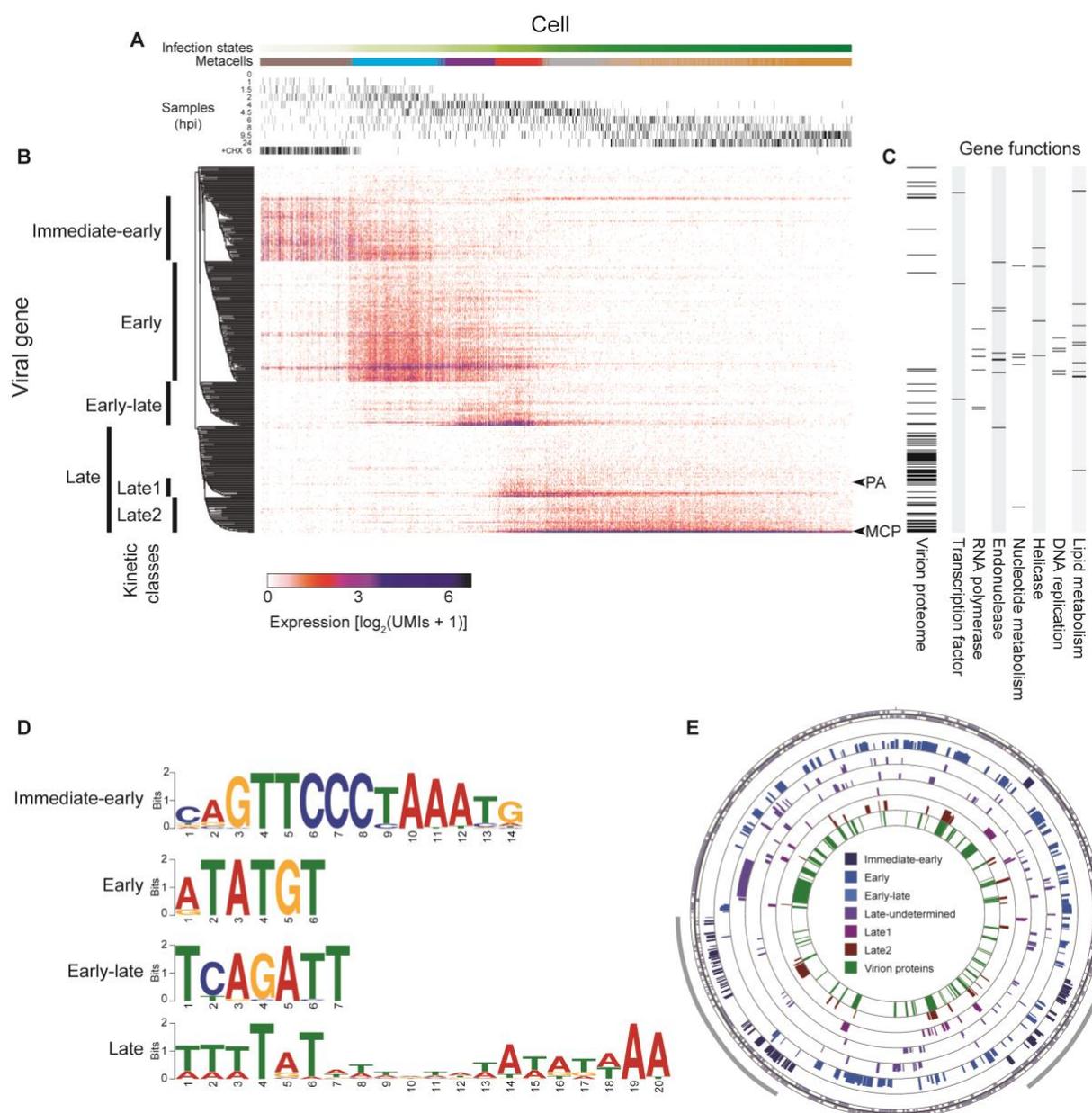
118
119 **Fig. 2 The continuum of infection states across individual cells.** (A) 2D projection of 2,072
120 single infected cells (dots) with active viral expression (≥ 10 viral UMIs) that constitute seven
121 metacells (in different colors). The cells form a continuum of infection states that can be
122 quantified by the relative distance to the initial infection state. (B) Distribution of cells from
123 each sampling time point along the infection continuum, with light gray dots representing
124 cells from the other time points. Numbers are mean \pm SD of the infection states as defined in
125 A. CHX: cycloheximide treatment before infection. The infection index (A) of an individual
126 cell was calculated as the ratio of its distance from the upper-left origin to the distance
127 between the origin and the lower-right end. A color gradient was painted based on the scale of
128 the infection indices between 0 and 1.

129
130 By re-ordering cells based on the newly defined infection states, we found that viral
131 genes were largely organized in five discrete kinetic classes of genes herein termed
132 immediate-early (73 genes), early (136), early-late (49), late 1 (22) and late 2 (40), which
133 formed partially overlapping sequential waves of expression (Fig. 3 and Table S2). The early

134 class encompasses most of the genes involved in information processing and metabolism,
135 including DNA replication (e.g., DNA polymerase delta subunit and proliferating cell nuclear
136 antigen) and lipid metabolism such as the unique EhV-encoded sphingolipid biosynthesis (13,
137 22–24) (Fig. 3C). The early genes also encode four RNA polymerase II subunits (RPB1,
138 RPB2, RPB3, and RPB6), which probably enable transcription independent of the host
139 machinery. In addition, two small RNA polymerase II subunits (RPB5 and RPB10) are
140 encoded in the early-late class. The late 1 and late 2 gene products include the packaging
141 ATPase for packing viral DNA into virions and the major capsid protein, the main component
142 of the virion capsid, respectively. In addition, our proteomic analysis showed that most
143 proteins encoded by the two late classes are integral components of EhV virions (Fig. 3C and
144 Table S3).

145 To characterize the regulatory mechanisms underlying the viral gene kinetic classes,
146 we looked for enriched sequence motifs *de novo* in the promoter regions (± 100 bp of the first
147 base of the start codon) of genes belonging to different kinetic classes (Fig. 3D and Table S4).
148 The most enriched motif in the promoter region of immediate-early genes ends mostly at the
149 ATG start codon (positions 12-14) and was previously identified as a putative promoter
150 element associated with initial expression in another EhV strain (25). We also revealed short,
151 highly enriched elements, including a 6-bp motif around the ATG of early genes and a 7-bp
152 one mostly upstream of the start codon in early-late genes. An AT-rich 20-bp motif was
153 detected, with a degenerated center, usually immediately upstream of the start codon of late 1,
154 late 2, and other late genes. We further found that these putative promoter elements are found
155 to be conserved across EhV genomes (fig. S4), which share most of their gene contents (fig.
156 S5 and Table S5). Intriguingly, the late motif bears striking similarity to the late promoter
157 element of mimivirus, a distantly related NCLDV that infects amoebae, which is comprised
158 by two AT-rich decamers separated by a degenerated tetramer and is also present in the
159 genome of the mimivirus virophage (11). Although genes in most kinetic classes are scattered

160 throughout the EhV201 genome, the immediate-early genes are mainly concentrated in two
 161 genomic regions (Fig. 3E), as also seen in the strain EhV86 (25). Such organization could
 162 facilitate the rapid activation of the immediate-early class once the host cell is infected,
 163 allowing transcription initiation of the entire cascade of the various kinetic classes. Taken
 164 together, our findings of newly defined viral kinetic classes are strongly supported by the
 165 discovery of potential regulatory mechanisms that are conserved across EhV genomes.
 166
 167



169 **Fig. 3 Viral kinetic classes, predicted promoter elements, and genomic organizations. (A)**

170 Cells (columns) are ordered by the pseudo-temporal scale (infection states) as defined in Fig.

171 2, with hours post infection (hpi) of each culture sample marked for each cell. **(B)** Viral genes

172 (rows) are divided into kinetic classes based on the hierarchical clustering of their expression

173 profiles across the individual cells. **(C)** Presence in the EhV201 virion proteome, as well as

174 categories of predicted functions, are indicated for each gene. PA: packing ATPase; MCP:

175 major capsid protein. **(D)** Enriched motifs in the promoter regions of viral genes in the newly

176 defined kinetic classes. **(E)** Distribution of genes in each kinetic class (circle) on the

177 assembled genome of EhV201. The outermost circle indicates the direction of transcription.

178 The heights of genes in each kinetic class are proportional to their log₂-transformed

179 expression levels. The innermost circle marks the presence in the proteome of EhV201

180 virions. Immediate-early genes are mostly clustered in two genomic regions (gray arcs).

181

182 Finally, our single-cell dual RNA-seq approach allows quantification of not only

183 transcripts encoded by the viral genome, but also respective host transcriptomes including

184 those of the mitochondria and chloroplasts (Fig. 4). As depicted in Fig. 1, there was a sharp

185 shift of the cellular transcriptome from host-encoded to virus-encoded transcripts. By

186 mapping the expression levels across infection states, we found differential shutdown of

187 nucleus-encoded and organelle-encoded genes (Fig. 4A). Whereas nuclear transcripts were

188 rapidly shut down at the onset of active viral expression, levels of the mitochondrial and

189 chloroplast transcripts were maintained. Moreover, the two types of organelles show distinct

190 expression patterns, with mitochondrial transcripts persisting at lower, yet relatively stable

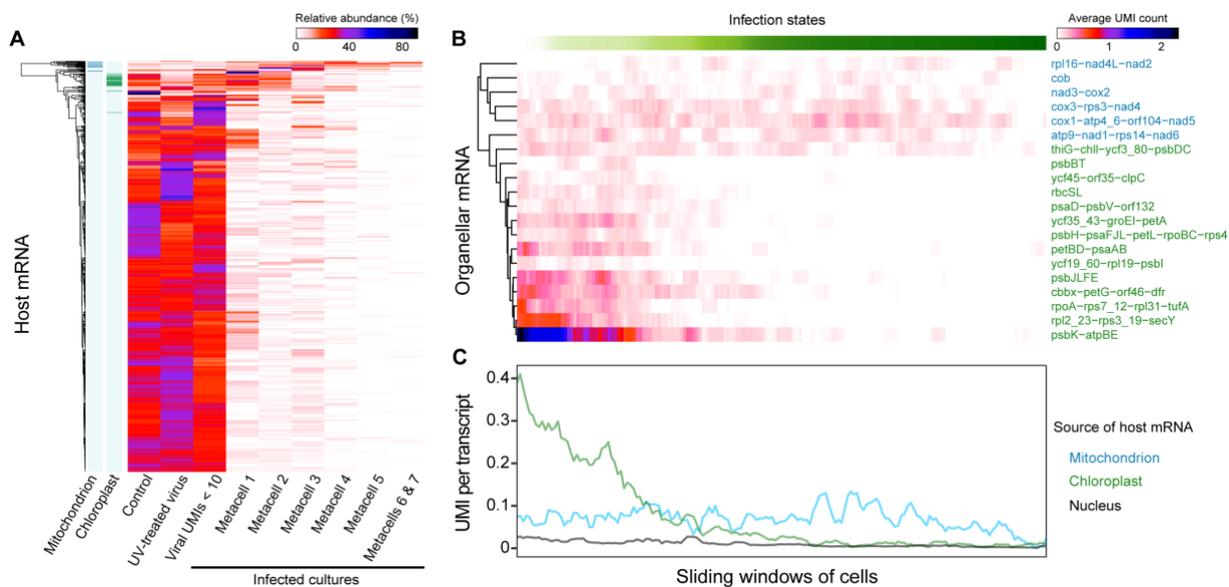
191 levels throughout the infection progression and chloroplast transcript levels that were high in

192 the beginning but then soon declined (Fig. 4B and 4C). The delayed shutdown of chloroplast-

193 encoded genes suggests the need for active chloroplasts in early stages of infection, which is

194 in agreement with the recent finding that photosynthetic electron transport remained active

195 during the initial phase of viral infection (26). In contrast to chloroplasts, the mitochondrial
 196 transcripts were maintained throughout the infection progression, consistent with the high
 197 requirement for energy and redox equivalents at all stages of infection.
 198



199
 200 **Fig. 4 Highly resolved dynamics of virus-induced shutoff of host mRNA.** (A) Levels of
 201 host transcripts across categories of cells. Single cells from cultures infected with EhV are
 202 divided based on overall viral UMI counts and metacell grouping (Fig. 2A). Protein-encoding
 203 genes with at least 100 UMIs were included. For each gene, the average UMI count was
 204 calculated for each category, and the sum of the average values of all categories was
 205 normalized to 100. Transcripts were then hierarchically clustered based on the normalized
 206 values. Organelle-encoded transcripts are indicated on the left. (B and C) Abundance of
 207 organellar transcripts along the infection progression. Cells with at least 10 viral UMIs are
 208 ordered according to the infection states as determined by the MetaCell analysis (Figs. 2 and
 209 3). Sliding window averages were calculated for each transcript (B) or mean of all transcripts
 210 (C) with a window size of 50 cells and steps of 10.

211
 212 Single-cell RNA-seq has proven to be a very powerful tool for comprehensive
 213 understanding of complex host-virus dynamics, disentangling cells from different sampling

214 times and infection states. Using this approach, we dissected the transcriptomic dynamics and
215 regulation of viral infection in a globally important microalga and showed rapid selective
216 shutdown of host genes induced by an NCLDV. The comprehensive high-resolution mapping
217 of infection states and kinetic classes allowed identification of conserved regulatory
218 sequences which likely orchestrate the sequential expression programs during the viral life
219 cycle. The applicability of scRNA-seq provides new resolution to track host-virus dynamics
220 and holds great potential for *in situ* quantification of active viral infection in natural
221 populations as well as identification of additional ecologically important host-virus systems.
222 This in turn will facilitate the assessment of the impact of viruses on the function and
223 composition of microbial food webs in the marine environment.

224

225

226 **References and Notes:**

227

- 228 1. J. L. Van Etten, R. H. Meints, *Annu. Rev. Microbiol.* **53**, 447–94 (1999).
- 229 2. N. Yutin, Y. I. Wolf, D. Raoult, E. V Koonin, *Viol. J.* **6**, 223 (2009).
- 230 3. S. W. Wilhelm, S. R. Coy, E. R. Gann, M. Moniruzzaman, J. M. A. Stough, *PLOS*
231 *Pathog.* **12**, e1005752 (2016).
- 232 4. F. Schulz *et al.*, *Science (80-.)*. **356**, 82–85 (2017).
- 233 5. Q. Carradec *et al.*, *Nat. Commun.* **9** (2018), doi:10.1038/s41467-017-02342-1.
- 234 6. A. C. Gregory *et al.*, *Cell.* **177**, 1–15 (2019).
- 235 7. C. W. Brown, J. A. Yoder, *J. Geophys. Res.* **99** (1994), p. 7467.
- 236 8. W. H. Wilson *et al.*, *J. Mar. Biol. Assoc. United Kingdom*, 369–377 (2002).
- 237 9. C. P. Laber *et al.*, *Nat. Microbiol.* **3**, 537–547 (2018).
- 238 10. U. Sheyn *et al.*, *ISME J.* (2018), pp. 1–10.
- 239 11. M. Legendre *et al.*, *Genome Res.* **20**, 664–674 (2010).

- 240 12. G. Blanc *et al.*, *PLoS One*. **9**, 1–10 (2014).
- 241 13. S. Rosenwasser *et al.*, *Plant Cell*. **26**, 2689–2707 (2014).
- 242 14. A. A. Kolodziejczyk, J. K. Kim, V. Svensson, J. C. Marioni, S. A. Teichmann, *Mol.*
243 *Cell*. **58**, 610–620 (2015).
- 244 15. A. Tanay, A. Regev, *Nature*. **541**, 331–338 (2017).
- 245 16. F. Zanini *et al.*, *Proc. Natl. Acad. Sci.* **115**, E12363–E12369 (2018).
- 246 17. A. B. Russell, C. Trapnell, J. D. Bloom, *Elife*. **7**, e32303 (2018).
- 247 18. U. Sheyn, S. Rosenwasser, S. Ben-Dor, Z. Porat, A. Vardi, *ISME J.*, 1–13 (2016).
- 248 19. H. Keren-Shaul *et al.*, *Nat. Protoc.*, 1 (2019).
- 249 20. Y. Baran *et al.*, 1–43 (2018).
- 250 21. T. Schneider-Poetsch *et al.*, *Nat. Chem. Biol.* **6**, 209–217 (2010).
- 251 22. W. H. Wilson *et al.*, *Science (80-.)*. **309**, 1090–1092 (2005).
- 252 23. A. Vardi *et al.*, *Science (80-.)*. **326**, 861–865 (2009).
- 253 24. C. Ziv *et al.*, *Proc. Natl. Acad. Sci.*, 201523168 (2016).
- 254 25. M. J. Allen *et al.*, *J. Virol.* **80**, 7699–7705 (2006).
- 255 26. K. Thamatrakoln *et al.*, *New Phytol.* **221**, 1289–1302 (2019).
- 256 27. M. D. Keller, R. C. Seluin, W. Claus, R. R. L. Guillard, *J. Phycol.* **23**, 633–638 (1987).
- 257 28. M. J. Allen, J. Martinez-Martinez, D. C. Schroeder, P. J. Somerfield, W. H. Wilson,
258 *Environ. Microbiol.* **9**, 971–982 (2007).
- 259 29. G. Schleyer *et al.*, *Nat. Microbiol.* **4**, 527–538 (2019).
- 260 30. J. I. Nissimov *et al.*, *J. Virol.* **86**, 2380–2381 (2012).
- 261 31. D. Schatz *et al.*, *Nat. Microbiol.* **2**, 1485–1492 (2017).
- 262 32. C. Camacho *et al.*, *BMC Bioinformatics.* **10**, 421 (2009).
- 263 33. M. V Sanchez-puerta, T. R. Bachvaroff, C. F. Delwiche, *DNA Res.* **156**, 151–156
264 (2005).
- 265 34. M. V Sanchez-puerta, *DNA Res.* **11**, 1–10 (2004).

- 266 35. N. Barak-Gavish *et al.*, *Sci. Adv.* (2018).
- 267 36. B. A. Read *et al.*, *Nature*. **499**, 209–13 (2013).
- 268 37. M. V Han, C. M. Zmasek, *BMC Bioinformatics*. **10**, 356 (2009).
- 269 38. Y. Mirzakhanyan, P. D. Gershon, *Microbiol. Mol. Biol. Rev.* **81**, e00010-17 (2017).
- 270 39. E. Assarsson *et al.*, *Proc. Natl. Acad. Sci. U. S. A.* **105**, 2140–5 (2008).
- 271 40. *Nucleic Acids Res.* **44**, D7–D19 (2016).
- 272 41. E. Afgan *et al.*, *Nucleic Acids Res.* (2018), doi:10.1093/nar/gky379.
- 273 42. T. L. Bailey, N. Williams, C. Misleh, W. W. Li, *Nucleic Acids Res.* **34**, 369–373
- 274 (2006).
- 275 43. T. L. Bailey, *Bioinformatics* (2011), doi:10.1093/bioinformatics/btr261.
- 276 44. J. R. Grant, A. S. Arantes, P. Stothard, *BMC Genomics* (2012), doi:10.1186/1471-
- 277 2164-13-202.
- 278 45. J. I. Nissimov *et al.*, *Viruses*. **9**, 52 (2017).

279

280

281 **Acknowledgements**

282 We thank N. Stern-Ginossar, A. Nachshon, M. Shnayder, and E. Koh for help with MARS-

283 seq and cycloheximide experiments; R. Avraham, D. Hoffman, S. H. Avivi, G. Rosenberg, A.

284 Solomon, and Advanced Sequencing Technologies Unit of the Life Science Core Facilities

285 for assistance with Illumina sequencing; Y. Levin and M. Kupervaser from the de Botton

286 Institute for Protein Profiling of the Nancy and Stephen Grand Israel National Center for

287 Personalized Medicine for assistance with proteomic analyses; and all Vardi lab members for

288 discussion. **Funding:** This work was supported by European Research Council (ERC) CoG

289 (VIROCELLSPHERE grant # 681715; A.V.) and EMBO Long-Term Fellowship (ALTF

290 1172-2016; C.K.). **Author contributions:** C.K., U.S., S.R. and A.V. designed the study, with

291 input from A.S.-P. and A.T.; U.S. and C.K. prepared the cell cultures; C.K. and U.S.

292 performed MARS-seq with assistance of A.S.-P.; C.K. analyzed the MARS-seq data with
293 assistance of A.S.-P. and A.T.; D.S. performed the virion proteome experiment; S.B.-D.
294 analyzed the motifs and compared the EhV genomes, with input from C.K.; C.K. prepared the
295 figures, with input from U.S., S.B.-D., S.R. and A.V.; C.K., S.R. and A.V. wrote the
296 manuscript, with input from all authors; A.V. supervised the study. **Competing interests:** All
297 authors declare that they have no competing interests. **Data and materials availability:** All
298 data was deposited in Gene Expression Omnibus with the accession number GSE135429.

299

300 **Supplementary materials:**

301 Materials and Methods

302 Figures S1-S5

303 Tables S1-S5

304 References (27-45)

305

306

Supplementary Materials for

Tracking infection dynamics at single cell resolution reveals highly resolved expression programs of a large marine virus

Chuan Ku, Uri Sheyn, Arnau Sebé-Pedrós, Shifra Ben-Dor, Daniella Schatz, Amos Tanay, Shilo Rosenwasser*, Assaf Vardi*

Correspondence to: shilo.rosenwaser@mail.huji.ac.il; assaf.vardi@weizmann.ac.il

This PDF file includes:

Materials and Methods

Figs. S1 to S5

Other Supplementary Materials for this manuscript include the following:

Tables S1 to S5

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 **Materials and Methods**

46 Algal culture and viral stock

47 The coccolithophore *Emiliana huxleyi* strain CCMP2090 was grown in K/2 (27)
48 filtered seawater medium at 18 °C with an irradiance of 100 $\mu\text{mol photons m}^{-2}\text{s}^{-1}$ and a 16:8 h
49 light-dark cycle. The viral strain EhV201 (28) was propagated on CCMP2090 culture until the
50 culture lysed to clearance, and viral lysate was filtered using 0.45- μm PVDF Stericup-HV
51 vacuum filtration system (Merck). We used the plaque assay as previously described (29) to
52 quantify infectious virions, and filtered viral lysate generally contained $\sim 10^8$ infectious virions
53 per ml. Cultures at the exponential phase ($\sim 10^6$ cells ml^{-1}) were used for viral infection at a
54 high multiplicity of infection (5 infectious virions to 1 cell) at 2 h after the light cycle began.
55 Algal cells and virion particles were enumerated using an Eclipse Flow Cytometer Analyzer
56 (iCyt) during an infection time course (fig. S1). UV-inactivated virus was obtained by
57 exposing viral lysate in a plastic petri dish to UV in a transilluminator system for 15 minutes.
58 For inhibition of protein synthesis, cycloheximide was added to culture for a final
59 concentration of 1 $\mu\text{g}/\text{mL}$ immediately before viral infection.

60

61 Single-cell isolation by fluorescence-activated cell sorting

62 Single cells were isolated into individual wells in 384-well plates on a FACSAria II or
63 III sorter (BD) using a 488-nm laser for excitation and a nozzle size of 100 μm . Cells were
64 identified based on chlorophyll red autofluorescence (663–737 nm) and separated under the
65 “Single Cell” mode for maximal purity.

66

67 Massively parallel single-cell RNA sequencing

68 The scRNA-seq method is based on the MARS-seq2.0 protocol (19). Before sorting,
69 each 384-well plate was prepared using a Bravo liquid handling platform (Agilent) to transfer
70 into each well 2 μl lysis buffer containing 2 or 4 nM reverse transcription (RT) primers with
71 cellular and molecular (UMI) barcodes. Plates with sorted cells were immediately frozen on
72 dry ice and later stored at -80 °C. After thawing, the plates were heated to 95 °C for 3 min to
73 evaporate the lysis buffer, cooled, added RT reagent mix, and placed in a Labcycler
74 (SensoQuest) for RT reaction. The unused primers were then degraded with Exonuclease I
75 (NEB). Complementary DNA from each set of wells with 384 (one pool per plate) or 192
76 (two pools per plate) unique cellular barcodes were then combined into pools, which

77 underwent second strand synthesis and *in vitro* transcription to amplify the sequences linearly.
78 The RNA products were fragmented and ligated with single-stranded DNA adapters
79 containing pool-specific barcodes. A second round of RT was carried out, followed by
80 polymerase chain reaction with primers containing Illumina adapters to construct DNA
81 libraries for paired-end sequencing on a NextSeq or MiniSeq machine (Illumina).

82

83 Read processing and mapping to reference sequences

84 The fastq files were processed using the analytical pipeline of MARS-seq2.0 (19)
85 which mapped the reads to viral and host reference sequences and demultiplexed them based
86 on the pool, cellular, and molecular barcodes. At least 4 reads per UMI were obtained per cell.
87 For the virus, the predicted coding sequences (CDS) in the EhV201 genome sequence
88 (JF974311) (30) were used as reference. For the host, an integrated transcriptome reference of
89 *E. huxleyi* (31) was used, where both forward and reverse (labeled with “_rev”) strands of
90 each sequence were included for mapping because the CDS directionality was often
91 unknown. In addition to nuclear sequences, the transcriptome reference contains chloroplast
92 and mitochondrial transcripts, which were identified by BLAST (32) searches against the
93 respective organellar genome sequences (33, 34).

94

95 Viral and host transcript abundance

96 To avoid empty wells or those with low transcript capture during library preparation,
97 wells with fewer than 10 UMIs were removed for all downstream analyses. Wells with more
98 than 1,000 viral transcripts were also removed to prevent wells with doublet or multiplet cells
99 from confounding single-cell infection state analyses. The vast majority (> 99%) of control or
100 infected *E. huxleyi* cells had fewer than 1,000 total UMIs, which reflects their low RNA
101 content due to the small cell size (3-5 μm in diameter for naked cells (35)).

102

103 Cliff diagrams of viral and host mRNA

104 A cliff diagram was plotted for each time point post infection (Fig. 1D). The total viral
105 mRNA transcript counts were calculated by summing up all UMIs assigned to virus-encoded
106 CDS. The total host mRNA counts were all host UMIs except for those assigned to the three
107 transcript sequences corresponding to nuclear, mitochondrial, and chloroplast ribosomal RNA
108 genes. For each sampling time point, the cells were ordered by the ratio of host to all

109 (host+virus) mRNA from high to low. Cells with fewer than 20 total UMIs were excluded
110 from the plots to avoid sampling biases.

111

112 Doublet assay and in silico simulations

113 The EhV201 viral genome has a much lower GC content (40.46%) (30) than the *E.*
114 *huxleyi* host genome (65%) (36). This difference in GC content and other factors might bias
115 the capture and amplification of viral and host transcripts during MARS-seq and might lead to
116 the either-or dominance of host and viral mRNA (Fig. 1D). To test if the presence of viral
117 mRNA has effects on the quantification of host mRNA transcripts and vice versa, we
118 conducted a doublet assay. Half of the doublet assay plate consisted of single cells sorted
119 from infected culture at 8-hpi and the other half plate consisted of two cells per well of which
120 one cell was sorted from infected culture at 8 hpi and the other cells from control culture
121 (Table S1). Cliff diagrams were then plotted for each half of the plate (fig. S2A-B). If the
122 viral mRNA does not interfere with host mRNA quantification, wells with cells isolated from
123 both the infected and the control cultures should show host and viral mRNA UMI counts that
124 are additive, as if the two cellular transcriptomes were sequenced separately and merged
125 together. To simulate this *in silico*, we randomly sampled sequenced single-cell
126 transcriptomes from the control culture plate with the closest timepoint to 8 hpi (Table S1)
127 using the *sample* function in R, and they were combined one-by-one with those from the 8-hpi
128 half plate with single cells. Cliff diagrams were plotted for each of the eight simulations (fig.
129 S2C-J).

130

131 MetaCell analysis of infected cells

132 The MetaCell package is an unbiased approach for cell grouping based on single-cell
133 expression profiles without enforcing any global structure (20). To account for the small
134 number of UMIs per cell and to use more viral genes for grouping infected cells, we
135 employed more inclusive criteria for gene marker selection ($T_{tot}=20$, $T_{top3}=1$, $T_{szcor}=-$
136 0.01 , $T_{vm}=0.2$, and $T_{niche}=0.05$). A total of 179 marker genes were used constructing *k*-
137 nearest neighbours graphs with $K=150$, followed by co-clustering with bootstrapping based
138 on 1,000 iterations of resampling 75% of the cells and an approximated target minimum
139 metacell size of 80. Unbalanced edges were filtered with $K=40$ and $\alpha=3$.

140 After plotting single cells in a 2D projection and identifying the directionality of
141 infection progression (Fig. 2), we quantified the infection states of individual cells based the
142 relative distance to the beginning of the progression (i.e., upper-left). The *ad hoc* origin was
143 defined as having x- and y-coordinates of the leftmost and the highest cells, respectively, and
144 the *ad hoc* end as having those of the rightmost and the lowest cells, respectively. The
145 infection index of an individual cell was calculated as the ratio of its distance from the origin
146 to the distance between the origin and the end. A color gradient was painted based on the
147 scale of the infection indices between 0 and 1.

148

149 Hierarchical clustering of genes

150 We used hierarchical clustering to group viral genes with similar expression patterns
151 across infected cells. Viral genes with fewer than 20 UMIs in total across all the cells with
152 active viral expression were excluded to avoid spurious clusters or groupings that are not well
153 supported. The pheatmap package in R was used for hierarchical clustering (hclust) of UMI
154 counts in log scale, using Pearson correlation as distance measure and the “average”
155 agglomeration method. The cluster tree of viral genes was manually curated by branch
156 swapping using Archaeopteryx (37) to order the genes by kinetic classes and by expression
157 levels.

158

159 Annotations of EhV201 genome

160 A comprehensive annotation table was prepared for EhV201 genes by ordering them
161 according to the manually curated cluster tree (Fig. 3B). The sequence descriptions were
162 based on the GenBank record (JF974311), with modifications based on relevant publications
163 (38, 39), BLAST searches against the nr database (40), and the Nucleo-Cytoplasmic Virus
164 Orthologous Groups (NCVOGs) (2).

165

166 Virion proteome composition

167 Purified EhV201 virion samples were lysed using an 8 M urea buffer and subjected to
168 tryptic digestion followed by LC-MS analysis on a Q Exactive HF instrument (Thermo). We
169 used Byonic version 2.10.5 (Protein Metrics) as the software to process the data and to search
170 against the viral CDS database (Table S3). In total 76 viral CDSs had at least 10 peptide-
171 spectrum matches and were shown as present in virion proteome in Figure 3.

172

173 Sequence motif analyses of viral genomes

174 A 200-bp region surrounding the predicted start codon ATG (A of ATG is base 101)
175 was extracted for each CDS in the EhV201 genome using the Extract Genomic DNA tool of
176 Galaxy (41). Promoter motif discovery was performed on the positive strand of these regions
177 for genes in different kinetic classes using MEME version 5.0.5 (42), Genomatix Genome
178 Analyzer CoreSearch (Intrexon Bioinformatics), and Improbizer
179 (<https://users.soe.ucsc.edu/~kent/improbizer/improbizer.html>), with the number of motifs to
180 detect set to 4. While all three programs agreed on the motifs detected, the results of the MEME
181 analyses are used in this study, where the most common motifs of lengths 6-20 are shown (Fig.
182 3D, fig. S4, and Table S4). To test if the short 6-bp sequence motif in early genes is significantly
183 enriched, a differential analysis was performed for early genes against all late (late 1, late 2 and
184 late-undetermined combined) using DREME (43) (Table S4).

185

186 Circular visualization of viral genes

187 The UMI counts of individual viral genes were summed over all cells showing active
188 viral expression (≥ 10 viral UMIs). The \log_2 values of these total UMI counts were visualized
189 using CCT (44), with a ring for each kinetic class and the height proportional to the log-
190 transformed expression value. Another ring shows the presence/absence of each gene product
191 in the virion proteome (Fig. 3E).

192

193 Comparative analyses of other EhV genomes

194 The best BLAST hits in the CCT analyses were used to find the homologous proteins
195 of the various kinetic classes in other EhV genomes. DNA surrounding the ATG of those
196 genes was extracted for each strain with the GetFastaBed tool in Galaxy (41). Promoter
197 analyses were run on the sequences of the promoter regions of three additional strains
198 (EhV207, EhV86 and EhV202) individually as described above and with all four together
199 (fig. S4 and Table S4). These three strains were chosen as representatives of three major
200 clades of EhV strains (45). A CCT plot was drawn to show the sequence similarity between
201 EhV201 and all 12 other EhV strains with whole genome sequences (fig. S5).

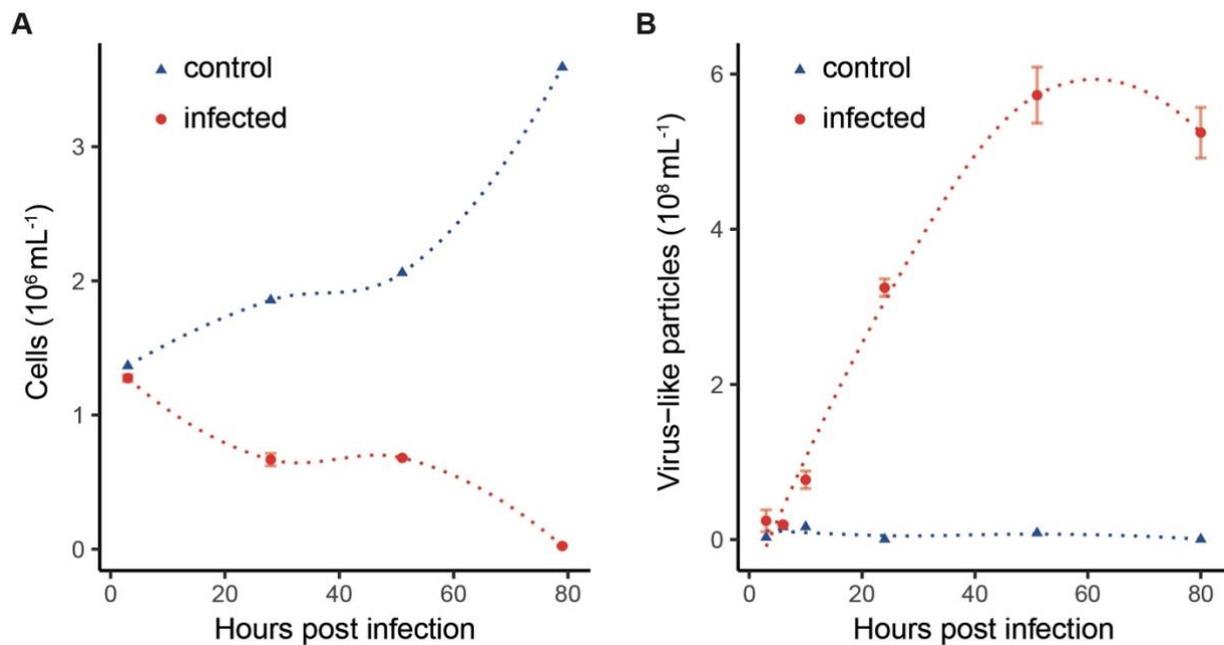
202

203 Nuclear and organellar transcripts

204 Host gene expression patterns were analyzed for different categories of single cells,
205 including control culture, mock infected culture, infected culture with fewer than 10 viral
206 UMIs, and different metacells of infected cells with at least 10 viral UMIs. Transcript
207 sequences with fewer than 100 total UMIs across all the cells were excluded to avoid lowly
208 expressed ones, leaving 720 unique transcripts in total (Fig. 4A). The average UMI count was
209 calculated for each transcript in each category, and the sum of the average values of all
210 categories was normalized to 100. Hierarchical clustering of host transcripts was done using
211 heatmap as described above based on the normalized values. The sliding windows (window
212 size = 50 cells, step = 10 cells) for organellar transcripts were generated using the rollapply
213 function of the zoo package in R.

214

215



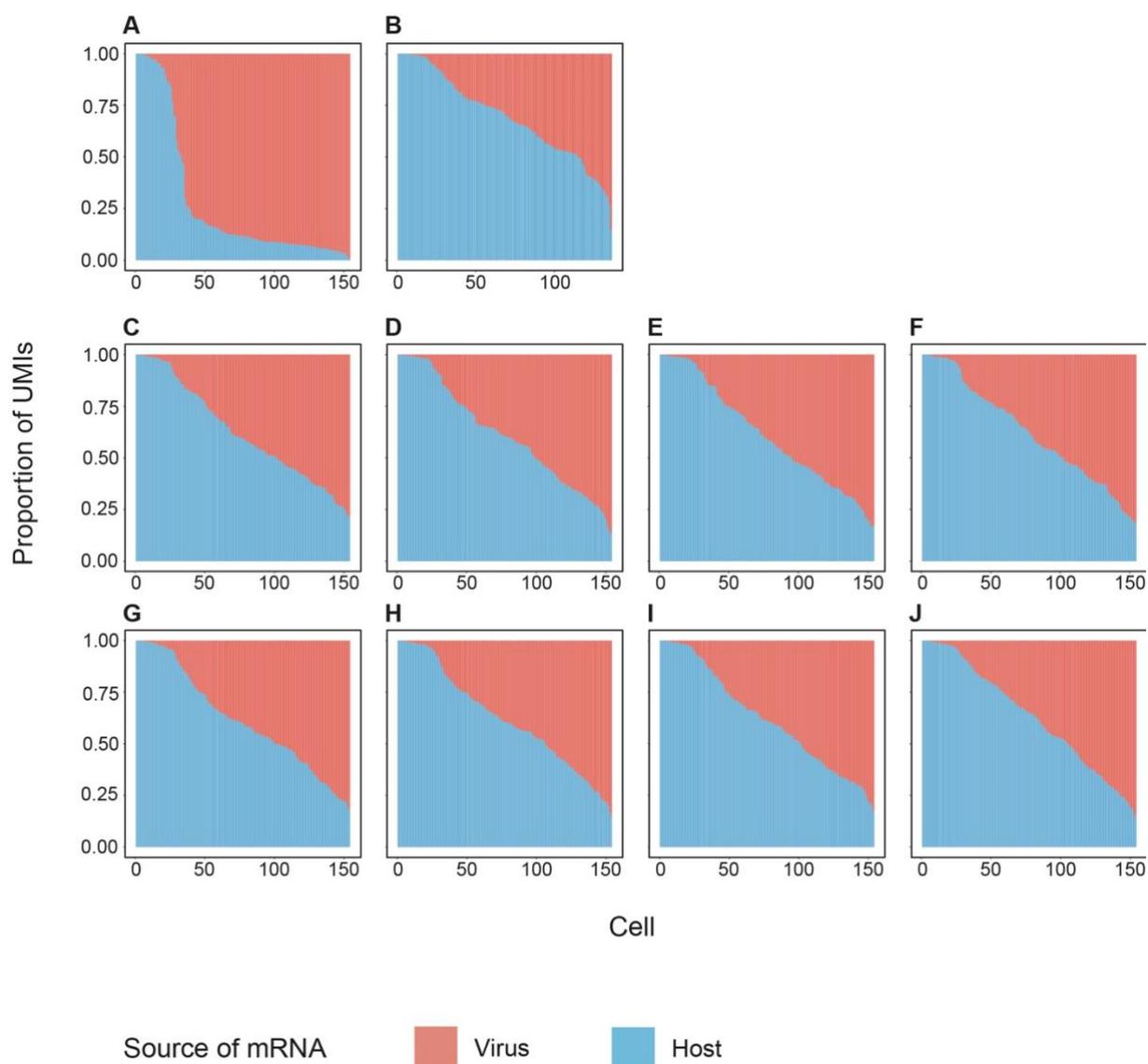
216

217 **Fig. S1.**

218 **Abundance of *E. huxleyi* cells and extracellular EhV201 virions during a time course of**
219 **infection.** (A) Numbers of *E. huxleyi* cells in control and infected cultures. (B) Numbers of
220 virus-like particles cells in control and infected cultures. A line was plotted for each treatment
221 using the locally estimated scatterplot smoothing method (span = 1) implemented in the
222 ggplot geom_smooth function in R. The infected culture had two biological replicates, and
223 mean \pm SD is shown.

224

225



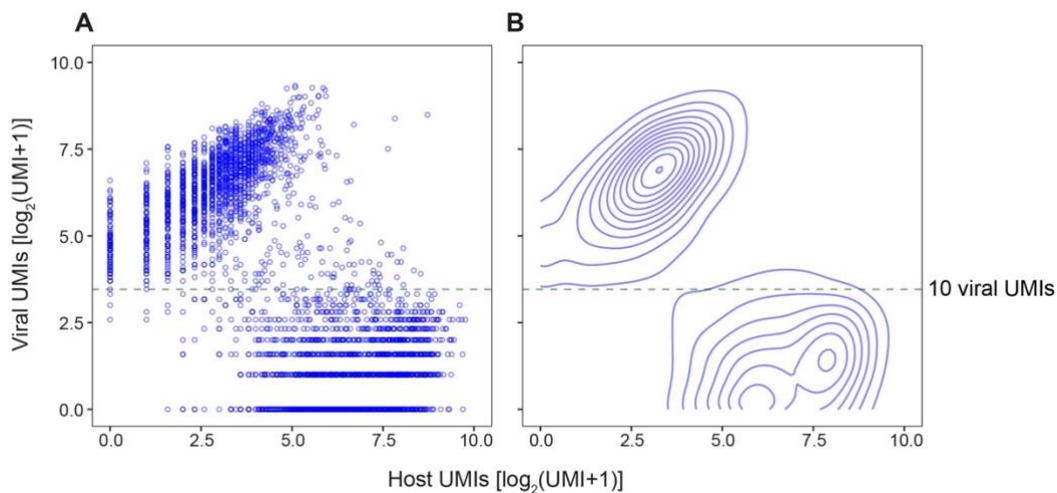
227 **Fig. S2.**

228 **The presence of viral mRNA does not affect the capture and sequencing of host mRNA**
229 **transcripts by MARS-seq.** To test if the shutdown of host mRNA (Fig. 1D) is not due to a
230 possible bias of MARS-seq toward viral transcripts given the differences in GC% content and
231 other properties between viral and host genes, a doublet assay was conducted. Cells at 8 hpi
232 were sorted and viral and host transcript quantifications were compared between wells with a
233 single cell from infected culture (half a 384-well plate) and wells with one from infected and
234 one from control cultures (the other half of the plate) (Table S1). Cliff diagrams are used to
235 represent viral and host proportions of the mRNA pool in single cells (Fig. 1D). (A) Wells
236 with one cell sorted from the infected culture at 8 hpi. (B) Wells with one cell sorted from the
237 infected culture and one from the control culture at 8 hpi. (C–J) Cliff diagrams based on four
238 *in silico* simulations, where the same number of cells from control plate 10 hpi were sampled
239 and randomly paired with the single cells from the infected culture at 8 hpi (i.e., A), show

240 patterns similar to the experimental results (**B**). This observation indicates that the presence of
241 viral molecules did not affect the quantification of host transcripts and that the absolute counts
242 of these transcripts are additive. For the doublet experiment (**B**) or simulations (**C–J**), a “cliff”
243 (sharp change in mRNA ratios) is absent; instead, a “slope” (gradual change) is observed
244 because of random pairings of infected and control culture cells of different cell size and
245 transcript content.

246

247



248

249 **Fig. S3.**

250 **Abundance of viral and host mRNA across single cells in EhV-infected culture of *E.***

251 ***huxleyi*.** (A) For each cell from infected samples and with total UMIs between 20 and 1,000,

252 the numbers of viral and host mRNA UMIs are plotted. (B) A contour plot of A based on 2D

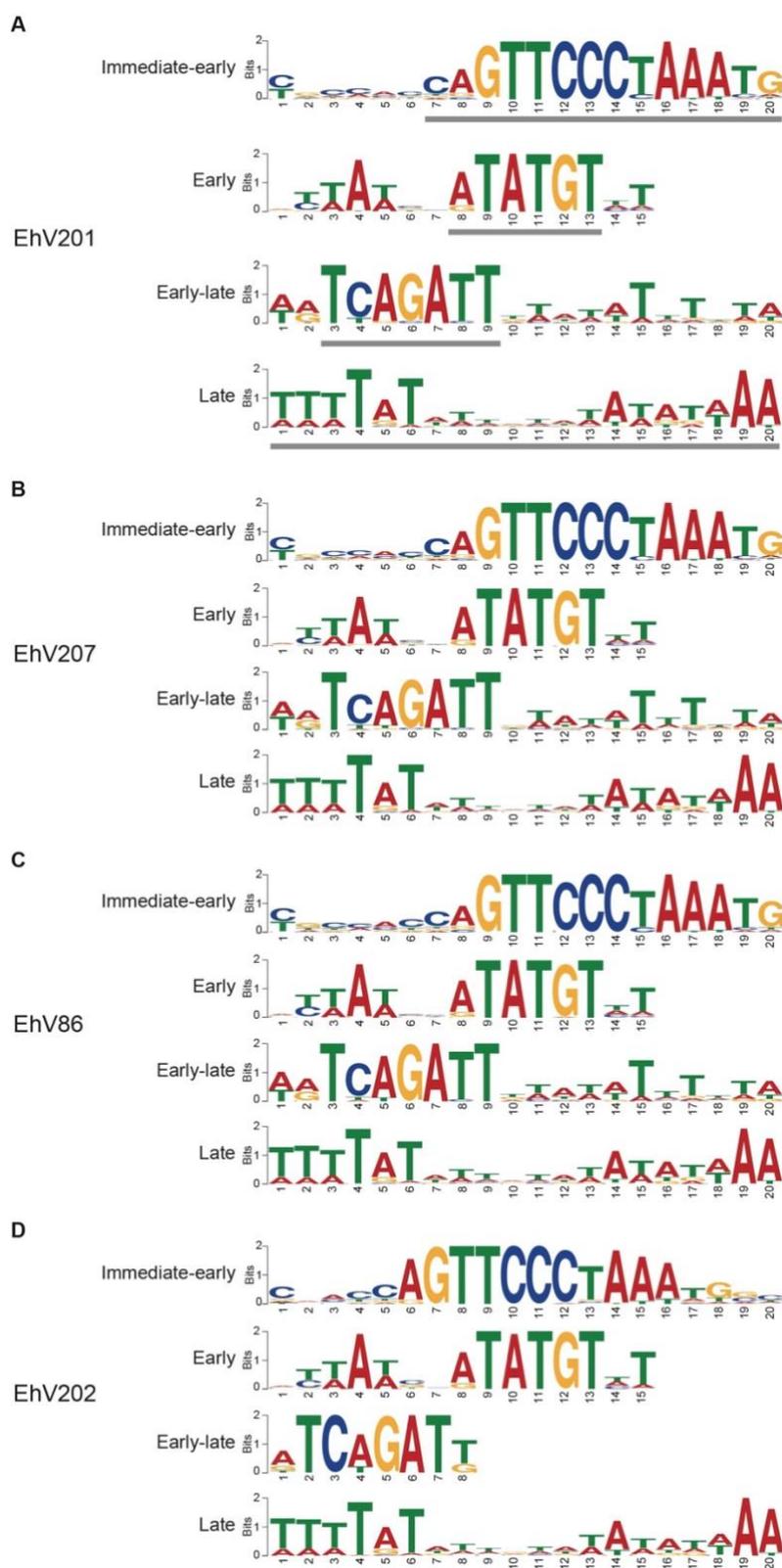
253 kernel density estimation using the `geom_density_2d` function of `ggplot2`. A threshold of 10

254 viral UMIs filters out wells with background or noise levels of viral UMIs and separates cells

255 with and without host shutdown.

256

257



258

259 **Fig. S4.**

260 **Promoter elements specific to EhV201 kinetic classes are conserved across EhV strains.**

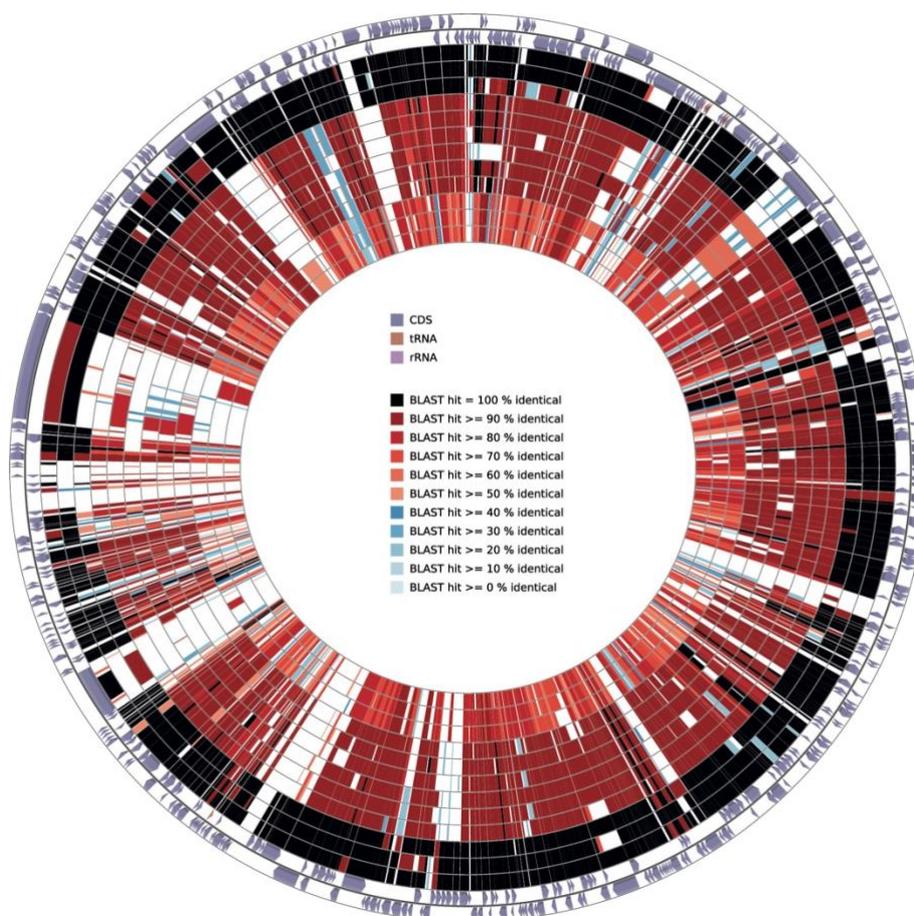
261 (A) Enriched sequence motifs (6-20 bp) found in the promoter regions (± 100 bp of the first
262 base of the start codon) of genes in different kinetic classes of EhV201. The gray bars indicate
263 the highly enriched motifs shown in Fig. 3D. (B-D) Enriched sequence motifs in the promoter

264 regions of genes in three other EhV strains that are homologous to the EhV201 genes in
265 different kinetic classes. The three strains represent the three major clades of EhV strains (45)
266 (Table S4) and are ordered with increasing divergence from EhV201. The numbers on the x-
267 axis indicate the positions within each motif.

268

269

270



271

272 **Fig. S5.**

273 **The EhV201 genes are largely conserved in 12 other EhV genomes.** In this similarity
274 circular plot, the two outermost circles indicates the direction of transcription and types of
275 genes. The other circles show the BLAST similarity of protein sequences between EhV201
276 and the following strains (from outer to inner): EhV207, EhV203, EhV208, EhV88, EhV86,
277 EhV84, EhV99B1, EhV164, EhV145, EhV202, EhV156, and EhV18. See also Table S5.

278

279